

1 **Robust enzyme kinetics prediction through pairwise relative learning**

2 Xiongwen Li^{1,2}, Zhengkai Li³, Jiawei Zou^{1,2}, Wenjie Chen^{1,2}, Shujia Liu^{1,2}, Ke Wu^{1,2}, Jiahao
3 Luo^{1,2}, Yu Chen⁴, Feiran Li^{1,2,*}

4
5 1 Institute of Biopharmaceutical and Health Engineering, Tsinghua Shenzhen International
6 Graduate School, Tsinghua University, Shenzhen 518055, China

7 2 Key Laboratory for Industrial Biocatalysis, Ministry of Education, Institute of Biochemical
8 Engineering, Department of Chemical Engineering, Tsinghua University, Beijing 100084,
9 China

10 3 Department of Computer Science, Boston University, Boston, MA 02215, USA

11 4 State Key Laboratory of Quantitative Synthetic Biology, Shenzhen Institute of Synthetic
12 Biology, Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences,
13 Shenzhen 518055, China

14
15 * Corresponding author:

16 Feiran Li, Email: feiranli@sz.tsinghua.edu.cn

18 **Abstract**

19 Enzyme kinetics prediction is a fundamental challenge, and numerous computational methods
20 have been developed. However, most methods rely on absolute regression across heterogeneous
21 datasets, which are noisy and prone to regression-to-the-mean bias, limiting accurate
22 identification of highly active enzymes. Here we present DeltaKcat, a siamese neural network
23 framework that reformulates kinetics prediction as a pairwise relative learning problem. By
24 learning pairwise differences within consistent experimental contexts, DeltaKcat improves
25 accuracy, robustness, and generalization over state-of-the-art methods. This framework enables
26 DeltaKcat to (i) mitigate inter-study variability, (ii) generalize to unseen enzymes and substrates,
27 and (iii) prioritize highly active enzymes. We demonstrate these capabilities through systematic
28 benchmarking and external adenylate kinase dataset validation. We anticipate that this approach
29 will be broadly useful for enzyme discovery and engineering, and more generally for learning
30 from noisy biological data.

31

32 Introduction

33 Enzymes are the fundamental functional units of life, governing the rates and efficiencies of
34 biochemical reactions¹⁻⁵. Their activity is quantitatively described by enzyme kinetics, which
35 determines catalytic efficiency⁶, including the turnover number (k_{cat}), Michaelis constant (K_m),
36 and catalytic efficiency (k_{cat}/K_m). However, despite the more than 230 million protein sequences
37 in UniProt⁷, fewer than 0.1% are linked to experimentally measured kinetic parameters in
38 databases such as BRENDA⁸ and SABIO-RK⁹, highlighting a substantial gap between enzyme
39 sequence availability and quantitative characterization of enzyme function¹⁰. Consequently,
40 enzyme kinetics characterization increasingly relies on computational methods¹¹⁻¹⁴. Recent
41 machine learning approaches predict enzyme kinetic parameters by jointly modeling enzyme
42 sequences and substrates to predict the kinetic values¹⁵⁻¹⁹. These methods have evolved from
43 traditional machine learning to deep learning frameworks¹⁵, and more recently to large
44 pretrained protein and molecular models that further improve predictive performance^{18, 20-23}.

45 However, two fundamental challenges critically limit current approaches and must be explicitly
46 addressed. First, the enzyme-substrate space is extremely sparse and unevenly sampled, with
47 many enzyme classes represented by only a small number of experimentally characterized
48 instances, severely limiting generalization to unseen sequences and substrates^{20, 24-27}. Second,
49 kinetic measurements exhibit substantial cross-study variability, where reported values for the
50 same enzyme-substrate system can differ by orders of magnitude due to differences in
51 experimental conditions and protocols, introducing strong systematic noise into training data^{28,}
52 ²⁹. Together, these limitations make direct prediction of absolute kinetic parameters challenging
53 for real-world applications such as enzyme discovery and engineering³⁰. In particular, models
54 trained on heterogeneous datasets tend to capture dataset-specific biases rather than true
55 biophysical relationships, limiting their ability to extrapolate beyond curated benchmarks³¹.
56 Therefore, overcoming these limitations requires a fundamental reformulation of enzyme
57 kinetics prediction beyond absolute value estimation toward learning invariant functional
58 relationships.

59 To address these challenges, we propose DeltaKcat, which reformulates enzyme kinetics
60 prediction as a pairwise relative learning problem, where the model learns kinetic differences
61 between paired enzyme-substrate systems rather than predicting absolute values independently.
62 By comparing enzyme pairs within a consistent experimental context (e.g., within the same
63 study), this formulation effectively reduces cross-study systematic bias and the impact of
64 measurement inconsistency. Importantly, it shifts the learning objective from dataset-dependent
65 values to functional relationships between sequence and substrate changes. This pairwise
66 formulation also increases the effective size of training data, as each study yields multiple valid
67 comparisons, substantially increasing usable data without additional experiments. DeltaKcat is
68 trained on a high-quality, large-scale dataset integrated from BRENDA, SABIO-RK, and
69 literature-mined resources. We further introduce a reference-based reconstruction strategy to
70 recover absolute kinetic parameters from predicted relative differences, enabling benchmarking
71 and evaluation across multiple challenging settings, including low-sequence similarity, cold-
72 start scenarios, and cross-family generalization. Overall, this work establishes a general
73 framework for enzyme kinetics prediction based on pairwise relative functional learning and
74 provides a principled approach for robust prediction of noisy dataset.

75

76 Results

77 Overview of DeltaKcat

78 We first constructed a comprehensive dataset of enzyme kinetic parameters by integrating k_{cat} ,
79 K_m , and k_{cat}/K_m measurements from the BRENDA⁸ and SABIO-RK⁹ databases, and further
80 expanding coverage using large language model-based literature mined resources, including
81 EnzyExtraDB³² and Enzyme Co-Scientist³³. After removing entries lacking substrate simplified
82 molecular input line entry system (SMILES) and duplicate records, the curated dataset contains
83 135,072 k_{cat} , 208,176 K_m , and 110,624 k_{cat}/K_m measurements, spanning 35,076 unique protein
84 sequences and 13,705 distinct substrates (Supplementary Fig. 1). To reduce cross-study
85 variability, we reorganized the data into paired representations of enzyme-substrate complexes
86 within the same study. By fixing either the enzyme or the substrate within each pair, this design
87 enables the model to disentangle sequence-driven and substrate-driven effects under a
88 consistent experimental context. By separating these two pair types, we could explicitly
89 distinguish kinetic changes caused by enzyme sequence variation from those caused by
90 substrate variation. These capture (i) variation across enzyme sequences while shared substrates,
91 and (ii) variation across substrates with shared enzymes, yielding 117,873/92,187/51,748 and
92 142,909/160,331/98,669 paired samples for k_{cat} , K_m , and k_{cat}/K_m , respectively (Supplementary
93 Fig. 2).

