# A symmetry-preserving and transferable representation for learning the Kohn-Sham density matrix

Liwei Zhang[1], Patrizia Mazzeo[2], Michele Nottoli[3], Edoardo Cignoni[2], Lorenzo Cupellini[2], and Benjamin Stamm[3]

[1]Institut für Geometrie und Praktische Mathematik, RWTH Aachen University, Templergraben 55, 52062 Aachen, Germany

[2]Dipartimento di Chimica e Chimica Industriale, Università di Pisa, 56124 Pisa, Italy

[3]Universität Stuttgart, Institute of Applied Analysis and Numerical Simulation, Pfaffenwaldring 57, 70569 Stuttgart, Germany

## Abstract

The Kohn-Sham (KS) density matrix is one of the most essential properties in KS density functional theory (DFT), from which many other physical properties of interest can be derived. In this work, we present a parameterized representation for learning the mapping from a molecular configuration to its corresponding density matrix using the Atomic Cluster Expansion (ACE) framework, which preserves the physical symmetries of the mapping, including isometric equivariance and Grassmannianity. Trained on several typical molecules, the proposed representation is shown to be systematically improvable with the increase of the model parameters and is transferable to molecules that are not part of and even more complex than those in the training set. The models generated by the proposed approach are illustrated as being able to generate reasonable predictions of the density matrix to either accelerate the DFT calculations or to provide approximations to some properties of the molecules.

## 1 Introduction

Computational chemistry often deals with many quantum mechanical calculations repeated on the same system or on similar systems. Examples are molecular dynamics (MD) simulations, repeated calculations on a statistical sampling, geometry optimizations, or even scans along some interesting coordinate. In all these cases, the results of already performed calculations can be used to fit a machine learning model able to predict energies and properties of subsequent calculations[1].

In the context of quantum chemistry, machine learning models have been used to fit properties, for example, the energy[2–7] and atomic forces[8–10], or to predict more fundamental quantities like the Hamiltonian[11–16] and the wavefunction[17, 18]. Machine learning models have also been used directly as interatomic potentials for molecular dynamics simulations of a variety of systems[19–23].

Among the more fundamental quantities, various methods have been proposed to fit the electronic density matrix. These either target the electronic density in real space[24–35], or they target the corresponding electronic density matrix in a basis[3, 36–40]. Fitting the electronic density is a powerful strategy, as the density can then be directly used to compute different observables. Other strategies often need to train an ad hoc model for each property of interest, but multiple properties are often required for answering a scientific question.

While being less general than fitting in real space, fitting the density in a suitable basis removes any projection error and removes the barrier between the predicted density and the quantum chemistry package of choice, which can be used to compute the properties of interest. Fitting the electronic density matrix provides two additional advantages. First, in the context of Hartree-Fock or density functional theory (DFT), an electronic density matrix can simply be used as an initial guess for the upcoming self-consistent field (SCF) calculation, instead of directly using it to access properties. This hybrid approach represents a middle ground between a full SCF calculation and directly using the density to access the properties: it retains the full accuracy of a normal SCF procedure, but at a reduced computational cost[40]. The better the guess, the more efficient is the full accuracy model. Second, given a predicted electronic density matrix $D$, it is possible to assemble the corresponding Fock / Kohn-Sham matrix $F(D)$, and the commutator $FD - DF$ provides a measure of how accurate the prediction is, thus providing the opportunity to either discard low quality predictions or mark the data points with the worst predictions, which is useful in active learning strategies[41].

However, the required mapping from the molecular configurations (coordinates and atomic numbers) to the corresponding density matrices is in general complicated and of high dimensionality, and therefore difficult to learn. The fitting problem becomes treatable by introducing appropriate molecular descriptors, which take into account physical knowledge such as invariance or equivariance of the target property. In this way, the required design space can be reduced. More specifically, the descriptors are functions of the molecular parameters satisfying a series of requisites: they are desired to be injective (exactly or approximately) and economical to compute, and should capture the aforementioned symmetries of the target property. Depending on the order in which the various invariances are introduced, different classes of descriptors are obtained. A possible strategy is to compute translationally and rotationally invariant functions of the coordinates, and only then introducing the permutational invariance. Examples of descriptors of this kind are permutationally invariant polynomials (PIPs)[42] and their variant atomic permutationally invariant polynomials (aPIPs)[43]. Alternatively, it is possible to compute functions that are permutationally and translationally invariant, thereafter enforcing the rotational invariance. This is the strategy followed by the smooth overlap of atomic positions (SOAP)[44], by the atomic cluster expansions (ACE)[45] and by the Behler-Parrinello descriptors[46, 47]. The ACE descriptors are of particular interest as they include, in principle, many-body terms of arbitrarily high order, and are cheaper to compute than other alternatives[48]. Notably, it has also been generalized to capture the equivariant properties[12, 49].

In this contribution, we propose a strategy that combines the strengths of the equivariant ACE descriptors with the flexibility of fitting the electronic density matrix in a basis, which respects the intrinsic properties of the density matrix. Specifically, the electronic density matrix is approximated with a linear regression in an ACE basis, similarly to the previous work on self-consistent Hamiltonians from one of the authors[12]. The strategy is used to train both specific models (that is, trained on a single molecule) and unified models (trained on multiple molecules). The resulting models are systematically improvable and, in the case of unified models, also transferrable to unseen molecules, provided that they share some chemical similarity with the training set. Both the specific and unified models can be used to reduce the number of SCF iterations or to directly predict the properties of interest.

## 2 Methods

### 2.1 Density Matrix

Let $\mathbf{R} = \{(Z_I, \boldsymbol{r}_I)\}_{I=1}^{N_{\mathrm{at}}} := \{\sigma_I\}_{I=1}^{N_{\mathrm{at}}}$ be a molecular configuration consisting of $N_{\mathrm{at}}$ atoms and $N$ (valence) electron-pairs, with $Z_I \in \mathbb{N}$ and $\boldsymbol{r}_I \in \mathbb{R}^3$ characterizing the atomic number and the position of the $I$-th atom, respectively. The union of all $Z_I$ characterizes the different elements in this given system, whose cardinality will be denoted by $n$. Other properties of the atoms can potentially also be included in $\sigma_I$ but that would go beyond the scope of this paper. If an $N_g(\geq N)$ dimensional discretization space is adopted, in which the orbitals are approximated, then the corresponding KS equation will read as

$$F_\mathbf{R}[D_\mathbf{R}]C_\mathbf{R} = S_\mathbf{R}C_\mathbf{R}E_\mathbf{R}.$$

where $F_\mathbf{R}$ and $S_\mathbf{R} \in \mathbb{R}^{N_g \times N_g}$ are the discretized KS operator (Hamiltonian) and the overlap matrix respectively, $C_\mathbf{R} \in \mathbb{R}^{N_g \times N}$ represents the coefficients of the orbitals in a given basis, $D_\mathbf{R} = C_\mathbf{R}C_\mathbf{R}^T$ is the density matrix, the main object of this paper, and $E_\mathbf{R}$, a diagonal matrix of order $N$, which contains the $N$ corresponding eigenvalues of the system (sorted in ascending order). Without loss of generality, we assume that $S_\mathbf{R} = I_{N_g}$, by adopting the Löwdin orthonormalization[50] if necessary. Under this setting, the above equation (2.1) can be rewritten as

$$F_\mathbf{R}[D_\mathbf{R}]C_\mathbf{R} = C_\mathbf{R}E_\mathbf{R}. \tag{1}$$

If $C_\mathbf{R}$ is chosen to be orthonormal, then $D_\mathbf{R}$ should lie in the following manifold

$$\mathcal{G}_{N_g}^N := \{D \in \mathbb{R}^{N_g \times N_g} : D^2 = D^T = D, \mathrm{tr}(D) = N\}, \tag{2}$$

which is equivalent to an $(N_g, N)$-Grassmann manifold, hence our notation.

In the context of linear combinations of atomic orbitals (LCAO), the discretization space is spanned by a set of atomic orbitals $\{\phi_{I\alpha}\}_{I \in \{1,\ldots,N_{\mathrm{at}}\}, \alpha \in \mathcal{I}_{Z_I}}$ where $\mathcal{I}_{Z_I}$ is the index set of the atomic orbitals centered at the $I$-th atom, depending only on the atomic number $Z_I$. The density matrix, consequently, has elements that are invariant under translations of the system or permutations of the index of the atoms, and is equivariant under rotations and reflections of the whole system. It can be divided into several subblocks that have respective symmetries and can be learned independently. In Figure 1, we take $C_3H_4O$ as an example to illustrate the block structure of the density matrix. Note that a similar strategy is used in Ref. [12] is here extended to a case with multiple different elements. The detailed derivation of the block-wise equivariance of the density matrix can be found in Appendix A. For simplicity, we omit the subscript $\mathbf{R}$ in $D_\mathbf{R}$ hereafter when no ambiguity is introduced.

As can be seen from Figure 1(a), there are two types of blocks appearing in the density matrix, the diagonal blocks and the off-diagonal ones. Depending on contexts, they are also called *onsite* and *offsite*, or *homo-* and *hetero-orbital*, respectively. We will use the terms *onsite* and *offsite* throughout this paper. For the targeted systems having $n$ different elements, there exist, $n(n + 3)/2$ matrix-valued functions which correspond to interactions of distinct elements (there can be a lack thereof, for instance, when the system has only one oxygen atom, there is no O-O interaction). As such, the block of the density matrix corresponding to the $I$-th and $J$-th atom is of the form

$$D_{IJ} = \begin{cases} D_{Z_I}(\mathbf{R}_I), & I = J, \\ D_{(Z_I, Z_J)}(\mathbf{R}_{IJ}), & I \neq J, \ Z_I \leq Z_J, \\ D_{(Z_J, Z_I)}(\mathbf{R}_{JI})^T, & I \neq J, \ Z_I > Z_J, \end{cases} \tag{3}$$
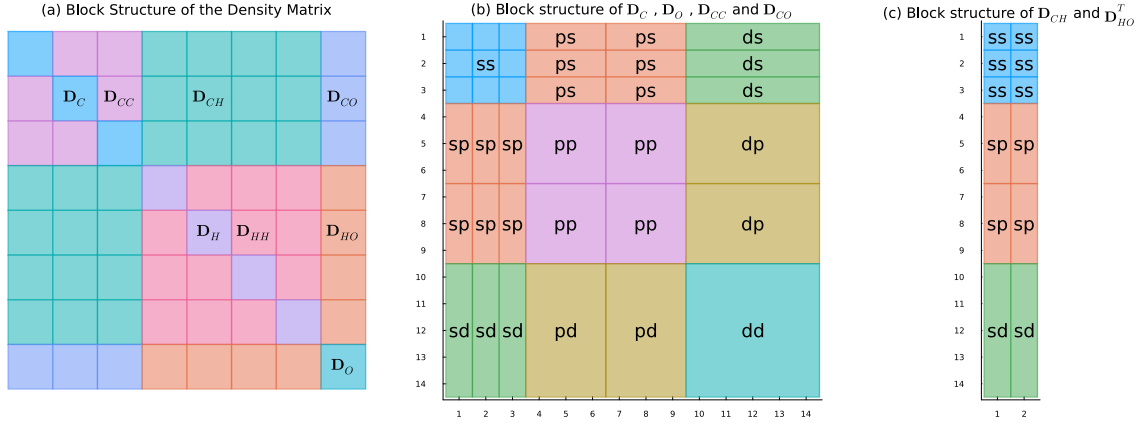
Figure 1: Block structure of the density matrix of the $C_3H_4O$ molecule, where the atomic basis 6-31G(d) is used (the block structure of $D_{HH}$ is omitted as it is just a 2 by 2 matrix consisting of 4 ss-blocks).

where $\mathbf{R}_I$ and $\mathbf{R}_{IJ}$ are global configurations, translated in order to be centered at the $I$-th atom or at some specific point of the $(I, J)$-th bond (the line segment that connects the two atoms) respectively. Note that the *Hermitian-ness* of the density matrix (*i.e.*, $D = D^T$) is used in the last line of (3). To unify the notations, we sometimes use the convention that $\mathbf{R}_{II} = \mathbf{R}_I$.

In addition, each of such matrix-valued functions can be further divided into completely independent sub-blocks corresponding to the atomic orbitals, as shown in Figure 1(b) and 1(c), *i.e.*,

$$D_{Z_I}(\mathbf{R}_I) = [D_{Z_I}^{\alpha\beta}(\mathbf{R}_I)]_{\alpha,\beta\in\mathcal{I}_{Z_I}},$$
$$D_{(Z_I,Z_J)}(\mathbf{R}_{IJ}) = [D_{(Z_I,Z_J)}^{\alpha\beta}(\mathbf{R}_{IJ})]_{\alpha\in\mathcal{I}_{Z_I},\beta\in\mathcal{I}_{Z_J}}. \tag{4}$$

Our target then becomes the unified functionals $D_{Z_i/(Z_i,Z_j)}^{\alpha\beta}$ for various atomic numbers $Z_i$, $Z_j$, which have their distinct isometric equivariance and can be dealt with separately. Such structure of the density matrix forms the foundation of the transferability and parallelizability of the proposed method. We refer readers to Appendix A for a detailed discussion.

In practice, it is commonly believed that only atoms near the central atoms make substantial contributions to the corresponding part of the density matrix (also known as the nearsightedness of the object). As a result, certain cutoff strategies are often used when constructing the input atomic environment. We illustrate our truncation strategy in Figure 2, where the particles in the red and blue spheres form the *onsite* and *offsite* environments $\mathbf{R}_I$ and $\mathbf{R}_{IJ}$, respectively. The rigorous definitions of the truncated $\mathbf{R}_I$ and $\mathbf{R}_{IJ}$ are provided in Appendix A.

Here, we assume that atoms of the same element are discretized by the same set of bases. However, the method proposed in this work can potentially be extended to the more general setting where atoms of the same element are assigned different atomic orbital basis functions, simply by artificially treating them as having different atom types.

## 2.2  Representation of the density matrix

One of the goals of this paper is to provide a faithful representation of the density matrix, respecting its inherent physical symmetries as much as possible to facilitate its learning. To this end, we adopt the equivariant ACE descriptors[12, 45] to approximate the functions $D_\bullet$, where the symbol $\bullet$ can be one of the indices appearing in the right hand side of (4).

For each function $D_\bullet$, there exists a set of ACE bases $\{\mathcal{B}_{\bullet,v}\}_v$ as functions of the local environments
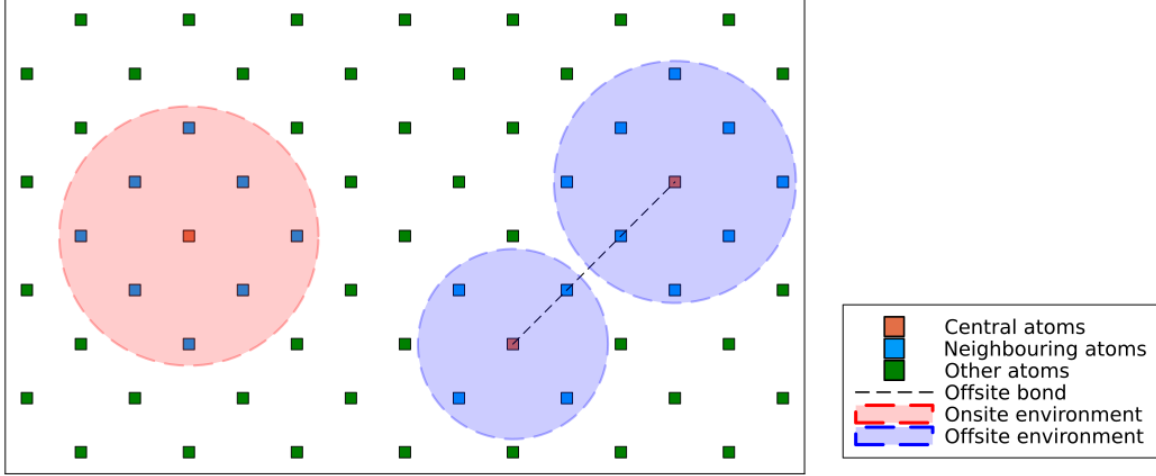
Figure 2: Illustration of the local atomic environments. Atoms in the red sphere form an *onsite* local configuration, with the radius of the sphere being the onsite cutoff. Atoms in the blue spheres form the *offsite* local configuration, with the radii of the sphere being the offsite cutoffs that can be chosen differently from the *onsite* one.

$\mathbf{R}_I$ or $\mathbf{R}_{IJ}$, which has the same isometric equivariance as $D_\bullet$ and asymptotically spans the function space to which $D_\bullet$ belongs[51]. The size of the basis is determined merely by two parameters: (i) the correlation order $\nu$, which corresponds to the body order in physics (up to a constant, precisely, it is 1 and 2 less for *onsite* and *offsite*, respectively) and (ii) the maximum polynomial degree $d_{\max}$, indicating the resolution of the one-particle basis. Additional details about these two parameters and the corresponding ACE basis, as well as the definition of the one-particle basis, can be found in Appendix B. Given the basis $\{\mathcal{B}_{\bullet,v}\}_v$, we can approximate each function $D_\bullet$ by a linear combination of $\{\mathcal{B}_{\bullet,v}\}_v$:

$$D_\bullet \approx \sum_v c_{\bullet,v} \mathcal{B}_{\bullet,v}. \tag{5}$$

To predict the corresponding sub-block of the density matrix, the only thing left now is to estimate the coefficient $c_{\bullet,v}$ for all possible indices $\bullet$.

## 2.3 Parameter estimation

Suppose that a dataset is given of the form $\{(\mathbf{R}^{(k)}, D^{(k)})\}_k$, where $k$ is the index of the data point, $\mathbf{R}^{(k)}$ is the $k$-th (global) molecular configuration and $D^{(k)}$ is the corresponding density matrix. The dataset can first be transformed into sets of local atomic clusters and their corresponding portions of the density matrix, according to the atomic number of each atom in $\mathbf{R}^{(k)}$, as

$$\{(\mathbf{R}_{IJ}^{(k)}, D_\bullet^{(k)})\}_{k,I,J},$$

where the subscript $\bullet$ has the same meaning as that in the preceding subsection. These sets are then used to train the coefficients of the corresponding models (5) independently. One of the most direct ways to estimate the coefficients is through a least squares approach, that is, they are determined by minimizing

$$L(\boldsymbol{c}_\bullet) = \sum_{k,I,J} \|D_\bullet^{(k)} - \sum_v c_{\bullet,v} \mathcal{B}_{\bullet,v}(\mathbf{R}_{IJ}^{(k)})\|^2 + \lambda \|\Gamma_{\mathcal{B}_\bullet}\|^2, \tag{6}$$

where $\boldsymbol{c}_\bullet = \{c_{\bullet,v}\}_v$, $\Gamma_{\mathcal{B}_\bullet}$ refers to some Tikhonov regularizer that can be customized with respect to

the basis $\mathcal{B}_\bullet$, and $\lambda$ is a regularization parameter. Throughout our experiments, we use $\lambda = 10^{-4}$ and choose $\Gamma$ to be the smooth prior given in Ref. [12]. Once an (approximate) minimizer is found, it is possible to provide an approximation of the ground state density matrix through (5) for any given configuration $\mathbf{R}$ as long as its chemical composition of elements does not go beyond that of the training set.

