1 SCREEN: A Graph-based Contrastive Learning Tool to Infer

2 Catalytic Residues and Assess Enzyme Mutations

- 3 Tong Pan^{1,2}, Yue Bi^{1,2}, Xiaoyu Wang^{1,2}, Ying Zhang^{1,3}, Geoffrey I. Webb⁴, Robin B. Gasser^{5,*}, Lukasz
- 4 Kurgan^{6,*}, Jiangning Song^{1,2,7,*}

5 ¹Monash Biomedicine Discovery Institute and Department of Biochemistry and Molecular Biology, Monash

- 6 University, Clayton, VIC 3800, Australia
- 7 ²Monash Biomedicine Discovery Institute-Wenzhou Medical University Alliance in Clinical and Experimental
- 8 Biomedicine, Monash University, Clayton, VIC 3800, Australia
- 9 ³School of Computer Science and Engineering, Nanjing University of Science and Technology, Nanjing
- 10 210094, China
- ⁴Department of Data Science and Artificial Intelligence, Monash University, Clayton, VIC 3800, Australia
- 12 ⁵Department of Veterinary Biosciences, Melbourne Veterinary School, The University of Melbourne,
- 13 Parkville, VIC 3010, Australia
- ⁶Department of Computer Science, Virginia Commonwealth University, Richmond 23284, USA
- ¹⁵ ⁷Key Laboratory of Clinical Laboratory Diagnosis and Translational Research of Zhejiang Province,
- 16 Department of Clinical Laboratory, The First Affiliated Hospital of Wenzhou Medical University, Wenzhou
- 17 325015, China
- 18 * Corresponding authors.
- 19 Email: <u>robinbg@unimelb.edu.au</u> (Gasser RB), <u>lkurgan@vcu.edu</u> (Kurgan L),
 20 jiangning.song@monash.edu (Song J).
- 21 Running Title: Pan T et al / SCREEN for Catalytic Residue Prediction
- 22 The counts of words (8540 words), references (64), tables (1) and figures (6), supplementary
- 23 figures (15), supplementary tables (3), The counts of letters in article title (90 characters),
- running title (44 characters), count of keywords (5 words) and words in Abstract (181 words).

25 Abstract

- 26 The accurate identification of catalytic residues contributes to our understanding of enzyme
- 27 functions in biological processes and pathways. The increasing number of protein sequences
- 28 necessitates computational tools for the automated prediction of catalytic residues in enzymes.

29 Here, we introduce SCREEN, a graph neural network for the high-throughput prediction of 30 catalytic residues via the integration of enzyme functional and structural information. SCREEN 31 constructs residue representations based on spatial arrangements and incorporates enzyme 32 function priors into such representations through contrastive learning. We demonstrate that 33 SCREEN (i) consistently outperforms currently-available predictors; (ii) provides accurate results when applied to inferred enzyme structures; and (iii) generalizes well to enzymes 34 35 dissimilar from those in the training set. We also show that the putative catalytic residues 36 predicted by SCREEN mimic key structural and biophysical characteristics of native catalytic residues. Moreover, using experimental data sets, we show that SCREEN's predictions can be 37 used to distinguish residues with a high mutation tolerance from those likely to cause functional 38 loss when mutated, indicating that this tool might be used to infer disease-associated mutations. 39 **SCREEN** https://github.com/BioColLab/SCREEN 40 is publicly available at and https://ngdc.cncb.ac.cn/biocode/tool/7580. 41

42

Keywords: Catalytic residue; Enzyme structure; Evolutionary conservation; Graph neural
network; Contrastive learning

45

46 Introduction

Enzymes are critical for a wide range of diverse biochemical, molecular and physiological 47 48 processes and pathways which sustain life [1]. The extraordinary catalytic proficiency of an enzyme is often intricately orchestrated by a selected set of amino acids within its active site(s), 49 referred to as the catalytic residues [2]. These spatially proximate catalytic residues can be 50 engaged in critical interactions with substrate molecules, catalyzing chemical reactions and 51 52 ensuring the catalytic efficiency and specificity of enzymes [3]. Catalytic residues often exhibit 53 conservation across species, particularly those within the same taxonomic groups [4], such that 54 mutations in catalytic sites can affect enzyme function(s), potentially inducing the onset of diseases, such as cancers and metabolic disorders [5]. For instance, mutations of catalytic 55 56 residues in the tumor suppressor phosphatase and tensin homolog (PTEN) have been shown to 57 culminate in various forms of cancers, such as glioblastoma multiforme, melanoma and breast 58 cancer [6]. Mutations in the catalytic sites of enzymes, such as CYP2C9, which are responsible 59 for the biotransformation of small molecule drugs, can impact individual drug responses and 60 potentially increase the risk of metabolic disorders [7].

61 The number for which enzymes with detailed catalytic residue annotations available in the 62 Mechanism and Catalytic Site Atlas (M-CSA) database [8] is substantially lower than the vast 63 number of enzyme sequences in the UniProt database [9] and enzyme structures in repositories 64 such as the Protein Data Bank (PDB) [10]. This gap relates primarily to the absence of highthroughput methods for identifying catalytic residues. Traditionally, the active sites in enzymes 65 have been established using site-directed mutagenesis and biochemical assays, providing insight 66 into the corresponding kinetic and thermodynamic parameters [11]. However, these laboratory 67 methods are relatively low throughput, time-consuming and labor-intensive, thereby restricting 68 analyses to small numbers of residues and constraining "down-stream" applications such as the 69 design of novel enzymes and inhibitors [12]. Furthermore, although M-CSA offers information 70 on enzyme functions, its coverage of the function space is not comprehensive, particularly for 71 oxidoreductases and translocases (with only 45.6% and 30% coverage, respectively), potentially 72 73 attributable to curation backlog and/or limited functional data/information [13]. There is significant demand for an *in silico* approach for the reliable and reproducible identification of 74 catalytic residues in enzymes from sequence and/or structure data to enable the exploration of 75 76 enzyme functions and accelerate biomolecular design.

A significant effort has been directed towards developing computational methods for the 77 identification of catalytic residues in enzymes [14]. These methods include homology-based 78 79 approaches, machine learning-based, and deep learning-based approaches. Homology-based 80 methods annotate catalytic residues by comparing the query enzyme's sequence or structural 81 similarity to that of the target enzymes with experimentally validated residues. Wallace et al. firstly introduced TESS [15], a program that uses a geometric hashing algorithm to identify 82 enzyme active sites by aligning the query enzyme to a structural template database. Mistry et al. 83 later proposed a sequence-based method that transfers previously verified catalytic residues to 84 85 other chains within the same Pfam family using a set of strict rules [16]. However, these 86 Homology-based methods are constrained by the availability of reliable templates. Compared 87 with homology-based methods, machine learning-based, and deep learning-based methods 88 utilize the "ground truth" annotations of catalytic residues from curated public databases to train 89 predictive models, which, in turn, can be applied to infer catalytic residues in most unknown enzymes [17]. For one of the earliest tools, Gutteridge et al. [18] trained a neural network model 90

91 to identify catalytic residues using enzyme structure- and sequence-derived features as inputs.
92 Subsequently developed predictors often employed support vector machine (SVM) models or
93 random forest models to a collection of manually curated features [19–23]. Interestingly, Chea
94 et al. [24] did not use a machine learning model but, instead, predicted catalytic residues using
95 statistical scores calculated from a network representation of protein structure and solvent
96 accessibility.