94
95 We developed DeltaKcat, a siamese neural network model that predicts relative enzyme kinetic
96 differences between enzyme-substrate complexes (Fig. 1). The model takes as input a pair of
97 enzyme-substrate complexes and learns the difference in their kinetic properties. For each
98 complex, the protein sequence and substrate molecule are encoded using pretrained models
99 (ProtT5³⁴ for proteins, and MolT5³⁵ for substrates), complemented by molecular access system
100 (MACCS) fingerprints, and integrated via a gated cross-attention module to capture enzyme-
101 substrate interactions. The resulting representations from the two complexes are then combined
102 through explicit differencing, together with their individual features, and fed into a multilayer
103 perceptron (MLP) to predict relative k_{cat} , K_m , or k_{cat}/K_m . This formulation shifts the learning
104 objective from absolute values to invariant functional relationships, improving robustness to
105 data heterogeneity.

106

107 Prediction of relative enzyme kinetics

108 DeltaKcat shows strong predictive performance across relative kinetics prediction on both
109 sequence-diversity and substrate-specificity tasks, achieving high correlations across all three
110 relative kinetic prediction targets (Δk_{cat} , ΔK_m , and $\Delta k_{cat}/K_m$), with Pearson correlation
111 coefficients (PCC) reaching up to 0.80 (Fig. 2a). Considering the inherent noise and cross-study
112 variability in enzyme kinetic measurements, these correlation levels indicate that the model is
113 able to recover substantial functional signal from highly heterogeneous data. Notably, we
114 observe a clear task-dependent pattern in predictive performance. Δk_{cat} is most accurately
115 predicted in the sequence-diversity task, whereas ΔK_m is best predicted in the substrate-
116 specificity task (Fig. 2a), suggesting that the model is more sensitive to sequence-driven
117 variations in catalytic turnover and substrate-driven variations in binding affinity. To further
118 assess generalization beyond random-split settings, we further evaluated DeltaKcat under a
119 more challenging unseen-sequence or unseen-substrate scenario (cold-start). While

120 performance is lower under this setting, ΔK_m remains strongly predictable in the unseen-
121 substrate setting, achieving a PCC of 0.74, highlighting robust generalization to novel
122 substrates (Fig. 2b), indicating that the model retains robust predictive capability when
123 extrapolating to previously unobserved enzyme-substrate combinations.

124

125 We next examined performance across different sequence contexts. DeltaKcat achieves higher
126 accuracy when comparing mutants than when comparing distinct wild-type sequences
127 (Supplementary Fig. 3), suggesting that kinetic changes induced by local sequence
128 perturbations are more predictable than those arising from larger evolutionary divergence. The
129 model showed consistently similar performance across all enzyme commission (EC) classes,
130 indicating robust generalization across diverse functional categories (Supplementary Fig. 4).
131 Beyond regression accuracy, DeltaKcat also reliably captures the direction of kinetic changes
132 with an accuracy of 0.71, indicating that the model learns meaningful relative sequence-
133 function relationships (Fig. 2c for k_{cat} , Supplementary Fig. 5 for K_m and k_{cat}/K_m). All results are
134 consistent across three independent runs (Fig. 2a-c), demonstrating the stability of the model
135 across different data splits. Ablation studies showed that each component of DeltaKcat
136 contributes to performance, with module-specific effects observed across sequence-diversity
137 and substrate-specificity tasks (Supplementary Fig. 6).

138

139 **Absolute kinetic parameter prediction on the sequence-diversity task**

140 While relative kinetic parameters are central to enzyme engineering, absolute kinetic
141 parameters are also essential for systems-level applications. We therefore extended DeltaKcat
142 to absolute value prediction using a reference-based inference approach, in which DeltaKcat's
143 relative outputs are anchored with experimentally characterized values to obtain k_{cat} , K_m , and
144 k_{cat}/K_m values (Methods). Under the random split, DeltaKcat achieves R^2 values of 0.79 for k_{cat} ,
145 0.80 for K_m , and 0.84 for k_{cat}/K_m (Supplementary Fig. 7). Under the cold-start split, DeltaKcat
146 also achieves strong performance across all three tasks. For k_{cat} , it reaches an R^2 of 0.71 and a
147 PCC of 0.85, outperforming all baseline methods, including DLKcat, CataPro, and UniKP (R^2
148 = 0.14 to 0.48) (Fig. 3a, Supplementary Fig. 8). Comparable gains are observed for K_m and
149 k_{cat}/K_m , with consistent improvements in both R^2 and PCC (Fig. 3a for R^2 , Supplementary Fig.
150 8 for PCC). As cold-start evaluation better reflects realistic application scenarios, we therefore
151 focus on the cold-start results in the following analyses.

152

153 To assess performance under realistic experimental conditions, we evaluated DeltaKcat within
154 individual publications. In this publication-level setting, kinetic parameters are measured in the
155 same laboratory under comparable experimental protocols. Under this setting, DeltaKcat
156 achieves a median PCC of 0.78 for k_{cat} prediction, outperforming baseline models (PCCs of
157 0.11-0.36) (Fig. 3b). Across EC classes, the largest gains are observed in EC 1, EC 2, and EC
158 6, while consistently outperforming baseline models across the remaining classes (Fig. 3c).
159 Consistent trends are observed for K_m and k_{cat}/K_m prediction (Supplementary Fig. 9-10).

160

161 Importantly, DeltaKcat shows clear advantages in identifying highly active enzymes, where
162 existing methods are often biased toward the mean. It better captures extreme kinetic regimes
163 and consistently outperforms all baselines across top-ranked enzyme subsets (top 10-25%),

164 with the largest gains observed for the top 10% most active enzymes (Fig. 3d). Consistent
165 improvements are observed for K_m and k_{cat}/K_m (Supplementary Fig. 11). This improvement
166 arises because models trained via absolute regression are biased toward the central tendency of
167 the training distribution, leading to systematic underestimation of extreme values. In contrast,
168 modeling relative differences removes this dependency on absolute value distributions and
169 better preserves extreme kinetic signals.

170

171 We further evaluated generalization across sequence similarity levels. DeltaKcat maintains
172 strong performance even in low-similarity regimes, achieving PCC values of 0.76 and 0.78 in
173 the 0-40% and 40-60% identity groups for k_{cat} (Fig. 3e). Consistent performance is also
174 observed for K_m and k_{cat}/K_m , where DeltaKcat achieves PCC values above 0.85 and 0.73 in the
175 0-40% identity group (Supplementary Fig. 12). These results demonstrate robust extrapolation
176 capability across diverse sequence space of DeltaKcat.