## 2.4 Retraction

The construction of the ACE basis as well as our representation (5) ensure that the predicted density matrix $D$ has the desired isometric-equivariance. However, it does not guarantee that the prediction belongs to the Grassmann manifold (2), *i.e.* that it is a valid density matrix. To bring this restriction back, we introduce a retraction operator that maps $D$ to the manifold.

Since $D$ is a real symmetric matrix of size $N_g \times N_g$, its eigenvalue decomposition can be written as

$$D = U_D \Sigma_D U_D^T,$$

where $U_D \in \mathbb{R}^{N_g \times N_g}$ is unitary and $\Sigma_D$ is a diagonal matrix containing all the eigenvalues of $D$, sorted in descending order. The retraction is then defined as

$$\mathcal{P}(D) = U_D E_{N_g}^N U_D^T, \tag{7}$$

where

$$E_{N_g}^N = \begin{bmatrix} I_N & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} \in \mathbb{R}^{N_g \times N_g}.$$

We mention that after applying the retraction, $\mathcal{P}(D)$ remains isometric equivariant and is the nearest element on the Grassmann manifold to $D$. A proof of this statement, as well as the well-definedness of $\mathcal{P}$, is given in Appendix C.

The following schematic (Figure 3) summarizes the whole procedure of our density matrix prediction scheme, from which we can see that the whole procedure, apart from the retraction step, is essentially parallelizable, including training and prediction.
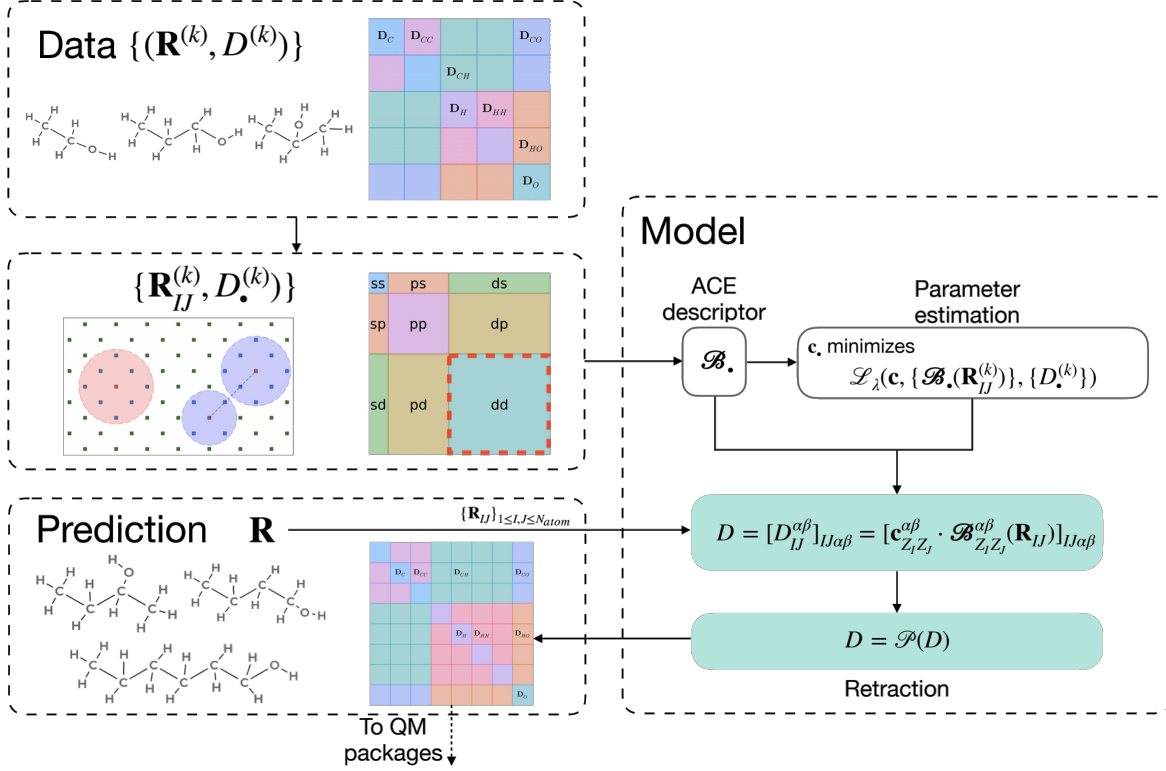
Figure 3: A schematic of the process of learning the density matrix described in this paper. Here, the loss function $\mathcal{L}_\lambda$ is defined as (6). The molecules for which the density matrices are predicted can be different from those in the training set. Finally, the predicted density matrix can be directly sent to some quantum chemistry packages of choice for further operations.

It is worth noting that the proposed scheme does not require specific data origination, but just requires a consistent atomic basis discretization across the data set (or more abstractly, it requires only the equivariance of the data). The resulting density matrix can be used in many different application scenarios (cf. Section 4). We remark that a similar strategy is used in [12] to fit the self-consistent Hamiltonian (*i.e.* Kohn-Sham) matrix for periodic crystal systems. We compared the aforementioned approach of learning the Hamiltonian matrix therein with the method proposed in this work, whose results can be found in Appendix D.

## 3  Results and discussion

Our approach of learning the density matrix was tested on various systems described with DFT. We designed tests of increasing complexity to validate the approach and test its performance. The various tests required both training and test datasets generated using a standard protocol, whose details can be found in Section 3.1. In general, the performance of the various models was assessed by computing the Root-Mean-Square-Error(RMSE) between the references ($\{D_{\text{ref}}^{(k)}\}_{k=1}^{N_{\text{data}}}$) and predicted density matrices ($\{D_{\text{pred}}^{(k)}\}_{k=1}^{N_{\text{data}}}$) as

$$\text{RMSE}_D = \sqrt{\frac{\sum_{k=1}^{N_{\text{data}}} \|D_{\text{ref}}^{(k)} - D_{\text{pred}}^{(k)}\|_F^2}{\sum_{i=1}^{N_{\text{data}}} \left(N_g^{(k)}\right)^2}}, \tag{8}$$

where $N_g^{(k)}$ is the size of the $k$-th density matrix.

| Molecule | N frames |
|---|---:|
| Acetaldehyde | 10000 |
| Acrolein | 10000 |
| Aniline | 10000 |
| o-Toluidine | 10000 |
| m-Toluidine | 10000 |
| Benzene | 100 |
| Toluene | 100 |
| Phenol | 100 |
| Benzaldehyde | 100 |
| p-Toluidine | 100 |
| 1-Propanol | 10000 |
| 1-Butanol | 10000 |
| 2-Butanol | 10000 |
| 1-Hexanol | 10000 |
| Ethanol | 100 |
| 2-Propanol | 100 |
| 2-Hexanol | 100 |
| 1-Heptanol | 100 |

Table 1: Datasets used in this work. Different molecules are grouped according to their chemical class: aldehydes, aromatics, alcohols. The level of theory is DFT $\omega$B97XD/6-31G(d).

## 3.1 Data preparation

Each dataset was prepared with the same protocol, consisting of a sampling step and a QM calculation step. In the sampling step, each molecule was optimized with DFT B3LYP/6-31G in water, treated with IEFPCM[52], and solvated with an octahedral box of TIP3P waters[53], extending up to 35 Å from the molecule. The solvent was then minimized while keeping the molecule fixed. Thereafter, the whole system was heated from 0 K to 100 K in a 5 ps NVT simulation and from 100 K to 310 K in a 100 ps NPT simulation. The QM/MM production simulation was then run for 150 ps in the NVT ensemble, using the Langevin thermostat. The molecule was treated at the DFTB3 level of theory[54] with 30b-3-1 parameters[55, 56]. Electrostatic interactions were treated with PME[57], using a 10 Å cutoff to divide the direct and reciprocal space. The first 50 ps of production trajectory were discarded. All simulations were performed with AMBER[58].

In the second step, QM DFT calculations were run on equally spaced frames along the trajectory, using the $\omega$B97XD/6-31G(d) level of theory, and enforcing the use of spherical atomic basis functions. The training-and-testing dataset comprises nine organic molecules featuring different functional groups, whereas the test-only datasets comprise nine similar but distinct molecules. For training-and-testing datasets, the calculations were run on 10000 frames, whereas for test only datasets, the calculations were run only on 100 frames. All the calculations for the datasets were performed using Gaussian 16[59].

The datasets were generated by storing the coordinates, the overlap matrices, the coefficient matrices, the Kohn-Sham matrices, as well as metadata such as the list of atoms and the calculation level in HDF5 binary files. An overview of the datasets is reported in Table 1. All the datasets are available in the corresponding archive for the sake of reproducibility.

## 3.2 Specific models

To assess the method we proposed, we first show that it generates systematically improvable results, so that we can refrain from fine-tuning the choice of model parameters (*i.e.*, correlation orders $\nu$ and maximum polynomial degrees $d_{\max}$). To this end, we show the results of molecule-specific models,

each of which is trained with the geometries of only one molecule. For each training molecule (those with 10000 available frames in Table 1) apart from acrolein and butanol, we use the first 3000 frames or a subset thereof for training, and test the resulting models on the rest among the 10000 frames. For the two exceptional molecules, we sampled alternatively/tertiarily from the first 6000/9000 frames, respectively. In any case, less than 30% of available data points are included in the training. More details on the selection of the dataset can be found in Table S1 and Figure S5 in the Supporting Information (SI). The training dataset is then used to train the models with $\nu = 2$, 3 and $4 \leq d_{\max} \leq 8$. For the local truncation, we choose the cutoffs (the radii of the three spheres displayed in Figure 2, which are chosen uniformly in this work regardless of the atomic types) to be $4.0\mathring{A}$ and $6.5\mathring{A}$, respectively. For the sake of simplicity, we only show the results of aniline and propanol. In the following Figure 4, the x-axis stands for the degree $d_{\max}$ used to generate the descriptors, whereas the y-axis is the element-wise RMSE of the predicted density matrix $D_{\mathrm{pred}}$ defined in (8), displayed in the logarithmic scale.
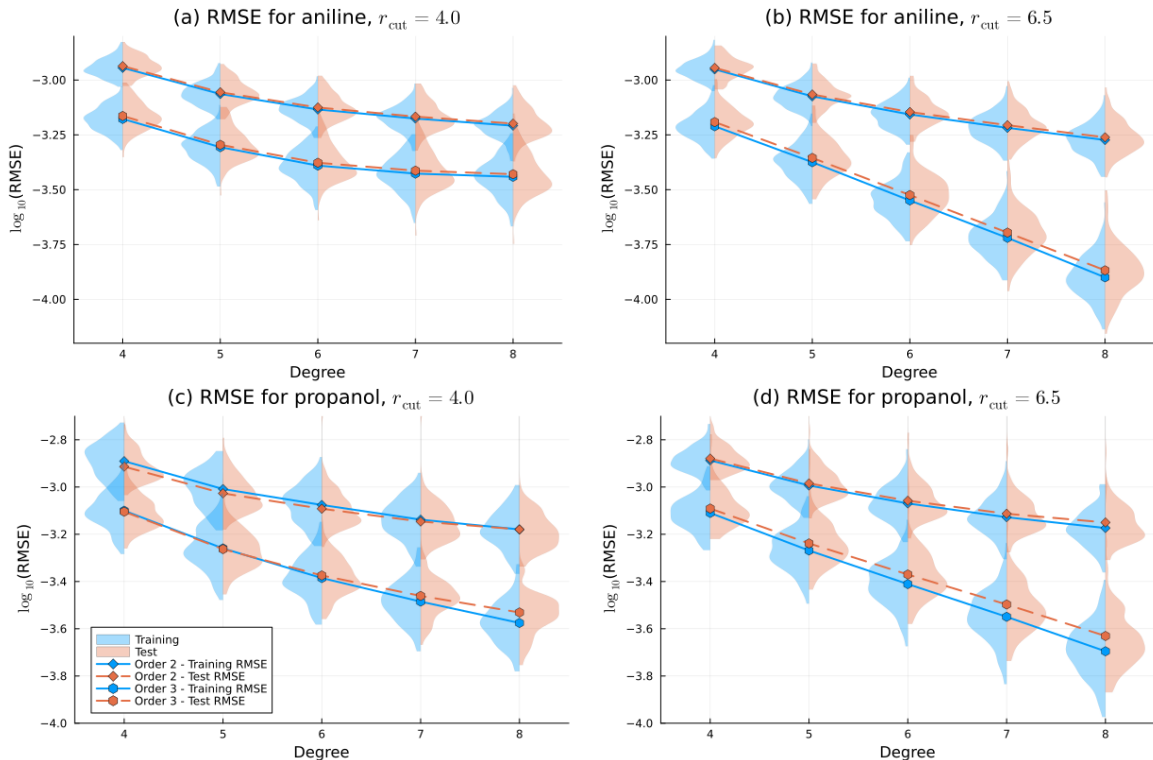


Figure 4: The RMSEs (8) of the predicted density matrices for: (a) aniline with $r_{\mathrm{cut}} = 4.0$, (b) aniline with $r_{\mathrm{cut}} = 6.5$, (c) propanol with $r_{\mathrm{cut}} = 4.0$ and (d) propanol with $r_{\mathrm{cut}} = 6.5$, obtained by the corresponding specific models with respect to different degrees $d_{\max}$ and for various orders $\nu$. The solid and dashed lines refer to the average training and test set errors, and the shaded areas show the distribution of the errors for the corresponding models.

As illustrated in Figure 4, the training and test set RMSEs align with each other nicely, and are nearly normally distributed, even with less than 30% of the data being involved in the training phase. The similarity of the training/test-set errors shows little overfitting in the training, validating the effects of the regularization we used (cf. equation (6)). Comparing the two truncation parameters, we find that the models trained with larger cutoff radii can give better results, but the differences are not too significant. It reaches the smallest average RMSE at around $10^{-4}$ and $2 \cdot 10^{-4}$ for Aniline and Propanol, respectively, which corresponds to a relative error in the whole density matrix of around 0.2% to 0.35% (note that $\|D\|_F^2 = N$ the number of electrons, so the relative error is simply given by

| Molecule | Specific Model | Unified Model | Unified Model-A |
|---|---|---|---|
| Acetaldehyde | $4.416 \cdot 10^{-5}$ | $3.278 \cdot 10^{-4}$ | $3.299 \cdot 10^{-4}$ |
| Acrolein | $2.514 \cdot 10^{-4}$ | $5.028 \cdot 10^{-4}$ | $5.081 \cdot 10^{-4}$ |
| Aniline | $4.300 \cdot 10^{-4}$ | $4.868 \cdot 10^{-4}$ | $4.876 \cdot 10^{-4}$ |
| o-Toluidine | $5.430 \cdot 10^{-4}$ | $5.962 \cdot 10^{-4}$ | $5.940 \cdot 10^{-4}$ |
| m-Toluidine | $5.384 \cdot 10^{-4}$ | $5.824 \cdot 10^{-4}$ | $5.822 \cdot 10^{-4}$ |
| Benzene* | - | $4.058 \cdot 10^{-4}$ | $3.646 \cdot 10^{-4}$ |
| Toluene* | - | $6.369 \cdot 10^{-4}$ | $5.980 \cdot 10^{-4}$ |
| Phenol** | - | $4.809 \cdot 10^{-3}$ | $6.770 \cdot 10^{-4}$ |
| Benzaldehyde** | - | $4.129 \cdot 10^{-3}$ | $1.201 \cdot 10^{-3}$ |
| p-Toluidine** | - | $2.840 \cdot 10^{-3}$ | $6.293 \cdot 10^{-4}$ |
| 1-Propanol | $3.049 \cdot 10^{-4}$ | $4.427 \cdot 10^{-4}$ | $4.444 \cdot 10^{-4}$ |
| 1-Butanol | $4.510 \cdot 10^{-4}$ | $4.921 \cdot 10^{-4}$ | $4.934 \cdot 10^{-4}$ |
| 2-Butanol | $9.173 \cdot 10^{-4}$ | $1.494 \cdot 10^{-3}$ | $1.509 \cdot 10^{-3}$ |
| 1-Hexanol | $1.031 \cdot 10^{-3}$ | $5.324 \cdot 10^{-4}$ | $5.314 \cdot 10^{-4}$ |
| Ethanol* | - | $9.644 \cdot 10^{-4}$ | $8.999 \cdot 10^{-4}$ |
| 2-Propanol* | - | $7.701 \cdot 10^{-4}$ | $7.480 \cdot 10^{-4}$ |
| 2-Hexanol* | - | $7.384 \cdot 10^{-4}$ | $7.353 \cdot 10^{-4}$ |
| 1-Heptanol* | - | $8.896 \cdot 10^{-4}$ | $8.980 \cdot 10^{-4}$ |

Table 2: The test set RMSEs obtained by the (3,8)-models trained on different datasets ($r_{\mathrm{cut}} = 4.0$). There is no specific model for test-only molecules, hence some dashes in the first column. In particular, the test molecules with a superscript $*$ are not included in the training process at all, and those with $**$ are involved in the training of the augmented model, Unified Model-A, with only 10 frames each included in training.

$\|D - D_{\mathrm{pred}}\|_F / \sqrt{N}$). Furthermore, the RMSE, in all cases, monotonically decreases with increasing ACE basis size, which implies that smaller errors can be expected simply by increasing either of the two model parameters.

## 3.3 Unified model

We then extend the training set to configurations from several distinct molecules. More specifically, the training set used in this subsection consists of a total of 2700 frames (300 frames evenly sampled from the first 9000 frames of each of the 9 training molecules). This extended training set is used to train the largest model mentioned above ($\nu = 3$, $d_{\max} = 8$) in order to obtain a unified model, which is then tested on the last 1000 configurations of the training molecules. We compare the average test-set RMSE obtained by the unified model with those obtained by the corresponding specific models of the same size, whose results can be found in the first two columns of Table 2. We note that to make the model more capable of capturing the similarity only of local chemical structures across different molecules, we choose $r_{\mathrm{cut}} = 4.0$ for the unified model. The results for the unified model with $r_{\mathrm{cut}} = 6.5$ are given in Table S1 in the SI, which indeed shows that a smaller cutoff is favorable in terms of the generalizability of the models.

As indicated in the first two columns of Table 2, the unified model overall achieves a performance

comparable to the specific models trained on each molecule independently. This indicates that the model is able to gather information from distinct molecules, and offers the advantage of predicting the density matrices for multiple systems within a single model. Whereas the RMSE of the unified model for acetaldehyde is slightly higher, the unified model performs even better than the specific one for 1-hexanol. This demonstrates that the models generated by the proposed method are able to predict the density matrices of diverse molecular systems, despite their inherent structural dissimilarities, as long as a certain number of their frames are included in the training set. We remark that the training set of the unified model comprises a slightly smaller number of configurations than that of the specific models, and was sampled without particular strategies.

For reference, we also trained a unified model for the alcohols, which is a simpler task compared to the unified one we introduce in this subsection. The results for the Unified Alcohols Model can be found in table S1 in the SI.

## 3.4 Transfer to other molecules

As a more challenging task, we directly use the unified model obtained in subsection 3.3 to predict the density matrices of some molecules beyond the training set, which may be larger or more complex. Specifically, we test the models on those molecules having 100 frames in Table 1. The corresponding average test set RMSEs are reported also in the second column of Table 2, from which we observe that the unified model can provide faithful predictions of the density matrices for benzene, toluene and all the alcohols, for which the obtained errors are similar to those within the training process. However, the unified model struggles with giving good predictions for phenol, benzaldehyde, and is less well in predicting the density matrices for p-toluidine. The poor prediction on these molecules can be attributed to the lack of information in the training set. Indeed, while the dataset includes both carbonyl compounds and alcohols, the chemical behaviour of these functional groups changes significantly when they are bonded to aromatic molecules. Additionally, when an aromatic molecule has two substituents, their effect depends on their relative positions. This explains why, despite the inclusion of o- and m-toluidine in the training set, the model struggles to accurately predict for p-toluidine.