97 Structural information of the enzymes has proven highly effective in addressing diverse 98 challenges, ranging from predicting enzyme function [25] to guiding enzyme engineering [26]. 99 However, most existing methods for catalytic residue prediction, which rely primarily on enzyme sequences or manually curated structural features, often struggle to capture the spatial 100 101 arrangements of amino acid residues, as catalytic residues tend to form spatial clusters. Moreover, 102 the intricate chemical reactions catalyzed by enzymes are typically driven by a small subset of catalytic residues. Despite their critical role in enzyme function, current methods often fail to 103 104 incorporate comprehensive enzyme functional data, thereby limiting the ability to fully explore the connections between catalytic residues, enzyme structure, and function. 105

106 We anticipate that recognizing patterns in the spatial arrangement of residues within enzyme 107 structures can substantially enhance the performance of catalytic residue prediction tools and deepen our understanding of enzyme function through the use of modern deep neural networks. 108 To this end, we propose here a deep learning-based solution, called SCREEN, for the accurate 109 110 prediction of catalytic residues in enzymes. SCREEN employs a graph neural network that 111 models the spatial arrangement of active sites in enzyme structures and combines data derived 112 from enzyme structure, sequence embedding, and evolutionary information obtained by using two complementary methods – BLAST (Basic Local Alignment Search Tool) [27] and HMMER 113 114 (sequence analysis tool using profile hidden Markov models) [28]. Moreover, we apply the contrastive learning framework to further enhance the predictive performance of SCREEN by 115 116 incorporating enzyme functional information.

117

118 Method

119 **Training and test datasets**

120 We curated a dataset comprising 1055 enzymes with annotated catalytic residues, which we used

121 to train and optimize our predictive model (Figure 1A). First, we collected data from M-CSA

122 database, which contains catalytic residue annotations and details about enzymic reaction mechanisms [29]. We combined the EF family dataset that includes catalytic residue annotations 123 124 for enzyems from different SCOP families [22]. We filtered the combined dataset by clustering 125 the proteins with the CD-HIT software at 40% sequence identity [30], then randomly selected 126 one protein from each cluster. This prevented overfitting the model into larger clusters of similar enzymes. Next, we collected enzyme three-dimensional structures from the Protein Data Bank 127 128 (PDB) [9]. We obtained the Enzyme Commission (EC) numbers [31] using the SIFTS database 129 [32]. We used the first-level EC numbers due to the relatively sparse/incomplete nature of these data at the lower levels [33]. We shuffled the dataset randomly and then divided it into training 130 (90%) and validation (10%) subsets (Figure S1). 131

We acquired five widely-used test datasets to conduct a comparative evaluation of our model 132 against existing tools. These test sets included the EF-fold dataset, EF-superfamily dataset, and 133 HA superfamily dataset, representing enzymes from every SCOP fold and superfamily 134 135 respectively. We also collected the PC test dataset [20], which was originally obtained from Catalytic Residue Dataset (CATRES) and represents proteins of the PIRSF protein groups [34]. 136 Finally, we obtained the NN test dataset [18], which comprises enzymes of six main classes. 137 138 Importantly, we excluded enzymes from these five test datasets from our training/validation dataset. Table S1 summarizes the training and the five test datasets. 139

140

141 **Overview of the SCREEN model**

142 SCREEN is a supervised deep learner that integrates information derived from atomic structures, 143 sequences, and evolutionary profiles. Specifically, SCREEN presents the input (enzyme structures) as graphs at the residue level, leveraging evolutionary information, sequence 144 145 embeddings – generated by a modern language model – and relevant structural characteristics, such as B-factors and solvent accessibility (Figure 1B) [35]. We correspondingly employed a 146 147 graph convolutional neural network to generate propensities for catalytic residues from these 148 inputs. The training process employed enzyme function information (Enzyme Commission 149 numbers) [36] via a contrastive learning framework, utilizing the Triplet Margin Loss function 150 to enable clustering enzymes of the same classes and separating enzymes from different classes 151 in a latent feature space (Figure 1C). This allowed our model to develop class-specific latent feature spaces, leading to improvements in predictive performance/capacity. Moreover, we 152

employed a dynamic training strategy, in which we initially trained the network using the
contrastive learning and then applied our model to accurately identify catalytic residues (Figures
S2 and S3).

156

157 Graph-based representation of protein structure

We represented the input enzyme structure composed of n residues $Enz = \{r_1, r_1, ..., r_n\}$ as an attributed graph encoded with evolutionary context, sequence embeddings and relevant structural characteristics, such as solvent accessibility and B-factors. Specifically, the graph representation G = (V, A, X) consists of three-dimensional enzyme structure by taking the residue set $V \subseteq Enz$ as graph nodes, the adjacency matrix A with $n \times n$ size that quantifies connectivity of nodes/residues, and the feature matrix $X \in \mathcal{R}^{n \times \theta}$.

The feature matrix $X = (X_L, X_G, X_A)$ covers the evolutionary, sequence and structural features. 164 165 The X_L descriptor quantifies evolutionary conservation by utilizing two complementary tools: PSI-BLAST, which is a heuristic algorithm that relies on the dynamic programming [27], and 166 167 HMMER that is based on the hidden Markov model (HMM) [28]. We run PSI-BLAST on the NCBI's non-redundant (nr) database, with three iterations and the E-value threshold of $< 10^{-3}$. 168 We normalize the output position-specific scoring matrix (PSSM) of size $n \times 20$ with the 169 sigmoid function: $\overline{x} = \frac{1}{1+e^{-x}}$. We use HMMER with the uniclust30 database [37] and default 170 parameters to generate the $n \times 30$ HMM matrix that we normalize to the [0, 1] range [28]. The 171 172 X_G descriptor captures sequence information computed by ProtT5 model, a deep learning language model that was pre-trained on 390 billion amino acids [38]. The enzyme sequence is 173 encoded into residue-level feature embeddings denoted as $X_G \in \mathcal{R}^{n \times h_1}$, where h_1 defaults to 174 175 1024. These vectors encapsulate information about individual residues that are adjacent in the sequence, and broader protein-level information. Lastly, the X_A descriptor encompasses several 176 177 key properties that are derived from the atomic-level data: atom types and atomic mass when excluding hydrogen atoms, **B**-factor, residue side-chain presence, the count of bonded hydrogen 178 179 atoms, ring membership, van der Waals radius, and solvent accessibility. Given that residues might have different numbers of atoms, we compute the average values across all atoms, 180 resulting in atomic descriptor $X_A \in \mathcal{R}^{n \times h_2}$, with $h_2 = 14$. 181

182

183 **Predictive model**

184 We designed the graph convolutional neural network (GCN) with three convolutional layers to 185 facilitate the propagation of feature embeddings for residues that share spatial proximity.

For a given graph defined by the adjacency matrix $A \in \{0, 1\}^{n \times n}$ and the feature matrix X =(X_L, X_G, X_A), our model produces residue-level representations $H^{(i)} \in \mathcal{R}^{n \times d_i}$ where d_i represents the embedding dimension for the *i*th convolutional layer.

189
$$H^{(i)} = GCN(A, [X_L, X_A])$$
 (1)

190 We refine residue representations through the process of neighbor aggregations as follows:

191
$$H^{(i)} = ReLU\left(\tilde{D}^{-\frac{1}{2}}(A+I_n)\tilde{D}^{-\frac{1}{2}}H^{(i-1)}W^{(i)}\right)$$
(2)

192
$$H^{(0)} = [X_L, X_A]$$
(3)

where $I_n \in \mathcal{R}^{n \times n}$ is the identity matrix, $\widetilde{D} \in \mathcal{R}^{n \times n}$ is the diagonal degree matrix with entries $D_{ii} = \sum_j (A + I_n)_{ij}$, $W^{(1)} \in \mathcal{R}^{\theta \times d_i}$ is the trainable weight matrix for the *i*th convolutional layer, ReLU denotes the Rectified Linear Unit activation function, and [] denotes the concatenation operation. The above architecture generates graph representation $X_E \in \mathcal{R}^{n \times d}$, where d =512 (Figure S4), which we combine using multilayer perception (MLP) network as follows:

199
$$X_E = ReLU\left(MLP([ReLU(MLP([H^1, H^2, H^3])), ReLU(MLP(X_G))])\right)$$
(4)

We employ three fully connected layers in the MLP network to reduce the feature space to the final output vector $Y \in \mathbb{R}^{n \times 2}$ that gives numeric propensities for putative catalytic residues.

