



Published in final edited form as:

J Chem Theory Comput. 2023 February 28; 19(4): 1261–1275. doi:10.1021/acs.jctc.2c01172.

QD π : A Quantum Deep Potential Interaction Model for Drug Discovery

Jinzhe Zeng,

Yujun Tao,

Timothy J. Giese,

Darrin M. York*

Laboratory for Biomolecular Simulation Research, Institute for Quantitative Biomedicine and Department of Chemistry and Chemical Biology, Rutgers University, Piscataway, NJ 08854, USA

Abstract

We report a QD π -v1.0 for modeling the internal energy of drug molecules containing H, C, N and O atoms. The QD π model is in the form of a quantum mechanical/machine learning potential correction (QM/MLP) that uses a fast 3rd-order self-consistent density-functional tight-binding (DFTB3/3OB) model that is corrected to a quantitatively high-level of accuracy through a deep-learning potential (DeepPot-SE). The model has the advantage that it is able to properly treat electrostatic interactions and handle changes in charge/protonation states. The model is trained against reference data computed at the ω B97X/6–31G* level (as in the ANI-1x data set) and compared to several other approximate semiempirical and machine learning potentials (ANI-1x, ANI-2x, DFTB3, MNDO/d, AM1, PM6, GFN1-xTB and GFN2-xTB). The QD π model is demonstrated to be accurate for a wide range of intra- and intermolecular interactions (despite its intended use as an internal energy model), and has shown to perform exceptionally well for relative protonation/deprotonation energies and tautomers. An example application to model reactions involved in RNA strand cleavage catalyzed by protein and nucleic acid enzymes illustrates the QD π has average errors less than 0.5 kcal/mol, whereas the other models compared have errors over an order of magnitude greater. Taken together, this makes QD π highly attractive as a potential force field model for drug discovery.

Graphical Abstract

*To whom correspondence should be addressed, Darrin.York@rutgers.edu.



1 Introduction

Computational methods that enable the prediction of the binding affinity and selectivity of small molecule drugs to protein or nucleic acid targets are essential tools for drug discovery.^{1–4} Among the most powerful of these methods are so-called “alchemical free energy” (AFE) simulations: physics-based approaches that strive to rigorously calculate the absolute and/or relative binding free energy (ABFE and RBEF, respectively) through atomistic simulations.⁵ The accuracy of such predictions depends critically on the quality and robustness of the underlying potential energy model from which atomic forces are derived.⁶

ABFE and RBEF simulations require the construction of thermodynamic cycles whereby ligands must dynamically sample phase space in “unbound” (aqueous solution) and “target-bound” (e.g., protein complexed) environments. Electrostatic interactions in these environments can differ substantially, and hence a desirable feature of the ligand potential energy model is the ability to explicitly polarize in order to electronically respond to these changes.⁷ Further, roughly 25% of potential drug molecules can exist in alternative tautomeric forms, and almost all of them can have multiple ionizable protonation states. These states are important as they are sensitive to their environment (e.g., pH, ionic conditions, aqueous versus membrane, etc.) and can change upon binding.^{8–10} In order to accommodate these changes, it is advantageous to have a “universal” potential energy model that is not restricted to a specific pre-determined bonding pattern or protonation state within the same simulation, unlike conventional molecular mechanical (MM) force fields (including polarizable force fields). High-level *ab initio* quantum mechanical (QM) models are universal in this sense, and also have been demonstrated to be robust and accurate,¹¹ but these methods require tremendous computational resources making them intractable for routine simulations. Approximate “semiempirical” QM models, on the other hand, are orders of magnitude faster and can be routinely applied in simulations where the QM region is limited to up to a few hundred atoms (which encompasses most drug molecules); however, these models typically do not have the quantitative accuracy that real-world drug discovery applications demand.¹²

An alternative approach is to develop machine learning (ML) potentials that are both fast and accurate within the scope of their training^{13–18}. To date, many such models have been developed and more continue to emerge^{19–35}. So-called “pure” ML potentials face many

challenges for use in free energy simulations. They must be able to model a wide range of intra- and intermolecular interactions,^{36–38} including relative conformational energies,³⁸ hydrogen bonding,³⁹ π stacking, London dispersion, and mixed interactions,^{40,41} in addition to different tautomers⁸ and protonation states^{10,42,43} as mentioned previously. The models must be able to distinguish variable electron number (charge) and spin (multiplicity). Finally, the models not only need to be trained to give back accurate energies and forces for the regions of configurational space expected to be sampled under relevant temperature and pressure conditions, they must also be trained to avoid inaccessible regions of configurational space.

Among the first and most widely recognized ML potentials are the ANI^{30,36,44,45} class of models that to a large degree formed the inspiration for the current work. These pure ML models are both robust and computationally efficient. The ANI models take as basic arguments the positions and identity of atoms in order to return an energy and through derivative relations, a set of forces. However, challenges remain for these models to distinguish different charge and/or spin states, and properly treat electrostatic interactions (although there has been recent progress to determine atomic charges³³). This is currently a serious limitation, as it has been estimated that up to 95% of drug molecules contain ionizable groups that can cause variations in charge state and greatly alter electrostatic interactions.⁴⁶

In the present work, we develop a Quantum Deep-learning Potential Interaction (QD π) model that uses a fast 3rd-order self-consistent density-functional tight-binding (DFTB3/3OB) model^{47,48} that is corrected to a quantitatively high-level of accuracy through a range-corrected deep-learning potential (DPRc).^{49,50} In this way, the QD π model developed here is the form of quantum mechanical/machine learning potential correction (QM/ -MLP).^{35,49–54} The use of DFTB3 as a robust QM base model has several important advantages. First, it provides a reasonable description of the conformational potential energy landscape, greatly reducing the requirement to explicitly train the MLP to avoid inaccessible high-energy regions. Second, DFTB3 uses polarizable atomic charge densities that are easily integrated into efficient particle-mesh Ewald⁵⁵ framework to capture long-range electrostatic interactions in condensed phase QM/MM⁵⁶ and QMFF^{57–59} simulations. Third, DFTB3 is able to model changes in charge, protonation and spin state in a size-consistent manner. The QD π is developed and validated with respect to a number of existing and new databases (DBs).^{36,37,39–43,60–66} Special emphasis is placed on developing a universal model that is able to quantitatively predict tautomers⁸ and protonation states.¹⁰ The present work develops the QD π for internal ligand energetics. This advance sets the stage for intermolecular interactions to be fully developed through quantum mechanical/molecular mechanical (QM/MM) -MLP. This would enable alchemical free energy simulations for drug discovery to be made using the QD π model through the use of indirect MM \rightarrow QD π free energy “book-ending” methods.^{67–70}

2 Methods

This paper brings together several facets in order to develop the QD π model for drug discovery. The first is the collection and curation of several existing molecular databases of

structures, energies and forces. The second is the generation of new data sets that fill needed gaps in training and/or testing data. Third, we develop new tools within DeePMD-kit^{71,72} that enable more general flexible forms of the loss function (including relative energies) used in training of the neural networks. Fourth, we create computational infrastructure for consistent comparison of a wide array of existing potential energy models. Each of these is described in detail below.

2.1 Preparation for data sets

The purpose of this first-generation QD π model is to create a highly robust universal potential that can accurately model drug-like molecules containing H, C, N and O atoms as relevant for binding to biological targets. Important properties for consideration include: relative conformational energies, a wide range of intermolecular interactions, as well as relative energies associated with different tautomers and protonation states. While ultimately this model can be extended to predict covalent binding (irreversible inhibition), the initial focus here is on non-covalent binding.

We prepared several data sets for training and (benchmark) testing of the QD π model. These are summarized in Table 1 and described in more detail below. As a general theme, we endeavored to be consistent with the ANI-1x data set that was generated using the ω B97X/6–31G* level of theory. Toward that end, where needed, we recalculated the energy and forces, and performed geometry optimizations at a consistent ω B97X/6–31G* level of theory⁷³ (all using Gaussian 16⁷⁴) in order to train evolving versions of the QD π model. The QD π model is trained to be a QM/ -MLP; i.e., a non-electronic DPRc “correction” to the DFTB3/3OB⁷⁵ QM model potential energy similar to previous work.^{49,50}

2.1.1 Broad data sets: ANI-1xm and COMP5m—These data sets contain a diverse range of bio and drug-like molecules at equilibrium and non-equilibrium conformations, and contain structures, potential energies and forces. Generally, previous chemical space data sets^{76–78} are usually derived from the GDB databases^{64,65,79} that contains billions of SMILES strings⁸⁰ for organic small molecules. Herein we use modified versions of the public ANI-1x³⁷ and COMP6³⁶ databases as follows.

ANI-1xm.: The ANI-1x data set is an open-source chemical space data set proposed by Smith *et al.*³⁷ that includes ω B97X/6–31G* energies and forces generated by diverse normal mode sampling (DNMS). We examined the ANI-1x data set and observed that the DNMS procedure would in some cases generate free radicals by breaking covalent bonds (which were still computed with a singlet spin state in the reference data set), and this led to problems in the QD π training (an example is provided in Section 1 of the Supporting Information). Thus, we curated a subset of the ANI-1x data to create a modified data set we refer to as ANI-1xm by analyzing and removing such predicted free radicals in addition to a few other select outlier molecules through the procedure described below.

As an example, the DFTB3/3OB base QM model is known to have rare anomalous outlier energies for some inorganic molecules such as cyanogen⁷⁵ that we did not consider as highly relevant for drug discovery. In other cases, it has been reported that some reference values

in ANI-1x data set are not reliable.⁵¹ We thus used the following outlier detection criteria to remove 50 points that satisfy the condition:

$$\frac{\|E_k - \bar{E}\|}{\sigma(E)} \geq 8, \quad (1)$$

where E_k is the energy difference between ω B97X/6–31G* and DFTB3, and \bar{E} and $\sigma(E)$ are the mean and standard deviation of the energy differences for all molecules with the same chemical formula. The threshold is taken from the TorchANI program.⁸¹ After curating the ANI-1xm data set in this way, we obtained a total of 2,641,429 points (a 46.7% reduction from the original ANI-1x data set).

COMP5m.: The COmprehensive Machine-learning Potential (COMP6)³⁶ benchmark is a chemical space data set that was built from six separate data sets: 1) the original S66×8 benchmark,^{40,41} 2) ANI-MD,³⁶ 3) GDB7to9,⁶⁴ 4) GDB10to13,⁶⁵ 4) Tripeptide,³⁶ and 5) DrugBank⁶⁶ data sets. We begin by extracting the S66×8 data set, which we will analyze separately, and we refer to the truncated data set as “COMP5”. The COMP5 data set, like ANI-1x, used the DNMS procedure which cleaved covalent bonds in some instances; therefore, we applied the same outlier detection procedure described above to arrive at a modified COMP5m data set containing 64,667 data points (a 35.9% reduction from the original COMP5 data set).

2.1.2 Intermolecular data sets.—Intermolecular data sets contain dimers at multiple separations. Each dimer was geometry optimized at the reference theoretical level and the intermonomeric distances without altering the internal geometries.⁴⁰ The ω B97X/6–31G* and DFTB3 energies and forces were evaluated at the reference geometries.

For the S66×8 data set and the HB375×10 data set, we compute the relative energies as follows:

$$\Delta E = E - E_{\min} \quad (2)$$

where E_{\min} is the minimum energy of the dimer (at the most favorable intermolecular separation). The QD π model was trained using all available dimer separations. Because most of the interaction energies are quite small, we report only the E of the most separated dimer configuration. The following intermolecular data sets are used in the current work.

S66×8.: The S66×8 data set⁴¹ is a nonbonded interaction data set containing 66 noncovalent pairs at 8 separations. The ω B97X/6–31G* energy and forces are directly taken from the COMP6 benchmark.³⁶ The 8 relative energies are computed for each of the 66 dimers.

HB375×10.: The HB375×10 data set³⁹ is an S66×8-like nonbonded interaction data set containing 375 hydrogen bonding pairs. We use geometry provided by the data set to compute the ω B97X/6–31G* energy and forces.

AEGIS:BP: The AEGIS:BP data set is a subset of 10 hydrogen-bonded nucleic acid base pairs (BPs) within the artificially expanded genetic information system (AEGIS)^{60,82}

database. The entire list of BPs is given in Figure S2 of the Supporting Information. The structures were generated by Open Babel and optimized at the ω B97X/6–31G* level.