177

178 Finally, we assessed the consistency of the relative-to-absolute conversion procedure. Predicted
179 absolute values show low variance across different anchor selections, with standard deviations
180 centered around 0.25 (Fig. 3f for k_{cat} , Supplementary Fig. 13 for K_m and k_{cat}/K_m). This indicates
181 that the inference procedure is stable and insensitive to the choice of reference enzyme-substrate
182 complex, supporting the reliability of the proposed conversion framework.

183

184 **Absolute value prediction on the substrate-specificity task**

185 We next evaluated DeltaKcat in the substrate-specificity setting by converting relative kinetic
186 predictions into absolute k_{cat} , K_m , and k_{cat}/K_m values. Despite substantially increased chemical
187 diversity, DeltaKcat consistently outperforms all baseline models. For k_{cat} , DeltaKcat achieves
188 an R^2 of 0.66 and a PCC of 0.83, exceeding CataPro ($R^2 = 0.42$) and DLKcat ($R^2 = 0.17$) (Fig.
189 4a). Similar improvements are observed for K_m ($R^2 = 0.47$, Fig. 4b) and k_{cat}/K_m ($R^2 = 0.53$ vs.
190 0.18 for the best baseline; Fig. 4c). Moreover, DeltaKcat also consistently outperforms baseline
191 models across publication-level analysis (Supplementary Fig. 14). To assess the generalization
192 of DeltaKcat to dissimilar substrates, we stratified test samples based on substrate similarity to
193 the training set. DeltaKcat maintains stable performance across all similarity ranges, including
194 the most challenging 0-80% interval, where it achieves PCCs of 0.80 (k_{cat}), 0.67 (K_m), and 0.77
195 (k_{cat}/K_m) (Fig. 4d-f). Overall, DeltaKcat generalizes well to substrate-driven kinetic variation
196 and uncharacterized substrates.

197

198 **Performance of DeltaKcat on an external adenylate kinase dataset**

199 To validate the real-world applicability of DeltaKcat, we tested it on an external adenylate
200 kinase (ADK) dataset that was completely excluded from training³⁶. ADKs are key enzymes in
201 energy metabolism that maintain intracellular adenylate homeostasis by catalyzing the
202 interconversion of ATP, ADP, and AMP (Fig. 5a)³⁷. This dataset contains 175 kinetic parameter
203 records derived from diverse ADK variants.

204

205 Under a zero-shot setting, we evaluated the model's utility for enzyme prioritization. Using the
206 top 5% of true k_{cat} values as a proxy for highly active enzymes, DeltaKcat successfully
207 identified one highly active enzyme within its top 5 candidates, whereas all baseline models

208 failed to retrieve any highly active enzyme within their top 10 candidates (Fig. 5 b, c). Across
209 the top 30 candidates, DeltaKcat recovers 56% of the true highly active enzymes, compared
210 with 22% for the best-performing baseline UniKP (Fig. 5b, c), highlighting its effectiveness for
211 practical enzyme discovery.

212

213 Under the zero-shot setting, DeltaKcat achieves a PCC of 0.40 and a Spearman correlation
214 coefficient (SCC) of 0.33 on this external ADK test set, representing the best performance
215 among all models evaluated, including UniKP (PCC = 0.22, SCC = 0.21) (Fig. 5d). Given the
216 complete absence of the ADK enzyme family from the training data, these findings indicate
217 that the paired comparative learning framework of DeltaKcat provides meaningful
218 generalization capabilities across unseen enzyme classes. Finally, we assessed whether limited
219 family-specific data can further improve performance. Using 20-140 training samples for fine-
220 tuning with a fixed held-out test set, performance increases steadily with data size, with a
221 substantial improvement beyond 60 samples. The fine-tuned model reached a PCC of 0.69,
222 corresponding to a 72% gain over the zero-shot setting (Fig. 5d, Supplementary Fig. 15), and
223 outperformed the previously reported ADK-specific model (ADK evolutionary). These results
224 demonstrate that DeltaKcat not only generalizes to completely unseen enzyme families, but can
225 also be rapidly adapted using small amounts of experimental data for improved predictive
226 accuracy and enzyme prioritization.

227

228 Discussion

229 Despite the rapid expansion of enzyme databases, the limited availability of high-quality kinetic
230 measurements remains a fundamental bottleneck for quantitative modeling and biocatalyst
231 design²⁴. Current approaches typically regress absolute k_{cat} , K_m and k_{cat}/K_m values from
232 heterogeneous experimental datasets^{20,21}, which introduces systematic variability and noise and
233 leads models to capture dataset-specific biases rather than generalizable enzyme-substrate
234 quantitative function relationships. In addition, standard random splits may overestimate
235 performance due to residual sequence or substrate similarity between training and test data.

236 To address these limitations, DeltaKcat reformulates enzyme kinetics prediction as a relational
237 learning problem. Instead of directly predicting absolute values, the model learns relative
238 kinetic differences between paired enzyme-substrate complexes under sequence or substrate
239 perturbations. This shifts learning from absolute scale fitting to invariant functional comparison,
240 improving robustness to experimental heterogeneity while focusing on perturbation-driven
241 signals.

242 Absolute kinetic parameters can then be reconstructed from relative predictions using a
243 reference-based inference strategy. Across diverse settings, including low sequence similarity,
244 low substrate similarity, and zero-shot validation on external ADK dataset, DeltaKcat
245 consistently outperforms state-of-the-art methods and maintains strong performance even on
246 unseen enzyme families, supporting its utility for enzyme discovery and functional annotation.
247 Despite these advantages, several limitations remain. First, although the dataset used in this
248 study is among the largest curated kinetic resources to date, the overall availability of
249 measurements remains limited, particularly for rare reaction types and underrepresented
250 enzyme classes. Second, key environmental factors such as pH, temperature, and assay
251 conditions, which are known to significantly influence catalytic activity, are not explicitly

252 incorporated into the current framework³⁸. Third, structural information of proteins and
253 enzyme-substrate complexes is not utilized, which may further improve model performance³⁹.
254 ⁴⁰.

255 More broadly, relational learning provides a general framework for biomolecular property
256 prediction under sparse, noisy, and heterogeneous data. This approach may facilitate enzyme
257 engineering, metabolic design, and functional annotation by enabling more reliable inference
258 in uncharacterized sequence and substrate spaces.

