

HEPOM: Using Graph Neural Networks for the Accelerated Predictions of Hydrolysis Free Energies in Different pH Conditions

Published as part of *Journal of Chemical Information and Modeling* special issue “Modeling Reactions from Chemical Theories to Machine Learning”.

Rishabh D. Guha,[◆] Santiago Vargas,[◆] Evan Walter Clark Spotte-Smith, Alexander Rizzolo Epstein, Maxwell Venetos, Ryan Kingsbury, Mingjian Wen, Samuel M. Blau, and Kristin A. Persson*



Cite This: <https://doi.org/10.1021/acs.jcim.4c02443>



Read Online

ACCESS |



Metrics & More

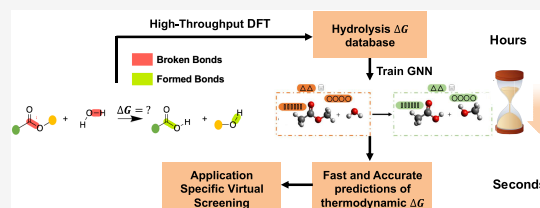


Article Recommendations



Supporting Information

ABSTRACT: Hydrolysis is a fundamental family of chemical reactions where water facilitates the cleavage of bonds. The process is ubiquitous in biological and chemical systems, owing to water’s remarkable versatility as a solvent. However, accurately predicting the feasibility of hydrolysis through computational techniques is a difficult task, as subtle changes in reactant structure like heteroatom substitutions or neighboring functional groups can influence the reaction outcome. Furthermore, hydrolysis is sensitive to the pH of the aqueous medium, and the same reaction can have different reaction properties at different pH conditions. In this work, we have combined reaction templates and high-throughput ab initio calculations to construct a diverse data set of hydrolysis free energies. The developed framework automatically identifies reaction centers, generates hydrolysis products, and utilizes a trained graph neural network (GNN) model to predict ΔG values for all potential hydrolysis reactions in a given molecule. The long-term goal of the work is to develop a data-driven, computational tool for high-throughput screening of pH-specific hydrolytic stability and the rapid prediction of reaction products, which can then be applied in a wide array of applications including chemical recycling of polymers and ion-conducting membranes for clean energy generation and storage.



1. INTRODUCTION

Water is one of the most essential molecules in chemistry, and yet, its unique properties make it notoriously difficult to characterize.^{1,2} The significant electronegativity differences between its oxygen and hydrogen atoms gives water a highly polar character that leads to its recognition as the “universal solvent”.^{3,4} Hydrolysis, or any reaction where water acts as both a reactant and the solvent medium,^{5,6} are a prevalent class of reactions across chemistry. Hydrolytic reactions are fundamental in biological^{7,8} and synthetic chemistry^{9,10} and play a critical role in various essential scientific processes and significant technological applications. These range from processes, such as human digestion,^{8,11} where enzymes facilitate the hydrolytic breakdown of complex macronutrients into simpler compounds, to the degradation of hazardous pollutants¹² and alternative plastic chemistries.¹³

At the molecular level, hydrolysis begins when a water molecule attacks specific sites on the reactant, initiating a sequence of bond cleavages and formations that lead to new product(s). The mechanism and the associated rate of this reaction is closely tied to the pH of the aqueous medium.^{14,15} The availability of protons (H^+) or hydroxide (OH^-) ions catalyzes the formation of charged species, which have markedly different reactivities compared to their neutral counterparts.^{13,16}

These ionized reactants can exhibit enhanced solubility^{17,18} by forming stronger hydrogen bonds with the solvent. Additionally, water can act as catalyst, facilitating ion transfer through the solvent and creating alternate reaction pathways with lower energy barriers.^{19,20} As a result, acid/base-catalyzed hydrolysis of the same reactant can have significantly different reaction rates compared to its neutral form, adding complexity to the study of these reactions.

Given activation barriers (ΔG^\ddagger), the experimental rate of a hydrolysis reaction can be directly correlated via the Eyring equation.^{16,21,22} This involves determining computationally intensive and difficult to find transition states for each individual reaction along the reaction coordinate of the potential energy surface (PES).^{16,23,24} In contrast, within a specific reaction family, the Bell–Evans–Polanyi principle (BEP)²⁵ can offer an approximate linear correlation between the thermodynamic Gibbs free energy change (ΔG_r) and the kinetic parameter

Received: January 2, 2025

Revised: February 25, 2025

Accepted: March 28, 2025

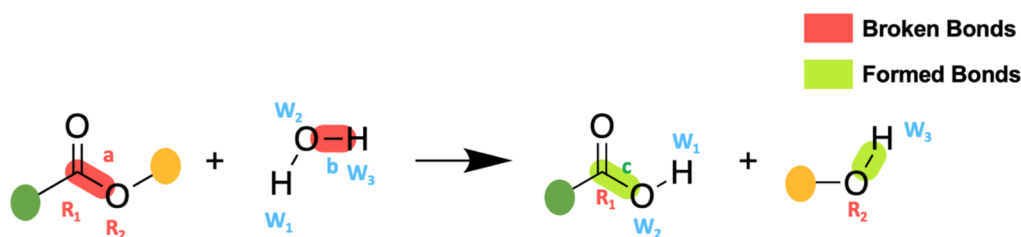


Figure 1. Set of bond cleavage and formations necessary to generate hydrolysis products for a representative ester molecule.

ΔG^\ddagger .^{26–28} In cases where it holds, BEP allows us to leverage the thermochemistry of products and reactants (ΔG_r) to approximate trends in the kinetic rates of the reaction. This opens the avenue for the development of a computational screening tool that can calculate the respective ΔG_r 's of all potential hydrolysis pathways and screen molecules for a specific hydrolysis-related application. Despite this, quantifying thermochemical quantities such as ΔG_r with high accuracy still requires DFT calculations with large basis sets and refined hybrid functionals at both reaction end points.^{29,30} Depending on the size of molecules, these calculations can take anywhere from several hours to days, particularly when employing implicit solvent models³¹ to approximate the contributions from the reaction environment.

Since computational cost is a severe bottleneck for any form of high-throughput screening, deep learning approaches have emerged as promising alternatives in the past decade, especially for tasks that involve the establishment of structure-to-property relationships.^{32,33} Recently, graph convolutions, which iteratively update node and edge features based on connectivity and local environment, have proven to be extremely effective in learning molecular^{34,35} and reaction representations.^{36,37} Despite these methodological advances, the largest roadblock to the development of an accurate model is typically the procurement of diverse, representative data. For instance, the model developed by Grambow et al.³³ was facilitated by a data set of 12,000 gas-phase reactions³⁸ sampled from a subset of molecules in the GDB-17 data set.³⁹ The bond dissociation energy (BDE) prediction framework developed by Wen et al.⁴⁰ was trained on a data set of over 60,000 homolytic and heterolytic bond dissociation reactions.⁴¹ In the realm of hydrolysis, no such comprehensive data set currently exists.

In this work, we have attempted to address these shortcomings by first developing a predictive framework based on reaction templates for different functional groups that can automatically generate hydrolysis products for multiple pathways in any molecule. This framework was then applied on a subset of the QM9⁴² and the Alchemy⁴³ databases to generate a database of over 65,000 hydrolysis reactions in an implicit aqueous solvation environment. For a given reactant molecule in the QM9 subset of the data, we have also generated corresponding hydroxylated and protonated states of the reactant molecule to approximate the effects of extreme pH on the ΔG_r of hydrolysis. In addition, the neutral fold of the data set was developed with reactants from the QM9 database and later augmented with the inclusion of larger reactant molecules from the Alchemy⁴³ data set. Combined, we provide a new data set hydrolysis reactions that encompasses thermodynamic properties at different charged (protonated/hydroxylated) states along with a large, exploratory set of neutral-pH reactions for analysis and model development.

We then proceeded to use this comprehensive data set to train a GNN model, which serves as a Hydrolysis Energy Predictor for

Organic Molecules (HEPOM). The model leverages the difference features of the atom (node), bond (edge), and global features between the products and the reactants to directly predict the DFT-calculated ΔG_r . The global reaction atom mapping allows the model to simultaneously track multiple bond dissociations and formations. For the neutral data set, the model achieved a low mean absolute error (MAE) of 1.73 kcal/mol on a diverse holdout set of hydrolysis reactions and it was also successful in outperforming a diverse set of benchmark models on the smaller and more complex protonated and hydroxylated data sets.

2. METHODS

2.1. Reaction Generation. We segmented the construction of our data set into four main parts: three derived from the QM9 data set (representing neutral, protonated, and hydroxylated reactions) and another, neutral reaction set from the Alchemy data set. Hydrolyzable molecules in QM9 were screened using RDKit⁴⁴ substructure matching for 20 standard, hydrolyzable functional groups (Figure S5b). These templates were adapted from the work by Tebes-Stevens et al.⁴⁵ and integrated into an automated framework to predict reaction products. For instance, in Figure 1, if an ester functional group was detected in a molecule, the reaction template would yield a carboxylic acid and an alcohol as the respective hydrolysis products. Bond 'a' in the reactant and bond 'b' in the water molecule were deleted with the RemoveBond functionality in RDKit. Then, AddBond was used to create bonds 'c' and 'd' between atoms R1-W2 and R2-W3 respectively, to yield a carboxylic acid and an alcohol as the respective products. Similar reaction templates were implemented for all functional groups. If multiple competing functional groups were identified in a single molecule, then the hydrolysis products were generated independently for each individual functional group (Schematic S1 in the Supporting Information) and added as separate reaction entries to the database.

Nitriles were treated differently: the reaction template yields amides and these can be further hydrolyzed into an amine and a carboxylic acid (Schematic S2 in the Supporting Information). As a result, the intermediate products of nitrile reactions served as reactants in additional hydrolysis reactions, thereby augmenting the data set. We generated a total of 16,264 hydrolysis reactions from the QM9 data set. An initial model was trained on 15,264 of these neutral reactions, while 1000 reactions were kept in an unseen holdout test set. Additional details regarding the product generation workflow has been included in Section S3 of the Supporting Information. The performance of this model is discussed in Section 3.2 (*vide infra*).

