

# Modern semiempirical electronic structure methods and machine learning potentials for drug discovery: Conformers, tautomers, and protonation states

Cite as: *J. Chem. Phys.* **158**, 124110 (2023); doi: [10.1063/5.0139281](https://doi.org/10.1063/5.0139281)

Submitted: 19 December 2022 • Accepted: 3 March 2023 •

Published Online: 22 March 2023



View Online



Export Citation



CrossMark

Jinzhe Zeng (曾晋哲),  Yujun Tao (陶玉君),  Timothy J. Giese,  and Darrin M. York<sup>a)</sup> 

## AFFILIATIONS

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

**Note:** This paper is part of the JCP Special Topic on Modern Semiempirical Electronic Structure Methods.

<sup>a)</sup>Author to whom correspondence should be addressed: [Darrin.York@rutgers.edu](mailto:Darrin.York@rutgers.edu)

## ABSTRACT

Modern semiempirical electronic structure methods have considerable promise in drug discovery as universal “force fields” that can reliably model biological and drug-like molecules, including alternative tautomers and protonation states. Herein, we compare the performance of several neglect of diatomic differential overlap-based semiempirical (MNDO/d, AM1, PM6, PM6-D3H4X, PM7, and ODM2), density-functional tight-binding based (DFTB3, DFTB/ChIMES, GFN1-xTB, and GFN2-xTB) models with pure machine learning potentials (ANI-1x and ANI-2x) and hybrid quantum mechanical/machine learning potentials (AIQM1 and QD $\pi$ ) for a wide range of data computed at a consistent  $\omega$ B97X/6-31G\* level of theory (as in the ANI-1x database). This data includes conformational energies, intermolecular interactions, tautomers, and protonation states. Additional comparisons are made to a set of natural and synthetic nucleic acids from the artificially expanded genetic information system that has important implications for the design of new biotechnology and therapeutics. Finally, we examine the acid/base chemistry relevant for RNA cleavage reactions catalyzed by small nucleolytic ribozymes, DNazymes, and ribonucleases. Overall, the hybrid quantum mechanical/machine learning potentials appear to be the most robust for these datasets, and the recently developed QD $\pi$  model performs exceptionally well, having especially high accuracy for tautomers and protonation states relevant to drug discovery.

Published under an exclusive license by AIP Publishing. <https://doi.org/10.1063/5.0139281>

## I. INTRODUCTION

Alchemical free energy (AFE) simulations<sup>1</sup> are widely used for the prediction of ligand-protein binding energies in drug discovery. These predictions are used to prioritize compounds for costly synthesis and testing in the lead optimization cycle.<sup>2</sup> The predictive capability of these methods relies critically on the accuracy of the force fields that are used.<sup>3</sup> For well-studied biological systems such as proteins<sup>4–6</sup> and common solvents such as water<sup>7–11</sup> and monovalent ions,<sup>12–15</sup> several molecular mechanical (MM) force fields<sup>16,17</sup> have been developed and have undergone extensive validation and revision based on comparison with a wide range of experiments. These force fields have evolved to become increasingly robust and reliable in long-time molecular dynamics simulations, despite the simplicity

of their functional forms. On the other hand, the “general” molecular mechanical force fields needed to model drug-like molecules that may not have ever been synthesized before are generally much less reliable. Moreover, conventional MM force fields are not “universal” in the sense that they use a pre-defined covalent bonding topology and are thus limited in their ability to model alternative tautomers and protonation states. This is important as 30% of the compounds in vendor databases and 21% of the compounds in drug databases have potential tautomers;<sup>18,19</sup> furthermore, it has been estimated that up to 95% of drug molecules contain ionizable groups<sup>18</sup> (~75% weak bases and ~20% weak acids<sup>20,21</sup>).

Modern semiempirical quantum mechanical (QM) electronic structure methods<sup>22,23</sup> provide an attractive alternative to the general MM force fields for drug discovery. The reason is that, unlike

a typical protein that may contain several thousands of atoms, ~79% of drugs are between 10 and 40 non-hydrogen atoms, and the vast majority are less than 100 non-hydrogen atoms.<sup>24</sup> This is of the size range where semiempirical QM methods are able to be used in combined quantum mechanical/molecular mechanical (QM/MM) simulations that include explicit MM representations of the entire protein and surrounding solvent bath under periodic boundary conditions.<sup>25–28</sup> Highly efficient [including parallel and graphics processing unit (GPU)-accelerated] implementations of semiempirical molecular orbital<sup>29</sup> and density-functional tight-binding<sup>30</sup> have been made and are available for molecular dynamics simulations. More importantly, in the context of AFE simulations, these QM/MM potentials can be efficiently integrated into thermodynamic cycles using indirect (or sometimes referred to as “book-ending” or “reference potential”) approaches<sup>31–35</sup> that apply an end-state MM  $\rightarrow$  QM free energy correction to a high-precision MM AFE simulation.

One potential caveat is the high level of accuracy required by drug discovery applications that seek to distinguish binding free energies at a resolution below  $k_B T$  (0.59 kcal/mol at 300 K).<sup>36–38</sup> This is extremely challenging for even the most advanced modern semiempirical QM methods. One path forward that appears promising is to use machine-learning potentials (MLPs) either as stand-alone alternative models,<sup>39–44</sup> or else to augment existing semiempirical QM methods.<sup>45–51</sup> We will refer to the former class as “pure MLPs” and the latter class as “QM/ $\Delta$ -MLPs”. MLPs have emerged as powerful tools to enable fast and accurate chemical models within the scope of their training.<sup>39,41–44</sup> Many such models have emerged for different applications,<sup>52–67</sup> although few, if any, have been used to their full potential in rigorous AFE simulations. Application of these models in drug discovery AFE simulations is challenging because they must: (1) make robust predictions for molecules within the relevant medicinal chemistry space that may have never been synthesized or characterized;<sup>68</sup> (2) model a wide range of intra- and intermolecular interactions, including relative conformational energies, hydrogen bonding,<sup>69</sup>  $\pi$  stacking,<sup>70,71</sup> London dispersion,<sup>72</sup> and mixed interactions; (3) quantitatively handle different tautomers,<sup>18,19,73</sup> and protonation states.<sup>21</sup> Currently, the ANI<sup>63,74–76</sup> class of models, and particularly the second generation ANI-2x,<sup>76</sup> have received widespread attention. A limitation of these models is that they were built for neutral molecules, and their functional forms do not explicitly account for total molecular charge or spin state. Consequently, they are not able to reliably predict the energetics of changing protonation states. This is a serious limitation, as it has been estimated that up to 95% of drug molecules contain ionizable groups.<sup>18</sup> Related to this, some of the pure MLPs did not initially treat long-ranged electrostatic interactions, although there have been efforts to remedy this.<sup>66</sup> Alternatively, there have been several recent efforts to develop new QM/ $\Delta$ -MLPs,<sup>45–51,77</sup> the most relevant in the current context being AIQM1,<sup>46</sup> which is based on the novel ODMx class of semiempirical models<sup>78</sup> and has recently been demonstrated to be robust for transition state optimizations.<sup>79</sup>

Very recently, we introduced a first-generation QM/ $\Delta$ -MLP for drug discovery.<sup>77</sup> The Quantum Deep-learning Potential Interaction (QD $\pi$ ) model uses a fast, robust third-order self-consistent density-functional tight binding (DFTB3/3OB) model<sup>80,81</sup> that is corrected to high-level accuracy through an MLP correction ( $\Delta$ -MLP) based

on our range-corrected deep-learning potential (DPRc)<sup>47,48</sup> as part of the DeePMD-kit<sup>82</sup> interfaced with AMBER.<sup>83</sup> The underlying DFTB3 model is able to capture long-range electrostatic interactions as well as changes in charge, protonation, and spin state. The intramolecular and short- to mid-range intermolecular interactions are made quantitatively accurate by training the DPRc model to correct the total energy and forces to match those of high-level *ab initio* methods.

In the present work, we compare the performance of several modern semiempirical QM, QM/ $\Delta$ -MLP, and pure MLP models against consistent reference data derived from databases relevant for drug discovery. Of particular focus in the present work is characterizing the ability of different potentials to accurately model intermolecular interactions, tautomers, and protonation states. Toward that end, we consider the dataset of natural and synthetic nucleic acids from the artificially expanded genetic information system (AEGIS)<sup>84–87</sup> that is being used for a wide range of biotechnology applications.<sup>88</sup> The system uses 12 different nucleobases in its genetic code, including the four canonical nucleobases found in DNA (adenine, cytosine, guanine, and thymine), in addition to eight synthetic nucleobases. These serve as good test systems as they contain complex covalent bonding and exhibit a rich set of tautomer forms, hydrogen bonded complexes, and alternative protonation states. The remainder of the article is organized as follows: Sec. II describes the computational details pertaining to the various semiempirical QM (MNDO/d,<sup>89</sup> AM1,<sup>90</sup> PM6,<sup>91</sup> PM6-D3H4X,<sup>92,93</sup> PM7,<sup>94</sup> ODM2,<sup>78</sup> DFTB3,<sup>95</sup> GFN1-xTB,<sup>96</sup> GFN2-xTB,<sup>97</sup> and DFTB/ChIMES<sup>98</sup>), MLP (ANI-1x<sup>74</sup> and ANI-2x<sup>76</sup>), and QM/ $\Delta$ -MLP (AIQM1<sup>46</sup> and most recently QD $\pi$ <sup>77</sup>) models, as well as the key modified databases (DBs) used as reference data at the  $\omega$ B97X/6-31G\*<sup>77,99</sup> level. Section III presents and analyzes data for a set of ten broad-spectrum databases for intermolecular interactions, tautomers, protonation states, and 2D conformational energy profiles. Further application is made to examine the performance of modern semiempirical QM, MLP, and QM/ $\Delta$ -MLPs against the AEGIS dataset.<sup>85,86</sup> Finally, the paper provides contextual examples of acid/base chemistry relevant for RNA cleavage reactions catalyzed by small nucleolytic ribozymes and ribonucleases.<sup>100</sup>

## II. METHODS

### A. Models compared in the current work

#### 1. Density-functional reference data

$\omega$ B97X/6-31G\*<sup>99</sup> was performed using Gaussian 16.<sup>101</sup> Reference energy and forces (including geometry optimizations, where needed) were performed at a consistent  $\omega$ B97X/6-31G\*<sup>99</sup> level of theory.

#### 2. NDDO-based semiempirical models

Semiempirical quantum mechanical (QM) models based on the neglect of diatomic differential overlap (NDDO) approximation enable the number of electron repulsion integrals to be drastically reduced and the single-particle density matrix to be decomposed into effective atom-centered atomic orbital products (and their resulting electrostatics represented as multipoles).<sup>102</sup> The NDDO approximation also eliminates the need to explicitly enforce orthogonalization of the molecular orbitals that normally would

be achieved by having an overlap matrix in the generalized eigenvalue equation. Consequently, this may lead to poor modeling of conformational energies and their barriers if left uncorrected. Much work has been performed to introduce orthogonalization corrections into the theoretical framework, which has resulted in the OMX class of methods.<sup>103–106</sup> In the current work, the following NDDO-based methods are considered: MNDO/d,<sup>89</sup> AM1,<sup>90</sup> and PM6<sup>91</sup> that were evaluated with the AMBER 20<sup>107</sup> SQM module;<sup>108</sup> the ODM2<sup>78</sup> method that was evaluated using the MNDO program<sup>109</sup> kindly provided by Dr. Axel Koslowski; and PM6-D3H4X<sup>92,93</sup> and PM7<sup>94</sup> that were performed using the MOPAC software.<sup>110</sup> PM6-D3H4X and PM7 correct PM6 using classical potentials and are often claimed to be the most suitable methodology for drug design among NDDO-based semiempirical models.<sup>111,112</sup>

### 3. DFTB-based semiempirical models

Density-functional tight binding methods offer an intriguing alternative to the NDDO-based semiempirical models. DFTB methods use an expansion of the energy<sup>113</sup> about a sum of neutral atom densities together with a two-center integral approximation to enable a framework for highly efficient calculations (speed is very comparable with NDDO-based methods). Unlike the NDDO-based methods, DFTB methods keep the overlap matrix in the generalized eigenvalue equation and thus explicitly deal with orbital orthogonalization. However, this complicates the decomposition of the density matrix, which now contains two-center products. Various density-matrix partition schemes can be used to map the density onto atomic centers such that an atom-centered (typically monopolar) representation can be made for the second-order electrostatic term in the expansion. The DFTB-based methods considered here include DFTB3<sup>95</sup> (3OB parameters<sup>114</sup>) that was performed using the AMBER 20<sup>107</sup> SQM module,<sup>108,115</sup> and GFN1-xTB,<sup>96</sup> GFN2-xTB,<sup>97</sup> and DFTB/ChIMES<sup>98</sup> (3OB parameters<sup>114</sup> and ChIMES parameters<sup>116</sup> kindly provided by Dr. Cong Huy Pham) models evaluated with the DFTB+ software.<sup>30</sup>

Compared to DFTB3 and GFN1-xTB, GFN2-xTB represents the first broadly parameterized tight-binding method, primarily designed for the fast calculation of structures and noncovalent interaction energies, to include electrostatic and exchange-correlation Hamiltonian terms up to second order in the multipole expansion.<sup>97</sup> In this way, the model takes into account anisotropic second order density fluctuation effects via short-range damped interactions of cumulative atomic multipole moments. DFTB/ChIMES,<sup>116</sup> on the other hand, leverages the relative simplicity of linear regression machine learning in the recently developed Chebyshev Interaction Model for Efficient Simulation (ChIMES) method.<sup>117</sup> Validation tests of DFTB/ChIMES demonstrate the model exhibits both transferability and extensibility and enables physical and chemical predictions with up to coupled-cluster accuracy.<sup>116</sup>

It should be noted that the use of machine learning methods to enhance DFTB models in one form or another is not new. Notable works along these lines, in addition to DFTB/ChIMES, include but are not limited to the ML-Hamiltonian approach of Yaron and co-workers,<sup>118</sup> the development of many-body potentials from deep tensor neural networks,<sup>119,120</sup> Gaussian process regression,<sup>121</sup> and unsupervised machine learning.<sup>122</sup>

### 4. Machine learning potentials (MLPs)

The pure machine learning potentials considered in this work produce energies and atomic forces of a molecule given the positions and elements. These potentials are quite fast compared with semiempirical QM models, and they have more favorable scaling properties. However, some initial pure MLPs were built for neutral molecules in singlet ground states, so they do not reliably model changes in charge state that occur with the addition or loss of electrons and/or protons. The latter of which is important for drug molecules that contain ionizable sites. The pure MLPs considered here include ANI-1x<sup>74</sup> and ANI-2x<sup>76</sup> models performed using the TorchANI software.<sup>123</sup> Both the ANI-1x and ANI-2x models use the ANI descriptor<sup>63</sup> with a cutoff radius of 6 Å and were trained against  $\omega$ B97X/6-31G\* with active learning cycles. The training data of ANI-1x only include energies, and the training data of ANI-2x include both energies and forces.

