1    **Deep learning based $k_{cat}$ prediction enables improved enzyme constrained model**

2    **reconstruction**

3

4    Feiran Li[1, #], Le Yuan[1, 2, #], Hongzhong Lu[1], Gang Li[1], Yu Chen[1], Martin K. M. Engqvist[1], Eduard

5    J Kerkhoven[1, 2], Jens Nielsen[1, 3, *]

6

7    1 Department of Biology and Biological Engineering, Chalmers University of Technology,

8    Kemivägen 10, SE-412 96 Gothenburg, Sweden

9    2 Novo Nordisk Foundation Center for Biosustainability, Chalmers University of Technology

10   , Kemivägen 10, SE-412 96, Gothenburg, Sweden

11   3 BioInnovation Institute, Ole Måløes Vej 3, DK2200 Copenhagen N, Denmark

12

13   # These authors contributed equally to this work: Feiran Li, Le Yuan.

14   * Correspondence to: nielsenj@chalmers.se

15

16

17

18

19

20

21

22

23

**Abstract**

Enzyme turnover numbers ($k_{cat}$ values) are key parameters to understand cell metabolism, proteome allocation and physiological diversity, but experimentally measured $k_{cat}$ data are sparse and noisy. Here we provide a deep learning approach to predict $k_{cat}$ values for metabolic enzymes in a high-throughput manner with the input of substrate structures and protein sequences. Our approach can capture $k_{cat}$ changes for mutated enzymes and identify amino acid residues with great impact on $k_{cat}$ values. Furthermore, we applied the approach to predict genome scale $k_{cat}$ values for over 300 yeast species, demonstrating that the predicted $k_{cat}$ values are consistent with current evolutional understanding. Additionally, we designed an automatic pipeline using the predicted $k_{cat}$ values to parameterize enzyme-constrained genome scale metabolic models (ecGEMs) facilitated by a Bayesian approach, which outperformed the default ecGEMs in predicting phenotypes and proteomes and enabled to explain phenotype differences among yeast species. The deep learning $k_{cat}$ prediction approach and automatic ecGEM construction pipeline would thus be a valuable tool to uncover the global trend of enzyme kinetics and physiological diversity, and to further elucidate cell metabolism on a large scale.

**Key words**: genome scale metabolic modelling, enzyme constraints, turnover rates, $k_{cat}$ values, deep learning, Bayesian approach.

**Introduction**

Enzyme turnover number ($k_{cat}$), which defines the maximum chemical conversion rate of a reaction, is a critical parameter for understanding metabolism, proteome allocation, growth and physiology of a certain organism[1–3]. There are large collections of $k_{cat}$ values available in the enzyme databases

2

47    BRENDA[4] and SABIO-RK[5], which are, however, still scarce compared to the variety of existing

48    organisms and metabolic enzymes, largely due to the lack of high-throughput methods for $k_{cat}$

49    measurements. Additionally, the experimentally measured $k_{cat}$ values have considerable

50    variabilities due to varying assay conditions such as pH, cofactor availability and experimental

51    methods[6]. Altogether, the sparse collection and considerable noise limit the usage of $k_{cat}$ data for

52    global analysis and may mask the enzyme evolution trend.

53

54    In particular, enzyme-constrained genome scale metabolic models (ecGEMs), where the whole-

55    cell metabolic network is constrained by enzyme catalytic capacities and thereby able to accurately

56    simulate maximum growth ability, metabolic shifts and proteome allocations, rely heavily on

57    genome scale $k_{cat}$ values[2,7]. Even for well-studied organisms, the $k_{cat}$ coverage is far less than

58    complete[8–10]. When data are missing, ecGEMs usually use assumed $k_{cat}$ values from similar

59    reactions or adopt available $k_{cat}$ values from other organisms, which could cause model predictions

60    deviating from experimental observations[7]. Thus, there is a clear requirement for obtaining a large

61    scale of $k_{cat}$ values to improve the model accuracy and get more reliable simulations for delicate

62    phenotypes[11].

63

64    Previously, machine learning has been used to predict $k_{cat}$ values based on features such as average

65    metabolic flux and the catalytic sites obtained from protein structures[9]. Due to the requirement of

66    feature data and absolute proteome data in the training dataset, this approach was only applied to

67    the most well-studied bacterium *Escherichia coli*, thus limiting its usage for large scale prediction

68    of $k_{cat}$ values for multiple organisms. In contrast, deep learning does not rely on feature selection

69    and has been applied and shown great performance in modeling chemical space[12], gene

3

70   expression[13], enzyme related parameters such as enzyme affinity[14], and enzyme commission

71   numbers (EC numbers)[15].

72

73   Inspired by these efforts, we developed a deep learning model and demonstrated its capability for

74   large scale prediction of $k_{cat}$ values, as well as for identifying key amino acid residues that affect

75   these predictions. We showcased the predictive power of the deep learning model by predicting

76   genome scale $k_{cat}$ profiles for 343 yeast/fungi species, accounting for more than 300,000 enzymes

77   and 3,000 substrates. The predicted $k_{cat}$ profiles enabled reconstruction of 343 ecGEMs for the

78   yeast/fungi species through an automatic Bayesian based pipeline, which can accurately simulate

79   growth phenotype among yeast species and identify the phenotype related key enzymes.

80

81   **Results**

82   **Construction of a deep learning framework for $k_{cat}$ prediction**

83   A deep learning framework was developed by combining a graph neural network (GNN)

84   for substrates and a convolutional neural network (CNN) for proteins (Fig. 1). In this framework,

85   substrates were represented as molecular graphs converted from SMILES (the simplified

86   molecular-input line-entry system) and protein sequences were split into overlapping n-gram

87   amino acids. To train the neural network, we generated a comprehensive dataset from the

88   BRENDA[4] and the SABIO-RK database[5]. Several rounds of data preprocessing and cleaning were

89   performed to filter out incomplete entries with missing information and redundant entries across

90   databases, to ensure that the dataset contains unique entries with substrate name, substrate SMILES,

91   EC number, protein sequence, organism name and $k_{cat}$ value information. The final dataset

92   contained 16,838 unique entries catalyzed by 7,822 unique protein sequences from 851 organisms

4

93    and converting 2,672 unique substrates (Supplementary Figure 1-2). This dataset was

94    randomly split into training, validation and test dataset by 80%, 10%, and 10%, respectively.