The broader goal of this work is to develop a framework capable of enumerating potential hydrolysis pathways for a wide range of molecules and predicting the thermodynamic free

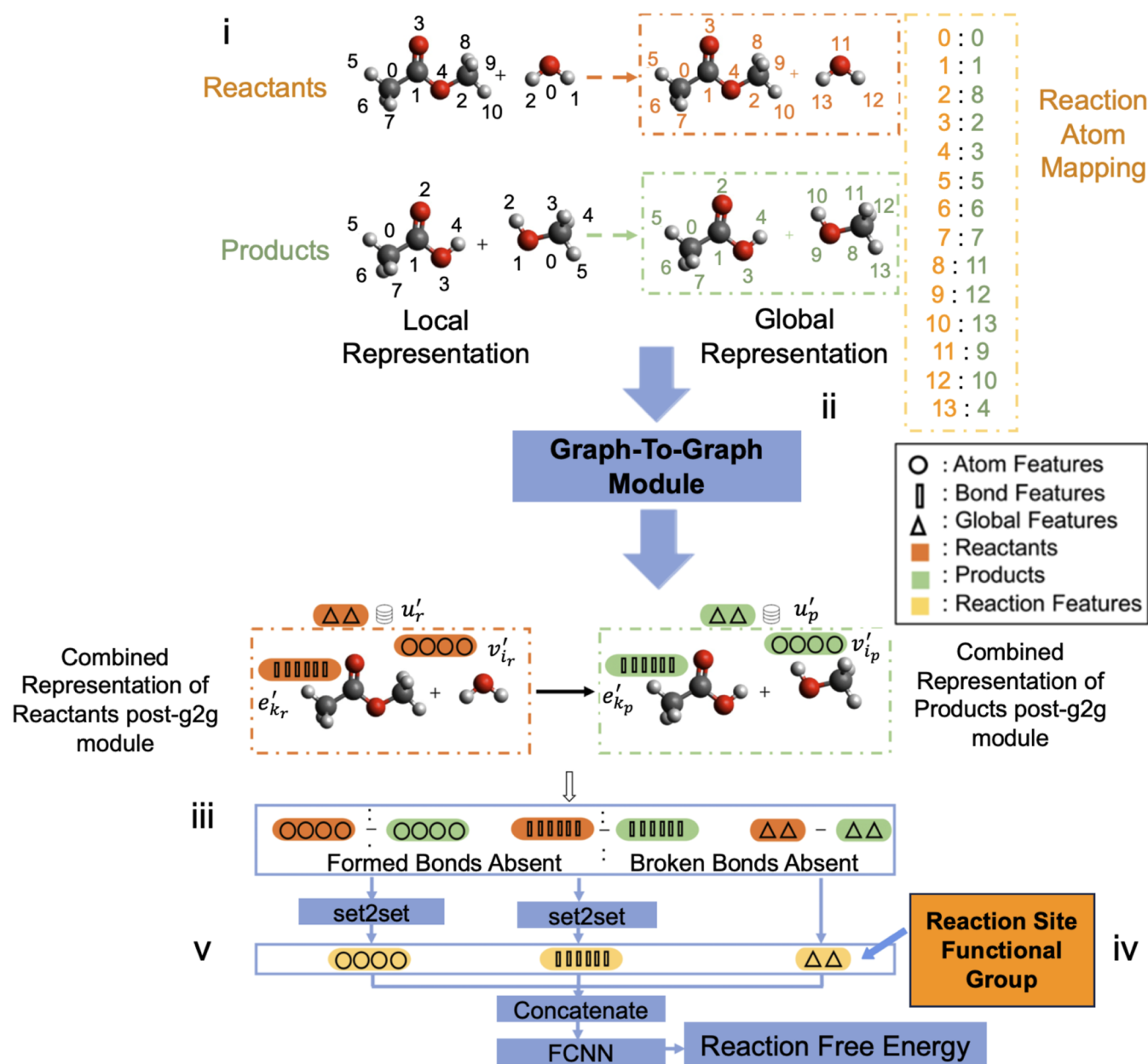


Figure 2. GNN Architecture: The user inputs atom-mapped sets of reactants and products (i) which undergo message-passing and update steps (ii). Using the user-specified mappings, these updated features are mapped to a global reaction graph (iii) where functional groups are the reaction site is added as a global feature (iv). Embeddings of bond and atom features plus global features directly serve as the fixed-size vector used in a conventional dense neural network for property prediction.

energies of these pathways with high accuracy. With this in mind, we screened molecules from the Alchemy data set⁴³ to generate data set entries representing larger, more complex reactions. The Alchemy data set includes molecules with up to 14 heavy atoms, though we only considered molecules with more than 10 heavy atoms. Neutral reactions generated from the Alchemy data set were added to the original data set, resulting in a combined data set of **41,006** reactions. Of these, **2800** reactions were filtered out and added to the previous QM9 test set, creating an unseen holdout test set of **3800** reactions (approximately 10% of the training set size). In the following sections of the manuscript, we refer to the original data set as the QM9 data set and the expanded data set as the QM9 + Alchemy data set.

While the data generated above are useful and novel, they are examples of hydrolysis in a neutral reaction medium. However, hydrolysis is often catalyzed in an acidic or basic reaction medium. For example, the hydrolysis rate of amides in a neutral medium is negligible, even after heating, but amidic hydrolysis proceeds at a moderate rate in an acidic or basic medium,^{46,47} forming carboxylic acid and an amine. Consequently, we explored whether this framework could extend to broader reaction conditions, potentially serving as a screening tool for identifying molecules amenable to pH-specific hydrolysis. It is important to clarify that under these reaction conditions, hydrolysis is initiated by the protonation or hydroxylation of the reacting functional group,^{20,48} and the overall reaction rate is heavily influenced by the pK_a values of the functional groups.^{49,50} Generating a diverse data set in high throughput

for acid- or base-catalyzed hydrolysis while accounting for pK_a was intractable with our current data generation scheme. Therefore, in this work, we focused our efforts on developing a unified model that can predict the differences in hydrolysis reaction-free energies for the same functional group under different pH conditions. The data sets generated for neutral or basic pH assume that the reaction medium is at an extreme pH, i.e., if a functional group can protonate or hydroxylate, it will. An alternative approach could involve applying a separate ML model which predicts the pK_a and identifies the most probable site for protonation or hydroxylation at a specific pH,^{51–53} but this is beyond the scope of the present study.

We separated extreme pH hydrolysis reactions into two reactions schemes. For an acidic medium, the reacting functional group was assumed to be protonated at the most electron-rich atom site (e.g., the carbonyl oxygen in an ester or amide, or the nitrogen atom in a nitrile). Conversely, for a basic pH, the functional group's relevant electrophilic site—such as the carbonyl carbon in an ester or amide or the ring carbon in the epoxide—was hydroxylated. The acidic pH reaction was then executed between the protonated reactant and two water molecules to maintain reaction stoichiometry. A representative example elucidating the differences in the hydrolysis reaction in acidic and neutral pH for a hydrolyzing carbamate molecule is demonstrated in Figure S4a,b of the Supporting Information. The extra water molecule on the reactant side absorbs the proton to generate hydronium as one of the reaction products. This approach circumvents the erroneous DFT-calculated energies of an isolated proton in an implicit solvent medium.⁵⁴ In the case of basic pH, the hydroxylated reactant decomposes into the reaction products and a hydroxide ion. For these two data sets, we focused on the QM9 molecules to limit the scope of computations, yielding a protonated data set of 11,323 reactions and a hydroxylated data set of 16,732 reactions. Holdout test sets consistent with the neutral QM9 data set were also extracted before model training. Since the protonated reactants have a +1 charge and the hydroxylated reactants a −1 charge, we will refer to these data sets as QM9⁺ and QM9[−], respectively, in the subsequent sections.

2.2. Density-Functional Theory. QChem (version 5 or 6)⁵⁵ was used to perform all the DFT calculations necessary to generate the data set. A specialized frequency-flattening optimization (FFOpt) workflow, originally developed by Spotte-Smith et al.⁴¹ and currently implemented in atomate⁵⁶ was used to optimize the reactant and product structures to a true minima and also obtain thermochemical quantities from the vibrational frequencies. The workflow iteratively performs successive geometry optimizations and frequency calculations until there are either none or a single negligible negative frequency (<15 cm^{−1}). This approach ensures that the optimized structure is a true local minimum of the PES and not a saddle point. Moreover, the workflow parses the necessary enthalpy and entropy terms from the QChem frequency output document for the free energy calculations. For all the DFT calculations, we used the range-separated metaGGA hybrid functional, ω B97M-V,⁵⁷ which employs the VV10 dispersion correction,⁵⁸ to improve the noncovalent interactions. The def2-SVPD basis set⁵⁹ was employed for the FFOpt workflow and the solvation effects were implicitly accounted for with the water SMD solvent model.¹⁷ The electronic energies of the optimized structures were refined with single-point calculations using a larger def2-QZVPPD basis set.⁵⁹

2.3. Model Architecture. The GNN model (Figure 2) is based on the previous BonDNet architecture.⁴⁰ Here we briefly discuss that model before highlighting our key departures and how these differences are key to working with our data sets here.

The original algorithm uses gated graph convolutional (GatedGC) layers to propagate initial node features within the graphs of individual species on both sides of a reaction. While GatedGC layers were used widely for structure-to-property models in chemistry and materials science,^{60,61} BonDNet improved on these previous implementations by integrating update and message-passing equations between global nodes and atom/bond type nodes; this allows for the treatment of species of different charges and provides a framework to include molecular-level features. Similar to other graph neural networks, more distant graph relationships are treated by iteratively stacking several (typically 2–4 layers) GatedGC layers. With updated species graphs, reaction graphs are built to hold reaction feature differences—atom and bond nodes are mapped to each other on both sides of a reaction and their features are subtracted between corresponding atoms/bonds. Broken bonds are represented by zero vectors in this scheme. Here BonDNet implemented a custom set2set⁶² global pooling feature to map updated reaction graphs to fixed-sized vectors. These vectors are passed to fully connected layers for property prediction.