### 5. Combined semiempirical quantum mechanical and machine-learning potentials (QM/ $\Delta$ -MLPs)

An attractive alternative to either semiempirical QM or pure MLPs is to combine the strengths of both into a combined QM/ $\Delta$ -MLP. In this way, it builds off of a fast and robust semiempirical QM that inherently can accommodate changes in electronic charge and spin states while using MLPs to greatly enhance the accuracy across a broad spectrum of chemical environments. The QM/ $\Delta$ -MLPs considered here include the QD $\pi$ <sup>77</sup> model, which is based on DFTB3/3OB<sup>95,107,108,114,115</sup> and the deep-learning potential available in DeePMD-kit,<sup>82,83</sup> and the AIQM1@DFT\*<sup>46</sup> model, which is based on an ODM2<sup>78,109</sup> model (which includes the D4 dispersion correction<sup>124</sup>) and a trained neural network correction using TorchANI.<sup>123</sup> The MLP component of QD $\pi$  uses the DeepPot-SE descriptor<sup>61</sup> with a cutoff radius of 6 Å and was trained against  $\omega$ B97X/6-31G\* energies and forces for 241 M steps; the MLP part of AIQM1@DFT\* uses the ANI descriptor<sup>63</sup> with a cutoff radius of 6 Å and was trained against  $\omega$ B97X/def2-TZVPP energies and forces for 1000 epochs.<sup>46</sup>

All geometry optimizations using semiempirical QM, MLP, or QM/ $\Delta$ -MLP models were performed using the Limited-memory Broyden–Fletcher–Goldfarb–Shanno (LBFGS) algorithm<sup>125</sup> in the ASE<sup>126</sup> package. Relaxed 2D torsion profiles were made using the same method described in Ref. 77.

### B. Databases and reference data used in the current work

The reference data used in the current work includes the modified ANI-1x,<sup>74,77,127</sup> the modified COMP5,<sup>74,77,128–130</sup> S66x8,<sup>74,131,132</sup> HB375x10,<sup>77,133</sup> TautoBase (TB),<sup>77,134,135</sup> amino acids (AAs) and nucleic acids (NAs),<sup>77,136</sup> PA26 and TAUT15,<sup>77,137</sup> RegioSQM20,<sup>77,138</sup> and the artificially expanded genetic information system (AEGIS).<sup>77,84–86</sup> All reference data were computed (or re-computed<sup>77</sup>) at the  $\omega$ B97X/6-31G\* level of theory (consistent with the most extensive ANI-1x and COMP5 databases).

Among all reference data, the ANI-1x (or modified version) dataset was used to parameterize DFTB/ChIMES, ANI-1x, ANI-2x, QD $\pi$ , and AIQM1; S66x8 was used to parameterize PM6-D3H4X, ANI-2x, and QD $\pi$ ; and TB, AA, NA, PA, and AEGIS were used for the training of QD $\pi$ .

### III. RESULTS AND DISCUSSION

The focus of the current article is on comparing modern semiempirical electronic structure methods and machine learning potentials with respect to their ability to accurately model conformers, tautomers, and protonation states of biological and drug-like molecules. These methods have potential impact for drug discovery owing to their efficiency and robustness.

#### A. Comparison of broad-spectrum databases

Important properties for consideration include relative conformational energies, a wide range of intermolecular interactions, as well as tautomeric and protonation state relative energies. The QD $\pi$  model was trained with the same reference theory level as ANI-2x<sup>76</sup> ( $\omega$ B97X/6-31G<sup>\*</sup>) and considered a number of DBs that encompass conformational energies (ANI-1x, COMP5), intermolecular interactions (S66x8, HB375x10), tautomers (TautoBase, Taut15), and protonation state relative energies (AA, NA, PA26, and RegioSQM20) that are described in detail elsewhere.<sup>77</sup> A comparison of 11 semiempirical quantum and machine learning models is compared against ten databases in Table I.

#### 1. Conformational energy datasets

With respect to the diverse conformational energy datasets (ANI-1x<sup>74,127</sup> and COMP5<sup>74,77,128–130</sup>), the mean absolute errors (maEs) in the forces are smallest for the MLP and  $\Delta$ -MLP potentials (QD $\pi$ , ANI-2x, AIQM1, and ANI-1x), and the QD $\pi$  model performs the overall best (maE values of 1.16 and 1.14 kcal/mol/Å for the ANI-1x and COMP5 datasets, respectively). This is likely due to the fact that the ANI-1x dataset was an integral part of the training of these models. In general, the DFTB models (GFN1-xTB, GFN2-xTB, and DFTB3/3OB) have lower force errors with respect to the reference  $\omega$ B97X/6-31G<sup>\*</sup> values (maE values range from 4.69 to 7.58 and 3.68 to 5.46 kcal/mol/Å for ANI-1x and COMP5, respectively), whereas the NDDO-based methods have considerably larger errors (maE values range from 11.98 to 15.14 and 8.54 to 12.13 kcal/mol/Å for ANI-1x and COMP5, respectively), with PM7 performing the best of the NDDO methods.

GFN2-xTB, and DFTB3/3OB) have lower force errors with respect to the reference  $\omega$ B97X/6-31G<sup>\*</sup> values (maE values range from 4.69 to 7.58 and 3.68 to 5.46 kcal/mol/Å for ANI-1x and COMP5, respectively), whereas the NDDO-based methods have considerably larger errors (maE values range from 11.98 to 15.14 and 8.54 to 12.13 kcal/mol/Å for ANI-1x and COMP5, respectively), with PM7 performing the best of the NDDO methods.

#### 2. Intermolecular interaction datasets

With respect to intermolecular interaction DBs (S66x8<sup>74,131,132</sup> and HB375x10<sup>133</sup>), several models have  $\Delta$ E values below 1 kcal/mol on average (QD $\pi$ , AIQM1, GFN1-xTB, GFN2-xTB, PM6-D3H4X, and PM7), with QD $\pi$  and AIQM1 having exceptional agreement with the reference data: QD $\pi$  has maE values of 0.13 and 0.44 kcal/mol, and AIQM1 has maE values of 0.57 and 0.71 kcal/mol for S66x8 and HB375x10, respectively. The ANI-2x model has excellent maE values for S66x8 (maE 0.37 kcal/mol) but does not perform quite as well for the HB375x10 DB (maE 1.40 kcal/mol). The DFTB3, DFTB/ChIMES, ODM2, and PM6 methods perform similarly with  $\Delta$ E maE values that range from 1.14 to 1.72 (S66x8) and 1.17 to 1.36 (HB375x10) kcal/mol for these DBs. The MNDO/d method has the largest  $\Delta$ E errors (6.67–9.36 kcal/mol), stemming from known limitations in the core–core interactions that particularly affect hydrogen bonding, which the empirical modified core–core repulsions in AM1 were designed in part to partially alleviate (AM1 maE values range from 2.17 to 2.57 kcal/mol).

#### 3. Tautomer datasets

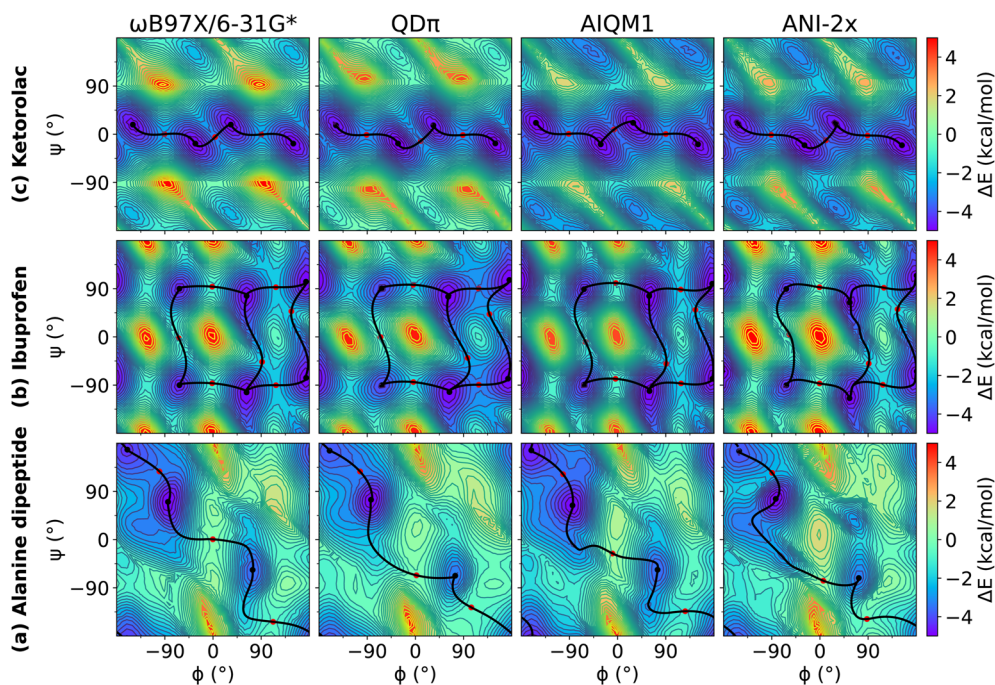
With respect to the tautomer databases, TautoBase<sup>134,135</sup> (TB) and Taut15<sup>137</sup> (T15), only the QD $\pi$  model achieves  $\Delta$ E errors less

TABLE I. Mean absolute errors for different datasets used for training and testing of the QD $\pi$  model.<sup>a</sup> Boldface denotes a vector.

	ANI-1x		S66	TB	AA	NA	PA	COMP5		HB	T15	SQM
	E	F	$\Delta$ E	$\Delta$ E	$\Delta$ E	$\Delta$ E	$\Delta$ E	E	F	$\Delta$ E	$\Delta$ E	$\Delta$ E
QD $\pi$	0.83	1.16	0.13	0.82	0.09	0.17	0.39	1.48	1.14	0.44	0.70	2.53
AIQM1	...	3.10	0.57	2.07	7.30	4.71	5.06	...	2.59	0.71	1.37	2.75
ANI-1x	1.48	4.48	1.41	1.73	86.95	52.68	43.02	1.96	3.72	1.25	1.63	16.85
ANI-2x	1.07	2.11	0.37	1.76	70.52	52.48	23.80	1.67	1.86	1.40	1.00	13.64
GFN2-xTB	...	5.81	0.74	5.68	5.77	8.45	7.35	...	4.33	0.85	2.84	4.12
GFN1-xTB	...	4.69	0.77	5.23	5.00	11.73	4.43	...	3.68	0.87	5.32	4.10
DFTB3	...	7.58	1.14	5.45	8.63	10.85	12.54	...	5.46	1.17	3.65	4.59
DFTB/ChIMES	...	4.82	1.72	5.04	9.47	9.70	12.87	...	4.14	1.36	3.00	6.70
ODM2	...	12.80	1.24	3.37	9.13	5.26	6.04	...	9.97	1.29	3.64	3.99
PM6	...	12.96	1.19	4.90	11.23	11.03	17.84	...	9.33	1.24	5.60	5.30
PM6-D3H4X	...	13.60	0.63	5.44	9.67	11.72	7.78	...	10.27	0.84	6.16	6.61
PM7	...	11.98	0.84	4.34	7.24	10.72	10.11	...	8.54	1.00	3.74	5.93
AM1	...	14.95	2.17	5.01	4.43	7.32	13.51	...	12.13	2.57	3.99	4.13
MNDO/d	...	15.14	6.67	9.69	11.71	11.29	13.07	...	11.52	9.36	7.78	5.18

<sup>a</sup> Mean absolute errors in the energy (E, kcal/mol), forces (F, kcal/mol/Å), and  $\Delta$ E for ANI-1x,<sup>74,127</sup> S66x8 (S66),<sup>74,131,132</sup> TautoBase (TB),<sup>134,135</sup> amino acid and nucleic acid proton affinities (AA and NA),<sup>136</sup> PA26 (PA),<sup>137</sup> COMP5,<sup>74,128–130</sup> HB375x10 (HB),<sup>133</sup> Taut15 (T15),<sup>137</sup> and RegioSQM20 (SQM)<sup>138</sup> databases. The datasets on the right were not part of the QD $\pi$  training.

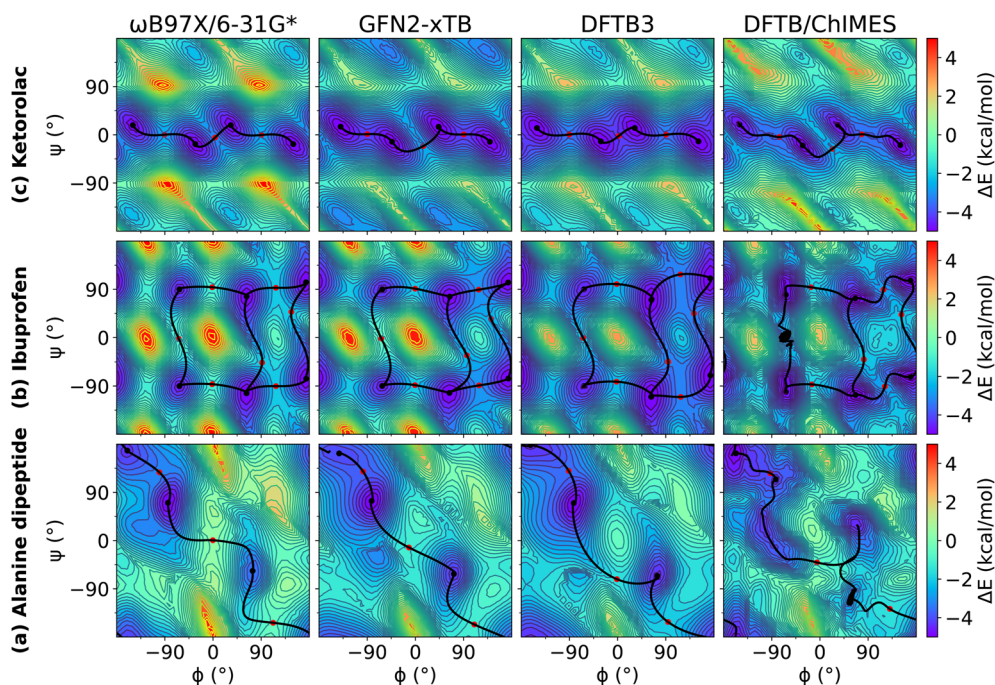




**FIG. 1.** Relaxed 2D torsion profiles for (a) alanine dipeptide; (b) ibuprofen; and (c) ketorolac. Each molecule was computed using  $\omega$ B97X/6-31G\*, QD $\pi$ , AIQM1, and ANI-2x, respectively. The reference level of theory is  $\omega$ B97X/6-31G\*. The color bars represent the potential energy (with respect to the minimum energy) in kcal/mol.

than 1 kcal/mol (maE values of 0.82 and 0.70 kcal/mol for TB and T15, respectively). The AIQM1 and ANI models perform admirably with errors generally below 2 kcal/mol (maE values range from 1.73 to 2.07 and 1.00 to 1.37 kcal/mol for TB and T15, respectively).