202

203 Contrastive learning

We used contrastive learning with Triplet Margin Loss to craft enzyme representations that improve the catalytic residue predictions. Using the graph representation $X_E \in \mathcal{R}^{n \times d}$, we employed average aggregation across residues, generating a fixed-sized sequence representation vector $Z \in \mathbb{R}^{1 \times h_3}$, with h_3 set to 1024. During training in each epoch, we iteratively refined every sequence representation vector and computed enzyme class cluster centres. When training with a query enzyme \mathbf{z}_a , we selected the enzyme cluster centre embedding from the same enzyme class as the positive sample z_p , and randomly sampled another cluster centre from a different enzyme class as the negative sample z_n , which resulted in the following Triplet Margin Loss function:

213

$$\mathcal{L}^{TM} = \| z_a - z_p \|_2 - \| z_a - z_n \|_2 + \alpha$$
(5)

where we set the margin α to the default value of 1. This loss function minimizes the Euclidean distance between enzyme representations belonging to the same main enzyme class while maximizing the distance between those form different main enzyme classes. We implemented a dynamic training strategy, where we performed contrastive learning for enzyme classification during early training epochs, and gradually shifted towards the default training that converges to produce accurate propensities for catalytic residues.

220

221 The multiplexed assays of variant effects data analyses

222 We gathered the multiplexed assays of variant effects (MAVE) measurements for four enzymes, 223 which provided insight into the impact of a broad collection of substitutions on both enzyme function [6,39]. We categorized the missense variants of PTEN into two main groups: functional 224 225 or inactive, regardless of their effect on abundance. The classification thresholds for the scores 226 generated by each MAVE were guided by an established methodology [40]. Specifically, we 227 used a minimal number of Gaussians (three) to ensure a reliable fit to the variant score 228 distributions, and the intersection point between the first and last Gaussian served as the classification cut-off. Adopting this binary classification approach allowed us to categorize 229 variants into two classes: (1) Wild Type-Like (WTL), variants characterized by high activity; 230 231 (2) Functional Loss (FL), variants assigned that exhibit low activity.

232

233 **Results**

234 SCREEN accurately predicts catalytic residues

We collected five commonly-used test datasets to comparatively assess SCREEN against eight current solutions. These test datasets included the EF superfamily and EF fold datasets [22], the HA superfamily dataset [24], the NN dataset [18], and the PC dataset [20]. We compared the results from SCREEN with those of a conventional sequence-based method, CRpred [19], and six tools that employed different predictive models based on enzyme structures. These tools 240 included a neural network-based approach [18], three SVM-based methods [20], a random 241 forests-based PREvaIL [23], and a statistical approach [24]. We compared the precision, recall, 242 and F1 score to evaluate the predictions of catalytic residues, referencing the reported 243 performance from the original paper (Table 1, Figure S5). We also evaluated the recently 244 proposed graph-based method, AEGAN [41], by retraining and testing it on the same training and test sets, with the hyperparameter for negative sample size set to 20. We showed that 245 246 SCREEN consistently outperformed the eight tools for the five test datasets. Compared to these 247 methods, SCREEN achieved a higher F1 score across all five test datasets. The high F1 scores achieved by SCREEN were coupled with balanced and high values of precision and recall that 248 ranged between 61.0 and 69.3 and between 61.2 and 82.0, respectively. We also quantified and 249 250 compared two other popular metrics, the area under the receiver operating curve (AUC) and the 251 area under the precision-recall curve (AUPR). Figure 2A reveals that SCREEN achieved substantially higher AUC and AUPR scores when compared with the latest structure-based 252 253 (PREvaIL) and the sequence-based (CRpred) tools, except for the EF superfamily and EF fold 254 dataset, where AUC and AUPRC values were comparable.

Using SCREEN, we also measured Best-F1 score and AUPR values for specific enzyme types (Figure S6). Particularly, for hydrolases, which represent the largest portion of the training dataset (313 of 1055; 29.7%), SCREEN obtained notable consistency across the five test datasets, with Best-F1 scores ranging from 0.63 to 0.78, AUPR values ranging from 0.56 to 0.68. For the isomerase data, the predictive performance was particularly high, with the Best-F1 score exceeding 0.74 and AUPR surpassing 0.80, even though these enzymes represent only a small portion of the training set (87 of 1055; 8.2%).

262

The use of enzyme structure information in SCREEN markedly improves the catalytic residue prediction

The catalytic residues typically tend to form cohesive clusters within the three-dimensional enzyme structures. Thus, we systematically investigated the spatial distribution of residues in enzyme structures by measuring the Euclidian distances to the nearest catalytic residue for both catalytic and non-catalytic residues in individual protein sequences. Figure 2B shows there was a clear difference in the distribution of Euclidian distances, with the catalytic residues peaking at ~ 6Å, and most non-catalytic residues exceeding 15Å, supporting that catalytic residues form cohesive active enzymatic sites. We investigated whether SCREEN could reconstruct the same spatial distributions for enzyme catalytic residues. A performance assessment of SCREEN using five test datasets (Figure 2C, Figure S7) showed that the majority of catalytic residues predicted grouped together in the structure and that the distance values were consistent among the datasets, with median values being ~ 6Å. These results agree with findings presented in Figure 2B, implying SCREEN accurately captures the spatial distribution of predicted catalytic residues.

These findings suggest the use of structure information in SCREEN likely results in predictive performance improvements. We further investigated whether SCREEN could improve results compared with a "baseline model" that excludes structure-based inputs and replaces the graph convolutional neural network (GCN) with a sequence-based convolutional neural network (CNN). The comparison of SCREEN with the baseline using five test datasets employing the Best-F1 score, AUPR, and AUC metrics (Figure 2D) showed that the use of graph network led to a marked improvement in predictive performance.

284 To assess predictions, we selected enzymes from the EF superfamily dataset ranked around the 10%, 50%, and 90% percentiles based on the Best-F1 score. We plotted catalytic residues 285 predicted by the structure-based SCREEN and the sequence-based baseline model against 286 287 ground-truth catalytic residues within enzyme structures. In the top 10%, for carboxylic ester hydrolase (PDB ID:1LE6), SCREEN accurately predicted all catalytic residues, whereas the 288 sequence-based baseline introduced three false positives (Figure 2E). In the median ranking, for 289 290 casein kinase-1 (PDB ID:1CSN), SCREEN successfully identified Asp-131, Lys-133, and Thr-291 181 as key residues, but the baseline model failed to identify Thr-181 (Figure 2F). In the 90% 292 ranking, for dihydropteroate synthase (PDB ID:1AD1), SCREEN identified Arg-239 and gave one false positive, whereas the baseline misidentified two non-essential residues (Figure 2G). 293 294 These findings indicate that SCREEN has a superior performance compared with the sequence-295 based baseline model. This supports our design and, in particular, the use of the graph network 296 and enzyme structure as a key input. We also investigated various graph convolution types, 297 including the extensively employed Graph Convolutional Layer (GCN), Graph Attention (GAT), 298 and Graph Isomorphism Network (GIN), but none of the models outperformed another using the 299 same test sets and metrics (Figure S8).

300

301 SCREEN accurately predicts catalytic residues using structure models

302 Although we showed SCREEN's predictive performance benefits from the use of enzyme structure data, this information is often missing for many proteins/enzymes. Recent advances in 303 304 protein structure prediction, like the AlphaFold algorithm [42-44], make it possible to accurately 305 predict protein structure from sequences and to use such structure models as the input to 306 SCREEN. We evaluated whether the use of predicted structures rather than experimentally 307 determined structures would alter SCREEN's predictive performance (Table S2). We generated C_{α} - C_{α} contact maps for enzymes based on experimental protein structures sourced from PDB, 308 as well as putative structures from AlphaFold. We tested SCREEN's performance by using each 309 310 of these two sets of contact maps, comparing it against a sequence-based baseline model devoid of structural information. Figure 3A showed that SCREEN consistently benefits from utilizing 311 putative enzyme structures (with Best-F1 = 0.702, AUPR = 0.644, and AUC = 0.985) compared 312 313 with sequence data alone (with Best-F1 = 0.691, AUPR = 0.624, and AUC = 0.973).

