

1

2

3

4

5

6 **MAHOMES II: A webserver for predicting if a metal binding site is enzymatic**

7

8 *Authors: Ryan Feehan<sup>1</sup>, Matthew Copeland<sup>1</sup>, Meghan W. Franklin<sup>1,\*</sup>, Joanna S. G.*

9 *Slusky<sup>1,2,†</sup>*

10

11 Affiliations:

12 <sup>1</sup>Center for Computational Biology, The University of Kansas, 2030 Becker Dr.,

13 Lawrence, KS 66047.<sup>2</sup>Department of Molecular Biosciences, The University of Kansas,

14 1200 Sunnyside Ave. Lawrence KS 66045-3101

15 \*Currently at Generate Biomedicines, Somerville, MA, United States.

16 †To whom correspondence should be addressed: [slusky@ku.edu](mailto:slusky@ku.edu)

17

18 *Abstract:*

19 Recent advances have enabled high-quality computationally generated structures for  
20 proteins with no solved crystal structures. However, protein function data remains  
21 largely limited to experimental methods and homology mapping. Since structure  
22 determines function, it is natural that methods capable of using computationally  
23 generated structures for functional annotations need to be advanced. Our laboratory  
24 recently developed a method to distinguish between metalloenzyme and non-enzyme  
25 sites. Here we report improvements to this method by upgrading our physicochemical  
26 features to alleviate the need for structures with sub-angstrom precision and using  
27 machine learning to reduce training data labeling error. Our improved classifier identifies  
28 protein bound metal sites as enzymatic or non-enzymatic with 94% precision and 92%  
29 recall. We demonstrate that both adjustments increased predictive performance and  
30 reliability on sites with sub-angstrom variations. We constructed a set of predicted  
31 metalloprotein structures with no solved crystal structures and no detectable homology  
32 to our training data. Our model had an accuracy of 90 - 97.5% depending on the quality  
33 of the predicted structures included in our test. Finally, we found the physicochemical  
34 trends that drove this model's successful performance were local protein density,  
35 second shell ionizable residue burial, and the pocket's accessibility to the site. We  
36 anticipate that our model's ability to correctly identify catalytic metal sites could enable  
37 identification of new enzymatic mechanisms and improve *de novo* metalloenzyme  
38 design success rates.

39

40 *Keywords:* Enzymes, Metalloenzymes, Metalloproteins, Machine Learning

41  
42 *Significance statement:* Identification of enzyme active sites on proteins with unsolved  
43 crystallographic structures can accelerate discovery of novel biochemical reactions,  
44 which can impact healthcare, industrial processes, and environmental remediation. Our  
45 lab has developed an ML tool for predicting sites on computationally generated protein  
46 structures as enzymatic and non-enzymatic. We have made our tool available on a  
47 webserver, allowing the scientific community to rapidly search previously unknown  
48 protein function space.

49

50 *Abbreviations footnote:*

51 ML = machine learning

52 RBF = Radial Basis Function

53 CV = cross validation

54 MAHOMES = metal activity heuristic of metalloprotein and enzyme sites

55 DROPP = distribution overlap of a physicochemical property

56 PDB = Protein Data Bank

57 PDE = probability density estimate

58 pLDDT = predicted Local Distance Difference Test

59 MCC = Mathews correlation coefficient

60 TNR = true negative rate

61 TN = true negative

62 TP = true positive

63 FN = false negative

64 FP = false positive

## 65 1. Introduction

66 Enzymes are biological catalysts that are known to lower activation energy for over  
67 8,000 reactions (McDonald, Boyce, and Tipton 2009). Furthermore, enzymes can  
68 increase reaction rates by factors of up to  $10^{17}$ -fold (Wolfenden, Ridgway, and Young  
69 1998). Enzymes are also becoming increasingly prevalent in industrial processes due to  
70 their greener chemistry (Sheldon and Woodley 2018). Despite the importance and  
71 extent of enzymatic research, a reproducible physicochemical basis of catalysis remains  
72 elusive. This unknown limits *de novo* enzyme design or even reliable identification of  
73 enzyme active sites from structure.

74 We have recently used protein structure-based machine learning (ML) to distinguish  
75 between very similar sites, metalloenzyme active sites and protein sites that bind metals  
76 without any enzyme activity (Feehan, Franklin, and Slusky 2021). Our model, metal  
77 activity heuristic of metalloproteins and enzyme sites (MAHOMES), uses an extra-trees  
78 algorithm to achieve better performance metrics than available enzyme function  
79 predictors. We attribute the classifier's success to training on structural physicochemical  
80 properties of similar sites. By training on negative sites that were also in pockets and  
81 also coordinated metals, rather than on all other sites on the protein, our classifier was  
82 able to assign feature importance based on characteristics that were particular to  
83 enzyme activity.

84 Using protein structure-based features enabled MAHOMES to focus on learning  
85 physicochemical properties related to catalysis but it relied on structurally determined  
86 proteins for its input. The PDB only has ~200,000 solved protein structures (Burley et al.

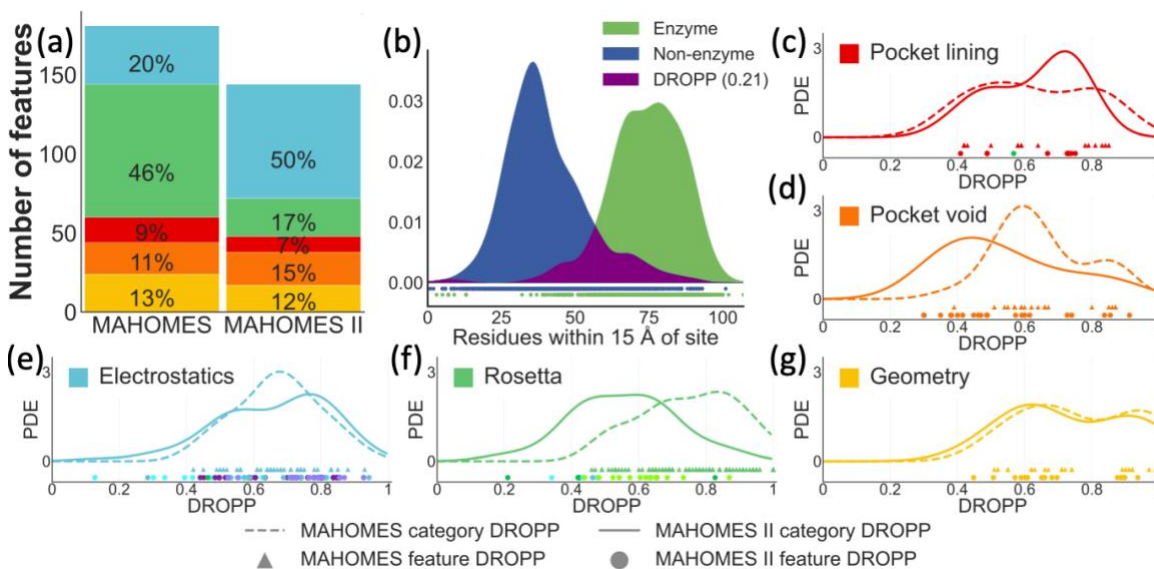
87 2019) thereby limiting MAHOMES utility. Recently, the ML tool AlphaFold2 generated  
88 quality protein structure predictions (Jumper et al. 2021) and now two hundred million  
89 predicted structures are available for download from AlphaFoldDB (Tunyasuvunakool et  
90 al. 2021). However, it remained unclear if these structures could be used for identifying  
91 catalytic sites. This concern was compounded by the finding that a relatively low  
92 percentage of AlphaFold models have a high enough confidence to be recommended  
93 for characterizing binding sites (Thornton, Laskowski, and Borkakoti 2021).

94 To test usage of the computationally generated structures, we updated the calculation  
95 methods used for several of MAHOMES features to reduce the need for sub-angstrom  
96 accuracy. Then, we use the new features when cross validating ML models to reduce  
97 labeling error in our training data labels. The improved features and training data were  
98 used with a variety of ML classifiers and techniques. We found that both the feature  
99 improvement and reduced labeling error led to increased performance for ML models.  
100 Our best ML model, MAHOMES II, outperformed its predecessor on our holdout test-set  
101 with 94% precision and 92% recall. Furthermore, MAHOMES II's predictions were more  
102 reliable for different input structures of the same site with sub-angstrom differences. We  
103 evaluated MAHOMES II on a new set of predicted metalloprotein structures, where it  
104 scored 97.5% accuracy on high confidence structures. Finally, we examined the  
105 features that MAHOMES II found to be the most important for making successful  
106 enzyme or non-enzyme predictions and found a preference for features describing  
107 enzyme sites to be densely packed, have buried second shell residues, and pockets  
108 that were highly accessible to the metal. MAHOMES II can be accessed online

109 (<https://mahomes.ku.edu>), allowing easy use for the scientific community, regardless of  
 110 computational expertise.

## 111 2. Results

**Figure 1: Feature category input space and DROPP** (a) The number of features used by MAHOMES and MAHOMES II for each feature category: blue for electrostatics, green for Rosetta energy terms, red for pocket lining, orange for pocket void, and yellow for coordination geometry. (b) Example of DROPP calculation for the number of residues within 15 Å of the site feature. Kernel density estimators for the feature's dataset values of enzymes (green) and non-enzymes (blue) are made and the overlapping region (purple) is calculated to give the features DROPP. (c-g) Comparison of MAHOMES and MAHOMES II DROPP probability density estimate (PDE) for each feature category: (c) pocket lining category, (d) pocket void category, (e) electrostatics category, (f) Rosetta energy terms category, and (g) geometry category. Dotted lines and triangles represent MAHOMES features. Solid lines and circles represent features used by MAHOMES II. Circles are colored by feature groups shown in figure 4b.



### 112 2.1 New and improved feature calculations

113 To transform metal sites into input for ML algorithms, we identified features belonging to  
 114 five categories which have previously been linked to enzymatic activity— coordination  
 115 geometry, electrostatics, pocket lining, pocket void, and Rosetta energy terms (Figure 1  
 116 a and b).

117 We use a metric to quantify how much a feature's values are similar between enzymatic  
118 and non-enzymatic sites. This metric, DROPP (distribution overlap of a physicochemical  
119 property) identifies how similar the distribution is between enzymatic and non-enzymatic  
120 sites(Figure 1b)(see methods). DROPP was previously found to be lower for features  
121 that are more important for predicting enzyme sites.(Feehan, Franklin, and Slusky  
122 2021). Therefore, when trying to improve our features, we used DROPP as indicator of  
123 feature improvement (Figure 1 c-g). We made efforts to improve all feature classes  
124 except coordination geometry, though some improvements were more successful than  
125 others.

126 *Electrostatics features expansion:* The most important feature used by the original  
127 MAHOMES model was an electrostatic feature (Feehan, Franklin, and Slusky 2021) ,  
128 which was the mean second moment of the of the theoretical titration curve's first  
129 derivative for ionizable residues in the second shell (3.5-9Å). We modeled this feature  
130 after the THEMATIC calculations, which have been used to identify catalytic residues  
131 due to their deviations from Henderson Hasselback titration behavior (Somarowthu,  
132 Yang, et al. 2011; Tong et al. 2009; Ko et al. 2005).