2.1.3 Tautomerization data sets.—The tautomerization data sets described below are used to evaluate the relative energy of the tautomeric configurations. The relative energy between A and B is the difference between their total energies (E_A and E_B).

$$\Delta E = E_A - E_B \quad (3)$$

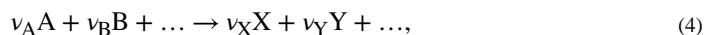
Each tautomer is optimized from the initial geometry at the ω B97X/6–31G* level. The relative energies and force corrections between ω B97X/6–31G* and DFTB3 at the reference geometries are used in the neural network training. When tabulating the results to compare different methods, we report the mean absolute errors (maE) and root mean squared errors (rmsE) of the relative energies; these values evaluate the energies upon geometry optimizing the structures with each method.

Tautobase.: The Tautobase data set⁶¹ is a broad tautomer data set. A subset of the Tautobase data set was constructed by Wieder et al.⁶² that includes 354 tautomer pairs with C, H, O, and N elements. Each pair involves the relocation/bonding of a hydrogen atom. The initial geometry is generated by Open Babel⁸³ and optimized at the ω B97X/6–31G* level.

TAUT15.: TAUT15 is a tautomer data set from GMTK55 database⁴² containing 15 relative energies. The initial geometry is provided by the data set and further optimized at the ω B97X/6–31G* level.

AEGIS:TAUT.: The AEGIS:TAUT data set is a subset of 25 tautomeric (TAUT) equilibria extracted from the AEGIS⁶⁰ database. The entire list of TAUT is given in Figure S3 of the Supporting Information. The initial geometry is provided by the data set and further optimized at the ω B97X/6–31G* level.

2.1.4 Protonation energy data sets.—These data sets are intended to be reflective of titratable sites in biological and ligand/drug-like molecules. We depart from public data sets that provide SMILES strings of absolute deprotonation energies $AH \rightarrow A^- + H^+$; however, we train and test the methods against relative protonation energies $AH + B^- \rightarrow A^- + BH$. Each reactant or product is optimized. For a generic chemical reaction



the relative energy for a reaction (denoted E_{rxn}) is defined as the potential energy difference between total reactants and total products:

$$\Delta E_{\text{rxn}} = \sum_{\text{products,p}} v_p E - \sum_{\text{reactants,r}} v_r E. \quad (5)$$

where v_p and v_r are the stoichiometric coefficients of each product p and reactant r.

Amino acids model compounds (AAMC): The AAMC data set contains 21 O-H and N-H bond-containing molecules (OHNH) with deprotonation energies, including amino acids model compounds from Ref. 43. The entire list of compounds is given in Section 5.1 of the Supporting Information. The initial geometries were generated by Open Babel⁸³ and optimized at the ω B97X/6–31G* level.

Nucleic acid model compounds (NAMC): The NAMC data set contains 53 deprotonation energies of nucleic acid (DNA and RNA bases) model compounds introduced in Ref. 43. The entire list of compounds is given in Section 5.2 of the Supporting Information. The initial geometry is generated by Open Babel⁸³ and optimized at the ω B97X/6–31G* level.

PA26: The PA26 data set is a subset of the GMTKN55 database⁴² containing 26 adiabatic proton affinities. The initial geometry is provided by the data set and optimized at the ω B97X/6–31G* level.

RegioSQM20: We selected a subset of the RegioSQM20⁶³ database containing C, H, O, and N elements. The subset was randomly divided into the training and test sets. Some outliers were removed using the procedure described above for the ANI-1x and COMP5m data sets. Then, there are 544 and 25 deprotonation energies in the training and test sets, respectively. The entire list of compounds is given in Section 5.3 of the Supporting Information. The initial geometry is generated by Open Babel⁸³ (or RDKit⁸⁴ for some compounds to get better structures) and optimized at the ω B97X/6–31G* level.

2.2 QD π (v1.0)

In this work, we develop a general QD π model as a -MLP correction⁵⁴ to DFTB3/3OB. The correction is parametrized to reproduce target energies and forces for closed-shell bio and drug-like organic molecules and ions composed of C, H, O, and N elements. The QD π energy is the sum of DFTB3 and neural network potential (NNP) model energies:

$$E_{\text{QD}\pi} = E_{\text{DFTB3}} + E_{\text{NNP}}. \quad (6)$$

where E_{DFTB3} is the DFTB3/3OB energy, and E_{NNP} is the -MLP correction using the Deep Potential — Smooth Edition (DeepPot-SE) functional form.²⁸ DFTB3 was chosen as the base model because it is robust, internally consistent, and has been reported³⁶ to have better overall accuracy for the ANI-1x data set compared to PM6.

2.2.1 DeepPot-SE—The QD π model parametrizes a Deep Potential (DP) using the DeepPot-SE descriptor²⁸ used as a -MLP correction. The functional form of the DP and DeepPot-SE descriptor has been previously described⁸⁵ and additional details can be found in Section 2 of the Supporting Information. DeepPot-SE is a popular descriptor implemented in the DeepPMD-kit package^{71,72} which has seen use in over 100 works⁸⁶ since its proposal in 2018. It also serves as the foundation for the DPRc⁴⁹ -MLP. The DPRc correction includes corrections for QM/MM interactions. Although the present work does not involve QM/MM interactions, the common framework between the DPRc and DP potentials affords the opportunity to extend the QD π model to QM/MM applications using the DPRc potential.

A recent work¹⁸ has compared the theories of different NNPs, including ANI-1 and DeepPot-SE. NNPs that use atomic-centered symmetry functions (ACSFs),¹⁹ such as ANI-1, have fixed descriptors that must be determined before training. In contrast, the descriptors used in the SchNet²⁴ and DeepPot-SE NNPs are trained to improve accuracy. It was found that NNPs with trainable descriptors require more computational effort to train because the descriptor includes additional parameters (and thus the additional parameter gradients must be evaluated during training).⁸⁷ To address this issue, a model compression scheme has been introduced that can freeze and compress the DeepPot-SE descriptor to improve performance (either during or after training).⁸⁷ In this work, we apply this model compression scheme in the latter part of training (see below).

2.2.2 Relative Energy loss—In this work, we introduce a new component to the loss function to train relative energies. It is common for ML training algorithms to update the neural network parameters using a subset (a “batch”) of the available training data. A “loss function”, L , is evaluated using the data contained within the batch, the neural network parameters are updated, and a new batch is created for the next optimization step by randomly selecting another subset of data.⁸⁸ In the past, a batch consisted of molecules whose total energies and forces are to be trained. In which case, the loss function consisted of two components: errors arising from the total energies L_E and forces L_f .

$$L = p_E L_E + p_f L_f \quad (7)$$

In the present work, we allow relative energies to be included within a batch.

$$L = p_E L_E + p_{\Delta E} L_{\Delta E} + p_f L_f \quad (8)$$

where L_E is the relative energy loss.

$$L_{\Delta E} = \frac{1}{\mathcal{B}_{\Delta E}} \sum_{k=1}^{\mathcal{B}_{\Delta E}} \frac{1}{N_k} (\Delta E_k - \Delta E_k^*)^2 \quad (9)$$

where $\mathcal{B}_{\Delta E}$ is the number of relative energies within the batch. N_k is the total number of atoms for system k (the sum of all product and reactant atoms in the case of reaction energy). E_k and ΔE_k^* are the model and reference relative energies, respectively. L_E and L_f are defined in the same way.

$$L_E = \frac{1}{\mathcal{B}_E} \sum_{k=1}^{\mathcal{B}_E} \frac{1}{N_k} (E_k - E_k^*)^2 \quad (10)$$

$$L_f = \frac{1}{\mathcal{B}_f} \sum_{k=1}^{\mathcal{B}_f} \frac{1}{3N_k} (f_k - f_k^*)^2 \quad (11)$$

where \mathcal{B}_E and \mathcal{B}_f are the number of relative energies and forces within the batch, respectively. E_k and E_k^* are the model and reference energy components, and f_k and f_k^* are the model

and reference force components, respectively. p_E , $p_{\Delta E}$, and p_f are weights assigned to energy, relative energy, and force contributions to the loss function. The weights are linearly updated with the learning rate, α :

$$p_E(t) = p_E^0 \frac{\alpha(t)}{\alpha(0)} + p_E^\infty \left(1 - \frac{\alpha(t)}{\alpha(0)}\right) \quad (12)$$

$$p_{\Delta E}(t) = p_{\Delta E}^0 \frac{\alpha(t)}{\alpha(0)} + p_{\Delta E}^\infty \left(1 - \frac{\alpha(t)}{\alpha(0)}\right) \quad (13)$$

$$p_f(t) = p_f^0 \frac{\alpha(t)}{\alpha(0)} + p_f^\infty \left(1 - \frac{\alpha(t)}{\alpha(0)}\right) \quad (14)$$

The learning rate decays exponentially with the training step, t .

$$\alpha(t) = \alpha_0 \lambda^{t/\tau} \quad (15)$$

where α_0 is the initial learning rate, λ is the decay rate, and τ is the decay steps. If energy, relative energy, or force is not available in a batch, p_E , $p_{\Delta E}$, or p_f will be set to zero to disable the corresponding loss contribution. In this work, we set $p_E^0 = 2$, $p_E^\infty = 20$, $p_{\Delta E}^0 = 2$, $p_{\Delta E}^\infty = 20$, $p_f^0 = 100$, $p_f^\infty = 0.1$, $\alpha_0 = 0.0001$, $\lambda = 0.99$, and $\tau = 400$. It's worth mentioning that direct training to the relative energies typically will not improve the accuracy of the absolute atomic energies (the energy of an atom in a vacuum). We have implemented the relative energy loss contributions into the DeePMD-kit package.^{71,72}

2.2.3 Training process—The QD π model was trained using the DeePMD-kit software package. As shown in Table 2, we performed 6 training iterations with different data sets and training properties. In the first two iterations, the model was trained to reproduce the total energies and forces of the molecules contained within the relative protonation energy data set. After the first iteration, the DP Compress⁸⁷ algorithm was applied to freeze the model descriptor \mathcal{D} to improve performance. All subsequent iterations restart the training from the previous iteration. Starting from iteration 3, the loss function was changed to train against relative energies rather than molecular total energies. All of the tables in the main text show results only for the QD π -v1.0, but the Supporting Information contains extended tables that have results for each version to compare.

2.3 Energy/force calculation and geometry optimizations

This section describes the various potentials compared, in addition to the basic methods used for performing geometry optimizations. Additional details for how relaxed 2D potential energy surface scans are provided in section 3 of the Supporting Information.

ω B97X/6–31G*.—We used Gaussian 16⁷⁴ to evaluate ω B97X/6–31G* energies and perform geometry optimizations.⁷³

Semiempirical methods.—The AMBER 20⁸⁹ SQM module⁹⁰ was used to perform geometry optimizations and evaluate the energies and forces of several semiempirical models, including DFTB3^{91,92} (3OB parameters⁷⁵), MNDO/d,⁹³ AM1,⁹⁴ and PM6.⁹⁵ The DFTB+⁹⁶ package was used to validate DFTB3 results from AMBER.

The DFTB+⁹⁶ package was used to calculate GFN1-xTB⁹⁷ and GFN2-xTB⁹⁸ energies and forces. For these models, the ASE package⁹⁹ was used to optimize the geometries with the Limited-memory Broyden-Fletcher-Goldfarb-Shanno (LBFGS) algorithm.¹⁰⁰

QD π models.—The QD π energy is the sum of the DFTB3 and the DP contributions. The DP contribution was directly evaluated within the DeePMD-kit program.⁷¹ The ASE package⁹⁹ was used to optimize the QD π structures.

ANI models.—The TorchANI program⁸¹ was used to provide energies and forces of the ANI-1x³⁶ and ANI-2x⁴⁵ models. Each ANI model consists of 8 independent parameter sets. We only use the first model (index 0) for benchmarking. It has been suggested that using an average of multiple models will improve the accuracy⁴⁴ but at additional computational cost.⁸¹ We performed the geometry optimizations with the ASE package.⁹⁹

3 Results and Discussion

For consistency, all the reference data used to train and test the QD π model was performed at the ω B97X/6–31G* level, as in the ANI databases^{30,36,44,45} used to train the ANI-1x and ANI-2x models. In making the comparison with the target reference data, we report various error metrics such as mean absolute and root-mean-square errors (maEs and rmsEs, respectively). Other models compared in this work have been trained to different reference data (levels of theory and data sets). Hence, it should not be concluded that deviation of these other models from the reference data used in this work implies they are necessarily less accurate in the theoretical electronic structure limit (which cannot be practically obtained for any of the data considered here). Thus, a comparison of results from other models is not meant to be critical, but rather provide a broader context with respect to variation from a well-defined reference. In our view, the real litmus test for a drug discovery force field is the accuracy of binding free energy predictions from rigorous and precise AFE simulations,^{5,101} which is beyond the scope of this work. The current work focuses on developing a robust and internally consistent internal energy model from drug-like molecules that provide a foundation from which to extend to accurate intermolecular interactions as a -MLP to the QM/MM potential.