95

96    **Deep learning model performance for $k_{cat}$ prediction**

97    We first evaluated the effects of different model hyperparameters on deep learning performance

98    using learning curves (Supplementary Figure 3). Note that 2-radius subgraphs and 3-gram amino

99    acids used to extract the substrate and protein vectors can considerably improve the deep learning

100    performance compared with other tested hyperparameter settings (Supplementary Figure

101    3a). When investigating the effect of vector dimensionality, we found that more highly

102    dimensional vectors used for substrates and proteins led to somewhat better performance

103    (Supplementary Figure 3b). Then, Additionally, the model performed much better when the

104    number of time steps/layers in GNN/CNN is 2 or 3 (Supplementary Figure 3c). With the settled

105    parameters (r-radius is 2, n-gram is 3, vector dimensionality is 20, number of time steps in GNN

106    is 3, and number of layers in CNN is 3), the training dataset was used to train the deep learning

107    model. We observed that the Root Mean Square Error (RMSE) of $k_{cat}$ prediction in the validation

108    and test datasets gradually decreased with increasing epoch (Fig. 2a), where the number of epochs

109    represents iterations of the dataset passing through the neural network. A final deep learning model

110    was trained and stored for further use, when the RMSE was 0.99 and 1.06 for the validation and

111    test datasets, respectively, signifying that the predicted and measured $k_{cat}$ values were overall

112    within one order of magnitude (Fig. 2a). As a result, the deep learning model showed a

113    high predictive accuracy on the original whole dataset and test dataset (Fig. 2b for whole dataset,

114    Pearson's r = 0.88; Supplementary Figure 4a for test dataset, Pearson's r = 0.71; Supplementary

115    Figure 4b for test dataset with substrates and enzymes that were not present in the training dataset,

5

116    Pearson's r = 0.70). To facilitate the further usage of our deep learning prediction tool, we also

117    supplied a user-friendly example for $k_{cat}$ prediction in our GitHub repository with the input of

118    substrate                        and                        protein                        sequence

119    (https://github.com/SysBioChalmers/DLKcat/tree/master/DeeplearningApproach/Code/example).

120

121    Besides, we investigated whether the deep learning model can identify the preferred substrates for

122    promiscuous enzymes. We classified substrates with the highest $k_{cat}$ value for promiscuous

123    enzymes as preferred substrates, and substrates with the lowest one as the alternative substrates,

124    then through comparing the predicted $k_{cat}$ values for preferred substrates and alternative substrates

125    (Fig. 2c), we found that our deep learning model are able to predict that the enzymes do indeed

126    have a higher $k_{cat}$ for the preferred substrates (median value = 6.45 /s) compared with alternative

127    substrates (median value = 1.49 /s) (P value < 1e-10, for promiscuous enzymes in all dataset),

128    which validates the predictive power of our deep learning model in identifying the preferred

129    substrates. The same trend was identified using the prediction for promiscuous enzymes in our test

130    dataset (Supplementary Figure 4c, P value = 0.009).

131

132    To explore the metabolic contexts for all wildtype enzymes in the original dataset, we mapped

133    these enzymes to four modules on the basis of categorization in KEGG database[16]: primary-CE

134    (enzymes involved in carbohydrate and energy metabolism), primary-AFN (amino acid, fatty acids

135    and nucleotide metabolism), intermediate (metabolism of common biomass components such as

136    cofactors) and secondary metabolism (condition specific metabolism or metabolism related to low

137    concentration metabolites) (Supplementary Table 1). Enzymes associated with primary-CE

138    metabolism on average exhibited a higher predicted $k_{cat}$ value than those of primary-AFN,

6

139    secondary and intermediate metabolism (Fig. 2d), which is in accordance with the previous finding

140    that enzyme-substrate pairs from central carbon metabolism tend to have relatively higher

141    $k_{cat}$ values than secondary and intermediate metabolism[6].

142

143    **Prediction and interpretation of $k_{cat}$ of mutated enzymes**

144    While the deep learning model displays an overall good performance for predicting $k_{cat}$ values (Fig.

145    2b), we next explored whether the model could capture more details such as the effects of amino

146    acid substitutions on $k_{cat}$ values of individual enzymes. To this end, we divided the original

147    annotated dataset into two categories: one including wildtype enzymes and the other mutated

148    enzymes with amino acid substitutions. In these two splits the median $k_{cat}$ value of mutant enzymes

149    is lower than that for wildtype enzymes (Supplementary Figure 5a). We found that the deep

150    learning model is a good predictor of $k_{cat}$ values for both wildtype enzymes (Fig. 3a for the whole

151    dataset, Pearson's r = 0.87; Supplementary Figure 5b for the test dataset, Pearson's r = 0.65) and

152    mutated enzymes (Fig. 3b for the whole dataset, Pearson's r = 0.90; Supplementary Figure 5c for

153    the test dataset, Pearson's r = 0.78). Next, several well-studied enzyme-substrate pairs were

154    collected from literature and original dataset from BRENDA[4] and SABIO-RK[5] where each

155    enzyme-substrate pair had $k_{cat}$ values reported for at least 25 unique amino acid substitutions

156    (Supplementary Table 2). The $k_{cat}$ values predicted by the deep learning model correlated very well

157    with the reported experimental $k_{cat}$ values (Pearson's r = 0.94; Fig. 3c). We subsequently divided

158    the entries for each enzyme-substrate pair into two groups based on their experimentally measured

159    $k_{cat}$ values: (i) within 0.5-2.0 fold change of the wildtype $k_{cat}$ value ('wildtype-like $k_{cat}$'); or (ii) less

160    than 0.5 fold change of the wildtype $k_{cat}$ value ('decreased $k_{cat}$'). Scarcity of mutated enzymes with

161    $k_{cat}$ values over 2-fold of wildtype $k_{cat}$ precluded defining the 'increased $k_{cat}$' group[17,18]. Using deep

7

162    learning predicted $k_{cat}$ values, we validated that the enzymes from the 'decreased $k_{cat}$' group indeed

163    showed significantly lower $k_{cat}$ values compared to those of enzymes from 'wildtype-like $k_{cat}$'

164    group for all of the enzyme-substrate pairs (Fig. 3d). The deep learning model is thereby able to

165    capture the effects of small changes in protein sequences on activities of individual enzymes.

