

NeuralSCF: Neural network self-consistent fields for density functional theory

Feitong Song¹ and Ji Feng^{1,2,*}

¹International Center for Quantum Materials, School of Physics, Peking University, Beijing 100871, China

²Hefei National Laboratory, Hefei 230088, China

Kohn-Sham density functional theory (KS-DFT) has found widespread application in accurate electronic structure calculations. However, it can be computationally demanding especially for large-scale simulations, motivating recent efforts toward its machine-learning (ML) acceleration. We propose a neural network self-consistent fields (NeuralSCF) framework that establishes the Kohn-Sham density map as a deep learning objective, which encodes the mechanics of the Kohn-Sham equations. Modeling this map with an SE(3)-equivariant graph transformer, NeuralSCF emulates the Kohn-Sham self-consistent iterations to obtain electron densities, from which other properties can be derived. NeuralSCF achieves state-of-the-art accuracy in electron density prediction and derived properties, featuring exceptional zero-shot generalization to a remarkable range of out-of-distribution systems. NeuralSCF reveals that learning from KS-DFT’s intrinsic mechanics significantly enhances the model’s accuracy and transferability, offering a promising stepping stone for accelerating electronic structure calculations through mechanics learning.

I. Introduction

Density functional theory (DFT) [1] is an *ab initio* computational method for investigating the electronic structure of matter. The Kohn-Sham (KS) formulation of DFT [2], by introducing reasonable approximations to the exchange-correlation (XC) functional, has brought this theory into practical use and established it as arguably the most widely used electronic structure method. Despite the balanced accuracy-efficiency trade-off offered by KS-DFT, its computational cost remains a bottleneck for large-scale simulations or high-throughput computations. As a potential solution, orbital-free DFT (OF-DFT) [3] theoretically features linear scaling with system size, but its practical use has been hindered by the lack of an accurate kinetic energy (KE) functional approximation. To address these challenges, recent developments have incorporated machine learning (ML) into DFT [4], where computationally intensive steps are bypassed with surrogate ML models, enabling dramatic acceleration without sacrificing significant accuracy.

We can categorize prior efforts toward ML acceleration of KS-DFT into three major paradigms based on the steps being bypassed. In the first paradigm, *property learning*, a machine-learning model predicts a specific property directly from atomic configurations, replacing the DFT pipeline en bloc in a purely data-driven fashion. A common and useful subset of this paradigm is machine-learned interatomic potentials (MLIPs) [5–10], which predict the potential energy surface (PES) of a system and derive the corresponding force field, enabling efficient molecular dynamics (MD) simulations. The second paradigm, *electronic structure learning*, gets around the iterative process of solving the Kohn-Sham equations by predicting quantities that fully describe the ground-state electronic structure, such as electron density [11–21] or the quantum Hamiltonian [15, 22–27]. This approach allows for the derivation of any properties obtainable through standard KS-DFT, at the cost of a single additional Kohn-Sham diagonalization step.

The third paradigm is what we shall refer to as *mechanics learning*. While retaining the traditional DFT workflow, a mechanics learning approach achieves acceleration by learning the inner workings of the theory. Current research in this paradigm primarily targets the KE functional $T_s[\rho]$ [28–33] or the XC functional $E_{xc}[\rho]$ [34–37]. Machine-learned KE functionals have enabled efficient OF-DFT calculations [3] where the energy functional is directly minimized without introducing Kohn-Sham orbitals. Recent progress [33] has successfully extended machine-learned OF-DFT to large-scale datasets of real-world molecules, demonstrating accuracy comparable with KS-DFT and outstanding extrapolative ability to larger systems. On the other hand, machine-learned XC functionals have mainly aimed to elevate KS-DFT’s accuracy to higher-level methods rather than accelerate it. Unlike the KE functionals used in OF-DFT, there currently lacks a suitable learning objective in the Kohn-Sham formalism that both captures its intrinsic mechanics and is effective for ML acceleration, which is desirable in view of KS-DFT’s wide application.

In this work, we establish that the *Kohn-Sham density map* is a suitable primary objective for the mechanics learning, which leads us to develop a novel framework called neural network self-consistent fields (NeuralSCF). As is well-known (also illustrated in Fig. 1a), the KS-DFT is a self-consistent field (SCF) theory; that is an input electron density ρ_{in} generates a Kohn-Sham effective potential $v_{\text{KS}}[\rho_{\text{in}}]$, which in turn yields an output electron density ρ_{out} upon solving the corresponding Kohn-Sham equations. This defines the Kohn-Sham density map, $\mathcal{F}_{\text{KS}} : \rho_{\text{in}} \mapsto \rho_{\text{out}}$, which encodes the core mechanics of the Kohn-Sham equations. The process of achieving self-consistency in KS-DFT amounts to locating the fixed point of \mathcal{F}_{KS} such that $\rho^* = \mathcal{F}_{\text{KS}}[\rho^*]$, where ρ^* is precisely the ground-state electron density.

Though originating from the Kohn-Sham equations, \mathcal{F}_{KS} can be seen as a universal orbital-free map that can be equivalently formulated as the Euler-Lagrange

equation of the kinetic energy functional (Section V A):

$$-\left. \frac{\delta T_s[\rho]}{\delta \rho} \right|_{\rho_{\text{out}}} = v_{\text{KS}}[\rho_{\text{in}}]. \quad (1)$$

This orbital-free nature allows us to circumvent the expensive Kohn-Sham equations with a machine-learned Kohn-Sham density map that directly operates on electron densities, referred to as the NeuralSCF density map hereafter. By representing the electron density as expansion coefficients under an atom-centered Gaussian basis, we build the NeuralSCF density map with a SE(3)-equivariant graph transformer, designed to strictly preserve the symmetries of this map between two scalar fields. For effective training, we develop a two-stage strategy (Section II B) utilizing SCF trajectories from KS-DFT, a rich data source that has been unexploited in property learning or electronic structure learning. Once trained, the NeuralSCF density map can be used in the same manner as the Kohn-Sham density map, performing self-consistent iterations to solve for its fixed point, i.e. the predicted ground-state electron density. From this prediction, electronic properties can eventually be derived by reconstructing the Kohn-Sham Hamiltonian and performing a single diagonalization.

We demonstrate NeuralSCF’s effectiveness by experimenting on a diverse range of molecular datasets. For comparison, we introduce an end-to-end density predictor with the same architecture, representing a state-of-the-art model of its kind. NeuralSCF shows consistently lower error than its end-to-end counterpart on standard molecular datasets in self-consistent electron density and derived properties, oftentimes by significant margins. Notably, on QM9 [38], NeuralSCF achieves a threefold lower error in electron density prediction over previous efforts, with errors in derived energy two orders of magnitude lower than the chemical accuracy threshold. We further highlight NeuralSCF’s ability to generalize to out-of-distribution samples in a zero-shot setting, including off-equilibrium geometries, bond rotation, and non-covalent systems. In these tests, NeuralSCF still achieves chemical accuracy while significantly outperforming its end-to-end counterpart. These results demonstrate NeuralSCF as an accurate, robust, and transferable framework for DFT acceleration, revealing the strong generalization capabilities of mechanics-based ML models and paving the way for their broader use in electronic structure calculations.

II. The NeuralSCF framework

NeuralSCF is a deep learning framework modeling the Kohn-Sham density map with a neural network, designed to predict the ground-state electron density by emulating SCF iterations, as demonstrated in Fig. 1a. In NeuralSCF, the all-electron density is expanded by an atom-centered Gaussian basis set $\{\chi_p\}$ as $\rho(\mathbf{r}) = \sum_{p=1}^{N_{\text{aux}}} d_p \chi_p(\mathbf{r})$, where N_{aux} represents the number of basis functions and d_p

are the density coefficients, as illustrated in Fig. 1b. The label p is the collection of $\{inlm\}$, with i denoting the atom index, n, ℓ, m denoting the principal, angular, and magnetic quantum numbers, respectively. This basis set $\{\chi_p\}$ is commonly referred to as the auxiliary basis, distinguishing it from the atomic orbital basis $\{\phi_\mu\}$ used to expand the Kohn-Sham orbitals.

In NeuralSCF, we introduce a SE(3)-equivariant graph transformer to model the Kohn-Sham density map (Section II A). This neural network receives two inputs: input density coefficients \mathbf{d}_{in} , and the atomic configuration $\mathcal{X} = \{(\mathbf{R}_i, Z_i)\}$, where \mathbf{R}_i denotes atomic coordinates and Z_i is the atomic number. The output is a new set of density coefficients:

$$\hat{\mathbf{d}}_{\text{out}} = f_\theta(\mathbf{d}_{\text{in}}; \mathcal{X}), \quad (2)$$

with f_θ denoting the NeuralSCF density map with trainable parameters θ . NeuralSCF’s final prediction of ground-state electron density $\hat{\mathbf{d}}^*$ is defined as the fixed point of f_θ , satisfying $\hat{\mathbf{d}}^* = f_\theta(\hat{\mathbf{d}}^*; \mathcal{X})$, which is coherent with deep equilibrium models [39] whose output is defined as the fixed point of a single neural network layer.

The inference workflow of NeuralSCF, i.e. finding the neural network’s fixed point $\hat{\mathbf{d}}^*$, closely resembles the SCF procedure in standard KS-DFT. The model starts with a superposition of atomic densities (SAD) [40] initial guess $\mathbf{d}_{\text{in}}^{(0)}$. In the t -th self-consistent iteration, the current density coefficients $\mathbf{d}_{\text{in}}^{(t)}$ are first passed through the model to obtain $\hat{\mathbf{d}}_{\text{out}}^{(t)} = f_\theta(\mathbf{d}_{\text{in}}^{(t)}; \mathcal{X})$. Instead of directly using $\hat{\mathbf{d}}_{\text{out}}^{(t)}$ as the next input, the coefficients are mixed with previous density coefficients to obtain the new input $\mathbf{d}_{\text{in}}^{(t+1)}$. This technique, known as density mixing, is essential to SCF processes to stabilize and accelerate fixed-point convergence. We introduce a modified version of Pulay mixing for NeuralSCF that operates on density coefficients, detailed in Section V D. Finally, with a highly accurate prediction of the self-consistent electron density, the complete electronic structure can be derived from a single extra Kohn-Sham iteration (Section V G).