3.1 Performance for internal (intramolecular) potential energy

The majority of data used to train the QD π model was from the ANI-1xm data set using energies and forces, and validated against the COMP5m data set. The ANI-1x model was trained to the energies only, whereas the ANI-2x was trained to both energy and forces. Comparison of the force component errors for these models are illustrated in Figure 1, and error results for QD π and all of the established comparison models are summarized in Table 3. Not surprisingly, the ANI-2x model performs best among the established models compared, having mean absolute errors (maE) of 1.07 kcal/mol and 2.11 kcal/(mol·Å) for

the energy and forces, respectively. Moreover, these errors are transferable to COMP5m, for which neither the ANI nor QD π models were trained (maE of 1.67 kcal/mol and 1.86 kcal/(mol·Å) for energy and force, respectively). The other models ranged in maE in forces of 4.69–15.14 kcal/(mol·Å) for the ANI-1xm and 3.68–12.13 kcal/(mol·Å) for the COMP5m (energies were not compared for the semiempirical models as the zero of total energy uses a different reference than for ω B97X/6–31G*). The worst models overall are the NDDO-based semiempirical models (MNDO/d, AM1 and PM6), which require additional fixes such as orthogonalization corrections¹⁰² or other empirical terms¹⁰³ in order to accurately reproduce relative conformational energies. The tight-binding models that explicitly account for orthogonalization through the inclusion of an overlap matrix in the generalized eigenvalue problem perform generally better, with the GFN models slightly out-performing DFTB3/3OB. The QD π performs exceptionally well on both the ANI-1xm training and COMP5m testing data sets, having maE of 0.83 and 1.48 kcal/mol for the energy, respectively, and 1.16 and 1.14 kcal/(mol·Å) for the forces, respectively.

In order to illustrate the degree to which the QD π model can reproduce conformational energy landscapes, we examined relaxed 2D torsion profiles for three systems: the alanine dipeptide, and the drug molecules ibuprofen and ketorolac illustrated in Fig. 2. Figure 3 compares the potential energy surface for 2D torsion scans at the ω B97X/6–31G*, QD π , ANI-2x, GFN2-xTB and DFTB3 levels. All of the models qualitatively predict the correct trends. Overall, QD π and ANI-2x are quite similar and have the closest agreement with ω B97X/6–31G*. While none of the models is able to reproduce conformational energy barriers below 1 kcal/mol in all cases (see Table S4 of the Supporting Information), the DFTB3 model errors are the largest and most systematic in their under-estimation of the barriers between minima.

3.2 Performance for intermolecular interactions

Despite the focus being on training a QD π internal energy model, we felt it important to include training and testing data to intermolecular interaction DBs (S66 \times 8^{36,40,41} and HB375 \times 10³⁹) as some large, drug-like molecules can form similar interactions (e.g., intramolecular hydrogen bonds). Figure 4 and Table 4 compares intermolecular interactions (E values) for S66 \times 8. Overall QD π has the smallest maE (0.30 kcal/mol) relative to the other models that ranged in maE from 0.59 kcal/mol (ANI-2x) to 3.59 kcal/mol for DFTB3. A more detailed breakdown of the errors into hydrogen bonding (HB), π stacking, London dispersion (LD) and mixed interactions are provided in Table 4. The QD π has maE values that range from 0.21–0.41 kcal/mol (the largest is the LD subset), whereas ANI-2x ranges from 0.40–0.91 kcal/mol (the largest being for the HB subset). Closer examination of the more extensive HB735 data set indicates that QD π has the lowest maE (0.44 kcal/mol), whereas ANI-2x, GFN2-xTB, and DFTB3 have maE values of 1.40, 0.85, and 1.17 kcal/mol, respectively. In the case of ANI-2x, the distribution of errors has a considerably larger variance as indicated by the rmsE value of 3.84 which is over 2.5 times larger than the maE.

3.3 Performance for tautomers

Figure 5 compares the tautomer energies for the Tautobase⁶¹ and TAUT15⁴² and AEGIS:TAUT data sets. The QD π model is the only model that has quantitative (less than 1

kcal/mol) accuracy for tautomer energies (maE values that range from 0.70 to 0.82 kcal/mol for data subsets listed in Table 5 and 0.79 kcal/mol overall). Of the other models, ANI-2x performs the best but still has a maE roughly twice as large (maE 1.63 kcal/mol, with data subset values that range from 1.00–1.76 kcal/mol). The GFN2-xTB and DFTB3 models have maE values exceeding 4.5 kcal/mol overall with the largest contributions coming from the Tautobase data subset (maE values for GFN2-xTB and DFTB3 both approximately 5.5 kcal/mol). Both QD π and ANI-2x have a high linear correlation with ω B97X/6–31G* (0.99 and 0.95, respectively), whereas GFN-xTB and DFTB3 have moderate correlation (0.69 and 0.55). Hence, the QD π model is the only model that consistently provides tautomer energies with errors below 1 kcal/mol relative to the reference.

3.4 Performance for relative protonation states

Figure 6 compares the differences in protonation energies for a series of amino acid and nucleic acid model compounds,⁴³ as well as molecules in the PA26⁴² data set. Changes in charge state are notoriously challenging for minimal valence basis models,¹⁰⁴ and are particularly problematic for the ANI-2x model that cannot distinguish molecules from ions and thus breaks down. Both the GFN2-xTB and DFTB3 models have a high correlation (0.99) with the ω B97X/6–31G*, owing mainly to the large range of values that are clustered into two sets (0–75 and 125–225 kcal/mol), and have large maE values are 7.78 and 10.24 kcal/mol, respectively. The QD π model performs exceptionally well with maE of 0.17 kcal/mol and almost perfect correlation.

A closer examination of the data subsets is listed in Table 6. The maE values between AAMC, NAMC and PA26 data subsets have ranges 0.09–0.39 (QD π), 5.77–8.45 (GFN2-xTB), 8.63–12.54 (DFTB3) and 23.8–70.5 kcal/mol (ANI-2x). Examination of errors from the RegioSQM* data set that was not used in training reveals more similar maE values: 2.53 (QD π), 4.12 (GFN2-xTB), 4.59 (DFTB3) and 13.6 kcal/mol (ANI-2x). The reason for this may be due in part to the design of the database to predict the regioselectivities of electrophilic aromatic substitution reactions from the calculation of proton affinities, but this is not fully clear.

3.5 Example application: acid/base reactions important in enzyme catalyzed RNA cleavage

Although the current QD π model has been designed with the intent ultimately for applications to drug discovery, it is nonetheless instructive to consider well-studied examples where protonation/deprotonation events are of biological significance. One such example presents itself in the chemistry of RNA strand cleavage^{105–107} that is catalyzed by protein ribonucleases^{108–110} as well as small self-cleaving ribozymes^{111,112} and several artificially engineered DNAzymes.^{113–116} In this reaction, the 2'-OH group of an RNA nucleotide is activated by deprotonation by a general base. The resulting activated nucleophile makes an in-line attack to the adjacent phosphorus of the scissile phosphate, and the reaction proceeds through a pentacovalent transition state followed by departure of the 5'-O leaving group that is facilitated by donation of a proton from a general acid. In the case of protein ribonucleases, the general base and acid is thought to be carried out by active site histidine residues,^{108–110} although early work speculated that a functionally important lysine residue

might also be a plausible candidate.¹¹⁷ In the case of small self-cleaving ribozymes and DNAzymes, general acid-base catalysis is carried out by nucleobases, metal ions or in some cases assisted by the 2'-OH of the ribose sugar moiety. Considering the nucleobase candidates, the general base is often an active site guanine (at the N1 position), whereas the general acid can be either a cytosine (at the N3 position), or else an adenine (either N1 or N3 positions).

We hence consider the energetics of reactions that involve the relative protonation/deprotonation of the general acid and base models with respect to the 2'OH nucleophile (a secondary alcohol modeled as isopropanol, iPrOH) and the 5'OH leaving group (a primary alcohol modeled as ethanol, EtOH). These are listed in Table 7. The model reactions where iPrOH is deprotonated to form iPrO⁻ is meant to represent a model system for general base activation of the 2'-OH nucleophile, whereas the reactions where EtO⁻ is protonated to form EtOH is meant to represent a model system for general acid stabilization of the 5'-O⁻ leaving group. The values shown in Table 7 show the relative energetics of the non-interacting molecular and ionic reaction species. In the gas phase, the formation of neutral molecules from non-interacting ions from neutral molecules is highly exothermic, although for interacting systems in an enzyme environment the differences are expected to be much smaller. Nonetheless, the inherent energetics associated with the relative protonation/deprotonation events still is a major factor that regulates reactivity. For the protein enzyme model reactions (top block, Table 7), the QD π performs exceptionally well with errors all less than 0.5 kcal/mol, whereas the DFTB3 and GFN2-xTB errors range from -11.33 to 11.72 and -7.02 to 9.66 kcal/mol, respectively. For the nucleobase reaction models (middle block, Table 7), QD π has larger errors that range from -1.11 to 1.25 kcal/mol, whereas the corresponding ranges are -8.63 to 16.00 and -2.69 to 11.40 for DFTB3 and GFN2-xTB, respectively. The relative errors between nucleobase general acids and base (bottom block, Table 7), again the QD π is in very close agreement with the reference values (maximum error -0.27 kcal/mol), where the DFTB3 and GFN2-xTB models have considerably larger errors (maximum errors 6.99 and 6.07 kcal/mol, respectively). Overall the average errors for the QD π is below 0.5 kcal/mol. As mentioned earlier, the ANI-2x model breaks down for system that have varying charge, with average errors of over 100 kcal/mol. These results provide an example in a biological context that emphasizes the importance of modeling relative protonation/deprotonation events with quantitative accuracy. For drug discovery, these will be especially important, as it has been estimated that over 95% of drug molecules have ionizable sites, many of which may potentially change upon binding to a biological targets.

3.6 Current limitations and future directions

A key aspect of this project was to create a first-generation potential energy model trained against broad data all computed at the same level of theory (and where possible, even using the same software package). At the time this project was initiated, the largest such data set was the ANI-1x DB³⁷ at the ω B97X/6-31G* level that only contained compounds with elements H, C, N and O. This chemical space is incomplete, as many drug molecules contain phosphorus, sulfur and halogen atoms, and some contain metal ions.¹¹⁸⁻¹²⁰ The ANI-2x model was extended to include S, F and Cl,⁴⁵ but the full data set, including the

important reference energies and forces at the ω B97X/6–31G* level, to our knowledge, has not yet become publicly available. Currently, there are a number of recent data sets that include compounds that contain phosphorus, sulfur and halogens at various levels of theory^{121–125} as well as metal ions.¹¹⁸ Among them, only the SPICE data set¹²⁴ includes forces at the ω B97M-D3BJ/def2-TZVPPD level and currently includes over 420K phosphorus, 520K sulfur, 750K halogen and 8K metal-containing structures.

Hence, current limitations of the QD π -v1.0 model include restricted chemical space (molecules containing H, C, N and O) and the ω B97X/6–31G* reference level of theory. The ω B97X/6–31G* reference level, like the DFTB3/3OB base QM model, lacks dispersion corrections and also does not include counterpoise corrections and complete basis set extrapolations that are important for intermolecular interactions. Further, this reference level of theory is not ideal for all molecular properties, including ionization energies and in some cases proton affinities of anions that may be sensitive to inclusion of diffuse basis functions. So while it is important to start with an established and consistent reference level of theory and chemical scope, ultimately as higher-level data sets become more complete and made publicly available, QD π and other machine learning potentials can continue to evolve.

The next step of future work will involve developing an intermolecular QM/MM interaction potential as a new range-corrected deep-learning potential.^{49,50} The full (internal and intermolecular interaction) QD π model is designed to be a correction to the QM/MM potential energy using DFTB3/3OB and the latest AMBER FF19SB for proteins,¹²⁶ OL3/OL15 for nucleic acids,^{127–129} OPC model for water^{130,131} and 12-6-4 ion models.^{132–134} Once the intermolecular interaction component of the QD π model has been developed and validated in alchemical free energy simulations,⁵ next steps will be to extend the chemical space of drug molecules to include P, S, F and Cl atoms. With GPU acceleration, the QD π is typically less than 10% overhead relative to a traditional QM/MM energy/force evaluation with DFTB3/3OB.⁴⁹ With this design, the QD π could in principle also be used to modify the internal energy of protein residues and/or solvent molecules in contact with the drug, but this would incur greater cost as the size of the QM region grows larger. Should treatment of these surrounding residues with a neural network correction potential be deemed necessary, an alternative strategy would be to extend the model such that it can directly correct the MM potential rather than the QM/MM potential.