166

167    To investigate which subsequence or amino acid residues dominate enzyme activity, we applied a

168    neural attention mechanism to back-trace important signals from an output of the neural network

169    toward its input[19]. This approach can assign attention weights to each amino acid residue, which

170    then quantitatively describes its importance for the predicted enzyme activity, where higher

171    attention weight signifies higher importance. By this method, we calculated the attention weights

172    for all residues of the *Homo sapiens* enzyme purine nucleoside phosphorylase (PNP) with inosine

173    as substrate, as rich mutation data is available for this enzyme-substrate pair[20] (Fig. 3e,

174    Supplementary Table 3). Subsequently situating the mutations from the 'wildtype-like $k_{cat}$' and

175    'decreased $k_{cat}$' groups (Fig. 3e) exhibit that mutations from the latter have significantly higher

176    attention weights (Fig. 3f, *P* value = 0.0014, Supplementary Table 4). Mutating amino acid

177    residues with higher attention weights is seemingly having a more substantial effect on enzyme

178    catalytic activity.

179

180    **$k_{cat}$ prediction for metabolic enzyme-substrate pairs in 343 yeast/fungi species**

181    There are reconstructed GEMs for 332 yeast species plus 11 outgroup fungi[21], but among these

182    only 14 GEMs were expanded with enzyme-constraints (ecGEMs) due to limited available $k_{cat}$

183    data[2,21]. Thus, we applied the deep learning model to populate enzyme-constrained genome scale

184    metabolic models (ecGEMs). As our developed deep learning model allows prediction of almost

8

185    all $k_{cat}$ values for metabolic enzymes against any substrates for any species except the pair with

186    generic substrates which does not have SMILES information, this enabled generation of ecGEMs

187    for all 343 yeast/fungi species. By using the metabolite and enzyme information extracted from

188    the 343 GEMs[21] as the input of the deep learning model for $k_{cat}$ prediction (Supplementary Figure

189    6), we predicted $k_{cat}$ values for around three million protein-substrate pairs in 343 yeast/fungi

190    species.

191

192    By inspecting the global trend for the predicted $k_{cat}$ values, we firstly found that yeast and fungal

193    enzymes from primary-CE metabolism have on average the highest $k_{cat}$ value compared with

194    enzymes from primary-AFN, secondary and intermediate metabolism (Supplementary Figure 7a),

195    which is consistent with the global trend of all enzymes (Fig. 2c) and literature report[6]. Secondly,

196    we found that specialist enzymes (with narrow substrate specificity) have higher $k_{cat}$ values

197    compared with generalist (promiscuous enzymes) that each catalyze more than one reaction in the

198    model (Supplementary Figure 7b). This is aligned with the hypothesis that ancestral enzymes that

199    exhibit broad substrate specificity and low catalytic efficiency improve their $k_{cat}$ when they evolve

200    to be a specialist through processes of mutation, gene duplication and horizonal gene transfer.

201    Consistent with reports for *E. coli*[22], this observation also holds for fungi. Thirdly, we investigated

202    whether sequence conservation trends with $k_{cat}$ values. The ratio of non-synonymous over

203    synonymous substitutions, denoted as dN/dS, is commonly used to detect proteins undergoing

204    adaptation[23]. Conserved enzymes with a lower dN/dS have significantly higher $k_{cat}$ values

205    compared with relatively lesser conserved enzymes (with high dN/dS), implying that conserved

206    yeast/fungi enzymes under evolutionary pressure are adapted to have higher $k_{cat}$ values

207    (Supplementary Figure 7c).

9

208

**Bayesian approach for 343 ecGEMs reconstruction**

210   Using the predicted $k_{cat}$ values for 343 yeast/fungi species we generated 343 DL-ecGEMs

211   (ecGEMs parameterized with $k_{cat}$ values derived from deep learning model prediction). Since the

212   training data for the deep learning model were primarily measured *in vitro*, this implies that also

213   *in vitro* $k_{cat}$ values are predicted by the deep learning model, which is undesired as *in vitro* $k_{cat}$

214   values can be considerably different from their *in vivo* counterparts[24]. To resolve these

215   uncertainties, we adopted a Bayesian genome scale modeling approach, which has been

216   successfully applied to resolve temperature dependence of yeast metabolism by quantifying and

217   reducing uncertainties in model parameters[25]. Here, we used predicted $k_{cat}$ values as mean values

218   for *Prior* distribution and used experimentally measured phenotypes to update it to *Posterior*. The

219   experimental data on yeast/fungi species were collected from literature, collating 445 entries on

220   growth data for 76 species with 16 carbon sources (Supplementary Figure 8, Supplementary Table

221   5). A sequential Monte Carlo based approximate Bayesian computation (SMC-ABC) approach[25]

222   was implemented to sample the $k_{cat}$ (Methods). The ecGEMs parameterized with the mean values

223   of sampled *Posterior* $k_{cat}$ values were hereafter represented as *Posterior*-mean-DL-ecGEMs.

224

225   To test the generality of this SMC-ABC approach and monitor the training process, we first applied

226   this method to ecGEM of *S. cerevisiae*, which has the most abundant experimental data. The

227   experimental phenotype datasets for *S. cerevisiae* were split into training (50%) and test datasets

228   (50%). The training dataset was used to update the *Prior*, which would then be tested on the test

229   dataset after each generation. RMSE between the experimental measurement and prediction for

230   the test dataset was reduced proportionally with the training dataset. After 30 generations, RMSE

10

231    for the training dataset was 0.5 and for the test dataset was 1, which demonstrates the generalization

232    of the SMC-ABC approach (Supplementary Figure 9).

233

234    The Bayesian learning process for *S. cerevisiae* and *Y. lipolytica* are shown as examples (Fig. 4 &

235    Supplementary Figure 10). We calculated RMSE values between measurements and predictions

236    for batch and chemostat growth of *S. cerevisiae* and *Y. lipolytica* under different carbon sources.

237    After several generations, the ecGEMs parameterized with sampled *Posterior* $k_{cat}$ achieved with a

238    RMSE lower than 0.5 (Fig. 4a & Supplementary Figure 10a), which can accurately describe the

239    experimental observations. For instance, the *S. cerevisiae* ecGEM with *Posterior* mean $k_{cat}$ values

240    captures the metabolic shift at increasing growth rate (Fig. 4b)—known as the Crabtree effect[26]—

241    while *Y. lipolytica* respires at its maximum growth rate (Supplementary Figure 10b). When

242    exploring which parameters were updated during the Bayesian process, a principal component

243    analysis (PCA) for all 9,500 generated $k_{cat}$ sets (95 generations with 100 sets each) showed a

244    gradual move from the *Prior* distribution to the distinct *Posterior* distribution (Fig. 4c for *S.*

245    *cerevisiae*). The similar gradual move was also observed for *Y. lipolytica* (Supplementary Figure

246    10c). By comparing the variances of the deep learning and sampled *Posterior* $k_{cat}$ datasets, we

247    found that the Bayesian training process mostly affected variance but not mean predicted $k_{cat}$ values

248    (Fig. 4d-e). For *S. cerevisiae*, 2,644 enzyme-substrate pairs reduced their $k_{cat}$ variance (Šidák adj.