A. Network architecture

Preserving symmetries is a core principle in the design of machine-learning models for physical systems. This means that when the input quantities to the model undergo a spatial translation or a proper rotation, the output quantities should transform accordingly. This property, known as SE(3)-equivariance, can be strictly achieved thanks to recent advances in model architecture design [41]. In NeuralSCF, we employ a spherical SE(3)-equivariant graph neural network [42], one of the most general and expressive equivariant frameworks.

A spherical EGNN represents the system as an atomistic graph, and uses spherical tensors for node and edge features to ensure the outputs are connected to the inputs

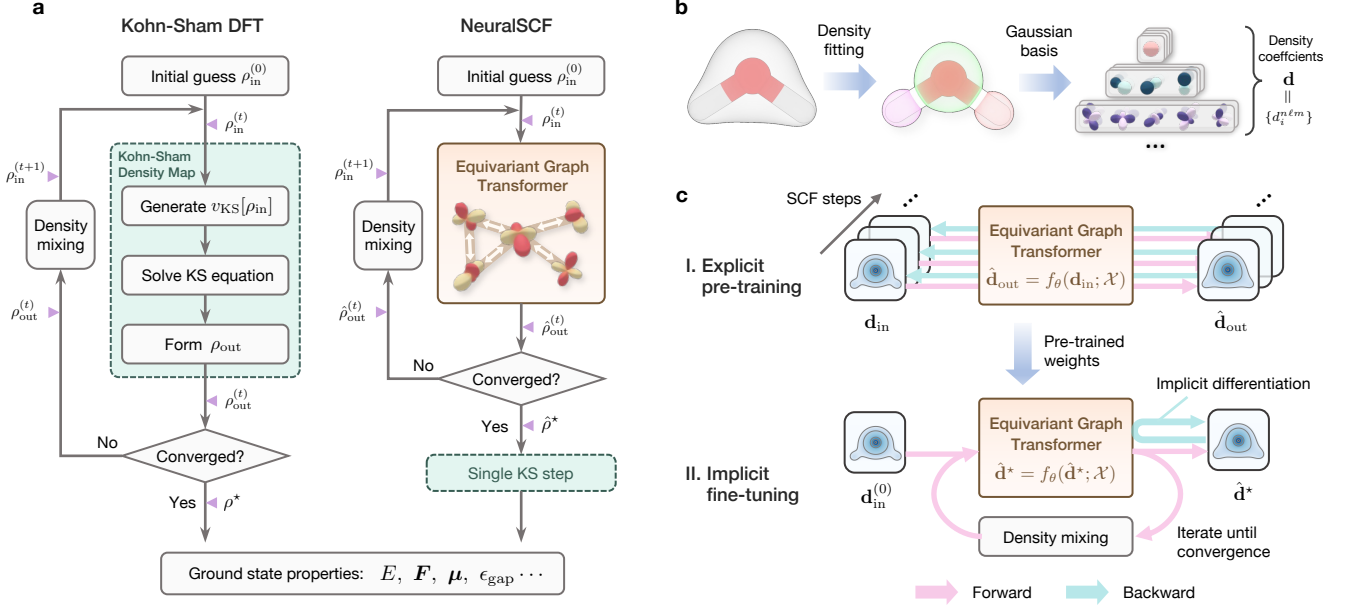


FIG. 1. **Overview of NeuralSCF.** (a) A comparison of the workflow of standard Kohn-Sham DFT and NeuralSCF. NeuralSCF models the Kohn-Sham density map using an equivariant graph transformer. NeuralSCF’s prediction of the ground-state electron density is defined by its fixed point, solved through self-consistent iterations aided by a density mixing scheme. Finally, ground state can be obtained from the predicted density with an extra Kohn-Sham step. (b) The electron density is represented by the expansion coefficients under a set of atom-centered Gaussian basis functions, which can be decomposed into atom-wise spherical tensors. (c) The two-stage training strategy of NeuralSCF. The explicit pre-training stage learns the Kohn-Sham density map from SCF trajectory data, while the implicit fine-tuning stage further aligns the model’s fixed point with the self-consistent electron density via implicit differentiation.

by a sequence of learnable equivariant operations. Specifically, an ℓ -degree ($\ell = 0, 1, 2, \dots$) irreducible spherical tensor $\mathbf{x}^{(\ell)} \in \mathbb{R}^{2\ell+1}$ carries an irreducible representation of $\text{SO}(3)$, and transforms under a proper rotation $\mathbf{R} \in \text{SO}(3)$ as:

$$(\mathbf{x}')_m^{(\ell)} = \sum_{m'=-\ell}^{\ell} \mathcal{D}_{mm'}^{(\ell)}(\mathbf{R}) x_{m'}^{(\ell)}, \quad m = -\ell, \dots, +\ell, \quad (3)$$

where $\mathcal{D}^{(\ell)}(\mathbf{R})$ is the Wigner D-matrix of degree ℓ . In NeuralSCF, density coefficients $\mathbf{d} = \{d_i^{n\ell m}\}$ are expansion coefficients under basis functions whose angular parts are spherical harmonics. Thus, these coefficients can be naturally decomposed into atom-wise spherical tensors $\{\mathbf{d}_i^{n\ell}\}$ of degree $\ell \geq 0$, which seamlessly integrate into spherical EGNNs.

The network architecture of the NeuralSCF density map $\hat{\mathbf{d}}_{\text{out}} = f_{\theta}(\mathbf{d}_{\text{in}}; \mathcal{X})$ is depicted in Fig. 2a. First, an atom-wise density encoder generates the initial node features from original inputs, as illustrated in Fig. 2c. In the density encoder, atom-wise density coefficients \mathbf{d}_i are rescaled and then transformed by an atom-type-specific equivariant linear layer into a homogeneous density feature. Meanwhile, the atomic number Z_i and the local environment $\{\mathbf{r}_{ij}\}$ are embedded into a geometry feature,

which is added to the density feature to form the initial node features.

The atomistic graph is then passed through a series of identical message-passing blocks, each with independent and trainable parameters, where node features are updated by aggregating information from neighboring nodes. We adopt a state-of-the-art equivariant transformer architecture from EquiformerV2 [43] for the message-passing blocks, which incorporates the widely successful self-attention mechanism [44] with the equivariant framework. Another key architectural improvement is the replacement of standard $\text{SO}(3)$ convolution with $\text{SO}(2)$ convolution [45], reducing the cost of spherical EGNNs’ computation bottleneck from $O(\ell_{\text{max}}^6)$ to that of the $\text{SO}(2)$ convolution, $O(\ell_{\text{max}}^3)$. This advancement enables us to augment the node representations to higher angular degrees ($\ell_{\text{max}} = 5$) without shrinking the channel width. Higher degree representations prove crucial for accurately modeling the electron density in NeuralSCF since density coefficients are themselves high-degree spherical tensors. Finally, the output node features are passed through an equivariant layer normalization module. They are then decoded by a density decoder into atom-wise output density coefficients, as shown in Fig. 2d.

To demonstrate the capabilities of the NeuralSCF framework, we implement an end-to-end density predictor,

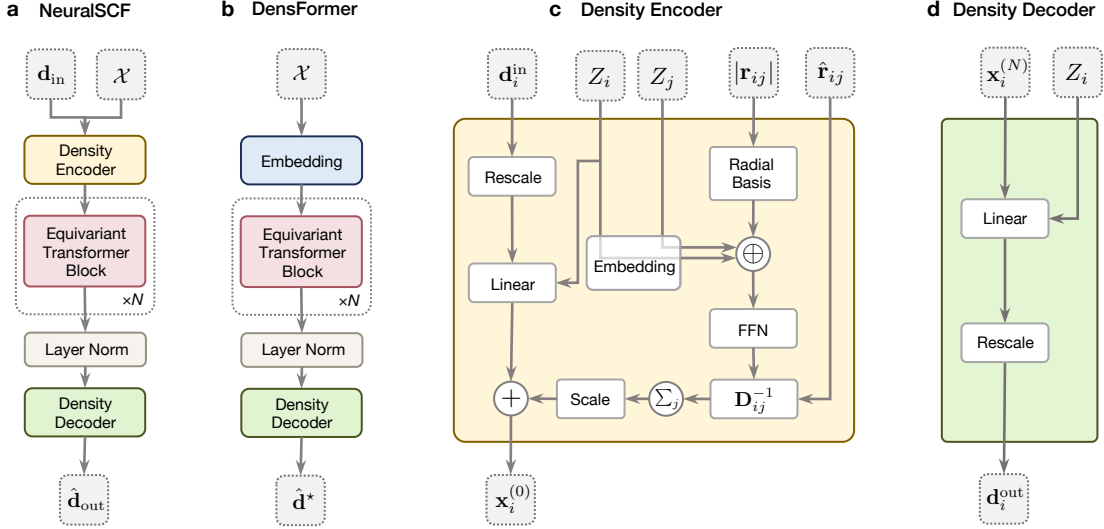


FIG. 2. **Network architecture.** (a) Network architecture of the proposed NeuralSCF density map $\hat{\mathbf{d}}_{\text{out}} = f_{\theta}(\mathbf{d}_{\text{in}}; \mathcal{X})$. (b) Network architecture of DensFormer, an end-to-end baseline model $\hat{\mathbf{d}} = g_{\theta}(\mathcal{X})$ that directly predicts self-consistent density from atomic configurations, sharing the same architecture as the NeuralSCF density map except for the input layer. (c) The density encoder transforms input atom-wise density coefficients, atom type, and the local environment into homogeneous input node features $\mathbf{x}_i^{(0)}$. Here, “FFN” stands for feed forward network, “ \oplus ” denotes concatenation, and “+” denotes element-wise addition. (d) The density decoder transforms output node features $\mathbf{x}_i^{(N)}$ into output atom-wise density coefficients.

DensFormer (**Density Transformer**), for comparison to rule out the contribution from architectural improvements. DensFormer directly predicts the self-consistent density coefficients from atomic configuration, conceptually similar to the equivariant GNN model proposed by Rackers *et al.* [18]. For a fair comparison, DensFormer shares the same equivariant transformer architecture as NeuralSCF in all experiments, except that the input layer is replaced by an atomic number embedding layer (Fig. 2b), resulting in approximately a 1% difference in the total number of parameters. Detailed descriptions of the neural network modules are available in Section V H and supplementary materials.