133 To improve our enzyme activity predictions and further our understanding of  
134 electrostatic properties responsible for catalytic activity, we expanded our electrostatic  
135 features category from 37 features to 152 features (Sup. Figure S1). To further  
136 investigate the success of electrostatics in MAHOMES, we added the Z-score  
137 calculations which are used by THEMATICs to measure the relative deviation of the  
138 theoretical titration curve's first derivative's second, third, and fourth moments (Ko et al.  
139 2005). We also added variables output by the generalized Born program we use for



140 generating theoretical titration curves, BLUUES (Fogolari et al. 2012b). Moreover, since  
141 catalytic residues often show interesting shifts in pKa (Pérez-Cañadillas et al. 1998;  
142 Bate and Warwicker 2004), we added features for the pKa shift from ideal amino acid  
143 values, the pKa shift due to desolvation, the pKa shift due to the interaction with other  
144 charges in the molecule with all titratable sites in their neutral state, and the pKa shift  
145 due to the interaction between titratable sites. Additional added features in the  
146 electrostatic category are the generalized Born atomic radii, a solvation exposure  
147 parameter, and solvation energies. After removing redundant features (see methods),  
148 72 electrostatic features are used by MAHOMES II (Figure 1a). Six of these new  
149 features had lower DROPP than any of the 37 previously used electrostatic features  
150 (Figure 1e).

151 *Rosetta features reduction:* In contrast to the expansion of the electrostatic feature  
152 space, we reduced the Rosetta feature space while also improving the features and  
153 improving our model reproducibility. We calculated Rosetta features in MAHOMES  
154 based on spheres with defined radii from the center of the site. In benchmarking  
155 MAHOMES, we found that sub-angstrom differences between relaxed structures of the  
156 same site caused large shifts in Rosetta feature values. To prevent sub-angstrom  
157 differences from significantly changing calculated feature values for the same site, we  
158 switched to a radial basis function (RBF) calculation for the Rosetta energy term  
159 features. The RBF calculation uses distance to weight each residue's influence on the  
160 calculated feature, which prevents subtle changes in the structure from having a  
161 significant impact on the calculated value. The RBF Rosetta energy terms category

162 decreased DROPP (Figure 1f) and the number of features used by the category (Figure  
163 1a).

164 *Pocket void and pocket lining improvements:* Our previous method, Rosetta pocket  
165 measure (Johnson and Karanicolas 2013), did not detect surface pockets for 645  
166 dataset sites, therefore 19% of the MAHOMES training data was missing values for  
167 pocket void and pocket lining features. GHECOM (Kawabata 2019, 2010), a tool that  
168 uses mathematical morphology for finding multi-scale pockets on protein surfaces,  
169 generated pocket for 99.5% of the dataset sites. To improve the quality of training data,  
170 we removed dataset sites that did not have pocket. Additionally, we added various  
171 pocket descriptors, including output features from GHECOM which describe the  
172 pocket's shallowness and size rank relative other pockets on the structure. Ultimately,  
173 the pocket output by GHECOM lowered the DROPP for features in both the pocket  
174 lining and pocket categories (Figure 1c and 1d).

## 175 2.2 Reduced training data labeling error

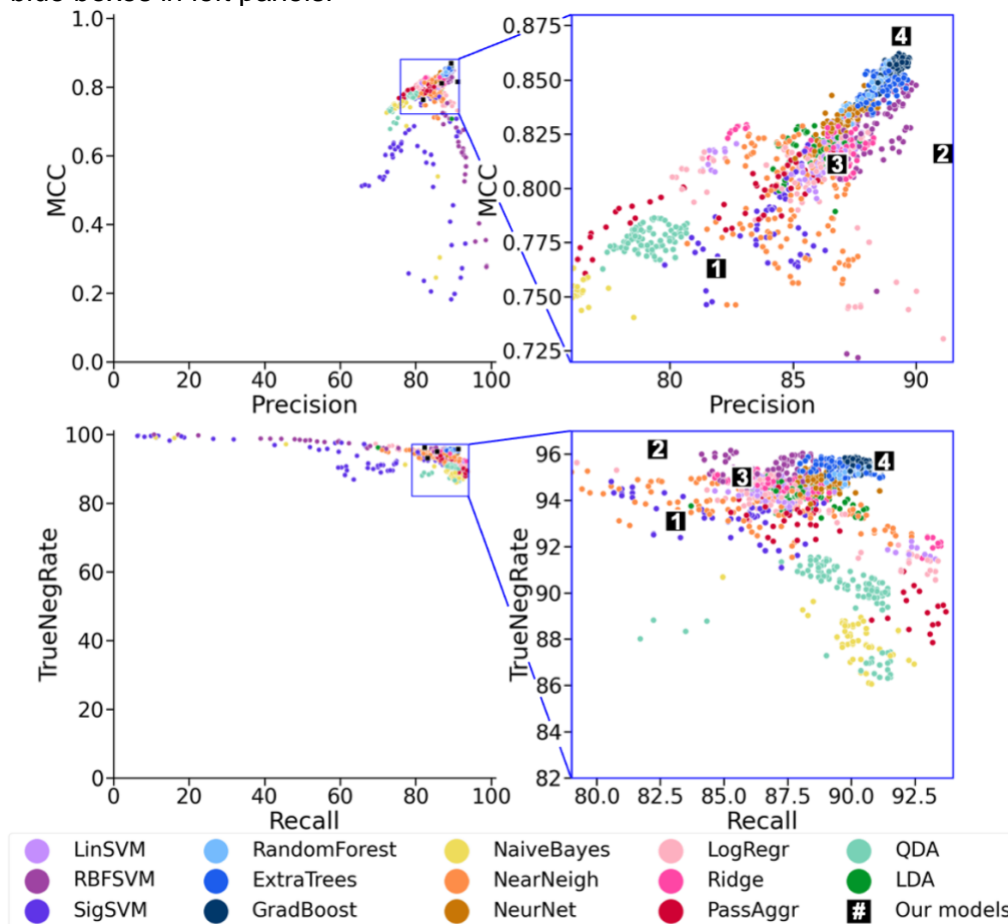
176 Using manual validation, we previously estimated that ~6% of our non-catalytic sites are  
177 mislabeled and that ~0% of our catalytic sites were mislabeled (Feehan, Franklin, and  
178 Slusky 2021). When using cross validation to evaluate newer (intermediate) iterations of  
179 MAHOMES, we found seemingly-incorrect predictions were often actually the sites our  
180 data generation pipeline mislabeled. We therefore intentionally used ML to hunt for  
181 mislabeled sites in our dataset via cross-validation.

182 Cross-validation is an ML method that leaves out a fraction of the dataset during training  
183 so that it can be used to assess the model's predictive performance. The left-out

184 fraction is iterated over the entire dataset, meaning a model makes predictions for each  
185 site in the training dataset. We manually examined non-enzymatic dataset sites that  
186 were predicted to be enzymatic during cross validation (see methods for more details).  
187 Because manual inspection during work on MAHOMES of 50 random dataset enzyme  
188 sites revealed an ~0% enzyme labeling error (Feehan, Franklin, and Slusky 2021), we  
189 did not examine enzymatic sites that were predicted to be non-enzymatic.

190 We used the available literature (structure publications, RCSB (Burley et al. 2019), and  
191 UniProt (UniProt 2019)) to investigate 225 sites that were previously labeled non-  
192 enzymatic but were classified during this cross validation as enzymatic. 94 of those  
193 sites had definitive literature support of catalytic activity (mislabeled) and 26 PDBs were  
194 removed from our set due to inconclusive evidence. Our previous estimate of 6%  
195 mislabeled non-catalytic sites implied approximately 158 mislabeled sites in the dataset.  
196 Therefore, we estimate that finding 94 mislabeled sites reduces our site mislabeling by  
197 60%.

**Figure 2: Cross-validation performance by algorithm.** Each dot represents one of the 1,792 models assessed in this work. The dots are colored to represent the type of ML algorithm the model uses: support vector machines = purples, decision-tree ensemble methods = blues, linear models = reds, discriminant analysis=greens, naive Bayes = yellow, nearest neighbor = orange, and neural network = brown. Better performing classifiers should have higher precision, Mathews correlation coefficient (MCC), true negative rate (TNR), and recall, meaning better classifiers will be close to the upper right corner. The black boxes with numbers show CV performance of: (1) the previously reported MAHOMES, (2) MAHOMES recalculated with the updated data labels, (3) MAHOMES retrained on updated labels, and (4) MAHOMES II – updated labels, trained on updated labels, and using new features. Right panels are zoomed in views of blue boxes in left panels.



## 198 2.3 ML model assessment

199 We generated 1,792 different ML models (Figure 2) using the following steps: feature  
 200 standardization, feature selection, and fourteen ML classification algorithms using one  
 201 of four optimization scoring terms. Since ML algorithms require or are greatly aided by

202 standardization of feature values in order to make comparable scales between the  
203 values of different features, we tested four different standardization techniques (see  
204 methods). Additionally, large numbers of input features can be detrimental to certain ML  
205 algorithms. To decrease the number of features with minimal information loss, we  
206 identified four feature subsets each using a different cut off to remove correlated  
207 features (see methods). In total we tried six feature sets (four low correlation subsets, all  
208 features, and a manually curated set) The six standardized feature sets were then used  
209 as inputs to ML classification algorithms which include: linear regression, decision-tree  
210 ensemble methods, support vector machines , nearest neighbors, Bayesian  
211 classification, and simple neural networks (see methods).

212 Selecting the best variation of the ML algorithm on the same data used to access a  
213 model can inflate performance metrics. To avoid inflated model assessment metrics, we  
214 used nested cross validation using an inner loop and an outer loop. During the inner  
215 loop, the ML algorithm was fine-tuned for a particular scalar and feature set using one of  
216 four different scoring metrics— accuracy, precision, Matthews correlation coefficient  
217 (MCC)(Matthews 1975), or a multi-score combination of accuracy, MCC, and Jaccard  
218 index(Jaccard 1907). Among our hyperparameter search space, each of the top three  
219 ranking ML algorithm variations were used to make models that were accessed using  
220 the outer loop. In total, we attempted 4,032 machine learning combinations (14  
221 algorithms x 6 feature sets x 4 standardization techniques x 4 optimization terms x 3 top  
222 algorithm variations). Due to convergence during model optimization, this process  
223 resulted in 1,792 different ML models.

224 The vast majority of all attempted ML models in this study outperformed the previous  
225 reported MAHOMES cross validation metrics (Figure 2, black box 1) because the  
226 training set was substantially corrected for all the new models. In order to make a more  
227 fair comparison between MAHOMES and the new models, we re-calculated MAHOMES  
228 cross validation performance metrics using the corrected enzyme/non-enzyme labels  
229 (Figure 2, black box 2). The number of MAHOMES cross validation false positives  
230 dropped from 182 to 90, which increased the precision by nearly 10% (Figure 2, top  
231 row) but the rest of the performance metrics remained far below those of our new  
232 models.

233 To assess if the increase in performance was purely due to corrected data labels, we  
234 assessed an intermediate model, which retrained MAHOMES using the corrected data  
235 but using the old MAHOMES features (Figure 2, black box 3). Despite an increase in  
236 recall, the retrained MAHOMES still identified significantly fewer enzyme sites than  
237 similar ML models that used the new features (Figure 2, CV blue). Thus, our ML  
238 benefitted from the improvement of both the quality of training labels and the improved  
239 features.

#### 240 2.4 MAHOMES II performance

241 To evaluate if these metrics are inflated from overtraining despite cross validation, we  
242 also predicted sites in an updated hold-out test set. In addition, we developed a new set  
243 derived from the hold-out test set to evaluate the reliability of the models. This set, the  
244 T-metal-sites<sub>10</sub>, includes ten different minimized structures for each site in T-metal-  
245 sites. The sub-angstrom variations for each site allowed us to calculate two

246 reproducibility metrics. First, we calculated the divergence frequency (equation 4,  
247 methods), which is the percent of test-set sites that received both an enzyme and non-  
248 enzyme prediction. Then, we calculated the divergence score (equation 5, methods), a  
249 measurement of the severity of divergent predictions. The divergence score ranges  
250 from 0 (the site receives the same prediction for every structure) to 1 (the site is  
251 predicted enzymatic for half of the structures and non-enzymatic for the other half).