249    one-tailed F-test $P$ value < 0.01), while only 146 pairs changed their mean predicted $k_{cat}$ (Šidák adj.

250    Welch's t test $P$ value < 0.01). For the non-conventional yeast *Y. lipolytica*, the value is 2,721 and

251    159 (Supplementary Figure 10d-e). Consequentially, the sampled *Posterior* $k_{cat}$ has a strong

252    correlation with the deep learning predicted $k_{cat}$ (Pearson's r = 0.83, for *S. cerevisiae*, Fig. 4f;

253    Pearson's r = 0.83, for *Y. lipolytica*, Supplementary Fig. S10f).

11

254

**Deep learning and Bayesian approach improve ecGEMs quality**

256     We subsequently generated *Posterior*-mean-ecGEMs from corresponding DL-ecGEMs for all the

257     343 yeast/fungi species. For comparison, we also built ecGEMs for the same species with a

258     classical $k_{cat}$ parameterization strategy that queried the BRENDA[4] and SABIO-RK[5] databases to

259     assign measured $k_{cat}$ values to enzyme/reaction pair in the model[2,27]. In case of missing data, certain

260     flexibility was introduced by matching the $k_{cat}$ value to other substrates, organisms, or even

261     introducing wild cards in the EC number. This approach is how ecGEMs are routinely

262     parameterized with $k_{cat}$ values, and the resulting models are hereafter referred to as Classical-

263     ecGEMs. The Classical-ecGEMs yielded $k_{cat}$ values for ca. 40% of enzymes included in the model

264     and generated enzymatic constraints for ca. 60% of the enzyme annotated reactions, while DL-

265     ecGEMs and their derived *Posterior*-mean-ecGEMs covered $k_{cat}$ values for ca. 80% of enzymes

266     and defined enzymatic constraints for ca. 90% of enzymatic reactions (Fig. 5a-b). While Classical-

267     ecGEMs have fewer assigned $k_{cat}$ values, their reconstruction pipeline also relies heavily on correct

268     enzyme EC number annotations and available measured $k_{cat}$ values in the databases, contrasting

269     with the DL-ecGEM reconstruction that relies only on protein sequences and substrate SMILES

270     while resulting in a higher coverage. The missing prediction for DL-ecGEMs and derived

271     *Posterior*-mean-ecGEMs are due to the missing $k_{cat}$ prediction for generic substrates which does

272     not have SMILES information.

273

274     The *Posterior*-mean-ecGEMs and DL-ecGEMs do not only have improved $k_{cat}$ coverage but also

275     outperform Classical-ecGEMs in the prediction of exchange rates (Fig. 5c) and are able to predict

276     maximum growth rates in line with the experimentally measured maximum growth rates under

12

277    different carbon sources and oxygen availabilities (Fig. 5d & more detailed Supplementary Figure

278    11). Moreover, we used the three types of models to predict required protein abundances and

279    compared this with published quantitative proteomics data from three species with different carbon

280    sources, culture mode and medium setup (Supplementary Table 6). Proteome predictions from

281    *Posterior*-mean-ecGEMs had the lowest RMSE, while DL-ecGEMs already reduced the RMSE

282    by 30% when compared to Classical-ecGEMs (Fig. 5e). Combined, this showed that not only the

283    increased $k_{cat}$ coverage but also the Bayesian learning approach contributed to ecGEMs that are

284    better representations of the 343 fungi/yeast species.

285

286    **$k_{cat}$ profile comparison enables to identify phenotype-related enzyme**

287    The predicted $k_{cat}$ values were furthermore able to distinguish between Crabtree positive and

288    negative yeast species. There is much interest in understanding the presence of the Crabtree

289    phenotype among yeast species[28,29], and a model of *S. cerevisiae* energy metabolism has been used

290    to interpret this phenotype by comparing protein efficiency, i.e. ATP produced per protein mass

291    per time, in its two energy-producing pathways. It was postulated that the Crabtree effect is related

292    to the high yield (HY) pathway (containing Embden–Meyerhof–Parnas (EMP) pathway,

293    tricarboxylic acid (TCA) cycle and electron transport chain (ETC)) having a lower protein

294    efficiency than the low yield (LY) pathway (containing EMP plus ethanol formation) (Fig. 6a)[1].

295    We here used the *Posterior*-mean-ecGEMs of 102 yeast species (of which 25 are Crabtree positive

296    and 77 are negative with experimental reported phenotype) to similarly calculate protein

297    efficiencies of HY and LY pathways. Of the 102 species we simulated, 89% follow the same trend

298    that Crabtree positive species have a higher LY efficiency while negative species have a higher

299    HY efficiency compared with its LY efficiency, which suggests that Crabtree positive yeast

13

300    species are more protein efficient using the LY pathway than the HY pathway for producing the

301    same amount of ATP (Supplementary Table 7). For five commonly studied species the results are

302    shown in Fig. 6b, and even though ATP yields in their HY pathways may be different in these

303    species, primarily due to the presence of Complex I, they still follow the same trend

304    (Supplementary Table 7). Inconsistencies in strains where the HY/LY protein efficiency ratio did

305    not trend with the Crabtree effects might be due to additional regulation not considered in

306    ecGEMs[30].

307

308    With the predicted genome scale $k_{cat}$ profiles for yeast species, we can investigate whether key

309    enzymes show significant different $k_{cat}$ among 25 Crabtree positive and 77 negative species. Of

310    the enzymes in the energy-producing pathways, only pyruvate kinase, citrate synthase, fumarase

311    and phosphoglucose isomerase had significantly different $k_{cat}$ values (Fig. 6c). Since fumarase and

312    phosphoglucose isomerase can operate in reversible direction, it is hard to explain the kinetic effect

313    towards the Crabtree effect. Thus, we would not further discuss the impact of these two enzymes

314    on the Crabtree effect. The $k_{cat}$ values of pyruvate kinase were higher in Crabtree positive species

315    compared to negative species (*P* value = 0.009 for deep learning predicted $k_{cat}$ values, Fig. 6c).