B. Two-stage training of NeuralSCF

Explicit pre-training. In the explicit pre-training stage, the NeuralSCF density map $\hat{\mathbf{d}}_{\text{out}} = f_{\theta}(\mathbf{d}_{\text{in}}; \mathcal{X})$ is trained to approximate the Kohn-Sham density map \mathcal{F}_{KS} . This is accomplished by explicitly learning from non-self-consistent density pairs $(\rho_{\text{in}}, \rho_{\text{out}})$ sampled from the SCF trajectory of KS-DFT calculations, as illustrated in Fig. 1c. Density pairs from different SCF iteration steps of the same structure are treated as independent and equal-weighted samples. We define the training loss as the L^2 -difference between the predicted density $\hat{\rho}_{\text{out}}$ and the ground truth ρ_{out} :

$$\mathcal{L}_{\text{exp}}(\theta) = \mathcal{L}^2(\rho_{\text{out}}, \hat{\rho}_{\text{out}}(\rho_{\text{in}}, \mathcal{X}; \theta)), \quad (4)$$

where $\mathcal{L}^2(\rho, \hat{\rho}) \equiv \int d^3\mathbf{r} |\rho(\mathbf{r}) - \hat{\rho}(\mathbf{r})|^2$. The L^2 -metric has an analytic expression under the density coefficients representation (Section V C), which can be optimized via the standard back-propagation algorithm. By leveraging SCF trajectory data, which is not utilized in property learning or electronic structure learning, NeuralSCF learns the internal mechanics of the Kohn-Sham equations without additional computational cost for data generation.

Implicit fine-tuning. After completing explicit pre-training, NeuralSCF can already predict the ground-state electron density through self-consistent iterations. Nonetheless, the prediction accuracy would be suboptimal as the training objective Eq. (4) lacks emphasis on the accuracy of the fixed point. To further enhance the model’s accuracy, we introduce an implicit fine-tuning stage where the model’s fixed point is calibrated to match the self-consistent electron density ρ^* , as depicted in Fig. 1c.

The forward pass of the fine-tuning stage is identical to the inference process, wherein the model solves for the fixed point $\hat{\mathbf{d}}^*$ iteratively. The training loss is defined as the L^2 -difference between the predicted self-consistent density $\hat{\rho}^*$ and the ground truth ρ^* :

$$\mathcal{L}_{\text{imp}}(\theta) = \mathcal{L}^2(\rho^*, \hat{\rho}^*(\mathcal{X}; \theta)). \quad (5)$$

The backward pass, however, is more challenging as the model’s prediction is now implicitly defined by the self-consistent equation $\hat{\mathbf{d}}^* = f_{\theta}(\hat{\mathbf{d}}^*; \mathcal{X})$. To compute the gradient of the loss function, one possible approach is to differentiate through all forward self-consistent iterations. Yet, this approach becomes computationally

expensive and numerically unstable as the number of forward SCF steps increases. Fortunately, the gradient can be alternatively calculated by applying the implicit function theorem at the fixed point [35, 39, 46], which solely requires the value of the fixed point and thus eliminates the need for tracking gradients during forward iterations. The implementation details of implicit differentiation in NeuralSCF are provided in Section V E.

III. Benchmark and performance

In this Section, we apply NeuralSCF to a range of molecular datasets drawn from literature, to assess the its ability to predict electron density and derived properties. As summarized in Table I, these datasets include equilibrium organic molecules, MD geometries, bond-rotated molecules, and non-covalent organic dimers.

A. Accurate prediction of electron density and derived properties

QM9. We first evaluate NeuralSCF’s compositional space generalizability on the QM9 dataset [38], a chemically diverse dataset commonly recognized as the gold standard for benchmarking atomistic machine-learning models. QM9 comprises 134k stable organic molecules with up to nine heavy atoms (C, N, O, F) in optimized equilibrium geometries. We partition it into 5,000 molecules for validation, 10,000 for testing, and the remainder for training. For explicit pre-training, we sample only density pairs from the first 8 SCF iterations, as changes in electron density beyond this point are generally negligible compared to the model’s precision.

We begin with the direct output from the models—the self-consistent electron density $\hat{\rho}^*$. Following common practice [12, 13, 17–19], we report the normalized mean absolute error (NMAE) ε_{ρ^*} , i.e. L^1 -error normalized by the number of electrons N_e , as the accuracy metric for electron density prediction:

$$\varepsilon_{\rho^*} = \frac{\int d^3\mathbf{r} |\rho^*(\mathbf{r}) - \hat{\rho}^*(\mathbf{r})|}{\int d^3\mathbf{r} \rho^*(\mathbf{r})} = \frac{\mathcal{L}^1(\rho^*, \hat{\rho}^*)}{N_e}, \quad (6)$$

NeuralSCF, after only the explicit pre-training stage, can already converge to a fixed point robustly during self-consistent inference, achieving an average $\varepsilon_{\rho^*} = 0.197\%$ on the QM9 test set. This accuracy already matches previous state-of-the-art density predictors, such as OrbNet-Equi’s 0.206% [17] and ChargE3Net’s 0.196% [20]. After the implicit fine-tuning stage, NeuralSCF further improves the NMAE to 0.064%, surpassing DensFormer’s average $\varepsilon_{\rho^*} = 0.070\%$ and highlighting a threefold improvement over previous efforts. Additionally, we examined outliers in the error distribution as indicators of the model’s robustness, marked in the histogram of self-consistent density ε_{ρ^*} (Fig. 3a). NeuralSCF exhibits a shorter tail

of outliers compared to DensFormer, with the largest outlier among 10,000 test molecules at 0.39% for NeuralSCF versus 0.53% for DensFormer. We also calculate dipole moments of the QM9 test set directly from the predicted electron densities, where NeuralSCF achieves a mean absolute error (MAE) of 0.035 D compared to DensFormer’s 0.042 D (Fig. 3b). Although NeuralSCF’s accuracy on dipole moment does not surpass that of the state-of-the-art property predictor reported in a recent benchmark [43], it remains comparable with common property predictors, even though predicting the full density distribution is much more challenging than predicting a single scalar value.

Next, we assess the accuracy of properties derived from the predicted densities by performing an additional Kohn-Sham diagonalization. In terms of total energy, NeuralSCF achieves a remarkably low mean absolute error (MAE) of 0.010 kcal mol⁻¹, surpassing DensFormer’s 0.017 kcal mol⁻¹ and being two orders of magnitude lower than the typical chemical accuracy threshold of 1 kcal mol⁻¹, as presented in Fig. 3c. This exceptionally low error in total energy stems from both NeuralSCF’s accurate prediction of electron density and a fundamental property of the total energy functional, whose deviation around the ground state is of second order in the density deviation [47], i.e. $E[\rho + \delta\rho] - E[\rho] = \mathcal{O}(\delta\rho^2)$. Surprisingly, even accounting for this second-order property, the gap in total energy MAE between NeuralSCF and DensFormer is much larger than expected from their difference in ε_{ρ^*} , suggesting that NeuralSCF’s predicted density may be more accurate than ε_{ρ^*} alone reflects. For reference, the state-of-the-art property predictor reported an energy MAE of 0.135 kcal mol⁻¹, and the KE-functional-based M-OFDFT [33] reported 0.93 kcal mol⁻¹. As another comparison, a previous attempt to derive total energy from predicted electron density [19] reported an energy MAE of 1.57 kcal mol⁻¹ on QM9, with the model trained on only 6% of the QM9 training set. For total energy, NeuralSCF exhibits a considerably shorter tail of error outliers than DensFormer, with its largest outlier being only 1.96 kcal mol⁻¹ compared to 5.93 kcal mol⁻¹ for DensFormer. Additionally, we derive the HOMO-LUMO gap, a quantity known to be notoriously difficult to predict accurately due to its nonlocal nature [19]. As shown in Fig. 3d, NeuralSCF achieves a MAE of 7.3 meV, outperforming DensFormer’s 9.2 meV and the state-of-the-art property predictor’s 29 meV.

MD datasets. We further benchmark NeuralSCF on datasets of molecular dynamics trajectories to examine its configurational space generalizability. We select two datasets: ethanol, a small molecule from the MD17 dataset [5], and Ac-Ala3-NHMe, a tripeptide comprising 42 atoms from the MD22 dataset [10]. For each dataset, we randomly sample 1,000 snapshots for training and 500 each for validation and testing from the full trajectory.

We observe that NeuralSCF outperforms DensFormer by a large margin on both datasets in terms of self-