4 Conclusion

We report a QD π -v1.0 for modeling the internal energy of drug molecules. The development of this model required: 1) collection and curation of several existing molecular databases of structures, energies and forces; 2) generation of new data sets at the ω B97X/6–31G* level that fill needed gaps in training and/or testing data; 3) development of new tools within DeePMD-kit that enable more general flexible forms of the loss function used in training of the neural networks; and 4) creation of computational infrastructure for consistent comparison of a wide array of existing potential energy models. The QD π model has the advantage that it is able to properly treat electrostatic interactions and handle changes in charge/protonation states. The QD π model is demonstrated to be accurate for a wide range of intra- and intermolecular interactions (despite its intended use as an internal energy

model), and has shown to perform exceptionally well for relative protonation/deprotonation energies and tautomers. Comparison with several other approximate semiempirical and machine learning potentials (ANI-1x, ANI-2x, DFTB3, MNDO/d, AM1, PM6, GFN1-xTB and GFN2-xTB) indicates QD π agrees much more closely with training and testing data at the reference ω B97X/6-31G* level. An example application to model reactions involved in RNA strand cleavage catalyzed by protein and nucleic acid enzymes further illustrates the QD π accuracy in a biological context. This makes QD π highly attractive as a potential force field model for drug discovery.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

The authors thank Dr. Han Wang for the code review of DeePMD-kit. The authors are grateful for the financial support provided by the National Institutes of Health (No. GM107485 to DMY). Computational resources were provided by the Office of Advanced Research Computing (OARC) at Rutgers, The State University of New Jersey; the Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant ACI-1548562 (supercomputer Expanse at SDSC through allocation CHE190067); and by the Texas Advanced Computing Center (TACC) at the University of Texas at Austin (supercomputer Longhorn through allocation CHE20002).

Data Availability

QD π -v1.0 is openly available in our GitLab repository at <https://gitlab.com/RutgersLBSR/qdpi>.

References

- (1). Jorgensen WL Efficient drug lead discovery and optimization. *Acc. Chem. Res* 2009, 42, 724–733. [PubMed: 19317443]
- (2). Abel R; Wang L; Harder ED; Berne BJ; Friesner RA Advancing Drug Discovery through Enhanced Free Energy Calculations. *Acc. Chem. Res* 2017, 50, 1625–1632. [PubMed: 28677954]
- (3). Cournia Z; Allen B; Sherman W Relative Binding Free Energy Calculations in Drug Discovery: Recent Advances and Practical Considerations. *J. Chem. Inf. Model* 2017, 57, 2911–2937. [PubMed: 29243483]
- (4). Cournia Z; Allen BK; Beuming T; Pearlman DA; Radak BK; Sherman W Rigorous Free Energy Simulations in Virtual Screening. *J. Chem. Inf. Model* 2020, 60, 4153–4169. [PubMed: 32539386]
- (5). Lee T-S; Allen BK; Giese TJ; Guo Z; Li P; Lin C; T. D. M. Jr.; Pearlman DA; Radak BK; Tao Y; Tsai H-C; Xu H; Sherman W; York DM Alchemical Binding Free Energy Calculations in AMBER20: Advances and Best Practices for Drug Discovery. *J. Chem. Inf. Model* 2020, 60, 5595–5623. [PubMed: 32936637]
- (6). Cole DJ; Horton JT; Nelson L; Kurdekar V The future of force fields in computer-aided drug design. *Future Med. Chem* 2019, 11, 2359–2363. [PubMed: 31544529]
- (7). Vanommeslaeghe K; MacKerell A Jr. {CHARMM} additive and polarizable force fields for biophysics and computer-aided drug design. *Biochim. Biophys. Acta* 2015, 1850, 861–871, Recent developments of molecular dynamics. [PubMed: 25149274]
- (8). Martin YC Let's not forget tautomers. *J. Comput.-Aided Mol. Des* 2009, 23, 693–704. [PubMed: 19842045]

- (9). Milletti F; Storchi L; Sforna G; Cross S; Cruciani G Tautomer Enumeration and Stability Prediction for Virtual Screening on Large Chemical Databases. *J. Chem. Inf. Model* 2009, 49, 68–75. [PubMed: 19123923]
- (10). Manallack DT The pK(a) Distribution of Drugs: Application to Drug Discovery. *Perspect Medicin Chem.* 2007, 1, 25–38. [PubMed: 19812734]
- (11). Rathore RS; Sumakanth M; Siva Reddy M; Reddanna P; Rao AA; Erion MD; Reddy MR Advances in Binding Free Energies Calculations: QM/MM-Based Free Energy Perturbation Method for Drug Design. *Curr. Pharm. Des* 2013, 19, 4674–4686. [PubMed: 23260025]
- (12). Kříž K; ezá J Benchmarking of Semiempirical Quantum-Mechanical Methods on Systems Relevant to Computer-Aided Drug Design. *J. Chem. Inf. Model* 2020, 60, 1453–1460. [PubMed: 32062970]
- (13). Behler J Perspective: Machine learning potentials for atomistic simulations. *J. Chem. Phys* 2016, 145, 170901. [PubMed: 27825224]
- (14). Butler KT; Davies DW; Cartwright H; Isayev O; Walsh A Machine learning for molecular and materials science. *Nature* 2018, 559, 547–555. [PubMed: 30046072]
- (15). Noé F; Tkatchenko A; Müller K-R; Clementi C Machine Learning for Molecular Simulation. *Annu. Rev. Phys. Chem* 2020, 71, 361–390. [PubMed: 32092281]
- (16). Pinheiro M Jr; Ge F; Ferré N; Dral PO; Barbatti M Choosing the right molecular machine learning potential. *Chem. Sci* 2021, 12, 14396–14413. [PubMed: 34880991]
- (17). Manzhos S; Carrington T Jr Neural Network Potential Energy Surfaces for Small Molecules and Reactions. *Chem. Rev* 2021, 121, 10187–10217. [PubMed: 33021368]
- (18). Zeng J; Cao L; Zhu T Neural network potentials. In *Quantum Chemistry in the Age of Machine Learning*; Dral PO, Ed.; Elsevier, 2022; Chapter 12, pp 279–294.
- (19). Behler J; Parrinello M Generalized Neural-Network Representation of High-Dimensional Potential-Energy Surfaces. *Phys. Rev. Lett* 2007, 98, 146401–146404. [PubMed: 17501293]
- (20). Bartók AP; Payne MC; Kondor R; Csányi G Gaussian Approximation Potentials: The Accuracy of Quantum Mechanics, without the Electrons. *Phys. Rev. Lett* 2010, 104, 136403. [PubMed: 20481899]
- (21). Behler J Atom-centered Symmetry Functions for Constructing High-dimensional Neural Network Potentials. *J. Chem. Phys* 2011, 134, 074106. [PubMed: 21341827]
- (22). Gastegger M; Schwiedrzik L; Bittermann M; Berzsenyi F; Marquetand P wACSF—Weighted atom-centered symmetry functions as descriptors in machine learning potentials. *J. Chem. Phys* 2018, 148, 241709. [PubMed: 29960372]
- (23). Chmiela S; Tkatchenko A; Sauceda HE; Poltavsky I; Schütt KT; Müller K-R Machine learning of accurate energy-conserving molecular force fields. *Sci. Adv* 2017, 3, 1603015.
- (24). Schütt KT; Arbabzadah F; Chmiela S; Müller KR; Tkatchenko A Quantum-chemical insights from deep tensor neural networks. *Nat. Commun* 2017, 8, 13890. [PubMed: 28067221]
- (25). Schütt K; Sauceda H; Kindermans P; Tkatchenko A; Müller K SchNet - A Deep Learning Architecture for Molecules and Materials. *J. Chem. Phys* 2018, 148, 241722. [PubMed: 29960322]
- (26). Chen X; Jørgensen MS; Li J; Hammer B Atomic Energies from a Convolutional Neural Network. *J. Chem. Theory Comput* 2018, 14, 3933–3942. [PubMed: 29812930]
- (27). Zhang L; Han J; Wang H; Car R; E W Deep potential molecular dynamics: a scalable model with the accuracy of quantum mechanics. *Phys. Rev. Lett* 2018, 120, 143001. [PubMed: 29694129]
- (28). Zhang L; Han J; Wang H; Saidi W; Car R; E W End-to-end Symmetry Preserving Inter-atomic Potential Energy Model for Finite and Extended Systems. In *Advances in Neural Information Processing Systems 31*; Bengio S, Wallach H, Larochelle H, Grauman K, Cesa-Bianchi N, Garnett R, Eds.; Curran Associates, Inc., 2018; pp 4436–4446.
- (29). Zhang Y; Hu C; Jiang B Embedded Atom Neural Network Potentials: Efficient and Accurate Machine Learning with a Physically Inspired Representation. *J. Phys. Chem. Lett* 2019, 10, 4962–4967. [PubMed: 31397157]
- (30). Smith JS; Isayev O; Roitberg AE ANI-1: an extensible neural network potential with DFT accuracy at force field computational cost. *Chem. Sci* 2017, 8, 3192–3203. [PubMed: 28507695]

- (31). Unke O; Meuwly M PhysNet: A Neural Network for Predicting Energies, Forces, Dipole Moments, and Partial Charges. *J. Chem. Theory Comput* 2019, 15, 3678–3693. [PubMed: 31042390]
- (32). Glick ZL; Metcalf DP; Koutsoukas A; Spronk SA; Cheney DL; Sherrill CD AP-Net: An atomic-pairwise neural network for smooth and transferable interaction potentials. *J. Chem. Phys* 2020, 153, 044112. [PubMed: 32752707]
- (33). Zubatiuk T; Isayev O Development of Multimodal Machine Learning Potentials: Toward a Physics-Aware Artificial Intelligence. *Acc. Chem. Res* 2021, 54, 1575–1585. [PubMed: 33715355]
- (34). Khajehpasha ER; Finkler JA; Kühne TD; Ghasemi SA CENT2: Improved charge equilibration via neural network technique. *Phys. Rev. B* 2022, 105, 144106.
- (35). Pan X; Yang J; Van R; Epifanovsky E; Ho J; Huang J; Pu J; Mei Y; Nam K; Shao Y Machine-Learning-Assisted Free Energy Simulation of Solution-Phase and Enzyme Reactions. *J. Chem. Theory Comput* 2021, 17, 5745–5758. [PubMed: 34468138]
- (36). Smith JS; Nebgen B; Lubbers N; Isayev O; Roitberg AE Less is more: Sampling chemical space with active learning. *J. Chem. Phys* 2018, 148, 241733–241743. [PubMed: 29960353]
- (37). Smith JS; Zubatyuk R; Nebgen B; Lubbers N; Barros K; Roitberg AE; Isayev O; Tretiak S The ANI-1ccx and ANI-1x data sets, coupled-cluster and density functional theory properties for molecules. *Sci Data* 2020, 7, 134. [PubMed: 32358545]
- (38). Rai BK; Sresht V; Yang Q; Unwalla R; Tu M; Mathiowetz AM; Bakken GA TorsionNet: A Deep Neural Network to Rapidly Predict Small-Molecule Torsional Energy Profiles with the Accuracy of Quantum Mechanics. *J. Chem. Inf. Model* 2022, 62, 785–800. [PubMed: 35119861]
- (39). ezá J Non-Covalent Interactions Atlas Benchmark Data Sets: Hydrogen Bonding. *J. Chem. Theory Comput* 2020, 16, 2355–2368. [PubMed: 32149503]
- (40). Goerigk L; Kruse H; Grimme S Benchmarking Density Functional Methods against the S66 and S66×8 Datasets for Non-Covalent Interactions. *Chem. Phys. Chem* 2011, 12, 3421–33. [PubMed: 22113958]
- (41). Brauer B; Kesharwani MK; Kozuch S; Martin JML The S66×8 benchmark for noncovalent interactions revisited: explicitly correlated ab initio methods and density functional theory. *Phys. Chem. Chem. Phys* 2016, 18, 20905–25. [PubMed: 26950084]
- (42). Goerigk L; Hansen A; Bauer C; Ehrlich S; Najibi A; Grimme S A look at the density functional theory zoo with the advanced GMTKN55 database for general main group thermochemistry, kinetics and noncovalent interactions. *Phys. Chem. Chem. Phys* 2017, 19, 32184–32215. [PubMed: 29110012]
- (43). Moser A; Range K; York DM Accurate Proton Affinity and Gas-Phase Basicity Values for Molecules Important in Biocatalysis. *J. Phys. Chem. B* 2010, 114, 13911–13921. [PubMed: 20942500]
- (44). Smith J; Nebgen B; Zubatyuk R; Lubbers N; Devereux C; Barros K; Tretiak S; Isayev O; Roitberg A Approaching Coupled Cluster Accuracy with a General-purpose Neural Network Potential Through Transfer Learning. *Nat. Commun* 2019, 10, 2903. [PubMed: 31263102]
- (45). Devereux C; Smith JS; Huddleston KK; Barros K; Zubatyuk R; Isayev O; Roitberg AE Extending the Applicability of the ANI Deep Learning Molecular Potential to Sulfur and Halogens. *J. Chem. Theory Comput* 2020, 16, 4192–4202. [PubMed: 32543858]
- (46). Martin YC Experimental and pK prediction aspects of tautomerism of drug-like molecules. *Drug Discov. Today. Technol* 2018, 27, 59–64. [PubMed: 30103864]
- (47). Yang Y; Yu H; York DM; Cui Q; Elstner M Extension of the Self-Consistent-Charge Density-Functional Tight-Binding Method: Third-Order Expansion of the Density Functional Theory Total Energy and Introduction of a Modified Effective Coulomb Interaction. *J. Phys. Chem. A* 2007, 111, 10861–10873. [PubMed: 17914769]
- (48). Gaus M; Lu X; Elstner M; Cui Q Parameterization of DFTB3/3OB for Sulfur and Phosphorus for Chemical and Biological Applications. *J. Chem. Theory Comput* 2014, 10, 1518–1537. [PubMed: 24803865]
- (49). Zeng J; Giese TJ; Ekesan ; York DM Development of Range-Corrected Deep Learning Potentials for Fast, Accurate Quantum Mechanical/Molecular Mechanical Simulations of