314 To further assess SCREEN's denoising power on predicted structure error, we evaluated the model performance employing AlphaFold-predicted structures with varying quality. Specifically, 315 316 we quantified the quality of predicted structures employing root-mean-square deviation (RMSD) metric as compared with experimental solved structures (RMSD = 0). Figure 3B revealed that 317 318 SCREEN's performance was better employing AlphaFold-derived structures than using 319 sequence data alone. SCREEN also out-performed currently-employed tools CRHunter and PREVAIL across the entire RMSD range, achieving the Best-F1 score of > 0.6 using AlphaFold-320 321 predicted structures, contrasting average Best-F1 scores of 0.45 and 0.26 for CRHunter and 322 PREvaIL, respectively.

323 SCREEN performed relatively well using predicted structures (Figure 3A and B), irrespective 324 of the quality of predictions via AlphaFold. Here, we used two examples to illustrate SCREEN's 325 ability to accurately identify catalytic residues even in relatively low-quality predicted structures (Figure 3C). For the human calcineurin heterodimer (PDB ID: 1AUI, Chain A) [45], where the 326 327 AlphaFold-predicted structure had an RMSD score of up to 6.76 Å, SCREEN successfully identified all 10 catalytic residues. Similarly, SCREEN accurately identified catalytic residues 328 (Ser-53, Pro-54, and Asp-96) in the PVUII DNA methyltransferase (PDB ID: 1BOO, Chain A) 329 [46], for which the structure predicted had an RMSD score of 4.75 Å. Taken together, these 330 331 results suggested that SCREEN can accurately predict catalytic residues from AlphaFold-332 predicted enzyme structures, which might be attributed to the robustness of the input features and how they are represented in the graph network model.

334

335 SCREEN predicts catalytic residues in "previously-unseen" enzymes

336 We investigated SCREEN's ability to generate accurate predictions for "previously unseen" 337 enzymes. To this end, we categorized enzymes in the five test sets into three distinct groups 338 based on their sequence identities, namely $\leq 30\%$ (low), 30 to 70% (moderate) and > 70% (high), 339 to enzymes in the training dataset. Figure 3D showed SCREEN's Best-F1 scores across three sequence identity ranges, *i.e.*, \leq 30% (low), 30 to 70% (moderate), and > 70% (high), for each 340 341 of the five test datasets. We showed that SCREEN consistently outperformed the sequence-based baseline model for each of the five datasets and the three identity ranges. Importantly, 342 343 SCREEN's predictions were accurate also for enzymes with limited sequence identities to those 344 in the training datasets, achieving the Best-F1 scores of 0.725, 0.738, 0.652, 0.718, and 0.794. These predictions were significantly better than those achieved using the existing tools, such as 345 346 CRHunter (Best-F1 values of 0.328, 0.315, 0.450, 0.365, and 0.140, respectively) and PREVAIL (Best-F1 of 0.264, 0.261, 0.263, 0.264, and 0.263, respectively). We also used CATH [47] 347 assignments to evaluate the model's robustness on enzymes sharing no homologous 348 349 superfamilies with those in the training dataset, as shown in Table S3 and Figure S9.

350

351 Model training for enzyme classes improves the prediction of catalytic residues

352 We investigated whether the training of the graph network model with distinct enzyme function 353 implications would improve predictive performance, considering that the small subset of 354 catalytic residues contributes to the intricate functions of enzymes. We used enzyme class information (first-level EC numbers) to refine enzyme representation through a contrastive 355 learning framework during the training process. This led to a separation of latent feature spaces 356 in SCREEN's deep network model for different types of enzymes. We displayed these latent 357 358 feature spaces among different enzyme classes employing t-distributed stochastic neighbor 359 embedding (t-SNE) [48]. We showed that predictive performance (metrics: F1, AUPR, and AUC) 360 of SCREEN was enhanced compared to when enzyme function was not incorporated for each of the data sets (EF fold, HA superfamily, EF superfamily, NN, and PC) (Figure 4A, Figure S10) 361 362 and that SCREEN was able to group enzymes with similar functions together and separating enzymes with distinct functions (Figure 4C). Taken together, these findings indicate that the 363

enzyme function incorporation through contrastive learning during the training process improves
predictive performance as SCREEN can differentiate catalytic from non-catalytic residues for
different distinct types of enzymes (Figures S11 and S12).

367

368 SCREEN can capture selected features of catalytic residues

We analyzed catalytic residues predicted by SCREEN, in order to investigate whether they 369 370 possess structural and biophysical characteristics expected for enzymes. To better understand the 371 relevance of the features learned by SCREEN, we initially displayed the general chemical properties (including hydrophobicity, charge, and hydrogen bonds), along with low-dimensional 372 projections of residue-level representations (Figure 4B). We observed that charged amino acids 373 374 dominated in the catalytic residues predicted (Figure 4D), consistent with previous findings showing that electrostatic filtering has a marked effect on enzyme substrate selection [49]. 375 Moreover, catalytic residues were inferred to be more rigid (structurally) than non-catalytic 376 377 residues (based on low vs. high B-factor values; see Figure 4E), which accords with a previous 378 study of native catalytic residues [50]. Fewer hydrophobic residues were associated with 379 catalytic residues (Figure 4D), which is consistent with limited solvent accessibility (Figure 4F) 380 and suggests substrate avoidance in substrate-enzyme interactions [51]. Collectively, these results show that SCREEN can capture key features that typify native catalytic residues in 381 382 distinct classes of enzymes.

383

384 Linking catalytic residues to enzyme function and structure

Based on these catalytic residues, we further analyzed the sequence-structure-function relationship of enzymes to gain deeper insights into enzymes' catalytic mechanisms. We categorized enzymes according to their catalytic functions defined by third-level EC numbers and then by (complete) fourth-level EC numbers (which link to substrates). For enzyme clusters sharing the same catalytic function/mechanism, we assessed structural similarity by their TMscores among cluster members [52], selected enzymes from individual clusters and mapped the catalytic residue predictions to respective three-dimensional structures (**Figure 5**).

The results indicated that same catalytic motif may represent distinct enzymes that serve different functions linked to diverse ligands or substrates. Figure 5A shows both 4hydroxyproline betaine 2-epimerase (PDB ID: 4h2h, Chain A) and L-Ala-D/L-Glu epimerase (PDB ID: 1TKK, Chain A) belong to the same superfamily and share a common catalytic motif
(KDEDK) [53], but they differ in their substrate specificity: 4-hydroxyproline betaine 2epimerase facilitates the 2-epimerization of trans-4-hydroxy-L-proline betaine (tHyp-B) to cis4-hydroxy-D-proline betaine (cHyp-B), whereas L-Ala-D/L-Glu epimerase catalyzes the
reversible epimerization of L-Ala-D-Glu to L-Ala-L-Glu.

400 Enzymes in a particular family can have similar structures and functions, despite undergoing 401 sequence divergence through evolution. A compelling illustration emerges when we studied two 402 related enzymes, the glutamate racemase (PDB ID: 1b73, Chain A) [54] and aspartate racemase 403 (PDB ID: 1jfl, Chain A) [55] (Figure 5A). Despite enabling similar reactions via the same mechanism, their catalytic residues are significantly different but have analogous tertiary 404 405 structures. Another example relates to protein-tyrosine-phosphatase non-receptor class (PDB ID: 1ytw, Chain A) [56] and protein-tyrosine-phosphatase non-receptor type 1 (PDB ID: 1bzc, Chain 406 A) [57] (Figure 5B). Here, although both enzymes are tyrosine phosphatases and catalyze the 407 same reaction to remove phosphoryl groups from tyrosine residues in proteins, their respective 408 catalytic residues are distinctly different. 409

The shared structural arrangements of catalytic residues can be associated with functional similarity. Figure 5B shows phosphatase 5 (PDB ID: 1S95, Chain A) and phosphatase 2B (PDB ID: 1aui, Chain A) exhibite significant structural similarity and share catalytic residues (motif: DHDDRNHHRH) pertaining to serine/threonine phosphatase function(s), characterized by executing a "nucleophilic assault" on the phosphorus atom within a phosphorylated serine or threonine residue [45].