To test whether the weaker performance of the unified model on the three molecules is caused by the limitation of the method itself or just by the training set, we design an augmented training set, which consists of the data points of the previous unified training set, and 10 frames from each of the three molecules, evenly sampled from the first 90 frames (2730 configurations in total). We train the above (3,8)-model with the augmented training set, and obtain a new model Unified Model-A, with the suffix "A" indicating that it is trained with an augmented set. The augmented unified model is again used to predict the density matrices for all the involved molecules, and the corresponding test set RMSEs are listed in the third column of Table 2. The results show that the augmented unified model achieves a higher accuracy for the three molecules with previously critical accuracy, which is similar in magnitude to that of the training molecules, while maintaining a comparable effectiveness for the other systems involved.

The RMSE results presented in this section suggest that the models generated by the proposed method can be uniformly refined simply by increasing the two model parameters. In addition, the proposed unified models can be transferred to the molecules that are not known at the training stage, provided some similarity in the chemical geometries. The performance of the generated models is mainly limited by the design of the training set, rather than the representation itself.

# 4    Applications

In this section, we showcase how the predicted density matrices can be used in some specific scenarios, as extended quality tests for the proposed models.

## 4.1    Accelerating the SCF iterations

A natural application to try is to use the predictions as the initial guesses of the SCF procedure, as it is no matter what not of full accuracy. For each test geometry, we use the proposed models to predict the density matrix and provide it to Gaussian as an initial guess. For these calculations, we used the development version of the Gaussian suite of programs[60]. Communication with Gaussian is possible thanks to the GauOpen open-source library[61]. We compared the number of iterations required to achieve convergence with all our models and with the default guess available in Gaussian. Table 3 reports the average of iterations obtained with the default guess, specific models with $r_{\mathrm{cut}} = 4.0$ and unified models for SCF convergence tolerance $10^{-6}$. The same results for the convergence levels $10^{-7}$ and $10^{-8}$ are reported in Table S4 in the SI. The value within parentheses indicates the percentage of reduction with respect to the default guess. The table shows that, as expected, the specific model achieves the highest reduction in the number of iterations for each molecule, with a maximum for acetaldehyde, where a 44% reduction is observed. We also compared the performance of the specific models with $r_{\mathrm{cut}} = 4.0$ and $r_{\mathrm{cut}} = 6.5$, finding no significant differences (see Table S2 in the SI). On average, the specific models allow us to save three iterations (30%). Moving to the unified model, we observe that for the majority of the molecules in the training set, the predicted density is comparable to that obtained with the respective specific model, with a slightly greater loss of accuracy for the two carbonyl molecules. For what concerns the out-of-sample molecules, the model exhibits good transferability for alcohols, achieving comparable results for both known and unknown molecules. Conversely, the predictions for phenol, p-toluidine, and benzaldehyde were particularly poor, even falling below the accuracy of the default guess. However, as demonstrated by the RMSEs, including just 10 frames in the training set for these three molecules enhances the performance and reduces the number of iterations by around two (Unified Model-A).

It is worth mentioning that the computational time for predicting a density matrix for a given configuration using our model is almost negligible compared to a single SCF iteration. For example, it takes about 112 ms to obtain a predicted density matrix for a propanol molecule using the unified model in a single thread, whereas a single SCF iteration, even carried out on 6 threads, takes an average of 626 ms. Hence, the percentage of reduction is almost exactly the acceleration that we gain.

## 4.2    Predictions of physical properties

The predicted density matrices can also be directly used to derive physical properties of interest, obtaining satisfactory predictions. This was achieved by providing the predicted density matrix as a guess and forcing Gaussian to stop the SCF procedure after a single iteration. Figure 5 illustrates the error in energy, Mulliken charges, dipole moment, and forces with respect to the results obtained from the corresponding converged density matrix. The plot compares the errors obtained using the density matrix predicted with Unified Model-A (pink) and the default guess available in Gaussian (blue), for which we also ran a single SCF iteration for consistency. The same plots for specific models, unified alcohols model and unified model are reported in the SI (see Figures S1, S2 and S3).

Averaging over all molecules, we obtain a mean absolute error (MAE) of 2.7 kcal/mol for energy and 6.5 kcal $\cdot$ mol$^{-1}$ $\cdot$ Å$^{-1}$ for forces. Although they do not achieve chemical accuracy (1 kcal/mol for energies and 1 kcal $\cdot$ mol$^{-1}$ $\cdot$ Å$^{-1}$ for forces) except for the two aldehydes, the predictions can be

| Molecule | Default guess | Specific Model | Unified Model | Unified Model-A |
|---|---|---|---|---|
| Acetaldehyde | $9.4 \pm 0.1$ | $5.2 \pm 0.1$ ($\sim 44\%$) | $7.5 \pm 0.1$ ($\sim 20\%$) | $7.5 \pm 0.1$ ($\sim 20\%$) |
| Acrolein | $10.5 \pm 0.1$ | $7.5 \pm 0.2$ ($\sim 29\%$) | $8.3 \pm 0.2$ ($\sim 21\%$) | $8.3 \pm 0.2$ ($\sim 21\%$) |
| Aniline | $9.9 \pm 0.1$ | $7.2 \pm 0.1$ ($\sim 28\%$) | $7.5 \pm 0.1$ ($\sim 24\%$) | $7.5 \pm 0.1$ ($\sim 24\%$) |
| o-Toluidine | $10.0 \pm 0.0$ | $7.6 \pm 0.1$ ($\sim 24\%$) | $7.8 \pm 0.1$ ($\sim 22\%$) | $7.8 \pm 0.1$ ($\sim 22\%$) |
| m-Toluidine | $10.0 \pm 0.0$ | $7.3 \pm 0.1$ ($\sim 27\%$) | $7.6 \pm 0.1$ ($\sim 24\%$) | $7.6 \pm 0.1$ ($\sim 24\%$) |
| Benzene* | $9.0 \pm 0.0$ | - | $7.4 \pm 0.1$ ($\sim 18\%$) | $7.4 \pm 0.1$ ($\sim 18\%$) |
| Toluene* | $9.0 \pm 0.0$ | - | $8.0 \pm 0.0$ ($\sim 11\%$) | $7.9 \pm 0.1$ ($\sim 12\%$) |
| Phenol** | $9.9 \pm 0.1$ | - | $10.0 \pm 0.1$ ($\sim$ -1%) | $8.1 \pm 0.1$ ($\sim 18\%$) |
| Benzaldehyde** | $10.7 \pm 0.1$ | - | $10.5 \pm 0.1$ ($\sim 2\%$) | $9.0 \pm 0.0$ ($\sim 16\%$) |
| p-Toluidine** | $10.0 \pm 0.0$ | - | $8.3 \pm 0.1$ ($\sim 16\%$) | $7.8 \pm 0.1$ ($\sim 21\%$) |
| 1-Propanol | $9.0 \pm 0.0$ | $6.0 \pm 0.1$ ($\sim 33\%$) | $6.5 \pm 0.1$ ($\sim 27\%$) | $6.5 \pm 0.1$ ($\sim 28\%$) |
| 1-Butanol | $9.0 \pm 0.0$ | $6.4 \pm 0.1$ ($\sim 29\%$) | $6.6 \pm 0.1$ ($\sim 27\%$) | $6.5 \pm 0.1$ ($\sim 27\%$) |
| 2-Butanol | $9.0 \pm 0.0$ | $7.2 \pm 0.1$ ($\sim 20\%$) | $7.0 \pm 0.1$ ($\sim 22\%$) | $7.1 \pm 0.1$ ($\sim 22\%$) |
| 1-Hexanol | $9.0 \pm 0.0$ | $6.4 \pm 0.1$ ($\sim 29\%$) | $6.5 \pm 0.1$ ($\sim 28\%$) | $6.5 \pm 0.1$ ($\sim 28\%$) |
| Ethanol* | $9.1 \pm 0.0$ | - | $7.2 \pm 0.1$ ($\sim 21\%$) | $7.1 \pm 0.1$ ($\sim 21\%$) |
| 2-Propanol* | $9.0 \pm 0.0$ | - | $7.1 \pm 0.1$ ($\sim 21\%$) | $7.0 \pm 0.0$ ($\sim 22\%$) |
| 2-Hexanol* | $9.0 \pm 0.0$ | - | $7.0 \pm 0.1$ ($\sim 23\%$) | $7.0 \pm 0.1$ ($\sim 22\%$) |
| 1-Heptanol* | $9.0 \pm 0.0$ | - | $6.6 \pm 0.1$ ($\sim 27\%$) | $6.6 \pm 0.1$ ($\sim 27\%$) |

Table 3: Average number of SCF iterations obtained by the (3,8)-models trained with different datasets ($r_{\text{cut}} = 4.0$). The values reported within parentheses indicate the percentage of reduction with respect to the default Gaussian guess. The test molecules with a superscript $*$ are not included in the training process at all, and those with $**$ are involved in the training of Unified Model-A, with only 10 frames each included.

considered as qualitatively correct results in most of the cases. Typically, the average errors of the properties derived from the predicted density matrix for all the involved molecules are 1 to 3 orders of magnitude smaller than those from the Gaussian default guesses. This trend holds consistently for both the aldehyde and aromatic families. Despite the existence of some outliers for the alcohols, especially 2-butanol, which also turned out to be the one having the largest test set RMSE within the training molecules (unified models), they are only rare occurrences, as indicated by the error distribution shown in the violin plots. We expect that this can be resolved by adjusting the training set to include the structures corresponding to the outliers. This result also suggests that one may need to give the alcohol family more weight in the training. It is therefore likely that a better selection of training points, obtained for example by active learning approaches (see *e.g.* Ref. [41] and the references therein), will give more stable errors. In Section 4.3, we provide a potential way to determine whether a given prediction should be disregarded or whether the corresponding structure should be included in the training set to improve model performance.

## 4.3   Commutator and errors

As a last application, we use the predicted density matrices to compute the corresponding KS matrix $F = F(D)$, and check how well the commutator condition $FD = DF$ is fulfilled. Indeed, when convergence is achieved, $FD = DF$ must hold exactly. Therefore, the norm of $FD - DF$ is a residue and can serve as a physical parameter to evaluate the accuracy of the prediction. In Figure 6, we
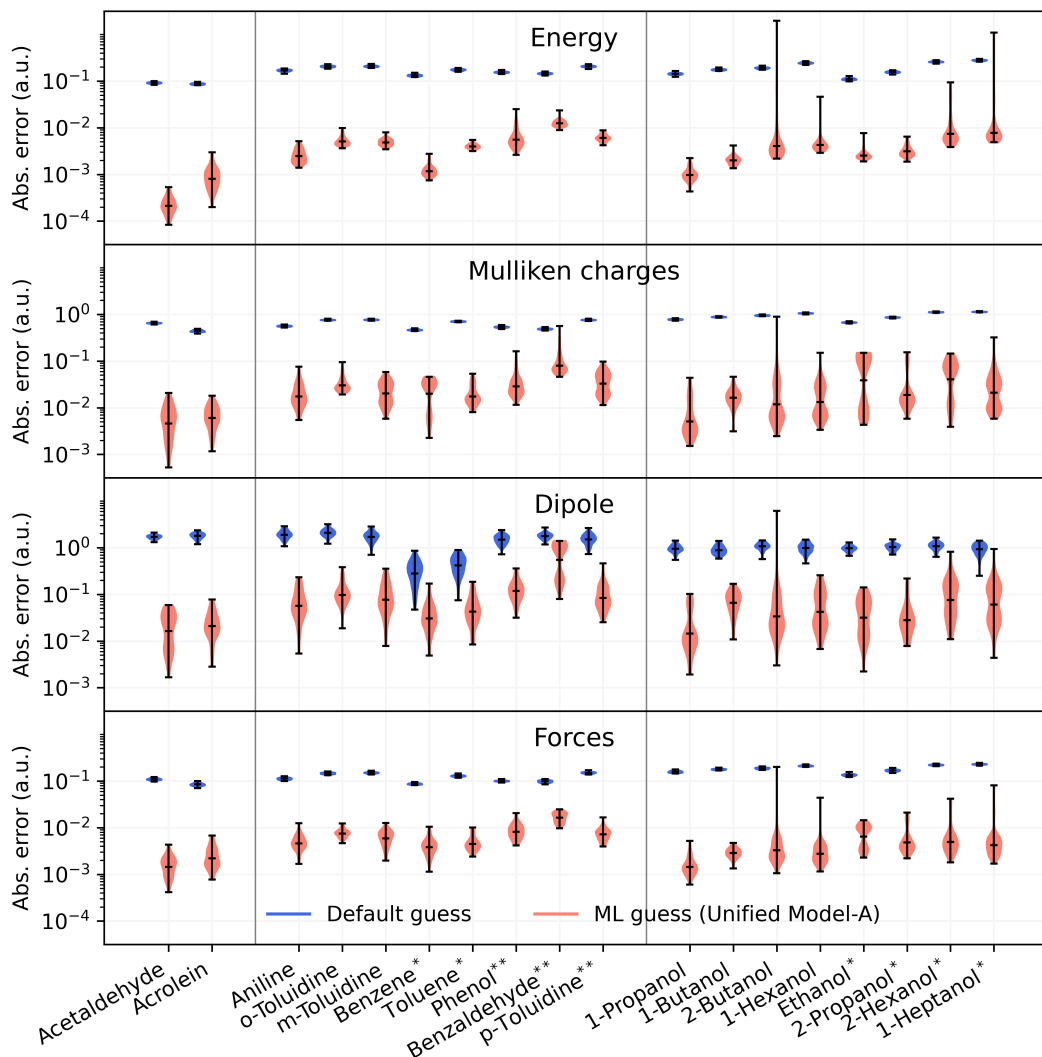
Figure 5: Plot of the error (in logarithmic scale) for energy, Mulliken charges, dipole moment and forces after a single SCF cycle. The blue line represents the default guess provided by Gaussian, while the pink line corresponds to the density matrix predicted using Unified Model-A. The test molecules with a superscript $*$ are not included in the training process at all, and those with $**$ are involved in the training of Unified Model-A, with only 10 frames each included.

present the relationship between the commutator violation error, measured in the Frobenius norm, and the relative error in the predicted energy. It turns out that there is an empirical algebraic relation observed between the two errors. Similar plots for other properties are provided in the SI (Figure S4), which also demonstrate positive correlations while the trend is less clear compared to that for the energies. Thus, we can use the commutator error to determine whether to disregard a prediction, without accessing the real physical properties of interest. From another perspective, we can also use the commutator error as an indicator of which frames to include in the training process in an active learning framework. As shown in the previous section, it is indeed the design of the training set that limits the accuracy of the proposed method, and this is therefore one of our immediate future works.
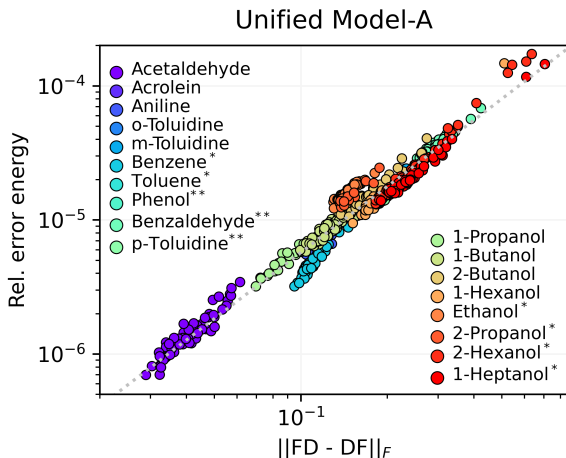
14

Figure 6: Plot of the Frobenius norm of the commutator between $F$ and $D$ versus the relative error in energy obtained after a single SCF cycle, using the density matrix predicted by the Unified Model-A as a guess. The test molecules with a superscript $*$ are not included in the training process at all, and those with $**$ are involved in the training of Unified Model-A, with only 10 frames each included.

## 5 Conclusions

We presented a simple yet powerful regression model for learning the ground-state density matrix of arbitrary molecules in an atom-centered basis set. Our model exploits the flexibility and the favorable symmetry properties of the equivariant ACE descriptors, which represent a natural set of features to represent the density matrix. The resulting model can be improved systematically by increasing the size of the ACE basis (order and degree) or by tuning the training samples. More importantly, our model can learn the relationship between molecular geometry and density matrix using information from multiple distinct molecules. This opens the possibility of building unified models for molecules with similar local structure, in contrast to other approaches[3, 40]. As a consequence, our model is transferable to unseen molecules, provided that they have a local chemical structure similar to the ones in the training set.

A model generating fast predictions for the density matrices provides, first of all, a way to accelerate KS-DFT calculations by generating a better starting guess for self-consistent iterations. Besides the straightforward application to *ab initio* molecular dynamics or geometry optimizations, the possibility of extrapolating predictions to unseen molecules provides the opportunity to accelerate KS-DFT calculations on many different molecules without molecule-specific training. The guesses generated by our unified model allow saving about 20% of the SCF iterations compared to the default guess in most of the cases.

Secondly, the predicted density matrix can be used in a quantum chemistry code to directly compute multiple properties. We have tested how well this model predicts energy, Mulliken atomic charges, molecular dipole, and atomic forces. The density matrix predictions give significantly better estimates than the standard guess for all these properties, and especially for the energy, even for unseen molecules.

Our model still presents some limitations. First, although the reduction in SCF iterations is significant, a substantial improvement would be necessary to consistently accelerate KS-DFT calculations by at least a factor of two. Second, the energies and forces predicted by our model still do not reach chemical accuracy, which would be needed for direct applications. Nonetheless, all models show significant room for improvement, both in the model flexibility and in the choice of the training set, making our strategy promising for both applications, especially thanks to the observed systematic improvability

and rigorous methodology.

Overall, our results show that learning the density matrix from descriptors encoding the correct symmetry features represents a promising strategy towards more complete and transferable ML models. In particular, the density matrix represents the solution of the KS-DFT equations and thus gives direct access to numerous properties with one single ML model. Further, the model turns out to be transferable to unseen molecular structures, which is a central stepping stone towards this development.

## Acknowledgements

## A    Equivariance of the density matrix

In this appendix, we discuss the equivariance of the density matrix and the rationality of decomposing it in the way that is illustrated in Figure 1.