Our implementation extends pooling by integrating a diverse set of global pooling functions, including set2set,⁶² WeightedMeanPooling, Self-attention pooling,⁶³ and Mean Pooling. This diverse set of global pooling functions was intended to provide a more comprehensive toolkit of architectures across different data set sizes. Previous benchmarks showed set2set layers did not always outperform simpler MeanPooling approaches.⁶⁴

In this implementation, the reaction mapping is altered from the original BonDNet to a global reaction graph that is constructed between the union set of bonds in products and reactants. Originally, BonDNet used the product graph as a scaffold and subtracted reactant features from corresponding nodes in the product graph. This limited the model to only being applicable for $A \rightarrow B$ and $A \rightarrow B + C$ type reactions with a single bond dissociation. The previous framework could not interpret a hydrolysis reaction that involves at least two elementary bond dissociation and formation reactions. Our algorithm builds a global reaction graph by taking the union set of atoms and bonds in products and reactants and uses this to build a graph structure with bonds from each side of reaction. The approach allows us to precompute global graphs, including mappings and descriptors, for offline preprocessing prior to training. This change allows for an arbitrary number of bond changes, simultaneous breaking and forming, to be treated by the model (Figure 2). In addition, we are able to generalize our model to any number of species on either side of the reaction—a feature critical for hydrolysis where no reaction can fit BonDNet's original implementation. This update is important for others looking to use our new model for reaction-property prediction as now the model is completely flexible to any number of bond changes and number of species.

For the task of hydrolysis, where we have a consistent reaction framework, we incorporated a one-hot encoding of functional group identity⁴⁵ into the global feature nodes. This encoding provides a simple, yet effective, descriptor that captures the reaction site of hydrolysis reactions alongside the more distant features generated by stacked message-passing layers. This is a particularly attractive feature, as sequential stacking of message-passing layers rapidly increases compute time and can lead to

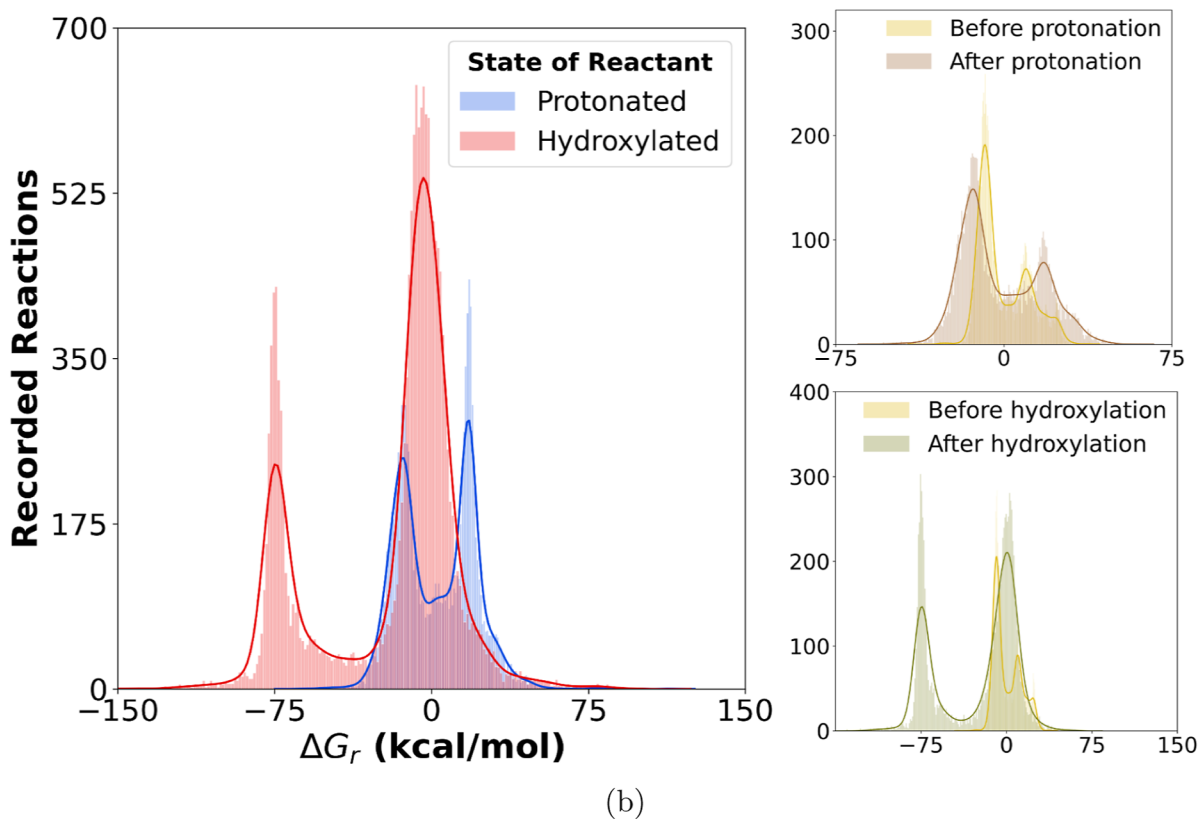
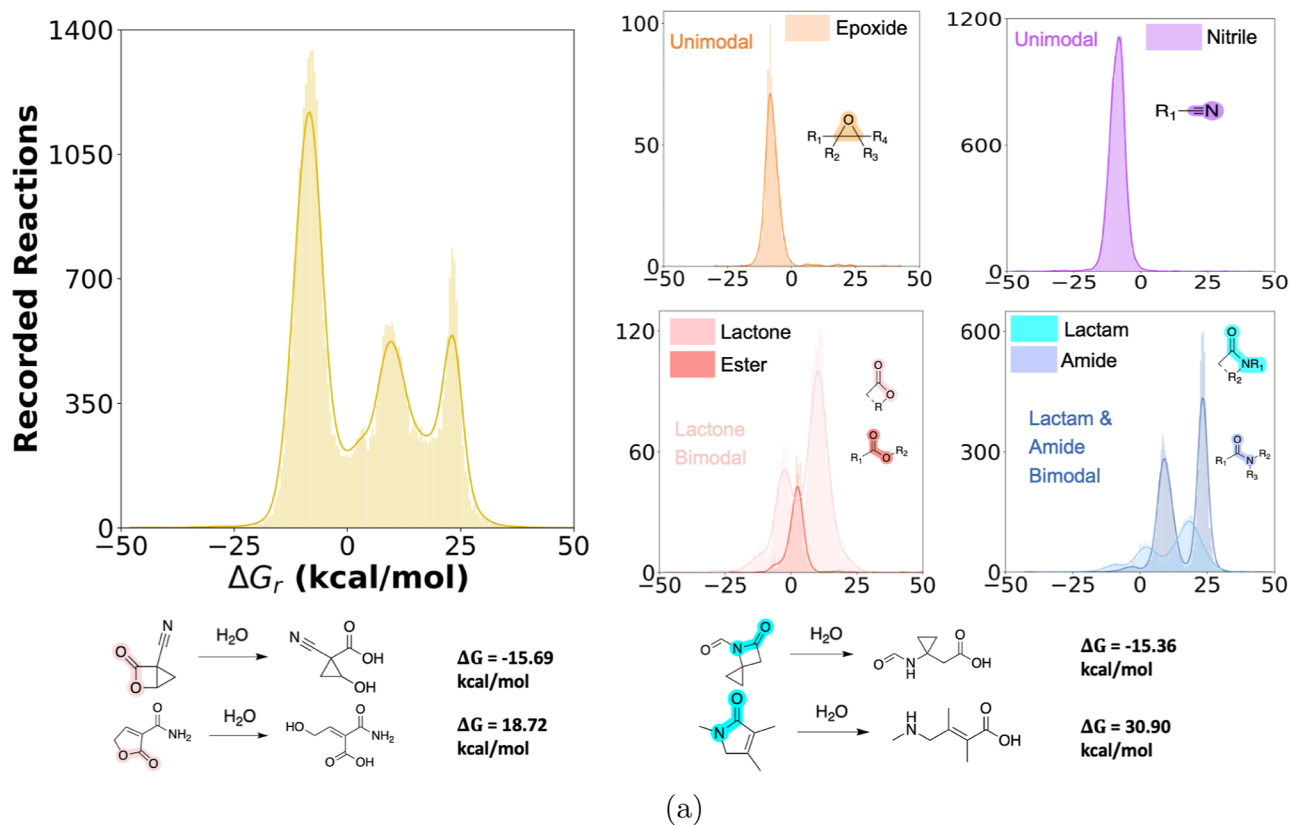


Figure 3. (a) Distribution of hydrolysis ΔG_r for the QM9 + Alchemy data set. The associated subfigures show the ΔG_r distributions for specific functional groups (b) Distribution of hydrolysis ΔG_r for the QM9⁺ and QM9⁻ data sets. The associated subfigures show the shift in the ΔG_r distributions for a subset of common reactants in the three data sets. The x and y axes on all the subfigures represent the same metrics as the main figures.

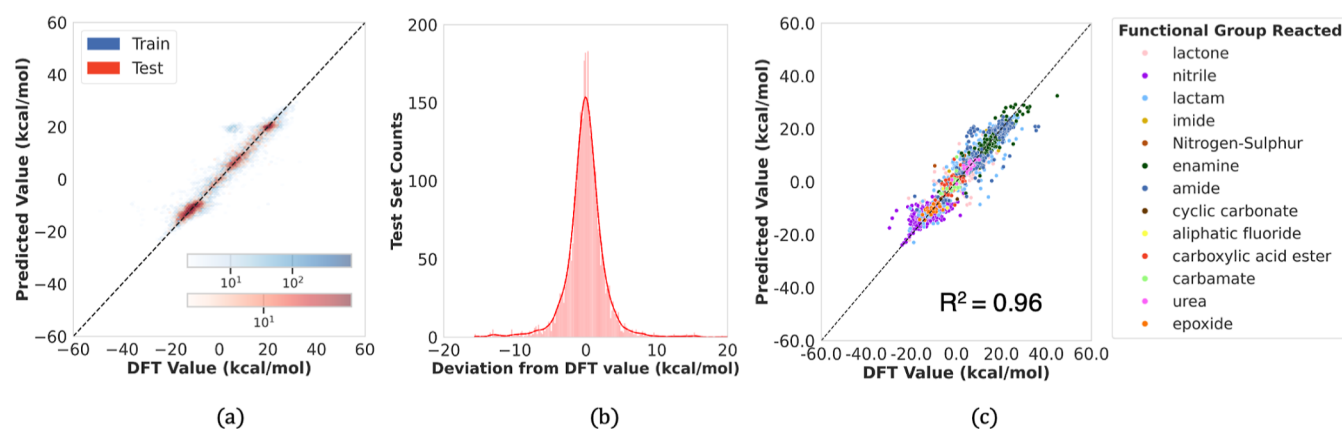


Figure 4. Performance of HEPOM on the QM9 + Alchemy data set. (a) ΔG_r predicted by HEPOM versus DFT reference labels for the train and test sets; (b) histogram of the prediction errors; (c) parity plot for the holdout test set segregated on the basis of hydrolyzed functional group.

problems such as oversmoothing.^{65,66} Our heterograph architecture, by including global nodes connected to every atom and bond, allows for distal information passing while avoiding such issues. We also implemented a host of computational features such as multiGPU compatibility, a pytorch-lightning implementation, and added support for preprocessing data sets. These latter features are vital for our data sets while also providing the community with a flexible, general GNN architecture for reaction-property prediction at scale.