316    This aligns with a report that increasing pyruvate kinase activity in the Crabtree positive species

317    *Schizosaccharomyces pombe* would increase its fermentation ratio, decrease the growth

318    dependence on respiration and provide resistance to growth inhibiting effects of antimycin A,

319    which inhibits the respiratory complex III[31]. Citrate synthase catalyzes the first and rate-limiting

320    step of the TCA cycle[32], condensing acetyl-coenzyme A and oxaloacetate to form citrate. We found

321    that the $k_{cat}$ of citrate synthase of Crabtree negative species are higher than the Crabtree positive

322    (*P* value = 0.008), which would benefit metabolic flux from entering the TCA cycle (Fig. 6a &

14

323    6c). This is consistent with [13]C-metabolic flux analysis results, which showed that Crabtree

324    negative species have higher TCA flux than Crabtree positive species[33,34].

325

326    **Discussion**

327    The diversity of biochemical reactions and organisms makes it difficult to generate genome scale

328    $k_{cat}$ profiles. Here we presented a deep learning model to predict $k_{cat}$ values of all metabolic

329    enzymes against all substrates, only requiring substrate SMILES and protein sequences of the

330    enzymes as input, simplifying the feature selection process required for the previous machine

331    learning model[9]. This deep learning approach can therefore be used as a versatile $k_{cat}$ prediction

332    tool for any species as long as protein sequence and substrate SMILES are available.

333

334    Another advantage of the deep learning model is that it can capture $k_{cat}$ changes towards precise

335    single amino acid substitutions. As amino acid substitution is a powerful technique in the enzyme

336    evolution field and is routinely used to probe the enzyme catalytic mechanism[35,36], it is valuable

337    that attention weight calculation with our deep learning model can identify which amino acid

338    residues have a major impact on the enzyme activity. Particularly, most amino acid substitution

339    experiments performed mutagenesis in the substrate binding site region, since it is hypothesized

340    that the binding region would have a high impact towards the catalytic activity. However, the

341    profound impact remote regions can have towards the catalytic activity has been reported[37,38]. Here,

342    we found high attention weights for the inosine binding region of human PNP enzyme, while also

343    identifying various non-binding residue sites with high attention weight that deserve further

344    validation. In total, our deep learning model is able to predict amino acid substitutions that can

345    impact $k_{cat}$ values and thereby serve as part of the protein engineering toolbox[39].

15

346

347     The deep learning model is able to predict genome scale $k_{cat}$ profiles for any species. Phenotype

348     related key enzymes can be identified through comparison of $k_{cat}$ values across groups with diverse

349     phenotypes, as done here to identify pyruvate kinase and citrate synthase as Crabtree-effect related

350     enzymes. This approach can as well be applied to identify phenotype related enzymes in other

351     species or even compare among species from different phylogenetic domains. Besides that, global

352     trends in enzyme evolution such as among generalist and specialist enzymes, can be analyzed.

353

354     On the other hand, predicted genome scale $k_{cat}$ profiles can facilitate the reconstruction of enzyme-

355     constrained models of metabolism. Deep learning predicted $k_{cat}$ proved to be a more comprehensive

356     but still practical alternative to matching *in vitro* $k_{cat}$ values from BRENDA[4] and SABIO-RK[5]

357     database as is common in Classical-ecGEMs[2,27,40]. Besides the limitation of the EC number

358     annotation for less studied species, $k_{cat}$ values measured for the well-studied species are also far

359     away from completeness (Supplementary Figure 1c). For the well-studied species *S. cerevisiae*,

360     only 47 $k_{cat}$ values are fully matched with proteins and substrates in the GEM, while other $k_{cat}$

361     values are mostly from fuzzy matching with other substrates, organisms, or even introducing wild

362     cards in the EC number[2], which also can introduce considerable uncertainty in the reconstructed

363     Classical-ecGEMs. In the earlier published ecGEM reconstruction, a lot of manual work is

364     required to ensure the functionality of Classical-ecGEMs[2]. Compared with the Classical-ecGEM

365     reconstruction, DL-ecGEMs is fully automatic, with reduced uncertainty, significantly increased

366     enzyme coverage and $k_{cat}$ coverage for enzymatic reactions and have a more reliable proteome

367     prediction. If there are available experimental growth data, then the ecGEM reconstruction can be

368     further improved through a Bayesian approach. Here, we showed that *Posterior*-mean-ecGEMs

16

369    are more accurate representatives for their phenotypes and the proteome predictions are also

370    improved, which illustrates how functional ecGEMs can be automatically reconstructed.

371

372    In conclusion, we showed how a deep learning approach yields realistic $k_{cat}$ which can be used to

373    direct future genetic engineering, understand enzyme evolution, reconstruct ecGEMs that can be

374    used to simulate metabolic flux and phenotype prediction. Besides that, we envision many other

375    possible uses of this deep learning based $k_{cat}$ prediction tool such as a novel tool in genome mining

376    and Genome-Wide Association Studies (GWAS) analysis. We also envision this automatic

377    Bayesian ecGEM reconstruction pipeline for further usage in ecGEMs reconstruction, for omics

378    data incorporation and analysis.

379

380    **Method and materials**

381    **Preparation of the dataset for deep learning model development**

382    The dataset used for deep learning model construction was extracted from the BRENDA[4] and

383    SABIO-RK database[5] on 10 July 2020 by customized scripts via Application Programming

384    Interface (API). We generated a comprehensive dataset including the substrate name, organism

385    information, Enzyme Commission number (EC number), protein ID (UniProt ID), enzyme type,

386    and $k_{cat}$ values. Besides, substrate SMILES (Simplified Molecular Input Line Entry System), a

387    string notation to represent the substrate structure, was extracted using substrate name to query the

388    PubChem compound database[41], which is the largest database of chemical compound information

389    and is easy to access[42]. As different substrates usually have various synonyms in different database

390    and GEMs, we used a customized Python-based script to ensure that the same canonical SMILES

17

391    could be output for the same substrates with various synonyms, which is essential to help filter

392    redundant entries obtained from different databases (Supplementary Figure 2).