Given a configuration $\mathbf{R}$, and a set of atomic orbitals $\{\phi_{I\alpha}\}_{I,\alpha}$ with which (1) is discretized, where $I$ ranges from 1 to $N_{\text{at}}$ and $\alpha \in \mathcal{I}_{Z_I}$, an index set depending only on $Z_I$. It is straightforward to see that $N_g = \sum_{I=1}^{N_{\text{at}}} \#\mathcal{I}_{Z_I}$ with $\#A$ denoting the cardinality of the set $A$. After the discretization, the solution of (1) can be used to approximate the KS orbitals, as

$$\varphi_k(\boldsymbol{r}; \mathbf{R}) = \sum_{(I,\alpha)} c_{k,(I,\alpha)}[\mathbf{R}]\phi_{I\alpha}(\boldsymbol{r}), \quad k = 1, 2, \ldots, N.$$

where the collection of all $\{c_{k,(I,\alpha)}[\mathbf{R}]\}$ form exactly the eigenvectors $C_{\mathbf{R}}$ in (1). Define the density matrix operator by

$$D(\boldsymbol{r}, \boldsymbol{r}'; \mathbf{R}) = \sum_{k=1}^{N} \varphi_k(\boldsymbol{r}; \mathbf{R})^* \varphi_k(\boldsymbol{r}'; \mathbf{R}),$$

and denote its diagonal by

$$D(\boldsymbol{r}; \mathbf{R}) = D(\boldsymbol{r}, \boldsymbol{r}; \mathbf{R}).$$

Then, the elements of the (discretized) density matrix are given by

$$
\begin{aligned}
\left(D_{\mathbf{R}}\right)_{IJ\alpha\beta} &= \int_{\mathbb{R}^3} \phi_{I\alpha}(\boldsymbol{r})^* D(\boldsymbol{r}; \mathbf{R})\phi_{J\beta}(\boldsymbol{r})d\boldsymbol{r} \\
&= \int_{\mathbb{R}^3} \phi_\alpha(\boldsymbol{r} - \boldsymbol{r}_I)^* D(\boldsymbol{r}; \mathbf{R})\phi_\beta(\boldsymbol{r} - \boldsymbol{r}_J)d\boldsymbol{r} \\
&= \int_{\mathbb{R}^3} \phi_\alpha(\boldsymbol{r} - \boldsymbol{r}_I)^* \tilde{D}(\boldsymbol{r}; \mathbf{R}_{IJ})\phi_\beta(\boldsymbol{r} - \boldsymbol{r}_J)d\boldsymbol{r}, \quad 1 \le I, J \le N_{\text{at}}, \ \alpha \in \mathcal{I}_{Z_I}, \ \beta \in \mathcal{I}_{Z_J}.
\end{aligned}
$$

In the last line, the whole configuration $\mathbf{R}$ is shifted to be centred at a proper position, which will be elaborated in more details in Appendix B and $\tilde{D}$ is simply the translated density matrix operator. Consequently, the $(I, J)$-th and the $(I', J')$-th blocks of the density matrix $D_{\mathbf{R}}$ share the same function form as long as $Z_I = Z_{I'}$ and $Z_J = Z_{J'}$, hence gives the unified form (9). In addition, we follow an

extremely similar discussion as in Ref. [12] to obtain the equivariance of the density matrix. More specifically, if

$$\mathcal{I}_{Z_I} = \{(n,l,m)\}_{\substack{l \in \{0,1,\ldots,l_I\}, \\ n \in \{0,1,\ldots,n_l\}, \\ m \in \{-l,-l+1,\ldots,l-1,l\}}}$$

where the index $l$ stands for the angular moment, indicating the type of atomic orbital being used (s,p,d orbitals etc.), $n_l$ denotes the number of orbitals of the type $l$, and $m$ is the standard angular index corresponding to $l$, then there holds

$$\left(D_{Q\mathbf{R}}\right)_{IJ} = \boldsymbol{\mathcal{D}}_I(Q)\left(D_\mathbf{R}\right)_{IJ}\boldsymbol{\mathcal{D}}_J(Q)^*, \quad \forall Q \in \mathrm{O}(3). \tag{9}$$

Here $\boldsymbol{\mathcal{D}}_\bullet(Q) = \mathrm{Diag}(\{\mathcal{D}^l(Q)\}_{l \in \mathcal{L}_\bullet})$ with the ordered tuple

$$\mathcal{L}_\bullet = [\![\underbrace{0,\ldots,0}_{n_0},\underbrace{1,\ldots,1}_{n_1},\ldots,\underbrace{l_\bullet,\ldots,l_\bullet}_{n_{l_\bullet}}]\!],$$

and $\bullet$ standing for $I$ or $J$. Moving one step forward, we have (cf. Figure 1(b) and (c))

$$\left(D_{Q\mathbf{R}}\right)_{IJ}^{ll'} = \mathcal{D}^l(Q)\left(D_\mathbf{R}\right)_{IJ}^{ll'}\mathcal{D}^{l'}(Q)^*, \quad \forall (l,l') \in \mathcal{L}_I \times \mathcal{L}_J, \ Q \in \mathrm{O}(3), \tag{10}$$

which is the smallest unit of describing the isometric symmetry of the density matrix.

# B   The construction of the ACE basis

Despite the existence of some other possible ways of construction, all the equivariant ACE bases can be obtained through the following procedure, following Ref. [12]:

1-particle basis $\rightarrow$ density projection $\rightarrow$ $\nu$-body correlation $\rightarrow$ symmetrization.

The main differences that distinguish the ACE bases are the construction of the 1-particle basis and the symmetrization, while the latter remains the same throughout this paper (cf. (9) or (10)). The focus is thus on the definition of the 1-particle basis. The 1-particle basis is, as its name suggests, a function applied to a single particle $\sigma$, which we will specify in detail for both the *onsite* and *offsite* cases, respectively, for completeness. Suppose the configuration is given by $\mathbf{R} = \{\sigma_I\}_{I=1,2,\ldots,N_{\mathrm{at}}}$, with $\sigma_I = (Z_I, \boldsymbol{r}_I)$.

***Onsite*** **basis:** the onsite environment $\mathbf{R}_I$ centering at the $I$-th atom is defined as

$$\mathbf{R}_I = \{(Z_K, Z_I, \boldsymbol{r}_{KI})\} =: \{\sigma_K\}_{K=1,2,\ldots,N_{\mathrm{at}}, K \neq I},$$

where $\boldsymbol{r}_{KI} = \boldsymbol{r}_K - \boldsymbol{r}_I$. The second variable in $\sigma_K$ indicates the center of the system. Given a particle $\sigma_K = (Z_K, Z_I, \boldsymbol{r}) \in \mathbf{R}_I$, we define the 1-particle basis for the *onsite* $Z_I$ model as

$$\phi_v^{\mathrm{on}}(\sigma_K) := \phi_{Znlm}^{\mathrm{on}}(\sigma_K) := \delta_{ZZ_K}P_{nl}(r)Y_{lm}(\hat{\boldsymbol{r}})f_{\mathrm{cut}}(r)$$

where $\delta_{ZZ_K}$ is the kronecker delta, $\boldsymbol{r} = r \cdot \hat{\boldsymbol{r}}$ with $r = |\boldsymbol{r}|$, and we have identified the composite index $v \equiv (Znlm)$ with $Z$ standing for some atomic number. The radial cutoff function $f_{\mathrm{cut}}(r)$ ensures that only the nearby atoms are taken into account (cf Figure 2), which may take different forms. In this work, we choose

$$f_{\text{cut}}(r; r_{\text{cut},Z_I}) = \begin{cases} (r^2/r_{\text{cut},Z_I}^2 - 1)^2, & r \leq r_{\text{cut},Z_I}, \\ 0, & r > r_{\text{cut},Z_I}. \end{cases} \tag{11}$$

Here, the subscript $Z_I$ indicates that the cutoff radius can be made element-dependent, and will sometimes be neglected for simplicity.

***Offsite*** **basis:** For the *offsite* interactions, we define the offsite local environment as

$$\mathbf{R}_{IJ} = \{\sigma_{IJ}\} \cup \{\sigma_K\}_{K=1,2,\ldots,N_{\text{at}}, K \neq I,J},$$

where $\sigma_{IJ} = \{(Z_I, Z_J), \boldsymbol{r}_{IJ}\}$, $\sigma_K = \{Z_K, (Z_I, Z_J), \boldsymbol{r}_{IJ,K}\}$ and

$$\boldsymbol{r}_{IJ,K} = \boldsymbol{r}_K - \boldsymbol{r}_{IJ,\theta}, \text{ with } \boldsymbol{r}_{IJ,\theta} = \boldsymbol{r}_J + \theta(\boldsymbol{r}_I - \boldsymbol{r}_J), \ \theta \in [0,1].$$

The 1-particle basis is then defined as

$$\phi_{nlm}^{\text{b}}(\sigma_{IJ}) = P_{nl}(r_{IJ}) Y_{lm}(\hat{\boldsymbol{r}}_{IJ}) f_{\text{cut}}^{\text{b}}(r_{IJ}),$$
$$\phi_{Znlm}^{\text{e}}(\sigma_K) = \delta_{ZZ_K} P_{nl}(r_{IJ,K}) Y_{lm}(\hat{\boldsymbol{r}}_{IJ,K}) f_{\text{cut}}^{\text{e}}(\boldsymbol{r}_{IJ,K}; \boldsymbol{r}_{IJ}).$$

Here the cutoff function for the bond can be given analogous to those for the *onsite* basis, as

$$f_{\text{cut}}^{\text{b}}(r; r_{\text{bond},Z_I Z_J}) = f_{\text{cut}}(r; r_{\text{bond},Z_I Z_J}),$$

with $f_{\text{cut}}$ being defined in (11). On the other hand, the cutoff function $f_{\text{cut}}^{\text{e}}$ for the environmental atom is defined as

$$f_{\text{cut}}^{\text{e}}(\boldsymbol{r}; r_{\text{cut},Z_I}, r_{\text{cut},Z_J}) = f_{\text{cut}}\big(|\boldsymbol{r} + (1-\theta)(\boldsymbol{r}_J - \boldsymbol{r}_I)|; r_{\text{cut},Z_I}\big) + f_{\text{cut}}\big(|\boldsymbol{r} + \theta(\boldsymbol{r}_I - \boldsymbol{r}_J)|; r_{\text{cut},Z_J}\big),$$

where $f_{\text{cut}}$ is defined in (11), $\boldsymbol{r}$ denotes the $\boldsymbol{r}_{IJ,K}$ element in $\sigma_K$ and $r_{\text{cut},Z_I/Z_J}$, the cutoff radii of the two spheres around both the $I$-th and the $J$-th atom. Throughout this work, we choose $\theta = 0.5$.

Given the one-particle basis, we can form the density projection and the $\nu$-correlations for the onsite case as

$$A_v(\mathbf{R}_I) := \sum_{\sigma \in \mathbf{R}_I} \phi_v(\sigma),$$
$$\boldsymbol{A_v}(\mathbf{R}_I) := \prod_{t=1}^{\nu} A_{v_t}(\mathbf{R}_I) \qquad \text{for } \boldsymbol{v} = (v_1, \ldots, v_\nu), \ \nu = 1, 2, \ldots,$$

and for the offsite,

$$A_v(\mathbf{R}_{IJ}) := \sum_{K \neq I,J} \phi_v^{\text{e}}(\sigma_K),$$
$$\boldsymbol{A_v}(\mathbf{R}_{IJ}) := \phi_{v^0}^{\text{b}}(\sigma_{IJ}) \cdot \prod_{t=1}^{\nu} A_{v^t}(\mathbf{R}_{IJ}) \qquad \text{for } \boldsymbol{v} = (v_0, \ldots, v_\nu), \ \nu = 1, 2, \ldots.$$

Finally, we perform the symmetrization over O(3), by leveraging an averaged integral

$$\mathcal{B}_{\boldsymbol{v},a}(\mathbf{R}_\bullet) = \fint_{\text{O}(3)} D(Q)(\boldsymbol{A_v}(Q\mathbf{R}_\bullet)E_a)D(Q)^* dQ,$$

18

where $\{E_a\}_a$ forms a canonical basis of the matrix space where the density matrices $D$, or a sub-block thereof, lies in.

By design, the $\mathcal{B}_{\boldsymbol{v},a}$ bases have exactly the same equivariance as (the subblocks) of the density matrix $D$. To simplify the notation, we absorbed the index $a$ into $\boldsymbol{v}$ in the main text.

## C    Properties of the retraction

In this section, we discuss the retraction operator $\mathcal{P}$ and some of its important properties. Denoting

$$S_{N_g} = \{D \in \mathbb{R}^{N_g \times N_g} : D^T = D, \lambda_{D,N} > \lambda_{D,N+1}\},$$

where $\{\lambda_{D,N}\}_{N=1}^{N_g}$ represents the eigenvalues of $D$, sorted descendingly, then we have the following proposition.

**Proposition 1.** *Let* $\mathcal{P} : S_{N_g} \to \mathcal{G}_{N_g}^N$ *be the retraction operator defined in* (7), *then:*

*(1.a) for* $D_{\mathbf{R}}$ *defined in* (9), *there holds*

$$\mathcal{P}(D_{Q\mathbf{R}}) = \mathcal{D}(Q)\mathcal{P}(D_{\mathbf{R}})\mathcal{D}(Q)^*, \ \forall Q \in \mathrm{O}(3);$$

*(1.b) for any* $D \in S_{N_g}$,

$$\mathcal{P}(D) = \underset{\tilde{D} \in \mathcal{G}_{N_g}^N}{\arg\min} \|\tilde{D} - D\|_F.$$

*Proof.* We first justify that $\mathcal{P}$ is well-defined. Let $D \in S_{N_g}$, then its eigenvalue decomposition can be written as $D = U\Sigma U^T$, where the unitary matrix $U = [u_1, u_2, \ldots, u_{N_g}]$ consists of $N_g$ orthonormal eigenvectors of $D$, $\Sigma$ is a diagonal matrix containing the eigenvalues of $D$, sorted decreasingly. By a straight calculation, we have

$$U E_{N_g}^N U^T = \sum_{i=1}^N u_i u_i^T.$$

Although $U$ may be non-unique, the result in the RHS of the above equality will not be influenced by the order and signs of the first $N$ orthonormal eigenvectors of $D$, since $\lambda_{D,N} > \lambda_{D,N+1}$. This shows the well-definedness of $\mathcal{P}$, *i.e.*, the image of $D$ via $\mathcal{P}$ is unique regardless of how the eigenvalue decomposition is performed.

Now we can move on to consider (1.a). Assume $D_{\mathbf{R}} = U_{D_{\mathbf{R}}} \Sigma_{D_{\mathbf{R}}} U_{D_{\mathbf{R}}}^T$, then $\mathcal{P}(D_{\mathbf{R}}) = U_{D_{\mathbf{R}}} E_{N_g}^N U_{D_{\mathbf{R}}}^T$. By (9), we have that for all $Q \in \mathrm{O}(3)$

$$
\begin{aligned}
D_{Q\mathbf{R}} &= \mathcal{D}(Q) D_{\mathbf{R}} \mathcal{D}(Q)^*, \\
&= \mathcal{D}(Q) U_{D_{\mathbf{R}}} \Sigma_{D_{\mathbf{R}}} U_{D_{\mathbf{R}}}^T \mathcal{D}(Q)^*,
\end{aligned}
$$

which means that $D_{Q\mathbf{R}}$ has an eigenvalue decomposition as above. As a result,

$$
\begin{aligned}
\mathcal{P}(D_{Q\mathbf{R}}) &= \mathcal{D}(Q) U_{D_{\mathbf{R}}} E_{N_g}^N U_{D_{\mathbf{R}}}^T \mathcal{D}(Q)^*, \\
&= \mathcal{D}(Q) \mathcal{P}(D_{\mathbf{R}}) \mathcal{D}(Q)^*, \quad \forall Q \in \mathrm{O}(3).
\end{aligned}
$$

This proves (1.a).

As of (1.b), we again suppose $D = U\Sigma U^T \in S_{N_g}$, with $U$ defined as above. We claim that

$$\mathcal{P}(D) = U E_{N_g}^N U^T = \underset{G \in \mathcal{G}_{N_g}^N}{\arg\min} \|G - D\|_F = \underset{\{P E_{N_g}^N P^T : \, P \in \mathcal{O}(N_g)\}}{\arg\min} \|P E_{N_g}^N P^T - D\|_F.$$

Note that the last equality above is based on

$$\mathcal{G}_{N_g}^N = \{PE_{N_g}^N P^T : \ P \in \mathcal{O}(N_g)\}.$$

We now estimate

$$\|PE_{N_g}^N P^T - D\|_F^2 = \|PE_{N_g}^N P^T - U\Sigma U^T\|_F^2 = \mathrm{tr}(PE_{N_g}^N P^T + U\Sigma^2 U^T - 2PE_{N_g}^N P^T U\Sigma U^T). \quad (12)$$

To minimize (12), we just need to maximize $\mathrm{tr}(PE_{N_g}^N P^T U\Sigma U^T)$, since the trace of the first two terms in the right hand side is a constant, thus

$$
\begin{aligned}
\max_{P \in \mathcal{O}(N_g)} \mathrm{tr}(PE_{N_g}^N P^T U\Sigma U^T) &= \max_{P \in \mathcal{O}(N_g)} \mathrm{tr}(U^T PE_{N_g}^N P^T U\Sigma) \\
&= \max_{P \in \mathcal{O}(N_g)} \mathrm{tr}(PE_{N_g}^N P^T \Sigma) \\
&= \max_{G \in \mathcal{G}_{N_g}^N} \mathrm{tr}(G\Sigma) = \max_{G \in \mathcal{G}_{N_g}^N} \sum_{i=1}^{N_g} \sigma_i G_{ii}.
\end{aligned}
$$

For any $G_0 \in \mathcal{G}_{N_g}^N$, there exists $P_0 \in \mathcal{O}(N_g)$ such that $G_0 = P_0 E_{N_g}^N P_0^T$. Consequently,

$$0 \le G_{0,ii} = \sum_{t=1}^{N} P_{0,it}^2 \le \sum_{t=1}^{N_g} P_{0,it}^2 = 1.$$

In addition, we have $\mathrm{tr}(G_0) = N$. Hence $\frac{1}{N}\sum_{i=1}^{N_g} \sigma_i G_{0,ii} = \sum_{i=1}^{N_g} \sigma_i \frac{G_{0,ii}}{N}$ becomes a convex combinition of $\{\sigma_i\}_{i=1}^{N_g}$, which achieves its maximum at $\sum_{i=1}^{N} \sigma_i$ when $G_0$ is chosen to be $E_{N_g}^N$. This completes the proof. $\qquad\square$

## D    Comparison of fitting the density matrix and the KS matrix

The KS matrix $F_{\mathbf{R}}$ has the same structure and symmetry as the density matrix $D_{\mathbf{R}}$ and the learning of it has been more broadly studied compared to that of the density matrix. In this appendix, we compare the fitting of the two objects using the same method. To make our comparison meaningful and fair, we first fix an ACE basis $\mathbf{B}_{\nu,d}$, where $\nu$ and $d$ stand for the correlation order and polynomial degree that define the basis (in addition, the cutoffs are also fixed but is not explicitly written here for the sake of simplicity). That said, the only thing different for the two targeting models (for the KS matrix and the density matrix) is their coefficients. We then pick the same data points, which have the form $\{(\mathbf{R}^{(k)}, F_{\mathbf{R}^{(k)}}, D_{\mathbf{R}^{(k)}})\}_{k=1}^{N_{\mathrm{data}}}$. By solving the least squares problems (6) with the Kohn-Sham matrices and density matrices, respectively, we obtain $\mathbf{c}_{\nu,d,H}$ and $\mathbf{c}_{\nu,d,D}$ for the two objects. Then we have two routines to get the predicted density matrix of a given configuration $\mathbf{R}$, which lies in the desired manifold (2).