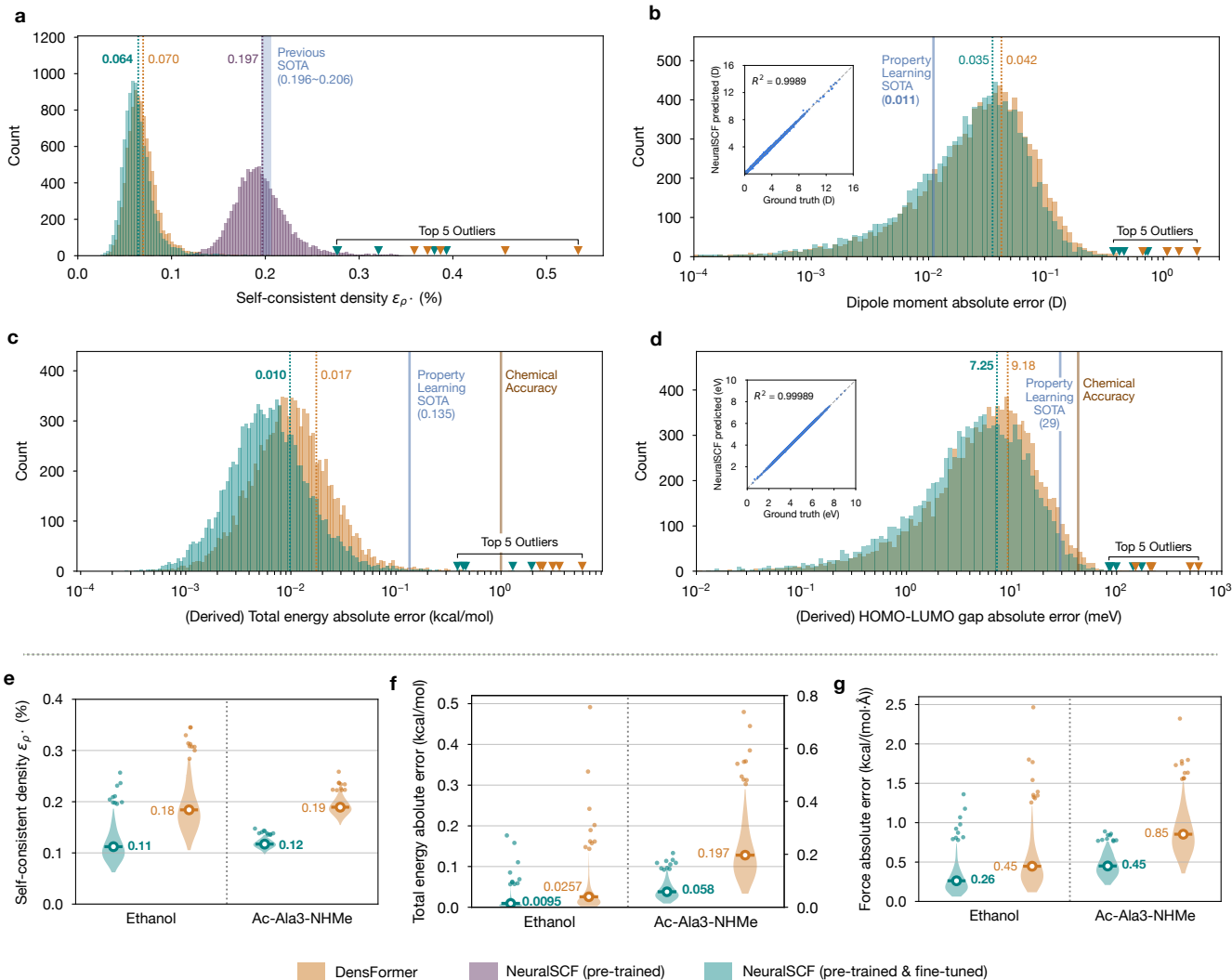


FIG. 3. **Results on QM9 and MD datasets.** (a-d) Histograms of absolute errors for self-consistent electron density, dipole moment, derived total energy, and derived HOMO-LUMO gap on QM9. The top 5 error outliers for NeuralSCF and its end-to-end density predictor counterpart DensFormer are marked with their respective colors. Vertical dotted lines indicate the mean absolute error (MAE) for each quantity. Property learning SOTAs are directly taken from a recent benchmark [43] as a reference. (e-g) Violin plots of absolute errors for self-consistent electron density, derived total energy and forces on ethanol (MD17) and Ac-Ala3-NHMe (MD22), with outliers shown as jittered scatter points. MAEs are indicated for each distribution.

consistent density, total energy, and forces, with their error distributions and MAEs summarized in Fig. 3e-g. Despite being trained on only 1,000 samples, NeuralSCF accurately predicts electron density on both datasets with $\epsilon_{\rho^*} \sim 0.1\%$, and the derived energy remains up to two orders of magnitude below the chemical accuracy threshold. This large performance gap between NeuralSCF and DensFormer on these datasets further indicates NeuralSCF’s advantage in low-data regime, likely due to the additional knowledge learned from SCF trajectories.

B. Zero-shot out-of-distribution generalization

Neural networks often perform well in generalizing to unseen samples but struggle with those far outside the training distribution, a phenomenon known as the out-of-distribution (OOD) problem [4]. This issue can be particularly pronounced in scientific applications, where new discoveries often lie outside the existing data distribution. The prevailing approach to address this is scaling up data and model sizes, which has proven effective in improving overall generalization. For instance, recent efforts in developing universal machine-learned interatomic potentials (MLIPs) across the periodic table [48–52] have demonstrated impressive zero-shot generalization to downstream

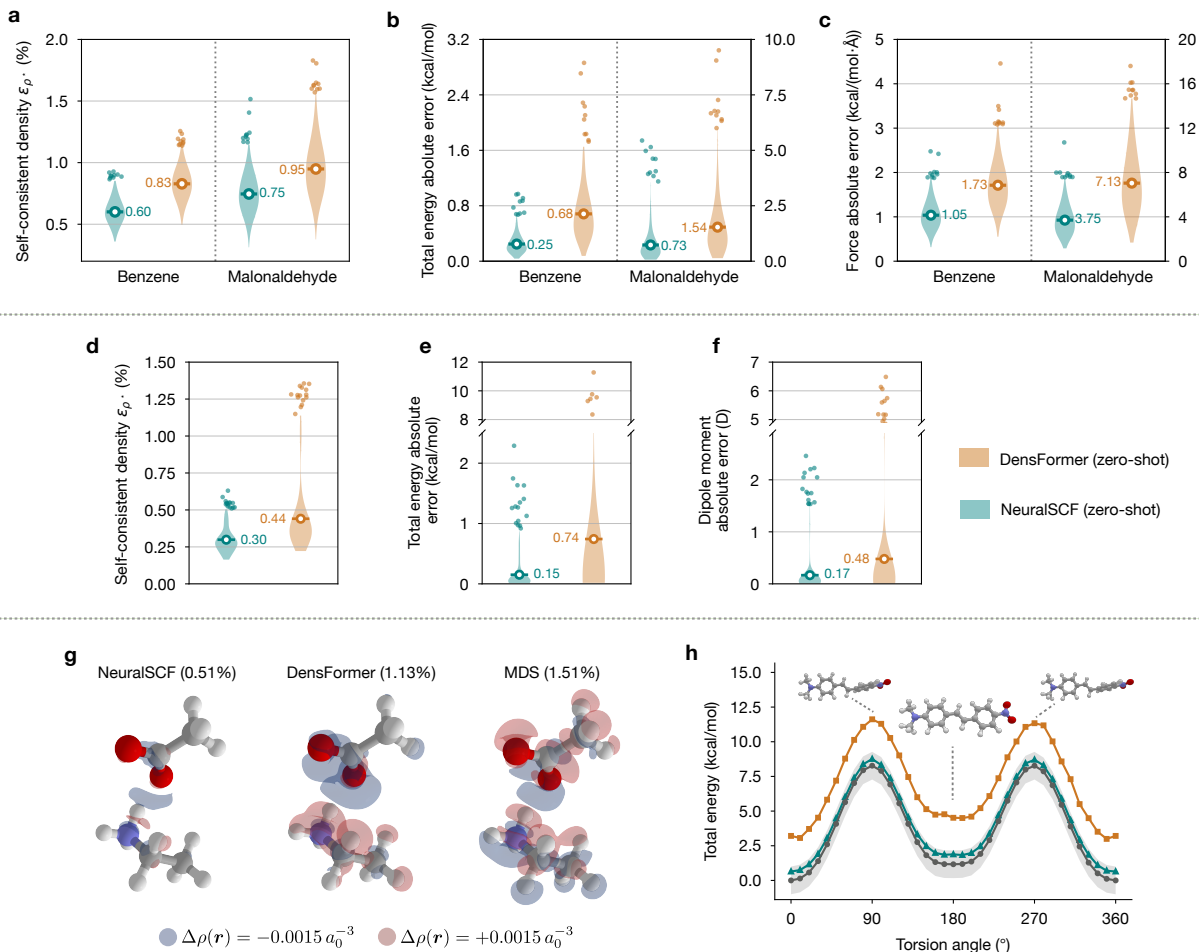


FIG. 4. **Zero-shot generalization to a wide range of datasets.** (a-c) Violin plots of absolute errors for self-consistent electron density, derived total energy and forces on MD17’s ethanol and malonaldehyde. (d-f) Violin plots of absolute errors for self-consistent electron density, derived total energy and dipole moment on BFDb-SSI. (g) Visualization of predicted density errors on the Glu⁻/Lys⁺ system, a challenging example from the BFDb-SSI dataset featuring significant charge transfer. (h) The true (dark grey) and predicted torsion energy profiles of 4-dimethylamino-4'-nitrostilbene. Shaded regions represent error within chemical accuracy 1 kcal mol⁻¹.

tasks. However, it remains uncertain whether more data fundamentally improves OOD generalization or simply pushes the OOD boundary further away. Another approach to improving OOD generalization is embedding physical principles into the model as *a priori* knowledge, which themselves are perfectly transferable. A notable example is the strict preservation of Euclidean symmetry, which has become a *de facto* standard in modern atomistic ML models due to the accuracy and data efficiency it provides. In the context of DFT surrogate models, Zhang *et al.* [33] recently showed that a machine-learned energy functional significantly outperforms its end-to-end energy predictor counterpart in extrapolating to larger systems, suggesting that choosing the underlying physical mechanics as the learning objective can qualitatively improve OOD generalization.

To further validate that learning mechanics learning is effective for enhancing OOD generalization, we conduct

comprehensive tests to evaluate NeuralSCF’s zero-shot OOD performance. Trained solely on the QM9 dataset, NeuralSCF is directly tested across a broad range of OOD systems, including off-equilibrium geometries, bond rotation, and non-covalent systems. Despite their modest performance gap on the in-distribution QM9 test set, NeuralSCF consistently demonstrates exceptional OOD generalization compared to its end-to-end counterpart, DensFormer, while maintaining chemical accuracy in terms of total energy.

Off-equilibrium geometries. As a test of configurational space generalization, we evaluate NeuralSCF, trained solely on QM9’s equilibrium geometries, directly on off-equilibrium geometries from the MD17 dataset. This type of generalization is typically not expected, as most previous atomistic ML models are limited to either chemically diverse equilibrium geometries or varied con-

figurations of a single system. Although a few models can handle both regimes [48–52], they all require training on large datasets spanning both compositional and configurational spaces.

We select benzene and malonaldehyde from MD17, randomly sampling 1,000 snapshots of each from the full MD trajectories for the zero-shot test set. As before, we report error distributions and MAEs on self-consistent density, total energy, and forces, as presented in Fig. 4a-c. Despite the notable performance degradation resulting from QM9’s complete absence of off-equilibrium sampling, NeuralSCF still maintains chemical accuracy and substantially outperforms DensFormer in terms of MAEs and the shorter tail of its error distribution for all three properties.

Non-covalent interactions. Non-covalent interactions play a crucial role in various chemical and biological processes and are essential components in molecular structures. More prevalent in larger systems or molecular assemblies, these interactions are underrepresented in QM9, which only contains small monomers. To evaluate NeuralSCF’s extrapolation to non-covalent systems, we conduct a zero-shot test on the BFDb-SSI dataset (side-chain side-chain interaction subset of the BioFragment database) [53], which includes 3,380 dimers representing a wide range of non-covalent interaction types. From a subset of 2,291 samples with no more than 25 atoms and no sulfur [13, 16], we further filter 1,796 neutral systems and randomly select 500 for testing.