- Chemical Reactions in Solution. *J. Chem. Theory Comput* 2021, 17, 6993–7009. [PubMed: 34644071]
- (50). Giese TJ; Zeng J; Ekesan ; York DM Combined QM/MM, Machine Learning Path Integral Approach to Compute Free Energy Profiles and Kinetic Isotope Effects in RNA Cleavage Reactions. *J. Chem. Theory Comput* 2022, 18, 4304–4317. [PubMed: 35709391]
- (51). Zheng P; Zubatyuk R; Wu W; Isayev O; Dral PO Artificial intelligence-enhanced quantum chemical method with broad applicability. *Nat. Commun* 2021, 12, 7022. [PubMed: 34857738]
- (52). Gómez-Flores CL; Maag D; Kansari M; Vuong V-Q; Irle S; Gräter F; Kuba T; Elstner M Accurate Free Energies for Complex Condensed-Phase Reactions Using an Artificial Neural Network Corrected DFTB/MM Methodology. *J. Chem. Theory Comput* 2022, 18, 1213–1226. [PubMed: 34978438]
- (53). Böser J; Kuba T; Elstner M; Maag D Reduction pathway of glutaredoxin 1 investigated with QM/MM molecular dynamics using a neural network correction. *J. Chem. Phys* 2022, 157, 154104. [PubMed: 36272777]
- (54). Dral PO; Zubatyuk T; Xue B-X Learning from multiple quantum chemical methods: γ -learning, transfer learning, co-kriging, and beyond. In *Quantum Chemistry in the Age of Machine Learning*; Dral PO, Ed.; Elsevier, 2022; Chapter 21, pp 491–507.
- (55). Darden T; York D; Pedersen L Particle mesh Ewald: An $N \log(N)$ method for Ewald sums in large systems. *J. Chem. Phys* 1993, 98, 10089–10092.
- (56). Nam K; Gao J; York DM An efficient linear-scaling Ewald method for long-range electrostatic interactions in combined QM/MM calculations. *J. Chem. Theory Comput* 2005, 1, 2–13. [PubMed: 26641110]
- (57). Giese TJ; Panteva MT; Chen H; York DM Multipolar Ewald methods, 1: Theory, accuracy, and performance. *J. Chem. Theory Comput* 2015, 11, 436–450. [PubMed: 25691829]
- (58). Giese TJ; Panteva MT; Chen H; York DM Multipolar Ewald methods, 2: Applications using a quantum mechanical force field. *J. Chem. Theory Comput* 2015, 11, 451–461. [PubMed: 25691830]
- (59). Giese TJ; York DM Ambient-Potential Composite Ewald Method for ab Initio Quantum Mechanical/Molecular Mechanical Molecular Dynamics Simulation. *J. Chem. Theory Comput* 2016, 12, 2611–2632. [PubMed: 27171914]
- (60). Eberlein L; Beierlein FR; van Eikema Hommes NJR; Radadiya A; Heil J; Benner SA; Clark T; Kast SM; Richards NGJ Tautomeric Equilibria of Nucleobases in the Hachimoji Expanded Genetic Alphabet. *J. Chem. Theory Comput* 2020, 16, 2766–2777. [PubMed: 32125859]
- (61). Wahl O; Sander T Tautobase: An Open Tautomer Database. *J. Chem. Inf. Model* 2020, 60, 1085–1089. [PubMed: 31967818]
- (62). Wieder M; Fass J; Chodera JD Fitting quantum machine learning potentials to experimental free energy data: predicting tautomer ratios in solution. *Chem. Sci* 2021, 12, 11364–11381. [PubMed: 34567495]
- (63). Ree N; Göller AH; Jensen JH RegioSQM20: improved prediction of the regio-selectivity of electrophilic aromatic substitutions. *J. Cheminform* 2021, 13, 10. [PubMed: 33579374]
- (64). Fink T; Bruggesser H; Reymond J-L Virtual Exploration of the Small-Molecule Chemical Universe below 160 Daltons. *Angew. Chem. Int. Ed* 2005, 44, 1504–8.
- (65). Blum LC; Reymond J-L 970 Million Druglike Small Molecules for Virtual Screening in the Chemical Universe Database GDB-13. *J. Am. Chem. Soc* 2009, 131, 8732–8733. [PubMed: 19505099]
- (66). Law V; Knox C; Djoumbou Y; Jewison T; Guo AC; Liu Y; Maciejewski A; Arndt D; Wilson M; Neveu V; Tang A; Gabriel G; Ly C; Adamjee S; Dame ZT; Han B; Zhou Y; Wishart DS DrugBank 4.0: shedding new light on drug metabolism. *Nucleic Acids Res.* 2014, 42, 1091–1097.
- (67). Giese TJ; York DM Development of a Robust Indirect Approach for MM \rightarrow QM Free Energy Calculations That Combines Force-Matched Reference Potential and Bennett's Acceptance Ratio Methods. *J. Chem. Theory Comput* 2019, 15, 5543–5562. [PubMed: 31507179]
- (68). König G; Mei Y; Pickard FC 4th; Simmonett AC; Miller BT; Herbert JM; Woodcock HL; Brooks BR; Shao Y Computation of Hydration Free Energies Using the Multiple Environment Single

- System Quantum Mechanical/Molecular Mechanical Method. *J. Chem. Theory Comput* 2016, 12, 332–44. [PubMed: 26613419]
- (69). Kearns FL; Hudson PS; Boresch S; Woodcock HL Methods for Efficiently and Accurately Computing Quantum Mechanical Free Energies for Enzyme Catalysis. *Methods Enzymol.* 2016, 577, 75–104. [PubMed: 27498635]
- (70). Schöller A; Kearns F; Woodcock HL; Boresch S Optimizing the Calculation of Free Energy Differences in Nonequilibrium Work SQM/MM Switching Simulations. *J. Phys. Chem. B* 2022, 126, 2798–2811. [PubMed: 35404610]
- (71). Wang H; Zhang L; Han J; E W DeePMD-kit: A deep learning package for many-body potential energy representation and molecular dynamics. *Comput. Phys. Commun* 2018, 228, 178–184.
- (72). Liang W; Zeng J; York DM; Zhang L; Wang H In *A Practical Guide to Recent Advances in Multiscale Modelling and Simulation for Biomolecules*; Wang Y, Ed.; AIP Publishing, 2023; Chapter Learning DeePMD-kit: A guide to building Deep Potential models.
- (73). Chai J-D; Head-Gordon M Systematic optimization of long-range corrected hybrid density functionals. *J. Chem. Phys* 2008, 128, 084106. [PubMed: 18315032]
- (74). Frisch MJ; Trucks GW; Schlegel HB; Scuseria GE; Robb MA; Cheeseman JR; Scalmani G; Barone V; Petersson GA; Nakatsuji H; Li X; Caricato M; Marenich AV; Bloino J; Janesko BG; Gomperts R; Mennucci B; Hratchian HP; Ortiz JV; Izmaylov AF; Sonnenberg JL; Williams-Young D; Ding F; Lipparini F; Egidi F; Goings J; Peng B; Petrone A; Henderson T; Ranasinghe D; Zakrzewski VG; Gao J; Rega N; Zheng G; Liang W; Hada M; Ehara M; Toyota K; Fukuda R; Hasegawa J; Ishida M; Nakajima T; Honda Y; Kitao O; Nakai H; Vreven T; Throssell K; Montgomery JA Jr.; Peralta JE; Ogliaro F; Bearpark MJ; Heyd JJ; Brothers EN; Kudin KN; Staroverov VN; Keith TA; Kobayashi R; Normand J; Raghavachari K; Rendell AP; Burant JC; Iyengar SS; Tomasi J; Cossi M; Millam JM; Klene M; Adamo C; Cammi R; Ochterski JW; Martin RL; Morokuma K; Farkas O; Foresman JB; Fox DJ Gaussian~16 Revision A.03. 2016; Gaussian Inc. Wallingford CT.
- (75). Gaus M; Goez A; Elstner M Parametrization and Benchmark of DFTB3 for Organic Molecules. *J. Chem. Theory Comput* 2013, 9, 338–354. [PubMed: 26589037]
- (76). Ramakrishnan R; Dral PO; Rupp M; von Lilienfeld OA Quantum chemistry structures and properties of 134 kilo molecules. *Sci. Data* 2014, 1, 140022. [PubMed: 25977779]
- (77). Rupp M; Tkatchenko A; Müller K-R; von Lilienfeld O Fast and Accurate Modeling of Molecular Atomization Energies with Machine Learning. *Phys. Rev. Lett* 2012, 108, 058301. [PubMed: 22400967]
- (78). Smith JS; Isayev O; Roitberg AE ANI-1, A data set of 20 million calculated off-equilibrium conformations for organic molecules. *Sci. Data* 2017, 4, 170193–170200. [PubMed: 29257127]
- (79). Ruddigkeit L; van Deursen R; Blum LC; Reymond J-L Enumeration of 166 Billion Organic Small Molecules in the Chemical Universe Database GDB-17. *J. Chem. Inf. Model* 2012, 52, 2864–2875. [PubMed: 23088335]
- (80). Weininger D SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *J. Chem. Inf. Comput. Sci* 1988, 28, 31–36.
- (81). Gao X; Ramezanghorbani F; Isayev O; Smith JS; Roitberg AE TorchANI: A Free and Open Source PyTorch-Based Deep Learning Implementation of the ANI Neural Network Potentials. *J. Chem. Inf. Model* 2020, 60, 3408–3415. [PubMed: 32568524]
- (82). Biondi E; Benner SA Artificially Expanded Genetic Information Systems for New Aptamer Technologies. *Biomedicines* 2018, 6, 53. [PubMed: 29747381]
- (83). O’Boyle NM; Banck M; James CA; Morley C; Vandermeersch T; Hutchison GR Open Babel: An open chemical toolbox. *J Cheminform* 2011, 3, 33. [PubMed: 21982300]
- (84). <http://www.rdkit.org>, RDKit: Open-source cheminformatics.
- (85). Zeng J; Cao L; Xu M; Zhu T; Zhang JZH Complex reaction processes in combustion unraveled by neural network-based molecular dynamics simulation. *Nat. Commun* 2020, 11, 5713. [PubMed: 33177517]
- (86). Wen T; Zhang L; Wang H; E W; Srolovitz DJ Deep potentials for materials science. *Mater. Futures* 2022, 1, 022601.