The remainder of the DFTB-based methods have maE values in excess of 5 kcal/mol for TB and similar values for the AM1, PM6, and PM6-D3H4X methods. The ODM2 method makes a notable improvement with reduced errors relative to the other NDDO-based



**FIG. 2.** Relaxed 2D torsion profiles for (a) alanine dipeptide; (b) ibuprofen; and (c) ketorolac. Each molecule was computed using  $\omega$ B97X/6-31G\*, GFN2-xTB, DFTB3, and DFTB/ChIMES, respectively. The reference level of theory is  $\omega$ B97X/6-31G\*. The color bars represent the potential energy (with respect to the minimum energy) in kcal/mol.



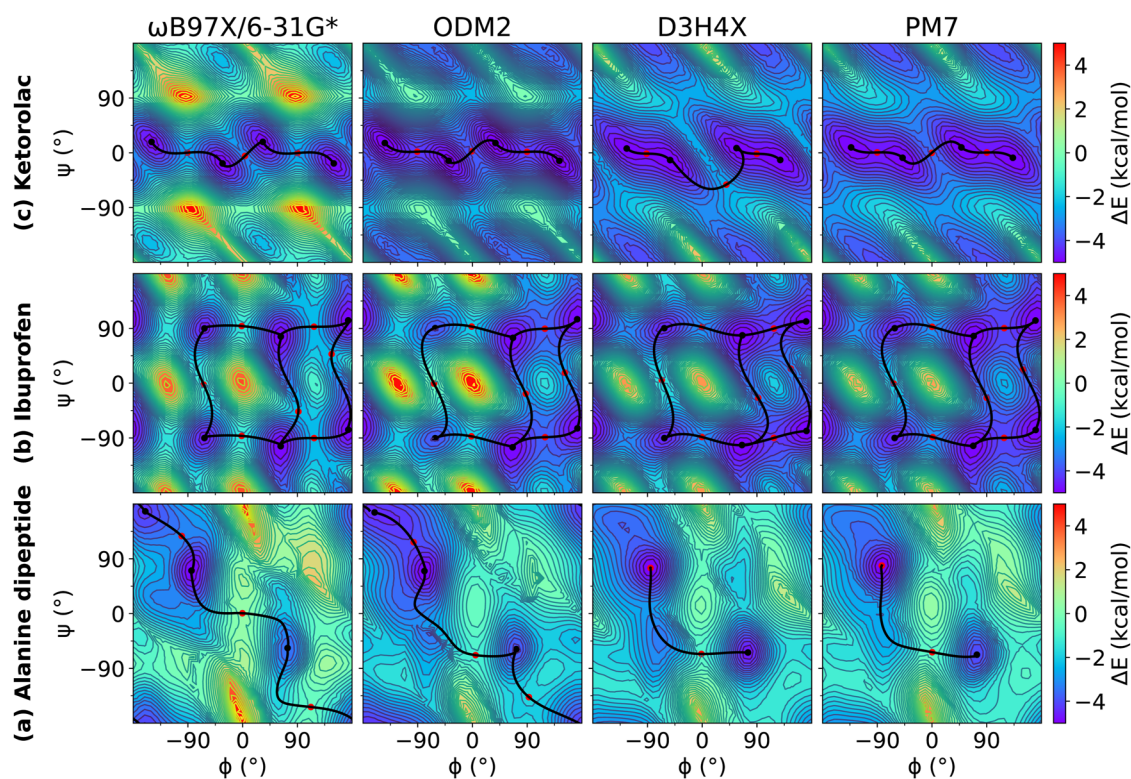
methods (maE values of 3.37 and 3.64 kcal/mol for TB and T15, respectively). The MNDO/d method overall performs the worst with maE values for TB and T15 exceeding 9 kcal/mol.

#### 4. Relative protonation datasets

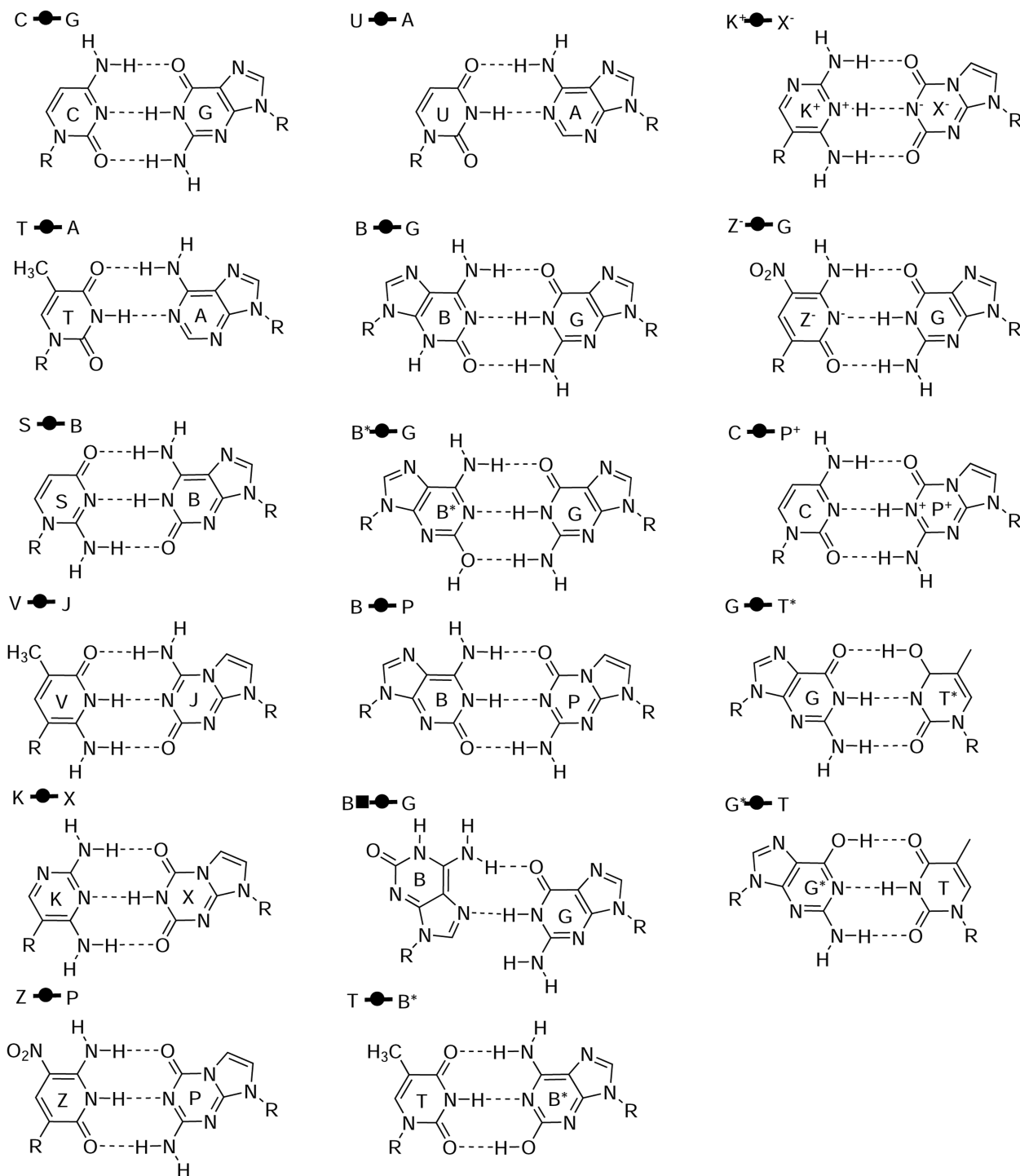
The relative protonation datasets include amino and nucleic acid model compounds<sup>136</sup> (AA and NA) as well as more general proton affinity (PA26<sup>137</sup>) datasets and a subset of the RegioSQM20<sup>138</sup> (SQM) database containing C, H, O, and N elements. The latter involves many relative protonation energies not related to ionizable sites in biological or drug-like molecules, and hence may be of less relevance for drug discovery. For the AA, NA, and PA26 datasets, the QD $\pi$  model stands alone with respect to having very high accuracy in relative deprotonation energies (maE values range from 0.09 to 0.39 kcal/mol). The next best models are AIQM1 (maE 4.71–7.30 kcal/mol). The other semiempirical QM models exhibit much larger ranges: GFN2-xTB (5.77–8.45 kcal/mol), GFN1-xTB (5.00–11.73 kcal/mol), DFTB3 (8.63–12.54 kcal/mol), DFTB/ChIMES (9.47–12.87), ODM2 (5.26–9.13 kcal/mol), PM6 (11.03–17.84 kcal/mol), PM6-D3H4X (7.78–11.72), PM7 (7.24–10.72), AM1 (4.43–13.51 kcal/mol), and MNDO/d (11.29–13.07 kcal/mol). With respect to the SQM dataset, again the QD $\pi$  and AIQM1 models perform best (maE values of 2.53 and 2.75 kcal/mol, respectively), and the remaining semiempirical QM models perform similarly with maE values that range from

3.99 to 6.70 kcal/mol. The pure MLP models (ANI-1x and ANI-2x) break down with respect to their ability to predict relative protonation/deprotonation energies, as these potentials were designed for neutral molecules.

Overall, the QD $\pi$  model performs exceptionally well across all datasets. The AIQM1 model is also impressive in this regard, with the exception of the protonation/deprotonation energies, where AIQM1 has larger errors for the AA, NA, and PA datasets. Clearly, the QM/ $\Delta$ -MLP form, using DFTB3 or ODM2 as a QM base model, considerably enhances the accuracy across all datasets listed in Table I. The pure MLP models, and particularly ANI-2x, generally perform better than the semiempirical QM models, with the exception of protonation/deprotonation energies, where the model gives very large errors. Of the semiempirical QM models, the DFTB-based methods have smaller force errors than the NDDO-based models. The GFN1-xTB, GFN2-xTB, PM6-D3H4X, and PM7 models perform well for intermolecular interactions, slightly better than the DFTB3, DFTB/ChIMES, and ODM2 models. All of the semiempirical QM models are fairly comparable in modeling tautomer energy differences (with the exception of MNDO/d, which is less accurate), with ODM2 performing best over a broad range of data. For protonation/deprotonation energies, however, there is no clear trend with the semiempirical QM potentials—they all deviate from the reference data with  $\Delta E$  maE values exceeding 8 kcal/mol for at least one of the datasets (AA, NA, PA, or SQM).



**FIG. 3.** Relaxed 2D torsion profiles for (a) alanine dipeptide; (b) ibuprofen; and (c) ketorolac. Each molecule was computed using  $\omega$ B97X/6-31G\*, ODM2, PM6-D3H4X, and PM7, respectively. The reference level of theory is  $\omega$ B97X/6-31G\*. The color bars represent the potential energy (with respect to the minimum energy) in kcal/mol.



**FIG. 4.** Structures for the artificially expanded genetic information system (AEGIS) base pair dataset<sup>77,85,86</sup> with Leontis and Westhof symbols used for the classification of nucleic acid base pairs.<sup>139-141</sup>

In the remainder of the article, we focus on comparisons to the most recent modern semiempirical QM (DFTB3, DFTB/ChIMES, GFN2-xTB, ODM2, PM6-D3H4X, PM7), MLP (ANI-2x), and QM/ $\Delta$ -MLP (QD $\pi$  and AIQM1) models.

## B. Comparison of 2D conformation energy profiles


















We examined relaxed 2D torsion profiles for three systems: the alanine dipeptide and the drug molecules ibuprofen and ketorolac, as illustrated in Figs. 1–3. These figures compare 2D torsion profiles at the  $\omega$ B97X/6-31G\* reference level with the QM/ $\Delta$ -MLP/pure MLP models QD $\pi$ , AIQM1, and ANI-2x (Fig. 1), DFTB-based GFN2-xTB, DFTB3 and DFTB/ChIMES (Fig. 2), and NDDO-based ODM2, PM3-D3H4X, and PM7 (Fig. 3) models. The relative energy values for the stationary points are provided in Table S1 of the supplementary material. All of the models qualitatively predict the correct trends. A modest exception occurs with PM6-D3H4X and PM7, which do not predict a pronounced minimum in the  $\beta$  region ( $\sim 180/180$ ) of the  $\phi/\psi$  map (Fig. 3). Overall, the QD $\pi$  and AIQM1 models have the closest agreement with  $\omega$ B97X/6-31G\*, with the

ANI-2x model only slightly worse. The GFN2-xTB, DFTB3, and ODM2 semiempirical QM models tend to systematically underestimate the conformational barriers between minima (Table S1 in the supplementary material). The largest errors that occur for the QD $\pi$  model are for the transition states in the ibuprofen example, which like the semiempirical QM models, are systematically underestimated.