416 Although enzymes catalyzing the same reactions often exhibit marked sequence and/or structural similarity, exceptions exist where structurally dissimilar enzymes enable similar 417 418 reactions via the same mechanism. This is expected since enzymes facilitate numerous reactions 419 using a finite set of building blocks in their residues, resulting in multiple enzymes inevitably 420 sharing components of their catalytic mechanisms. Here, we showed that non-homologous 421 proteins, protein phosphatase 5 (PDB ID: 1S95, Chain A) and dual-specificity phosphatase (PDB 422 ID: 1d5r, Chain A) employ distinct structural motifs to execute the same reaction that dephosphorylates a phosphoprotein substrate (Figure 5B, Figure S13). 423

424

425 Associating mutations with the SCREEN-predicted catalytic pockets

Here, we studied the patterns of mutations in the context of their proximity to the catalytic
pockets predicted by SCREEN. We collected the multiplexed assays of variant effects data,
probing mutation effect(s) on the functions of four different enzymes, namely PTEN tumor
suppressor (PDB ID: 1d5r, Chain A) [58], human cytochrome P450 CYP2C9 (PDB ID: 1og5,
Chain A) [59], NUDT15 (PDB ID: 5lpg, Chain A) [60], and *Escherichia coli* TEM1 betalactamase (PDB ID: 1btl, Chain A) [61], encompassing a total of 15,665 variants across 1,343
residues.

433 To function effectively, enzymes must be present at sufficiently high levels and have suitable catalytic residues in the active sites [62]; however, mutations can affect both of these aspects, 434 435 potentially resulting in impaired enzymatic function. We used the MAVE data [63] for two residue classes: (1) wild type-like (WTL) residues that exhibit high functional tolerance to 436 mutations, whose most missense mutations do not adversely impact enzyme function; and (2) 437 functional loss (FL) residues that are prone to mutations that either decrease abundance (e.g., 438 439 unstable structures) and/or impair function, leading to diminished enzyme activity (Figure S14). We systematically analyzed key characteristics of mutations in the context of the predicted 440 catalytic residues. Specifically, we applied an additional tree-structured model to SCREEN 441 442 (Figure 6A). The Euclidean distance values to the closest catalytic residue predicted by 443 SCREEN combined with solvent accessibility, which, as expected, were inferred to vary according to residue type (WTL or FL), allowing to differentiate among different mutation 444 445 groups. Figure 6B illustrates the catalytic residues predicted by SCREEN along with residue 446 mutation type predictions. We performed five-fold cross-validation on all MAVE data; our 447 results revealed an average accuracy of 0.70 and 0.84, and precision of 0.58 and 0.88 on the 448 validation data and the entire dataset, respectively (Figure 6C). We found distinct spatial 449 distribution patterns for the WTL and FL residues based on their Euclidean distances from the putative catalytic residues (Figure 6D). Our result aligns well with experimental data showing 450 451 that FL residues are relatively close to the catalytic site, while WTL residues are distributed 452 throughout the structure (Figure S15). This result indicates that SCREEN can be useful to 453 establish the impact of mutations on enzyme structure and function and provides a tool to guide the identification of disease-associated mutations in enzymes. 454

455

456 **Discussion**

457 Enzymes can catalyze a broad set of chemical reactions using a limited set of catalytic residues 458 [64]. Identifying these residues allows us to understand how existing enzymes function at the 459 molecular level and to design new ones. In this work, we hypothesize that the structural 460 organization of catalytic residues in spatial space, along with their generally high evolutionary 461 conservation, collectively contributes to catalytic residue identification. To this end, we conceptualized, designed, and assessed SCREEN, a structure-based graph network that uses 462 463 functional priors through contrastive learning and combines structure-, sequence-, and 464 evolutionary profile-based representations to accurately predict catalytic residues in enzymes.

Comparative empirical assessments using five commonly-utilized test datasets and seven 465 currently-available (published) predictors revealed that SCREEN (1) accurately predicts 466 catalytic residues in known and computationally-modeled enzymes; (2) outperforms current 467 tools; and (3) generalizes well to enzymes that have limited similarity to enzymes used to train 468 the model, suggesting that SCREEN is applicable to currently-unknown enzymes. Incorporating 469 470 enzyme function as a prior could improve the prediction of catalytic residues by enhancing the consistency of enzyme's latent representations based on their functions. However, we did not 471 472 explore enzyme function prediction in depth, e.g., addressing questions such as can we also 473 utilize this model to solve enzyme function classification tasks? A more comprehensive understanding of the spatial distribution across diverse enzyme functions and the sequence-474 structure-function relationship could be achieved by analyzing a larger sample of enzymes that 475 476 thoroughly covers the EC space. Additionally, our analysis is limited by treating enzymes as 477 independent units. Enzymatic reactions involve multiple residues, substrates, and cofactors 478 interacting across various chemical steps [13]; as such, an integrated analysis would be necessary 479 for a more comprehensive understanding of the catalytic chemical activity. Further, future research should focus on determining the level of confidence that can be assigned to model 480 predictions of catalytic residues, as well as exploring the techniques that can effectively assess 481 482 this confidence.

We demonstrate that SCREEN could infer key structural and biophysical features, including amino acid charge, solvent accessibility and structural rigidity, of predicted and known catalytic residues. Further, We undertook sequence-structure-function analyses to link catalytic residues to enzyme structure and function. Our analyses revealed that while enzymes that catalyze identical reactions often display significant sequence and/or structural similarity, exceptions 488 arise wherein dissimilar sequences and/or structures can catalyze reactions via the same 489 mechanism. In addition, using multiplexed enzyme mutation data, we showed that SCREEN 490 could infer the tolerance of individual catalytic residues to mutations and, thus, predict which 491 mutations in catalytic residues likely lead to the functional loss of an enzyme. Taken together, 492 SCREEN should provide a useful tool for the reliable prediction of catalytic residues to support 493 studies of known and unknown enzyme groups/classes as well as enable *in silico* investigations 494 of diseases linked to mutations.

495

496 **Conclusion**

497 SCREEN is an efficient and robust method for high-throughput prediction of catalytic residues 498 by integrating enzyme functional and structural information. We demonstrate SCREEN's 499 effectiveness and robustness across various widely used datasets, illustrating that the predicted 500 putative catalytic residues closely align with the key structural and biophysical characteristics of native catalytic residues. Furthermore, we performed sequence-structure-function analyses to 501 502 establish connections between catalytic residues and enzyme structure and function. This highlights SCREEN's potential for reliably predicting catalytic residues in both known and 503 504 unknown enzyme groups/classes, thereby supporting studies of the molecular mechanisms underlying enzyme functions. 505

506

507 Code availability

All source data needed to evaluate the conclusions in the paper can be found at <u>https://huggingface.co/datasets/Biocollab/SCREEN/tree/main</u>. All source code is available at <u>https://github.com/BioColLab/SCREEN and https://ngdc.cncb.ac.cn/biocode/tool/7580</u>.