Fig. 4d-f presents the zero-shot error distribution of NeuralSCF and DensFormer for self-consistent density, dipole moment, and derived total energy, highlighting NeuralSCF’s outstanding advantage in zero-shot generalization to non-covalent systems. For self-consistent density, NeuralSCF achieves an average $\varepsilon_{\rho^*} = 0.30\%$ compared to DensFormer’s 0.45%, with a significantly shorter tail of outliers. Remarkably, NeuralSCF’s zero-shot accuracy matches SA-GPR’s 0.29% [13] and is comparable to the more recent model OrbNet-Equi’s 0.19% [16], despite both being trained on 2,000 samples from the same dataset. For dipole moment and total energy, the performance difference is even more pronounced. NeuralSCF achieves an energy MAE of $0.15 \text{ kcal mol}^{-1}$ compared to DensFormer’s $0.74 \text{ kcal mol}^{-1}$, and a dipole moment MAE of 0.17 D compared to 0.48 D . This substantial performance gap in total energy and dipole moment, despite the relatively smaller difference in ε_{ρ^*} , indicates that NeuralSCF might more accurately capture charge transfer and redistribution caused by non-covalent interactions.

To further access NeuralSCF’s generalizability to systems with significant charge transfer, we examine a representative strongly interacting glutamic acid–lysine system whose Glu⁻/Lys⁺ salt bridge is essential for the helix stabilization in short peptides [54], following Qiao *et al.* [16]. For this system, NeuralSCF predicts the self-consistent density with an $\varepsilon_{\rho^*} = 0.50\%$, considerably lower than DensFormer’s 1.13%. As a simple baseline, monomer density superposition (MDS)–superpositioning independently

calculated DFT monomer densities–yields an $\varepsilon_{\rho^*} = 1.51\%$. The density errors are visualized in Fig. 4g using the $\Delta\rho(\mathbf{r}) = \pm 0.015 a_0^{-3}$ isosurfaces, showing that NeuralSCF accurately predicts the charge transfer between the Glu⁻ and Lys⁺ moieties.

Bond rotation. We further demonstrate NeuralSCF’s OOD configurational space generalization by investigating bond rotation. Following Chmiela *et al.* [10], we study 4-dimethylamino-4’-nitrostilbene, a donor-bridge-acceptor-type molecule consisting of two substituted phenyl rings connected by an ethylene bridge, which is rotated around the single bond between the acceptor and the ethylene bridge. This test is highly extrapolative beyond the QM9 dataset due to this molecule’s significantly larger size (20 heavy atoms) and additional rotational degree of freedom absent in QM9. Nevertheless, NeuralSCF predicts its electron density with an $\varepsilon_{\rho^*} = 0.28\%$ averaged over the full rotation, while accurately reproducing the torsion energy profile within chemical accuracy with an energy MAE of $0.51 \text{ kcal mol}^{-1}$, as presented in Fig. 4h. In comparison, DensFormer predicts the electron density with an average $\varepsilon_{\rho^*} = 0.43\%$, while its energy profile deviates significantly from the reference DFT curve with an MAE of $3.21 \text{ kcal mol}^{-1}$, failing to reproduce the correct minimum at 0° or the symmetry of the profile.

IV. Summary and outlook

To sum up, we have presented NeuralSCF, a novel deep-learning framework to accelerate KS-DFT by learning from its core mechanics. We establish the Kohn-Sham density map as the learning objective, enabling the model to learn the mechanics of the Kohn-Sham equations from previously underutilized SCF trajectory data. NeuralSCF predicts the electron density by emulating self-consistent iterations with the learned density map, and derives other properties from the predicted density.

NeuralSCF achieves state-of-the-art performance on benchmarks including QM9 and MD datasets, surpassing previous density predictor models by a large margin and outperforming its end-to-end counterpart. Moreover, NeuralSCF demonstrates exceptional zero-shot generalization to various out-of-distribution systems, including off-equilibrium geometries, bond rotation, and non-covalent interactions, considerably outmatching its end-to-end counterpart while maintaining chemical accuracy. These results mark a significant milestone, showing that mechanics-based models can achieve leading performance in both in-distribution and out-of-distribution generalization, indicating a potential path toward universal electronic structure models. The robustness and strong extrapolative performance of NeuralSCF also suggest its potential in accelerating high-throughput DFT data generation, offering near-DFT convergence predictions for a wide range of unseen systems and significantly reducing the DFT cost to as low as a single Kohn-Sham step.

As a general framework leveraging the self-consistent nature of KS-DFT, NeuralSCF opens up a number of new possibilities for future extensions and developments. While our current experiments have focused on neutral, closed-shell molecules, this framework can apply to closed-shell charged molecules and holds promise for future generalization to spin-polarized systems. NeuralSCF is also extensible to periodic systems, either using the current density coefficients representation or the grid-based representation more commonly used for plane wave basis, with an appropriate backbone model such as convolutional networks. While we utilize the orbital-free nature of the Kohn-Sham density map by using density coefficients, NeuralSCF can naturally adapt to density matrix representations, making it compatible with hybrid functionals that include the orbital-dependent exchange functional. Moreover, the self-consistent quantity in NeuralSCF is not limited to electron density; the KS effective potential or the KS Hamiltonian can also be subject to Kohn-Sham maps, where similar generalization improvements can be naturally envisioned.

V. Methods

A. Orbital-free definition of the Kohn-Sham density map

We initially define the Kohn-Sham density map $\rho_{\text{out}} = \mathcal{F}_{\text{KS}}[\rho_{\text{in}}]$ by the spin-unpolarized Kohn-Sham equations:

$$\left(-\frac{1}{2}\nabla^2 + v_{\text{KS}}[\rho_{\text{in}}](\mathbf{r})\right) \psi_i(\mathbf{r}) = \epsilon_i \psi_i(\mathbf{r}), \quad (7)$$

$$\rho_{\text{out}}(\mathbf{r}) = 2 \sum_{i=1}^{N_e/2} |\psi_i(\mathbf{r})|^2, \quad (8)$$

where $\{\psi_i(\mathbf{r})\}$ are the Kohn-Sham orbitals. Fixing ρ_{in} , $v_{\text{KS}}[\rho_{\text{in}}]$ can be treated as a regular single-particle potential. The energy functional of the non-interacting auxiliary system defined by Eq. (7) takes the form:

$$E_s[\rho; \rho_{\text{in}}] = T_s[\rho] + \int d^3\mathbf{r} \rho(\mathbf{r}) v_{\text{KS}}[\rho_{\text{in}}](\mathbf{r}) \quad (9)$$

By definition (8), ρ_{out} is the ground-state density of this auxiliary system, and should therefore minimize E_s according to the Hohenberg-Kohn theorems [1]:

$$0 = \left. \frac{\delta E_s[\rho; \rho_{\text{in}}]}{\delta \rho} \right|_{\rho_{\text{out}}} = \left. \frac{\delta T_s[\rho]}{\delta \rho} \right|_{\rho_{\text{out}}} + v_{\text{KS}}[\rho_{\text{in}}]. \quad (10)$$

This shows that the Kohn-Sham density map can be equivalently defined in an orbital-free manner (1).

B. Density fitting

In KS-DFT calculations under atomic orbital basis sets, the electron density is represented with the density matrix

\mathbf{D} as:

$$\rho(\mathbf{r}) = \sum_{\mu=1}^{N_{\text{ao}}} \sum_{\nu=1}^{N_{\text{ao}}} D_{\mu\nu} \phi_{\mu}(\mathbf{r}) \phi_{\nu}(\mathbf{r}). \quad (11)$$

The density matrix \mathbf{D} can be projected to density coefficients \mathbf{d} within the context of density fitting, a technique originally introduced to accelerate the calculation of electron repulsion integrals (ERIs). Let $\rho_{\mathbf{d}}$ and $\rho_{\mathbf{D}}$ represent the electron density generated by density coefficients \mathbf{d} and the density matrix \mathbf{D} , respectively. The density fitting problem is solved by minimizing the fitting residual $(\rho_{\mathbf{d}} - \rho_{\mathbf{D}} | \rho_{\mathbf{d}} - \rho_{\mathbf{D}})$, where $(\cdot | \cdot)$ denotes an inner product of two real space functions defined by the 2-center Coulomb integral:

$$(\phi_1 | \phi_2) \equiv \int d^3\mathbf{r}_1 \int d^3\mathbf{r}_2 \frac{\phi_1(\mathbf{r}_1) \phi_2(\mathbf{r}_2)}{|\mathbf{r}_1 - \mathbf{r}_2|}. \quad (12)$$

This optimization problem can be solved analytically [55], resulting in a linear transformation from the density matrix to density coefficients:

$$d_p = \sum_{q=1}^{N_{\text{aux}}} \sum_{\mu, \nu=1}^{N_{\text{ao}}} ((\mathbf{J}^{\text{X}})^{-1})_{pq} J_{\mu\nu;p} D_{\mu\nu}, \quad (13)$$

where $J_{pq}^{\text{X}} = (\chi_p | \chi_q)$, and $J_{\mu\nu;p} = (\phi_{\mu} \phi_{\nu} | \chi_p)$ is the 3-center Coulomb integral.

C. Efficient evaluation of density metrics

L^2 -metric. Under the density coefficients representation, the L^2 -metric of electron density has the following analytical expression:

$$\mathcal{L}^2(\rho, \hat{\rho}) = \Delta \mathbf{d}^T \mathbf{S}^{\text{X}} \Delta \mathbf{d}, \quad (14)$$

where $\Delta \mathbf{d} = \mathbf{d} - \hat{\mathbf{d}}$ is the difference between density coefficients and their predicted values, and \mathbf{S}^{X} is the overlap matrix of the auxiliary basis, defined by $S_{pq}^{\text{X}} = (\chi_p | \chi_q) = \int d^3\mathbf{r} \chi_p(\mathbf{r})^* \chi_q(\mathbf{r})$.