## C. Comparison of hydrogen bond complex energies for natural and artificial nucleic acids

The natural and modified nucleic acids exhibit a wide range of canonical and non-canonical hydrogen bonded base pairs, including some that involve non-standard tautomer forms and protonation states. The base pairs considered in the AEGIS dataset<sup>77,85,86</sup> are illustrated in Fig. 4. This dataset represents a rich set of hydrogen bonding interactions between endocyclic and exocyclic amines and carbonyl and hydroxyl functional groups. The results are listed in

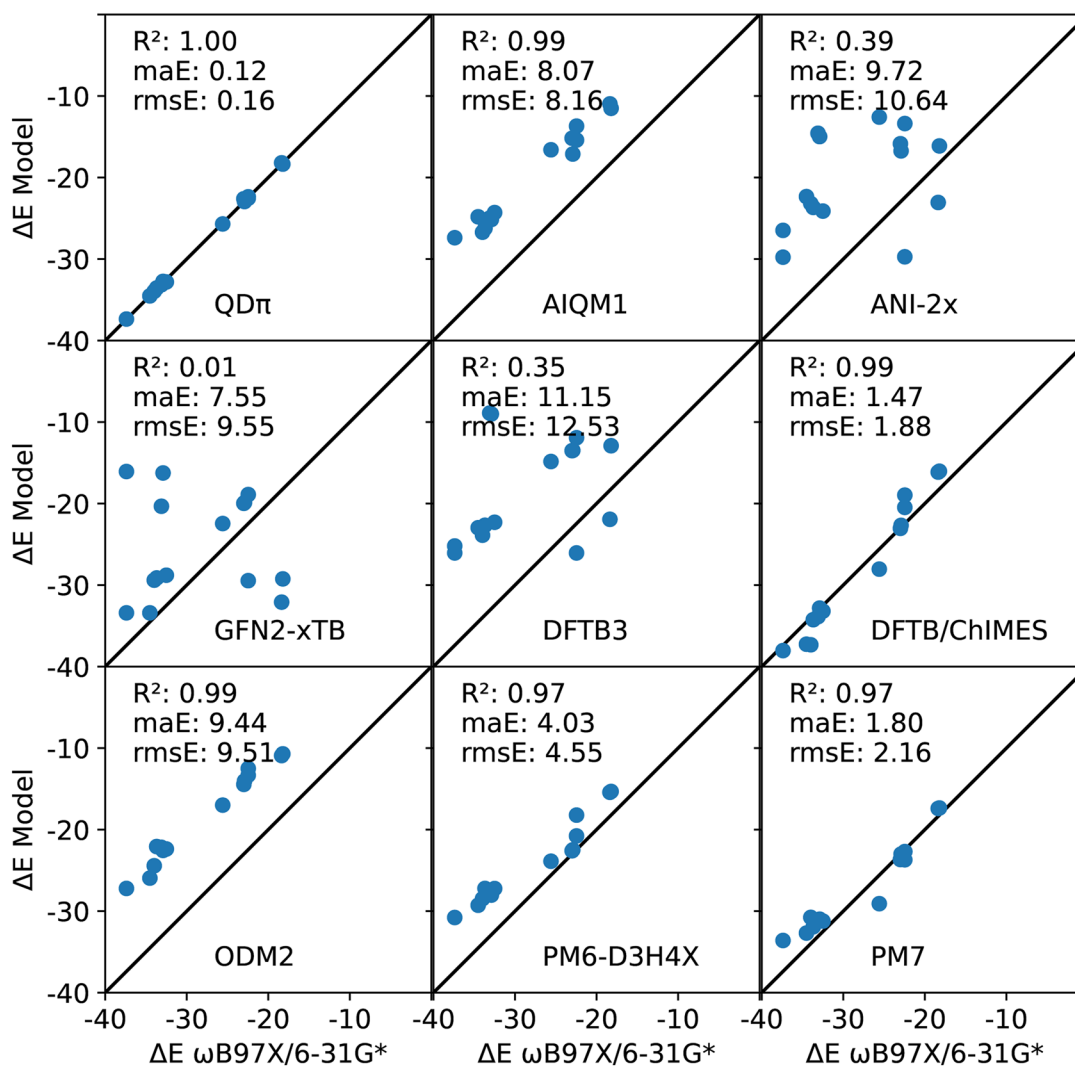
**Table II.** Hydrogen bond complex energies from  $\omega$ B97X/6-31G\* and model errors (kcal/mol) for the artificially expanded genetic information system (AEGIS) base pair dataset<sup>77,85,86</sup> with Leontis and Westhof symbols used for the classification of nucleic acid base pairs,<sup>139–141</sup> including complexes that involve alternative tautomers and protonation states.<sup>a</sup>

Complex	$\omega$ B97X $\Delta E$	QM/ $\Delta$ -MLP or MLP			DFTB			NDDO		
		QD $\pi$ Err	AIQM1 Err	ANI-2x Err	GFN2 Err	DFTB3 Err	ChIMES Err	ODM2 Err	D3H4X Err	PM7 Err
C  G	=32.90	0.16	7.75	9.84	3.66	10.98	0.08	10.35	4.87	1.92
T  A	=18.22	=0.14	6.71	3.65	2.15	9.33	2.19	7.51	2.90	0.88
U  A	=18.36	0.17	7.39	3.39	2.12	9.35	2.22	7.45	2.94	0.96
S  B	=37.40	0.03	10.01	7.67	3.99	11.34	=0.64	10.18	6.62	3.80
V  J	=34.52	=0.00	9.70	8.04	2.43	9.32	=2.73	8.57	5.25	1.83
K  X	=22.46	=0.08	7.06	6.34	2.13	9.56	1.99	9.14	1.69	=1.24
Z  P	=33.11	=0.05	8.13	10.77	3.65	10.14	=0.76	10.93	5.50	1.83
B  G	=32.50	=0.31	8.20	8.38	3.72	10.21	=0.69	10.13	5.28	1.28
B  P	=33.68	0.13	7.45	10.02	4.56	11.02	=0.56	11.61	6.47	1.77
B  G	=22.46	0.09	8.76	9.08	3.56	10.54	3.50	9.95	4.24	=0.24
G  T*	=33.99	0.03	7.28	10.79	4.59	10.11	=3.34	9.56	5.51	3.22
G*  T	=23.00	0.39	7.82	7.14	3.03	9.49	=0.04	8.55	0.41	=0.67
B*  G	=25.59	=0.11	8.99	13.00	3.14	10.75	=2.45	8.60	1.71	=3.49
T  B*	=22.92	=0.03	5.80	6.18	3.04	9.44	0.24	8.88	0.43	=0.08
K <sup>+</sup>  X <sup>-</sup>	=144.48	0.10	12.53	79.73	14.30	18.59	8.26	15.79	15.52	15.29
Z <sup>-</sup>  G	=43.33	=0.07	12.36	11.58	6.69	15.35	3.42	12.16	4.25	1.90
C  P <sup>+</sup>	=47.17	0.04	6.77	24.45	3.38	13.41	0.87	9.41	7.35	4.82
maE	...	0.11	8.46	14.17	4.25	11.22	1.99	10.08	4.87	2.77
rmsE	...	0.15	8.66	22.51	5.09	11.49	2.83	10.26	5.97	4.44

<sup>a</sup>Models and datasets are described in Sec. II. An illustration of each of the complexes is provided in Fig. 4. Complexes include adenine (A), cytosine (C), guanine (G), thymine (T), uracil (U), isoguanine (B), isocytosine (S), 6-amino-5-nitropyridin-2-one (Z), 2-aminoimidazo[1,2-a][1,3,5]triazin-4(1H)-one (P), imidazo[1,2-a]-1,3,5-triazine-2(8H)-4(3H)-dione (X), 2,4-diaminopyrimidine (K), 4-aminoimidazo[1,2-a][1,3,5]triazin-2(8H)-one (J), and 6-amino-3-methylpyridin-2(1H)-one (V).<sup>86,142</sup> The "\*" symbol refers to tautomeric form, and the "+" and "-" symbols refer to the positive and negative charge.



● AEGIS Hydrogen Bonded Base Pair Dataset



**FIG. 5.** Relation between hydrogen bond complex energies calculated by  $\omega$ B97X/6-31G\* and QD $\pi$ , AIQM1, ANI-2x, GFN2-xTB, DFTB3, DFTB/ChIMES, ODM2, PM6-D3H4X, and PM7 for the artificially expanded genetic information system (AEGIS) base pair dataset,<sup>86</sup> including complexes that involve alternative tautomers and protonation states. Illustrations of each of the complexes are provided in Fig. 4. The three base pairs that involve ionized nucleobases are excluded from the regression as they have much larger binding energy values that would artificially skew the correlation.

Table II, and the neutral base pairs are illustrated in Fig. 5. Overall, the QD $\pi$  model gives excellent agreement with the  $\omega$ B97X/6-31G\* reference level over the entire set with  $\Delta E$  maE of 0.11 kcal/mol and a maximum error of 0.39 kcal/mol for G\*  $\bullet$  T. The DFTB/ChIMES model has the next lowest error (maE 1.99 kcal/mol), followed by PM7 (maE 2.77 kcal/mol), GFN2-xTB (maE 4.25 kcal/mol), and PM6-D3H4X (maE 4.87 kcal/mol). The remainder of the models have maE values in excess of 8 kcal/mol. The ANI-2x model has a large maE value (14.17 kcal/mol), but the errors are dominated by base pairs involving ionized nucleobases

that range from 11.58 to 79.73 kcal/mol, whereas the range of errors for neutral base pairs is much smaller (3.39–13.00 kcal/mol; maE of the neutral subset is 9.72 kcal/mol).

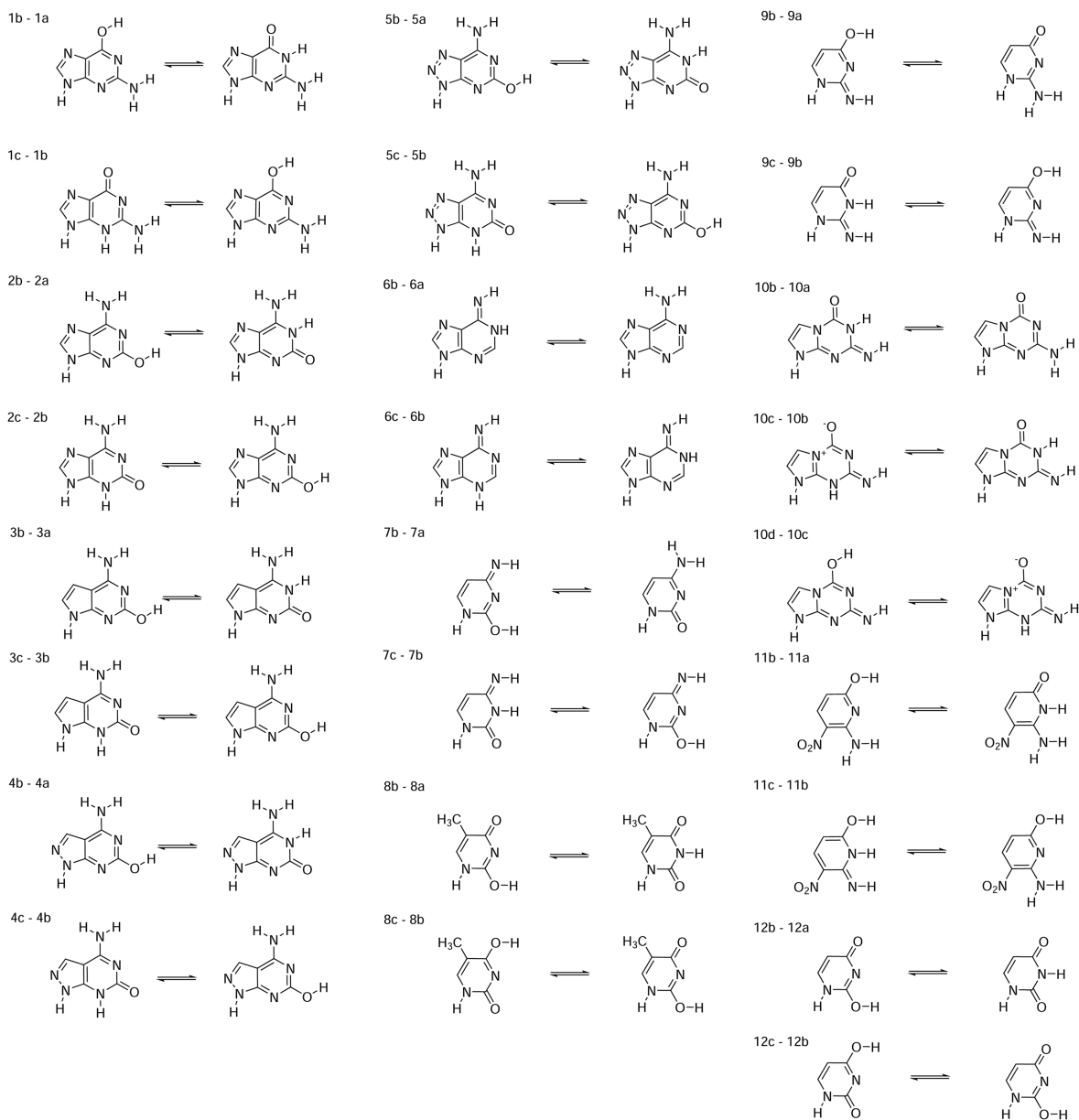
Examination of the correlation of hydrogen complex energies for neutral nucleobases reveals that QD $\pi$  has the highest correlation ( $R^2$  value of 0.999), followed by DFTB/ChIMES, AIQM1, and ODM2 with  $R^2$  values of 0.99. Whereas DFTB/ChIMES is well aligned with the reference data, the ODM2 and related AIQM1 models have values that have been systematically shifted to lower  $\Delta E$  values. Both PM7 and PM6-D3H4X models show

impressive correlation ( $R^2$  values of 0.97) and low maE values (1.80 and 4.03 kcal/mol, respectively) for complexes of these neutral nucleobases.

#### D. Comparison of tautomer energies for natural and artificial nucleic acids

The artificially expanded genetic information system (AEGIS) dataset also exhibits a rich set of tautomeric forms that have been

extensively studied with computational methods.<sup>77,85–87</sup> These tautomeric pairs are illustrated in Fig. 6, and their  $\Delta E$  values are listed in Table III and illustrated in Fig. 7. Overall, both QD $\pi$  and AIQM1 give excellent agreement with the  $\omega$ B97X/6-31G<sup>+</sup> reference values, with  $\Delta E$  maE values of 0.71 and 0.77 kcal/mol, respectively, and high correlation ( $R^2$  value of 0.99). The ANI-2x is the next most accurate, but with errors roughly twice as large (maE 1.41 kcal/mol) and ( $R^2$  value of 0.97). The DFTB/ChIMES and GFN2-xTB models have considerably higher errors (maE values of 2.20 and



**FIG. 6.** Structures for the artificially expanded genetic information system (AEGIS) tautomer dataset.<sup>77,85</sup> Guanine derivatives (1–5, 2: nucleobase code B), codes 6: A, 7: C, 8: T, 9: S, 10: P, and 11: Z.