511

512 **CRediT author statement**

Tong Pan: Conceptualization, Methodology, Investigation, Writing – original draft. Yue Bi:
Visualization. Xiaoyu Wang: Visualization. Ying Zhang: Visualization. Geoffrey I. Webb:
Methodology, Supervision. Robin B. Gasser: Investigation, Supervision, Writing – review &
editing. Lukasz Kurgan: Methodology, Investigation, Supervision, Writing – review & editing.
Jiangning Song: Conceptualization, Methodology, Investigation, Supervision, Writing – review

- 518 & editing. All authors have read and approved the final manuscript.
- 519

520 Competing interests

- 521 All authors declare they have no competing interests.
- 522

523 Supplementary material

524 Supplementary material is available at *Genomics, Proteomics & Bioinformatics* online 525 (https://doi.org/10.1093/gpbinl/qzaxxxx).

526

527 Acknowledgments

- 528 Financial support to RBG and JS from the Australian Research Council (Grant No. ARC
- 529 LP220200614) is gratefully acknowledged.
- 530

531 **ORCID**

- 532 0009-0003-6676-5727 (Tong Pan)
- 533 0009-0001-6220-7475 (Yue Bi)
- 534 0000-0003-4444-6197 (Xiaoyu Wang)
- 535 0000-0003-1792-0121 (Ying Zhang)
- 536 0000-0001-9963-5169 (Geoffrey I. Webb)
- 537 0000-0002-4423-1690 (Robin B. Gasser)
- 538 0000-0002-7749-0314 (Lukasz Kurgan)
- 539 0000-0001-8031-9086 (Jiangning Song)
- 540

541 **Reference**

- 542 [1] Benkovic SJ, Hammes-Schiffer S. A perspective on enzyme catalysis. Science 543 2003;301:1196–202.
- 544 [2] O'maille PE, Malone A, Dellas N, Andes Hess Jr B, Smentek L, Sheehan I, et al. Quantitative
- 545 exploration of the catalytic landscape separating divergent plant sesquiterpene synthases. Nat
- 546 Chem Biol 2008;4:617–23.
- 547 [3] Zhou J, Yan W, Hu G, Shen B. Amino acid network for prediction of catalytic residues in
- 548 enzymes: a comparison survey. Curr Protein Pept Sci 2016;17:41–51.

- 549 [4] Makarova KS, Wolf YI, Iranzo J, Shmakov SA, Alkhnbashi OS, Brouns SJ, et al.
- Evolutionary classification of CRISPR-Cas systems: a burst of class 2 and derived variants. Nat
 Rev Microbiol 2020;18:67–83.
- 552 [5] Adzhubei IA, Schmidt S, Peshkin L, Ramensky VE, Gerasimova A, Bork P, et al. A method 553 and server for predicting damaging missense mutations. Nat Methods 2010;7:248–9.
- 54 [6] Mighell TL, Evans-Dutson S, O'Roak BJ. A saturation mutagenesis approach to
- understanding PTEN lipid phosphatase activity and genotype-phenotype relationships. Am J
- 556 Hum Genet 2018;102:943–55.
- 557 [7] Markin CJ, Mokhtari DA, Sunden F, Appel MJ, Akiva E, Longwell SA, et al. Revealing 558 enzyme functional architecture via high-throughput microfluidic enzyme kinetics. Science 559 2021:373:eabf8761.
- 560 [8] Porter CT, Bartlett GJ, Thornton JM. The Catalytic Site Atlas: a resource of catalytic sites 561 and residues identified in enzymes using structural data. Nucleic Acids Res 2004;32:D129–D33.
- and residues identified in enzymes using structural data. Nucleic Acids Res 2004;32:D129–D33.
 [9] Berman HM, Westbrook J, Feng Z, Gilliland G, Bhat TN, Weissig H, et al. The protein data
- 563 bank. Nucleic Acids Res 2000;28:235–42.
- 564 [10] Consortium U. UniProt: a worldwide hub of protein knowledge. Nucleic Acids Res
- 565 2019;47:D506–D15.
- 566 [11] Storici F, Lewis LK, Resnick MA. In vivo site-directed mutagenesis using oligonucleotides.
- 567 Nat Biotechnol 2001;19:773–6.
- 568 [12] Süel GM, Lockless SW, Wall MA, Ranganathan R. Evolutionarily conserved networks of
 569 residues mediate allosteric communication in proteins. Nat Struct Biol 2003;10:59–69.
- 570 [13] Ribeiro AJ, Tyzack JD, Borkakoti N, Holliday GL, Thornton JM. A global analysis of 571 function and conservation of catalytic residues in enzymes. J Biol Chem 2020;295:314–24.
- 572 [14] Radivojac P, Clark WT, Oron TR, Schnoes AM, Wittkop T, Sokolov A, et al. A large-scale
- evaluation of computational protein function prediction. Nat Methods 2013;10:221–7.
- 574 [15] Wallace AC, Borkakoti N, Thornton JM. TESS: a geometric hashing algorithm for deriving
- 3D coordinate templates for searching structural databases. Application to enzyme active sites.
 Protein Sci 1997;6:2308–23.
- 577 [16] Mistry J, Bateman A, Finn RD. Predicting active site residue annotations in the Pfam 578 database. BMC Bioinformatics 2007;8:1–14.
- [17] Zhang T, Zhang H, Chen K, Shen S, Ruan J, Kurgan L. Accurate sequence-based prediction
 of catalytic residues. Bioinformatics 2008;24:2329–38.
- 581 [18] Gutteridge A, Bartlett GJ, Thornton JM. Using a neural network and spatial clustering to 582 predict the location of active sites in enzymes. J Mol Biol 2003;330:719–34.
- 582 predict the location of active sites in enzymes. J Mor Biol 2005,550.719–54.
- [19] Zhang T, Zhang H, Chen K, Shen S, Ruan J, Kurgan L. Accurate sequence-based prediction
 of catalytic residues. Bioinformatics 2008;24:2329–38.
- 585 [20] Petrova NV, Wu CH. Prediction of catalytic residues using Support Vector Machine with 586 selected protein sequence and structural properties. BMC Bioinformatics 2006;7:1–12.
- [21] Sun J, Wang J, Xiong D, Hu J, Liu R. CRHunter: integrating multifaceted information to
 predict catalytic residues in enzymes. Sci Rep 2016;6:34044.
- [22] Youn E, Peters B, Radivojac P, Mooney SD. Evaluation of features for catalytic residue
 prediction in novel folds. Protein Sci 2007;16:216–26.
- 591 [23] Song J, Li F, Takemoto K, Haffari G, Akutsu T, Chou K-C, et al. PREvaIL, an integrative
- 592 approach for inferring catalytic residues using sequence, structural, and network features in a 593 machine-learning framework. J Theor Biol 2018;443:125–37.