259

260 **Method**

261 **Dataset preparation**

262 Experimentally measured k_{cat} , K_m , and k_{cat}/K_m values were collected from the BRENDA and
263 SABIO-RK⁹ open-source databases in October 2024. Each data sample included the EC number,
264 protein sequence, organism information, substrate name, substrate SMILES, PubMed ID,
265 UniProt ID, temperature, pH, modification information, and enzyme kinetic parameter values.
266 The SMILES string of the substrate was queried from the PubChem Compound Database⁴¹
267 using the substrate name or international chemical identifiers (InChI). Enzyme sequences are
268 primarily retrieved from the UniProt database⁶ using UniProt ID. For entries lacking UniProt
269 IDs, sequences are obtained from the UniProt database based on the corresponding EC numbers
270 and species information. Mutant enzyme sequences were generated by applying the amino acid
271 substitutions recorded in the dataset to the corresponding wild-type sequences. We further
272 incorporated data from EnzyExtraDB and Enzyme Co-Scientist databases, two databases
273 constructed using large language models to extract kinetic information from the literature. The
274 same set of data fields was retained across all sources to ensure consistency. The final dataset
275 contains 135072 unique k_{cat} entries, 208176 K_m entries, and 110624 k_{cat}/K_m entries.

276

277 **Relative enzyme kinetic parameter datasets preparation**

278 Two types of relative enzyme kinetic parameter datasets were constructed based on specific
279 prediction tasks, the sequence-diversity task and the substrate-specificity task. In both datasets,
280 pairs of enzyme-substrate complexes were constructed under the constraint that samples
281 originated from the same publication, ensuring consistent experimental conditions within each
282 pair. For each pair, the relative kinetic parameter was calculated as the \log_2 -transformed ratio
283 between the kinetic value of the first complex and that of the second complex. For the sequence-
284 diversity dataset, we retained only records in which at least two distinct protein sequences were
285 associated with the same substrate within the same study. A cold-start split was then applied: a
286 subset of protein sequences was held out as test sequences. Training pairs were formed
287 exclusively from training sequences sharing the same substrate and study, while test pairs were
288 formed by pairing complexes involving test sequences with any other complex under the same
289 constraints. This design ensures that no test sequence appears during training. For the substrate-
290 specificity task, we retained only records in which at least two distinct substrates were
291 associated with the same protein sequence within the same study. A cold-start split was applied
292 by holding out a subset of substrates as test substrates. Test pairs were formed by pairing
293 complexes containing test substrates with other complexes from the same study, while training
294 pairs were formed solely from complexes involving training substrates. This prevents any test
295 substrate from appearing in the training set.

296

297 Evaluation metrics

298 In this work, the coefficient of determination (R^2), Pearson's correlation coefficient (PCC),
299 Spearman correlation coefficient (SCC), and root mean squared error (RMSE) was used to
300 assess prediction accuracy. Recall@K was used to assess top-K retrieval performance.

$$301 R^2 = 1 - \frac{\sum_{i=1}^n (m_i - z_i)^2}{\sum_{i=1}^n (m_i - \bar{m}_i)^2}$$

302

$$303 PCC = 1 - \frac{\sum_{i=1}^n (z_i - \bar{z}_i)(m_i - \bar{m}_i)}{\sqrt{\sum_{i=1}^n (z_i - \bar{z}_i)^2} \cdot \sqrt{\sum_{i=1}^n (m_i - \bar{m}_i)^2}}$$

$$304 SCC = 1 - \frac{6 \sum_{i=1}^n d_i^2}{n(n^2 - 1)}$$

$$305 RMSE = 1 - \sqrt{\frac{1}{n} \sum_{i=1}^n (m_i - z_i)^2}$$

$$306 Recall@K = \frac{\sum_{i=1}^K r_i}{R}$$

307 where m_i is the value of the experimentally measured enzyme kinetic parameter, z_i denotes the
308 value of the model-predicted enzyme kinetic parameter. d_i denotes the rank difference between
309 the predicted and experimental values for a given sample. \bar{m}_i denotes the mean value of the
310 experimentally measured enzyme kinetic parameter, \bar{z}_i denotes the mean value of the model-
311 predicted enzyme kinetic parameter, and n denotes the number of samples. K represents the
312 predetermined number of top ranked candidates being evaluated. r_i is an indicator variable that
313 equals 1 if the candidate at rank i is a true positive and 0 if it is a false positive. R represents the
314 total number of actual positive samples present in the entire dataset.

315

316 Absolute kinetic parameter inference from relative predictions

317 To extend DeltaKcat from relative prediction to absolute kinetic parameter estimation, we
318 applied a reference-based inference approach that converts predicted relative ratios into
319 absolute values. For each test sample, we first identified whether a valid reference complex was
320 available in the training set. Specifically, we retained only those comparative pairs in which
321 one enzyme-substrate complex was present in the training set and had an experimentally
322 measured kinetic parameter. This known complex was treated as the reference anchor, whereas
323 the other complex was regarded as the target whose absolute parameter was to be inferred.
324 DeltaKcat predicts the relative kinetic relationship between the target complex and the
325 reference complex. We then combined the predicted relative ratio with the experimentally
326 determined kinetic value of the reference anchor to infer the absolute kinetic parameter of the
327 target complex. When multiple valid training-set anchors were available for the same target
328 complex, we performed inference independently for each anchor and then averaged the
329 resulting absolute values to obtain the final prediction. Using this approach, DeltaKcat projects
330 comparative predictions into the absolute parameter space through experimentally reference
331 points.

332

333 External ADK dataset and evaluation

334 The external ADK dataset contains 175 kinetic parameter records collected from diverse ADK
335 variants. For zero-shot evaluation, we applied the relative prediction framework without using
336 any ADK data for training. Five ADK enzymes were randomly selected as reference enzymes
337 with known experimental kinetic parameters, and the remaining enzymes were treated as target
338 enzymes. Each target enzyme was paired with all five references to obtain relative predictions.
339 These predictions were then converted to absolute kinetic parameters using the corresponding
340 reference values, and the five inferred values were averaged to obtain the final prediction for
341 each target enzyme. For fine-tuning experiments, 30 ADK samples were randomly selected as
342 a fixed test set. From the remaining samples, training subsets of different sizes were constructed,
343 ranging from 20 to 140 samples. DeltaKcat was fine-tuned on each training subset and
344 evaluated on the fixed test set.

345

346 Encoding for enzyme sequence, structure and substrates

347 We used pre-trained language models to characterize proteins and substrates, aiming to obtain
348 more expressive and transferable encodings. Enzyme sequences were encoded as residue-level
349 embeddings from ProtT5-XL-UniRef, yielding an $L \times 1024$ tensor (L is the sequence length). To
350 control computation and memory, sequences longer than 1024 residues were truncated at
351 residue 1024. Substrates were represented with two features: a 768-dimensional continuous
352 vector from MolT5 applied to normalized SMILES, and a 167-bit MACCS key fingerprint. The
353 MACCS keys were generated using the RDKit package.