393

394    For the BRENDA database[4], 69,140 entries could be found after downloading and simply

395    processing the accessible data, including 46,417 entries with wildtype enzymes and 22,723 entries

396    with mutated enzymes according to the classification of enzyme type. All these entries contain the

397    required information regarding substrate name, organism, EC number, UniProt ID, enzyme type

398    and $k_{cat}$ value. Then we removed duplicates in the entries, and if there are multiple reported

399    measurements for the same enzyme, we only used the maximum value. For the SABIO-RK

400    database[5], the same data cleaning process was performed. Besides that, we removed the entries

401    with non-standard units for $k_{cat}$ values, such as $s^{-1}*g^{-1}$, $mol*s^{-1}*g^{-1}$, J/mol, etc. All

402    $k_{cat}$ values were converted to the unit in $s^{-1}$. Available SMILES for substrates were obtained via

403    the API of the PubChem database[41]. Then we combined the dataset extracted from BRENDA

404    database and the SABIO-RK database. Due to high overlap between these two databases, 48,659

405    unique entries could be found after data cleaning by merging the entries with the same substrate

406    name, EC number, organism, enzyme type and $k_{cat}$ value for both databases, and all of the entries

407    have specific substrate SMILES information. Besides the similar approach to keep the maximal

408    values for the multiple measurement, duplicates caused by different synonyms usage in these two

409    databases are filtered using the canonical SMILES. Next, protein sequences are queried with two

410    methods, for entries with UniProt ID information, the amino acid sequences could be obtained via

411    the API of the UniProt database[43]; for entries without UniProt ID, the amino acid sequences were

412    acquired from the UniProt database[43] and the BRENDA database[4] based on their EC number and

413    organism information. After that, the sequences of those entries with wildtype enzymes were

18

414    mapped directly and the sequences of those entries with mutated enzymes were changed according

415    to the mutated sites. Finally, 16,838 entries (including 9,411 entries with wildtype enzymes and

416    7,427 entries with mutated enzymes) were left as the high-quality dataset for deep learning model

417    construction. Detailed numbers for the data cleaning can be found in Supplementary Figure 2. Data

418    availability:

419    https://github.com/SysBioChalmers/DLKcat/tree/master/DeeplearningApproach/Data/database

420

421    **Construction of the deep learning pipeline**

422    In this work, we developed an approach for *in vitro* $k_{cat}$ value prediction by combining a graph

423    neural network (GNN) for substrates and a convolutional neural network (CNN) for proteins. The

424    integration of GNN and CNN can be naturally used to handle pairs of data with different structures,

425    i.e., molecular graphs and protein sequences. In this approach, substrates are represented as

426    molecular graphs where the vertices are atoms, the edges are chemical bonds, and proteins are

427    represented as sequences in which the characters are amino acids.

428

429    For substrates, there are just a few types of chemical atoms (e.g., carbon and hydrogen) and

430    chemical bonds (e.g., single bond and double bond). To obtain more learning parameters, we

431    employed r-radius subgraphs to get the vector representations, which are induced by the

432    neighboring vertices and edges within radius r from a vertex[44]. Firstly, substrate SMILES was

433    converted to a molecular graph using RDKit (https://www.rdkit.org). Given a substrate graph, the

434    GNN can update each atom vector and its neighboring atom vectors transformed by the neural

435    network via a non-linear function, e.g., ReLU[45]. Besides, two transitions were developed in the

436    GNN, including vertex transitions and edge transitions. The aim of transitions is to ensure that the

19

437    local information of vertices and edges is propagated in the graph by iterating the process and

438    summing neighboring embeddings. And the final output of the GNN is a set of real-valued

439    molecular vector representations for substrates.

440

441    Similarly, by using the CNN to scan protein sequences, we can obtain low-dimensional vector

442    representations for protein sequences transformed by the neural network via a non-linear function,

443    e.g., ReLU. To apply the CNN to proteins, we defined 'words' in protein sequence and split a

444    protein sequence into an overlapping n-gram (n = 1, 2, 3) amino acids[46]. In this work, to avoid

445    low-frequency words in the learning representations, relatively smaller n-gram number of 1, 2 or

446    3 was set. Also, other important parameters of the neural networks (CNN & GNN) were set as

447    follows: number of layers in CNN: 2, 3 or 4; number of time steps in GNN: 2, 3 or 4; window size:

448    11 (fixed); r-radius: 0, 1 or 2; vector dimensionality: 5, 10 or 20. These different settings were

449    explored based on R Squared ($R^2$) in Equation 1 during the hypermeter tuning to find which

450    hyperparameter is better for improving the deep learning performance. And finally, we used the

451    optimal hyperparameters to train our deep learning model.

$$R^2 = 1 - \frac{\sum_{i=1}^{n}(y_{ie}-y_{ip})^2}{\sum_{i=1}^{n}(y_{ie}-\bar{y})^2} \quad (1)$$

452

453    where $y_{ip}$ is the predicted $k_{cat}$ value, $y_{ie}$ is the experimental $k_{cat}$ value, n is the total number of

454    validation dataset.

455

456    After the acquisition of the substrate molecular vector representations and the protein sequence

457    vector representations, we concatenated them together and an output vector ($k_{cat}$ value) to train the

458    deep learning framework. During the training process, all the datasets were shuffled at the first

459    step, and then were randomly split into training dataset, validation dataset and test dataset at the

20

460    ratio of 80%:10%:10%. Given a set of substrate-protein pairs and the $k_{cat}$ values in the training

461    dataset, the aim of training process is to minimize its loss function. The best model was chosen

462    according to the minimal Root Mean Square Error (RMSE) in Equation 2 on the validation dataset

463    with the least spread between training dataset and validation dataset. For building and training

464    models, the PyTorch v1.4.0 software package was utilized and accessed using the python interface

465    under CUDA/10.1.243.

466    $$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_{ip} - y_{ie})^2} \quad (2)$$

467    where $y_{ip}$ is the predicted $k_{cat}$ value, $y_{ie}$ is the experimental $k_{cat}$ value, n is the total number of dataset

468    (validation dataset or test dataset).

469

470    **Analysis of experimental and deep learning-based $k_{cat}$ values across different metabolic**

471    **contexts**

472    According to the classification of metabolic pathways, metabolic contexts were mainly divided

473    into four different subsystems: primary metabolism-CE (carbohydrate and energy), involving the

474    main carbon and energy metabolism, e.g., glycolysis/gluconeogenesis, TCA cycle, pentose

475    phosphate pathway, etc; primary metabolism-AFN (amino acids, fatty acids, and nucleotides);

476    intermediate metabolism, related to the biosynthesis and degradation of cellular components, such

477    as coenzymes and cofactors; and secondary metabolism, associated with metabolites that are

478    produced in specific cells or tissues, e.g., flavonoid biosynthesis, caffeine metabolism etc[6]. To

479    explore the metabolic subsystems for all of the wildtype enzymes in the experimental dataset, the

480    module in KEGG database[16] was utilized to assign metabolic pathways for enzyme-substrate pairs

481    by linking the detailed metabolic pathway in KEGG API with EC number annotated in each

482    enzyme-substrate pair. Detailed classification can be found in Supplementary Table 1. Using the

483      trained deep learning model, the predicted $k_{cat}$ values were generated for all the enzyme-substrate

484      pairs. The relationship between these predicted $k_{cat}$ values and various metabolic contexts was

485      further analyzed, which was compared with the trends of the annotated experimental results.