- 594 [24] Chea E, Livesay DR. How accurate and statistically robust are catalytic site predictions 595 based on closeness centrality? BMC Bioinformatics 2007;8:1–14.
- 596 [25] Pan T, Li C, Bi Y, Wang Z, Gasser RB, Purcell AW, et al. PFresGO: an attention 597 mechanism-based deep-learning approach for protein annotation by integrating gene ontology 598 inter-relationships. Bioinformatics 2023;39:btad094.
- [26] Chen K, Arnold FH. Engineering new catalytic activities in enzymes. Nat Catal 2020;3:203–
 13.
- [27] Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST
- and PSI-BLAST: a new generation of protein database search programs. Nucleic Acids Res1997;25:3389–402.
- [28] Potter SC, Luciani A, Eddy SR, Park Y, Lopez R, Finn RD. HMMER web server: 2018
 update. Nucleic Acids Res 2018;46:W200–W4.
- 606 [29] Ribeiro AJM, Holliday GL, Furnham N, Tyzack JD, Ferris K, Thornton JM. Mechanism
- and Catalytic Site Atlas (M-CSA): a database of enzyme reaction mechanisms and active sites.
 Nucleic Acids Res 2018;46:D618–D23.
- 609 [30] Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of protein 610 or nucleotide sequences. Bioinformatics 2006;22:1658–9.
- 611 [31] Bairoch A. The ENZYME database in 2000. Nucleic Acids Res 2000;28:304–5.
- [32] Velankar S, Dana JM, Jacobsen J, Van Ginkel G, Gane PJ, Luo J, et al. SIFTS: structure
- 613 integration with function, taxonomy and sequences resource. Nucleic Acids Res 2012;41:D483–
 614 D9.
- [33] Schomburg I, Chang A, Ebeling C, Gremse M, Heldt C, Huhn G, et al. BRENDA, the
- enzyme database: updates and major new developments. Nucleic Acids Res 2004;32:D431–D3.
- 617 [34] Bartlett GJ, Porter CT, Borkakoti N, Thornton JM. Analysis of catalytic residues in enzyme
- 618 active sites. J Mol Biol 2002;324:105–21.
- 619 [35] Yip KM, Fischer N, Paknia E, Chari A, Stark H. Atomic-resolution protein structure 620 determination by cryo-EM. Nature 2020;587:157–61.
- [36] Yu T, Cui H, Li JC, Luo Y, Jiang G, Zhao H. Enzyme function prediction using contrastive
 learning. Science 2023;379:1358–63.
- [37] Mirdita M, Von Den Driesch L, Galiez C, Martin MJ, Söding J, Steinegger M. Uniclust
- databases of clustered and deeply annotated protein sequences and alignments. Nucleic Acids
 Res 2017;45:D170–D6.
- [38] Elnaggar A, Heinzinger M, Dallago C, Rehawi G, Wang Y, Jones L, et al. Prottrans: Toward
- 627 understanding the language of life through self-supervised learning. IEEE Trans Pattern Anal
- 628 Mach Intell 2021;44:7112–27.
- [39] Matreyek KA, Starita LM, Stephany JJ, Martin B, Chiasson MA, Gray VE, et al. Multiplex
- assessment of protein variant abundance by massively parallel sequencing. Nat Genet2018;50:874–82.
- [40] Cagiada M, Johansson KE, Valanciute A, Nielsen SV, Hartmann-Petersen R, Yang JJ, et al.
- 633 Understanding the origins of loss of protein function by analyzing the effects of thousands of634 variants on activity and abundance. Mol Biol Evol 2021;38:3235–46.
- 635 [41] Shen X, Zhang S, Long J, Chen C, Wang M, Cui Z, et al. A Highly Sensitive Model Based
- on Graph Neural Networks for Enzyme Key Catalytic Residue Prediction. J Chem Inf Model
 2023;63:4277–90.
- 638 [42] Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, et al. Applying and
- 639 improving AlphaFold at CASP14. Proteins 2021;89:1711–21.

- [43] Elofsson A. Progress at protein structure prediction, as seen in CASP15. Curr Opin Struct
 Biol 2023;80:102594.
- [44] Simpkin AJ, Mesdaghi S, Sanchez Rodriguez F, Elliott L, Murphy DL, Kryshtafovych A,
 et al. Tertiary structure assessment at CASP15. Proteins 2023;91:1616–35.
- 644 [45] Kissinger CR, Parge HE, Knighton DR, Lewis CT, Pelletier LA, Tempczyk A, et al. Crystal
- 645 structures of human calcineurin and the human FKBP12-FK506-calcineurin complex. Nature
- 646 1995;378:641–4.
- [46] Gong W, O'Gara M, Blumenthal RM, Cheng X. Structure of pvu II DNA-(cytosine N4)
- 648 methyltransferase, an example of domain permutation and protein fold assignment. Nucleic 649 Acids Res 1997;25:2702–15.
- 650 [47] Pearl FM, Bennett C, Bray JE, Harrison AP, Martin N, Shepherd A, et al. The CATH
- database: an extended protein family resource for structural and functional genomics. Nucleic
 Acids Res 2003;31:452–5.
- 653 [48] Van der Maaten L, Hinton G. Visualizing data using t-SNE. J Mach Learn Res 654 2008;9:2579–2605.
- [49] Shah NH, Wang Q, Yan Q, Karandur D, Kadlecek TA, Fallahee IR, et al. An electrostatic
- selection mechanism controls sequential kinase signaling downstream of the T cell receptor.Elife 2016;5:e20105.
- [50] Bartlett GJ, Porter CT, Borkakoti N, Thornton JM. Analysis of catalytic residues in enzyme
 active sites. J Mol Biol 2002;324:105–21.
- 660 [51] Johnson JL, Yaron TM, Huntsman EM, Kerelsky A, Song J, Regev A, et al. An atlas of 661 substrate specificities for the human serine/threonine kinome. Nature 2023;613:759–66.
- [52] Koehler Leman J, Szczerbiak P, Renfrew PD, Gligorijevic V, Berenberg D, Vatanen T, et
- al. Sequence-structure-function relationships in the microbial protein universe. Nat Commun2023;14:2351.
- [53] Zhao S, Kumar R, Sakai A, Vetting MW, Wood BM, Brown S, et al. Discovery of new
 enzymes and metabolic pathways by using structure and genome context. Nature 2013;502:698–
 702.
- 668 [54] Hwang KY, Cho C-S, Kim SS, Sung H-C, Yu YG, Cho Y. Structure and mechanism of 669 glutamate racemase from Aquifex pyrophilus. Nat Struct Biol 1999;6:422–6.
- 670 [55] Liu L, Iwata K, Kita A, Kawarabayasi Y, Yohda M, Miki K. Crystal structure of aspartate
- racemase from Pyrococcus horikoshii OT3 and its implications for molecular mechanism ofPLP-independent racemization. J Mol Biol 2002;319:479–89.
- [56] Fauman EB, Yuvaniyama C, Schubert HL, Stuckey JA, Saper MA. The X-ray crystal
 structures of Yersinia tyrosine phosphatase with bound tungstate and nitrate: mechanistic
- 675 implications. J Biol Chem 1996;271:18780–8.
- [57] Groves MR, Yao Z-J, Roller PP, Burke TR, Barford D. Structural basis for inhibition of the
 protein tyrosine phosphatase 1B by phosphotyrosine peptide mimetics. Biochemistry
 1998;37:17773–83.
- [58] Furnari FB, Huang HS, Cavenee WK. The phosphoinositol phosphatase activity of PTEN
 mediates a serum-sensitive G1 growth arrest in glioma cells. Cancer Res 1998;58:5002–8.
- [59] Amorosi CJ, Chiasson MA, McDonald MG, Wong LH, Sitko KA, Boyle G, et al. Massively
- parallel characterization of CYP2C9 variant enzyme activity and abundance. Am J Hum Genet
- 683 2021;108:1735–51.

[60] Moriyama T, Nishii R, Perez-Andreu V, Yang W, Klussmann FA, Zhao X, et al. NUDT15
polymorphisms alter thiopurine metabolism and hematopoietic toxicity. Nat Genet 2016;48:367–
73.