**TABLE III.** Tautomerization energies from  $\omega$ B97X/6-31G\* and model errors (kcal/mol) for the artificially expanded genetic information system (AEGIS) tautomer dataset.<sup>a</sup>

Tautomer pair	$\omega$ B97X $\Delta E$	QM/ $\Delta$ -MLP or MLP			DFTB			NDDO		
		QD $\pi$ Err	AIQM1 Err	ANI-2x Err	GFN2 Err	DFTB3 Err	ChIMES Err	ODM2 Err	D3H4X Err	PM7 Err
1b-1a	2.43	=0.36	=1.07	=1.35	=2.75*	=6.25*	=2.38	3.48	6.60	2.67
1c-1b	17.39	0.36	0.40	0.83	=0.72	2.68	2.58	=6.65	=11.01	=7.46
2b-2a	=5.41	=1.12	0.13	=0.29	2.81	=4.39	=1.02	6.31*	12.62*	8.85*
2c-2b	4.95	1.22	=0.74	2.30	=2.61	4.74	1.54	=7.50*	=14.46*	=10.85*
3b-3a	=6.32	=0.28	=0.16	0.89	3.53	=3.90	=0.77	6.39*	12.95*	9.24*
3c-3b	4.05	=0.15	=0.15	3.05	=2.48	4.44	1.02	=7.07*	=13.56*	=9.97*
4b-4a	=6.81	0.69	0.37	0.58	3.60	=4.06	=0.94	7.06*	13.49*	10.01*
4c-4b	2.76	0.06	=0.82	1.14	=1.90	5.47	1.98	=6.48*	=12.55*	=8.93*
5b-5a	=6.12	=0.41	0.66	=0.28	3.07	=4.60	=1.43	7.34*	13.36*	9.71*
5c-5b	3.23	=0.57	=0.57	0.49	=1.92	5.97	2.56	=6.59*	=13.03*	=9.41*
6b-6a	12.24	=0.01	=0.20	=2.56	=1.85	=3.22	=3.63	3.76	=4.94	=0.92
6c-6b	20.12	=0.10	=0.35	2.63	=5.80	=3.10	=0.26	=4.60	=8.47	=8.48
7b-7a	20.15	1.40	=0.63	=1.39	=4.87	=5.71	=5.32	5.63	=0.67	=0.43
7c-7b	=19.17	=0.96	0.94	=0.22	4.79	6.93	1.64	=0.25	=0.01	2.86
8b-8a	21.32	0.84	=1.79	=0.08	=5.43	=10.25	=2.12	1.57	3.71	1.24
8c-8b	=6.16	0.31	=0.43	1.10	=1.52	0.18	1.28	=1.34	=3.27	=3.05
9b-9a	5.42	0.77	=0.89	=0.45	=5.39	=7.41*	=4.94	4.64	=3.31	=3.53
9c-9b	=10.02	=0.74	0.40	0.29	3.20	6.10	1.29	=0.76	=1.38	1.91
10b-10a	7.95	2.11	=0.57	=2.50	=1.83	0.77	=3.26	3.28	=3.84	=1.36
10c-10b	22.20	=2.92	=2.92	=5.08	=6.40	=15.21	0.78	=9.14	=11.36	=9.26
10d-10c	4.01	0.04	1.54	2.86	=1.71	5.13	=3.27	9.90	12.36	7.71
11b-11a	=0.86	=0.63	=0.19	=0.36	=0.47	=4.04	=0.58	3.42*	7.69*	4.68*
11c-11b	24.35	=0.84	=0.90	=4.13	=4.19	=5.43	=6.39	1.67	=8.47	=6.13
12b-12a	22.00	0.41	=1.80	0.14	=5.51	=10.37	=2.16	1.34	3.73	1.31
12c-12b	=7.79	0.54	0.54	0.29	=0.62	0.87	1.78	=0.98	=2.70	=2.59
maE	...	0.71	0.77	1.41	3.16	5.25	2.20	4.69	7.98	5.70
rmsE	...	0.97	1.00	1.94	3.59	6.12	2.67	5.42	9.28	6.72

<sup>a</sup>Models and datasets are described in Sec. II. Illustrations of each of the tautomerization reactions are provided in Fig. 7. Errors corresponding to wrong prediction of more stable tautomer are indicated by an asterisk (\*).

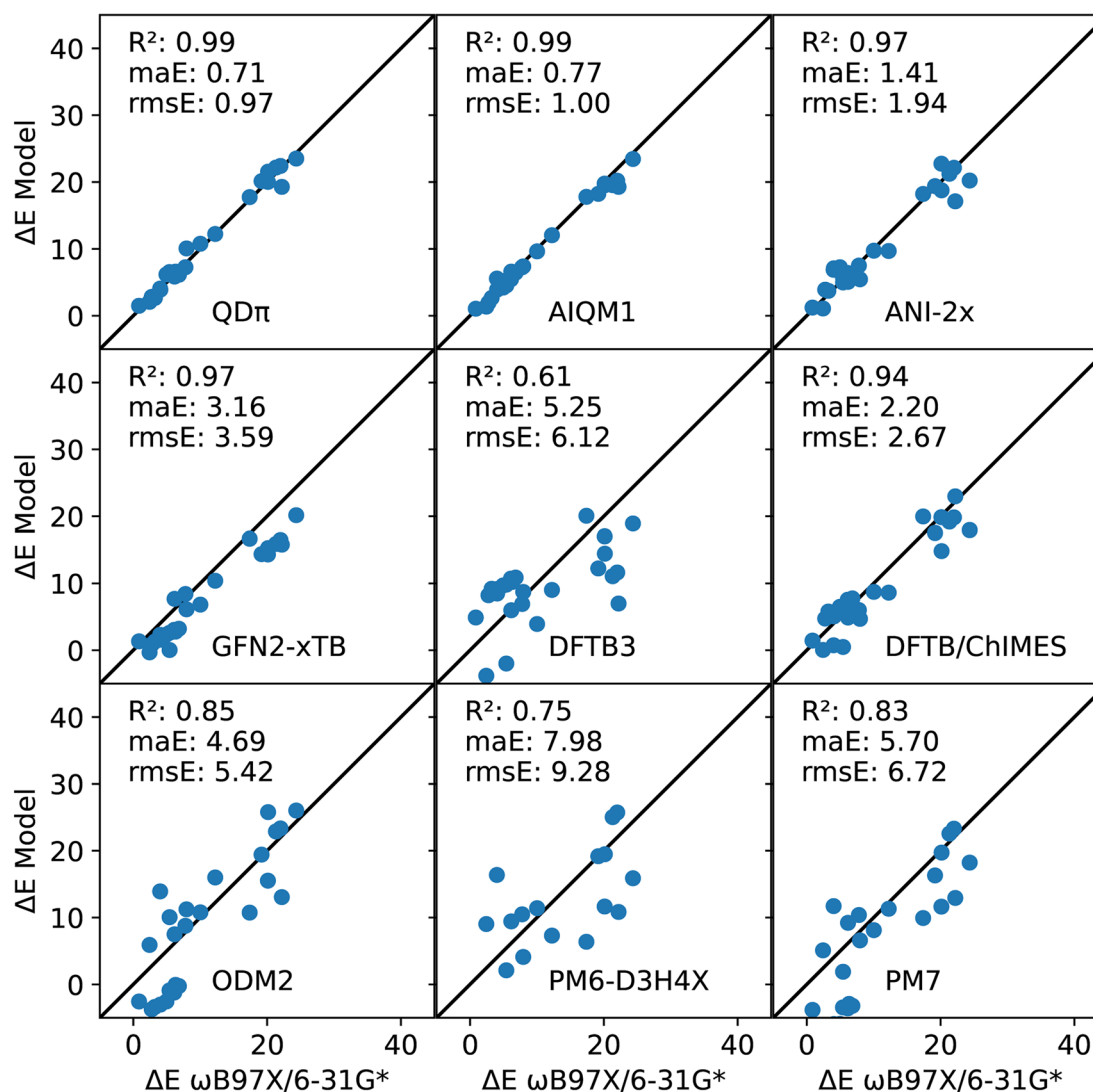
3.16 kcal/mol, respectively), but maintain high correlation with the reference values ( $R^2$  value of 0.97), whereas DFTB3 and ODM2 have larger errors (maE values of 5.25 and 4.69 kcal/mol, respectively), and lower correlation ( $R^2$  values of 0.61 and 0.85, respectively). The largest errors occur for PM7 and PM6-D3H4X (maE values of 5.70 and 7.98 kcal/mol, respectively).

It has been estimated that 30% of the compounds in vendor databases and 21% of the compounds in drug databases have potential tautomers.<sup>18,19</sup> For drug discovery applications, it is thus vitally important to be able to model alternative tautomer forms, discern which forms are relevant for ligand-protein binding, and if binding induces a change in tautomer state, to quantitatively determine the tautomerization energy contribution to binding with sub-kcal/mol accuracy. In some cases, the semiempirical QM models incorrectly predict the lowest energy tautomer (one case for GFN2-xTB, two cases for DFTB3, and nine cases each for ODM2, PM6-D3H4X, and PM7). For the models compared here, only QD $\pi$  and AIQM1 are able to achieve the requisite accuracy for quantitative prediction of ligand-protein binding applications.

## E. Comparison of protonation energies for common general acids and bases

Modeling protonation states is important for drug discovery applications as it has been estimated that up to 95% of drug molecules contain ionizable groups<sup>18</sup> (~75% weak bases and ~20% weak acids<sup>20,21</sup>), and protonation states can sometimes change upon ligand binding. Hence, quantitatively accurate modeling of protonation/deprotonation events at these ionizable sites is critical to obtain high predictive capability. As an illustrative set of examples, we examine simple model systems that mimic the acid/base chemistry associated with RNA cleavage reactions catalyzed by small nucleolytic RNA enzymes (ribozymes) and protein enzymes (ribonucleases).<sup>100</sup> In these reactions, the 2'OH of an RNA nucleotide, modeled by the secondary alcohol isopropanol (iPrOH), becomes activated (deprotonated) by a general base that in ribozymes is often an ionized (deprotonated) guanine residue (G:N $^-$ ), and in RNase A<sup>143-145</sup> is generally believed to be a histidine (His:N $^e$ ) although it has been speculated that a neutral lysine (Lys:NH $_2$ ) might also be

## ● AEGIS Tautomer Dataset



**FIG. 7.** Relation between tautomerization energies calculated by  $\omega$ B97X/6-31G\* and QD $\pi$ , AIQM1, ANI-2x, GFN2-xTB, DFTB3, DFTB/ChIMES, ODM2, PM6-D3H4X, and PM7 for the artificially expanded genetic information system (AEGIS) tautomer dataset. An illustration of each of the complexes is provided in Fig. 6. In the regression plot shown, the sign convention (direction of the tautomer reaction) is chosen such that the reference  $\Delta E$  value is positive (this is performed to circumvent “spreading out” of the data and artificially inflating the correlation).

capable. The activated nucleophile then attacks the scissile phosphate, passing through a pentavalent transition state, followed by the departure of the 5' O leaving group (modeled by the primary alkoxide ethoxide ( $\text{EtO}^-$ ) with the assistance of a general acid that in ribozymes can be either a protonated adenine at the N1 or N3 positions ( $\text{A:N}_1\text{H}^+$  and  $\text{A:N}_3\text{H}^+$ , respectively) or an ionized (protonated) cytosine ( $\text{C:N}_3\text{H}^+$ ), and in RNase A is a protonated histidine ( $\text{His:N}_\epsilon\text{H}^+$ ).

Table IV lists relative protonation/deprotonation reactions that model general acid/base events in RNA cleavage reactions.<sup>100</sup> Overall, QD $\pi$  performs extremely well, with  $\Delta E$  maE of 0.50 kcal/mol. Of the semiempirical QM methods, GFN2-xTB is the least inaccurate (maE value of 5.94 kcal/mol), followed by PM7 (6.97 kcal/mol), with other models notably higher (maE values ranging from 9.12 to 14.67 kcal/mol). As mentioned earlier, the ANI-2x model was not designed to handle ions; it produces errors on the order of



**TABLE IV.** Selected relative protonation/deprotonation energies from  $\omega$ B97X/6-31G\* and model error (kcal/mol) relevant to acid/base catalysis in RNA cleavage reactions.<sup>a</sup>

Protonation pair	$\omega$ B97X $\Delta E$	QM/ $\Delta$ -MLP or MLP			DFTB			NDDO		
		QD $\pi$ Err	AIQM1 Err	ANI-2x Err	GFN2 Err	DFTB3 Err	ChIMES Err	ODM2 Err	D3H4X Err	PM7 Err
[Lys:NH <sub>2</sub> ,iPrOH]	167.76	0.00	=0.64	=115.04	0.04	6.11	=-6.87	=15.24	=13.15	=10.15
[His:N <sub><math>\epsilon</math></sub> ,iPrOH]	158.33	0.08	=9.22	=126.62	=7.02	=11.33	=17.47	=18.74	=12.96	-6.19
[EtO <sup>-</sup> ,His:N <sub><math>\epsilon</math></sub> H <sup>+</sup> ]	=160.25	=0.02	12.82	137.70	9.66	11.71	18.09	21.92	12.24	5.36
[G:N <sub>1</sub> ,iPrOH]	43.06	=1.11	=1.15	=28.62	-2.69	=8.63	=11.17	=10.84	0.89	5.45
[EtO <sup>-</sup> ,A:N <sub>1</sub> H <sup>+</sup> ]	=165.06	1.25	12.94	137.24	10.02	15.21	23.15	20.74	8.07	1.25
[EtO <sup>-</sup> ,A:N <sub>3</sub> H <sup>+</sup> ]	=190.89	1.21	12.88	143.42	11.40	16.00	24.17	19.43	17.97	8.39
[EtO <sup>-</sup> ,C:N <sub>3</sub> H <sup>+</sup> ]	=160.33	0.89	12.78	145.20	6.58	4.66	16.93	20.19	7.03	2.22
[G:N <sub>1</sub> ,A:N <sub>1</sub> H <sup>+</sup> ]	=120.07	0.08	8.19	97.55	4.69	6.20	11.36	6.73	9.69	7.53
[G:N <sub>1</sub> ,A:N <sub>3</sub> H <sup>+</sup> ]	=145.91	0.04	8.12	103.73	6.07	6.99	12.38	5.41	19.58	14.67
[G:N <sub>1</sub> ,C:N <sub>3</sub> H <sup>+</sup> ]	=115.34	=0.27	8.03	105.50	1.25	=4.35	5.14	6.17	8.64	8.51
maE	...	0.50	8.68	114.06	5.94	9.12	14.67	14.54	11.02	6.97
rmsE	...	0.72	9.72	118.72	6.96	9.96	15.87	15.84	12.17	7.88

<sup>a</sup>Models and datasets are described in Sec. II. Protonation pairs are written in the general form as follows: [B, A]: B + A  $\rightarrow$  BH<sup>+</sup> + A<sup>-</sup>, or [B<sup>-</sup>, AH<sup>+</sup>]: B<sup>-</sup> + AH<sup>+</sup>  $\rightarrow$  BH + A. Here, B/BH<sup>+</sup> and B<sup>-</sup>/BH are the base/conjugate acid pairs and A/A<sup>-</sup> and AH<sup>+</sup>/A are the acid/conjugate base pairs. These are model systems for general acid and base catalysis in RNA cleavage reactions by small nucleolytic ribozymes and ribonucleases.<sup>100</sup> The molecules indicated are isopropanol (iPrOH), ethoxide (EtO<sup>-</sup>), neutral lysine (Lys:NH<sub>2</sub>), neutral histidine (His:N <sub>$\epsilon$</sub> ), protonated histidine (His:N <sub>$\epsilon$</sub> H<sup>+</sup>), deprotonated guanine at the N1 position (G:N<sub>1</sub><sup>-</sup>), protonated adenine at the N1 and N3 positions (A:N<sub>1</sub>H<sup>+</sup> and A:N<sub>3</sub>H<sup>+</sup>), and protonated cytosine at the N3 position (C:N<sub>3</sub>H<sup>+</sup>).