486

487      **Interpretation of the reasoning of deep learning with neural attention mechanism**

488      To interpretate which subsequences or residue sites are more important for the substrate, the neural

489      attention mechanism was employed by assigning attention weights to the subsequences[19]. A higher

490      attention weight of one residue means that residue is more important for the enzyme activity

491      towards the specific substrate. Such attention weights were modeled based on the output of the

492      neural network.

493 $$C = \left\{ c_1^{(t)}, c_2^{(t)}, c_3^{(t)}, ..., c_n^{(t)} \right\} \quad (3)$$

494 $$h_{substrate} = f(W_{inter} y_{substrate} + b) \quad (4)$$

495 $$h_i = f(W_{inter} c_i + b) \quad (5)$$

496 $$\alpha_i = \sigma(h^T_{substrate} h_i) \quad (6)$$

497      where C is a set of hidden vectors for the protein sequence, $c_1^{(t)}$ to $c_n^{(t)}$ are the sub-hidden vectors

498      for the split subsequences, $y_{substrate}$ is the substrate molecular vector, $W_{inter}$ and b are the weight

499      matrix and the bias vector in the neural network, respectively, f is a non-linear activation function

500      (e.g., ReLU), $\alpha_i$ is the final attention weight value.

501

502      For a defined protein, it could be split into overlapping n-gram amino acids and calculated as a set

503      of hidden vectors in Equation 3. Given a substrate molecular vector $y_{substrate}$ and a set of protein

504      hidden vectors, the substrate embeddings ($h_{substrate}$) and subsequence embeddings ($h_i$) could be

505      output based on the neural network as shown in Equation 4 and Equation 5. By considering the

22

506    embeddings of y$_{substrate}$, the attention weight value for each subsequence was accessible in Equation

507    6, which represents the importance signals of the protein subsequence towards the enzyme activity

508    for a certain substrate.

509

510    **Prediction of $k_{cat}$ values for 343 yeast/fungi species**

511    The GEMs of 343 yeast/fungi species were downloaded from the GitHub repository[21]. For each

512    model, all reversible enzymatic reactions were split to forward and backward reactions. Reactions

513    catalyzed by isoenzymes were also split to multiple reactions with one enzyme complex for each

514    reaction. Substrates were extracted from the model and mapped to MetaNetX database to get

515    SMILES structure using corresponding annotated MetaNet IDs for metabolites[47]. Protein IDs for

516    the enzymes were from the model.grRules. Since there are around 200 yeast species are newly

517    sequenced[48] and are not included in the UniProt database[43], protein sequences were queried by the

518    protein ID in the protein fasta file for each species (Supplementary Dataset). Reaction IDs,

519    substrate names, substrate SMILES and protein IDs were combined as the input file for the deep

520    learning $k_{cat}$ prediction model.

521

522    **Analysis of $k_{cat}$ values and dN/dS for 343 yeast/fungi species**

523    In a previous study, the genomes of 343 yeast/fungi species combined with comprehensive genome

524    annotations were publicly available[48]. The gene-level dN/dS of gene sequences for pairs of

525    orthologous genes from the 343 species were calculated with yn00 from PAML v4.7[49]. For this

526    computational framework, the input is the single-copy ortholog groups (OGs), and the output is

527    the gene-level dN/dS values extracted from the PAML output files. By mapping the predicted $k_{cat}$

528    values with the gene-level dN/dS values via the bridge of protein ID, a global analysis was

23

529    performed between the $k_{cat}$ values and the dN/dS values for 343 yeast/fungi species across the

530    outgroup (11 fungal species) together with 12 major clades divided by the genus-level phylogeny

531    for 332 yeast species.

532

533    **ecGEM reconstruction**

534    ecGEMs are reconstructed by adding enzymatic constraints (Equation 7) into the basic constraints

535    of basic GEMs.

536    $$v_j \leq k_{cat}^{i,j} * [E_i] \ (7)$$

537    where $v_j$ stands for the metabolic flux (mmol/gDW/h) of the reaction $j$, $[E_i]$ stands for the enzyme

538    concentration for the enzyme $i$ that catalyzes reaction $j$ and $k_{cat}^{i,j}$ is the catalytic turnover number

539    for the enzyme catalyzing reaction $j$. This constraint is applied to all enzymatic reactions with

540    available $k_{cat}$ values.

541

542    We used two formats of ecGEMs in the reconstruction process: we adopted the sMOMENT[27]

543    format in the Bayesian modeling process to speed up the $k_{cat}$ mapping process and linear problem

544    construction in the SMC-ABC search; while in the model evaluation and final format, we used the

545    GECKO format to compile all $k_{cat}$ values in the model S matrix which would be compatible with

546    all developed GECKO functions[2,50]. There is a developed customized function

547    convertToGeckoModel to facilitate the conversion for these two formats.

548

549    Classical-ecGEM reconstruction queries $k_{cat}$ values from BRENDA database by matching the EC

550    number, which is heavily relied on the database EC number annotation for the specific species[2,27].

551    Since more than 200 out of 343 yeast/fungi species are not annotated in UniProt[43] and KEGG[16],

24

552  EC numbers for orthologs annotated in *S. cerevisiae* were borrowed to facilitate Classical-ecGEM

553  reconstruction process for all these 343 species. The $k_{cat}$ extraction process used the criteria from

554  the process 13 in the reconstruction methods of the reference[40].

555

556  DL-ecGEM reconstruction extracts all $k_{cat}$ values from the deep learning predicted file. To assign

557  $k_{cat}$ value for each metabolic reaction, we follow the criteria below 1) $k_{cat}$ values predicted for

558  currency metabolites such as $H_2O$, $H^+$ were excluded; 2) If there are multiple substrates in the

559  reaction, maximum values among substrates were kept; 3) If multiple subunits exist in the enzyme

560  complex, we used the maximum values among all subunits to represent the $k_{cat}$ for the complex.