3. RESULTS AND DISCUSSION

3.1. Data Set Overview. As detailed in Section 2.1, our hydrolysis database, in its current form, comprises a total of 68,761 reactions, making it the largest molecular database for hydrolysis reactions. Among these, the QM9 + Alchemy data set contains 41,006 reactions with reactant molecules in their neutral state, while the remaining reactions are approximately evenly split between the QM9⁺ (protonated) and QM9⁻ data sets, representing acidic and basic pH conditions, respectively.

The QM9 + Alchemy data set contains reactants with up to 12 heavy atoms. The distribution of reactants based on the number of heavy elements is illustrated in Figure S5c of the Supporting Information. For the charged subsets (QM9⁺ and QM9⁻), the reactants are restricted to a maximum of 9 heavy atoms.

The number of hydrolyzed products varies depending on the reacting functional group, with reactions yielding one, two, and, in some cases (e.g., urea and carbamates), three products. Figure S5a of the Supporting Information visualizes the distribution of reactions based on the number of products generated, and Figure S5b illustrates the distribution across different hydrolyzed functional groups.

The hydrolysis ΔG_r distribution for the QM9 + Alchemy data set is presented in Figure 3a, where three peaks are observed: two distinct peaks in the endergonic region ($\Delta G_r > 0$) and one larger peak in the exergonic regime ($\Delta G_r < 0$). Interestingly, the ΔG_r distribution in Figure 3a is almost perfectly balanced with 20,547 reactions (50.11%) of the neutral reactions falling within the endergonic regime.

Further analysis across different functional groups reveals some interesting insights. Distributions of epoxides, nitriles, and esters exhibit unimodal energy distributions, while cyclic esters and cyclic amides (e.g., lactones and lactams) are bimodal. Sampling random lactone and lactam reactions (Figure 3a) from the endergonic and exergonic regimes indicates that cyclic

structures with strained rings have more favorable thermodynamic hydrolysis pathways, whereas stable five-membered rings are more resistant to hydrolysis.^{67,68} The hydrolysis of amides shows a distinctly bimodal nature, with both peaks centered in the endergonic regime, consistent with the established trend of the thermodynamic infeasibility of amide hydrolysis in a neutral reaction medium.⁶⁹

The energy distribution for the protonated (QM9⁺) and hydroxylated (QM9⁻) data sets is shown in Figure 3b. It is evident that the ΔG_r distribution for hydroxylated reactants shifts strongly toward the exergonic regime with greater than 70% of the reactions with a thermodynamically exergonic hydrolytic pathway. The shift in the protonated data set is more subtle; however, when comparing the corresponding slices of the same reactions in the neutral and protonated states (Figure 3b), it is clear that the distribution broadens after protonation.

The following section will discuss our model's performance on these different data sets and how it compares to existing benchmarks.

3.2. Model Performance—Neutral Data Set. In the initial round of model training, we utilized a data set of hydrolysis reactions generated exclusively from reactants extracted from the QM9 database. Despite a modest training set of 15,264 reactions, the model performed well on the holdout test set of 1000 reactions across ten different functional groups. The mean absolute error (MAE) was 2.44 kcal/mol. Further details on the model's performance with this smaller data set are provided in Section S6 of the Supporting Information.

Compared to other studies using molecular GNNs for property prediction,^{35,40,70} this training data set is relatively small. However, to the best of our knowledge, there are no publicly available data sets specifically for hydrolysis reactions. To evaluate the impact of additional training data, we curated 24,742 more reactions from the Alchemy database, focusing on molecules with 10, 11, or 12 heavy atoms. This expansion also increased the variety of hydrolyzing functional groups from 10 to 13. The expanded data set was randomly split into training and testing sets at roughly a 9:1 ratio, resulting in a holdout test set of 3800 reactions. As shown in Figure 4a, the model generalized effectively on this test set, with the MAE improving to 1.73 kcal/mol. The parity plot comparing model predictions with DFT labels for the test set, shown in Figure 4c, demonstrates a high coefficient of determination (R^2) of 0.96. The distribution of deviations between model predictions and DFT labels,

illustrated in Figure 4b, indicates that errors are closely centered around a mean of zero kcal/mol. A detailed breakdown of these errors is provided in Table 1, showing that only 53 out of 3800

Table 1. Error Distribution for the QM9 + Alchemy Holdout Test Set

absolute error (kcal/mol)	counts
<2	2674
>2 and <5	893
>5 and <10	180
>10	53

reactions had prediction errors exceeding 10 kcal/mol, thus suggesting the model is suitable for screening purposes in this regime. Section S7 of the Supporting Information includes examples of five such outlier reactions, which often feature unusual structures with multiple strained rings, potentially contributing to the larger prediction errors. Another important aspect of evaluating the model's applicability is its capability to classify the overall thermodynamic feasibility of hydrolysis reactions as either exergonic or endergonic, based on the DFT labels. The model correctly classifies 97.1% of the reactions in the test set, demonstrating its strong predictive power in distinguishing between these two thermodynamic outcomes. Among the 117 misclassified reactions, a significant proportion (72) had DFT-calculated ΔG_r values close to zero, highlighting the inherent difficulty of correctly classifying these reactions.

Table 2 summarizes MAE statistics by functional group, revealing that no specific functional group performs poorly.

Table 2. MAE Statistics Based on Functional Group Hydrolyzed

functional group	MAE (kcal/mol)
lactone	2.198
nitrile	1.433
lactam	2.408
imide	2.498
nitrogen–sulfur cleavage	2.176
enamine	2.098
amide	1.724
cyclic carbonate	1.252
aliphatic fluoride	1.022
carboxylic acid ester	1.436
carbamate	1.591
urea	1.091
epoxide	1.571
overall average	1.731

Notably, functional groups like lactones and lactams, which exhibit a broader range of ΔG_r as shown in Figure 3a, tend to have higher MAEs. The higher MAE for imides may be attributed to their lower representation in the database (165 out of 41,006 reactions). However, interestingly, some functional groups with lower representation, such as aliphatic fluoride (191) and cyclic carbonate (116), show lower MAEs compared to the model average.

To assess our model's performance relative to other reaction-property prediction algorithms, we benchmarked it against several models. As detailed in Section 2.3, our model is highly generalizable and capable of handling reactions with varying numbers of bond changes—a feature not commonly found in

reaction-property algorithms. This limitation narrowed the range of models suitable for benchmarking. We tested a simple reactant-only graph neural network with both atom and bond features, incorporating standard cheminformatics features such as bond degree, element identity, atomic weight, ring inclusion, and hybridization, as well as global features such as the number of atoms and bonds, molecular weight, and one-hot encoding for the hydrolyzing functional group and charge. Additionally, we evaluated an XGBoost model with Morgan Fingerprints and Chemprop,⁷⁰ another modern algorithm. Both the XGBoost and Chemprop models were tuned using Bayesian optimization before final testing. The performance of these models is summarized in Table 3. Although Chemprop performs

Table 3. Performance Comparison against Benchmark Models for QM9 + Alchemy Dataset^a

model	test MAE (kcal/mol)	test RMSE (kcal/mol)
mean	12.745	14.670
reactant GNN (atom)	4.008	5.429
reactant GNN (atom + bond)	3.445	4.875
XGB + morgan	2.448	3.705
chemprop	2.257	3.528
HEPOM	1.731	2.674

^aWe also include a benchmark to trivially guessing the mean of the training set.

competitively (MAE: 2.25 kcal/mol vs 1.73 kcal/mol), our model outperforms all benchmarked models in terms of performance metrics on the holdout test set. The performance of individual benchmark models is shown in Figure 5.

3.3. Model Embeddings—Neutral Data Set. Visualizing the feature space provides insight into the underlying patterns the model learns during training. To analyze the learned representations for the trained model, we extracted high-dimensional difference feature vectors for each test set reaction before they are implemented into the fully connected layer for prediction. These vectors were then reduced to a two-dimensional (2D) space using the Uniform Manifold Approximation and Projection (UMAP) method.⁷¹ The evolution of these 2D embeddings at different epochs during training is shown in Figure S8 of the Supporting Information.

Initially, the embeddings are loosely clustered based on the functional groups of the hydrolyzing reactants, as expected. However, as training progresses, these clusters become tighter and more defined, reflecting not only functional groups but also other underlying chemical similarities not explicitly known to the model. Figure 6 illustrates the final 2D representations of the feature vectors for the test set, each tagged with its respective hydrolyzing functional group. In addition to clustering by functional groups, a clear distinction emerges between the embeddings of uniprotect reactions and those of biproduct and triproduct reactions. Uniprotect reactions predominantly cluster on one side of the feature vector space, while reactions yielding more than one product aggregate oppositely. For uniprotect reactions, the model forms a distinct cluster for cyclic functional groups (lactams, lactones, and imides). This suggests that the model identifies additional common features beyond functional group similarity, such as ring-opening during hydrolysis.