100 kcal/mol. The AIQM1 model is greatly improved with respect to ANI-2x and ODM2 (the base QM model). The QD $\pi$   $\Delta E$  maE value is dominated by large positive errors involving the ethoxide and protonated nucleobases (0.89–1.25 kcal/mol). The ethoxide anion is a primary alkoxide that is only marginally stable in the gas phase, and thus especially challenging. The QD $\pi$  model is by far the most accurate for protonation/deprotonation energies. It is a promising candidate for use in drug discovery applications.

#### IV. CONCLUSION

We have compared the performance of several NDDO-based (MNDO/d, AM1, PM6, PM6-D3H4X, PM7, and ODM2) and density-functional tight-binding based (DFTB3, DFTB/ChIMES, GFN1-xTB, and GFN2-xTB) semiempirical models with pure machine learning potentials (ANI-1x and ANI-2x) and hybrid quantum mechanical/machine learning potentials (AIQM1 and QD $\pi$ ). We examine broad datasets computed at a consistent  $\omega$ B97X/6-31G\* level of theory that includes conformational energies, intermolecular interactions, tautomers, and protonation states. The methods were further compared against the AEGIS dataset and acid/base chemistry relevant for RNA cleavage reactions catalyzed by small nucleolytic ribozymes and ribonucleases. Overall, the recently developed QD $\pi$  model performs exceptionally well across all datasets, with especially high accuracy for tautomers and protonation states relevant to drug discovery. The AIQM1 model also has impressive performance in many cases, including tautomerization energies. All other methods examined have various strengths and weaknesses, but none have the broad range of quantitative accuracy

of the QD $\pi$  model for the data examined. Taken together, this suggests that QM/ $\Delta$ -MLPs such as QD $\pi$  and AIQM1 have considerable promise as universal force fields for drug discovery applications.

#### SUPPLEMENTARY MATERIAL

See the [supplementary material](#) for relative energies for the minima and transition states of the alanine dipeptide, ibuprofen, and ketorolac.

#### ACKNOWLEDGMENTS

The authors are grateful for the financial support provided by the National Institutes of Health (Grant No. GM107485 to D.M.Y.) and the National Science Foundation (Grant No. 2209718 to D.M.Y.). 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 the National Science Foundation under Grant No. ACI-1548562 (supercomputer Expanse at SDSC through allocation under Grant No. CHE190067), and by the Texas Advanced Computing Center (TACC) at the University of Texas at Austin (supercomputer Longhorn through allocation under Grant No. CHE20002).

#### AUTHOR DECLARATIONS

##### Conflict of Interest

The authors have no conflicts to disclose.

## Author Contributions

**Jinze Zeng** (曾晋哲): Data curation (equal); Software (equal); Visualization (equal); Writing – original draft (equal); Writing – review & editing (equal). **Yujun Tao** (陶玉君): Data curation (equal); Visualization (equal); Writing – review & editing (equal). **Timothy J. Giese**: Writing – review & editing (equal). **Darrin M. York**: Formal analysis (equal); Funding acquisition (equal); Project administration (equal); Resources (equal); Supervision (equal); Writing – original draft (equal); Writing – review & editing (equal).

## DATA AVAILABILITY

QD $\pi$ -v1.0 is openly available in our GitLab repository at <https://gitlab.com/RutgersLBSR/qdpi>, which was previously released.<sup>77</sup> The data that support the findings of this study are available from the corresponding author upon reasonable request.

## REFERENCES

- 1 T.-S. Lee, B. K. Allen, T. J. Giese, Z. Guo, P. Li, C. Lin, T. D. McGee, Jr., D. A. Pearlman, B. K. Radak, Y. Tao, H.-C. Tsai, H. Xu, W. Sherman, and D. M. York, “Alchemical binding free energy calculations in AMBER20: Advances and best practices for drug discovery,” *J. Chem. Inf. Model.* **60**, 5595–5623 (2020).
- 2 W. L. Jorgensen, “Efficient drug lead discovery and optimization,” *Acc. Chem. Res.* **42**, 724–733 (2009).
- 3 D. J. Cole, J. T. Horton, L. Nelson, and V. Kurdekar, “The future of force fields in computer-aided drug design,” *Future Med. Chem.* **11**, 2359–2363 (2019).
- 4 A. D. MacKerell, Jr., “Empirical force fields for biological macromolecules: Overview and issues,” *J. Comput. Chem.* **25**, 1584–1604 (2004).
- 5 K. Lindorff-Larsen, P. Maragakis, S. Piana, M. P. Eastwood, R. O. Dror, and D. E. Shaw, “Systematic validation of protein force fields against experimental data,” *PLoS One* **7**, e32131 (2012).
- 6 S. Piana, J. L. Klepeis, and D. E. Shaw, “Assessing the accuracy of physical models used in protein-folding simulations: Quantitative evidence from long molecular dynamics simulations,” *Curr. Opin. Struct. Biol.* **24**, 98–105 (2014).
- 7 W. L. Jorgensen, J. Chandrasekhar, J. D. Madura, R. W. Impey, and M. L. Klein, “Comparison of simple potential functions for simulating liquid water,” *J. Chem. Phys.* **79**, 926–935 (1983).
- 8 H. W. Horn, W. C. Swope, J. W. Pitera, J. D. Madura, T. J. Dick, G. L. Hura, and T. Head-Gordon, “Development of an improved four-site water model for biomolecular simulations: TIP4P-Ew,” *J. Chem. Phys.* **120**, 9665–9678 (2004).
- 9 Y. Wu, H. L. Tepper, and G. A. Voth, “Flexible simple point-charge water model with improved liquid-state properties,” *J. Chem. Phys.* **124**, 024503 (2006).
- 10 S. Izadi, R. Anandakrishnan, and A. V. Onufriev, “Building water models: A different approach,” *J. Phys. Chem. Lett.* **5**, 3863–3871 (2014).
- 11 S. Izadi and A. V. Onufriev, “Accuracy limit of rigid 3-point water models,” *J. Chem. Phys.* **145**, 074501–074510 (2016).
- 12 I. S. Joung and T. E. Cheatham III, “Molecular dynamics simulations of the dynamic and energetic properties of alkali and halide ions using water-model-specific ion parameters,” *J. Phys. Chem. B* **113**, 13279–13290 (2009).
- 13 P. Li, B. P. Roberts, D. K. Chakravorty, and K. M. Merz, Jr., “Rational design of particle mesh Ewald compatible Lennard-Jones parameters for +2 metal cations in explicit solvent,” *J. Chem. Theory Comput.* **9**, 2733–2748 (2013).
- 14 P. Li and K. M. Merz, Jr., “Taking into account the ion-induced dipole interaction in the nonbonded model of ions,” *J. Chem. Theory Comput.* **10**, 289–297 (2014).
- 15 P. Li and K. M. Merz, “Metal ion modeling using classical mechanics,” *Chem. Rev.* **117**, 1564–1686 (2017).
- 16 P. E. M. Lopes, O. Guvench, and A. D. MacKerell, Jr., “Current status of protein force fields for molecular dynamics simulations,” *Methods Mol. Biol.* **1215**, 47–71 (2015).
- 17 C. Tian, K. Kasavajhala, K. A. A. Belfon, L. Raguette, H. Huang, A. N. Miguels, J. Bickel, Y. Wang, J. Pincay, Q. Wu, and C. Simmerling, “ff19SB: Amino-acid-specific protein backbone parameters trained against quantum mechanics energy surfaces in solution,” *J. Chem. Theory Comput.* **16**, 528–552 (2020).
- 18 Y. C. Martin, “Experimental and pK<sub>a</sub> prediction aspects of tautomerism of drug-like molecules,” *Drug Discovery Today: Technol.* **27**, 59–64 (2018).
- 19 F. Milletti and A. Vulpetti, “Tautomer preference in PDB complexes and its impact on structure-based drug discovery,” *J. Chem. Inf. Model.* **50**, 1062–1074 (2010).
- 20 J. I. Wells, *Pharmaceutical Preformulation: The Physicochemical Properties of Drug Substances* (E. Horwood, Chichester, UK, 1988).
- 21 C. D. Navo and G. Jiménez-Osés, “Computer prediction of pK<sub>a</sub> values in small molecules and proteins,” *ACS Med. Chem. Lett.* **12**, 1624–1628 (2021).
- 22 W. Thiel, “Semiempirical quantum–chemical methods,” *Wiley Interdiscip. Rev.: Comput. Mol. Sci.* **4**, 145–157 (2014).
- 23 K. Kříž and J. Řezáč, “Benchmarking of semiempirical quantum-mechanical methods on systems relevant to computer-aided drug design,” *J. Chem. Inf. Model.* **60**, 1453–1460 (2020).
- 24 V. Khanna and S. Ranganathan, “Physicochemical property space distribution among human metabolites, drugs and toxins,” *BMC Bioinf.* **10**, S10 (2009).
- 25 T. Darden, D. York, and L. Pedersen, “Particle mesh Ewald: An N log(N) method for Ewald sums in large systems,” *J. Chem. Phys.* **98**, 10089–10092 (1993).
- 26 K. Nam, J. Gao, and D. M. York, “An efficient linear-scaling Ewald method for long-range electrostatic interactions in combined QM/MM calculations,” *J. Chem. Theory Comput.* **1**, 2–13 (2005).
- 27 T. J. Giese, M. T. Panteva, H. Chen, and D. M. York, “Multipolar Ewald methods, 2: Applications using a quantum mechanical force field,” *J. Chem. Theory Comput.* **11**, 451–461 (2015).
- 28 T. J. Giese and D. M. York, “Ambient-potential composite Ewald method for *ab initio* quantum mechanical/molecular mechanical molecular dynamics simulation,” *J. Chem. Theory Comput.* **12**, 2611–2632 (2016).
- 29 J. T. Margraf, M. Hennemann, and T. Clark, “EMPIRE: A highly parallel semiempirical molecular orbital program: 3: Born-Oppenheimer molecular dynamics,” *J. Mol. Model.* **26**, 43 (2020).
- 30 B. Hourahine, B. Aradi, V. Blum, F. Bonafé, A. Buccheri, C. Camacho, C. Cevallos, M. Y. Deshayé, T. Dumitrică, A. Dominguez, S. Ehlert, M. Elstner, T. van der Heide, J. Hermann, S. Irlé, J. J. Kranz, C. Köhler, T. Kowalczyk, T. Kubař, I. S. Lee, V. Lutsker, R. J. Maurer, S. K. Min, I. Mitchell, C. Negre, T. A. Niehaus, A. M. N. Niklasson, A. J. Page, A. Pecchia, G. Penazzi, M. P. Persson, J. Řezáč, C. G. Sánchez, M. Sternberg, M. Stöhr, F. Stuckenberg, A. Tkatchenko, V. W.-z. Yu, and T. Frauenheim, “DFTB+, a software package for efficient approximate density functional theory based atomistic simulations,” *J. Chem. Phys.* **152**, 124101 (2020).
- 31 T. J. Giese and D. M. York, “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.* **15**, 5543–5562 (2019).
- 32 F. L. Kearns, P. S. Hudson, H. L. Woodcock, and S. Boresch, “Computing converged free energy differences between levels of theory via nonequilibrium work methods: Challenges and opportunities,” *J. Comput. Chem.* **38**, 1376–1388 (2017).
- 33 S. Boresch and H. L. Woodcock, “Convergence of single-step free energy perturbation,” *Mol. Phys.* **115**, 1200–1213 (2017).
- 34 P. S. Hudson, F. Aviat, R. Meana-Pañeda, L. Warrensford, B. C. Pollard, S. Prasad, M. R. Jones, H. L. Woodcock, and B. R. Brooks, “Obtaining QM/MM binding free energies in the SAMPL8 drugs of abuse challenge: Indirect approaches,” *J. Comput.-Aided Mol. Des.* **36**, 263–277 (2022).
- 35 A. Schöller, F. Kearns, H. L. Woodcock, and S. Boresch, “Optimizing the calculation of free energy differences in nonequilibrium work SQM/MM switching simulations,” *J. Phys. Chem. B* **126**, 2798–2811 (2022).
- 36 G. Arya, “Models for recovering the energy landscape of conformational transitions from single-molecule pulling experiments,” *Mol. Simul.* **42**, 1102–1115 (2016).
- 37 A. N. Naganathan, U. Doshi, and V. Muñoz, “Protein folding kinetics: Barrier effects in chemical and thermal denaturation experiments,” *J. Am. Chem. Soc.* **129**, 5673–5682 (2007).