354

355 DeltaKcat model and architecture

356 DeltaKcat was developed to predict relative enzyme kinetic parameters between two enzyme-
357 substrate complexes, including relative k_{cat} , relative K_m , and relative k_{cat}/K_m . Instead of
358 predicting the absolute kinetic value of each complex independently, the model takes a pair of
359 enzyme-substrate complexes as input and directly learns the kinetic difference between them.
360 For each complex $i \in [1, 2]$, the protein sequence and substrate SMILES were encoded
361 separately using pretrained language models. The protein sequence was represented by ProtT5
362 as a residue-level embedding matrix.

$$363 P_i \in \mathbb{R}^{L_i \times d_p}$$

364 where L_i is the sequence length and $d_p=1024$. The substrate SMILES was represented by
365 MolT5 as a molecular embedding.

$$366 S_i \in \mathbb{R}^{d_s}$$

367 where $d_s=768$. To capture enzyme-substrate interactions, DeltaKcat used a gated cross-
368 attention module. The substrate representation was used to query the protein sequence
369 representation, enabling the model to identify enzyme features most relevant to the
370 corresponding substrate. A gating unit was added after the softmax operation in the attention
371 layer to introduce an additional nonlinear modulation step. This design improved the flexibility
372 of attention allocation and strengthened the representation of enzyme-substrate interactions.

373 The gated cross-attention module was applied independently to the two input complexes,
374 generating two interaction representations, H_1 and H_2 . A difference representation was then
375 constructed as

$$376 H_{\Delta} = H_1 - H_2$$

377 The final pair representation was obtained by concatenating the two complex-specific
378 interaction vectors and their difference

$$379 \quad H_{pair} = H_1 || H_2 || H_{\Delta}$$

380 where || denotes vector concatenation. This representation was fed into a multilayer perceptron

$$381 \quad y = MLP(H_{pair})$$

382 In this way, DeltaKcat learns relative kinetic changes directly from paired enzyme-substrate
383 complexes.

384

385 **Data availability**

386 All relevant data supporting the key findings of this study are available within the article and
387 its Supplementary Information files. All the data analyzed in this study is publicly available
388 from either public database, including BRENDA (<https://www.brenda-enzymes.org/>), SABIO-
389 RK (<https://sabiork.h-its.org/>), UniProt (<https://www.uniprot.org/>), PubChem
390 (<https://www.uniprot.org/>) databases or supplementary datasets of referenced articles
391 (<https://github.com/ChemBioHTP/EnzyExtract>, [https://github.com/JackKuo666/LLM-
392 BioDataExtractor](https://github.com/JackKuo666/LLM-BioDataExtractor)). All data used in the paper can be found in the GitHub repository:
393 <https://github.com/LiLabTsinghua/DeltaKcat>.

394

395 **Code availability**

396 In order to facilitate additional utilization, we have made available all of the codes and thorough
397 instructions in our GitHub repository located at <https://github.com/LiLabTsinghua/DeltaKcat>.

398

399 **Acknowledgements**

400 We acknowledge financial support from National Key R&D Program of China
401 (2024YFA0920300), Shenzhen Medical Research Fund (A2403013) and the National Natural
402 Science Foundation of China (22478223 and 62532017).

403

404 **Author Contributions**

405 F.L., Y.C. and X.L. and Z.L. designed the research. X.L. and Z.L. performed the research. X.
406 L., Z. L., J. Z., S. L., W. C., K. W., J. L., Y. C., and F. L. analyzed the data. F.L. and X.L. wrote
407 the paper. All authors approved the final paper.