First, with $\mathbf{B}_{\nu,d}$ and $\mathbf{c}_{\nu,d,D}$, we obtain directly a feasible approximation of the density matrix

$$\tilde{D}_{\mathbf{R}} = \mathcal{P}(\mathbf{c}_{\nu,d,D} \cdot \mathbf{B}_{\nu,d}(\mathbf{R})),$$

where $\mathcal{P}$ is the retraction operator defined in (7). Alternatively, we can construct

$$\tilde{F}_{\mathbf{R}} = \mathbf{c}_{\nu,d,F} \cdot \mathbf{B}_{\nu,d}(\mathbf{R}),$$

| Molecule | Density Matrix | | KS Matrix | |
|---|---|---|---|---|
| | Specific Model | Unified Model | Specific Model | Unified Model |
| Acetaldehyde | $4.416 \cdot 10^{-5}$ | $3.278 \cdot 10^{-4}$ | $2.369 \cdot 10^{-5}$ | $2.137 \cdot 10^{-4}$ |
| Acrolein | $2.514 \cdot 10^{-4}$ | $5.028 \cdot 10^{-4}$ | $1.688 \cdot 10^{-4}$ | $4.052 \cdot 10^{-4}$ |
| Aniline | $4.300 \cdot 10^{-4}$ | $4.868 \cdot 10^{-4}$ | $5.400 \cdot 10^{-4}$ | $1.316 \cdot 10^{-3}$ |
| o-Toluidine | $5.430 \cdot 10^{-4}$ | $5.962 \cdot 10^{-4}$ | $4.760 \cdot 10^{-4}$ | $3.806 \cdot 10^{-3}$ |
| m-Toluidine | $5.384 \cdot 10^{-4}$ | $5.824 \cdot 10^{-4}$ | $5.401 \cdot 10^{-4}$ | $4.115 \cdot 10^{-3}$ |
| Benzene* | - | $4.058 \cdot 10^{-4}$ | - | $4.333 \cdot 10^{-4}$ |
| Toluene* | - | $6.369 \cdot 10^{-4}$ | - | $4.893 \cdot 10^{-4}$ |
| Phenol** | - | $4.809 \cdot 10^{-3}$ | - | $3.204 \cdot 10^{-2}$ |
| Benzaldehyde** | - | $4.129 \cdot 10^{-3}$ | - | $3.101 \cdot 10^{-2}$ |
| p-Toluidine** | - | $2.840 \cdot 10^{-3}$ | - | $1.922 \cdot 10^{-2}$ |
| 1-Propanol | $3.049 \cdot 10^{-4}$ | $4.427 \cdot 10^{-4}$ | $4.044 \cdot 10^{-4}$ | $3.468 \cdot 10^{-4}$ |
| 1-Butanol | $4.510 \cdot 10^{-4}$ | $4.921 \cdot 10^{-4}$ | $2.801 \cdot 10^{-4}$ | $3.579 \cdot 10^{-4}$ |
| 2-Butanol | $9.173 \cdot 10^{-4}$ | $1.494 \cdot 10^{-3}$ | $5.748 \cdot 10^{-4}$ | $1.709 \cdot 10^{-3}$ |
| 1-Hexanol | $1.031 \cdot 10^{-3}$ | $5.324 \cdot 10^{-4}$ | $3.264 \cdot 10^{-4}$ | $4.193 \cdot 10^{-4}$ |
| Ethanol* | - | $9.644 \cdot 10^{-4}$ | - | $1.133 \cdot 10^{-3}$ |
| 2-Propanol* | - | $7.701 \cdot 10^{-4}$ | - | $2.948 \cdot 10^{-3}$ |
| 2-Hexanol* | - | $7.384 \cdot 10^{-4}$ | - | $1.219 \cdot 10^{-3}$ |
| 1-Heptanol* | - | $8.896 \cdot 10^{-4}$ | - | $2.234 \cdot 10^{-3}$ |

Table 4: The average test set RMSEs on the density matrices obtained by the (3,8) specific models and unified model trained with the density matrices and the KS matrices ($r_{\mathrm{cut}} = 4.0$).

following by solving

$$\tilde{F}_{\mathbf{R}}\tilde{C}_{\mathbf{R}} = \tilde{C}_{\mathbf{R}}E_{\mathbf{R}}$$

for its N eigenvectors $\tilde{C}_{\mathbf{R}}$, and finally obtain

$$\bar{D}_{\mathbf{R}} = \tilde{C}_{\mathbf{R}}\tilde{C}_{\mathbf{R}}^{T}.$$

In table 4, we compare the element-wise test set RMSE in the predicted density matrices obtained by the two approaches above, including both the cases of specific models and the unified model. To be fair, we use exactly the same data points as those mentioned in Section 3 for training. It can be observed that fitting the KS matrix gives comparable, or even smaller test set RMSEs on the predicted density matrices with respect to the specific models, whereas it performs poorly on the unified model. This implies that the approach of fitting the KS matrix only by minimizing the element-wise error is either less transferable or requires more careful weighting for different systems. In any case, it is less robust than the same approach for the density matrix.

In addition, to see the performance of the specific models more clearly, we compared the properties derived from the specific models tailored to both the density matrices and the KS matrices. The errors on the properties are illustrated in Figure 7, from which we can see that although the approach of predicting the KS matrices provides smaller element-wise errors, it cannot guarantee that the derived properties are equally effective. Overall, the proposed approach of fitting the matrix element is more suitable for the density matrix rather than for the KS matrix, at least for the molecular systems that were mentioned in our experiments.
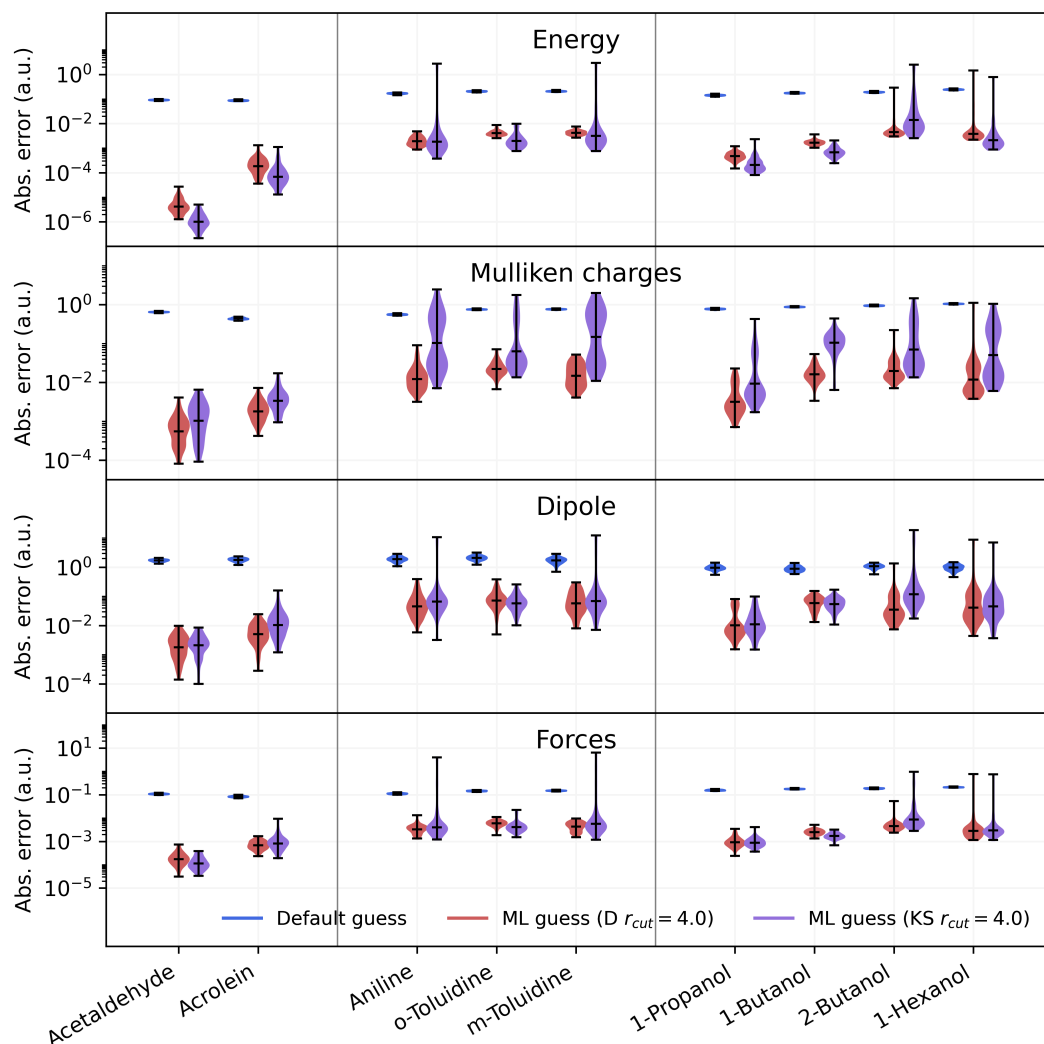
Figure 7: Plot of the error (in logarithmic scale) for energy, Mulliken charges, dipole moment and forces after a single SCF cycle. The blue line represents the default guess provided by Gaussian, and red and violet lines refer to the ML guess obtained with models trained on the density matrices and the KS matrices, respectively ($r_{\text{cut}} = 4.0$).

# References

[1] Nikita Fedik et al. "Extending machine learning beyond interatomic potentials for predicting molecular properties". In: *Nature Reviews Chemistry* 6.9 (Aug. 2022), pp. 653–672. ISSN: 2397-3358. DOI: 10.1038/s41570-022-00416-3.

[2] Zhuoran Qiao et al. "OrbNet: Deep learning for quantum chemistry using symmetry-adapted atomic-orbital features". In: *The Journal of Chemical Physics* 153.12 (Sept. 2020). ISSN: 1089-7690. DOI: 10.1063/5.0021955.

[3] Xuecheng Shao et al. "Machine learning electronic structure methods based on the one-electron reduced density matrix". In: *Nature Communications* 14.1 (Oct. 2023). ISSN: 2041-1723. DOI: 10.1038/s41467-023-41953-9.

[4]   Yixiao Chen et al. "Ground State Energy Functional with Hartree–Fock Efficiency and Chemical Accuracy". In: *The Journal of Physical Chemistry A* 124.35 (Aug. 2020), pp. 7155–7165. ISSN: 1520-5215. DOI: 10.1021/acs.jpca.0c03886.

[5]   Sebastian Dick and Marivi Fernandez-Serra. "Machine learning accurate exchange and correlation functionals of the electronic density". In: *Nature Communications* 11.1 (July 2020). ISSN: 2041-1723. DOI: 10.1038/s41467-020-17265-7.

[6]   Anders S. Christensen et al. "OrbNet Denali: A machine learning potential for biological and organic chemistry with semi-empirical cost and DFT accuracy". In: *The Journal of Chemical Physics* 155.20 (Nov. 2021). ISSN: 1089-7690. DOI: 10.1063/5.0061990.

[7]   Matthew Welborn, Lixue Cheng, and Thomas F. Miller. "Transferability in Machine Learning for Electronic Structure via the Molecular Orbital Basis". In: *Journal of Chemical Theory and Computation* 14.9 (July 2018), pp. 4772–4779. ISSN: 1549-9626. DOI: 10.1021/acs.jctc.8b00636.

[8]   Oliver T. Unke et al. "Machine Learning Force Fields". In: *Chemical Reviews* 121 (16 Aug. 2021), pp. 10142–10186. ISSN: 15206890. DOI: 10.1021/ACS.CHEMREV.0C01111.

[9]   Max Pinheiro et al. "Choosing the right molecular machine learning potential". In: *Chemical Science* 12 (43 Nov. 2021), pp. 14396–14413. ISSN: 2041-6539. DOI: 10.1039/D1SC03564A.

[10]  Oliver T. Unke et al. "SpookyNet: Learning force fields with electronic degrees of freedom and nonlocal effects". In: *Nature Communications 2021 12:1* 12 (1 Dec. 2021), pp. 1–14. ISSN: 2041-1723. DOI: 10.1038/s41467-021-27504-0.

[11]  K. T. Schütt et al. "Unifying machine learning and quantum chemistry with a deep neural network for molecular wavefunctions". In: *Nature Communications* 10.1 (Nov. 2019). ISSN: 2041-1723. DOI: 10.1038/s41467-019-12875-2.

[12]  Liwei Zhang et al. "Equivariant analytical mapping of first principles Hamiltonians to accurate and transferable materials models". In: *npj Computational Materials* 8.1 (July 2022). ISSN: 2057-3960. DOI: 10.1038/s41524-022-00843-2.

[13]  Jigyasa Nigam, Michael J. Willatt, and Michele Ceriotti. "Equivariant representations for molecular Hamiltonians and N-center atomic-scale properties". In: *The Journal of Chemical Physics* 156.1 (Jan. 2022). ISSN: 1089-7690. DOI: 10.1063/5.0072784.

[14]  He Li et al. "Deep-learning density functional theory Hamiltonian for efficient ab initio electronic-structure calculation". In: *Nature Computational Science* 2.6 (June 2022), pp. 367–377. ISSN: 2662-8457. DOI: 10.1038/s43588-022-00265-6.

[15]  Edoardo Cignoni et al. "Electronic Excited States from Physically Constrained Machine Learning". In: *ACS Central Science* 10.3 (Feb. 2024), pp. 637–648. ISSN: 2374-7951. DOI: 10.1021/acscentsci.3c01480.

[16]  Mohammad Shakiba and Alexey V. Akimov. "Machine-Learned Kohn–Sham Hamiltonian Mapping for Nonadiabatic Molecular Dynamics". In: *Journal of Chemical Theory and Computation* 20.8 (Apr. 2024), pp. 2992–3007. ISSN: 1549-9626. DOI: 10.1021/acs.jctc.4c00008.

[17]  Jan Hermann, Zeno Schätzle, and Frank Noé. "Deep-neural-network solution of the electronic Schrödinger equation". In: *Nature Chemistry* 12.10 (Sept. 2020), pp. 891–897. ISSN: 1755-4349. DOI: 10.1038/s41557-020-0544-y.

[18]  Xiang Li et al. "Fermionic neural network with effective core potential". In: *Physical Review Research* 4.1 (Jan. 2022). ISSN: 2643-1564. DOI: 10.1103/physrevresearch.4.013021.

[19] Stefan Chmiela et al. "sGDML: Constructing accurate and data efficient molecular force fields using machine learning". In: *Computer Physics Communications* 240 (July 2019), pp. 38–45. ISSN: 0010-4655. DOI: 10.1016/J.CPC.2019.02.007.

[20] Kirill Zinovjev et al. "emle-engine: A Flexible Electrostatic Machine Learning Embedding Package for Multiscale Molecular Dynamics Simulations". In: *Journal of Chemical Theory and Computation* 20 (11 June 2024), pp. 4514–4522. ISSN: 15499626. DOI: 10.1021/ACS.JCTC.4C00248.

[21] Raimondas Galvelis et al. "NNP/MM: Accelerating Molecular Dynamics Simulations with Machine Learning Potentials and Molecular Mechanics". In: *Journal of Chemical Information and Modeling* 63 (18 Sept. 2023), pp. 5701–5708. ISSN: 1549960X. DOI: 10.1021/ACS.JCIM.3C00773.

[22] Adil Kabylda et al. *Molecular Simulations with a Pretrained Neural Network and Universal Pairwise Force Fields*. Oct. 2024. DOI: 10.26434/chemrxiv-2024-bdfr0.

[23] Patrizia Mazzeo et al. "Electrostatic embedding machine learning for ground and excited state molecular dynamics of solvated molecules". In: *Digital Discovery* 3 (12 Dec. 2024), pp. 2560–2571. ISSN: 2635-098X. DOI: 10.1039/D4DD00295D.

[24] Felix Brockherde et al. "Bypassing the Kohn-Sham equations with machine learning". In: *Nature Communications* 8.1 (Oct. 2017). ISSN: 2041-1723. DOI: 10.1038/s41467-017-00839-3.

[25] John M. Alred et al. "Machine learning electron density in sulfur crosslinked carbon nanotubes". In: *Composites Science and Technology* 166 (Sept. 2018), pp. 3–9. ISSN: 0266-3538. DOI: 10.1016/j.compscitech.2018.03.035.

[26] Anand Chandrasekaran et al. "Solving the electronic structure problem with machine learning". In: *npj Computational Materials* 5.1 (Feb. 2019). ISSN: 2057-3960. DOI: 10.1038/s41524-019-0162-7.

[27] Sheng Gong et al. "Predicting charge density distribution of materials using a local-environment-based graph convolutional network". In: *Physical Review B* 100.18 (Nov. 2019). ISSN: 2469-9969. DOI: 10.1103/physrevb.100.184103.

[28] Alberto Fabrizio et al. "Electron density learning of non-covalent systems". In: *Chemical Science* 10.41 (2019), pp. 9424–9432. ISSN: 2041-6539. DOI: 10.1039/c9sc02696g.

[29] Bruno Cuevas-Zuviría and Luis F. Pacios. "Analytical Model of Electron Density and Its Machine Learning Inference". In: *Journal of Chemical Information and Modeling* 60.8 (Aug. 2020), pp. 3831–3842. ISSN: 1549-960X. DOI: 10.1021/acs.jcim.0c00197.

[30] Ralf Meyer, Manuel Weichselbaum, and Andreas W. Hauser. "Machine Learning Approaches toward Orbital-free Density Functional Theory: Simultaneous Training on the Kinetic Energy Density Functional and Its Functional Derivative". In: *Journal of Chemical Theory and Computation* 16.9 (Aug. 2020), pp. 5685–5694. ISSN: 1549-9626. DOI: 10.1021/acs.jctc.0c00580.

[31] J. A. Ellis et al. "Accelerating finite-temperature Kohn-Sham density functional theory with deep neural networks". In: *Physical Review B* 104.3 (July 2021). ISSN: 2469-9969. DOI: 10.1103/physrevb.104.035120.

[32] Bruno Cuevas-Zuviría and Luis F. Pacios. "Machine Learning of Analytical Electron Density in Large Molecules Through Message-Passing". In: *Journal of Chemical Information and Modeling* 61.6 (May 2021), pp. 2658–2666. ISSN: 1549-960X. DOI: 10.1021/acs.jcim.1c00227.

[33] Peter Bjørn Jørgensen and Arghya Bhowmik. "Equivariant graph neural networks for fast electron density estimation of molecules, liquids, and solids". In: *npj Computational Materials* 8.1 (Aug. 2022). ISSN: 2057-3960. DOI: 10.1038/s41524-022-00863-y.

[34] Bruno Focassio et al. "Linear Jacobi-Legendre expansion of the charge density for machine learning-accelerated electronic structure calculations". In: *npj Computational Materials* 9.1 (May 2023). ISSN: 2057-3960. DOI: 10.1038/s41524-023-01053-0.

[35] Ryong-Gyu Lee and Yong-Hoon Kim. "Convolutional network learning of self-consistent electron density via grid-projected atomic fingerprints". In: *npj Computational Materials* 10.1 (Oct. 2024). ISSN: 2057-3960. DOI: 10.1038/s41524-024-01433-0.

[36] Andrea Grisafi et al. "Transferable Machine-Learning Model of the Electron Density". In: *ACS Central Science* 5.1 (Dec. 2018), pp. 57–64. ISSN: 2374-7951. DOI: 10.1021/acscentsci.8b00551.

[37] Alan M. Lewis et al. "Learning Electron Densities in the Condensed Phase". In: *Journal of Chemical Theory and Computation* 17.11 (Oct. 2021), pp. 7203–7214. ISSN: 1549-9626. DOI: 10.1021/acs.jctc.1c00576.