Despite this convenient form, the requirement of an overlap matrix for each structure presents challenges to storage and computation. Given the near-sparse nature of the overlap matrix, a straightforward approximation is to remove its insignificant matrix elements. However, since \mathbf{S}^{X} often possesses numerous small positive eigenvalues, minor modifications can break its positive definiteness and lead to a diverged training loss. To address this issue, we execute a Cholesky decomposition $\mathbf{S}^{\text{X}} = (\mathbf{L}^{\text{X}})^T \mathbf{L}^{\text{X}}$, with \mathbf{L}^{X} being a unique lower triangular matrix. This decomposition reduces the L^2 -metric Eq. (14) to a squared vector norm:

$$\mathcal{L}^2(\rho, \hat{\rho}) = (\mathbf{L}^{\text{X}} \Delta \mathbf{d})^T (\mathbf{L}^{\text{X}} \Delta \mathbf{d}) = \|\mathbf{L}^{\text{X}} \Delta \mathbf{d}\|_2^2, \quad (15)$$

which assures positive definiteness and is even more efficient to compute. Empirical observations show that

Dataset	Characteristic	Trained on	As test set
QM9 [38]	Equilibrium small organic molecules	✓	ID
MD17 [5] (ethanol)	MD trajectory of small organic molecules	✓	ID
MD17 [5] (benzene, malonaldehyde)	MD trajectory of small organic molecules	✗	OOD
MD22 [10] (Ac-Ala3-NHMe)	MD trajectory of larger organic molecules	✓	ID
4-dimethylamino-4'-nitrostilbene [10]	Full rotation around a single bond	✗	OOD
BFDb-SSI [53]	Organic dimers with non-covalent interactions	✗	OOD

TABLE I. Summary of datasets used in this work. “ID” stands for in-distribution, and “OOD” stands for out-of-distribution. Datasets that were never trained on are used as out-of-distribution test sets for models trained solely on QM9.

the Cholesky factor \mathbf{L}^χ maintains sparsity. Its matrix elements are thus pruned, cast to 16-bit float precision and stored in sparse COO format.

L^1 -metric. The L^1 -metric, however, can only be evaluated via numerical integration. We generate molecular grids with PySCF using a preset grid level of 0. The numerical integration is then formulated as

$$\mathcal{L}^1(\rho, \hat{\rho}) = \mathbf{w}^T \text{abs}(\mathbf{C}\Delta\mathbf{d}). \quad (16)$$

Here, $\mathbf{w} \in \mathbb{R}^{N_{\text{grid}}}$ denotes the weights of the grid points with N_{grid} being the grid size. The collocation matrix $\mathbf{C} \in \mathbb{R}^{N_{\text{grid}} \times N_{\text{aux}}}$ is defined by $(\mathbf{C})_{jp} = \chi_p(\mathbf{r}_j)$, where \mathbf{r}_j denotes the j -th grid point. The collocation matrix \mathbf{C} is applied with the same set of approximations used for the Cholesky factor to save memory usage. These approximations are applied only to the training and validation sets, while the test set is evaluated with full precision. With all auxiliary quantities (\mathbf{L}^χ , \mathbf{w} , \mathbf{C}) precomputed, we implement batched versions of both L^1 and L^2 -metric with PyTorch [56], utilizing GPU for highly efficient sparse operations which significantly reduces the overhead caused by evaluating density metrics.

D. Pulay mixing of density coefficients

Pulay mixing [57], also known as direct inversion of the iterative subspace (DIIS) or Anderson’s method, is one of the most robust and efficient mixing schemes for SCF calculations [58]. In NeuralSCF, we introduce a modified version of Pulay mixing that operates on density coefficients to facilitate the convergence of the forward iterations.

Consistent with the original Pulay mixing, the input for the next iteration is a linear combination of output density coefficients from previous n iterations:

$$\mathbf{d}_{\text{in}}^{(t+1)} = \sum_{i=1}^n \alpha_i \mathbf{d}_{\text{out}}^{(t-i+1)}, \quad (17)$$

where $n = \max\{t, N\}$ is the history length with a cut-off N , and $\{\alpha_i\}$ are mixing coefficients to be determined. Defining the residue of the current iteration

as $\delta\mathbf{d}^{(t)} = \mathbf{d}_{\text{out}}^{(t)} - \mathbf{d}_{\text{in}}^{(t)}$, the residue of the next iteration is estimated as the linear combination of previous residues using the same set of mixing coefficients, i.e. $\widetilde{\delta\mathbf{d}}^{(t+1)} = \sum_{i=1}^n \alpha_i \delta\mathbf{d}^{(t-i+1)}$.

The mixing coefficients are then determined by minimizing the L^2 -norm of the estimated density residue $\widetilde{\delta\rho}^{(t+1)}(\mathbf{r}) = \sum_p \widetilde{\delta d}_p^{(t+1)} \chi_p(\mathbf{r})$. Note that we minimize the norm of the electron density instead of the vector norm of density coefficients, as the former is more physically relevant and enjoys significantly faster fixed-point convergence in practice. Utilizing the Cholesky factor \mathbf{L}^χ of the overlap matrix introduced in Section VC, this optimization problem can be formulated as

$$\min_{\{\alpha_i\}} \left\| \mathbf{L}^\chi \widetilde{\delta\mathbf{d}}^{(t+1)} \right\|_2^2, \text{ s.t. } \sum_{i=1}^n \alpha_i = 1. \quad (18)$$

The optimal mixing coefficients are the solutions to the following linear equations [57]:

$$\begin{bmatrix} \mathbf{0} & \mathbf{1}^T \\ \mathbf{1} & \mathbf{B} \end{bmatrix} \begin{bmatrix} \alpha_0 \\ \boldsymbol{\alpha} \end{bmatrix} = \begin{bmatrix} \mathbf{1} \\ \mathbf{0} \end{bmatrix}, \quad (19)$$

where $\mathbf{1}$, $\mathbf{0}$, $\boldsymbol{\alpha} \in \mathbb{R}^n$ denotes column vectors of ones, zeros, and the mixing coefficients, respectively. The matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ is a symmetric matrix with elements $B_{ij} = (\mathbf{L}^\chi \delta\mathbf{d}^{(t-i+1)})^T (\mathbf{L}^\chi \delta\mathbf{d}^{(t-j+1)})$.

Finally, the self-consistent iteration terminates at the T -th step if the relative L^2 -norm of the density residue falls below a preset threshold η_ρ :

$$\frac{\int d^3\mathbf{r} |\delta\rho^{(T)}(\mathbf{r})|^2}{\int d^3\mathbf{r} |\rho^{(T)}(\mathbf{r})|^2} = \frac{\|\mathbf{L}^\chi \delta\mathbf{d}^{(T)}\|_2^2}{\|\mathbf{L}^\chi \mathbf{d}^{(T)}\|_2^2} < \eta_\rho, \quad (20)$$

which is set to $\eta_\rho = 10^{-8}$ for all experiments. At test time where the overlap matrix may be unavailable, the vector norm of density coefficients can be alternatively used as the convergence criterion. This is equivalent to simply setting $\mathbf{L}^\chi = \mathbf{I}$ in all the above equations.

E. Implicit differentiation for fine-tuning

Following the derivations in deep equilibrium models [39], the gradient of the loss function in Eq. (5) to model pa-

rameters is given by implicit function theorem:

$$\frac{\partial \mathcal{L}_{\text{imp}}}{\partial \theta} = \frac{\partial \mathcal{L}_{\text{imp}}}{\partial \hat{\mathbf{d}}^*} (\mathbf{I} - \mathbf{J}_f^*)^{-1} \frac{\partial f_\theta(\hat{\mathbf{d}}^*; \mathcal{X})}{\partial \theta}, \quad (21)$$

where

$$\mathbf{J}_f^* = \left. \frac{\partial f_\theta(\mathbf{d}_{\text{in}}; \mathcal{X})}{\partial \mathbf{d}_{\text{in}}} \right|_{\mathbf{d}_{\text{in}}=\hat{\mathbf{d}}^*} \quad (22)$$

is the Jacobian matrix at the model’s fixed point. This matrix can be prohibitively large to compute directly, not to mention finding its inverse. Nevertheless, it is sufficient to evaluate the following product as a whole:

$$\mathbf{y}^T \equiv \frac{\partial \mathcal{L}_{\text{imp}}}{\partial \hat{\mathbf{d}}^*} (\mathbf{I} - \mathbf{J}_f^*)^{-1}. \quad (23)$$

By multiplying $(\mathbf{I} - \mathbf{J}_f^*)$ to both sides of Eq. (23), one can immediately derive another linear self-consistent equation about \mathbf{y}^T :

$$\mathbf{y}^T = \mathbf{y}^T \mathbf{J}_f^* + \frac{\partial \mathcal{L}_{\text{imp}}}{\partial \hat{\mathbf{d}}^*}, \quad (24)$$

where the vector-Jacobian product $\mathbf{y}^T \mathbf{J}_f^*$ can be evaluated using PyTorch’s automatic differentiation engine. Thereby, \mathbf{y}^T can be solved with an iterative fixed-point solver to obtain the gradient. Note that only the value of the fixed point $\hat{\mathbf{d}}^*$ is required for the gradient, irrespective of the trajectory of the fixed-point iterations. This allows the forward pass to operate efficiently without tracking gradients, resulting in constant memory usage.

In practice, we solve Eq. (24) using a simple linear mixing scheme with a mixing factor $0 < \lambda < 1$. Here, \mathbf{y}^T is initialized with $\lambda \frac{\partial \mathcal{L}_{\text{imp}}}{\partial \hat{\mathbf{d}}^*}$, and then updated iteratively as

$$\mathbf{y}^T := (1 - \lambda) \left(\mathbf{y}^T \mathbf{J}_f^* + \frac{\partial \mathcal{L}_{\text{imp}}}{\partial \hat{\mathbf{d}}^*} \right) + \lambda \mathbf{y}^T, \quad (25)$$

for a fixed number of iterations, M . Although this strategy does not guarantee an exact solution to Eq. (24), we empirically observe that it does not detract from the model’s final performance compared to using the exact gradient given by Anderson’s method. This observation is consistent with the finding that random noises in stochastic gradient descent serve as regularization and lead to better generalization than gradient descent [59]. We choose $M = 8$ and $\lambda = 0.4$ for all experiments, which consistently results in stable training dynamics while significantly reducing the computational cost.

F. Dataset preparation

KS-DFT calculations. We perform KS-DFT calculations on all datasets present in this study using the PySCF package [60] with the cc-pVTZ basis set and the PBE exchange-correlation functional [61]. Calculations are

based on the geometries provided in the original datasets except for 4-dimethylamino-4’-nitrostilbene, whose geometry is re-optimized at PBE/cc-pVTZ level. Density fitting is applied in DFT calculations to reduce memory usage. To minimize the error introduced by density fitting, we generate a large auxiliary basis set for cc-pVTZ using Basis Set Exchange’s [62] implementation of the AutoAux algorithm [63]. This auxiliary basis set, henceforth referred to as cc-pVTZ-AutoAux, is used for both DFT calculations and density coefficients representation. The SCF convergence threshold is set to 10^{-8} Hartree, and the DFT grid level to 3 for all calculations. We extract SCF trajectories of density matrix $\{(\mathbf{D}_{\text{in}}^{(t)}, \mathbf{D}_{\text{out}}^{(t)})\}$ from the DFT calculations with a custom callback function at the end of each SCF iteration. Since PySCF performs DIIS on the Hamiltonian instead of the density matrix, an additional diagonalization of the Hamiltonian is required within the callback function to obtain the corresponding \mathbf{D}_{out} , which only adds an insignificant overhead to the total computation time.