- (87). Lu D; Jiang W; Chen Y; Zhang L; Jia W; Wang H; Chen M DP Compress: A Model Compression Scheme for Generating Efficient Deep Potential Models. *J. Chem. Theory Comput* 2022, 18, 5559–5567. [PubMed: 35926122]
- (88). Goodfellow I; Bengio Y; Courville A *Deep Learning*; MIT Press, 2016; <http://www.deeplearningbook.org>.
- (89). Case DA; Belfon K; Ben-Shalom IY; Brozell SR; Cerutti DS; Cheatham TE III; Cruzeiro VWD; Darden TA; Duke RE; Giambasu G; ; Gilson MK; Gohlke H; Goetz AW; Harris R; Izadi S; Izmailov SA; Kasavajhala K; Kovalenko K; Krasny R; Kurtzman T; Lee T; Le-Grand S; Li P; Lin C; Liu J; Luchko T; Luo R; Man V; Merz K; Miao Y; Mikhailovskii O; Monard G; ; Nguyen H; Onufriev A; Pan F; Pantano S; Qi R; Roe DR; Roitberg A; Sagui C; Schott-Verdugo S; Shen J; Simmerling CL; Skrynnikov N; Smith J; Swails J; Walker RC; Wang J; Wilson RM; Wolf RM; Wu X; Xiong Y; Xue Y; York DM; Kollman PA AMBER 20. University of California, San Francisco: San Francisco, CA, 2020.
- (90). Walker RC; Crowley MF; Case DA The implementation of a fast and accurate QM/MM potential method in Amber. *J. Comput. Chem* 2008, 29, 1019–1031. [PubMed: 18072177]
- (91). Gaus M; Cui Q; Elstner M DFTB3: Extension of the self-consistent-charge density-functional tight-binding method (SCC-DFTB). *J. Chem. Theory Comput* 2011, 7, 931–948.
- (92). Seabra G; Walker RC; Elstner M; Case DA; Roitberg AE Implementation of the SCC-DFTB method for hybrid QM/MM simulations within the Amber molecular dynamics package. *J. Phys. Chem. A* 2007, 111, 5655–5664. [PubMed: 17521173]
- (93). Dewar MJS; Thiel W A semiempirical model for the two-center repulsion integrals in the NDDO approximation. *Theor. Chim. Acta* 1977, 46, 89–104.
- (94). Dewar MJS; Zoebisch E; Healy EF; Stewart JJP Development and use of quantum mechanical molecular models. 76. AM1: a new general purpose quantum mechanical molecular model. *J. Am. Chem. Soc* 1985, 107, 3902–3909.
- (95). Stewart JJP Optimization of parameters for semiempirical methods V: Modification of NDDO approximations and application to 70 elements. *J. Mol. Model* 2007, 13, 1173–1213. [PubMed: 17828561]
- (96). Hourahine B; Aradi B; Blum V; Bonafe F; Buccheri A; Camacho C; Cevallos C; Deshayes MY; Dumitrica T; Dominguez A; Ehlert S; Elstner M; van der Heide T; Hermann J; Irle S; Kranz JJ; Kohler C; Kowalczyk T; Kubar T; Lee IS; Lutsker V; Maurer RJ; Min SK; Mitchell I; Negre C; Niehaus TA; Niklasson AMN; Page AJ; Pecchia A; Penazzi G; Persson MP; Rezac J; Sanchez CG; Sternberg M; Stohr M; Stuckenberg F; Tkatchenko A.; z. Yu VW; Frauenheim T DFTB+, a software package for efficient approximate density functional theory based atomistic simulations. *J. Chem. Phys* 2020, 152, 124101. [PubMed: 32241125]
- (97). Grimme S; Bannwarth C; Shushkov P A Robust and Accurate Tight-Binding Quantum Chemical Method for Structures, Vibrational Frequencies, and Noncovalent Interactions of Large Molecular Systems Parametrized for All spd-Block Elements ($Z = 1-86$). *J. Chem. Theory Comput* 2017, 13, 1989–2009. [PubMed: 28418654]
- (98). Bannwarth C; Ehlert S; Grimme S GFN2-xTB—An Accurate and Broadly Parametrized Self-Consistent Tight-Binding Quantum Chemical Method with Multipole Electrostatics and Density-Dependent Dispersion Contributions. *J. Chem. Theory Comput* 2019, 15, 1652–1671. [PubMed: 30741547]
- (99). Hjorth Larsen A; Jørgen Mortensen J; Blomqvist J; Castelli IE; Christensen R; Dułak M; Friis J; Groves MN; Hammer B; Hargus C; Hermes ED; Jennings PC; Bjerre Jensen P; Kermode J; Kitchin JR; Leonhard Kolsbjerg E; Kubal J; Kaasbjerg K; Lysgaard S; Bergmann Maronsson J; Maxson T; Olsen T; Pastewka L; Peterson A; Rostgaard C; Schiøtz J; Schütt O; Strange M; Thygesen KS; Vegge T; Vilhelmsen L; Walter M; Zeng Z; Jacobsen KW The atomic simulation environment—a Python library for working with atoms. *J. Phys. Condens. Matter* 2017, 29, 273002. [PubMed: 28323250]
- (100). Liu DC; Nocedal J On the limited memory BFGS method for large scale optimization. *Mathematical Programming* 1989, 45, 503–528.
- (101). Courmia Z; Chipot C; Roux B; York DM; Sherman W In *Free Energy Methods in Drug Discovery—Introduction*; Armacost KA, Thompson DC, Eds.; ACS Symposium Series; 2021; Vol. 1397; pp 1–38.

- (102). Weber W; Thiel W Orthogonalization corrections for semiempirical methods. *Theor. Chem. Acc* 2000, 103, 495–506.
- (103). Huang M; Giese TJ; Lee T-S; York DM Improvement of DNA and RNA Sugar Pucker Profiles from Semiempirical Quantum Methods. *J. Chem. Theory Comput* 2014, 10, 1538–1545. [PubMed: 24803866]
- (104). Kuechler ER; Giese TJ; York DM Charge-dependent many-body exchange and dispersion interactions in combined QM/MM simulations. *J. Chem. Phys* 2015, 143, 234111. [PubMed: 26696050]
- (105). Perreault DM; Anslyn EV Unifying the Current Data on the Mechanism of Cleavage-Transesterification of RNA. *Angew. Chem. Int. Ed* 1997, 36, 432–450.
- (106). Emilsson GM; Nakamura S; Roth A; Breaker RR Ribozyme speed limits. *RNA* 2003, 9, 907–918. [PubMed: 12869701]
- (107). Bevilacqua PC; Harris ME; Piccirilli JA; Gaines C; Ganguly A; Kostenbader K; Ekesan ; York DM An Ontology for Facilitating Discussion of Catalytic Strategies of RNA-Cleaving Enzymes. *ACS Chem. Biol* 2019, 14, 1068–1076. [PubMed: 31095369]
- (108). Anslyn E; Breslow R On the mechanism of catalysis by ribonuclease: cleavage and isomerization of the dinucleotide UpU catalyzed by imidazole buffers. *J. Am. Chem. Soc* 1989, 111, 4473–4482.
- (109). Raines RT Ribonuclease A. *Chem. Rev* 1998, 98, 1045–1066. [PubMed: 11848924]
- (110). Gu H; Zhang S; Wong K-Y; Radak BK; Dissanayake T; Kellerman DL; Dai Q; Miyagi M; Anderson VE; York DM; Piccirilli JA; Harris ME Experimental and computational analysis of the transition state for ribonuclease A-catalyzed RNA 2'-*O*-transphosphorylation. *Proc. Natl. Acad. Sci. USA* 2013, 110, 13002–13007. [PubMed: 23878223]
- (111). Lilley DMJ Classification of the nucleolytic ribozymes based upon catalytic mechanism. *F1000 Res.* 2019, 8, 1462.
- (112). Gaines CS; Piccirilli JA; York DM The L-platform/L-scaffold framework: a blueprint for RNA-cleaving nucleic acid enzyme design. *RNA* 2020, 26, 111–125. [PubMed: 31776179]
- (113). Breaker RR; Joyce GF A DNA enzyme that cleaves RNA. *Chem. Biol* 1994, 1, 223–229. [PubMed: 9383394]
- (114). Liu H; Yu X; Chen Y; Zhang J; Wu B; Zheng L; Haruehanroengra P; Wang R; Li S; Lin J; Li J; Sheng J; Huang Z; Ma J; Gan J Crystal Structure of an RNA-Cleaving DNAzyme. *Nat. Commun* 2017, 8, 2006–2015. [PubMed: 29222499]
- (115). Ekesan ; York DM Dynamical ensemble of the active state and transition state mimic for the RNA-cleaving 8–17 DNAzyme in solution. *Nucleic Acids Res.* 2019, 47, 10282–10295. [PubMed: 31511899]
- (116). Borggräfe J; Victor J; Rosenbach H; Viegas A; Gertzen CGW; Wuebben C; Kovacs H; Gopalswamy M; Riesner D; Steger G; Schiemann O; Gohlke H; Span I; Eitzkorn M Time-resolved structural analysis of an RNA-cleaving DNA catalyst. *Nature* 2022, 601, 144–149. [PubMed: 34949858]
- (117). Formoso E; Matxain JM; Lopez X; York DM Molecular dynamics simulation of bovine pancreatic ribonuclease A-CpA and transition state-like complexes. *J. Phys. Chem. B* 2010, 114, 7371–7382. [PubMed: 20455590]
- (118). Brandstetter H; Grams F; Glitz D; Lang A; Huber R; Bode W; Krell HW; Engh RA The 1.8-Å Crystal Structure of a Matrix Metalloproteinase 8-Barbiturate Inhibitor Complex Reveals a Previously Unobserved Mechanism for Collagenase Substrate Recognition. *J. Biol. Chem* 2001, 276, 17405–12. [PubMed: 11278347]
- (119). Natesan S; Subramaniam R; Bergeron C; Balaz S Binding Affinity Prediction for Ligands and Receptors Forming Tautomers and Ionization Species: Inhibition of Mitogen-Activated Protein Kinase-Activated Protein Kinase 2 (MK2). *J. Med. Chem* 2012, 55, 2035–47. [PubMed: 22280316]
- (120). Soliva R; Gelpí JL; Almansa C; Virgili M; Orozco M Dissection of the Recognition Properties of p38 MAP Kinase. Determination of the Binding Mode of a New Pyridinyl-Heterocycle Inhibitor Family. *J. Med. Chem* 2007, 50, 283–93. [PubMed: 17228870]

- (121). Christensen AS; Sirumalla SK; Qiao Z; O'Connor MB; Smith DGA; Ding F; Bygrave PJ; Anandkumar A; Welborn M; Manby FR; Miller TF 3rd OrbNet Denali: A machine learning potential for biological and organic chemistry with semi-empirical cost and DFT accuracy. *J. Chem. Phys* 2021, 155, 204103. [PubMed: 34852495]
- (122). Zubatyuk R; Smith JS; Nebgen BT; Tretiak S; Isayev O Teaching a neural network to attach and detach electrons from molecules. *Nat. Commun* 2021, 12, 4870. [PubMed: 34381051]
- (123). Isert C; Atz K; Jiménez-Luna J; Schneider G QMugs, quantum mechanical properties of drug-like molecules. *Sci. Data* 2022, 9, 273. [PubMed: 35672335]
- (124). Eastman P; Behara PK; Dotson DL; Galvelis R; Herr JE; Horton JT; Mao Y; Chodera JD; Pritchard BP; Wang Y; De Fabritiis G; Markland TE SPICE, A Dataset of Drug-like Molecules and Peptides for Training Machine Learning Potentials. 2022; <https://arxiv.org/abs/2209.10702>.
- (125). Khrabrov K; Shenbin I; Ryabov A; Tsypin A; Telepov A; Alekseev A; Grishin A; Strashnov P; Zhilyaev P; Nikolenko S; Kadurin A nablaDFT: Large-Scale Conformational Energy and Hamiltonian Prediction benchmark and dataset. *Phys. Chem. Chem. Phys* 2022, 24, 25853–25863. [PubMed: 36279016]
- (126). Tian C; Kasavajhala K; Belfon KAA; Raguette L; Huang H; Miguez AN; Bickel J; Wang Y; Pincay J; Wu Q; Simmerling C ff19SB: Amino-Acid-Specific Protein Backbone Parameters Trained against Quantum Mechanics Energy Surfaces in Solution. *J. Chem. Theory Comput* 2020, 16, 528–552. [PubMed: 31714766]
- (127). Pérez A; Marchán I; Svozil D; Sponer J; Cheatham TE III; Lughton CA; Orozco M Refinement of the AMBER force field for nucleic acids: Improving the description of α/γ conformers. *Biophys. J* 2007, 92, 3817–3829. [PubMed: 17351000]
- (128). Zgarbová M; Otyepka M; Šponer J; Mládek A; Banáš P; Cheatham TE III; Jurek P Refinement of the Cornell et al. nucleic acids force field based on reference quantum chemical calculations of glycosidic torsion profiles. *J. Chem. Theory Comput* 2011, 7, 2886–2902. [PubMed: 21921995]
- (129). Bergonzo C; T. E. C. III Improved Force Field Parameters Lead to a Better Description of RNA Structure. *J. Chem. Theory Comput* 2015, 11, 3969–3972. [PubMed: 26575892]
- (130). Izadi S; Anandakrishnan R; Onufriev AV Building Water Models: A Different Approach. *J. Phys. Chem. Lett* 2014, 5, 3863–3871. [PubMed: 25400877]
- (131). Izadi S; Onufriev AV Accuracy limit of rigid 3-point water models. *J. Chem. Phys* 2016, 145, 074501–074510. [PubMed: 27544113]
- (132). Li P; Roberts BP; Chakravorty DK; Merz KM Jr. Rational design of Particle Mesh Ewald compatible Lennard-Jones parameters for +2 metal cations in explicit solvent. *J. Chem. Theory Comput* 2013, 9, 2733–2748. [PubMed: 23914143]
- (133). Li P; Merz KM Jr. Taking into account the ion-induced dipole interaction in the nonbonded model of ions. *J. Chem. Theory Comput* 2014, 10, 289–297. [PubMed: 24659926]
- (134). Li P; Merz KM Metal Ion Modeling Using Classical Mechanics. *Chem. Rev* 2017, 117, 1564–1686. [PubMed: 28045509]