- <sup>38</sup>J. Basran, S. Patel, M. J. Sutcliffe, and N. S. Scrutton, "Importance of barrier shape in enzyme-catalyzed reactions," *J. Biol. Chem.* **276**, 6234–6242 (2001).
- <sup>39</sup>J. Behler, "Perspective: Machine learning potentials for atomistic simulations," *J. Chem. Phys.* **145**, 170901 (2016).
- <sup>40</sup>K. T. Butler, D. W. Davies, H. Cartwright, O. Isayev, and A. Walsh, "Machine learning for molecular and materials science," *Nature* **559**, 547–555 (2018).
- <sup>41</sup>F. Noé, A. Tkatchenko, K.-R. Müller, and C. Clementi, "Machine learning for molecular simulation," *Annu. Rev. Phys. Chem.* **71**, 361–390 (2020).
- <sup>42</sup>M. Pinheiro, Jr., F. Ge, N. Ferré, P. O. Dral, and M. Barbatti, "Choosing the right molecular machine learning potential," *Chem. Sci.* **12**, 14396–14413 (2021).
- <sup>43</sup>S. Manzhos and T. Carrington, Jr., "Neural network potential energy surfaces for small molecules and reactions," *Chem. Rev.* **121**, 10187–10217 (2021).
- <sup>44</sup>J. Zeng, L. Cao, and T. Zhu, "Neural network potentials," in *Quantum Chemistry in the Age of Machine Learning*, edited by P. O. Dral (Elsevier, 2022) Chap. 12, pp. 279–294.
- <sup>45</sup>X. Pan, J. Yang, R. Van, E. Epifanovsky, J. Ho, J. Huang, J. Pu, Y. Mei, K. Nam, and Y. Shao, "Machine-learning-assisted free energy simulation of solution-phase and enzyme reactions," *J. Chem. Theory Comput.* **17**, 5745–5758 (2021).
- <sup>46</sup>P. Zheng, R. Zubatyuk, W. Wu, O. Isayev, and P. O. Dral, "Artificial intelligence-enhanced quantum chemical method with broad applicability," *Nat. Commun.* **12**, 7022 (2021).
- <sup>47</sup>J. Zeng, T. J. Giese, Ş. Ekesan, and D. M. York, "Development of range-corrected deep learning potentials for fast, accurate quantum mechanical/molecular mechanical simulations of chemical reactions in solution," *J. Chem. Theory Comput.* **17**, 6993–7009 (2021).
- <sup>48</sup>T. J. Giese, J. Zeng, Ş. Ekesan, and D. M. York, "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.* **18**, 4304–4317 (2022).
- <sup>49</sup>C. L. Gómez-Flores, D. Maag, M. Kansari, V.-Q. Vuong, S. Irle, F. Gräter, T. Kubař, and M. Elstner, "Accurate free energies for complex condensed-phase reactions using an artificial neural network corrected DFTB/MM methodology," *J. Chem. Theory Comput.* **18**, 1213–1226 (2022).
- <sup>50</sup>J. Böser, T. Kubař, M. Elstner, and D. Maag, "Reduction pathway of glutaredoxin 1 investigated with QM/MM molecular dynamics using a neural network correction," *J. Chem. Phys.* **157**, 154104 (2022).
- <sup>51</sup>P. O. Dral, T. Zubatiuk, and B.-X. Xue, "Learning from multiple quantum chemical methods:  $\Delta$ -learning, transfer learning, co-kriging, and beyond," in *Quantum Chemistry in the Age of Machine Learning*, edited by P. O. Dral (Elsevier, 2022) Chap. 21, pp. 491–507.
- <sup>52</sup>J. Behler and M. Parrinello, "Generalized neural-network representation of high-dimensional potential-energy surfaces," *Phys. Rev. Lett.* **98**, 146401–146404 (2007).
- <sup>53</sup>A. P. Bartók, M. C. Payne, R. Kondor, and G. Csányi, "Gaussian approximation potentials: The accuracy of quantum mechanics, without the electrons," *Phys. Rev. Lett.* **104**, 136403 (2010).
- <sup>54</sup>J. Behler, "Atom-centered symmetry functions for constructing high-dimensional neural network potentials," *J. Chem. Phys.* **134**, 074106 (2011).
- <sup>55</sup>M. Gastegger, L. Schwiedrzik, M. Bittermann, F. Berzsenyi, and P. Marquetand, "wACSF—Weighted atom-centered symmetry functions as descriptors in machine learning potentials," *J. Chem. Phys.* **148**, 241709 (2018).
- <sup>56</sup>S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller, "Machine learning of accurate energy-conserving molecular force fields," *Sci. Adv.* **3**, 1603015 (2017).
- <sup>57</sup>K. T. Schütt, F. Arbabzadah, S. Chmiela, K. R. Müller, and A. Tkatchenko, "Quantum-chemical insights from deep tensor neural networks," *Nat. Commun.* **8**, 13890 (2017).
- <sup>58</sup>K. T. Schütt, H. E. Sauceda, P.-J. Kindermans, A. Tkatchenko, and K.-R. Müller, "SchNet – A deep learning architecture for molecules and materials," *J. Chem. Phys.* **148**, 241722 (2018).
- <sup>59</sup>X. Chen, M. S. Jørgensen, J. Li, and B. Hammer, "Atomic energies from a convolutional neural network," *J. Chem. Theory Comput.* **14**, 3933–3942 (2018).
- <sup>60</sup>L. Zhang, J. Han, H. Wang, R. Car, and W. E, "Deep potential molecular dynamics: A scalable model with the accuracy of quantum mechanics," *Phys. Rev. Lett.* **120**, 143001 (2018).
- <sup>61</sup>L. Zhang, J. Han, H. Wang, W. Saidi, R. Car, and W. E, "End-to-end symmetry preserving inter-atomic potential energy model for finite and extended systems," in *Advances in Neural Information Processing Systems 31*, edited by S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett (Curran Associates, Inc., 2018), pp. 4436–4446.
- <sup>62</sup>Y. Zhang, C. Hu, and B. Jiang, "Embedded atom neural network potentials: Efficient and accurate machine learning with a physically inspired representation," *J. Phys. Chem. Lett.* **10**, 4962–4967 (2019).
- <sup>63</sup>J. S. Smith, O. Isayev, and A. E. Roitberg, "ANI-1: An extensible neural network potential with DFT accuracy at force field computational cost," *Chem. Sci.* **8**, 3192–3203 (2017).
- <sup>64</sup>O. Unke and M. Meuwly, "PhysNet: A neural network for predicting energies, forces, dipole moments, and partial charges," *J. Chem. Theory Comput.* **15**, 3678–3693 (2019).
- <sup>65</sup>Z. L. Glick, D. P. Metcalf, A. Koutsoukas, S. A. Spronk, D. L. Cheney, and C. D. Sherrill, "AP-Net: An atomic-pairwise neural network for smooth and transferable interaction potentials," *J. Chem. Phys.* **153**, 044112 (2020).
- <sup>66</sup>T. Zubatiuk and O. Isayev, "Development of multimodal machine learning potentials: Toward a physics-aware artificial intelligence," *Acc. Chem. Res.* **54**, 1575–1585 (2021).
- <sup>67</sup>E. R. Khajepasha, J. A. Finkler, T. D. Kühne, and S. A. Ghasemi, "CENT2: Improved charge equilibration via neural network technique," *Phys. Rev. B* **105**, 14261–14261 (2022).
- <sup>68</sup>Y. Hirano, N. Okimoto, S. Fujita, and M. Taiji, "Molecular dynamics study of conformational changes of tankyrase 2 binding subsites upon ligand binding," *ACS Omega* **6**, 17609–17620 (2021).
- <sup>69</sup>P. W. Kenny, "Hydrogen-bond donors in drug design," *J. Med. Chem.* **65**, 14261–14275 (2022).
- <sup>70</sup>H. Yuki, Y. Tanaka, M. Hata, H. Ishikawa, S. Neya, and T. Hoshino, "Implementation of  $\pi$ - $\pi$  interactions in molecular dynamics simulation," *J. Comput. Chem.* **28**, 1091–1099 (2007).
- <sup>71</sup>T. Chen, M. Li, and J. Liu, " $\pi$ - $\pi$  Stacking interaction: A nondestructive and facile means in material engineering for bioapplications," *Cryst. Growth Des.* **18**, 2765–2783 (2018).
- <sup>72</sup>M. Mohebifar, E. R. Johnson, and C. N. Rowley, "Evaluating force-field London dispersion coefficients using the exchange-hole dipole moment model," *J. Chem. Theory Comput.* **13**, 6146–6157 (2017).
- <sup>73</sup>L. I. Vazquez-Salazar, E. D. Boittier, O. T. Unke, and M. Meuwly, "Impact of the characteristics of quantum chemical databases on machine learning prediction of tautomerization energies," *J. Chem. Theory Comput.* **17**, 4769–4785 (2021).
- <sup>74</sup>J. S. Smith, B. Nebgen, N. Lubbers, O. Isayev, and A. E. Roitberg, "Less is more: Sampling chemical space with active learning," *J. Chem. Phys.* **148**, 241733–241743 (2018).
- <sup>75</sup>J. Smith, B. Nebgen, R. Zubatyuk, N. Lubbers, C. Devereux, K. Barros, S. Tretiak, O. Isayev, and A. Roitberg, "Approaching coupled cluster accuracy with a general-purpose neural network potential through transfer learning," *Nat. Commun.* **10**, 2903 (2019).
- <sup>76</sup>C. Devereux, J. S. Smith, K. K. Huddleston, K. Barros, R. Zubatyuk, O. Isayev, and A. E. Roitberg, "Extending the applicability of the ANI deep learning molecular potential to sulfur and halogens," *J. Chem. Theory Comput.* **16**, 4192–4202 (2020).
- <sup>77</sup>J. Z. amd Yujun Tao, T. J. Giese, and D. M. York, "QD $\pi$ : A quantum deep potential interaction model for drug discovery," *J. Chem. Theory Comput.* **19**, 1261–1275 (2023).
- <sup>78</sup>P. O. Dral, X. Wu, and W. Thiel, "Semiempirical quantum-chemical methods with orthogonalization and dispersion corrections," *J. Chem. Theory Comput.* **15**, 1743–1760 (2019).
- <sup>79</sup>Y. Chen, Y. Ou, P. Zheng, Y. Huang, F. Ge, and P. O. Dral, "Benchmark of general-purpose machine learning-based quantum mechanical method AIQM1 on reaction barrier heights," *J. Chem. Phys.* **158**, 074103 (2023).
- <sup>80</sup>Y. Yang, H. Yu, D. M. York, Q. Cui, and M. Elstner, "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* **111**, 10861–10873 (2007).