252 **Table 1.** ML model performance evaluations. Predictive performance of MAHOMES,  
253 retrained MAHOMES with corrected labels, and MAHOMES II on the holdout test-set, T-  
254 metal-sites10, and the quality AlphaFold set, which is the subset of generated sites with  
255 confidence scores recommended to characterize binding ( $pLDDT > 90$ ) within 15 Å of the  
256 metal. TNR is the true negative rate and MCC is the Mathews correlation coefficient.  
257 Descriptions of performance metric calculations in methods. \*Evaluation using T-metal-  
258 sites, which includes ten incorrectly labeled sites and eight undeterminable sites which  
259 were removed from T-metal-sites10.

ML model	Evaluation	Accuracy	Precision	Recall	TNR	MCC	div. freq.	div. score
MAHOMES*	T-metal-sites	94.2	92.2	90.1	96.2	0.87	-	-
Recalculated MAHOMES	T-metal-sites10	92.6	94.2	85.3	96.9	0.84	6.6	0.53
	Quality AlphaFold set	92.7	94.1	66.7	99.0	0.75	-	-
Retrained MAHOMES	T-metal-sites10	93.4	91.9	90.0	95.4	0.86	5.9	0.55
MAHOMES II	T-metal-sites10	94.9	94.1	92.1	96.6	0.89	5.0	0.47
	Quality AlphaFold set	97.6	95.7	91.7	99.0	0.92	-	-

260

261 For MAHOMES II, we selected a GradBoost model that used FeatureSet4. We selected  
262 that as our final model because of its high cross validation metrics. However,  
263 ExtraTrees models, which is the algorithm used by the previous MAHOMES model, had

264 the lowest divergence frequency (Figure S2). So, we further refined hyperparameters,  
265 which were too computationally expensive to optimize for GradBoost during our inner  
266 cross validation using GridSearch optimization to mimic those favored by the  
267 ExtraTrees models (supplemental methods MAHOMES II, fine tuning). The final  
268 MAHOMES II model had a cross validation MCC higher than any other ML model  
269 (Figure 2 black box 4). Though this could have indicated an overfit model, our final  
270 performance evaluation on the hold-out test set T-metal-sites10 (Table 1), which we  
271 only used for reproducibility metric calculations during optimization, fell within the  
272 projected performance of the CV assessment (Sup. Table S1), thereby supporting the  
273 dependability of both the CV assessment and final evaluation of MAHOMES II  
274 performance.

275 In addition to improved performance, our aim was to lessen the effects of sub-angstrom  
276 deviations in input structures. Retraining MAHOMES with corrected labels decreased  
277 the divergency frequency but increased the divergence score (Table 1). Updated  
278 features in MAHOMES II further decreased the divergence frequency *and* decreased  
279 the divergence score, demonstrating that our feature improvements were effective at  
280 improving reproducibility.

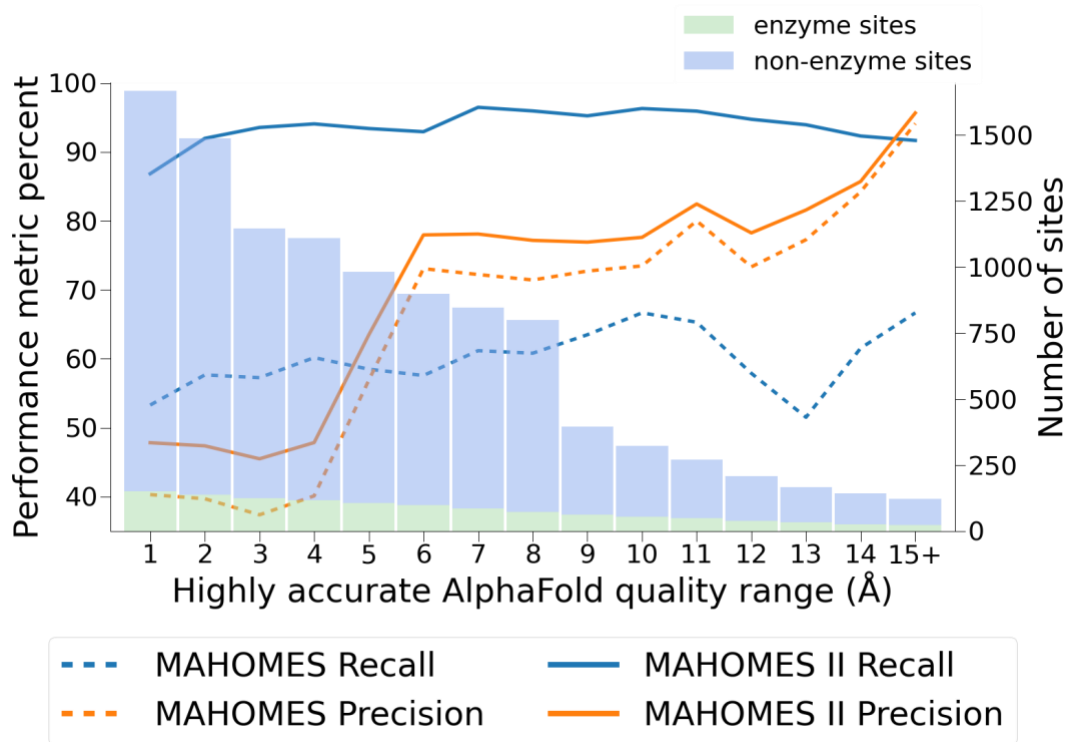
281 To test our hypothesis that upgraded features and improved training data can be used  
282 to successfully predict enzyme activity for predicted structures, we tested MAHOMES II  
283 on a set of AlphaFold generated structures (Tunyasuvunakool et al. 2021; Jumper et al.  
284 2021). However, AlphaFold generated structures do not have ligands such as metals.  
285 To create the AlphaFold set of protein structures with metals we queried UniProt  
286 (UniProt 2019) for proteins with known metal coordinating residues and no solved



287 crystal structure and filtered for metal ions that could be mapped to AlphaFoldDB  
288 structures. Benchmarking our metal method using dataset sites revealed sub-angstrom  
289 placement accuracy.

290 AlphaFold predictions have a confidence metric associated with each residue. The  
291 AlphaFold authors recommend using residues with high confidence (pLDDT >90) for  
292 characterizing binding sites. Very few sites in our AlphaFold set had high confidence for  
293 all residues used to calculate the MAHOMES II features (residues within 15 Å of the  
294 metal). So, we made MAHOMES II predictions on the entire non-homologous,  
295 metalloprotein AlphaFold set and made multiple performance evaluations by requiring

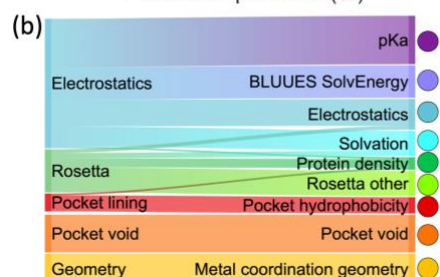
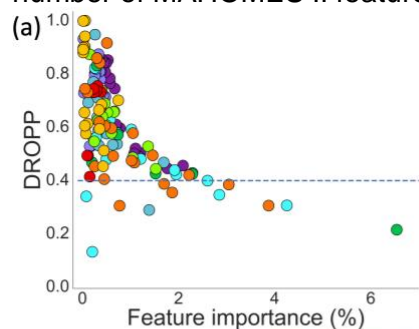
**Figure 3.** AlphaFold set performance evaluation. The number of enzyme (green bar) and non-enzyme (blue bar) AlphaFold set sites containing only highly confident residues within X Å, where X ranges from 1 to 15. The recall (blue lines) and precision (orange lines) of MAHOMES (dotted lines) and MAHOMES II (solid lines) are shown for the AlphaFold set sites at each cutoff.



296 residues within  $X \text{ \AA}$  of the site to be high quality (Figure 3, Table 1 and S2). For the

297 entire AlphaFold set, MAHOMES II was able to

**Figure 4.** Feature importance and DROPP. (a) Each dot represents MAHOMES II feature and is colored by physicochemical group. The y-axis is the feature's DROPP, or overlap between values for enzyme and non-enzyme dataset sites. The x-axis is feature importance for MAHOMES II, which is a measurement of the mean decrease of impurity by a feature during training. The blue dotted line represents the lowest feature DROPP from MAHOMES. (b) Sankey diagram of feature distribution between feature categories and feature groups, where width is representative of number of MAHOMES II features.



correctly identify 87% of the enzyme sites (recall) and 90% of the non-enzyme sites (true negative rate)(Sup Table S2). As we removed structures with low quality residue predictions close to the metal, MAHOMES II performance increases up to an accuracy of 97.5% an improvement even over our test-set metrics (Figure 3).

Interestingly, the coordinating residues' quality was not the most important, as MAHOMES II performance increases the most as the quality range increases from  $4 \text{ \AA}$  to  $6 \text{ \AA}$  (Figure 3). MAHOMES II enzyme recall (predicting which protein sites are catalytic) was very stable over all confidence regions and vastly outmatched the previous MAHOMES model (Figure 3).

## 2.5 Feature importance

314 Because MAHOMES II uses a decision-tree based  
315 gradient boosting algorithm, we can measure each feature's importance via the feature  
316 contribution to the decrease in impurity on the training data. As previously shown  
317 (Feehan, Franklin, and Slusky 2021), features with high importance had low DROPP

318 (overlap between enzyme site and non-enzyme site feature values). The five most  
 319 important features for MAHOMES II had lower DROPP than any MAHOMES feature  
 320 (Figure 4). However, the feature with the lowest DROPP, the minimum solvent exposure  
 321 parameter for outer sphere (3.5 – 9 Å) ionizable residues, was not important to  
 322 MAHOMES II—it ranked 116<sup>th</sup> in feature importance for MAHOMES II (Sup Table S3).  
 323 Hence, though quantitative differences, such as those measured by DROPP, can  
 324 indicate potentially important features, MAHOMES II is learning more than just these  
 325 numerical differences in order to successfully differentiate between enzyme and non-  
 326 enzyme sites.

327 **Table 2.** Feature group importance. Each feature group is described by its number of  
 328 included features (num total), the percent of MAHOMES II feature space accounted for  
 329 by the group, the total feature importance for all group features, the mean feature  
 330 importance of features in the group, the rank of the most important feature in the group,  
 331 and the mean DROPP of features in the group.

Feature group	Number of features	MAHOMES II feature space	Feature importance			DROPP mean
			total	mean	Max (rank)	
Local protein density	7	4.9%	14%	2.0%	6.5% (1)	0.48
solvation	14	9.7%	18%	1.3%	4.3% (2)	0.46
Pocket void	21	14.6%	22%	1.1%	3.9% (3)	0.54
pKa	27	18.8%	19%	0.7%	2.1% (10)	0.69
Rosetta	14	9.7%	9%	0.6%	1.5% (17)	0.62
Electrostatics	17	11.8%	9%	0.5%	1.7% (16)	0.63
Pocket hydrophobicity	9	6.2%	2%	0.2%	0.4% (75)	0.64
BLUUES SolvEnergy	18	12.5%	4%	0.2%	0.5% (59)	0.75
Metal coordination geometry	17	11.8%	3%	0.2%	0.7% (34)	0.74

332

333

334 Since the original feature categories were based on calculation method, we transitioned  
335 to feature groups (Fig. 1B) for analyzing which physicochemical properties were the  
336 most important for identifying catalytic activity. For example, the Coulombic electrostatic  
337 potential RBF feature had been in the Rosetta category but was a better fit for the  
338 electrostatic group. Due to differences in feature importance, we split features  
339 describing solvation into two groups. The BLUUES SolvEnergy group includes features  
340 calculated directly from the BLUUES solvation energy output. We placed other solvation  
341 related features in the Solvation group.