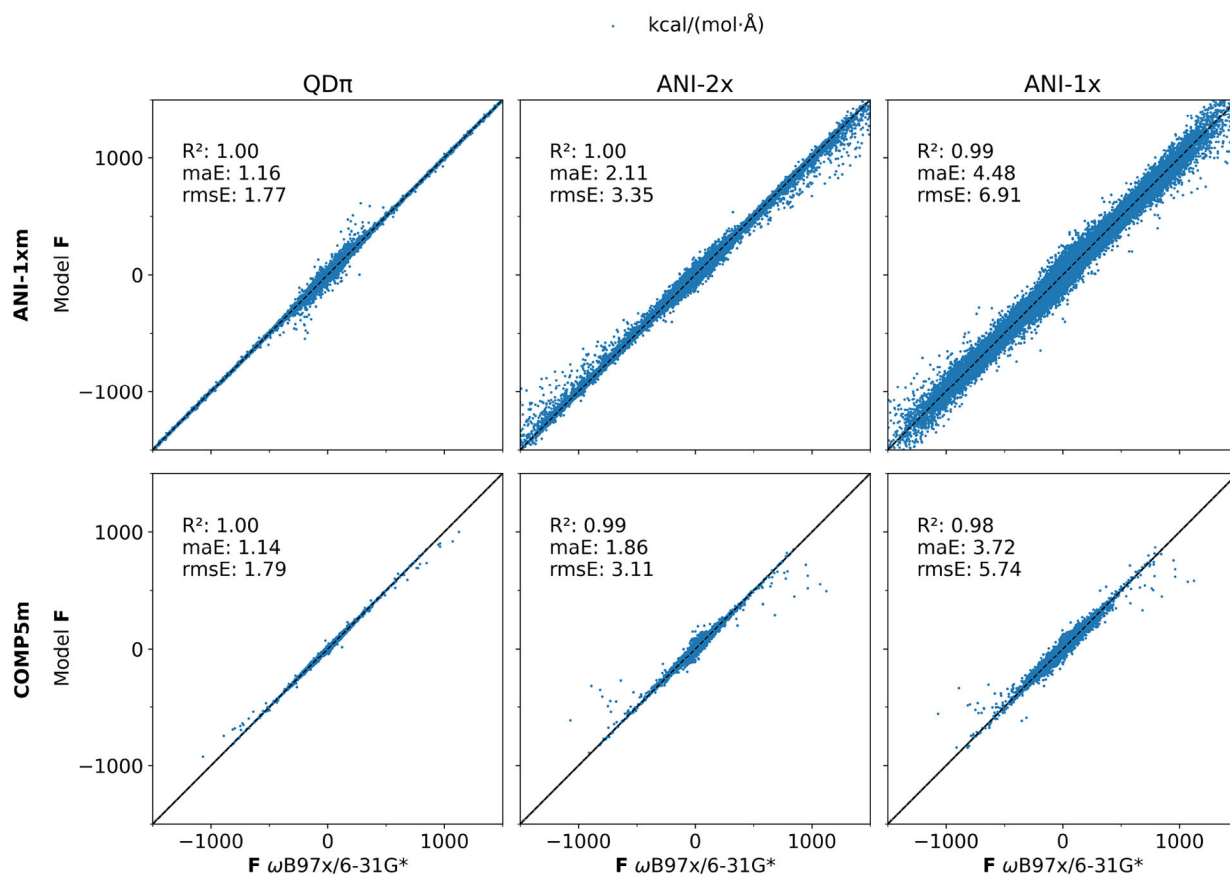


Figure 1: Relation between forces in kcal/(mol·Å) calculated by ω B97X/6-31G* and QD π , ANI-2x, and ANI-1x, respectively for the ANI-1xm and COMP5m data sets.

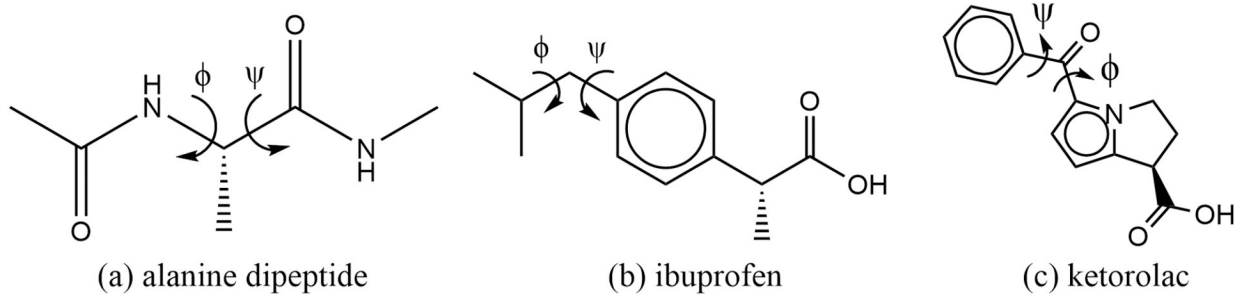


Figure 2:
Model (a) alanine dipeptide; (b) ibuprofen; (c) ketorolac. ϕ and ψ represent of the 2D torsion angles.

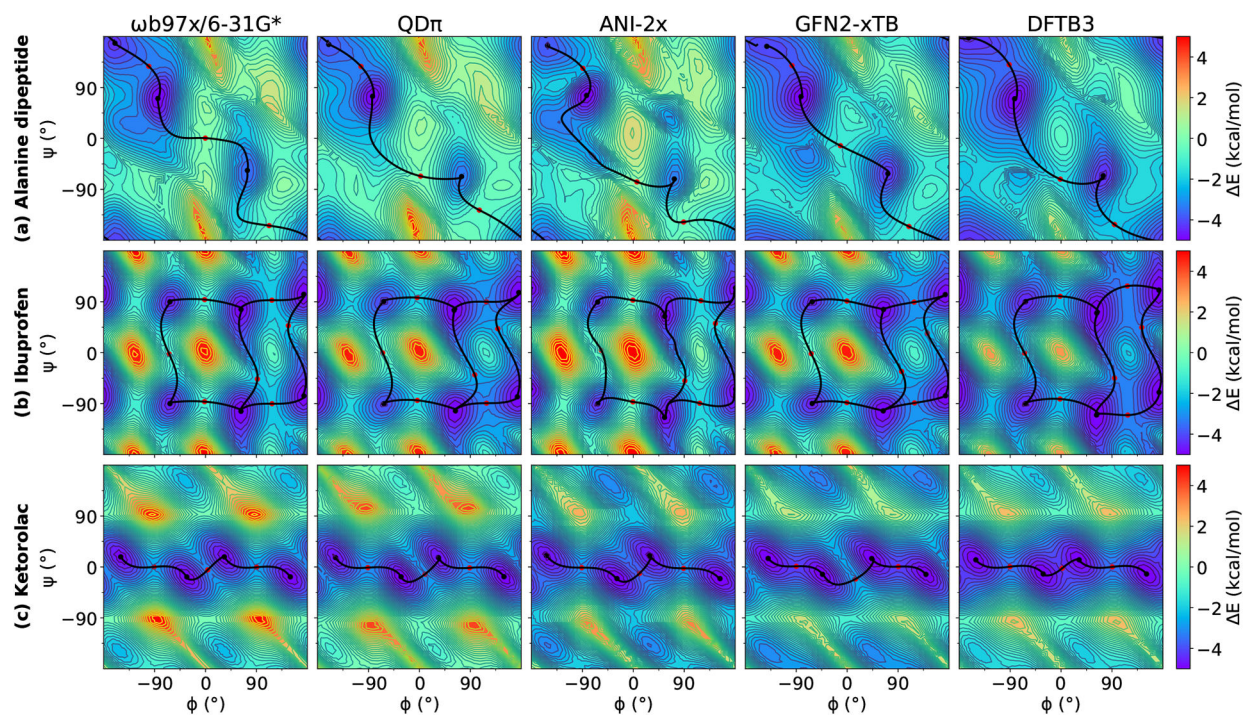


Figure 3:

Relaxed 2D torsion profiles for (a) alanine dipeptide; (b) ibuprofen; (c) ketorolac. Each molecule was computed using ω B97X/6-31G*, QD π , DFTB3, ANI/2x, and GFN2-xTB, respectively. ω B97X/6-31G* is the reference potential and the other potentials are compared with ω B97X/6-31G*. The color bars represent the potential energy (with respect to the minimum energy) ω B97X/6-31G* in kcal/mol. The black and red points represent the minima and the transition states, respectively, and the black curves represent the transition paths between minima.

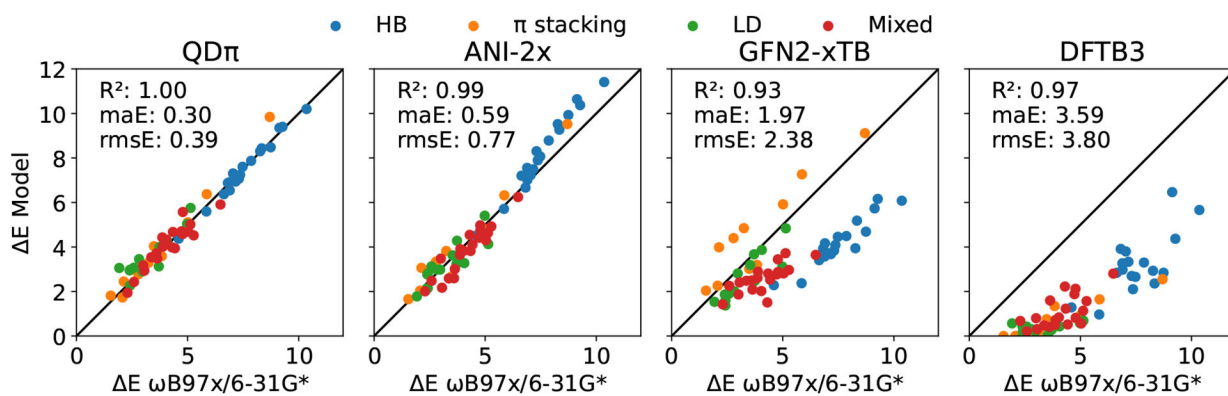


Figure 4: Relation between relative energies in kcal/mol calculated by $\omega B97X/6-31G^*$ and QD π , ANI-2x, GFN2-xTB, and DFTB3, respectively for the S66 \times 8 data set. Relative energies consist of the difference between the optimized structure and the structure with the furthest distance in each of the 66 dimer pairs.