Dataset preprocessing. The raw outputs from KS-DFT calculations are preprocessed to extract quantities essential for training NeuralSCF. For each structure, we extract its atomic configuration, input/output density coefficients at each SCF iteration (including the initial guess), pruned Cholesky factors, grid weights, and pruned collocation matrices. We retain the top 12% of nonzero matrix elements with the largest absolute values for Cholesky factors. For training and validation collocation matrices, we retain the top 6% of nonzero elements on QM9 and 10% on all other datasets. To facilitate efficient training, the atomic configurations, density coefficients, Cholesky factors, and grid weights are packed into a PyTorch Geometric’s [64] `InMemoryDataset` and loaded into shared CPU memory. Collocation matrices are stored on disk as an HDF5 file and are loaded into GPU memory in batches during the evaluation of the L^1 -error.

G. Deriving properties from predicted density

The dipole moment $\boldsymbol{\mu}$ is directly calculated from the predicted density coefficients as:

$$\boldsymbol{\mu} = \sum_i Z_i \mathbf{R}_i - \sum_p d_p \cdot \int d^3 \mathbf{r} \chi_p(\mathbf{r}) \mathbf{r}. \quad (26)$$

To obtain the total energy and other orbital-dependent properties such as the HOMO-LUMO gap, the Kohn-Sham Hamiltonian has to be constructed:

$$\mathbf{H}_{\text{KS}} = \mathbf{T} + \mathbf{V}_{\text{ext}} + \mathbf{V}_{\text{H}} + \mathbf{V}_{\text{xc}}. \quad (27)$$

Here, the kinetic energy operator \mathbf{T} and the external potential \mathbf{V}_{ext} are independent of electron density and are evaluated exactly under the atomic orbital basis using PySCF. The Coulomb matrix \mathbf{V}_{H} can be analytically

computed from the predicted density coefficients:

$$(\mathbf{V}_H)_{\mu\nu} = \sum_p (\phi_\mu \phi_\nu | \chi_p) \hat{d}_p. \quad (28)$$

For non-hybrid XC functionals, \mathbf{V}_{xc} is evaluated based on the electron density values on a numerical grid. We set PySCF’s grid level to 3 for all datasets. Finally, the Kohn-Sham Hamiltonian is diagonalized, from which the total energy and other ground-state electronic properties can be derived.

H. Details on neural network modules

In this section, we provide a detailed description of the neural network modules present in Fig. 2 except for the equivariant transformer block, which is elaborated in the supplementary materials. Full model details and training hyperparameters for each experiment are also provided in the supplementary materials.

Density coefficients rescaling. The scales of different components of the density coefficients often vary across multiple orders of magnitude, which causes difficulties in model training. To resolve this issue, we rescale the density coefficients equivariantly before feeding them into the neural network. For a given atom type Z , the atom-wise density coefficients are rescaled as:

$$d_Z^{n\ell m} \mapsto \frac{d_Z^{n\ell m} - \mu_Z^{n\ell}}{\sigma_Z^{n\ell}}. \quad (29)$$

For scalar ($\ell = 0$) components, $\mu_Z^{n0} = \langle d_i^{*n00} \rangle_{i \in Z}$ with $\langle \cdot \rangle_{i \in Z}$ denoting the average over all atoms of type Z in the training set, and $\sigma_Z^{n0} = \sqrt{\langle (d_i^{*n00} - \mu_Z^{n0})^2 \rangle_{i \in Z}}$. For $\ell > 0$ components, $\mu_Z^{n\ell}$ is set to zero to preserve equivariance, and $\sigma_Z^{n\ell} = \sqrt{\langle (d_i^{*n\ell m})^2 \rangle_{m, i \in Z}}$, where $\langle \cdot \rangle_{m, i \in A}$ denotes the average over all atoms of type Z in the training set and all orders $m = -\ell, \dots, \ell$ for a given ℓ . Density coefficients are rescaled by Eq. (29) in the density encoder and by its inverse transformation in the density decoder.

Equivariant linear layer. We denote a general spherical tensor, also known as an irreps feature [42], as $\mathbf{x} = \{x_{cm}^{(\ell)}\}$ with $1 \leq c \leq C_\ell$ being the channel index, i.e. the multiplicity of each irreducible representation. An equivariant linear layer transforms the input spherical tensor by mixing exclusively across channels:

$$(\mathbf{x}')_{cm}^{(\ell)} = \sum_{c'=1}^{C'_\ell} W_{cc'}^{(\ell)} x_{c'm}^{(\ell)} + b_c^{(0)} \delta_{\ell,0}, \quad (30)$$

where $W_{cc'}^{(\ell)}$ is the learnable weight, $\delta_{\ell,0}$ denotes the Kronecker delta, and $b_c^{(0)}$ is the learnable bias applied only to the scalar components for equivariance. In both the

density encoder and decoder, we employ an independent equivariant linear layer for each atom type to convert between the atom-wise density coefficients and a homogeneous spherical tensor feature irrespective of atom type.

Constructing initial geometry feature. As described in Section. II A, the density encoder creates an initial node feature by combining a geometry feature with a density feature. Here, we elaborate on the construction of the geometry feature (Fig. 2c). First, an edge distance embedding is generated by expanding the relative distances $|\mathbf{r}_{ij}|$ with a Gaussian radial basis [6]. Two atom embeddings are created independently for the source and target atoms i, j by passing the one-hot vectors representing their atom types through a scalar linear layer. The two atom embeddings are then concatenated with the edge distance embedding to form a single scalar vector, which is further transformed by a 2-layer feed-forward network with SiLU activation [65]. The resulting scalar vector, interpreted as the complete embedding of edge ij viewed in the local frame specified by $\hat{\mathbf{r}}_{ij}$, is then rotated back to the global frame by the inverse Wigner-D matrix $\mathcal{D}(\mathbf{R}_{ij})^{-1}$. Here, $\mathbf{R}_{ij} \in \text{SO}(3)$ and such that $\mathbf{R}_{ij} \cdot \hat{\mathbf{r}}_{ij} = (0, 0, 1)^T$. Finally, the geometry feature of atom i is obtained by summing up all the edge embeddings contributed by its neighbors j , then divided with the squared root of the average number of neighbors in the training set.

Equivariant layer normalization. We adopt the separable layer normalization introduced in EquiformerV2 [43], which is somewhat similar to the density coefficients rescaling method, where a spherical tensor is transformed as:

$$x_{cm}^{(\ell)} \mapsto \gamma_c^{(\ell)} \frac{x_{cm}^{(\ell)} - \mu^{(\ell)}}{\sigma^{(\ell)}} + \beta_c^{(0)} \delta_{\ell,0}. \quad (31)$$

Here, $\gamma_c^{(\ell)}$ and $\beta_c^{(0)}$ are learnable parameters. For scalar components, $\mu^{(0)} = \langle x_{c0}^{(0)} \rangle_c$ and $\sigma^{(0)} = \sqrt{\langle (x_{c0}^{(0)} - \mu^{(0)})^2 \rangle_c}$, where $\langle \cdot \rangle_c$ denotes the average across all channels. For non-scalar components, $\mu^{(\ell>0)} = 0$ and all degrees $\ell > 0$ share the same $\sigma^{(\ell>0)} = \sqrt{\langle (x_{cm}^{(\ell)})^2 \rangle_{\ell>0, c, m}}$, where $\langle \cdot \rangle_{\ell>0, c, m}$ denotes the average over all non-scalar degrees $\ell > 0$, all channels c , and all orders m .

Acknowledgments

We are grateful for insightful discussions with Yihao Lin, Qiangqiang Gu, Yi-Lun Liao, and Zhengyang Geng. We acknowledge the financial support from the National Natural Science Foundation of China (Grants No. 12274003, No. 11725415, and No. 11934001), the National Key R&D Program of China (Grants No. 2018YFA0305601 and No. 2021YFA1400100), and the Innovation Program for Quantum Science and Technology (Grant No. 2021ZD0302600).