3.4. Model Performance: QM9⁺ and QM9[−] Data Sets. In Section 2.1, we discussed that in most practical scenarios,

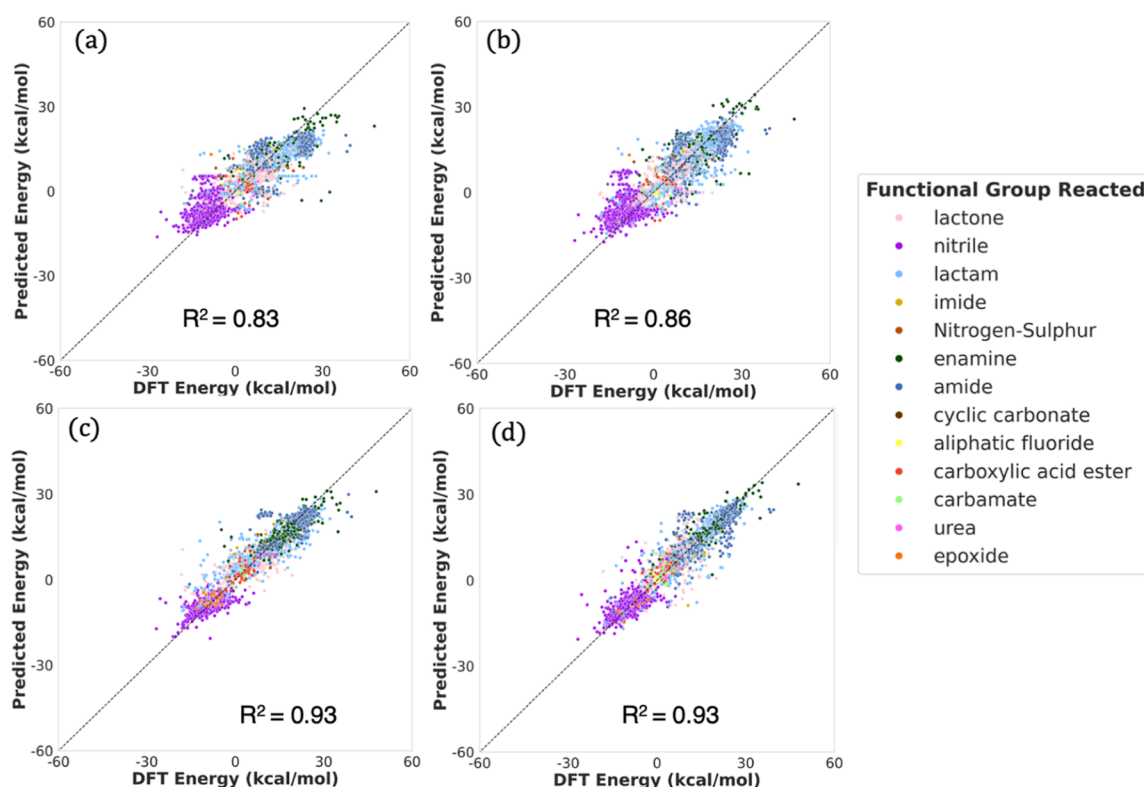


Figure 5. Parity plots for the performance of benchmark models on the holdout test of the QM9 + Alchemy data set. (a) Reactant only GNN—node features; (b) reactant only GNN—node + edge features; (c) XGBoost + Morgan fingerprints; (d) Chemprop.⁷⁰

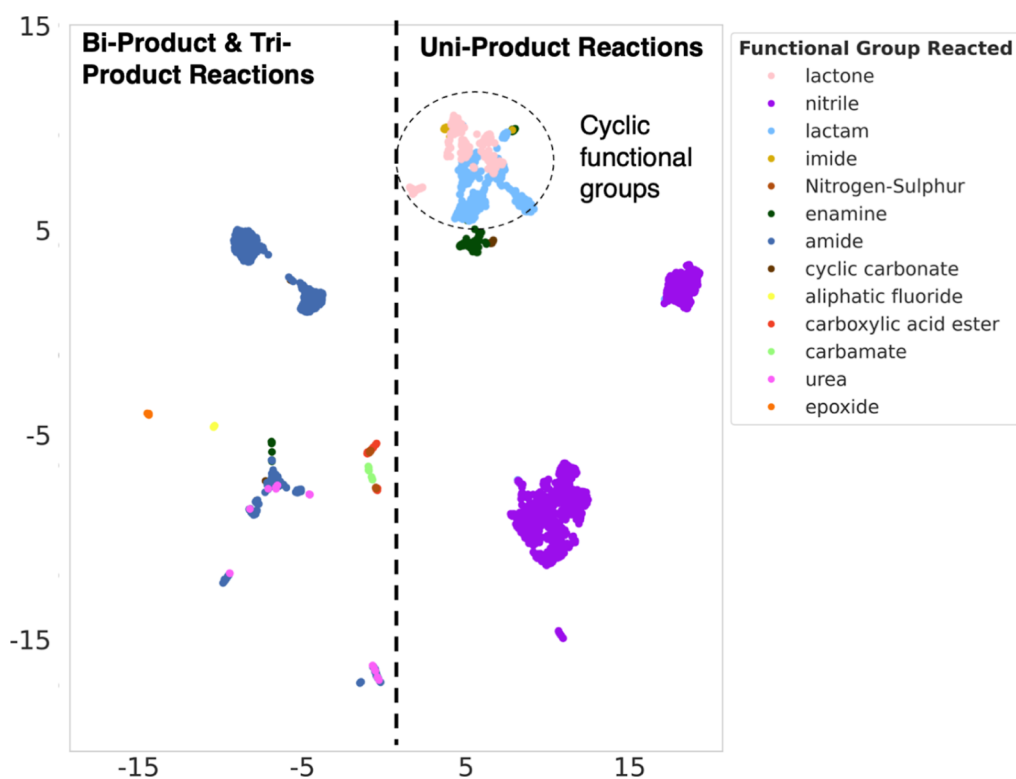


Figure 6. UMAP embeddings of the high-dimensional feature vectors representing the hydrolysis reactions into a two-dimensional space.

hydrolysis occurs in a reaction medium where acidic or basic pH expedites the reaction. To extend the model's applicability, we generated the QM9⁺ and QM9[−] data sets, which include protonated and hydroxylated reactants, respectively, to simulate

extreme pH conditions. In the current work, the protonation/hydroxylation was limited to only the QM9 molecules. Therefore, these charged data sets are considerably smaller than the neutral QM9 + Alchemy data set. The trained models'

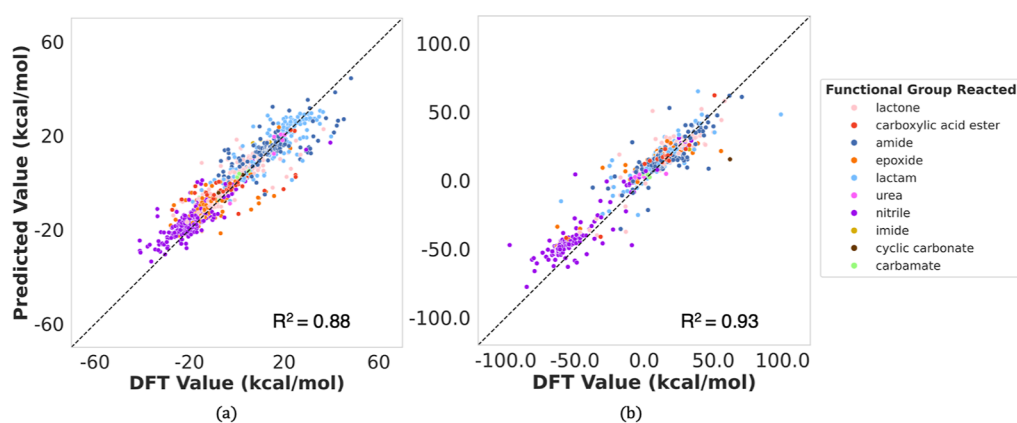


Figure 7. Test Set Performance on models trained with the (a) QM9⁺ protonated and (b) QM9[−] hydroxylated data sets.

performance on the holdout test sets is shown in Figure 7 and summarized in Table 4.

Table 4. Performance Comparison against Benchmark Models for QM9⁺ and QM9[−] Datasets, Best Model Is Bolded

model	QM9 ⁺ data set		QM9 [−] data set	
	test MAE (kcal/mol)	test RMSE (kcal/mol)	test MAE (kcal/mol)	test RMSE (kcal/mol)
mean	15.234	17.381	34.831	36.919
XGB + Morgan	5.394	8.195	8.687	14.375
chemprop	6.275	8.864	5.373	9.682
HEPOM	4.282	6.213	6.607	9.326

As expected, the model's performance on these data sets deteriorates, evidenced by the lower coefficient of determination (R^2) and higher mean absolute errors (MAEs) for both the QM9⁺ and QM9[−] test sets. This performance decrease is particularly pronounced for the hydroxylated model, which shows a relatively high MAE of 6.607 kcal/mol. However, it is important to contextualize this result by noting that the range of ΔG_r values in this data set is radically different, roughly spanning between −150 and 100 kcal/mol, compared to the QM9 + Alchemy data sets (−40 to 40 kcal/mol). Notably, the MAE value also corresponds to a strong R^2 of 0.93.

Given these differences, a fair comparison of model performance should be made with relevant benchmarks rather than the neutral data set. For these two data sets, we conducted hyperparameter optimization using XGBoost and Chemprop models. As shown in Table 4, our model outperforms the benchmarks in most metrics in this smaller and more complex data set. The exception is the mean absolute error (MAE) on the QM9[−] data set, where Chemprop slightly outperforms HEPOM. Interestingly, despite the higher MAE for HEPOM on the QM9[−] test set, it demonstrates greater robustness to outlier predictions, achieving a lower root mean squared error (RMSE) compared to Chemprop. Furthermore, HEPOM significantly outperforms Chemprop in correctly classifying the thermodynamic feasibility (endergonic vs exergonic) of reactions in the QM9[−] data set, achieving a classification accuracy of 95.5% compared to Chemprop's 84.8%. XGBoost has the lowest classification accuracy of 76.3% on the same test set. Additional details, including parity plots for the benchmarks and MAE statistics by functional group for these holdout test sets, are provided in Section S9 of the Supporting Information.