342 The three most important feature groups are local protein density, solvation, and pocket  
343 void (Table 2). Despite only making up 29% of MAHOMES II feature space, these  
344 groups account for 55% of what the model learned during training. These feature  
345 groups also have the lowest average DROPP. Using the DROPP plots for features in  
346 these groups, we identify specific subgroups that were fundamental to MAHOMES II  
347 distinguishing between enzyme and non-enzyme sites.

348 The local protein density feature group (seven features) has the highest average feature  
349 importance and includes Lennard-Jones energies and the number of residues within a  
350 certain distance of the site. This group includes the most important MAHOMES II  
351 feature, the number of residues within 15 Å (Figure 1B), which is more important than  
352 the 44 least important features combined (Table S3).

353 The next most important group, the solvation feature group (fourteen features), includes  
354 Rosetta solvation features and BLUUES generalized Born features. The second most  
355 important MAHOMES II feature is the average BLUUES solvent exposure parameter for

356 second shell (3.5-9Å) ionizable residues. The DROPP plot for this feature shows that  
357 most second shell ionizable residues are buried for enzyme sites and relatively exposed  
358 for non-enzyme sites (Figure S3b). This group also contains the sixth most important  
359 feature, the maximum second shell generalized Born radius, which measures an atom's  
360 shielding from high solvent dielectric (Figure S3f). Enzyme sites also have higher  
361 Rosetta solvation features that rank fifth, twelfth, and twenty-first in feature importance  
362 (Figure S3e, Table S3), which corresponds with the energetic cost associated with  
363 buried charged residues.

364 The third most important feature group is the pocket void group (twenty-one features).  
365 The pocket void group has features that describe the pockets' location, shape, and size.  
366 The third most important MAHOMES II feature describes the slice of the pocket closest  
367 to the metal as being larger for enzymes (Figure S3c). The fourth most important  
368 feature is the shortest distance between a metal and pocket grid point, which is smaller  
369 for enzyme sites (Figure S3d). These features combine to make a subgroup describing  
370 site accessible pockets.

### 371 **3. Discussion**

372 Our previous classifier, MAHOMES, outperformed available, alternative methods for  
373 classifying enzymes or non-enzymes. MAHOMES II, outperforms its predecessor with  
374 increased reliability thanks to both upgraded features and reduced training data error.  
375 MAHOMES II's performance generalize to new, unseen metalloproteins. Moreover,  
376 MAHOMES II learned physicochemical properties related to our current understanding  
377 of enzyme function.

### 378 3.1 MAHOMES II learned general enzyme activity

379 A key question of any classifier is if it has learned beyond its training, i.e. can it predict  
380 for examples it has never seen before. For MAHOMES II, training on solved, crystal  
381 metalloproteins structures could limit its performance to the 0.056% of proteins with  
382 experimentally determined structures (UniProt 2019). Our evaluation using the newly  
383 curated AlphaFold set finds that MAHOMES II generalizes to new enzyme reactions  
384 and even generalizes to very unrelated proteins.

385 Alternative tools that can be used to identify enzymatic activity (Zou et al. 2019; Kumar  
386 and Skolnick 2012) are less successful than MAHOMES II at predicting if our set of  
387 metalloproteins are catalytic (Table S4). Despite using ML, these enzymatic activity  
388 classifiers and catalytic residues predictors (Somarowthu and Ondrechen 2012; Song et  
389 al. 2018) rely on homology-based features causing their performance to not be  
390 transferable to catalysis more generally or be applicable for novel or designed enzymes.

391 To make our training data different enough from our testing data to facilitate  
392 generalizability, our training datasets and test-sets in both in this work and our previous  
393 work (Feehan, Franklin, and Slusky 2021) remove redundancy using local similarity. We  
394 only kept sites with dissimilar surrounding amino acid identities preventing training and  
395 evaluation of repeated sites among homologs and rare cases of similar active sites on  
396 different structural folds (Parasuram et al. 2016).

397 Using the AlphaFold data set, we determined that our model was extremely  
398 generalizable and was not implicitly using homology trends. The extensive quantity of  
399 AlphaFold structures and experimental data from UniProt for enzyme labeling (instead

400 of homology) allowed us to use a very high E-value of 1, i.e. only proteins with no  
401 evolutionary relationship, for creating our AlphaFold set. In comparison, only 17% of our  
402 previous test set, T-metal-sites10, sequences have no detectable homology to  
403 metalloproteins used for training MAHOMES II (E-value > 1). Furthermore, only seven  
404 of the 46 biochemical reactions included in the AlphaFold set are also included in the  
405 dataset used to train MAHOMES II. Despite the use of computationally generated  
406 structures and strict redundancy removal, MAHOMES II's 90-97.5% accuracy on the  
407 AlphaFold set was similar to its CV and T-metal-sites10 evaluations. Therefore, we  
408 believe our assessment of MAHOMES II performance will remain true for any natural  
409 metalloprotein structure uploaded by the community on the webserver, even if it is for a  
410 novel enzyme reaction. However, due to a lack of available structures, we remain  
411 uncertain if MAHOMES II performance transfers to *de novo* metalloproteins.

### 412 3.2 The less important first shell

413 Frequently, enzyme bioinformatics focuses on the active site's first shell, which is the  
414 residues interacting directly with substrate(s) or cofactor(s), such as metal ion(s) (Bartlett  
415 et al. 2002; Furnham et al. 2016; Ribeiro et al. 2018). The crucial roles played by first  
416 shell residues are well supported by conservation and experimental studies (Morley and  
417 Kazlauskas 2005; Ribeiro et al. 2020). MAHOMES II has 60 features that describe only  
418 first shell properties, covering coordination geometry, inner shell electrostatics (< 3.5 Å  
419 from metal), and pocket lining. Despite making up 42% of MAHOMES II's feature space,  
420 first shell features account for only 18% of feature importance. Since the same metal  
421 and coordinating residues are found to participate in enzyme and non-enzyme functions  
422 (Lee et al. 2019), it makes sense that first shell features are largely incapable of

423 differentiating enzyme and non-enzyme sites in metalloproteins since in both the first  
424 shell coordinates metals. Consequently, despite the well-known critical roles of the first  
425 shell, distinction between metallo-enzymes and metallo-proteins is driven by more  
426 distant physicochemical properties.

### 427 3.3 Comparing important MAHOMES II subgroups to current enzyme paradigms

428 The physicochemical features most important to MAHOMES II success can be  
429 considered as three groups/subgroups—1) local protein density, 2) second shell  
430 ionizable residue burial, and 3) site accessibility of pockets (in the pocket void feature  
431 group) – align with the current paradigm of the enzyme function, which also consists of  
432 three features: 1) local environment control of functional sites through control of water  
433 access, 2) networks of residue interactions spanning from functional residues, and 3)  
434 conformational dynamics (Mazmanian, Sargsyan, and Lim 2020; Agarwal 2019).

435 *Control of water access:* MAHOMES II captures local environmental control through  
436 water access with two of the important MAHOMES II feature subgroups: the local  
437 protein density group and site accessibility of pockets feature subgroup. The local  
438 protein density features, detect the dense packing of enzyme sites which protects them  
439 from high external dielectrics of bulk water, enhancing the local electrostatic effects from  
440 hydrogen-bonding and charge-charge interactions. The site accessibility of pockets  
441 subgroup identifies close pockets with large openings adjacent to the site that can  
442 enable access by individual water molecules, which commonly participate as  
443 nucleophiles, to form hydrogen-bonding networks, and to facilitate the release of



444 products. Hence, MAHOMES II can detect the control of water access to enzyme sites  
445 by combing local protein density and site accessibility of pockets.

446 *Networks of connected interactions:* Distal residues that interact with catalytic residues  
447 or as part of networks connecting to the catalytic site are essential for fine tuning and  
448 optimization of enzyme activity(Dudev et al. 2003; Somarowthu, Brodtkin, et al. 2011;  
449 Parasuram et al. 2018; Brodtkin et al. 2015; Tiwari et al. 2014; Coulther, Ko, and  
450 Ondrechen 2021; Coulther et al. 2021; Ngu et al. 2020). The MAHOMES II burial of  
451 ionizable residues feature subgroup differentiates enzyme sites based on buried second  
452 shell polar and charged residues, which would be a direct result of crucial coupled  
453 interactions that enhance enzyme activity. In addition, the MAHOMES II local protein  
454 density group uses the density of residues surrounding the active site to provide the  
455 most basic description of networks of interactions connected to enzyme sites with the  
456 potential to promote activity. The combination of these two subgroups therefore seems  
457 to accurately estimate connected networks.

458 *Conformational flexibility:* Although we did not design any MAHOMES II features to  
459 directly describe conformational dynamics, the final aspect of the current enzyme  
460 function paradigm, all of the three most important physicochemical subgroups describe  
461 properties that affect conformational stability. Local protein density describes tight  
462 packing that increases backbone hydrogen-bonding which increases stability and  
463 rigidity. Burial of charged residues amongst nonpolar sidechains makes for an  
464 energetically unfavorable conformation that will promote destabilization and flexibility.  
465 Moreover, interactions between charged sidechains will also increase or decrease  
466 stability of various active site conformations depending on the charges. Finally, the site

467 accessibility of pockets enables active site interactions with cofactors, substrates, and  
468 solvent that will change the flexibility or rigidity of an active site. Therefore, all three  
469 important subgroups contribute to the conformational changes required for enzyme  
470 activity, such as the shifting from the ground state to transition state(s).

### 471 3.4 Machine learning lessons for metalloenzyme design

472 Considering our training dataset covers all enzyme reaction types(Feehan, Franklin,  
473 and Slusky 2021), the physicochemical properties highlighted by MAHOMES II gives us  
474 insight for making better metalloenzyme designs. Their feature importance indicates a  
475 fundamental blueprint that is harnessed by a range of known catalysis mechanisms  
476 performed by nature. To this point, recent work exploring the functional space of non-  
477 metallo TIM barrel enzymes has also highlighted the importance of local atomic  
478 density(Lipsh-Sokolik et al. 2023). *De novo* enzyme designs on non-enzyme protein  
479 backbones could benefit from selecting densely surrounded positions with large pocket  
480 openings. Furthermore, lower solvation penalties for buried for ionizable residues might  
481 also help design active sites that more closely resemble those in nature. We anticipate  
482 that dense protein regions with buried ionizable residues can improve the success rate  
483 of designed enzymes and limit additional steps that are currently necessary to reach  
484 native enzyme reaction rates, such as directed evolution (Yang, Wu, and Arnold 2019).

## 485 **4. Conclusion**

486 Our ML classifier, MAHOMES II (<https://mahomes.ku.edu>), uses protein structure-based  
487 features describing the local site to distinguish between enzyme and non-enzyme metal  
488 ion sites on proteins with 94% precision and 92% recall. We demonstrated that

489 MAHOMES II can make quality predictions for computationally generated structures,  
490 which greatly expands its utility when combined with the structure prediction tool  
491 AlphaFold. Additionally, the similarity among performance metrics for our cross-  
492 validation, holdout test-set, and evolutionarily unrelated AlphaFold set supports that  
493 MAHOMES II evaluation is not bias, overfit, or the result of off-target learning. Finally,  
494 we were able to identify that MAHOMES II was making successful predictions due to its  
495 use of physicochemical features related to densely packed active sites, burial of second  
496 shell ionizable residues, and site accessible pockets.