- * jfeng11@pku.edu.cn
- [1] P. Hohenberg and W. Kohn, Inhomogeneous electron gas, *Phys. Rev.* **136**, B864 (1964).
 - [2] W. Kohn and L. J. Sham, Self-consistent equations including exchange and correlation effects, *Phys. Rev.* **140**, A1133 (1965).
 - [3] W. Mi, K. Luo, S. B. Trickey, and M. Pavanello, Orbital-free density functional theory: An attractive electronic structure method for large-scale first-principles simulations, *Chemical Reviews* **123**, 12039 (2023).
 - [4] X. Zhang *et al.*, Artificial intelligence for science in quantum, atomistic, and continuum systems (2023), [arXiv:2307.08423 \[cs.LG\]](https://arxiv.org/abs/2307.08423).
 - [5] S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller, Machine learning of accurate energy-conserving molecular force fields, *Science Advances* **3**, e1603015 (2017).
 - [6] K. Schütt *et al.*, SchNet: A continuous-filter convolutional neural network for modeling quantum interactions, in *Advances in Neural Information Processing Systems*, Vol. 30 (2017).
 - [7] L. Zhang, J. Han, H. Wang, R. Car, and W. E, Deep potential molecular dynamics: A scalable model with the accuracy of quantum mechanics, *Phys. Rev. Lett.* **120**, 143001 (2018).
 - [8] S. Batzner *et al.*, E(3)-equivariant graph neural networks for data-efficient and accurate interatomic potentials, *Nature Communications* **13**, 2453 (2022).
 - [9] I. Batatia, D. P. Kovacs, G. Simm, C. Ortner, and G. Csanyi, MACE: Higher order equivariant message passing neural networks for fast and accurate force fields, in *Advances in Neural Information Processing Systems*, Vol. 35 (2022) pp. 11423–11436.
 - [10] S. Chmiela *et al.*, Accurate global machine learning force fields for molecules with hundreds of atoms, *Science Advances* **9**, eadf0873 (2023).
 - [11] F. Brockherde, L. Vogt, L. Li, M. E. Tuckerman, K. Burke, and K.-R. Müller, Bypassing the Kohn-Sham equations with machine learning, *Nature Communications* **8**, 872 (2017).
 - [12] A. Grisafi, A. Fabrizio, B. Meyer, D. M. Wilkins, C. Corminboeuf, and M. Ceriotti, Transferable machine-learning model of the electron density, *ACS Central Science* **5**, 57 (2019), PMID: 30693325.
 - [13] A. Fabrizio, A. Grisafi, B. Meyer, M. Ceriotti, and C. Corminboeuf, Electron density learning of non-covalent systems, *Chemical science* **10**, 9424 (2019).
 - [14] L. Zepeda-Núñez, Y. Chen, J. Zhang, W. Jia, L. Zhang, and L. Lin, Deep Density: Circumventing the Kohn-Sham equations via symmetry preserving neural networks, *Journal of Computational Physics* **443**, 110523 (2021).
 - [15] O. Unke, M. Bogojeski, M. Gastegger, M. Geiger, T. Smidt, and K.-R. Müller, SE(3)-equivariant prediction of molecular wavefunctions and electronic densities, in *Advances in Neural Information Processing Systems*, Vol. 34 (2021) pp. 14434–14447.
 - [16] Z. Qiao, A. S. Christensen, M. Welborn, F. R. Manby, A. Anandkumar, and T. F. Miller, Informing geometric deep learning with electronic interactions to accelerate quantum chemistry, *Proceedings of the National Academy of Sciences* **119**, e2205221119 (2022).
 - [17] P. B. Jørgensen and A. Bhowmik, Equivariant graph neural networks for fast electron density estimation of molecules, liquids, and solids, *npj Computational Materials* **8**, 183 (2022).
 - [18] J. A. Rackers, L. Tecot, M. Geiger, and T. E. Smidt, A recipe for cracking the quantum scaling limit with machine learned electron densities, *Machine Learning: Science and Technology* **4**, 015027 (2023).
 - [19] A. Grisafi, A. M. Lewis, M. Rossi, and M. Ceriotti, Electronic-structure properties from atom-centered predictions of the electron density, *Journal of Chemical Theory and Computation* **19**, 4451 (2023), PMID: 36453538.
 - [20] T. Koker, K. Quigley, and L. Li, Higher order equivariant graph neural networks for charge density prediction, in *NeurIPS 2023 AI for Science Workshop* (2023).
 - [21] 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).
 - [22] Q. Gu, L. Zhang, and J. Feng, Neural network representation of electronic structure from ab initio molecular dynamics, *Science Bulletin* **67**, 29 (2022).
 - [23] J. Nigam, M. J. Willatt, and M. Ceriotti, Equivariant representations for molecular Hamiltonians and N-center atomic-scale properties, *The Journal of Chemical Physics* **156**, 014115 (2022).
 - [24] 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).
 - [25] Q. Gu, Z. Zhouyin, S. K. Pandey, P. Zhang, L. Zhang, and W. E, DeePTB: A deep learning-based tight-binding approach with *ab initio* accuracy (2023), [arXiv:2307.04638 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2307.04638).
 - [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] H. Yu, Z. Xu, X. Qian, X. Qian, and S. Ji, Efficient and equivariant graph networks for predicting quantum Hamiltonian, in *Proceedings of the 40th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 202 (PMLR, 2023) pp. 40412–40424.
 - [28] J. C. Snyder, M. Rupp, K. Hansen, K.-R. Müller, and K. Burke, Finding density functionals with machine learning, *Phys. Rev. Lett.* **108**, 253002 (2012).
 - [29] R. Meyer, M. Weichselbaum, and A. W. Hauser, Machine learning approaches toward orbital-free density functional theory: Simultaneous training on the kinetic energy density functional and its functional derivative, *Journal of Chemical Theory and Computation* **16**, 5685 (2020).
 - [30] F. Imoto, M. Imada, and A. Oshiyama, Order- N orbital-free density-functional calculations with machine learning of functional derivatives for semiconductors and metals, *Phys. Rev. Res.* **3**, 033198 (2021).
 - [31] K. Ryczko, S. J. Wetzel, R. G. Melko, and I. Tambllyn, Toward orbital-free density functional theory with small data sets and deep learning, *Journal of Chemical Theory and Computation* **18**, 1122 (2022).
 - [32] R. Remme, T. Kaczun, M. Scheurer, A. Dreuw, and F. A.

- Hamprecht, KineticNet: Deep learning a transferable kinetic energy functional for orbital-free density functional theory, *The Journal of Chemical Physics* **159**, 144113 (2023).
- [33] H. Zhang *et al.*, Overcoming the barrier of orbital-free density functional theory for molecular systems using deep learning, *Nature Computational Science* [10.1038/s43588-024-00605-8](https://doi.org/10.1038/s43588-024-00605-8) (2024).
- [34] L. Li, S. Hoyer, R. Pederson, R. Sun, E. D. Cubuk, P. Riley, and K. Burke, Kohn-Sham equations as regularizer: Building prior knowledge into machine-learned physics, *Phys. Rev. Lett.* **126**, 036401 (2021).
- [35] M. F. Kasim and S. M. Vinko, Learning the exchange-correlation functional from nature with fully differentiable density functional theory, *Phys. Rev. Lett.* **127**, 126403 (2021).
- [36] Y. Chen, L. Zhang, H. Wang, and W. E, DeePKS: A comprehensive data-driven approach toward chemically accurate density functional theory, *Journal of Chemical Theory and Computation* **17**, 170 (2021), pMID: 33296197.
- [37] J. Kirkpatrick *et al.*, Pushing the frontiers of density functionals by solving the fractional electron problem, *Science* **374**, 1385 (2021).
- [38] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. von Lilienfeld, Quantum chemistry structures and properties of 134 kilo molecules, *Scientific Data* **1**, 140022 (2014).
- [39] S. Bai, J. Z. Kolter, and V. Koltun, Deep equilibrium models, in *Advances in Neural Information Processing Systems*, Vol. 32 (2019).
- [40] J. H. Van Lenthe, R. Zwaans, H. J. J. Van Dam, and M. F. Guest, Starting SCF calculations by superposition of atomic densities, *Journal of Computational Chemistry* **27**, 926 (2006).
- [41] A. Duval *et al.*, A hitchhiker’s guide to geometric GNNs for 3D atomic systems (2024), [arXiv:2312.07511 \[cs.LG\]](https://arxiv.org/abs/2312.07511).
- [42] M. Geiger *et al.*, *Euclidean neural networks: e3nn* (2022).
- [43] Y.-L. Liao, B. M. Wood, A. Das, and T. Smidt, EquiformerV2: Improved equivariant transformer for scaling to higher-degree representations, in *The Twelfth International Conference on Learning Representations* (2024).
- [44] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. u. Kaiser, and I. Polosukhin, Attention is all you need, in *Advances in Neural Information Processing Systems*, Vol. 30 (2017).
- [45] S. Passaro and C. L. Zitnick, Reducing SO(3) convolutions to SO(2) for efficient equivariant GNNs, in *Proceedings of the 40th International Conference on Machine Learning*, Proceedings of Machine Learning Research, Vol. 202 (2023) pp. 27420–27438.
- [46] X. Zhang and G. K.-L. Chan, Differentiable quantum chemistry with PySCF for molecules and materials at the mean-field level and beyond, *The Journal of Chemical Physics* **157**, 204801 (2022).
- [47] E. Engel, *Density Functional Theory* (Springer, 2011) p. 63.
- [48] C. Chen and S. P. Ong, A universal graph deep learning interatomic potential for the periodic table, *Nature Computational Science* **2**, 718 (2022).
- [49] B. Deng, P. Zhong, K. Jun, J. Riebesell, K. Han, C. J. Bartel, and G. Ceder, CHGNet as a pretrained universal neural network potential for charge-informed atomistic modelling, *Nature Machine Intelligence* **5**, 1031 (2023).
- [50] 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).
- [51] D. Zhang *et al.*, DPA-2: Towards a universal large atomic model for molecular and material simulation (2023), [arXiv:2312.15492 \[physics.chem-ph\]](https://arxiv.org/abs/2312.15492).
- [52] H. Yang *et al.*, MatterSim: A deep learning atomistic model across elements, temperatures and pressures (2024), [arXiv:2405.04967 \[cond-mat.mtrl-sci\]](https://arxiv.org/abs/2405.04967).
- [53] L. A. Burns *et al.*, The BioFragment Database (BFDdb): An open-data platform for computational chemistry analysis of noncovalent interactions, *The Journal of Chemical Physics* **147**, 161727 (2017).
- [54] S. Marqusee and R. L. Baldwin, Helix stabilization by Glu...Lys+ salt bridges in short peptides of de novo design., *Proceedings of the National Academy of Sciences* **84**, 8898 (1987).
- [55] B. I. Dunlap, N. Rösch, and S. Trickey, Variational fitting methods for electronic structure calculations, *Molecular Physics* **108**, 3167 (2010).
- [56] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, *et al.*, PyTorch: An imperative style, high-performance deep learning library, in *Advances in Neural Information Processing Systems*, Vol. 32 (2019).
- [57] P. Pulay, Improved SCF convergence acceleration, *Journal of Computational Chemistry* **3**, 556 (1982).
- [58] N. D. Woods, M. C. Payne, and P. J. Hasnip, Computing the self-consistent field in Kohn–Sham density functional theory, *Journal of Physics: Condensed Matter* **31**, 453001 (2019).
- [59] S. J. Prince, *Understanding Deep Learning* (The MIT Press, 2023) p. 144.
- [60] Q. Sun *et al.*, PySCF: the Python-based simulations of chemistry framework, *WIREs Computational Molecular Science* **8**, e1340 (2018).
- [61] J. P. Perdew, K. Burke, and M. Ernzerhof, Generalized gradient approximation made simple, *Phys. Rev. Lett.* **77**, 3865 (1996).
- [62] B. P. Pritchard, D. Altarawy, B. Didier, T. D. Gibson, and T. L. Windus, New basis set exchange: An open, up-to-date resource for the molecular sciences community, *Journal of Chemical Information and Modeling* **59**, 4814 (2019), pMID: 31600445.
- [63] G. L. Stoychev, A. A. Auer, and F. Neese, Automatic generation of auxiliary basis sets, *Journal of Chemical Theory and Computation* **13**, 554 (2017).
- [64] M. Fey and J. E. Lenssen, Fast graph representation learning with PyTorch Geometric, in *ICLR Workshop on Representation Learning on Graphs and Manifolds* (2019).
- [65] S. Elfving, E. Uchibe, and K. Doya, Sigmoid-weighted linear units for neural network function approximation in reinforcement learning, *Neural Networks* **107**, 3 (2018), special issue on deep reinforcement learning.