[38] Ksenia R. Briling, Alberto Fabrizio, and Clemence Corminboeuf. "Impact of quantum-chemical metrics on the machine learning prediction of electron density". In: *The Journal of Chemical Physics* 155.2 (July 2021). ISSN: 1089-7690. DOI: 10.1063/5.0055393.

[39] Joshua A Rackers et al. "A recipe for cracking the quantum scaling limit with machine learned electron densities". In: *Machine Learning: Science and Technology* 4.1 (Feb. 2023), p. 015027. ISSN: 2632-2153. DOI: 10.1088/2632-2153/acb314.

[40] S. Hazra, U. Patil, and S. Sanvito. "Predicting the One-Particle Density Matrix with Machine Learning". In: *Journal of Chemical Theory and Computation* 20.11 (May 2024), pp. 4569–4578. ISSN: 1549-9626. DOI: 10.1021/acs.jctc.4c00042.

[41] Cas van der Oord et al. "Hyperactive learning for data-driven interatomic potentials". In: *npj Computational Materials* 9 (2023), p. 168. DOI: 10.1038/s41524-023-01104-6.

[42] Zhen Xie and Joel M. Bowman. "Permutationally Invariant Polynomial Basis for Molecular Energy Surface Fitting via Monomial Symmetrization". In: *Journal of Chemical Theory and Computation* 6.1 (Nov. 2009), pp. 26–34. ISSN: 1549-9626. DOI: 10.1021/ct9004917.

[43] Cas van der Oord et al. "Regularised atomic body-ordered permutation-invariant polynomials for the construction of interatomic potentials". In: *Machine Learning: Science and Technology* 1.1 (2020), p. 015004. DOI: 10.1088/2632-2153/ab527c.

[44] Albert P. Bartók, Risi Kondor, and Gábor Csányi. "On representing chemical environments". In: *Physical Review B* 87.18 (May 2013). ISSN: 1550-235X. DOI: 10.1103/physrevb.87.184115.

[45] Ralf Drautz. "Atomic cluster expansion for accurate and transferable interatomic potentials". In: *Physical Review B* 99.1 (Jan. 2019). ISSN: 2469-9969. DOI: 10.1103/physrevb.99.014104.

[46] Jörg Behler. "Atom-centered symmetry functions for constructing high-dimensional neural network potentials". In: *The Journal of Chemical Physics* 134.7 (Feb. 2011). ISSN: 1089-7690. DOI: 10.1063/1.3553717.

[47] Jörg Behler. "Four Generations of High-Dimensional Neural Network Potentials". In: *Chemical Reviews* 121.16 (Mar. 2021), pp. 10037–10072. ISSN: 1520-6890. DOI: 10.1021/acs.chemrev.0c00868.

[48] Yury Lysogorskiy et al. "Performant implementation of the atomic cluster expansion (PACE) and application to copper and silicon". In: *npj Computational Materials* 7.1 (2021), p. 97. DOI: 10.1038/s41524-021-00559-9.

[49] Ralf Drautz. "Atomic cluster expansion of scalar, vectorial, and tensorial properties including magnetism and charge transfer". In: *Physical Review B* 102.2 (2020), p. 024104. DOI: 10.1103/PhysRevB.102.024104.

[50]  Per-Olov Löwdin. "On the Non-Orthogonality Problem Connected with the Use of Atomic Wave Functions in the Theory of Molecules and Crystals". In: *The Journal of Chemical Physics* 18.3 (1956), pp. 365–375. DOI: 10.1063/1.1748067.

[51]  Geneviève Dusson et al. "Atomic Cluster Expansion: Completeness, Efficiency and Stability". In: *Journal of Computational Physics* 454 (2022), p. 110946. DOI: 10.1016/j.jcp.2022.110946.

[52]  E. Canc'es, B. Mennucci, and J. Tomasi. "A new integral equation formalism for the polarizable continuum model: Theoretical background and applications to isotropic and anisotropic dielectrics". In: *The Journal of Chemical Physics* 107.8 (Aug. 1997), pp. 3032–3041. ISSN: 1089-7690. DOI: 10.1063/1.474659.

[53]  Pekka Mark and Lennart Nilsson. "Structure and Dynamics of the TIP3P, SPC, and SPC/E Water Models at 298 K". In: *The Journal of Physical Chemistry A* 105.43 (Oct. 2001), pp. 9954–9960. ISSN: 1520-5215. DOI: 10.1021/jp003020w.

[54]  Michael Gaus, Qiang Cui, and Marcus Elstner. "DFTB3: Extension of the Self-Consistent-Charge Density-Functional Tight-Binding Method (SCC-DFTB)". In: *Journal of Chemical Theory and Computation* 7.4 (Mar. 2011), pp. 931–948. ISSN: 1549-9626. DOI: 10.1021/ct100684s.

[55]  Michael Gaus et al. "Parameterization of DFTB3/3OB for Sulfur and Phosphorus for Chemical and Biological Applications". In: *Journal of Chemical Theory and Computation* 10.4 (2014), pp. 1518–1537. DOI: 10.1021/ct401002w.

[56]  Xiya Lu et al. "Parametrization of DFTB3/3OB for Magnesium and Zinc for Chemical and Biological Applications". In: *The Journal of Physical Chemistry B* 119.3 (2015), pp. 1062–1082. DOI: 10.1021/jp506557r.

[57]  Tom Darden, Darrin York, and Lee Pedersen. "Particle mesh Ewald: An N·log(N) method for Ewald sums in large systems". In: *The Journal of Chemical Physics* 98.12 (June 1993), pp. 10089–10092. ISSN: 1089-7690. DOI: 10.1063/1.464397.

[58]  D. A. Case et al. *AMBER 2022*. University of California, San Francisco. 2022.

[59]  M. J. Frisch et al. *Gaussian˜16 Revision C.01*. Gaussian Inc. Wallingford CT. 2016.

[60]  M. J. Frisch et al. *Gaussian Development Version, Revision J.19*. Gaussian, Inc., Wallingford CT, 2020. 2020.

[61]  *GauOpen*. https://gaussian.com/interfacing/, Accessed 20 Feb. 2025.

# Supporting Information for "A symmetry-preserving and transferable representation for learning the Kohn-Sham density matrix"

Liwei Zhang[1], Patrizia Mazzeo[2], Michele Nottoli[3], Edoardo Cignoni[2], Lorenzo Cupellini[2], and Benjamin Stamm[3]

[1]Institut für Geometrie und Praktische Mathematik, RWTH Aachen University, Templergraben 55, 52062 Aachen, Germany

[2]Dipartimento di Chimica e Chimica Industriale, Università di Pisa, 56124 Pisa, Italy

[3]Universität Stuttgart, Institute of Applied Analysis and Numerical Simulation, Pfaffenwaldring 57, 70569 Stuttgart, Germany

# 1 Additional tables

| Molecule | Specific Model | | | | Unified Model | | | | Unified Model-A | | | Unified-Alcohol-Model | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | training frames | test frames | $r_{cut}=4.0$ | $r_{cut}=6.5$ | training frames | test frames | $r_{cut}=4.0$ | $r_{cut}=6.5$ | training frames | test frames | $r_{cut}=4.0$ | training frames | test frames | $r_{cut}=4.0$ |
| Acetaldehyde | 0:2999 | 3000:9999 | $4.416 \cdot 10^{-5}$ | $4.240 \cdot 10^{-5}$ | 0:30:8999 | 9000:9999 | $3.278 \cdot 10^{-4}$ | $3.861 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $3.299 \cdot 10^{-4}$ | - | - | - |
| Acrolein | 0:2:5999 | 6000:9999 | $2.514 \cdot 10^{-4}$ | $1.409 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.028 \cdot 10^{-4}$ | $4.277 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.081 \cdot 10^{-4}$ | - | - | - |
| Aniline | 0:2:2999 | 3000:9999 | $4.300 \cdot 10^{-4}$ | $1.634 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $4.868 \cdot 10^{-4}$ | $3.177 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $4.876 \cdot 10^{-4}$ | - | - | - |
| o-Toluidine | 0:2:2999 | 3000:9999 | $5.430 \cdot 10^{-4}$ | $2.738 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.962 \cdot 10^{-4}$ | $3.584 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.940 \cdot 10^{-4}$ | - | - | - |
| m-Toluidine | 0:2:2999 | 3000:9999 | $5.384 \cdot 10^{-4}$ | $2.596 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.824 \cdot 10^{-4}$ | $3.727 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.822 \cdot 10^{-4}$ | - | - | - |
| Benzene* | - | - | - | - | | 0:99 | $4.058 \cdot 10^{-4}$ | $8.638 \cdot 10^{-3}$ | | 0:99 | $3.646 \cdot 10^{-4}$ | - | - | - |
| Toluene* | - | - | - | - | | 0:99 | $6.369 \cdot 10^{-4}$ | $7.674 \cdot 10^{-3}$ | | 0:99 | $5.980 \cdot 10^{-4}$ | - | - | - |
| Phenol** | - | - | - | - | | 0:99 | $4.809 \cdot 10^{-3}$ | $1.795 \cdot 10^{-2}$ | 0:10:89 | others | $6.770 \cdot 10^{-4}$ | - | - | - |
| Benzaldehyde** | - | - | - | - | | 0:99 | $4.129 \cdot 10^{-3}$ | $1.997 \cdot 10^{-2}$ | 0:10:89 | others | $1.201 \cdot 10^{-3}$ | - | - | - |
| p-Toluidine** | - | - | - | - | | 0:99 | $2.840 \cdot 10^{-3}$ | $1.907 \cdot 10^{-3}$ | 0:10:89 | others | $6.293 \cdot 10^{-4}$ | - | - | - |
| Propanol | 0:2999 | 3000:9999 | $3.049 \cdot 10^{-4}$ | $2.421 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $4.427 \cdot 10^{-4}$ | $4.080 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $4.444 \cdot 10^{-4}$ | 0:24:8999 | 9000:9999 | $3.932 \cdot 10^{-4}$ |
| Butanol | 0:3:8999 | 9000:9999 | $4.510 \cdot 10^{-4}$ | $7.643 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $4.921 \cdot 10^{-4}$ | $1.273 \cdot 10^{-3}$ | 0:30:8999 | 9000:9999 | $4.934 \cdot 10^{-4}$ | 0:24:8999 | 9000:9999 | $4.477 \cdot 10^{-4}$ |
| 2-Butanol | 0:2999 | 3000:9999 | $9.173 \cdot 10^{-4}$ | $7.603 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $1.494 \cdot 10^{-3}$ | $6.581 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $1.509 \cdot 10^{-3}$ | 0:24:8999 | 9000:9999 | $7.598 \cdot 10^{-4}$ |
| Hexanol | 0:2:2999 | 3000:9999 | $1.031 \cdot 10^{-3}$ | $7.021 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.423 \cdot 10^{-4}$ | $9.614 \cdot 10^{-4}$ | 0:30:8999 | 9000:9999 | $5.314 \cdot 10^{-4}$ | 0:24:8999 | 9000:9999 | $4.602 \cdot 10^{-4}$ |
| Ethanol* | - | - | - | - | | 0:99 | $9.644 \cdot 10^{-4}$ | $1.021 \cdot 10^{-3}$ | | 0:99 | $8.999 \cdot 10^{-4}$ | | 0:99 | $5.556 \cdot 10^{-4}$ |
| 2-Propanol* | - | - | - | - | | 0:99 | $7.701 \cdot 10^{-4}$ | $7.573 \cdot 10^{-4}$ | | 0:99 | $7.480 \cdot 10^{-4}$ | | 0:99 | $4.686 \cdot 10^{-4}$ |
| 2-Hexanol* | - | - | - | - | | 0:99 | $7.384 \cdot 10^{-4}$ | $4.853 \cdot 10^{-3}$ | | 0:99 | $7.353 \cdot 10^{-4}$ | | 0:99 | $5.485 \cdot 10^{-4}$ |
| Heptanol* | - | - | - | - | | 0:99 | $8.896 \cdot 10^{-4}$ | $1.538 \cdot 10^{-3}$ | | 0:99 | $8.980 \cdot 10^{-4}$ | | 0:99 | $7.048 \cdot 10^{-4}$ |

Table S1: test set RMSEs for different molecules obtained by (3,8)-models trained with different datasets. In the table, when "-" appears, it means that the corresponding molecule is not included in the training process, and hence there is no dedicated model for these molecules.

3

| Molecule | Tol. | Default guess | Specific Model $r_{\text{cut}} = 4.0$ | $r_{\text{cut}} = 6.5$ |
|---|---|---|---|---|
| Acetaldehyde | $10^{-6}$ | $9.4 \pm 0.1$ | $5.2 \pm 0.1 \ (\sim 44\%)$ | $5.1 \pm 0.1 \ (\sim 46\%)$ |
| | $10^{-7}$ | $12.4 \pm 0.2$ | $7.8 \pm 0.1 \ (\sim 37\%)$ | $7.8 \pm 0.1 \ (\sim 37\%)$ |
| | $10^{-8}$ | $14.9 \pm 0.1$ | $9.8 \pm 0.1 \ (\sim 34\%)$ | $9.7 \pm 0.1 \ (\sim 35\%)$ |
| Acrolein | $10^{-6}$ | $10.5 \pm 0.1$ | $7.5 \pm 0.2 \ (\sim 29\%)$ | $6.3 \pm 0.2 \ (\sim 40\%)$ |
| | $10^{-7}$ | $13.4 \pm 0.1$ | $9.8 \pm 0.1 \ (\sim 26\%)$ | $9.1 \pm 0.1 \ (\sim 32\%)$ |
| | $10^{-8}$ | $16.0 \pm 0.2$ | $11.5 \pm 0.1 \ (\sim 28\%)$ | $11.2 \pm 0.1 \ (\sim 30\%)$ |
| Aniline | $10^{-6}$ | $9.9 \pm 0.1$ | $7.2 \pm 0.1 \ (\sim 28\%)$ | $6.5 \pm 0.1 \ (\sim 35\%)$ |
| | $10^{-7}$ | $12.5 \pm 0.1$ | $9.7 \pm 0.1 \ (\sim 23\%)$ | $8.7 \pm 0.1 \ (\sim 31\%)$ |
| | $10^{-8}$ | $14.7 \pm 0.1$ | $11.9 \pm 0.1 \ (\sim 19\%)$ | $11.1 \pm 0.1 \ (\sim 24\%)$ |
| o-Toluidine | $10^{-6}$ | $10.0 \pm 0.0$ | $7.6 \pm 0.1 \ (\sim 24\%)$ | $7.2 \pm 0.1 \ (\sim 28\%)$ |
| | $10^{-7}$ | $12.8 \pm 0.1$ | $9.9 \pm 0.1 \ (\sim 22\%)$ | $9.7 \pm 0.1 \ (\sim 24\%)$ |
| | $10^{-8}$ | $14.9 \pm 0.0$ | $12.3 \pm 0.1 \ (\sim 18\%)$ | $11.7 \pm 0.1 \ (\sim 22\%)$ |
| m-Toluidine | $10^{-6}$ | $10.0 \pm 0.0$ | $7.3 \pm 0.1 \ (\sim 27\%)$ | $7.1 \pm 0.1 \ (\sim 29\%)$ |
| | $10^{-7}$ | $12.7 \pm 0.1$ | $9.8 \pm 0.1 \ (\sim 23\%)$ | $9.5 \pm 0.1 \ (\sim 25\%)$ |
| | $10^{-8}$ | $14.8 \pm 0.1$ | $11.9 \pm 0.1 \ (\sim 19\%)$ | $11.7 \pm 0.1 \ (\sim 21\%)$ |
| 1-Propanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.0 \pm 0.1 \ (\sim 33\%)$ | $6.0 \pm 0.1 \ (\sim 33\%)$ |
| | $10^{-7}$ | $11.0 \pm 0.1$ | $8.0 \pm 0.0 \ (\sim 27\%)$ | $8.0 \pm 0.0 \ (\sim 27\%)$ |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $9.8 \pm 0.1 \ (\sim 25\%)$ | $9.8 \pm 0.1 \ (\sim 26\%)$ |
| 1-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.4 \pm 0.1 \ (\sim 29\%)$ | $6.4 \pm 0.1 \ (\sim 29\%)$ |
| | $10^{-7}$ | $11.0 \pm 0.0$ | $8.1 \pm 0.1 \ (\sim 27\%)$ | $8.2 \pm 0.1 \ (\sim 26\%)$ |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $10.1 \pm 0.1 \ (\sim 24\%)$ | $10.2 \pm 0.1 \ (\sim 23\%)$ |
| 2-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $7.2 \pm 0.1 \ (\sim 20\%)$ | $7.5 \pm 0.1 \ (\sim 16\%)$ |
| | $10^{-7}$ | $11.1 \pm 0.1$ | $9.0 \pm 0.1 \ (\sim 19\%)$ | $9.3 \pm 0.1 \ (\sim 16\%)$ |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $10.9 \pm 0.1 \ (\sim 18\%)$ | $11.0 \pm 0.1 \ (\sim 17\%)$ |
| 1-Hexanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.4 \pm 0.1 \ (\sim 29\%)$ | $6.6 \pm 0.1 \ (\sim 27\%)$ |
| | $10^{-7}$ | $10.9 \pm 0.1$ | $8.2 \pm 0.1 \ (\sim 24\%)$ | $8.3 \pm 0.1 \ (\sim 23\%)$ |
| | $10^{-8}$ | $13.0 \pm 0.0$ | $10.3 \pm 0.1 \ (\sim 21\%)$ | $10.4 \pm 0.1 \ (\sim 20\%)$ |

Table S2: Average of the number of iterations obtained using the predicted density matrix as guess for the calculation on the same geometry. The results refer to the (3,8)-models trained targeting the density matrix. The first column specifies the molecule, the second column indicates the convergence tolerance for the SCF procedure, the third column reports the average of the number of iterations obtained using the default guess available in the Gaussian program and in the last two columns we report the average number of iterations obtained with our predicted density matrix as guess, using the specific model for each molecule and with 4.0 and 6.5 radius of cutoff, respectively. The values reported in parentheses indicate the percentage of reduction with respect to the iterations required by the default guess to converge.