## 497 **5. Methods**

### 498 5.1 Metal ion dataset and T-metal-sites

499 The data developed to train and evaluate MAHOMES (Feehan, Franklin, and Slusky  
500 2021) is, to our knowledge, the largest non-redundant dataset of enzymatic and non-  
501 enzymatic labeled protein bound metal ions. Briefly, protein structures containing  
502 transition metals were filtered to remove poor quality structures and structures dissimilar  
503 to metalloenzymes. Metal ion sites bound to multiple chains were removed to avoid  
504 labeling partial enzyme active sites during our homology-based enzyme labeling.  
505 Metalloenzymes were identified using explicit enzymatic annotations and homology to  
506 entries in the Mechanism and Catalytic Site Atlas (M-CSA)(Ribeiro et al. 2018), a  
507 database of enzyme active sites. Alignment with M-CSA homolog structures was used  
508 to label metal sites on metalloenzymes as enzyme or non-enzyme. Metals on  
509 metalloproteins that lacked explicit enzymatic annotations and had no M-CSA homolog  
510 were labeled as non-enzyme sites. Finally, sequence and structural homology were  
511 used to remove redundancy. Sites on structures deposited in the Protein Data Bank  
512 (PDB)(Berman et al. 2000) prior to 2018 were placed in the dataset, which was used for  
513 training ML models, ML optimization, and model selection. Sites on structures deposited  
514 in the PDB in 2018 or later were placed in a holdout test-set, called T-metal-sites, which  
515 was used for final model evaluation.

516 All the metalloprotein structures in the dataset and T-metal-sites were relaxed using  
517 Rosetta (Conway et al. 2014) using a previously provided RosettaScript (Feehan,  
518 Franklin, and Slusky 2021). To remove loosely bound metals that are less likely to be

519 physiologically relevant, i.e. crystal artefacts, we removed 729 sites that moved more  
520 than 3 Å from the aligned crystal structure. We also removed 179 sites that failed  
521 MAHOMES feature calculations since this was commonly due to issues like lack of  
522 multiple coordinating atoms. New sites were defined as any metals within 5 Å of each  
523 other in a relaxed structure. Since the original dataset defined sites using the crystal  
524 structures, the revised set slightly differs in the number of sites. All sites containing a  
525 metal atom that was previously a part of an enzyme site were labeled enzyme. Any  
526 remaining site with a metal atom that was previously included in a non-enzymatic site  
527 was labeled non-enzymatic. All other sites on these structures were discarded. We  
528 found that MAHOMES was susceptible to making different predictions for the same site  
529 on different relaxed structures. To check model prediction reproducibility, we included  
530 the ten relaxed structures with the lowest Rosetta energy units for each metalloprotein  
531 in T-metal-sites, making T-metal-sites10. We repeated the labeling procedure for the T-  
532 metal-sites10 sites.

533 We removed sites that were flagged by our automated feature process as problematic  
534 or that had extreme outlier feature values (greater than ten standard deviations from the  
535 dataset mean). Manual examination of sites with incorrect ML predictions identified a  
536 significant number of incorrect non-enzyme labels by our pipeline for sites in both the  
537 dataset and T-metal-sites10 (Table S5). Sites found to actually be enzymatic were  
538 relabeled and sites with undeterminable enzyme activity were removed. At the end of  
539 this work, the MAHOMES II dataset contained 957 enzyme sites and 2,467 non-enzyme  
540 sites. The final T-metal-sites10 consisted of 1,895 enzyme entries and 3,277 non-

541 enzyme entries, which were representative of 189 enzyme sites and 328 non-enzyme  
542 sites.

## 543 5.2 Feature engineering

544 For machine learning input features, we calculated physicochemical properties for five  
545 categories – Rosetta energy terms, pocket void, pocket lining, electrostatics, and  
546 coordination geometry. A complete feature list with descriptions can be found in Table  
547 S6. For exact calculations, please see our available github code (Feehan et al. 2022).

### 548 5.2.1 Rosetta energy terms

549 Rosetta 3.13 was used to score all residues in a metalloprotein structure for all energy  
550 terms in the energy function *beta\_nov16* with all weights set to one (Alford et al. 2017).  
551 For each energy term,  $E$ , a squared inverse radial basis function (Eq. 1) was used to  
552 calculate the energy for a given site,  $\mathbf{s}$ .

$$553 \quad E(\mathbf{s}) = \sum_r \frac{E(r)}{d(r)^2} \quad (\text{Eq. 1})$$

554 where  $r$  is a residue with a distance  $d(\mathbf{r}) < 15 \text{ \AA}$  from the site center. We included the  
555 number of residues used for these calculations as a feature. Our Rosetta energy terms  
556 category included 25 features in total.

### 557 5.2.2 Pocket void terms

558 We used GHECOM (Kawabata 2019) – a program for detecting multiscale pockets via  
559 grid representations and probes – to identify all pockets for a given metalloprotein.  
560 Then, for each site on the metalloprotein, we identify all adjacent pockets – pockets with

561 a grid point within 5Å of the site's center. For sites with multiple adjacent pockets, we  
562 select the closest adjacent pocket with a volume greater than 100 Å<sup>3</sup> as its pocket. In all  
563 other cases, the GHECOM pocket with the closest grid point is selected.

564 We used the selected pocket to calculate the pocket void features previously used by  
565 MAHOMES (Feehan, Franklin, and Slusky 2021), which include volume, Euclidean and  
566 Manhattan distance between pocket centroid and site center, terms describing the size  
567 and shape of three pocket slices at the site center, pocket center, and midpoint of site  
568 center and pocket center.

569 We added pocket void features output by GHECOM for the relative rank among all the  
570 metalloprotein's detected pockets and terms that describe the pockets shallowness and  
571 openness of the pocket. Then, we rotated the pocket so that the z-axis runs from the  
572 protein centroid to the pockets centroid and calculate its height, max z – min z, and  
573 depth, the Euclidean distance between the grid points with the max z and min z. Finally,  
574 we calculate the site's height and depth in the pocket using the Euclidean distance  
575 between the site center and the max z grid point or min z grid point respectively.

### 576 5.2.3 Pocket lining

577 The selected GHECOM pocket was used to identify pocket adjacent residues (within 3.0  
578 Å). We identified pocket lining residues as pocket adjacent residues with a sidechain  
579 distance of less than 2.2 Å or where the sidechain was closer to the pocket than the  
580 backbone. The pocket lining residues were used to calculate the average, minimum,  
581 maximum, skew, and standard deviation of the hydrophobicity for both Eisenberg  
582 (Eisenberg et al. 1984) and Kyte-Doolittle (Kyte and Doolittle 1982). We also calculated

583 the sum of the pocket lining residues van der Waals sidechain volumes, the volume of  
584 the pocket without the lining residues' sidechains, and the percent of that volume  
585 occupied by the sidechains. Finally, the number of pocket lining residues and the  
586 number of backbone adjacent residues (pocket adjacent but not considered pocket  
587 lining) were included as features.

#### 588 5.2.4 Electrostatic terms

589 Our previous electrostatics features were based on the use of theoretical titration curves  
590 by THEMATICS (Ondrechen, Clifton, and Ringe 2001; Somarowthu, Yang, et al. 2011;  
591 Ko et al. 2005). To calculate these, we used the generalized Borne program, BLUUES  
592 (Fogolari et al. 2012a). For the first derivative of the theoretical titration curve, we  
593 calculated the second, third, and fourth moment of each ionizable residue using SciPy's  
594 (Virtanen et al. 2020) variation, skew, and kurtosis functions respectively. The mean,  
595 standard deviation, maximum, minimum, and range was calculated for two shells. The  
596 first, inner shell included ionizable residues within 3.5 Å of a site's metal atom(s). The  
597 second, outer shell included ionizable residues within 9 Å of a site's metal atom(s),  
598 excluding residues that are in the first shell. For each moment calculation, the Z-score  
599 was calculated (Eq. 2), where  $x$  is the residue's moment value,  $\mu$  is the moment's average  
600 for all of the structure's ionizable residues, and  $\sigma$  is the moment's standard deviation. We  
601 turned this into a site feature by counting the number of residues with a Z-score greater  
602 than 1.

$$603 \quad z = \frac{x - \mu}{\sigma} \quad (\text{Eq. 2})$$



604 All residues from both shells were used to calculate an environmental feature for each  
605 moment using a squared inverse radial basis function (Eq. 1). The number of residues  
606 used for the inner shell, outer shell, and environmental feature calculations were also  
607 saved to be used as features.

608 The inner shell, outer shell and environmental features were also calculated for additional  
609 BLUUES outputs, which included: generalized Born radius, residuals for the deviation of  
610 the theoretical titration curve from a sigmoidal curve, pKa shift from ideal amino acid  
611 values, pKa shift due to desolvation, pKa shift due to the interaction with other charges in  
612 the molecule with all titratable sites in their neutral state, pKa shift due to the interaction  
613 between titratable sites, solvation energies, and solvent exposure parameter. All of the  
614 structure's residues were sorted by BLUUES solvation energy and placed into five bins;  
615 destabilizing ranks were assigned from highest to lowest solvation energy and stabilizing  
616 ranks were assigned from lowest to highest solvation energy. The inner shell, outer shell  
617 and environmental features were calculated for both destabilizing and stabilizing ranks.  
618 Overall, there are 152 electrostatic features.

### 619 5.2.5 Coordination geometry terms

620 We used FindGeo (Andreini, Cavallaro, and Lorenzini 2012) to determine the  
621 coordination geometries of the site's metal atom(s). First, we record the total number of  
622 ligand N, O, S, and other atoms used as input for FindGeo (within 3.5 Å of any site  
623 metal). Then, we record if the metal atom(s) were identified as a regular, distorted, or  
624 irregular geometry. If the geometry is regular or distorted, we use the RMSD from the  
625 idealized geometry, the number of coordinating atoms for the geometry, and if it is

626 completely or partially filled. To prevent issues with ML algorithms and categorical  
627 features, the number of coordinating geometries are one hot key encoded, giving us a  
628 total of 20 coordination geometry features.

#### 629 5.2.6 Feature analysis

630 DROPP was calculated as previously described for feature similarity (Feehan, Franklin,  
631 and Slusky 2021). For discrete features, we used the Jaccard index between the  
632 proportions observed in the enzymatic and non-enzymatic sites (Jaccard 1907). For  
633 continuous features, we calculated the overlap of the kernel density estimators for the  
634 values of the enzymatic and non-enzymatic sites. To prevent extreme outliers from  
635 having large influence on DROPP, values greater than ten standard deviations the  
636 mean were ignored. Only dataset sites were used for calculating DROPP. The code for  
637 the DROPP calculation can be found in the MAHOMES II repository file  
638 FeatureCalculations/CalcFeatureDROPP.py (Feehan et al. 2022).

639 Any feature that had the same value calculated for the entire dataset was discarded,  
640 leaving a total of 250 features. To decrease feature space with minimal information loss,  
641 additional subsets of features were identified for ML input using maximum correlation  
642 cutoffs between features in the same category (Sup. Figure S1). For highly correlated  
643 features, the feature with the higher DROPP was removed. FeatureSet2, FeatureSet3,  
644 and FeatureSet5 used correlation cutoffs of 0.99, 0.9, and 0.75 respectively. Due to the  
645 dramatic increase of electrostatic features, FeatureSet4 used a correlation cutoff of 0.75  
646 for electrostatic features and 0.9 for other categories. Finally, we manually selected  
647 features for FeatureSet6.