3.5. Combined Model Training. For the neutral data set, we observed a significant improvement in model performance after incorporating additional data from the Alchemy data set. Given this result, it can be expected that the performance of the protonated (QM9⁺) and hydroxylated (QM9[−]) models would

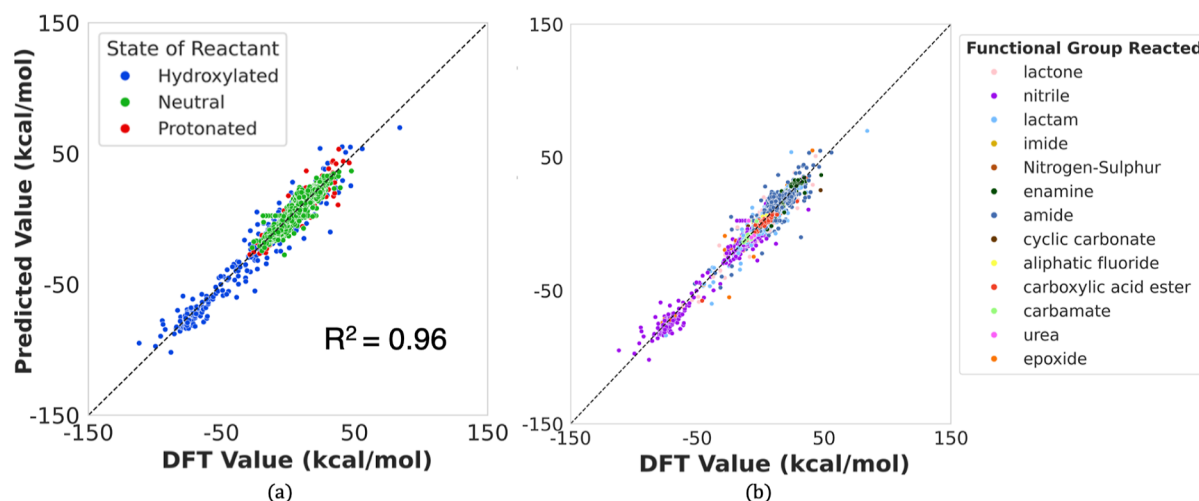


Figure 8. Test Set Performance on models trained with the combined data set depicted as (a) the charged state of the reactant and (b) the hydrolyzing functional group.

also benefit from more data. However, obtaining this would require another round of computationally intensive data curation. Instead, we chose to augment the data sets by combining all three data sets. We found that adding the neutral data significantly improved the performance of both charged models, as shown in Figure 8.

Importantly, this improvement was not limited to our models but was also observed in the benchmark models, particularly Chemprop. This suggests that the higher MAEs of the QM9⁺ and QM9⁻ models may be simply due to insufficient data. However, as seen in Table 5, the most pronounced improvement was in our HEPOM model, which allowed it to surpass the other benchmarks, including in the previously higher MAE for the QM9⁻ test set (Table 4).

Table 5. Performance Comparison against Benchmark Models the Combined Dataset, Best Model Is Bolded

model	test MAE (kcal/mol)	test RMSE (kcal/mol)
XGB + Morgan	9.998	14.756
chemprop	4.920	8.562
HEPOM	3.054	4.281

The parity plots in Figure 8a and the MAE split based on the reactant state, compiled in Table 6, show that the improved

Table 6. Our Model MAE Stratified on the Combined Test Set Based on the Reactant State

state of reactant	mean absolute error (kcal/mol)
hydroxylated	4.057
neutral	2.805
protonated	2.838

model performance for the charged data sets comes at the cost of a slightly higher MAE for the neutral test set. The corresponding parity plots for the benchmark models are included in Section S10 of the Supporting Information.

4. CONCLUSION

In this work, we combined reaction templates and high-throughput DFT calculations to generate a large and diverse data set of ΔG_r values of hydrolytic pathways for molecules selected from two popular molecular databases (QM9 and Alchemy). We then used this data set to train a custom message-passing GNN on the difference features of the products and reactants, resulting in a model capable of predicting the thermodynamic feasibility (ΔG_r) of hydrolysis reactions. The model demonstrates remarkable accuracy on the neutral data set of hydrolysis reactions and outperforms benchmark models on smaller, more complex data sets involving charged reactants, simulating extreme pH conditions. In addition, by combining all three of our data sets, we find that our model is able to reasonably predict across all three classes at once.

We believe that this model is valuable for high-throughput screening of molecules and automated chemical synthesis in various domains, including drug development, environmental chemistry, and chemical deconstruction. The comprehensive data set developed in this work also serves as a critical resource for training other machine learning models. In terms of the model, although this study focuses on hydrolysis, the model can be easily fine-tuned for any arbitrary reaction data sets with available reactant and product molecule graphs.

Training and holdout test sets for all models are publicly accessible via Figshare, and detailed information about the reactant and product molecules for the QM9 database is available via the MPCules⁷² interface, with future plans to integrate the Alchemy data set reactants and products as well. The code for training the model can be accessed at the GitHub repository.

■ ASSOCIATED CONTENT

Data Availability Statement

The training and holdout test sets for all models are publicly available via Figshare: https://figshare.com/articles/dataset/Hydrolysis_Datasets_from_HEPOM_paper_/27851130. The code for training the model is available on the HEPOM GitHub repository: <https://github.com/HEPOM/HEPOM>. Model hyperparameters for each of the trained models are available in Section S11 of the Supporting Information

Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jcim.4c02443>.

Additional details about data set generation, data set statistics, model performance against benchmarks, and model hyperparameters (PDF)

■ AUTHOR INFORMATION

Corresponding Author

Kristin A. Persson – Department of Materials Science and Engineering, University of California, Berkeley, California 94720, United States; Materials Science Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, United States; orcid.org/0000-0003-2495-5509; Email: kristinpersson@berkeley.edu

Authors

Rishabh D. Guha – Materials Science Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, United States; Present Address: Schrödinger Inc. 1540 Broadway, New York, NY 10024, USA; orcid.org/0000-0003-1977-6039

Santiago Vargas – Chemical Sciences Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, United States

Evan Walter Clark Spotte-Smith – Department of Materials Science and Engineering, Carnegie Mellon University, Pittsburgh, Pennsylvania 15213, United States; orcid.org/0000-0003-1554-197X

Alexander Rizzolo Epstein – Department of Materials Science and Engineering, University of California, Berkeley, California 94720, United States

Maxwell Venetos – Materials Science Division, Lawrence Berkeley National Laboratory, Berkeley, California 94720, United States; Department of Materials Science and Engineering, University of California, Berkeley, California 94720, United States

Ryan Kingsbury – Department of Civil and Environmental Engineering, Princeton University, Princeton, New Jersey 08544, United States; orcid.org/0000-0002-7168-3967

Mingjian Wen – Institute of Fundamental and Frontier Sciences, University of Electronic Science and Technology of China, Chengdu 610054, China; orcid.org/0000-0003-0013-575X

Samuel M. Blau – Energy Storage and Distributed Resources,
Lawrence Berkeley National Laboratory, Berkeley, California
94720, United States; orcid.org/0000-0003-3132-3032

Complete contact information is available at:
<https://pubs.acs.org/10.1021/acs.jcim.4c02443>

Author Contributions

◆R.D.G. and S.V. contributed equally to this work. Rishabh D. Guha: Conceptualization, data curation, formal analysis, investigation, methodology, software, visualization, writing—original draft, writing—review and editing, Santiago Vargas.: Software, methodology, formal analysis, investigation, visualization, writing—original draft, writing—review and editing, Evan Walter Clarke Spotte-Smith: Conceptualization, investigation, methodology, writing—review and editing, Alexander Rizzolo Epstein: Conceptualization, investigation, methodology, writing—review and editing, Maxwell Venetos: Conceptualization, investigation, methodology, writing—review and editing, Ryan Kingsbury: Conceptualization, methodology, funding acquisition, writing—review and editing, Mingjian Wen: Software, writing—review and editing, Samuel M. Blau: Conceptualization, methodology, funding acquisition, writing—review and editing, Kristin A. Persson: Conceptualization, methodology, funding acquisition, project administration, resources, supervision, writing—review and editing.

Notes

The authors declare no competing financial interest.

ACKNOWLEDGMENTS

This work was intellectually led by the Moore Foundation Grant, which is funded by the Gordon and Betty Moore Foundation, under Grant no. 10454. Additional support was provided by the U.S. Department of Energy, Office of Science, Office of Basic Energy Sciences, Materials Sciences and Engineering Division under contract no. DE-AC02-05-CH11231, Unlocking Chemical Circularity in Recycling by Controlling Polymer Reactivity across Scales Program CUP-LBL- Helms. Data for this study were produced using computational resources provided by Eagle and Swift high-performance computing (HPC) systems at the National Renewable Energy Laboratory and the Savio HPC cluster at the University of California, Berkeley. The GNN models were trained on the Eagle HPC. This research was supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Department of Energy Computational Science Graduate Fellowship under Award Number DE-SC0020347.

REFERENCES

- (1) Breynaert, E.; Houlléberghs, M.; Radhakrishnan, S.; Grübel, G.; Taulelle, F.; Martens, J. A. Water as a tuneable solvent: a perspective. *Chem. Soc. Rev.* **2020**, *49*, 2557–2569.
- (2) Idolor, O.; Guha, R. D.; Berkowitz, K.; Grace, L. An experimental study of the dynamic molecular state of transient moisture in damaged polymer composites. *Polym. Compos.* **2021**, *42*, 3391–3403.
- (3) Franks, F. *Water in Crystalline Hydrates Aqueous Solutions of Simple Nonelectrolytes*; Springer, 1973; pp 1–54.
- (4) Pohorille, A.; Pratt, L. R. Is water the universal solvent for life? *Origins Life Evol. Biospheres* **2012**, *42*, 405–409.
- (5) Butler, R. N.; Coyne, A. G. Water: Nature's Reaction Enforcer Comparative Effects for Organic Synthesis "In-Water" and "On-Water". *Chem. Rev.* **2010**, *110*, 6302–6337.