| Molecule | Tol. | Default guess | Specific Model | |
| | | | $r_{\text{cut}} = 4.0$ | $r_{\text{cut}} = 6.5$ |
|---|---|---|---|---|
| Acetaldehyde | $10^{-6}$ | $9.4 \pm 0.1$ | $4.8 \pm 0.1\ (\sim 49\%)$ | $4.7 \pm 0.1\ (\sim 50\%)$ |
| | $10^{-7}$ | $12.4 \pm 0.2$ | $7.3 \pm 0.2\ (\sim 41\%)$ | $7.3 \pm 0.2\ (\sim 41\%)$ |
| | $10^{-8}$ | $14.9 \pm 0.1$ | $9.3 \pm 0.1\ (\sim 38\%)$ | $9.2 \pm 0.1\ (\sim 38\%)$ |
| Acrolein | $10^{-6}$ | $10.5 \pm 0.1$ | $7.2 \pm 0.2\ (\sim 32\%)$ | $5.9 \pm 0.2\ (\sim 44\%)$ |
| | $10^{-7}$ | $13.4 \pm 0.1$ | $9.6 \pm 0.1\ (\sim 29\%)$ | $8.6 \pm 0.1\ (\sim 36\%)$ |
| | $10^{-8}$ | $16.0 \pm 0.2$ | $11.2 \pm 0.1\ (\sim 30\%)$ | $10.8 \pm 0.1\ (\sim 32\%)$ |
| Aniline | $10^{-6}$ | $9.9 \pm 0.1$ | $7.0 \pm 0.1\ (\sim 29\%)$ | $6.3 \pm 0.1\ (\sim 36\%)$ |
| | $10^{-7}$ | $12.5 \pm 0.1$ | $9.4 \pm 0.2\ (\sim 25\%)$ | $8.8 \pm 0.2\ (\sim 30\%)$ |
| | $10^{-8}$ | $14.7 \pm 0.1$ | $11.7 \pm 0.2\ (\sim 20\%)$ | $11.0 \pm 0.2\ (\sim 25\%)$ |
| o-Toluidine | $10^{-6}$ | $10.0 \pm 0.0$ | $7.0 \pm 0.1\ (\sim 30\%)$ | $6.9 \pm 0.1\ (\sim 31\%)$ |
| | $10^{-7}$ | $12.8 \pm 0.1$ | $9.6 \pm 0.1\ (\sim 25\%)$ | $9.4 \pm 0.1\ (\sim 27\%)$ |
| | $10^{-8}$ | $14.9 \pm 0.0$ | $11.7 \pm 0.1\ (\sim 21\%)$ | $11.4 \pm 0.1\ (\sim 24\%)$ |
| m-Toluidine | $10^{-6}$ | $10.0 \pm 0.0$ | $7.1 \pm 0.2\ (\sim 29\%)$ | $6.9 \pm 0.2\ (\sim 31\%)$ |
| | $10^{-7}$ | $12.7 \pm 0.1$ | $9.7 \pm 0.2\ (\sim 23\%)$ | $9.4 \pm 0.2\ (\sim 25\%)$ |
| | $10^{-8}$ | $14.8 \pm 0.1$ | $11.9 \pm 0.2\ (\sim 20\%)$ | $11.6 \pm 0.2\ (\sim 22\%)$ |
| 1-Propanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.0 \pm 0.0\ (\sim 33\%)$ | $5.9 \pm 0.0\ (\sim 34\%)$ |
| | $10^{-7}$ | $11.0 \pm 0.1$ | $8.0 \pm 0.1\ (\sim 27\%)$ | $7.8 \pm 0.1\ (\sim 29\%)$ |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $9.9 \pm 0.1\ (\sim 25\%)$ | $9.7 \pm 0.1\ (\sim 27\%)$ |
| 1-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.2 \pm 0.1\ (\sim 31\%)$ | $6.2 \pm 0.1\ (\sim 32\%)$ |
| | $10^{-7}$ | $11.0 \pm 0.0$ | $8.1 \pm 0.1\ (\sim 27\%)$ | $8.1 \pm 0.1\ (\sim 27\%)$ |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $10.1 \pm 0.1\ (\sim 24\%)$ | $10.0 \pm 0.1\ (\sim 24\%)$ |
| 2-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $7.3 \pm 0.1\ (\sim 18\%)$ | $7.1 \pm 0.1\ (\sim 21\%)$ |
| | $10^{-7}$ | $11.1 \pm 0.1$ | $9.3 \pm 0.1\ (\sim 17\%)$ | $9.1 \pm 0.1\ (\sim 18\%)$ |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $11.1 \pm 0.1\ (\sim 16\%)$ | $11.0 \pm 0.1\ (\sim 17\%)$ |
| 1-Hexanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.6 \pm 0.1\ (\sim 27\%)$ | $6.7 \pm 0.1\ (\sim 25\%)$ |
| | $10^{-7}$ | $10.9 \pm 0.1$ | $8.3 \pm 0.1\ (\sim 24\%)$ | $8.6 \pm 0.1\ (\sim 21\%)$ |
| | $10^{-8}$ | $13.0 \pm 0.0$ | $10.2 \pm 0.1\ (\sim 22\%)$ | $10.4 \pm 0.1\ (\sim 20\%)$ |

Table S3: Average of the number of iterations obtained using the predicted density matrix as guess for the calculation on the same geometry. The results refer to the (3,8)-models trained targeting the Kohn-Sham matrix. The first column specifies the molecule, the second column indicates the convergence tolerance for the SCF procedure, the third column reports the average of the number of iterations obtained using the default guess available in the Gaussian program and in the last two columns we report the average number of iterations obtained with our predicted density matrix as guess, using the specific model for each molecule and with 4.0 and 6.5 radius of cutoff, respectively. The values reported in parentheses indicate the percentage of reduction with respect to the iterations required by the default guess to converge.

| Molecule | Tol. | Default guess | Specific Model | Unified Model | Unified Model-A |
|---|---|---|---|---|---|
| Acetaldehyde | $10^{-6}$ | $9.4 \pm 0.1$ | $5.2 \pm 0.1$ ($\sim 44\%$) | $7.5 \pm 0.1$ ($\sim 20\%$) | $7.5 \pm 0.1$ ($\sim 20\%$) |
| | $10^{-7}$ | $12.4 \pm 0.2$ | $7.8 \pm 0.1$ ($\sim 37\%$) | $9.7 \pm 0.1$ ($\sim 22\%$) | $9.7 \pm 0.1$ ($\sim 22\%$) |
| | $10^{-8}$ | $14.9 \pm 0.1$ | $9.8 \pm 0.1$ ($\sim 34\%$) | $11.1 \pm 0.1$ ($\sim 26\%$) | $11.1 \pm 0.1$ ($\sim 25\%$) |
| Acrolein | $10^{-6}$ | $10.5 \pm 0.1$ | $7.5 \pm 0.2$ ($\sim 29\%$) | $8.3 \pm 0.2$ ($\sim 21\%$) | $8.3 \pm 0.2$ ($\sim 21\%$) |
| | $10^{-7}$ | $13.4 \pm 0.1$ | $9.8 \pm 0.1$ ($\sim 26\%$) | $10.5 \pm 0.1$ ($\sim 22\%$) | $10.5 \pm 0.1$ ($\sim 21\%$) |
| | $10^{-8}$ | $16.0 \pm 0.2$ | $11.5 \pm 0.1$ ($\sim 28\%$) | $12.0 \pm 0.1$ ($\sim 25\%$) | $12.1 \pm 0.1$ ($\sim 24\%$) |
| Aniline | $10^{-6}$ | $9.9 \pm 0.1$ | $7.2 \pm 0.1$ ($\sim 28\%$) | $7.5 \pm 0.1$ ($\sim 24\%$) | $7.5 \pm 0.1$ ($\sim 24\%$) |
| | $10^{-7}$ | $12.5 \pm 0.1$ | $9.7 \pm 0.1$ ($\sim 23\%$) | $9.9 \pm 0.1$ ($\sim 21\%$) | $9.9 \pm 0.1$ ($\sim 21\%$) |
| | $10^{-8}$ | $14.7 \pm 0.1$ | $11.9 \pm 0.1$ ($\sim 19\%$) | $12.1 \pm 0.1$ ($\sim 18\%$) | $12.1 \pm 0.1$ ($\sim 18\%$) |
| o-Toluidine | $10^{-6}$ | $10.0 \pm 0.0$ | $7.6 \pm 0.1$ ($\sim 24\%$) | $7.8 \pm 0.1$ ($\sim 22\%$) | $7.8 \pm 0.1$ ($\sim 22\%$) |
| | $10^{-7}$ | $12.8 \pm 0.1$ | $9.9 \pm 0.1$ ($\sim 22\%$) | $10.1 \pm 0.1$ ($\sim 21\%$) | $10.1 \pm 0.1$ ($\sim 21\%$) |
| | $10^{-8}$ | $14.9 \pm 0.0$ | $12.3 \pm 0.1$ ($\sim 18\%$) | $12.5 \pm 0.1$ ($\sim 16\%$) | $12.5 \pm 0.1$ ($\sim 16\%$) |
| m-Toluidine | $10^{-6}$ | $10.0 \pm 0.0$ | $7.3 \pm 0.1$ ($\sim 27\%$) | $7.6 \pm 0.1$ ($\sim 24\%$) | $7.6 \pm 0.1$ ($\sim 24\%$) |
| | $10^{-7}$ | $12.7 \pm 0.1$ | $9.8 \pm 0.1$ ($\sim 23\%$) | $9.9 \pm 0.1$ ($\sim 22\%$) | $9.9 \pm 0.1$ ($\sim 22\%$) |
| | $10^{-8}$ | $14.8 \pm 0.1$ | $11.9 \pm 0.1$ ($\sim 19\%$) | $12.1 \pm 0.1$ ($\sim 18\%$) | $12.1 \pm 0.1$ ($\sim 18\%$) |
| Benzene* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $7.4 \pm 0.1$ ($\sim 18\%$) | $7.4 \pm 0.1$ ($\sim 18\%$) |
| | $10^{-7}$ | $11.3 \pm 0.1$ | - | $9.7 \pm 0.1$ ($\sim 14\%$) | $9.7 \pm 0.1$ ($\sim 14\%$) |
| | $10^{-8}$ | $13.5 \pm 0.1$ | - | $11.7 \pm 0.1$ ($\sim 13\%$) | $11.7 \pm 0.1$ ($\sim 13\%$) |
| Toluene* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $8.0 \pm 0.0$ ($\sim 11\%$) | $7.9 \pm 0.1$ ($\sim 12\%$) |
| | $10^{-7}$ | $11.6 \pm 0.1$ | - | $10.1 \pm 0.1$ ($\sim 13\%$) | $10.0 \pm 0.0$ ($\sim 14\%$) |
| | $10^{-8}$ | $13.9 \pm 0.1$ | - | $12.4 \pm 0.1$ ($\sim 11\%$) | $12.3 \pm 0.1$ ($\sim 12\%$) |
| Phenol** | $10^{-6}$ | $9.9 \pm 0.1$ | - | $10.0 \pm 0.1$ ($\sim -1\%$) | $8.1 \pm 0.1$ ($\sim 18\%$) |
| | $10^{-7}$ | $12.2 \pm 0.1$ | - | $12.2 \pm 0.1$ ($\sim 1\%$) | $10.5 \pm 0.1$ ($\sim 14\%$) |
| | $10^{-8}$ | $14.6 \pm 0.1$ | - | $14.4 \pm 0.1$ ($\sim 2\%$) | $12.3 \pm 0.1$ ($\sim 15\%$) |
| Benzaldehyde** | $10^{-6}$ | $10.7 \pm 0.1$ | - | $10.5 \pm 0.1$ ($\sim 2\%$) | $9.0 \pm 0.0$ ($\sim 16\%$) |
| | $10^{-7}$ | $13.2 \pm 0.1$ | - | $12.9 \pm 0.0$ ($\sim 2\%$) | $11.5 \pm 0.1$ ($\sim 14\%$) |
| | $10^{-8}$ | $15.7 \pm 0.1$ | - | $15.1 \pm 0.1$ ($\sim 4\%$) | $14.2 \pm 0.1$ ($\sim 10\%$) |
| p-Toluidine** | $10^{-6}$ | $10.0 \pm 0.0$ | - | $8.3 \pm 0.1$ ($\sim 16\%$) | $7.8 \pm 0.1$ ($\sim 21\%$) |
| | $10^{-7}$ | $12.6 \pm 0.1$ | - | $11.1 \pm 0.1$ ($\sim 12\%$) | $10.2 \pm 0.1$ ($\sim 20\%$) |
| | $10^{-8}$ | $14.9 \pm 0.1$ | - | $13.3 \pm 0.1$ ($\sim 11\%$) | $12.6 \pm 0.1$ ($\sim 15\%$) |

**Continued on the next page**

Table S4: Average of the number of iterations obtained using the predicted density matrix as guess for the calculation on the same geometry. The results refer to the (3,8)-models trained targeting the density matrix ($r_{\text{cut}} = 4.0$). The first column specifies the molecule, the second column indicates the convergence tolerance for the SCF procedure, the third column reports the average of the number of iterations obtained using the default guess available in the Gaussian program and in the last three columns we report the average number of iterations obtained with our predicted density matrix as guess, using the specific models, the unified model and the Unified Model-A, respectively. The values reported in parentheses indicate the percentage of reduction with respect to the iterations required by the default guess to converge. The test molecules with a superscript * are not included in the training process at all, and those with ** are involved in the training of Unified Model-A, with only 10 frames each included.

| | | | | | |
|---|---|---|---|---|---|
| 1-Propanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.0 \pm 0.1$ ($\sim 33\%$) | $6.5 \pm 0.1$ ($\sim 27\%$) | $6.5 \pm 0.1$ ($\sim 28\%$) |
| | $10^{-7}$ | $11.0 \pm 0.1$ | $8.0 \pm 0.0$ ($\sim 27\%$) | $8.1 \pm 0.1$ ($\sim 26\%$) | $8.1 \pm 0.1$ ($\sim 26\%$) |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $9.8 \pm 0.1$ ($\sim 25\%$) | $10.1 \pm 0.1$ ($\sim 23\%$) | $10.1 \pm 0.1$ ($\sim 23\%$) |
| 1-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.4 \pm 0.1$ ($\sim 29\%$) | $6.6 \pm 0.1$ ($\sim 27\%$) | $6.5 \pm 0.1$ ($\sim 27\%$) |
| | $10^{-7}$ | $11.0 \pm 0.0$ | $8.1 \pm 0.1$ ($\sim 27\%$) | $8.2 \pm 0.1$ ($\sim 26\%$) | $8.2 \pm 0.1$ ($\sim 26\%$) |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $10.1 \pm 0.1$ ($\sim 24\%$) | $10.2 \pm 0.1$ ($\sim 23\%$) | $10.2 \pm 0.1$ ($\sim 22\%$) |
| 2-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $7.2 \pm 0.1$ ($\sim 20\%$) | $7.0 \pm 0.1$ ($\sim 22\%$) | $7.1 \pm 0.1$ ($\sim 22\%$) |
| | $10^{-7}$ | $11.1 \pm 0.1$ | $9.0 \pm 0.1$ ($\sim 19\%$) | $8.8 \pm 0.1$ ($\sim 21\%$) | $8.8 \pm 0.1$ ($\sim 21\%$) |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $10.9 \pm 0.1$ ($\sim 18\%$) | $10.8 \pm 0.1$ ($\sim 19\%$) | $10.8 \pm 0.1$ ($\sim 19\%$) |
| 1-Hexanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.4 \pm 0.1$ ($\sim 29\%$) | $6.5 \pm 0.1$ ($\sim 28\%$) | $6.5 \pm 0.1$ ($\sim 28\%$) |
| | $10^{-7}$ | $10.9 \pm 0.1$ | $8.2 \pm 0.1$ ($\sim 24\%$) | $8.1 \pm 0.1$ ($\sim 25\%$) | $8.2 \pm 0.1$ ($\sim 25\%$) |
| | $10^{-8}$ | $13.0 \pm 0.0$ | $10.3 \pm 0.1$ ($\sim 21\%$) | $10.2 \pm 0.1$ ($\sim 21\%$) | $10.2 \pm 0.1$ ($\sim 21\%$) |
| Ethanol* | $10^{-6}$ | $9.1 \pm 0.0$ | - | $7.2 \pm 0.1$ ($\sim 21\%$) | $7.1 \pm 0.1$ ($\sim 21\%$) |
| | $10^{-7}$ | $11.3 \pm 0.1$ | - | $9.0 \pm 0.0$ ($\sim 20\%$) | $9.0 \pm 0.0$ ($\sim 20\%$) |
| | $10^{-8}$ | $13.7 \pm 0.1$ | - | $10.9 \pm 0.0$ ($\sim 20\%$) | $10.8 \pm 0.1$ ($\sim 21\%$) |
| 2-Propanol* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $7.1 \pm 0.1$ ($\sim 21\%$) | $7.0 \pm 0.0$ ($\sim 22\%$) |
| | $10^{-7}$ | $11.3 \pm 0.1$ | - | $9.1 \pm 0.1$ ($\sim 20\%$) | $9.0 \pm 0.1$ ($\sim 20\%$) |
| | $10^{-8}$ | $13.7 \pm 0.1$ | - | $11.0 \pm 0.0$ ($\sim 19\%$) | $10.9 \pm 0.0$ ($\sim 20\%$) |
| 2-Hexanol* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $7.0 \pm 0.1$ ($\sim 23\%$) | $7.0 \pm 0.1$ ($\sim 22\%$) |
| | $10^{-7}$ | $11.0 \pm 0.0$ | - | $8.6 \pm 0.1$ ($\sim 22\%$) | $8.6 \pm 0.1$ ($\sim 22\%$) |
| | $10^{-8}$ | $13.0 \pm 0.0$ | - | $10.5 \pm 0.1$ ($\sim 19\%$) | $10.5 \pm 0.1$ ($\sim 19\%$) |
| 1-Heptanol* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $6.6 \pm 0.1$ ($\sim 27\%$) | $6.6 \pm 0.1$ ($\sim 27\%$) |
| | $10^{-7}$ | $10.9 \pm 0.1$ | - | $8.3 \pm 0.1$ ($\sim 24\%$) | $8.2 \pm 0.1$ ($\sim 24\%$) |
| | $10^{-8}$ | $13.0 \pm 0.0$ | - | $10.2 \pm 0.1$ ($\sim 21\%$) | $10.3 \pm 0.1$ ($\sim 21\%$) |

| Molecule | Tol. | Default guess | Specific Model | Unified Alcohols Model |
|----------|------|---------------|----------------|------------------------|
| 1-Propanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.0 \pm 0.1$ ($\sim 33\%$) | $6.4 \pm 0.1$ ($\sim 29\%$) |
| | $10^{-7}$ | $11.0 \pm 0.1$ | $8.0 \pm 0.0$ ($\sim 27\%$) | $8.1 \pm 0.1$ ($\sim 27\%$) |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $9.8 \pm 0.1$ ($\sim 25\%$) | $10.0 \pm 0.1$ ($\sim 24\%$) |
| 1-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.4 \pm 0.1$ ($\sim 29\%$) | $6.4 \pm 0.1$ ($\sim 29\%$) |
| | $10^{-7}$ | $11.0 \pm 0.0$ | $8.1 \pm 0.1$ ($\sim 27\%$) | $8.1 \pm 0.1$ ($\sim 26\%$) |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $10.1 \pm 0.1$ ($\sim 24\%$) | $10.2 \pm 0.1$ ($\sim 23\%$) |
| 2-Butanol | $10^{-6}$ | $9.0 \pm 0.0$ | $7.2 \pm 0.1$ ($\sim 20\%$) | $6.9 \pm 0.1$ ($\sim 24\%$) |
| | $10^{-7}$ | $11.1 \pm 0.1$ | $9.0 \pm 0.1$ ($\sim 19\%$) | $8.4 \pm 0.1$ ($\sim 24\%$) |
| | $10^{-8}$ | $13.2 \pm 0.1$ | $10.9 \pm 0.1$ ($\sim 18\%$) | $10.5 \pm 0.1$ ($\sim 21\%$) |
| 1-Hexanol | $10^{-6}$ | $9.0 \pm 0.0$ | $6.4 \pm 0.1$ ($\sim 29\%$) | $6.2 \pm 0.1$ ($\sim 31\%$) |
| | $10^{-7}$ | $10.9 \pm 0.1$ | $8.2 \pm 0.1$ ($\sim 24\%$) | $8.1 \pm 0.1$ ($\sim 26\%$) |
| | $10^{-8}$ | $13.0 \pm 0.0$ | $10.3 \pm 0.1$ ($\sim 21\%$) | $10.1 \pm 0.1$ ($\sim 22\%$) |
| Ethanol* | $10^{-6}$ | $9.1 \pm 0.0$ | - | $7.0 \pm 0.0$ ($\sim 23\%$) |
| | $10^{-7}$ | $11.3 \pm 0.1$ | - | $8.4 \pm 0.1$ ($\sim 26\%$) |
| | $10^{-8}$ | $13.7 \pm 0.1$ | - | $10.4 \pm 0.1$ ($\sim 24\%$) |
| 2-Propanol* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $7.0 \pm 0.0$ ($\sim 22\%$) |
| | $10^{-7}$ | $11.3 \pm 0.1$ | - | $8.4 \pm 0.1$ ($\sim 25\%$) |
| | $10^{-8}$ | $13.7 \pm 0.1$ | - | $10.5 \pm 0.1$ ($\sim 23\%$) |
| 2-Hexanol* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $6.7 \pm 0.1$ ($\sim 26\%$) |
| | $10^{-7}$ | $11.0 \pm 0.0$ | - | $8.4 \pm 0.1$ ($\sim 24\%$) |
| | $10^{-8}$ | $13.0 \pm 0.0$ | - | $10.5 \pm 0.1$ ($\sim 20\%$) |
| 1-Heptanol* | $10^{-6}$ | $9.0 \pm 0.0$ | - | $6.4 \pm 0.1$ ($\sim 29\%$) |
| | $10^{-7}$ | $10.9 \pm 0.1$ | - | $8.1 \pm 0.1$ ($\sim 26\%$) |
| | $10^{-8}$ | $13.0 \pm 0.0$ | - | $10.2 \pm 0.1$ ($\sim 22\%$) |