## 648 5.2.7 Preparation for ML

649 ML algorithms require or are greatly aided by standardization of feature values to take  
650 them close to zero and make comparable scales between the values of different  
651 features. We selected four different standardization techniques available in scikit-learn  
652 (Pedregosa et al. 2011) to use during model optimization and selection. The  
653 StandardScaler removes the mean and divides by the features standard deviation. The  
654 RobustScaler removes the median and divides by range between the 20th and 80th  
655 quantile to mitigate the effect of extreme outliers. We also examined uniform and  
656 gaussian QuantileTransformers which use non-linear transformations to map feature  
657 values to uniform or gaussian distributions respectively. All scalars include an imputer to  
658 fill missing feature values with the dataset sites' average feature values. MAHOMES II  
659 used the uniform QuantileTransformer.

## 660 5.3 Machine Learning

### 661 5.3.1 Classification performance metric calculations

662 To calculate predictive performance metrics, a prediction is counted as a true positive  
663 (TP) if it is an enzyme prediction for an enzyme labeled site. It is considered a false  
664 positive (FP) if it is an enzyme prediction for a non-enzyme labeled site. A true negative  
665 (TN) is a non-enzyme prediction for a non-enzyme labeled site. Finally, a false negative  
666 (FN) is a non-enzyme prediction for an enzyme labeled site. The TP, FP, TN, and FN  
667 counts are then used to calculate a model's accuracy, precision, recall, true negative  
668 rate (TNR), and Matthews correlation coefficient (MCC)(Matthews 1975).

669 
$$accuracy = \frac{TP+TN}{TP+FP+TN+FN} \quad (\text{Eq. 3})$$

670 
$$precision = \frac{TP}{TP+FP} \quad (\text{Eq. 4})$$

671 
$$recall = \frac{TP}{TP+FN} \quad (\text{Eq. 5})$$

672 
$$TNR = \frac{TN}{TN+FP} \quad (\text{Eq. 6})$$

673 
$$MCC = \frac{(TP \times TN) - (FP \times FN)}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}} \quad (\text{Eq. 7})$$

### 674 5.3.2 Optimization of ML model(s)

675 Since the best performing hyperparameters for ML classification algorithm can change  
676 depending on input feature subset and standardization technique, a nested cross  
677 validation (CV) strategy was used to find the optimal ML models. The outer CV used  
678 stratified k-fold and the inner loop used stratified shuffle split. During the inner loop,  
679 different hyperparameters sets were attempted and scored using one of four terms -  
680 accuracy, precision, MCC, or a multi-score combination of accuracy, MCC and Jaccard  
681 index. The hyperparameter sets were ranked according to the average of the scoring  
682 metric. Our multi-score optimization ranked each scoring metric and then averaged the  
683 rankings. Due to the potential for ties, the top three were selected as the optimized ML  
684 models. Depending on the convergence of the different scoring terms, ML algorithm,  
685 feature subset, and standardization could have between three and twelve optimized ML  
686 models. To reliably compare optimized ML models, we used the average performance  
687 metrics during stratified k-fold CV with ten repetitions during each iteration using  
688 different random state hyperparameters for the classifier when applicable.

### 689 5.3.3 Evaluating a model's reproducibility

690 For considering a model's reproducibility, different minimized structure inputs in T-  
691 metal-sites10 were used to make a set of predictions,  $\mathbf{p}$ , for the same site,  $s$ . The site's  
692 divergence,  $d(s)$ , was calculated using Equation 8, where  $p_i$  is either 1 (enzyme  
693 prediction) or -1 (non-enzyme prediction) and  $n$  is the number of minimized input  
694 versions for  $s$ .

$$695 \quad d(s) = 1 - \left| \frac{\sum_{i=1}^n p_i}{n} \right| \quad (\text{Eq. 8})$$

696 Therefore,  $d(s)$  ranges from 0 to 1, where 0 is a site with the same prediction for all  
697 minimized inputs and 1 is a site with five enzyme predictions and five non-enzyme  
698 predictions. Using the set of all sites in T-metal-sites10,  $\mathbf{T}$ , and the subset of divergent  
699 sites,  $\mathbf{D} = \{s \mid s \in \mathbf{T} \text{ and } d(s) > 0\}$ , we calculated our reliability metrics. The divergence  
700 frequency is the percent of sites in T-metal-sites10 that were divergent (Eq. 9).

$$701 \quad \text{divergence frequency} = \frac{n(\mathbf{D})}{n(\mathbf{T})} = \frac{n(\{s \mid s \in \mathbf{T} \text{ and } d(s) > 0\})}{n(\mathbf{T})} \quad (\text{Eq. 9})$$

702 The divergence score is the average site divergence of the divergent sites (Eq. 10).

$$703 \quad \text{divergence score} = \frac{\sum\{d(s) \mid s \in \mathbf{D}\}}{n(\mathbf{D})} = \frac{\sum d(s)\{s \mid s \in \mathbf{T} \text{ and } d(s) > 0\}}{n(\{s \mid s \in \mathbf{T} \text{ and } d(s) > 0\})} \quad (\text{Eq. 10})$$

704 Since only divergent sites are considered, the lowest divergence score is 0.2.

### 705 5.3.4 ML guided dataset manual annotation

706 We optimized decision tree-based algorithms (ExtraTrees, GradBoost, and  
707 RandomForest) using FeatureSet4 and selected an ExtraTrees model with high MCC  
708 and high recall. We manually checked non-enzyme sites that were predicted to be  
709 enzyme sites during the k-fold CV using available structure publications, RCSB(Burley  
710 et al. 2019), and UniProt(UniProt 2019). We changed the labels for 73 sites that were  
711 actually enzyme sites and removed all sites from 16 metalloproteins that were  
712 undeterminable, for reasons such as lack of publication and enzymatic homologs. We  
713 repeated the process without ExtraTrees models to minimize redundancy of  
714 mispredicted sites. We fixed 19 additional site labels found during the second iteration,  
715 removed all sites from eight undeterminable metalloproteins. We also removed two sites  
716 that were located on the edge of active sites, meaning they could be correctly labeled  
717 as both enzymatic and non-enzymatic.

### 718 5.3.5 Recalculated MAHOMES and Retrained MAHOMES

719 Since we were able to identify a significant amount of labeling error in the data used for  
720 previously reported MAHOMES performance evaluations, we made updated  
721 performance evaluations to enable more fair comparisons between MAHOMES and  
722 MAHOMES II. We recalculated the k-fold performance for the MAHOMES predictions  
723 using the corrected dataset labels for what was and wasn't an enzyme. Since we made  
724 new relaxed structures for T-metal-site10 (the MAHOMES II test set), new MAHOMES  
725 predictions were made for T-metal-sites10. Both variations of the test-set received  
726 nearly the same predictions, but the performance evaluation changes significantly due  
727 to reduced labeling error in T-metal-sites10.

728 The recalculated MAHOMES performance evaluations were still from a model that was  
729 trained using dataset sites which have since been identified as mislabeled. So, we  
730 made a retrained MAHOMES model, which differs from recalculated MAHOMES in two  
731 ways. The first difference is that the fixed dataset with updated labels and removed  
732 undeterminable sites were used during training. The number of enzyme sites increases  
733 by 10% when the dataset labels are fixed, which prevents under-sampling at a ratio of 3  
734 non-enzyme:1 enzyme site during training. So, the retrained MAHOMES model under-  
735 samples by randomly removing 10% of the enzyme sites, followed by random removal  
736 of non-enzyme sites until the ratio of training data is 3 non-enzyme:1 enzyme sites.  
737 Otherwise, retrained MAHOMES model uses the same methods as the recalculated  
738 MAHOMES model, including calculated feature values, algorithm, and optimized  
739 hyperparameter set.

#### 740 5.3.6 Model selection and performance evaluation

741 Despite the favorable performance of decision tree-based classifiers during work on  
742 MAHOMES(Feehan, Franklin, and Slusky 2021), we tested fourteen ML classification  
743 algorithms from Scikit-learn(Pedregosa et al. 2011) for MAHOMES II, which include:  
744 linear regression, decision-tree ensemble methods, support vector machines , nearest  
745 neighbors, Bayesian classification, and simple neural networks. We decided to attempt  
746 these various algorithms because decision tree ensemble-based classifiers are known  
747 to be robust against mislabeled data, large feature spaces, and outlier feature values.  
748 So, our upgraded features and reduced training label error does not affect decision tree  
749 ensemble-based algorithms as much as it affects alternative ML classification  
750 algorithms.

751 In total, we assessed 4,032 machine learning combinations (14 algorithms x 6 feature  
752 sets x 4 standardization techniques x 4 optimization terms x 3 top hyperparameter sets).  
753 However, we only ended up with 1,792 unique ML models due to convergence during  
754 model optimization. The specific code used for the algorithms, standardization  
755 techniques, and hyperparameters can be found in the MAHOMES II repository file  
756 MachineLearning/GeneralML.py (Feehan et al. 2022).

757 We selected a model that used a gradient boosting classifier with FeatureSet4 and a  
758 uniform QuantileTransformer because it had the highest recall for models with greater  
759 than 0.845 MCC and 88.5% precision. For MAHOMES II, we further refined this model  
760 to improve both its cross validation MCC, divergence frequency and divergence score  
761 by adjusting hyper-parameters that were too computationally expensive to optimize  
762 during our inner cross validation using GridSearch optimization (Sup. methods).

### 763 5.3.7 Feature importance

764 For algorithms using decision-tree classifiers, sci-kit learn has a built in feature  
765 importance output that measures the mean decrease in impurity that a feature was  
766 responsible for during training. The MAHOMES II feature importance is the average of  
767 feature importance output of the models trained during k-fold cross validation.

## 768 5.4 AlphaFold set

### 769 5.4.1 AlphaFold set generation

770 To make the AlphaFold set, we queried UniProt(UniProt 2019) for reviewed entries with  
771 no solved crystal structure, an AlphaFold model (as of February 15, 2022), and metal



772 binding data. Entries with no EC number and no catalytic activity annotation of any kind  
773 were labeled as non-enzyme. Entries with annotated with experimental catalytic activity  
774 were labeled as enzyme. Remaining unlabeled entries were removed.

775 Since each entry is a protein sequence at this point, we chose to remove homology with  
776 training and evaluating data next. We used PHMMER (Eddy 2011) to search each  
777 protein sequence against all sequences in the PDB as of May 21, 2020. Entries with  
778 detected homology, using an E-value  $< 1$ , to any protein sequence in the dataset or  
779 test-set were removed.

780 To go from sequence to the site level data, we retrieved all available metal binding data  
781 from UniProt for each of the remaining entries. We removed data for metals other than  
782 Copper, Iron, Magnesium, Manganese, Zinc, and Nickel. To ensure that the labels were  
783 accurate at the site level, we removed enzyme labeled entries that did not include  
784 'catalytic' annotations for the metal binding site. Due to automatic metal site annotations  
785 or lack of EC coverage, non-enzyme labeled entries with 'catalytic' metal binding  
786 annotations also had to be removed. Entries with only one or two listed metal binding  
787 residue(s) were removed. We did not relax or perform any additional structure  
788 minimization. The resulting AlphaFold set contains 1740 computationally generated  
789 structures with 1583 non-enzyme sites and 157 enzyme sites.

790 We placed metals in sites using the average coordinates of the atoms binding to the  
791 metals. For hydrophilic amino acids, we used the coordinate of the sidechain atom  
792 capable of binding a metal ion (N, O, or S). For amino acids with multiple sidechain  
793 atoms capable of binding metal residues (GLU, ASP, GLN, ASN), we used the average

794 of these atomic coordinates. For GLY, we used the average coordinate of the backbone  
795 N and O. Some entries also listed other non-polar amino acids as coordinating residues.  
796 We found the average coordinate of sidechain carbons worked best for placing these  
797 metals without any steric clash. Since some entries with multiple metal binding sites did  
798 not differentiate different bound metal sites, we removed any entries if the coordinating  
799 residues were more than 12 Å from each other. The resulting 1,675 metalloprotein  
800 structures and enzyme/non-enzyme labels are available on Zenodo (Feehan 2023).