- [61] Sirot D, Recule C, Chaibi E, Bret L, Croize J, Chanal-Claris C, et al. A complex mutant of
 TEM-1 beta-lactamase with mutations encountered in both IRT-4 and extended-spectrum TEM-
- 689 15, produced by an Escherichia coli clinical isolate. Antimicrob Agents Chemother 690 1997;41:1322–5.
- [62] Jack BR, Meyer AG, Echave J, Wilke CO. Functional sites induce long-range evolutionary
 constraints in enzymes. PLoS Biol 2016;14:e1002452.
- 693 [63] Cagiada M, Bottaro S, Lindemose S, Schenstrøm SM, Stein A, Hartmann-Petersen R, et al.
- Discovering functionally important sites in proteins. Nat Commun 2023;14:4175.
- [64] Kraut DA, Carroll KS, Herschlag D. Challenges in enzyme mechanism and energetics.
 Annu Rev Biochem 2003;72:517–71.
- 697
- 698

699 **Figure legends**

700 Figure 1 SCREEN – a predictor of catalytic residues in enzymes

- A. Data collection. We curate catalytic residue annotations from the M-CSA database, extract the corresponding enzyme structures from the RCSB PDB database, and retrieve corresponding enzyme function information using the SIFTS program [32]. **B.** Generation of inputs that include evolutionary profiles based on multiple sequence alignments (MSA), sequence embeddings that leverage a large-scale protein language model, and structural characteristics derived from the atomic structures. **C.** The architecture of SCREEN's predictive model.
- 708 Figure 2 Predictive performance of SCREEN

709 A. Comparison of the sequence-based CRpred and structure-based PREvaIL on the five test 710 datasets. B. The distribution of the Euclidean distances to the nearest catalytic residues for the 711 residues in enzymes from the M-CSA database [13]. C. The distribution of the Euclidean 712 distances to the nearest catalytic residue for the catalytic residues predicted by SCREEN. For 713 boxplots, the center line represents the median, top and bottom edges are the first and third 714 quartiles, respectively. **D.** Comparison of the structure-based SCREEN model (GCN encoder) 715 with a sequence-based baseline model (CNN encoder) on the five test datasets. E.-G. Three 716 examples, ranked approximately at 10% (E), 50% (F), and 90% (G) based on the Best-F1 scores 717 for catalytic residues predicted by SCREEN. We compare SCREEN with a sequence-based 718 baseline model (CNN encoder). The native catalytic residues are in green in a zoomed figure.

The correctly identified catalytic residues by SCREEN and baseline model (CNN encoder) are
marked in red, while misidentified residues are in gray.

721

Figure 3 Analysis of predictive performance when using putative enzyme structure and low similarity test proteins

724 A. Comparison of results produced by the SCREEN using the native structure, structure 725 predicted from sequence with AlphaFold, and sequence-based baseline predictor (CNN 726 encoder); the error bars represent standard deviation of the mean based on 10 independent runs. B. Comparison of results produced by the SCREEN using the native structure, AlphaFold-727 predicted structure, and sequence-based baseline predictor in the context of the quality of the 728 729 AlphaFold-predicted structure. The dashed horizontal lines represent the Best-F1 scores generated by CRHunter (blue line) and PREvaIL (yellow line). C. Examples of contact maps by 730 ground-truth of native enzyme structures (PDB) and AlphaFold-predicted enzyme structures, 731 732 and corresponding catalytic residues identified by SCREEN for enzymes 1aui-A and 1boo-A. D. 733 Evaluation for enzymes that share varying levels of sequence identify with the training proteins 734 on the five test datasets.

735

736 Figure 4 Analysis of putative catalytic residues generated by SCREEN

737 A. Comparison of the SCREEN models with its variant SCREEN_NoEC that does not 738 incorporate enzyme function through contrastive learning on the five test datasets. This figure 739 also shows results from CRpred and PREvaIL. B. The t-SNE based visualization of latent feature 740 spaces in the SCREEN model for residue characteristics, such as hydrophobicity (left), hydrogen bond types (medium) and charges (right). The "H-bond acceptor" denotes residues exclusively 741 742 serving as hydrogen bond acceptors without containing H-bond donor atoms. C. The t-SNE based visualization of latent feature spaces in the SCREEN model for different color-coded 743 744 enzyme classes. D.-F. Analysis of structural and biophysical characteristics, which include 745 biophysical properties of amino acids (D), B-factor (E) and solvent accessibility (F) for the 746 putative catalytic residues generated by SCREEN.

747

Figure 5 Diversity of enzyme structure and catalytic residues with the same catalytic mechanism

We examine enzymes that have the same catalytic mechanism: Isomerases acting on amino acids and derivatives with EC number 5.1.1 (**A**) and phosphoric monoester hydrolases with EC number 3.1.3 (**B**). We plot the TM-score as a measure of structural similarity as a heatmap, with larger numbers (more yellow) representing more similar structures. We also map the catalytic residues prediction by SCREEN onto the structures as well as enzyme reactions on the right.

755

756 Figure 6 Assessing the tolerance of catalytic residues to mutations

757 A. Architecture of the tree-structured model for characterizing mutated residues based on catalytic residue predictions by SCREEN. B. SCREEN identified catalytic residues of four 758 different enzymes: PTEN tumor suppressor (PDB ID: 1d5r, Chain A), human cytochrome P450 759 760 CYP2C9 (PDB ID: 10g5, Chain A), NUDT15 (PDB ID: 5lpg, Chain A), and Escherichia coli TEM1 beta-lactamase (PDB ID: 1btl, Chain A) (top). Residues within enzyme structures are 761 colored according to their predicted mutant class: blue corresponds to the Wild Type-like (WTL) 762 763 residues, while gray to the Functional Loss (FL) residues (bottom). C. Quality of distance-based classification of residues with different mutation classes, measured by accuracy and precision on 764 765 both the validation and entire datasets. **D.** Distribution of the Euclidean distances for residues of 766 different mutation classes.

767

768 Table 1 Comparison with current predictors of catalytic residues





Click here to access/download;Figure;Figure 2_Revised.pdf ±

D 0.85

_ون

2C



Ε

1le6: Chain A

	AUC	AUPR	F1
GCN	1.0	1.0	1.0
CNN	0.986	0.622	0.667

GCN Gly-28, His-46, Asp-91

CNN Phe-26, Gly-28, His-46, Asp-47, Tyr-65, Asp-91



0.8



G



4

F

1csn: Chain A					
	AUC	AUPR	F1		
GCN	0.949	0.662	0.750		
CNN	0.958	0.488	0.571		

GCN

Asp-131, Lys-133, Asp-135, Asn-136, Thr-181 CNN

Asp-131, Lys-133, Asp-135, Asn-136, Thr-181



	AUC	AUPR	F1
GCN	0.985	0.496	0.500
CNN	0.941	0.118	0.333

GCN Asp-84, Arg-239 CNN

Asp-84, Lys-203, Arg-239









Methods	Measurement	EF	EF fold	HA	NN	PC
	(%)	Superfamily dataset	dataset	superfamily dataset	dataset	dataset
Methods from Youn et al.[22]; Chea et al.[24]; Gutteridge et	Precision (Recall)	16.9 (53.9) ^a	17.1 (51.1) ^a	16.5 (29.3) ^b	56.0 (14.0) ^c	7.0 (90.0) ^d
al.[18]; Petrova and Wu [20]	F1	25.7	25.6	21.1	22.4	13.0
CRpred	Precision (Recall)	15.9 (52.1)	16.1 (48.0)	24.7 (49.7)	65.9 (18.0)	5.6 (84.5)
	F1	24.4	24.1	33	28.3	10.5
CRHunter	Precision (Recall)	21.5 (68.7)	21.0 (62.7)	33.2 (69.6)	76.4 (24.0)	7.6 (92.1)
	F1	32.8	31.5	45.0	36.5	14.0
PREvaIL	Precision (Recall)	17.0 (59.4)	17.0 (56.5)	17.0 (57.9)	58.9 (17.0)	17.0 (58.1)
	F1	26.4	26.1	26.3	26.4	26.3
AEGAN	Precision (Recall)	31.7 (85.7)	31.0 (83.7)	29.8 (85.7)	29.6 (83.9)	28.6 (86.6)
	F1	45.9	44.6	43.6	43.3	41.3
SCREEN (This study)	Precision (Recall)	61.9 (61.2)	61.0 (68.4)	69.3 (74.8)	68.5 (79.9)	67.6 (82.0)
	F1	61.5	64.5	72.0	73.8	74.1

1
 Table 1 Comparison with current predictors of catalytic residues

Note: ^a Model performance on EF superfamily and EF fold datasets by Youn et al. [22]; ^b Model 2

3

performance on HA family dataset by Chea et al. [24]; ^{*c*} Molde performance on NN dataset using the structure-based method without spatial clustering by Gutteridge et al. [18]; ^{*d*} Model 4

Performance on PC dataset by Petrova and Wu [20]. 5