Table S5: Average of the number of iterations obtained using the predicted density matrix as guess for the calculation on the same geometry. The results refer to the (3,8)-models trained targeting the density matrix. The first column specifies the molecule, the second column indicates the convergence tolerance for the SCF procedure, the third column reports the average of the number of iterations obtained using the default guess available in the Gaussian program and in the last two columns we report the average number of iterations obtained with our predicted density matrix as guess, using the specific models and the unified model trained on alcohols, respectively. The values reported in parentheses indicate the percentage of reduction with respect to the iterations required by the default guess to converge. The test molecules with a superscript * are not included in the training process at all.
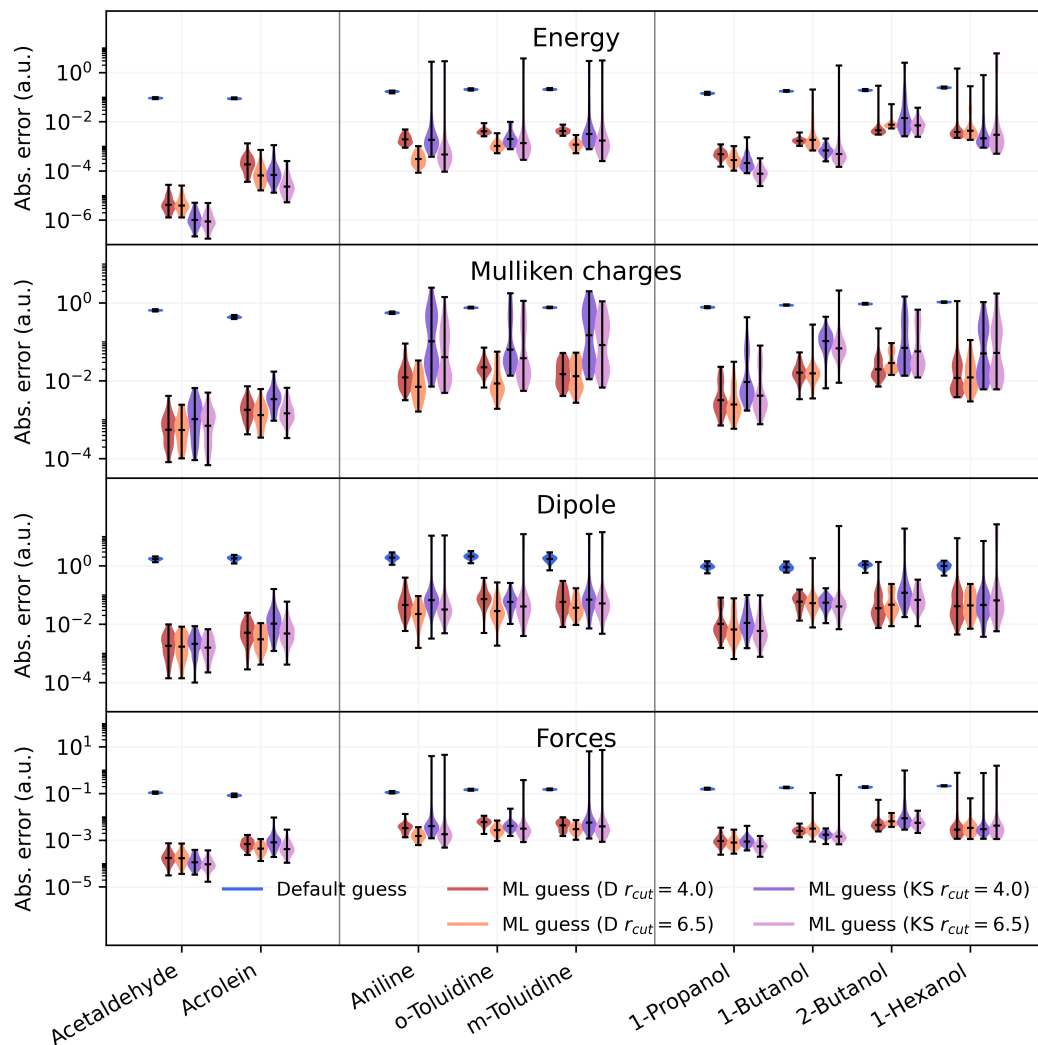
# 2 Additional plots



Figure S1: Plot of the error (in logarithmic scale) for energy, Mulliken charges, dipole moment and forces after a single SCF cycle. The blue line represents the default guess provided by Gaussian, red and orange lines refer to models trained on density matrices with cutoff of 4.0 and 6.5 Å, respectively, and violet and pink lines refer to models trained on Kohn-Sham matrices with cutoff of 4.0 and 6.5 Å, respectively.
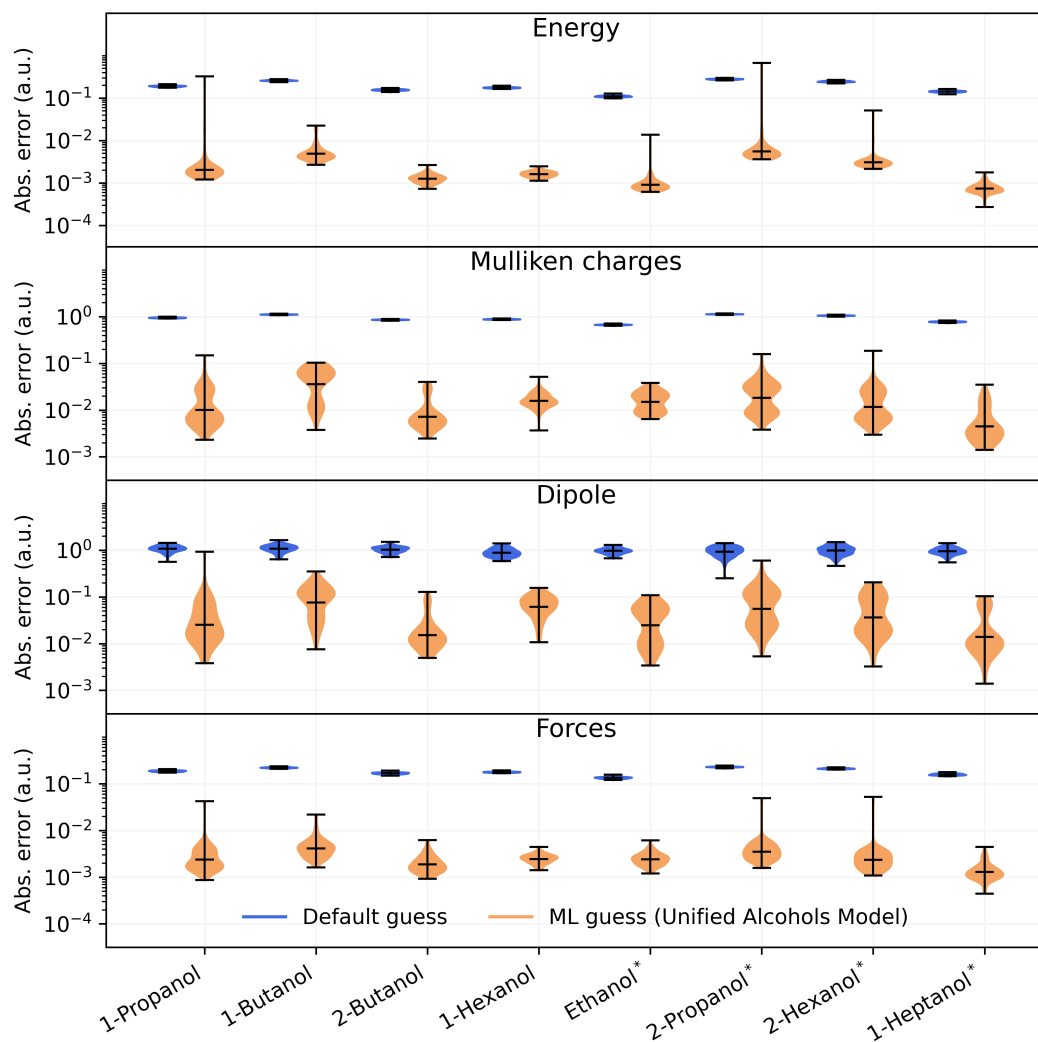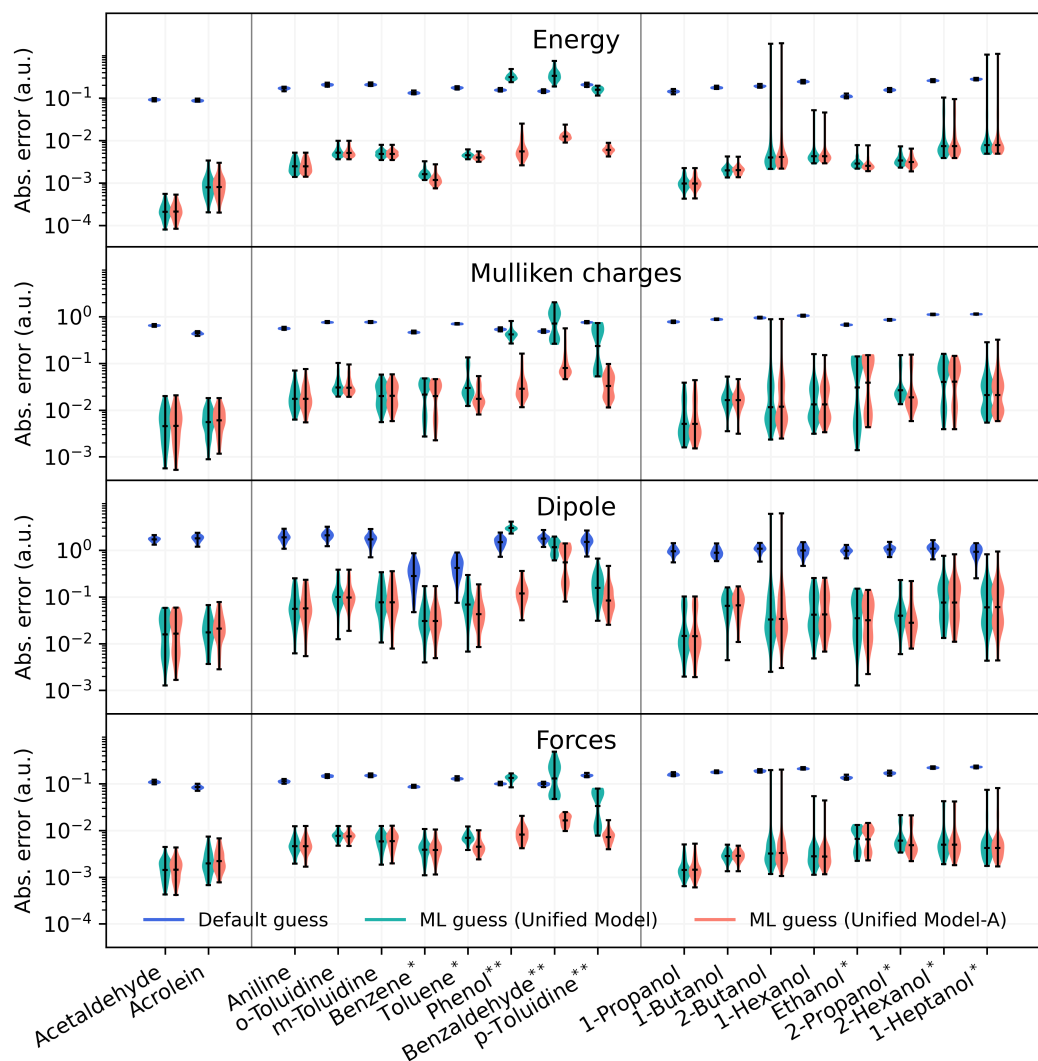
Figure S2: Plot of the error (in logarithmic scale) for energy, Mulliken charges, dipole moment and forces after a single SCF cycle. The blue line represents the default guess provided by Gaussian, while the orange line corresponds to the density matrix predicted using the unified alcohols model. The test molecules with a superscript ∗ are not included in the training process at all.

Figure S3: Plot of the error (in logarithmic scale) for energy, Mulliken charges, dipole moment and forces after a single SCF cycle. The blue line represents the default guess provided by Gaussian, while the sea-green and pink lines correspond to the density matrix predicted using the unified model and Unified Model-A, respectively. The test molecules with a superscript ∗ are not included in the training process at all, and those with ∗∗ are involved in the training of Unified Model-A, with only 10 frames each included.
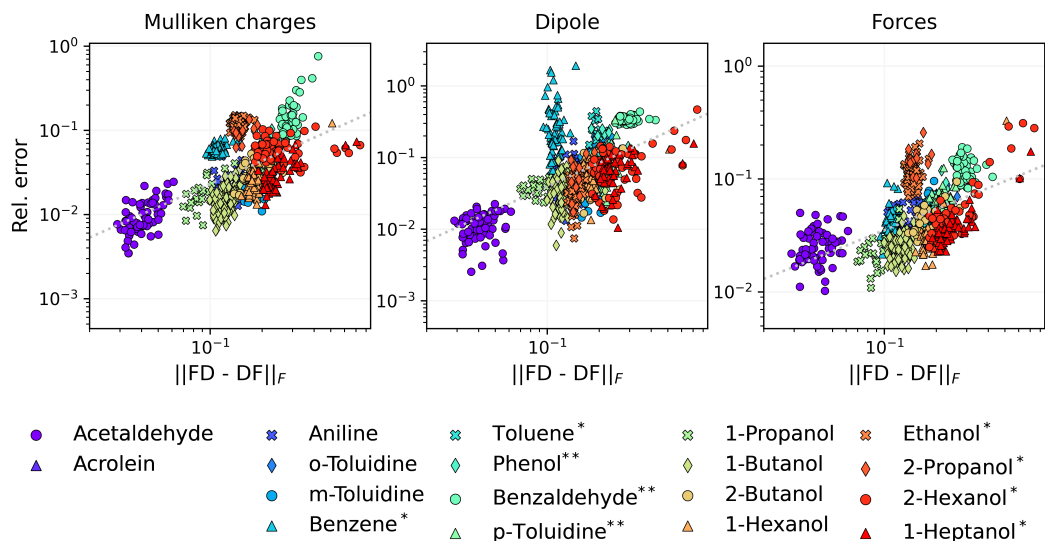
Figure S4: Plot of the Frobenius norm of the commutator between $F$ and $D$ versus the relative error in Mulliken charges, dipole moment and atomic forces, obtained after a single SCF cycle, using the density matrix predicted by Unified Model-A as a guess. The test molecules with a superscript $*$ are not included in the training process at all, and those with $**$ are involved in the training of Unified Model-A, with only 10 frames each included. As expected, for the dipole moment of benzene, which is close to zero, a clear correlation is not observed.
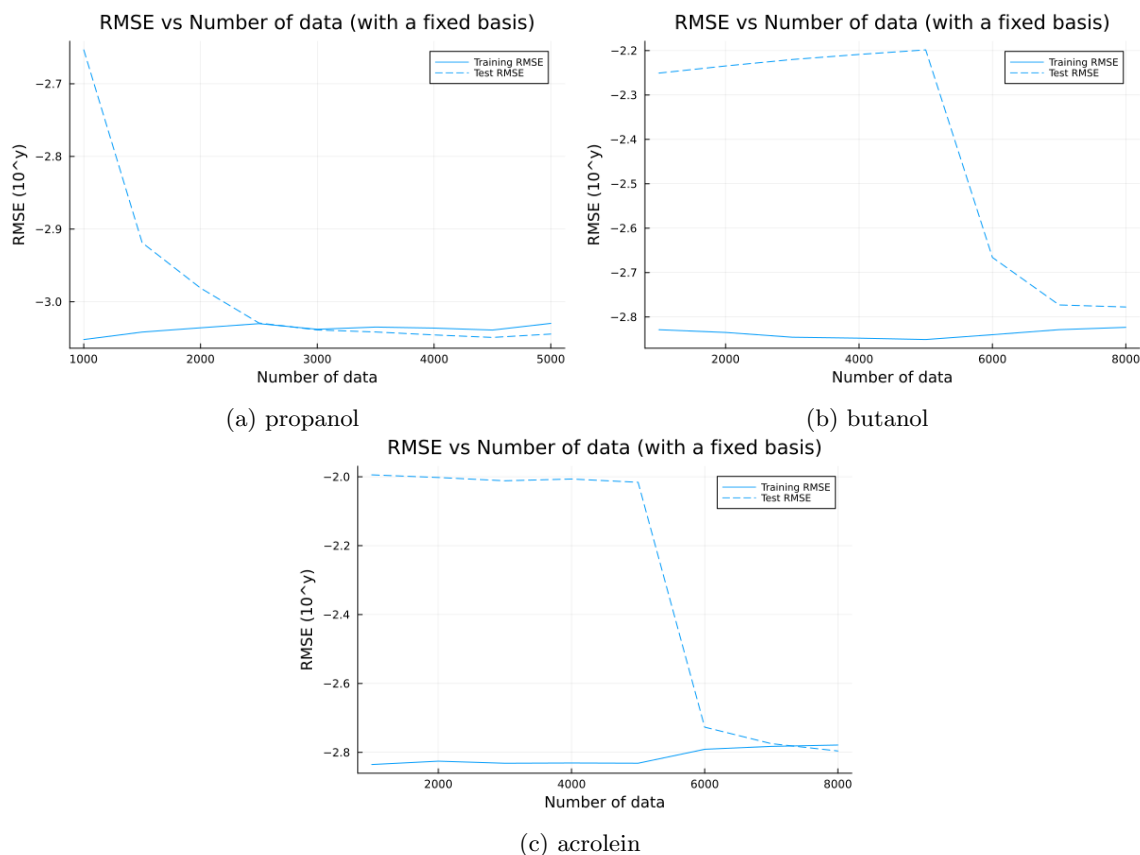


Figure S5: Plot of the relationship between the RMSEs and the number of training configurations for different molecules obtained with a model of moderate size. We pick the amount of training data points for the specific models for each molecule by finding a balance between the training and test set errors.