#### 801 5.4.2 AlphaFold set metal ion placement accuracy

802 We evaluated the accuracy of the AlphaFold set metal placement by adding a metal to  
803 the sites in our dataset and test-set using UniProt data and comparing it to the metal  
804 location in the relaxed crystal structures.

805 To find appropriate crystallographically-resolved sites to compare with the AlphaFold set  
806 metal placement, we retrieved available binding site data for 2,608 of the 2,643 UniProt  
807 entries in our dataset and test-set. However, only 1,207 entries included data for relevant  
808 metal binding sites -- CHEBI ids: 29105 (Zn<sup>2+</sup>), 29033 (Fe<sup>2+</sup>), 29034 (Fe<sup>3+</sup>), 29035  
809 (Mn<sup>2+</sup>), 29036 (Cu<sup>2+</sup>), 49552 (Cu<sup>+</sup>), 18420 (Mg<sup>2+</sup>), or 49786 (Ni<sup>2+</sup>). To create the  
810 dataset for benchmarking, we removed UniProt entries and PDBs with multiple sites.  
811 Also, entries with different metals in the PDB and UniProt binding data were removed. To  
812 accurately depict the placement of sites in our AlphaFold set, we removed sites with fewer  
813 than three coordinating residues. Finally, coordinating residues had to be among the  
814 previously described amino acid types, resolved in the PDB, and indexed with the same

815 numbers in UniProt and the PDB. These filtering steps resulted in a total of 103 sites  
816 remaining for benchmarking.

817 The final benchmark set consists of 103 successfully placed sites. The average distance  
818 between the placed metal and the metal in the relaxed crystal structure was 0.87 Å (Fig.  
819 S4). For comparison, the same 103 sites moved an average of 0.54 Å during minimization  
820 of the crystal structure with Rosetta relax. Moreover, 56% of sites were placed within 1 Å,  
821 96% were placed within 2 Å, and only one was more than 3 Å from its respective  
822 experimentally resolved location in the PDB structure.

823 The UniProt metal binding annotations were converted to binding site annotations  
824 (Coudert et al. 2023). This data conversion occurred after we created the AlphaFold set  
825 but before we benchmarked our metal binding site placement method. This conversion  
826 therefore required different data retrieval scripts for the benchmarking set than for the  
827 AlphaFold set.

## 828 5.5 webserver

829 The MAHOMES Web Server was implemented in Python 3 on the back end using the  
830 Flask framework with Jinja for templates in creating the HTML client-side interfaces.  
831 When a PDB file and email are submitted to MAHOMES, metadata about the job is  
832 stored in a JSON file. The information necessary to schedule the job for processing is  
833 placed into an SQLite3 database.

834

835 The job execution program, which is written in Python 3, monitors the SQLite database  
836 for new user submissions, and then handles executing the job, monitoring the job  
837 execution, and then sending an email to the user with a link to the results page.

838

839 The jobs and the web application are run on the Slusky Lab web server, which is a  
840 virtual machine running in the University of Kansas's enterprise data center. Running  
841 this service as a virtual machine has allowed us to scale up the hardware backing the  
842 instance as we have needed additional resources while working to control the long-term  
843 costs associated with running the MAHOMES service.

844

## 845 **6. Supplementary material description**

846 *Figure S1.* Comparison of feature sets.

847 *Figure S2.* ML models reproducibility.

848 *Figure S3.* Top 6 feature DROPP plots

849 *Figure S4.* Metal binding site placement accuracy

850 *Table S1.* Performance evaluations.

851 *Table S2.* AlphaFold set evaluations.

852 *Table S3.* MAHOMES II feature importance

853 *Table S4.* Comparison of MAHOMES II performance to similar tools that make  
854 enzymatic and non-enzymatic predictions

855 *Table S5. Manually annotated sites*

856 *Table S6. Feature details*

## 857 **7. Acknowledgements**

858 We gratefully acknowledge helpful feedback from members of the Slusky lab, especially  
859 Daniel Montezano and Samuel Lim. We also gratefully acknowledge NIGMS awards  
860 DP2GM128201 and P20GM103418 to JSGS, NSF award 2226804 to JSGS, and  
861 funding from the Scandinavian American Foundation to JSGS.

## 862 **8. Data and code availability**

863 The data and code used to train and evaluate MAHOMES II can be found at  
864 <https://github.com/SluskyLab/MAHOMES-II> (Feehan et al. 2022). The AlphaFold set  
865 metalloprotein structures and enzyme/non-enzyme labels can be downloaded from  
866 <https://doi.org/10.5281/zenodo.7703098> (Feehan 2023).

## 867 **8. Conflicts of interests**

868 The authors declare no conflicts of interests.

## 869 9. References

- 870 Agarwal, Pratul K. 2019. 'A Biophysical Perspective on Enzyme Catalysis', *Biochemistry*, 58:  
871 438-49.
- 872 Alford, Rebecca F., Andrew Leaver-Fay, Jeliaskov, Matthew J. O'Meara, Frank P.  
873 DiMaio, Hahnbeom Park, Maxim V. Shapovalov, P. Douglas Renfrew, Vikram K.  
874 Mulligan, Kalli Kappel, Jason W. Labonte, Michael S. Pacella, Richard Bonneau, Philip  
875 Bradley, Roland L. Dunbrack, Rhiju Das, David Baker, Brian Kuhlman, Tanja  
876 Kortemme, and Jeffrey J. Gray. 2017. 'The Rosetta All-Atom Energy Function for  
877 Macromolecular Modeling and Design', *Journal of Chemical Theory and Computation*,  
878 13: 3031-48.
- 879 Andreini, Claudia, Gabriele Cavallaro, and Serena Lorenzini. 2012. 'FindGeo: a tool for  
880 determining metal coordination geometry', *Bioinformatics*, 28: 1658-60.
- 881 Bartlett, Gail J., Craig T. Porter, Neera Borkakoti, and Janet M. Thornton. 2002. 'Analysis of  
882 Catalytic Residues in Enzyme Active Sites', *Journal of Molecular Biology*, 324: 105-21.
- 883 Bate, Paul, and Jim Warwicker. 2004. 'Enzyme/Non-enzyme Discrimination and Prediction of  
884 Enzyme Active Site Location Using Charge-based Methods', *Journal of Molecular  
885 Biology*, 340: 263-76.
- 886 Berman, Helen M., John Westbrook, Zukang Feng, Gary Gilliland, T. N. Bhat, Helge Weissig,  
887 Ilya N. Shindyalov, and Philip E. Bourne. 2000. 'The Protein Data Bank', *Nucleic Acids  
888 Research*, 28: 235-42.
- 889 Brodtkin, Heather R., Nicholas A. DeLateur, Srinivas Somarowthu, Caitlyn L. Mills, Walter R.  
890 Novak, Penny J. Beuning, Dagmar Ringe, and Mary Jo Ondrechen. 2015. 'Prediction of  
891 distal residue participation in enzyme catalysis', *Protein Science*, 24: 762-78.
- 892 Burley, Stephen K., Helen M. Berman, Charmi Bhikadiya, Chunxiao Bi, Li Chen, Luigi Di  
893 Costanzo, Cole Christie, Ken Dalenberg, Jose M. Duarte, Shuchismita Dutta, Zukang  
894 Feng, Sutapa Ghosh, David S. Goodsell, Rachel K. Green, Vladimir Guranovic, Dmytro  
895 Guzenko, Brian P. Hudson, Tara Kalro, Yuhe Liang, Robert Lowe, Harry Namkoong,  
896 Ezra Peisach, Irina Periskova, Andreas Prlic, Chris Randle, Alexander Rose, Peter Rose,  
897 Raul Sala, Monica Sekharan, Chenghua Shao, Lihua Tan, Yi-Ping Tao, Yana Valasatava,  
898 Maria Voigt, John Westbrook, Jesse Woo, Huanwang Yang, Jasmine Young, Marina  
899 Zhuravleva, and Christine Zardecki. 2019. 'RCSB Protein Data Bank: biological  
900 macromolecular structures enabling research and education in fundamental biology,  
901 biomedicine, biotechnology and energy', *Nucleic Acids Research*, 47: D464-D74.
- 902 Conway, Patrick, Michael D. Tyka, Frank DiMaio, David E. Konerding, and David Baker. 2014.  
903 'Relaxation of backbone bond geometry improves protein energy landscape modeling',  
904 *Protein science : a publication of the Protein Society*, 23: 47-55.
- 905 Coudert, Elisabeth, Sebastien Gehant, Edouard de Castro, Monica Pozzato, Delphine Baratin,  
906 Teresa Neto, Christian J. A. Sigrist, Nicole Redaschi, Alan Bridge, and Consortium The  
907 UniProt. 2023. 'Annotation of biologically relevant ligands in UniProtKB using ChEBI',  
908 *Bioinformatics*, 39: btac793.
- 909 Coulther, Timothy A., Jaeju Ko, and Mary Jo Ondrechen. 2021. 'Amino acid interactions that  
910 facilitate enzyme catalysis', *The Journal of Chemical Physics*, 154: 195101.
- 911 Coulther, Timothy A., Moritz Pott, Cathleen Zeymer, Donald Hilvert, and Mary Jo Ondrechen.  
912 2021. 'Analysis of electrostatic coupling throughout the laboratory evolution of a  
913 designed retroaldolase', *Protein Science*, 30: 1617-27.