- (6) Gorb, L.; Asensio, A.; Tuñón, I.; Ruiz-López, M. F. The mechanism of formamide hydrolysis in water from ab initio calculations and simulations. *Chem. - Eur. J.* **2005**, *11*, 6743–6753.
- (7) Arumugam, P.; Gruber, S.; Tanaka, K.; Haering, C. H.; Mechtler, K.; Nasmyth, K. ATP Hydrolysis Is Required for Cohesin's Association with Chromosomes. *Curr. Biol.* **2003**, *13*, 1941–1953.
- (8) Bashkin, J. K. Hydrolysis of phosphates, esters and related substrates by models of biological catalysts. *Curr. Opin. Chem. Biol.* **1999**, *3*, 752–758.
- (9) Blazek, J.; Gilbert, E. P. Effect of Enzymatic Hydrolysis on Native Starch Granule Structure. *Biomacromolecules* **2010**, *11*, 3275–3289.
- (10) Helms, B. A. Polydiketoenamides for a Circular Plastics Economy. *Acc. Chem. Res.* **2022**, *55*, 2753–2765.
- (11) Le Feunteun, S.; Verkempinck, S.; Floury, J.; Janssen, A.; Kondjoyan, A.; Marze, S.; Mirade, P.-S.; Pluschke, A.; Sicard, J.; Van Aken, G.; et al. Mathematical modelling of food hydrolysis during in vitro digestion: From single nutrient to complex foods in static and dynamic conditions. *Trends Food Sci. Technol.* **2021**, *116*, 870–883.
- (12) Meng, X.; Guo, Y.; Wang, Y.; Fan, S.; Wang, K.; Han, W. A systematic review of photolysis and hydrolysis degradation modes, degradation mechanisms, and identification methods of pesticides. *J. Chem.* **2022**, *2022*, 9552466.
- (13) Demarteau, J.; Epstein, A. R.; Christensen, P. R.; Abubekarov, M.; Wang, H.; Teat, S. J.; Seguin, T. J.; Chan, C. W.; Scown, C. D.; Russell, T. P.; Keasling, J. D.; Persson, K. A.; Helms, B. A. Circularity in mixed-plastic chemical recycling enabled by variable rates of polydiketoenamide hydrolysis. *Sci. Adv.* **2022**, *8*, No. eabp8823.
- (14) Olsson, E.; Menzel, C.; Johansson, C.; Andersson, R.; Koch, K.; Järnström, L. The effect of pH on hydrolysis, cross-linking and barrier properties of starch barriers containing citric acid. *Carbohydr. Polym.* **2013**, *98*, 1505–1513.
- (15) Mitchell, S. M.; Ullman, J. L.; Teel, A. L.; Watts, R. J. pH and temperature effects on the hydrolysis of three β -lactam antibiotics: Ampicillin, cefalotin and cefoxitin. *Sci. Total Environ.* **2014**, *466*, 547–555.
- (16) Epstein, A. R.; Demarteau, J.; Helms, B. A.; Persson, K. A. Variable Amine Spacing Determines Depolymerization Rate in Polydiketoenamides. *J. Am. Chem. Soc.* **2023**, *145*, 8082–8089.
- (17) Marenich, A. V.; Cramer, C. J.; Truhlar, D. G. Universal Solvation Model Based on Solute Electron Density and on a Continuum Model of the Solvent Defined by the Bulk Dielectric Constant and Atomic Surface Tensions. *J. Phys. Chem. B* **2009**, *113*, 6378–6396.
- (18) Bergström, C. A.; Luthman, K.; Artursson, P. Accuracy of calculated pH-dependent aqueous drug solubility. *Eur. J. Pharm. Sci.* **2004**, *22*, 387–398.
- (19) Grisuta, B.; Janssen, L. P. B. M.; Heeres, H. J. Kinetic Study on the Acid-Catalyzed Hydrolysis of Cellulose to Levulinic Acid. *Ind. Eng. Chem. Res.* **2007**, *46*, 1696–1708.
- (20) Carlson, D. L.; Than, K. D.; Roberts, A. L. Acid- and Base-Catalyzed Hydrolysis of Chloroacetamide Herbicides. *J. Agric. Food Chem.* **2006**, *54*, 4740–4750.
- (21) Jorner, K.; Brinck, T.; Norrby, P.-O.; Buttar, D. Machine learning meets mechanistic modelling for accurate prediction of experimental activation energies. *Chem. Sci.* **2021**, *12*, 1163–1175.
- (22) Hirao, H.; Que, L., Jr.; Nam, W.; Shaik, S. A Two-State Reactivity Rationale for Counterintuitive Axial Ligand Effects on the C H Activation Reactivity of Nonheme FeIV O Oxidants. *Chem. - Eur. J.* **2008**, *14*, 1740–1756.
- (23) Epstein, A. R.; Spotte-Smith, E. W. C.; Venetos, M. C.; Andriuc, O.; Persson, K. A. Assessing the Accuracy of Density Functional Approximations for Predicting Hydrolysis Reaction Kinetics. *J. Chem. Theory Comput.* **2023**, *19*, 3159–3171.
- (24) Malick, D. K.; Petersson, G. A.; Montgomery, J. A. Transition states for chemical reactions I. Geometry and classical barrier height. *J. Chem. Phys.* **1998**, *108*, S704–S713.
- (25) Evans, M. G.; Polanyi, M. Further considerations on the thermodynamics of chemical equilibria and reaction rates. *Trans. Faraday Soc.* **1936**, *32*, 1333.

- (26) Stuyver, T.; Coley, C. W. Machine Learning-Guided Computational Screening of New Candidate Reactions with High Bioorthogonal Click Potential. *Chem. - Eur. J.* **2023**, *29*, No. e202300387.
- (27) Zhou, S.; Nguyen, B. T.; Richard, J. P.; Kluger, R.; Gao, J. Origin of Free Energy Barriers of Decarboxylation and the Reverse Process of CO₂ Capture in Dimethylformamide and in Water. *J. Am. Chem. Soc.* **2021**, *143*, 137–141.
- (28) Lawson, K. E.; Dekle, J. K.; Adamczyk, A. J. Towards pharmaceutical protein stabilization: DFT and statistical learning studies on non-enzymatic peptide hydrolysis degradation mechanisms. *Comput. Theor. Chem.* **2022**, *1218*, 113938.
- (29) Mardirossian, N.; Head-Gordon, M. Thirty years of density functional theory in computational chemistry: an overview and extensive assessment of 200 density functionals. *Mol. Phys.* **2017**, *115*, 2315–2372.
- (30) Ribeiro, A. J. M.; Ramos, M. J.; Fernandes, P. A. Benchmarking of DFT Functionals for the Hydrolysis of Phosphodiester Bonds. *J. Chem. Theory Comput.* **2010**, *6*, 2281–2292.
- (31) Cramer, C. J.; Truhlar, D. G. A universal approach to solvation modeling. *Acc. Chem. Res.* **2008**, *41*, 760–768.
- (32) Gilmer, J.; Schoenholz, S. S.; Riley, P. F.; Vinyals, O.; Dahl, G. E. Neural Message Passing for Quantum Chemistry. *arXiv* **2017**, arXiv:1704.01212v2.
- (33) Grambow, C. A.; Pattanaik, L.; Green, W. H. Deep Learning of Activation Energies. *J. Phys. Chem. Lett.* **2020**, *11*, 2992–2997.
- (34) Kearnes, S.; McCloskey, K.; Berndl, M.; Pande, V.; Riley, P. Molecular graph convolutions: moving beyond fingerprints. *J. Comput.-Aided Mol. Des.* **2016**, *30*, 595–608.
- (35) Yang, K.; Swanson, K.; Jin, W.; Coley, C.; Eiden, P.; Gao, H.; Guzman-Perez, A.; Hopper, T.; Kelley, B.; Mathea, M.; Palmer, A.; Settels, V.; Jaakkola, T.; Jensen, K.; Barzilay, R. Analyzing Learned Molecular Representations for Property Prediction. *J. Chem. Inf. Model.* **2019**, *59*, 3370–3388.
- (36) Heid, E.; Green, W. H. Machine Learning of Reaction Properties via Learned Representations of the Condensed Graph of Reaction. *J. Chem. Inf. Model.* **2022**, *62*, 2101–2110.
- (37) Wen, M.; Blau, S. M.; Xie, X.; Dwaraknath, S.; Persson, K. A. Improving machine learning performance on small chemical reaction data with unsupervised contrastive pretraining. *Chem. Sci.* **2022**, *13*, 1446–1458.
- (38) Grambow, C. A.; Pattanaik, L.; Green, W. H. Reactants, products, and transition states of elementary chemical reactions based on quantum chemistry. *Sci. Data* **2020**, *7*, 137.
- (39) Ruddigkeit, L.; Van Deursen, R.; Blum, L. C.; 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.
- (40) Wen, M.; Blau, S. M.; Spotte-Smith, E. W. C.; Dwaraknath, S.; Persson, K. A. BondNet: a graph neural network for the prediction of bond dissociation energies for charged molecules. *Chem. Sci.* **2021**, *12*, 1858–1868.
- (41) Spotte-Smith, E. W. C.; Blau, S. M.; Xie, X.; Patel, H. D.; Wen, M.; Wood, B.; Dwaraknath, S.; Persson, K. A. Quantum chemical calculations of lithium-ion battery electrolyte and interphase species. *Sci. Data* **2021**, *8*, 203.
- (42) Ramakrishnan, R.; Dral, P. O.; Rupp, M.; von Lilienfeld, O. A. Quantum chemistry structures and properties of 134 kilo molecules. *Sci. Data* **2014**, *1*, 140022.
- (43) Chen, G.; Chen, P.; Hsieh, C.-Y.; Lee, C.-K.; Liao, B.; Liao, R.; Liu, W.; Qiu, J.; Sun, Q.; Tang, J.; et al. Alchemy: A quantum chemistry dataset for benchmarking ai models. *arXiv* **2019**, arXiv:1906.09427.
- (44) Landrum, G. *RDKit: A Software Suite for Cheminformatics, Computational Chemistry, and Predictive Modeling*, 2013.
- (45) Tebes-Stevens, C.; Patel, J. M.; Jones, W. J.; Weber, E. J. Prediction of hydrolysis products of organic chemicals under environmental pH conditions. *Environ. Sci. Technol.* **2017**, *51*, 5008–5016.
- (46) Schowen, R. L.; Jayaraman, H.; Kershner, L. Catalytic Efficiencies in Amide Hydrolysis. The Two-Step Mechanism I. *J. Am. Chem. Soc.* **1966**, *88*, 3373–3375.
- (47) Zahn, D. On the role of water in amide hydrolysis. *Eur. J. Org. Chem.* **2004**, *2004*, 4020–4023.
- (48) Jencks, W. P.; Carriuolo, J. General base catalysis of ester Hydrolysis I. *J. Am. Chem. Soc.* **1961**, *83*, 1743–1750.
- (49) Rupp, M.; Korner, R.; V Tetko, I. Predicting the pKa of small molecules. *Comb. Chem. High Throughput Screening* **2011**, *14*, 307–327.
- (50) Baba, T.; Matsui, T.; Kamiya, K.; Nakano, M.; Shigeta, Y. A density functional study on the pKa of small polyprotic molecules. *Int. J. Quantum Chem.* **2014**, *114*, 1128–1134.
- (51) Mansouri, K.; Cariello, N. F.; Korotcov, A.; Tkachenko, V.; Grulke, C. M.; Sprankle, C. S.; Allen, D.; Casey, W. M.; Kleinstreuer, N. C.; Williams, A. J. Open-source QSAR models for pKa prediction using multiple machine learning approaches. *J. Cheminf.* **2019**, *11*, 60.
- (52) Wu, J.; Kang, Y.; Pan, P.; Hou, T. Machine learning methods for pKa prediction of small molecules: Advances and challenges. *Drug Discovery Today* **2022**, *27*, 103372.
- (53) Yang, Q.; Li, Y.; Yang, J.-D.; Liu, Y.; Zhang, L.; Luo, S.; Cheng, J.-P. Holistic prediction of the pKa in diverse solvents based on a machine-learning approach. *Angew. Chem.* **2020**, *132*, 19444–19453.
- (54) Dutra, F. R.; Silva, C. d. S.; Custodio, R. On the Accuracy of the Direct Method to Calculate pKa from Electronic Structure Calculations. *J. Phys. Chem. A* **2021**, *125*, 65–73.
- (55) Epifanovsky, E.; Gilbert, A. T. B.; Feng, X.; Lee, J.; Mao, Y.; Mardirossian, N.; Pokhilko, P.; White, A. F.; Coons, M. P.; Dempwolff, A. L.; Gan, Z.; Hait, D.; Horn, P. R.; Jacobson, L. D.; Kaliman, I.; Kussmann, J.; Lange, A. W.; Lao, K. U.; Levine, D. S.; Liu, J.; McKenzie, S. C.; Morrison, A. F.; Nanda, K. D.; Plasser, F.; Rehn, D. R.; Vidal, M. L.; You, Z.-Q.; Zhu, Y.; Alam, B.; Albrecht, B. J.; Aldossary, A.; Alguire, E.; Andersen, J. H.; Athavale, V.; Barton, D.; Begam, K.; Behn, A.; Bellonzi, N.; Bernard, Y. A.; Berquist, E. J.; Burton, H. G. A.; Carreras, A.; Carter-Fenk, K.; Chakraborty, R.; Chien, A. D.; Closser, K. D.; Cofer-Shabica, V.; Dasgupta, S.; de Wergifosse, M.; Deng, J.; Diedenhofen, M.; Do, H.; Ehlert, S.; Fang, P.-T.; Fatehi, S.; Feng, Q.; Friedhoff, T.; Gayvert, J.; Ge, Q.; Gidofalvi, G.; Goldey, M.; Gomes, J.; González-Espinoza, C. E.; Gulania, S.; Gunina, A. O.; Hanson-Heine, M. W. D.; Harbach, P. H. P.; Hauser, A.; Herbst, M. F.; Hernández Vera, M.; Hodecker, M.; Holden, Z. C.; Houck, S.; Huang, X.; Hui, K.; Huynh, B. C.; Ivanov, M.; Jász, A.; Ji, H.; Jiang, H.; Kaduk, B.; Kähler, S.; Khistyayev, K.; Kim, J.; Kis, G.; Klunzinger, P.; Koczor-Benda, Z.; Koh, J. H.; Kosenkov, D.; Koulis, L.; Kowalczyk, T.; Krauter, C. M.; Kue, K.; Kunitsa, A.; Kus, T.; Ladžanski, I.; Landau, A.; Lawler, K. V.; Lefrançois, D.; Lehtola, S.; Li, R. R.; Li, Y.-P.; Liang, J.; Liebenthal, M.; Lin, H.-H.; Lin, Y.-S.; Liu, F.; Liu, K.-Y.; Loipersberger, M.; Luenser, A.; Manjanath, A.; Manohar, P.; Mansoor, E.; Manzer, S. F.; Mao, S.-P.; Marenich, A. V.; Markovich, T.; Mason, S.; Maurer, S. A.; McLaughlin, P. F.; Menger, M. F. S. J.; Mewes, J.-M.; Mewes, S. A.; Morgante, P.; Mullinax, J. W.; Oosterbaan, K. J.; Parag, G.; Paul, A. C.; Paul, S. K.; Pavšević, F.; Pei, Z.; Prager, S.; Proynov, E. I.; Rák, A.; Ramos-Cordoba, E.; Rana, B.; Rask, A. E.; Rettig, A.; Richard, R. M.; Rob, F.; Rossomme, E.; Scheele, T.; Scheurer, M.; Schneider, M.; Sergueev, N.; Sharada, S. M.; Skomorowski, W.; Small, D. W.; Stein, C. J.; Su, Y.-C.; Sundstrom, E. J.; Tao, Z.; Thirman, J.; Tornai, G. J.; Tsuchimochi, T.; Tubman, N. M.; Veccham, S. P.; Vydrov, O.; Wenzel, J.; Witte, J.; Yamada, A.; Yao, K.; Yeganeh, S.; Yost, S. R.; Zech, A.; Zhang, I. Y.; Zhang, X.; Zhang, Y.; Zuev, D.; Aspuru-Guzik, A.; Bell, A. T.; Besley, N. A.; Bravaya, K. B.; Brooks, B. R.; Casanova, D.; Chai, J.-D.; Coriani, S.; Cramer, C. J.; Cserey, G.; DePrince, A. E., III; DiStasio, R. A., Jr.; Dreuw, A.; Dunietz, B. D.; Furlani, T. R.; Goddard, W. A., III; Hammes-Schiffer, S.; Head-Gordon, T.; Hehre, W. J.; Hsu, C.-P.; Jagau, T.-C.; Jung, Y.; Klamt, A.; Kong, J.; Lambrecht, D. S.; Liang, W.; Mayhall, N. J.; McCurdy, C. W.; Neaton, J. B.; Ochsenfeld, C.; Parkhill, J. A.; Peverati, R.; Rassolov, V. A.; Shao, Y.; Slipchenko, L. V.; Stauch, T.; Steele, R. P.; Subotnik, J. E.; Thom, A. J. W.; Tkatchenko, A.; Truhlar, D. G.; Van Voorhis, T.; Wesolowski, T. A.; Whaley, K. B.; Woodcock, H. L., III; Zimmerman, P. M.; Faraji, S.; Gill, P. M. W.; Head-Gordon, M.; Herbert, J. M.; Krylov, A. I. Software for the frontiers of quantum chemistry: An overview of developments in the Q-Chem 5 package. *J. Chem. Phys.* **2021**, *155*, 084801.