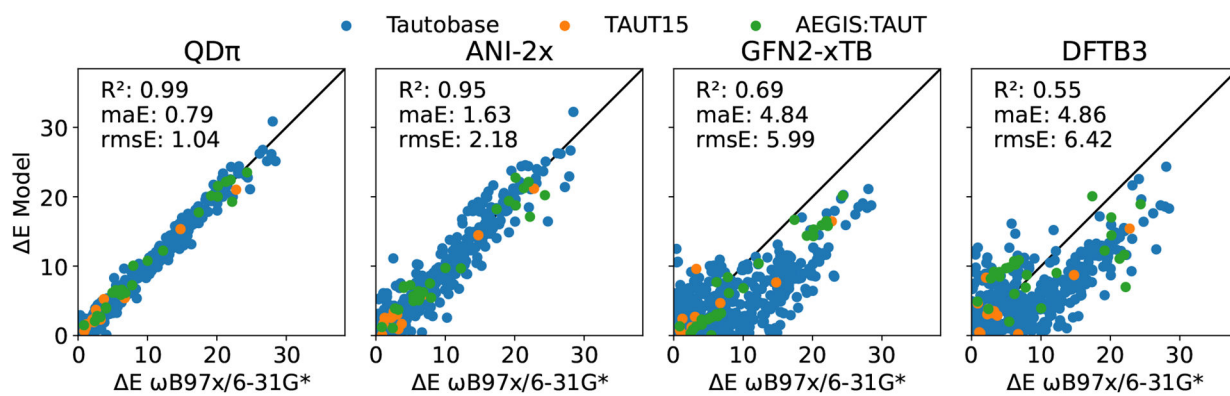


Figure 5:
Relation between tautomerization energies in kcal/mol calculated by ω B97X/6-31G* and QD π , ANI-2x, GFN2-xTB, and DFTB3, respectively, for TAUT15 data set and the artificially expanded genetic information system: Tautomer (AEGIS:TAUT) data set.

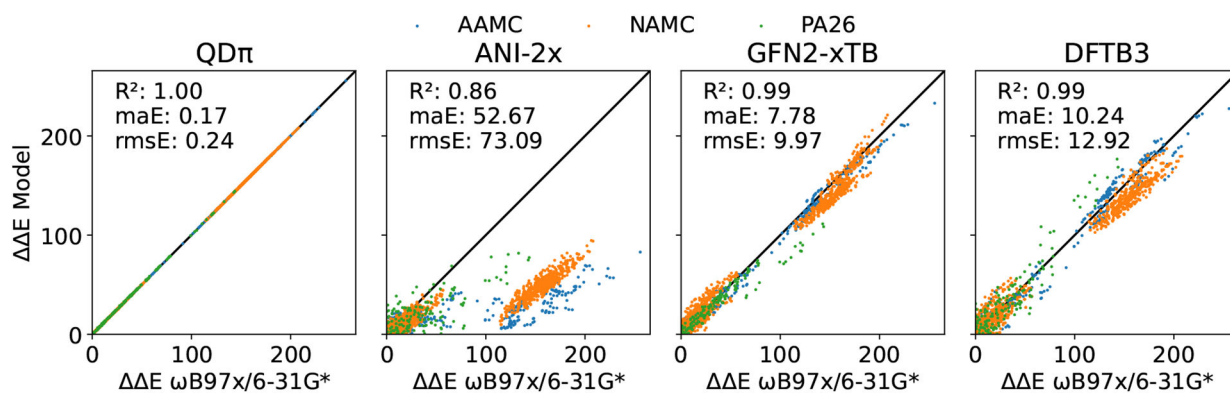


Figure 6:
Relation between of relative protonation energies ($AH + B^- \rightarrow A^- + BH$) in kcal/mol calculated by $\omega B97X/6-31G^*$ and QD π , ANI-2x, GFN2-xTB, and DFTB3, respectively, for AAMC, NAMC, and PA26 data sets.

Table 1:Data sets used in the current work.^a

Usage	data set	# of data points			Ref
		E	F	E	
	ANI-1x	4,956,005	227,101,443	...	36,37
Training	ANI-1xm	2,641,429	130,421,121	...	36,37
	S66×8	528	31,536	462	36,40,41
	HB375×10	3750	192,690	3375	39
	AEGIS:BP	32	1953	10	60
	AEGIS:TAUT	37	1668	25	60
	Tautobase	700	39,216	350	61,62
	AAMC	50	1527	25	43
	NAMC	68	3018	53	43
	PA26	34	1137	17	42
	RegioSQM20 (95%)	1088	84,576	544	63
Testing	COMP5m	64,667	5,215,848	...	36,64–66
	TAUT15	21	831	13	42
	RegioSQM20 (5%)	50	4023	25	63

^aData sets are described in the Methods section. In the current work, all reference DFT data is computed at the ω B97X/6–31G* level of theory for consistency.

Table 2:Data sets and neural network optimization steps used in training different QD π model versions.

Iter.	data set	Steps	Descriptor	Output model
1	ANI-1x (E,F)	10,000,000	Normal	
2	ANI-1xm (E,F)	60,000,000	Compressed	QD π v0.0
3	ANI-1xm (E,F), Tautobase (E, F), AAMC (E, F)	31,000,000	Compressed	QD π v0.1
4	ANI-1xm (E,F), Tautobase (E, F), AAMC (E, F), NAMC (E, F), S66 \times 8 (E, F)	50,000,000	Compressed	QD π v0.2
5	ANI-1xm (E,F), Tautobase (E, F), AAMC (E, F), NAMC (E, F), S66 \times 8 (E, F), PA26 (E, F), RegioSQM20 (95%) (E, F), HB375 \times 10 (E, F)	47,000,000	Compressed	QD π v0.3
6	ANI-1xm (E,F), Tautobase (E, F), AAMC (E, F), NAMC (E, F), S66 \times 8 (E, F), PA26 (E, F), RegioSQM20 (95%) (E, F), HB375 \times 10 (E, F), AEGIS:BP (E, F), AEGIS:TAUT (E, F)	43,800,000	Compressed	QD π v1.0

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

Table 3:

Mean absolute error (MAE) and root mean square error (rmsE) of energies in kcal/mol and forces in kcal/(mol·Å) for the ANI-1xm training and COMP5m testing (marked with an “*”) data sets (Table 1).^a

Model	ANI-1xm				COMP5m*			
	Energy		Force		Energy		Force	
	maE	rmsE	maE	rmsE	maE	rmsE	maE	rmsE
QD π v1.0	0.83	1.22	1.16	1.77	1.48	2.44	1.14	1.79
ANI-1x	1.48	2.07	4.48	6.91	1.96	3.33	3.72	5.74
ANI-2x	1.07	1.58	2.11	3.35	1.67	2.66	1.86	3.11
DFTB3	7.58	12.45	5.46	8.76
MNDO/d ^b	15.14	24.53	11.52	17.77
AM1	14.95	24.29	12.13	18.07
PM6	12.96	23.63	9.33	14.30
GFN1-xTB	4.69	7.02	3.68	5.40
GFN2-xTB	5.81	8.65	4.33	6.33

^a Models and data sets are described in the Methods section.

^b Some points (<0.1 %) fail to coverage and are removed.

Table 4:

Mean absolute error (maE) and root mean square error (rmsE) of relative energies (E) in kcal/mol for hydrogen bonding (HB), π stacking, London dispersion (LD) and mixed influence (Mixed) subsets of the S66 \times 8 and HB375 \times 10 training data sets (Table 1).^a

Model	S66 \times 8 E (subsets)								HB375 \times 10	
	HB		π stacking		LD		Mixed		E	
	maE	rmsE	maE	rmsE	maE	rmsE	maE	rmsE	maE	rmsE
QD π v1.0	0.21	0.26	0.35	0.47	0.41	0.51	0.31	0.38	0.90	1.18
ANI-1x	2.01	3.12	1.34	1.49	0.78	0.98	1.42	1.70	1.54	1.97
ANI-2x	0.91	1.11	0.49	0.57	0.40	0.48	0.40	0.49	1.62	3.49
DFTB3	4.64	4.75	3.14	3.39	2.93	3.07	3.06	3.15	4.12	4.34
MNDO/d	5.37	5.73	2.68	2.83	2.49	2.58	2.66	2.75	6.44	7.04
AM1	7.07	7.63	3.76	4.21	3.04	3.26	3.74	3.83	5.28	5.74
PM6	9.86	10.87	3.83	4.33	3.28	3.42	4.12	4.24	3.92	4.23
GFN1-xTB	3.31	3.45	0.88	1.09	0.90	1.03	1.95	2.03	2.52	2.68
GFN2-xTB	3.53	3.59	0.96	1.12	0.65	0.82	1.55	1.69	2.66	2.89

^aModels and data sets are described in the Methods section. Relative energies consist of the difference between the optimized structure and the structure with the furthest distance in each of dimer pairs.

Table 5:

Mean absolute error (maE) and root mean square error (rmsE) of energies and relative energies (E) in kcal/mol for the Tautobase/AEGIS:TAUT training and TAUT15 testing (marked with an “**”) data sets (Table 1).^a

Model	Tautobase		TAUT15*		AEGIS:TAUT	
	maE	rmsE	maE	rmsE	maE	rmsE
QD π v1.0	0.82	1.09	0.70	0.89	0.71	0.97
ANI-1x	1.73	2.42	1.63	1.83	1.54	2.08
ANI-2x	1.76	2.39	1.00	1.20	1.41	1.94
DFTB3	5.45	6.93	3.65	4.60	5.25	6.12
MNDO/d	9.69	11.35	7.78	9.55	8.20	9.43
AM1	5.01	6.34	3.99	5.85	3.88	4.66
PM6	4.90	6.12	5.60	7.11	7.39	8.85
GFN1-xTB	5.23	6.51	5.32	6.53	5.61	6.58
GFN2-xTB	5.68	6.81	2.84	3.62	3.16	3.59

^aModels and data sets are described in the Methods section.

Table 6:

Mean absolute error (maE) and root mean square error (rmsE) of energies and relative energies (E) in kcal/mol and forces in kcal/(mol·Å) for relative protonation energies for the AAMC, NAMC and PA26 training and RegioSQM* testing (validation) data sets (Table 1).^a

Model	AAMC		NAMC		PA26		RegioSQM*	
	E		E		E		E	
	maE	rmsE	maE	rmsE	maE	rmsE	maE	rmsE
QD π v1.0	0.09	0.14	0.17	0.22	0.39	0.49	2.53	3.19
ANI-1x	86.95	112.31	52.68	71.74	43.02	62.28	16.85	22.15
ANI-2x	70.52	89.39	52.48	72.33	23.80	30.54	13.64	17.24
DFTB3	8.63	11.12	10.85	13.33	12.54	15.84	4.59	5.74
MNDO/d	11.71	14.13	11.29	14.08	13.07	16.10	5.18	6.29
AM1	4.43	5.49	7.32	9.10	13.51	20.89	4.13	5.11
PM6	11.23	13.86	11.03	13.58	17.84	34.36	5.30	6.57
GFN1-xTB	5.00	6.07	11.73	35.39	4.43	5.39	4.10	5.07
GFN2-xTB	5.77	7.14	8.45	10.40	7.35	11.73	4.12	4.96

^aModels and data sets are described in the Methods section.

Table 7:

Selected relative protonation/deprotonation energies from ω B97X/6-31G* and model error (kcal/mol) relevant to acid/base catalysis in RNA cleavage reactions.^a

Protonation pair	ω B97X/6-31G*	QD π	DFTB3	ANI-2x	GFN2-xTB
	E	Err	Err	Err	Err
Lys:NH ₂ + iPrOH → Lys:NH ₃ ⁺ + iPrO ⁻	167.76	0.00	6.11	-115.04	0.04
His:N _e + iPrOH → His:N _e H ⁺ + iPrO ⁻	158.33	0.02	-11.33	-126.62	-7.02
His:N _e H ⁺ + EtO ⁻ → His:N _e + EtOH	-160.25	0.04	11.71	137.70	9.66
G:N ₁ ⁻ + iPrOH → G:N ₁ H + iPrO ⁻	43.06	-1.11	-8.63	-28.62	-2.69
A:N ₁ H ⁺ + EtO ⁻ → A:N ₁ + EtOH	-165.06	1.25	15.21	137.24	10.02
A:N ₃ H ⁺ + EtO ⁻ → A:N ₃ + EtOH	-190.89	1.21	16.00	143.42	11.40
C:N ₃ H ⁺ + EtO ⁻ → C:N ₃ + EtOH	-160.33	0.89	4.66	145.20	6.58
A:N ₁ H ⁺ + G:N ₁ ⁻ → A:N ₁ + G:N ₁ H	-120.07	0.08	6.20	97.55	4.69
A:N ₃ H ⁺ + G:N ₁ ⁻ → A:N ₃ + G:N ₁ H	-145.91	0.04	6.99	103.73	6.07
C:N ₃ H ⁺ + G:N ₁ ⁻ → C:N ₃ + G:N ₁ H	-115.34	-0.27	-4.35	105.50	1.25
maE	...	0.49	9.12	114.06	5.94
rmsE	...	0.71	9.96	118.72	6.96

^aModels and data sets are described in the Methods section. Shown are model reactions for protein enzymes (top block) and nucleic acid enzymes (middle block). Additionally, the relative acid and base protonation/deprotonation energies for different nucleobases are provided (bottom block).