408

409 **Competing Interests Statement**

410 The authors declare no competing interests.

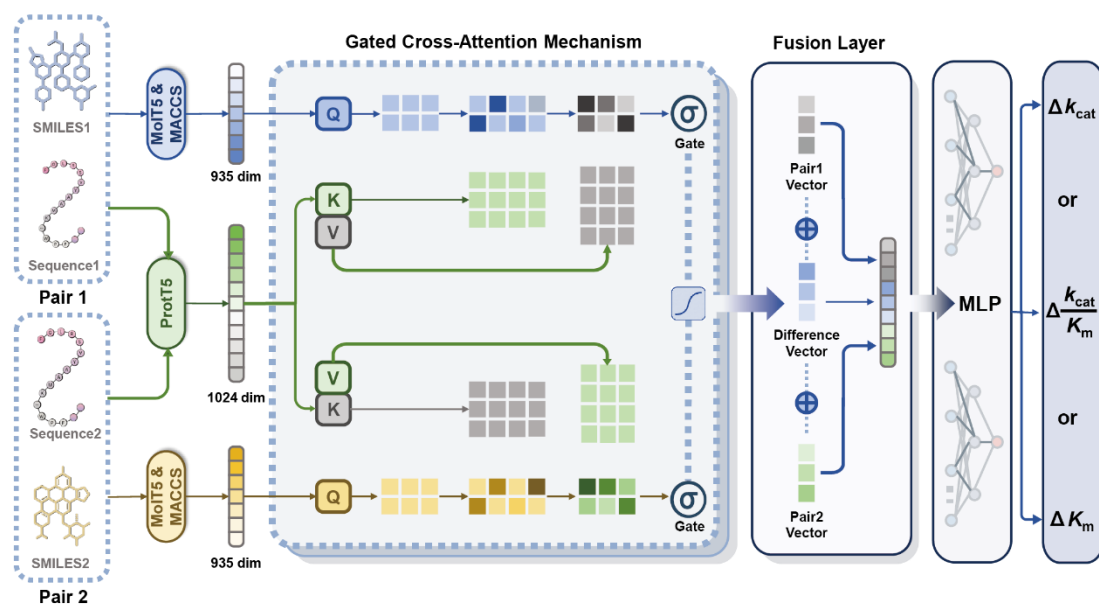
411

412

413 **Reference:**

- 414 1. Currin, A., Swainston, N., Day, P.J. & Kell, D.B. Synthetic biology for the directed evolution of
415 protein biocatalysts: navigating sequence space intelligently. *Chem Soc Rev* **44**, 1172-1239
416 (2015).
- 417 2. Riva, S. Laccases: blue enzymes for green chemistry. *Trends Biotechnol* **24**, 219-226 (2006).
- 418 3. Chen, Y. & Nielsen, J. Energy metabolism controls phenotypes by protein efficiency and
419 allocation. *Proc Natl Acad Sci U S A* **116**, 17592-17597 (2019).
- 420 4. Meghwanshi, G.K. et al. Enzymes for pharmaceutical and therapeutic applications. *Biotechnol*
421 *Appl Biochem* **67**, 586-601 (2020).
- 422 5. Raveendran, S. et al. Applications of Microbial Enzymes in Food Industry. *Food Technol*
423 *Biotechnol* **56**, 16-30 (2018).
- 424 6. Goldman, S., Das, R., Yang, K.K. & Coley, C.W. Machine learning modeling of family wide
425 enzyme-substrate specificity screens. *PLoS Comput Biol* **18**, e1009853 (2022).
- 426 7. Bateman, A. et al. UniProt: the Universal Protein Knowledgebase in 2025. *Nucleic Acids*
427 *Research* **52**, D609-D617 (2024).
- 428 8. Chang, A. et al. BRENDA, the ELIXIR core data resource in 2021: new developments and updates.
429 *Nucleic Acids Research* **49**, D498-D508 (2021).
- 430 9. Wittig, U., Rey, M., Weidemann, A., Kania, R. & Müller, W. SABIO-RK: an updated resource for
431 manually curated biochemical reaction kinetics. *Nucleic Acids Research* **46**, D656-D660 (2018).
- 432 10. Nilsson, A., Nielsen, J. & Palsson, B.O. Metabolic Models of Protein Allocation Call for the
433 Kinetome. *Cell Syst* **5**, 538-541 (2017).
- 434 11. Zhang, D. et al. Discovery of Toxin-Degrading Enzymes with Positive Unlabeled Deep Learning.
435 *ACS Catalysis* **14**, 3336-3348 (2024).
- 436 12. Cui, H. et al. Enzyme specificity prediction using cross-attention graph neural networks. *Nature*
437 **647**, 639-647 (2025).
- 438 13. Meier, J. et al. Language models enable zero-shot prediction of the effects of mutations on
439 protein function. *Advances in neural information processing systems* **34**, 29287-29303 (2021).
- 440 14. Lin, Z. et al. Evolutionary-scale prediction of atomic-level protein structure with a language
441 model. *Science* **379**, 1123-1130 (2023).
- 442 15. Li, F. et al. Deep learning-based k cat prediction enables improved enzyme-constrained model
443 reconstruction. *Nature Catalysis* **5**, 662-672 (2022).
- 444 16. Wang, J. et al. MPEK: a multitask deep learning framework based on pretrained language
445 models for enzymatic reaction kinetic parameters prediction. *Brief Bioinform* **25** (2024).
- 446 17. Hua, C. et al. Reactzyme: A benchmark for enzyme-reaction prediction. *Advances in Neural*
447 *Information Processing Systems* **37**, 26415-26442 (2024).
- 448 18. Kroll, A., Rousset, Y., Hu, X.P., Liebrand, N.A. & Lercher, M.J. Turnover number predictions for
449 kinetically uncharacterized enzymes using machine and deep learning. *Nat Commun* **14**, 4139
450 (2023).
- 451 19. Shen, X. et al. EITLEM-Kinetics: A deep-learning framework for kinetic parameter prediction of
452 mutant enzymes. *Chem Catalysis* (2024).
- 453 20. Wang, Z. et al. Robust enzyme discovery and engineering with deep learning using CataPro.
454 *Nat Commun* **16**, 2736 (2025).
- 455 21. Yu, H., Deng, H., He, J., Keasling, J.D. & Luo, X. UniKP: a unified framework for the prediction of
456 enzyme kinetic parameters. *Nat Commun* **14**, 8211 (2023).

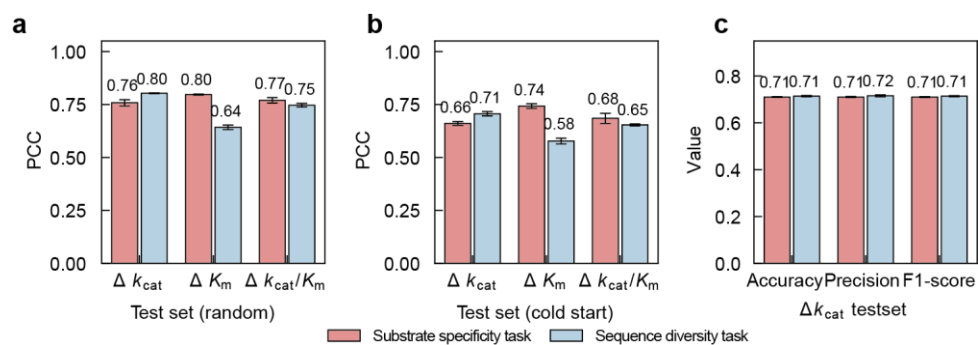
- 457 22. Rives, A. et al. Biological structure and function emerge from scaling unsupervised learning to
458 250 million protein sequences. *Proc Natl Acad Sci U S A* **118** (2021).
- 459 23. Honda, S., Shi, S. & Ueda, H.R. Smiles transformer: Pre-trained molecular fingerprint for low
460 data drug discovery. *arXiv preprint arXiv:1911.04738* (2019).
- 461 24. Sapoval, N. et al. Current progress and open challenges for applying deep learning across the
462 biosciences. *Nat Commun* **13**, 1728 (2022).
- 463 25. Lyu, B. et al. GotEnzymes2: expanding coverage of enzyme kinetics and thermal properties.
464 *Nucleic Acids Res* **54**, D583-D592 (2026).
- 465 26. Qiu, S., Saeed, H., Leonard, W., Li, F. & Yang, A. Machine learning for enzyme catalytic activity:
466 current progress and future horizons. *Brief Bioinform* **27** (2026).
- 467 27. Sagi, O. & Rokach, L. Ensemble learning: A survey. *Wiley interdisciplinary reviews: data mining
468 and knowledge discovery* **8**, e1249 (2018).
- 469 28. Jimenez-Luna, J. et al. DeltaDelta neural networks for lead optimization of small molecule
470 potency. *Chem Sci* **10**, 10911-10918 (2019).
- 471 29. Yu, J. et al. Computing the relative binding affinity of ligands based on a pairwise binding
472 comparison network. *Nat Comput Sci* **3**, 860-872 (2023).
- 473 30. Sajeevan, K.A. et al. Robust prediction of enzyme variant kinetics with RealKcat. *bioRxiv* (2025).
- 474 31. Malli, A., Vasyutyn, D. & Kim, J.R. Advances in Machine Learning Models for Predicting Enzyme
475 Kinetic Parameters. *J Chem Inf Model* **66**, 42-60 (2026).
- 476 32. Wei, G., Ran, X., Ai-Abssi, R. & Yang, Z. Finding the dark matter: Large language model-based
477 enzyme kinetic data extractor and its validation. *Protein Sci* **34**, e70251 (2025).
- 478 33. Jiang, J. et al. Enzyme Co-scientist: harnessing large language models for enzyme kinetic data
479 extraction from literature. *BioRxiv*, 2025.2003. 2003.641178 (2025).
- 480 34. Elnaggar, A. et al. ProtTrans: Toward Understanding the Language of Life Through Self-
481 Supervised Learning. *IEEE Trans Pattern Anal Mach Intell* **44**, 7112-7127 (2022).
- 482 35. Edwards, C. et al. in Proceedings of the 2022 Conference on Empirical Methods in Natural
483 Language Processing 375-413 (2022).
- 484 36. Muir, D.F. et al. Evolutionary-scale enzymology enables exploration of a rugged catalytic
485 landscape. *Science* **388**, eadu1058 (2025).
- 486 37. Kerns, S.J. et al. The energy landscape of adenylate kinase during catalysis. *Nat Struct Mol Biol*
487 **22**, 124-131 (2015).
- 488 38. Bar-Even, A. et al. The moderately efficient enzyme: evolutionary and physicochemical trends
489 shaping enzyme parameters. *Biochemistry* **50**, 4402-4410 (2011).
- 490 39. Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583-
491 589 (2021).
- 492 40. Song, Y. et al. Accurately predicting enzyme functions through geometric graph learning on
493 ESMFold-predicted structures. *Nat Commun* **15**, 8180 (2024).
- 494 41. Kim, S. et al. PubChem 2025 update. *Nucleic Acids Res* **53**, D1516-D1525 (2025).
- 495
- 496