- (56) Mathew, K.; Montoya, J. H.; Faghaninia, A.; Dwarkanath, S.; Aykol, M.; Tang, H.; Chu, I.-h.; Smidt, T.; Bocklund, B.; Horton, M.; Dagdelen, J.; Wood, B.; Liu, Z.-K.; Neaton, J.; Ong, S. P.; Persson, K.; Jain, A. Atomate: A high-level interface to generate, execute, and analyze computational materials science workflows. *Comput. Mater. Sci.* **2017**, *139*, 140–152.
- (57) Mardirossian, N.; Head-Gordon, M. B97M-V: A combinatorially optimized, range-separated hybrid, meta-GGA density functional with VV10 nonlocal correlation. *J. Chem. Phys.* **2016**, *144*, 214110.
- (58) Vydrov, O. A.; Van Voorhis, T. Nonlocal van der Waals density functional: The simpler the better. *J. Chem. Phys.* **2010**, *133*, 244103.
- (59) Hellweg, A.; Rappoport, D. Development of new auxiliary basis functions of the Karlsruhe segmented contracted basis sets including diffuse basis functions (def2-SVPD, def2-TZVPPD, and def2-QVPPD) for RI-MP2 and RI-CC calculations. *Phys. Chem. Chem. Phys.* **2015**, *17*, 1010–1017.
- (60) Bresson, X.; Laurent, T. Residual Gated Graph ConvNets. *arXiv* **2018**, arXiv:1711.07553v2.
- (61) Dwivedi, V. P.; Joshi, C. K.; Luu, A. T.; Laurent, T.; Bengio, Y.; Bresson, X. Benchmarking Graph Neural Networks. *arXiv* **2022**, arXiv:2003.00982v5.
- (62) Vinyals, O.; Bengio, S.; Kudlur, M. Order Matters: Sequence to sequence for sets. *arXiv* **2016**, arXiv:1511.06391v4.
- (63) Lee, J.; Lee, I.; Kang, J. Self-attention graph pooling. In *International Conference on Machine Learning*, 2019; pp 3734–3743.
- (64) Schweidtmann, A. M.; Rittig, J. G.; Weber, J. M.; Grohe, M.; Dahmen, M.; Leonhard, K.; Mitsos, A. Physical pooling functions in graph neural networks for molecular property prediction. *Comput. Chem. Eng.* **2023**, *172*, 108202.
- (65) Zhou, K.; Dong, Y.; Wang, K.; Lee, W. S.; Hooi, B.; Xu, H.; Feng, J. Understanding and Resolving Performance Degradation in Graph Convolutional Networks. *arXiv* **2021**, arXiv:2006.07107v3.
- (66) Rusch, T. K.; Bronstein, M. M.; Mishra, S. A Survey on Oversmoothing in Graph Neural Networks. *arXiv* **2023**, arXiv:2303.10993v1.
- (67) Hall, Jr. H.; Brandt, M.; Mason, R. Hydrolysis rates and mechanisms of cyclic monomers. *J. Am. Chem. Soc.* **1958**, *80*, 6420–6427.
- (68) Yang, J.-C.; Gorenstein, D. G. Contribution of ring strain and the stereoelectronic effect to the hydrolysis of cyclic five-membered ring phosphorus esters. *Tetrahedron* **1987**, *43*, 479–486.
- (69) Krug, J.; Popelier, P.; Bader, R. Theoretical study of neutral and of acid and base-promoted hydrolysis of formamide. *J. Phys. Chem.* **1992**, *96*, 7604–7616.
- (70) Heid, E.; Greenman, K. P.; Chung, Y.; Li, S.-C.; Graff, D. E.; Vermeire, F. H.; Wu, H.; Green, W. H.; McGill, C. J. Chemprop: A Machine Learning Package for Chemical Property Prediction. *J. Chem. Inf. Model.* **2023**, *64*, 9–17.
- (71) McInnes, L.; Healy, J.; Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv* **2020**, arXiv:1802.03426v3.
- (72) Spotte-Smith, E. W. C.; Cohen, O. A.; Blau, S. M.; Munro, J. M.; Yang, R.; Guha, R. D.; Patel, H. D.; Vijay, S.; Huck, P.; Kingsbury, R.; Horton, M. K.; Persson, K. A. A database of molecular properties integrated in the Materials Project. *Digital Discovery* **2023**, *2*, 1862–1882.