- 914 Dudev, Todor, Lin, Minko Dudev, and Carmay Lim. 2003. 'First–Second Shell Interactions in  
915 Metal Binding Sites in Proteins: A PDB Survey and DFT/CDM Calculations', *Journal of*  
916 *the American Chemical Society*, 125: 3168-80.
- 917 Eddy, Sean R. 2011. 'Accelerated Profile HMM Searches', *PLOS Computational Biology*, 7:  
918 e1002195.
- 919 Eisenberg, D., E. Schwarz, M. Komaromy, and R. Wall. 1984. 'Analysis of membrane and  
920 surface protein sequences with the hydrophobic moment plot', *Journal of Molecular*  
921 *Biology*, 179: 125-42.
- 922 Feehan, Ryan. 2023. 'Metalloprotein AlphaFold set with enzyme/non-enzyme labeled sites',  
923 *Zenodo*.
- 924 Feehan, Ryan, Matthew Copeland, Meghan W. Franklin, and Joanna S. G. Slusky. 2022.  
925 'SluskyLab/MAHOMES-II: v1.0.0', *Zenodo*.
- 926 Feehan, Ryan, Meghan W. Franklin, and Joanna S. G. Slusky. 2021. 'Machine learning  
927 differentiates enzymatic and non-enzymatic metals in proteins', *Nature Communications*,  
928 12: 3712.
- 929 Fogolari, F., A. Corazza, V. Yarra, A. Jalaru, P. Viglino, and G. Esposito. 2012a. 'Bluues: a  
930 program for the analysis of the electrostatic properties of proteins based on generalized  
931 Born radii', *BMC Bioinformatics*, 13 Suppl 4: S18.
- 932 Fogolari, Federico, Alessandra Corazza, Vijaylakshmi Yarra, Anusha Jalaru, Paolo Viglino, and  
933 Gennaro Esposito. 2012b. 'Bluues: a program for the analysis of the electrostatic  
934 properties of proteins based on generalized Born radii', *BMC Bioinformatics*, 13: S18.
- 935 Furnham, Nicholas, Natalie L. Dawson, Syed A. Rahman, Janet M. Thornton, and Christine A.  
936 Orengo. 2016. 'Large-Scale Analysis Exploring Evolution of Catalytic Machineries and  
937 Mechanisms in Enzyme Superfamilies', *Journal of Molecular Biology*, 428: 253-67.
- 938 Jaccard, Paul. 1907. 'La distribution de la flore dans la zone alpine', *Revue Générale des Sciences*  
939 *Pures et Appliquées*, 18: 961-7.
- 940 Johnson, David K., and John Karanicolas. 2013. 'Druggable Protein Interaction Sites Are More  
941 Predisposed to Surface Pocket Formation than the Rest of the Protein Surface', *PLOS*  
942 *Computational Biology*, 9: e1002951.
- 943 Jumper, John, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf  
944 Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Židek, Anna Potapenko,  
945 Alex Bridgland, Clemens Meyer, Simon A. A. Kohl, Andrew J. Ballard, Andrew Cowie,  
946 Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back,  
947 Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin Steinegger,  
948 Michalina Pacholska, Tamas Berghammer, Sebastian Bodenstern, David Silver, Oriol  
949 Vinyals, Andrew W. Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis.  
950 2021. 'Highly accurate protein structure prediction with AlphaFold', *Nature*, 596: 583-89.
- 951 Kawabata, Takeshi. 2010. 'Detection of multiscale pockets on protein surfaces using  
952 mathematical morphology', *Proteins: Structure, Function, and Bioinformatics*, 78: 1195-  
953 211.
- 954 ———. 2019. 'Detection of cave pockets in large molecules: Spaces into which internal probes  
955 can enter, but external probes from outside cannot', *Biophysics and physcobiology*, 16:  
956 391-406.
- 957 Ko, Jaeju, Leonel F. Murga, Pierrette André, Huyuan Yang, Mary Jo Ondrechen, Ronald J.  
958 Williams, Akochi Agunwamba, and David E. Budil. 2005. 'Statistical criteria for the

- 959 identification of protein active sites using theoretical microscopic titration curves',  
960 *Proteins: Structure, Function, and Bioinformatics*, 59: 183-95.
- 961 Kumar, Narendra, and Jeffrey Skolnick. 2012. 'EFICAZ2.5: application of a high-precision  
962 enzyme function predictor to 396 proteomes', *Bioinformatics (Oxford, England)*, 28:  
963 2687-88.
- 964 Kyte, Jack, and Russell F. Doolittle. 1982. 'A simple method for displaying the hydrophatic  
965 character of a protein', *Journal of Molecular Biology*, 157: 105-32.
- 966 Lee, Yu-Ming, Cédric Grauffel, Ting Chen, Karen Sargsyan, and Carmay Lim. 2019. 'Factors  
967 Governing the Different Functions of Zn<sup>2+</sup>-Sites with Identical Ligands in Proteins',  
968 *Journal of Chemical Information and Modeling*, 59: 3946-54.
- 969 Lipsh-Sokolik, R., O. Khersonsky, S. P. Schröder, C. de Boer, S. Y. Hoch, G. J. Davies, H. S.  
970 Overkleeft, and S. J. Fleishman. 2023. 'Combinatorial assembly and design of enzymes',  
971 *Science*, 379: 195-201.
- 972 Matthews, B. W. 1975. 'Comparison of the predicted and observed secondary structure of T4  
973 phage lysozyme', *Biochimica et Biophysica Acta (BBA) - Protein Structure*, 405: 442-51.
- 974 Mazmanian, Karine, Karen Sargsyan, and Carmay Lim. 2020. 'How the Local Environment of  
975 Functional Sites Regulates Protein Function', *Journal of the American Chemical Society*,  
976 142: 9861-71.
- 977 McDonald, Andrew G., Sinéad Boyce, and Keith F. Tipton. 2009. 'ExplorEnz: the primary  
978 source of the IUBMB enzyme list', *Nucleic Acids Research*, 37: D593-D97.
- 979 Morley, Krista L., and Romas J. Kazlauskas. 2005. 'Improving enzyme properties: when are  
980 closer mutations better?', *Trends in Biotechnology*, 23: 231-37.
- 981 Ngu, Lisa, Jenifer N. Winters, Kien Nguyen, Kevin E. Ramos, Nicholas A. DeLateur, Lee  
982 Makowski, Paul C. Whitford, Mary Jo Ondrechen, and Penny J. Beuning. 2020. 'Probing  
983 remote residues important for catalysis in Escherichia coli ornithine transcarbamoylase',  
984 *PLOS ONE*, 15: e0228487.
- 985 Ondrechen, Mary Jo, James G. Clifton, and Dagmar Ringe. 2001. 'THEMATICS: A simple  
986 computational predictor of enzyme function from structure', *Proceedings of the National  
987 Academy of Sciences*, 98: 12473-78.
- 988 Parasuram, Ramya, Timothy A. Coulther, Judith M. Hollander, Elise Keston-Smith, Mary Jo  
989 Ondrechen, and Penny J. Beuning. 2018. 'Prediction of Active Site and Distal Residues in  
990 E. coli DNA Polymerase III alpha Polymerase Activity', *Biochemistry*, 57: 1063-72.
- 991 Parasuram, Ramya, Caitlyn L. Mills, Zhouxi Wang, Saroja Somasundaram, Penny J. Beuning,  
992 and Mary Jo Ondrechen. 2016. 'Local structure based method for prediction of the  
993 biochemical function of proteins: Applications to glycoside hydrolases', *Methods*, 93: 51-  
994 63.
- 995 Pedregosa, Fabian, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion,  
996 Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, and Vincent Dubourg.  
997 2011. 'Scikit-learn: Machine learning in Python', *The Journal of Machine Learning  
998 Research*, 12: 2825-30.
- 999 Pérez-Cañadillas, José Manuel, Ramón Campos-Olivas, Javier Lacadena, Alvaro Martínez del  
1000 Pozo, José G. Gavilanes, Jorge Santoro, Manuel Rico, and Marta Bruix. 1998.  
1001 'Characterization of pKa Values and Titration Shifts in the Cytotoxic Ribonuclease  $\alpha$ -  
1002 Sarcin by NMR. Relationship between Electrostatic Interactions, Structure, and Catalytic  
1003 Function', *Biochemistry*, 37: 15865-76.



- 1004 Ribeiro, António J. M., Gemma L. Holliday, Nicholas Furnham, Jonathan D. Tyzack, Katherine  
1005 Ferris, and Janet M. Thornton. 2018. 'Mechanism and Catalytic Site Atlas (M-CSA): a  
1006 database of enzyme reaction mechanisms and active sites', *Nucleic Acids Research*, 46:  
1007 D618-D23.
- 1008 Ribeiro, António J. M., Jonathan D. Tyzack, Neera Borkakoti, Gemma L. Holliday, and Janet M.  
1009 Thornton. 2020. 'A global analysis of function and conservation of catalytic residues in  
1010 enzymes', *The Journal of biological chemistry*, 295: 314-24.
- 1011 Sheldon, Roger A., and John M. Woodley. 2018. 'Role of Biocatalysis in Sustainable Chemistry',  
1012 *Chemical Reviews*, 118: 801-38.
- 1013 Somarowthu, Srinivas, Heather R. Brodtkin, J. Alejandro D'Aquino, Dagmar Ringe, Mary Jo  
1014 Ondrechen, and Penny J. Beuning. 2011. 'A Tale of Two Isomerases: Compact versus  
1015 Extended Active Sites in Ketosteroid Isomerase and Phosphoglucose Isomerase',  
1016 *Biochemistry*, 50: 9283-95.
- 1017 Somarowthu, Srinivas, and Mary Jo Ondrechen. 2012. 'POOL server: machine learning  
1018 application for functional site prediction in proteins', *Bioinformatics (Oxford, England)*,  
1019 28: 2078-79.
- 1020 Somarowthu, Srinivas, Huyuan Yang, David G. C. Hildebrand, and Mary Jo Ondrechen. 2011.  
1021 'High-performance prediction of functional residues in proteins with machine learning  
1022 and computed input features', *Biopolymers*, 95: 390-400.
- 1023 Song, Jiangning, Fuyi Li, Kazuhiro Takemoto, Gholamreza Haffari, Tatsuya Akutsu, Kuo-Chen  
1024 Chou, and Geoffrey I. Webb. 2018. 'PREvail, an integrative approach for inferring  
1025 catalytic residues using sequence, structural, and network features in a machine-learning  
1026 framework', *Journal of Theoretical Biology*, 443: 125-37.
- 1027 Thornton, Janet M., Roman A. Laskowski, and Neera Borkakoti. 2021. 'AlphaFold heralds a  
1028 data-driven revolution in biology and medicine', *Nature Medicine*, 27: 1666-69.
- 1029 Tiwari, Manish Kumar, Vipin C. Kalia, Yun Chan Kang, and Jung-Kul Lee. 2014. 'Role of a  
1030 remote leucine residue in the catalytic function of polyol dehydrogenase', *Molecular  
1031 BioSystems*, 10: 3255-63.
- 1032 Tong, Wenxu, Ying Wei, Leonel F. Murga, Mary Jo Ondrechen, and Ronald J. Williams. 2009.  
1033 'Partial order optimum likelihood (POOL): maximum likelihood prediction of protein  
1034 active site residues using 3D Structure and sequence properties', *PLOS Computational  
1035 Biology*, 5: e1000266-e66.
- 1036 Tunyasuvunakool, Kathryn, Jonas Adler, Zachary Wu, Tim Green, Michal Zielinski, Augustin  
1037 Židek, Alex Bridgland, Andrew Cowie, Clemens Meyer, Agata Laydon, Sameer  
1038 Velankar, Gerard J. Kleywegt, Alex Bateman, Richard Evans, Alexander Pritzel, Michael  
1039 Figurnov, Olaf Ronneberger, Russ Bates, Simon A. A. Kohl, Anna Potapenko, Andrew J.  
1040 Ballard, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub Jain, Ellen Clancy,  
1041 David Reiman, Stig Petersen, Andrew W. Senior, Koray Kavukcuoglu, Ewan Birney,  
1042 Pushmeet Kohli, John Jumper, and Demis Hassabis. 2021. 'Highly accurate protein  
1043 structure prediction for the human proteome', *Nature*, 596: 590-96.
- 1044 UniProt, Consortium. 2019. 'UniProt: a worldwide hub of protein knowledge', *Nucleic Acids  
1045 Research*, 47: D506-D15.
- 1046 Virtanen, P., R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski,  
1047 P. Peterson, W. Weckesser, and J. Bright. 2020. 'Scipy 1.0: Fundamental Algorithms for  
1048 Scientific Computing in Python', *Nat. Methods*, 17: 261.

- 1049 Wolfenden, Richard, Caroline Ridgway, and Gregory Young. 1998. 'Spontaneous Hydrolysis of  
1050 Ionized Phosphate Monoesters and Diesters and the Proficiencies of Phosphatases and  
1051 Phosphodiesterases as Catalysts', *Journal of the American Chemical Society*, 120: 833-  
1052 34.
- 1053 Yang, Kevin K., Zachary Wu, and Frances H. Arnold. 2019. 'Machine-learning-guided directed  
1054 evolution for protein engineering', *Nature Methods*, 16: 687-94.
- 1055 Zou, Zhenzhen, Shuye Tian, Xin Gao, and Yu Li. 2019. 'mlDEEPre: Multi-Functional Enzyme  
1056 Function Prediction With Hierarchical Multi-Label Deep Learning', *Frontiers in  
1057 Genetics*, 9.  
1058