- <sup>81</sup>M. Gaus, X. Lu, M. Elstner, and Q. Cui, "Parameterization of DFTB3/3OB for sulfur and phosphorus for chemical and biological applications," *J. Chem. Theory Comput.* **10**, 1518–1537 (2014).
- <sup>82</sup>H. Wang, L. Zhang, J. Han, and W. E, "DeePMD-kit: A deep learning package for many-body potential energy representation and molecular dynamics," *Comput. Phys. Commun.* **228**, 178–184 (2018).
- <sup>83</sup>W. Liang, J. Zeng, D. M. York, L. Zhang, and H. Wang, "Learning DeePMD-kit: A guide to building deep potential models," in *A Practical Guide to Recent Advances in Multiscale Modeling and Simulation of Biomolecules*, edited by Y. Wang and R. Zhou (AIP Publishing, 2023), Chap. 6, pp. 6–1–6–20.
- <sup>84</sup>Z. Yang, D. Hutter, P. Sheng, A. M. Sismour, and S. A. Benner, "Artificially expanded genetic information system: A new base pair with an alternative hydrogen bonding pattern," *Nucleic Acids Res.* **34**, 6095–6101 (2006).
- <sup>85</sup>L. Eberlein, F. R. Beierlein, N. J. R. van Eikema Hommes, A. Radadiya, J. Heil, S. A. Benner, T. Clark, S. M. Kast, and N. G. J. Richards, "Tautomeric equilibria of nucleobases in the hachimoji expanded genetic alphabet," *J. Chem. Theory Comput.* **16**, 2766–2777 (2020).
- <sup>86</sup>E. Biondi and S. A. Benner, "Artificially expanded genetic information systems for new aptamer technologies," *Biomedicines* **6**, 53 (2018).
- <sup>87</sup>B. Behera, P. Das, and N. R. Jena, "Accurate base pair energies of artificially expanded genetic information systems (AEGIS): Clues for their mutagenic characteristics," *J. Phys. Chem. B* **123**, 6728–6739 (2019).
- <sup>88</sup>C. A. Jerome, S. Hoshika, K. M. Bradley, S. A. Benner, and E. Biondi, "In vitro evolution of ribonucleases from expanded genetic alphabets," *Proc. Natl. Acad. Sci. U. S. A.* **119**, e2208261119 (2022).
- <sup>89</sup>M. J. S. Dewar and W. Thiel, "A semiempirical model for the two-center repulsion integrals in the NDDO approximation," *Theor. Chim. Acta* **46**, 89–104 (1977).
- <sup>90</sup>M. J. S. Dewar, E. Zebisch, E. F. Healy, and J. J. P. Stewart, "Development and use of quantum mechanical molecular models. 76. AM1: A new general purpose quantum mechanical molecular model," *J. Am. Chem. Soc.* **107**, 3902–3909 (1985).
- <sup>91</sup>J. J. P. Stewart, "Optimization of parameters for semiempirical methods V: Modification of NDDO approximations and application to 70 elements," *J. Mol. Model.* **13**, 1173–1213 (2007).
- <sup>92</sup>J. Řezáč and P. Hobza, "Advanced corrections of hydrogen bonding and dispersion for semiempirical quantum mechanical methods," *J. Chem. Theory Comput.* **8**, 141–151 (2012).
- <sup>93</sup>J. Řezáč and P. Hobza, "A halogen-bonding correction for the semiempirical PM6 method," *Chem. Phys. Lett.* **506**, 286–289 (2011).
- <sup>94</sup>J. J. P. Stewart, "Optimization of parameters for semiempirical methods VI: More modifications to the NDDO approximations and re-optimization of parameters," *J. Mol. Model.* **19**, 1–32 (2013).
- <sup>95</sup>M. Gaus, Q. Cui, and M. Elstner, "DFTB3: Extension of the self-consistent-charge density-functional tight-binding method (SCC-DFTB)," *J. Chem. Theory Comput.* **7**, 931–948 (2011).
- <sup>96</sup>S. Grimme, C. Bannwarth, and P. Shushkov, "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.* **13**, 1989–2009 (2017).
- <sup>97</sup>C. Bannwarth, S. Ehlert, and S. Grimme, "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.* **15**, 1652–1671 (2019).
- <sup>98</sup>N. Goldman, K. E. Kweon, B. Sadigh, T. W. Heo, R. K. Lindsey, C. H. Pham, L. E. Fried, B. Aradi, K. Holliday, J. R. Jeffries, and B. C. Wood, "Semi-automated creation of density functional tight binding models through leveraging Chebyshev polynomial-based force fields," *J. Chem. Theory Comput.* **17**, 4435–4448 (2021).
- <sup>99</sup>J.-D. Chai and M. Head-Gordon, "Systematic optimization of long-range corrected hybrid density functionals," *J. Chem. Phys.* **128**, 084106 (2008).
- <sup>100</sup>P. C. Bevilacqua, M. E. Harris, J. A. Piccirilli, C. Gaines, A. Ganguly, K. Kostenbader, Š. Ekesan, and D. M. York, "An ontology for facilitating discussion of catalytic strategies of RNA-cleaving enzymes," *ACS Chem. Biol.* **14**, 1068–1076 (2019).
- <sup>101</sup>M. J. Frisch, G. W. Trucks, H. B. Schlegel, G. E. Scuseria, M. A. Robb, J. R. Cheeseman, G. Scalmani, V. Barone, G. A. Petersson, H. Nakatsuji, X. Li, M. Caricato, A. V. Marenich, J. Bloino, B. G. Janesko, R. Gomperts, B. Mennucci, H. P. Hratchian, J. V. Ortiz, A. F. Izmaylov, J. L. Sonnenberg, D. Williams-Young, F. Ding, F. Lipparini, F. Egidi, J. Goings, B. Peng, A. Petrone, T. Henderson, D. Ranasinghe, V. G. Zakrzewski, J. Gao, N. Rega, G. Zheng, W. Liang, M. Hada, M. Ehara, K. Toyota, R. Fukuda, J. Hasegawa, M. Ishida, T. Nakajima, Y. Honda, O. Kitao, H. Nakai, T. Vreven, K. Throssell, J. A. Montgomery, Jr., J. E. Peralta, F. Ogliaro, M. J. Bearpark, J. J. Heyd, E. N. Brothers, K. N. Kudin, V. N. Staroverov, T. A. Keith, R. Kobayashi, J. Normand, K. Raghavachari, A. P. Rendell, J. C. Burant, S. S. Iyengar, J. Tomasi, M. Cossi, J. M. Millam, M. Klene, C. Adamo, R. Cammi, J. W. Ochterski, R. L. Martin, K. Morokuma, O. Farkas, J. B. Foresman, and D. J. Fox, Gaussian 16 Revision A.03, Gaussian, Inc., Wallingford CT, 2016.
- <sup>102</sup>D. L. Beveridge, *Approximate Molecular Orbital Theory of Nuclear and Electron Magnetic Resonance Parameters* (Springer, Boston, 1977).
- <sup>103</sup>T. Tuttle and W. Thiel, "OMx-D: Semiempirical methods with orthogonalization and dispersion corrections. Implementation and biochemical application," *Phys. Chem. Chem. Phys.* **10**, 2159–2166 (2008).
- <sup>104</sup>D. Tuna, Y. Lu, A. Koslowski, and W. Thiel, "Semiempirical quantum-chemical orthogonalization-corrected methods: Benchmarks of electronically excited states," *J. Chem. Theory Comput.* **12**, 4400–4422 (2016).
- <sup>105</sup>P. O. Dral, X. Wu, L. Spörkel, A. Koslowski, W. Weber, R. Steiger, M. Scholten, and W. Thiel, "Semiempirical quantum-chemical orthogonalization-corrected methods: Theory, implementation, and parameters," *J. Chem. Theory Comput.* **12**, 1082–1096 (2016).
- <sup>106</sup>P. O. Dral, X. Wu, L. Spörkel, A. Koslowski, and W. Thiel, "Semiempirical quantum-chemical orthogonalization-corrected methods: Benchmarks for ground-state properties," *J. Chem. Theory Comput.* **12**, 1097–1120 (2016).
- <sup>107</sup>D. A. Case, K. Belfon, I. Y. Ben-Shalom, S. R. Brozell, D. S. Cerutti, T. E. Cheatham III, V. W. D. Cruzeiro, T. A. Darden, R. E. Duke, G. Giambasu, M. K. Gilson, H. Gohlke, A. W. Goetz, R. Harris, S. Izadi, S. A. Izmailov, K. Kasavajhala, K. Kovalenko, R. Krasny, T. Kurtzman, T. Lee, S. Le-Grand, P. Li, C. Lin, J. Liu, T. Luchko, R. Luo, V. Man, K. Merz, Y. Miao, O. Mikhailovskii, G. Monard, H. Nguyen, A. Onufriev, F. Pan, S. Pantano, R. Qi, D. R. Roe, A. Roitberg, C. Sagui, S. Schott-Verdugo, J. Shen, C. L. Simmerling, N. Skrynnikov, J. Smith, J. Swails, R. C. Walker, J. Wang, R. M. Wilson, R. M. Wolf, X. Wu, Y. Xiong, Y. Xue, D. M. York, and P. A. Kollman, AMBER 20, University of California, San Francisco, CA, 2020.
- <sup>108</sup>R. C. Walker, M. F. Crowley, and D. A. Case, "The implementation of a fast and accurate QM/MM potential method in Amber," *J. Comput. Chem.* **29**, 1019–1031 (2008).
- <sup>109</sup>M. Thiel, MNDO, Max-Planck-Institut für Kohlenforschung, Mülheim an der Ruhr, 2022.
- <sup>110</sup>J. J. P. Stewart, "MOPAC: A semiempirical molecular orbital program," *J. Comput.-Aided Mol. Des.* **4**, 1–105 (1990).
- <sup>111</sup>J. Hostaš, J. Řezáč, and P. Hobza, "On the performance of the semiempirical quantum mechanical PM6 and PM7 methods for noncovalent interactions," *Chem. Phys. Lett.* **568**, 161–166 (2013).
- <sup>112</sup>A. V. Sulimov, D. C. Kutov, E. V. Katkova, I. S. Ilin, and V. B. Sulimov, "New generation of docking programs: Supercomputer validation of force fields and quantum-chemical methods for docking," *J. Mol. Graphics Modell.* **78**, 139–147 (2017).
- <sup>113</sup>T. J. Giese and D. M. York, "Density-functional expansion methods: Grand challenges," *Theor. Chem. Acc.* **131**, 1145 (2012).
- <sup>114</sup>M. Gaus, A. Goez, and M. Elstner, "Parameterization and benchmark of DFTB3 for organic molecules," *J. Chem. Theory Comput.* **9**, 338–354 (2013).
- <sup>115</sup>G. Seabra, R. C. Walker, M. Elstner, D. A. Case, and A. E. Roitberg, "Implementation of the SCC-DFTB method for hybrid QM/MM simulations within the amber molecular dynamics package," *J. Phys. Chem. A* **111**, 5655–5664 (2007).
- <sup>116</sup>C. H. Pham, R. K. Lindsey, L. E. Fried, and N. Goldman, "High-accuracy semiempirical quantum models based on a minimal training set," *J. Phys. Chem. Lett.* **13**, 2934–2942 (2022).
- <sup>117</sup>R. K. Lindsey, L. E. Fried, and N. Goldman, "ChIMES: A force matched potential with explicit three-body interactions for molten carbon," *J. Chem. Theory Comput.* **13**, 6222–6229 (2017).



- <sup>118</sup>H. Li, C. Collins, M. Tanha, G. J. Gordon, and D. J. Yaron, "A density functional tight binding layer for deep learning of chemical Hamiltonians," *J. Chem. Theory Comput.* **14**, 5764–5776 (2018).
- <sup>119</sup>J. Zhu, V. Q. Vuong, B. G. Sumpter, and S. Irle, "Artificial neural network correction for density-functional tight-binding molecular dynamics simulations," *MRS Commun.* **9**, 867–873 (2019).
- <sup>120</sup>M. Stöhr, L. Medrano Sandonas, and A. Tkatchenko, "Accurate many-body repulsive potentials for density-functional tight binding from deep tensor neural networks," *J. Phys. Chem. Lett.* **11**, 6835–6843 (2020).
- <sup>121</sup>C. Panosetti, A. Engelmann, L. Nemeč, K. Reuter, and J. T. Margraf, "Learning to use the force: Fitting repulsive potentials in density-functional tight-binding with Gaussian process regression," *J. Chem. Theory Comput.* **16**, 2181–2191 (2020).
- <sup>122</sup>J. J. Kranz, M. Kubillus, R. Ramakrishnan, O. A. von Lilienfeld, and M. Elstner, "Generalized density-functional tight-binding repulsive potentials from unsupervised machine learning," *J. Chem. Theory Comput.* **14**, 2341–2352 (2018).
- <sup>123</sup>X. Gao, F. Ramezanghorbani, O. Isayev, J. S. Smith, and A. E. Roitberg, "TorchANI: A free and open source PyTorch-based deep learning implementation of the ANI neural network potentials," *J. Chem. Inf. Model.* **60**, 3408–3415 (2020).
- <sup>124</sup>E. Caldeweyher, C. Bannwarth, and S. Grimme, "Extension of the D3 dispersion coefficient model," *J. Chem. Phys.* **147**, 034112 (2017).
- <sup>125</sup>D. C. Liu and J. Nocedal, "On the limited memory BFGS method for large scale optimization," *Math. Program.* **45**, 503–528 (1989).
- <sup>126</sup>A. Hjorth Larsen, J. Jørgen Mortensen, J. Blomqvist, I. E. Castelli, R. Christensen, M. Dułak, J. Friis, M. N. Groves, B. Hammer, C. Hargus, E. D. Hermes, P. C. Jennings, P. Bjerre Jensen, J. Kermode, J. R. Kitchin, E. Leonhard Kolsbjerg, J. Kubal, K. Kaasbjerg, S. Lysgaard, J. Bergmann Maronsson, T. Maxson, T. Olsen, L. Pastewka, A. Peterson, C. Rostgaard, J. Schiøtz, O. Schütt, M. Strange, K. S. Thygesen, T. Vegge, L. Vilhelmsen, M. Walter, Z. Zeng, and K. W. Jacobsen, "The atomic simulation environment—A Python library for working with atoms," *J. Phys.: Condens. Matter* **29**, 273002 (2017).
- <sup>127</sup>J. S. Smith, R. Zubatyuk, B. Nebgen, N. Lubbers, K. Barros, A. E. Roitberg, O. Isayev, and S. Tretiak, "The ANI-1ccx and ANI-1x data sets, coupled-cluster and density functional theory properties for molecules," *Sci. Data* **7**, 134 (2020).
- <sup>128</sup>T. Fink, H. Bruggesser, and J.-L. Reymond, "Virtual exploration of the small-molecule chemical universe below 160 daltons," *Angew. Chem., Int. Ed.* **44**, 1504–1508 (2005).
- <sup>129</sup>L. C. Blum and J.-L. Reymond, "970 million druglike small molecules for virtual screening in the chemical universe database GDB-13," *J. Am. Chem. Soc.* **131**, 8732–8733 (2009).
- <sup>130</sup>V. Law, C. Knox, Y. Djoumbou, T. Jewison, A. C. Guo, Y. Liu, A. Maciejewski, D. Arndt, M. Wilson, V. Neveu, A. Tang, G. Gabriel, C. Ly, S. Adamjee, Z. T. Dame, B. Han, Y. Zhou, and D. S. Wishart, "DrugBank 4.0: Shedding new light on drug metabolism," *Nucleic Acids Res.* **42**, 1091–1097 (2014).
- <sup>131</sup>L. Goerigk, H. Kruse, and S. Grimme, "Benchmarking density functional methods against the S66 and S66x8 datasets for non-covalent interactions," *ChemPhysChem.* **12**, 3421–3433 (2011).
- <sup>132</sup>B. Brauer, M. K. Kesharwani, S. Kozuch, and J. M. L. Martin, "The S66x8 benchmark for noncovalent interactions revisited: Explicitly correlated *ab initio* methods and density functional theory," *Phys. Chem. Chem. Phys.* **18**, 20905–20925 (2016).
- <sup>133</sup>J. Řezáč, "Non-covalent interactions atlas benchmark data sets: Hydrogen bonding," *J. Chem. Theory Comput.* **16**, 2355–2368 (2020).
- <sup>134</sup>O. Wahl and T. Sander, "Tautobase: An open tautomer database," *J. Chem. Inf. Model.* **60**, 1085–1089 (2020).
- <sup>135</sup>M. Wieder, J. Fass, and J. D. Chodera, "Fitting quantum machine learning potentials to experimental free energy data: Predicting tautomer ratios in solution," *Chem. Sci.* **12**, 11364–11381 (2021).
- <sup>136</sup>A. Moser, K. Range, and D. M. York, "Accurate proton affinity and gas-phase basicity values for molecules important in biocatalysis," *J. Phys. Chem. B* **114**, 13911–13921 (2010).
- <sup>137</sup>L. Goerigk, A. Hansen, C. Bauer, S. Ehrlich, A. Najibi, and S. Grimme, "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.* **19**, 32184–32215 (2017).
- <sup>138</sup>N. Ree, A. H. Göller, and J. H. Jensen, "RegioSQM20: Improved prediction of the regioselectivity of electrophilic aromatic substitutions," *J. Cheminf.* **13**, 10 (2021).
- <sup>139</sup>N. B. Leontis and E. Westhof, "Geometric nomenclature and classification of RNA base pairs," *RNA* **7**, 499–512 (2001).
- <sup>140</sup>N. B. Leontis, J. Stombaugh, and E. Westhof, "The non-Watson–Crick base pairs and their associated isostericity matrices," *Nucleic Acids Res.* **30**, 3497–3531 (2002).
- <sup>141</sup>N. B. Leontis and E. Westhof, "Analysis of RNA motifs," *Curr. Opin. Struct. Biol.* **13**, 300–308 (2003).
- <sup>142</sup>I. Singh, M.-J. Kim, R. W. Molt, S. Hoshika, S. A. Benner, and M. M. Georgiadis, "Structure and biophysics for a six letter DNA alphabet that includes imidazo[1,2-a]-1,3,5-triazine-2(8H)-4(3H)-dione (X) and 2,4-diaminopyrimidine (K)," *ACS Synth. Biol.* **6**, 2118–2129 (2017).
- <sup>143</sup>R. T. Raines, "Ribonuclease A," *Chem. Rev.* **98**, 1045–1066 (1998).
- <sup>144</sup>H. Gu, S. Zhang, K.-Y. Wong, B. K. Radak, T. Dissanayake, D. L. Kellerman, Q. Dai, M. Miyagi, V. E. Anderson, D. M. York, J. A. Piccirilli, and M. E. Harris, "Experimental and computational analysis of the transition state for ribonuclease A-catalyzed RNA 2'-O-transphosphorylation," *Proc. Natl. Acad. Sci. U. S. A.* **110**, 13002–13007 (2013).
- <sup>145</sup>M. E. Harris, J. A. Piccirilli, and D. M. York, "Integration of kinetic isotope effect analyses to elucidate ribonuclease mechanism," *Biochim. Biophys. Acta, Proteins Proteomics* **1854**, 1801–1808 (2015).