497

498 **Fig. 1 DeltaKcat framework for predicting relative enzyme kinetic parameters.** DeltaKcat
499 takes a pair of enzyme-substrate complexes as input and predicts their relative k_{cat} , K_m , or $k_{cat} /$
500 K_m values as the output, respectively. Protein sequences are encoded by ProtT5, and substrates
501 are represented using MolT5 and MACCS features. Sequence and substrate embeddings are
502 integrated through a gated cross-attention module, in which substrate-derived queries attend to
503 protein-derived keys and values to capture enzyme-substrate interactions. A gating unit is
504 applied after the softmax step to refine attention allocation. The resulting interaction vectors for
505 the two complexes, together with their difference vector, are concatenated and passed to a
506 multilayer perceptron for prediction.

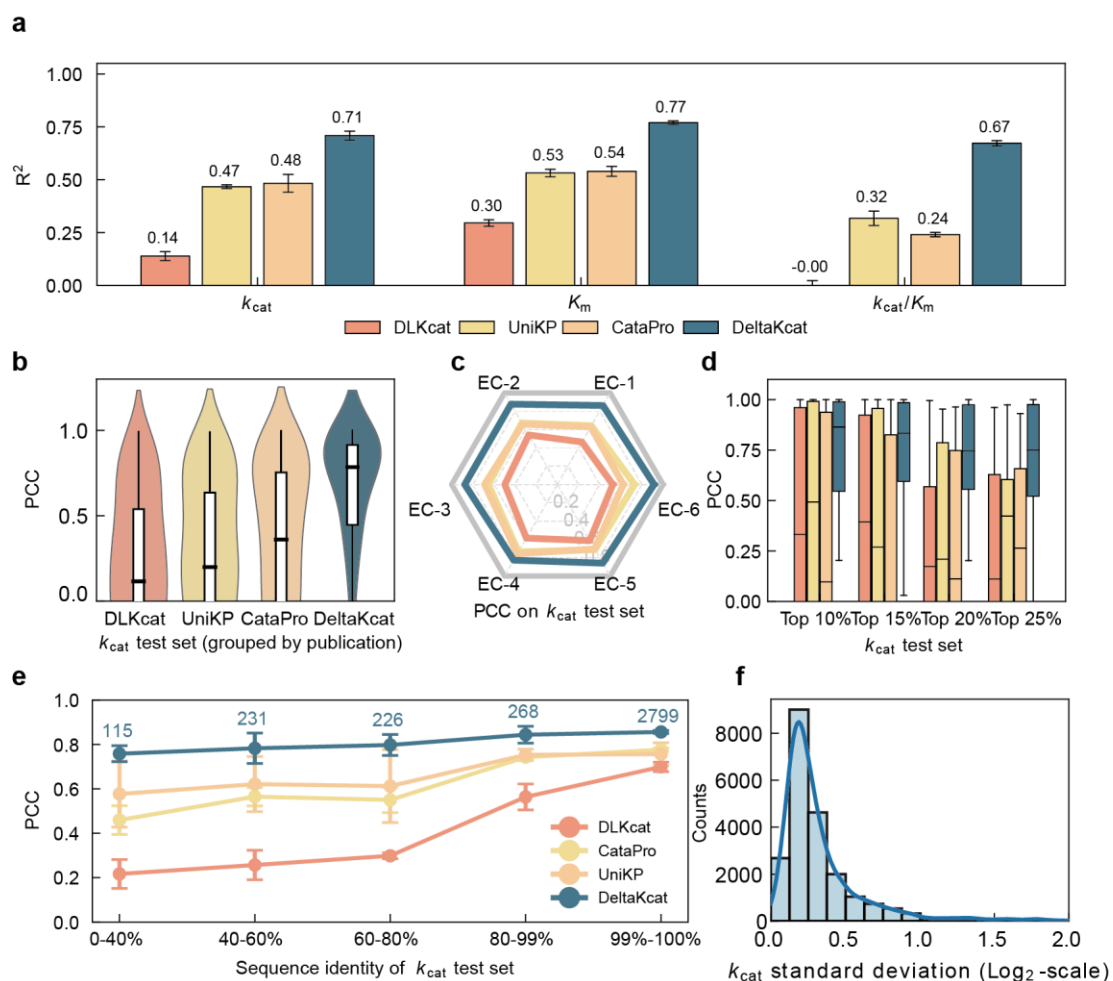
507



508

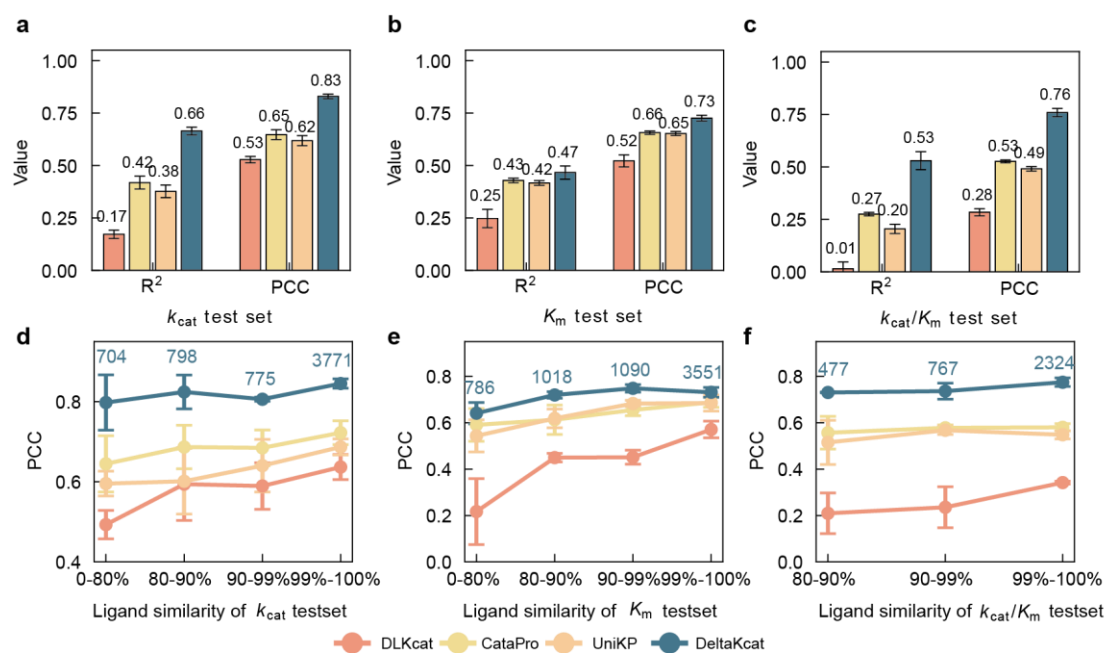
509 **Fig. 2 Prediction of relative enzyme kinetic parameters by DeltaKcat.** (a, b) Predictive
510 performance of DeltaKcat on the sequence-diversity dataset and substrate-specificity dataset
511 under random-split (a) and cold-start (b) settings for relative k_{cat} , relative K_m , and relative k_{cat}/K_m .
512 In the cold-start split, each test pair contained at least one protein sequence not observed during
513 training. (c) Classification performance of DeltaKcat for predicting the direction of kinetic
514 changes in the sequence-diversity and substrate-specificity tasks. The data in a-c are presented
515 as mean values, and the error bars represent the standard deviation (s.d.) across three
516 independent experiments using different random seeds.

517



518

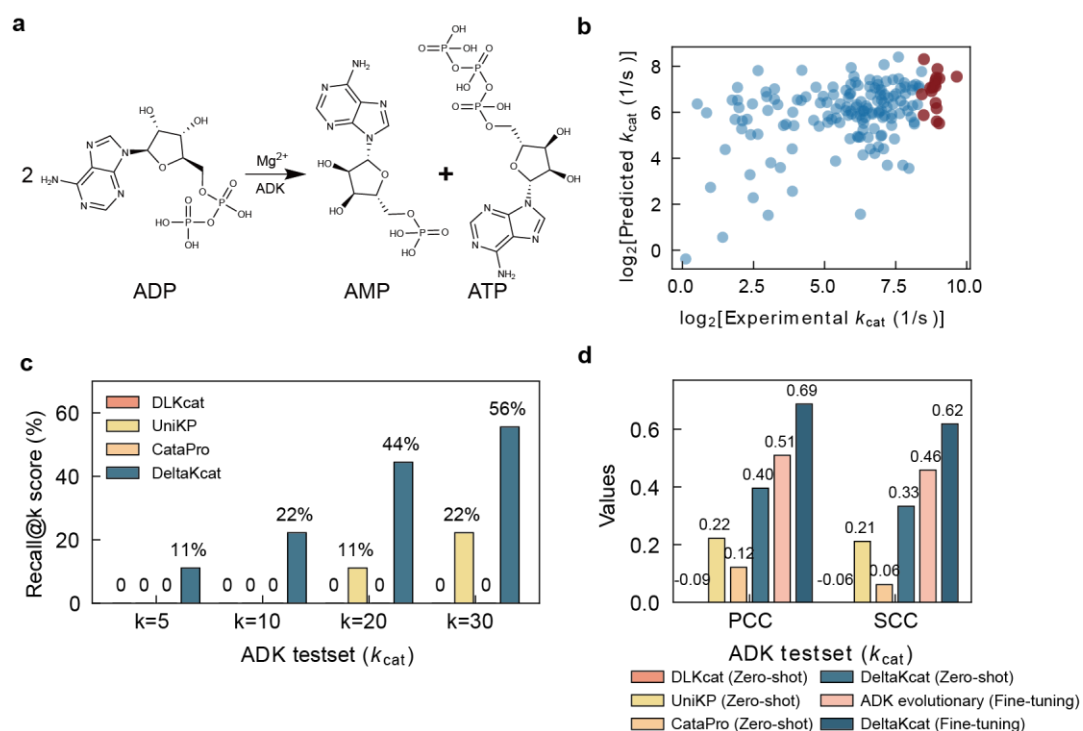
519 **Fig. 3 Absolute kinetic parameter prediction on the sequence-diversity dataset.** (a)
 520 Comparison of R^2 values for absolute parameter prediction. Performance of DeltaKcat versus
 521 DLKcat, UniKP, and CataPro baseline models on inferring absolute k_{cat} , K_m , and k_{cat}/K_m . (b)
 522 Publication-level Pearson correlation coefficients for absolute prediction across studies with
 523 more than five data points. (c) PCC results of different model based on EC classifications 1-6.
 524 (d) Evaluation of k_{cat} prediction (PCC) on subsets of the test set representing the top 10%, 15%,
 525 20%, and 25% of absolute k_{cat} values within each EC class. (e) PCC evaluation across varying
 526 thresholds of maximum sequence similarity to the training set. (f) Histogram showing the
 527 distribution of the standard deviation of multiple absolute k_{cat} estimates for the same test
 528 complex, derived from different training-set reference anchors. The solid blue line is a kernel
 529 density estimate. The data in a, e and f are presented as mean values, and the error bars
 530 represent the standard deviation (s.d.) across three independent experiments using different
 531 random seeds. In each box plot, the lower and upper boundaries of the box represent the first
 532 quartile (Q1) and the third quartile (Q3), respectively. The whiskers extend from the quartiles
 533 to the minimum and maximum values within 1.5 times the interquartile range. The black line
 534 inside the box represents the median.
 535



536

537 **Fig. 4 Absolute kinetic parameter prediction on the substrate-specificity dataset.** (a-c)
 538 Comparison of R^2 and PCC values for absolute parameter prediction. Performance of DeltaKcat
 539 versus DLKcat, UniKP, and CataPro baseline models on inferring absolute k_{cat} , K_m , and k_{cat}/K_m .
 540 (d-f) PCC evaluation across varying thresholds of maximum ligand similarity to the training
 541 set on k_{cat} , K_m , and k_{cat}/K_m . The data in a-f are presented as mean values, and the error bars
 542 represent the standard deviation (s.d.) across three independent experiments using different
 543 random seeds.

544



545

546 **Fig. 5 Evaluation of DeltaKcat on an external adenylate kinase (ADK) dataset.** (a)

547 Adenylate kinase catalyzes the reaction. (b) Scatter plot of predicted versus observed absolute

548 k_{cat} values for the external ADK dataset. Red points denote enzymes with observed k_{cat} values

549 in the top 5% of the dataset. (c) Recall of highly active enzymes among top-ranked candidates.

550 Highly active enzymes were defined as the top 5% by true k_{cat} . (d) Comparison of prediction

551 performance between zero-shot and fine-tuned models on the ADK dataset, evaluated by

552 Pearson and Spearman correlation coefficients.

553

554