561

562  *Posterior*-mean-ecGEM reconstruction uses mean values for accepted *Posterior* distribution. The

563  $k_{cat}$ values in the DL-ecGEMs combined with the RMSE (which is 1 in log10 scale) of the $k_{cat}$

564  prediction were used as mean values and variance to make the *Prior* distribution. Each $k_{cat}$ value

565  was described with a log normal distribution N($kcat_i$ , 1). This *Prior* iteratively morphs into a

566  *Posterior* through multiple generations[25]. For each generation, we sampled 128 $k_{cat}$ datasets within

567  the distribution, and 100 among those 128 datasets with smaller distance (see next section for the

568  SMC-ABC distance calculation) between phenotype measurements and predictions which can

569  better represent the phenotype were kept to make the distribution for the next generation. Until the

570  distance is lower than the cutoff (RMSE of 0.5), then we accepted the final distribution as *Posterior*

571  distubiton[25].

572

573  **SMC-ABC distance function**

25

574    Experimental growth data and related exchange rates in batch and chemostat conditions were

575    collected for yeast/fungi species, which are available at Supplementary Table 5. The distance

576    function was designed as RMSE between simulated and experimental values for maximal growth

577    simulations and exchange rates simulations. As for maximal growth simulation, the medium was

578    set in the model by allowing the free uptake of composition, and the objective function was set to

579    maximizing growth. The RMSE was calculated for the simulated and measured growth rates. For

580    the exchange rates simulation, the carbon source uptake rates were constrained based on

581    experimental measurements, and the objective function was also set to maximizing growth. The

582    RMSE was calculated for the simulated and measured exchange rates of all measured exo-

583    metabolites. All measured and simulated rates were normalized by the carbon numbers of the

584    corresponding metabolites before calculation of RMSE. The carbon number for biomass is 41

585    (mean value for the molecular wight of 1 Cmol biomass of yeast is ~24.42 g[51], the biomass equals

586    to 1000 mg). Note that if the substrate or byproduct does not contain any carbon such as $O_2$, then

587    the normalizing number is 1. Then the average RMSE of both simulations was used to represent

588    the distance. SMC-ABC search would stop once the RMSE reaches the accepted value or reaches

589    the maximum generation. The accepted value for the distance is set to be lower than 0.5 and the

590    maximum generation is set to be 150.

591

592    **Simulations with ecGEMs**

593    We performed different kinds of simulations using the ecGEMs including simulations of growth

594    and protein abundance. Different mediums and growth conditions were set to match the

595    experiment measurement condition, e.g., using xylose as the carbon source or anaerobic condition.

596    Since there are no measured total protein abundance in the biomass for all yeast/fungi species, we

26

597    used the protein content mass to serve as the total protein abundance for each species and used a

598    sigma factor of 0.5 to serve as the ratio of metabolic protein ratio in total protein abundance.

599

600    **Statistical tests for comparison between sampled *Prior* and *Posterior* dataset**

601    Sampled *Prior* and *Posterior* $k_{cat}$ datasets were compared for the difference in the mean values and

602    the variance. Welch's t test was used to test the significance for the mean values, while one-tailed

603    F-test was used for the reduced variances. The cutoff for the significance was set to 0.01 for the

604    adjusted *P* value corrected by the Šidák method.

605

606    **Proteome data collection**

607    All    collected    proteome    data    are    available    in    the    GitHub    repository

608    (https://github.com/SysBioChalmers/DLKcat/tree/master/BayesianApporach/Data/Proteome_ref.

609    xlsx). For relative proteome datasets, we normalized by the identical condition of the absolute

610    proteome data from the literature following the same method as[52,53]. Reference absolute datasets

611    for those relative proteome datasets were documented in the same file.

612

613    **Calculation of protein cost and efficiency**

614    To calculate the protein cost of the HY pathway, the glucose uptake rate was fixed at 1

615    mmol/gDW/h, and the non-growth associated maintenance energy (NGAM) reaction was

616    maximized. The total protein pool reaction was then minimized with fixing the NGAM reaction at

617    the maximized value. The minimized flux through the total protein pool reaction is the protein cost

618    of the HY pathway for converting one glucose to ATP. As for the protein cost calculation of LY

619    pathway, glucose uptake rate was fixed at 1 mmol/gDW/h, the ethanol production was maximized.

27

620    Then the ethanol exchange rate was fixed at the maximized value, and NGAM was maximized.

621    After that, NGAM was also fixed at the maximized value, and total protein pool was minimized

622    to calculate the protein cost for LY pathway. We also examined the flux distribution to ensure that

623    other energy producing pathways are all inactive during this simulation. Protein efficiency is

624    defined as the protein cost for producing one flux ATP in both pathways.

625

626    **Code and data availability**

627    To facilitate further usage, we provide all codes, example and detailed instruction in GitHub

628    repository: https://github.com/SysBioChalmers/DLKcat. Protein sequence fasta files, deep

629    learning predicted $k_{cat}$ values, classcial-ecGEMs, DL-ecGEMs and *Posterior*-mean-ecGEMs for

630    343 yeast/fungi species are available as Supplementary Dataset on the zenodo:

631    https://doi.org/10.5281/zenodo.5164210.

632

633    **Author contribution**

634    F.L, L.Y., H.L. and J.N. designed the research. F.L. and L.Y. performed the research. F.L, L.Y.,

635    Y.C., G.L., E.K. and J.N. analyzed the data. L.Y. and M.E. collected the $k_{cat}$ data. F.L, L.Y., H.L,

636    G.L., Y.C., M.E., E.K. and J.N. wrote the paper. All authors approved the final paper.

637

638    **Acknowledgement**

643    Computing (SNIC) at Chalmers Centre for Computational Science and Engineering (C3SE) and

644    High Performance Computing Center North (HPC2N), partially funded by the Swedish Research

645    Council through grant agreement no. 2018-05973.

646

647    **Competing interests**

648    The authors declare no competing interests.

649

650    **Reference:**

651    1.    Chen, Y. & Nielsen, J. Energy metabolism controls phenotypes by protein efficiency and

652          allocation. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 17592–17597 (2019).

653    2.    Sánchez, B. J. *et al.* Improving the phenotype predictions of a yeast genome-scale

654          metabolic model by incorporating enzymatic constraints. *Mol. Syst. Biol.* **13**, 935 (2017).

655    3.    Klumpp, S., Scott, M., Pedersen, S. & Hwa, T. Molecular crowding limits translation and

656          cell growth. *Proc. Natl. Acad. Sci. U. S. A.* **110**, 16754–16759 (2013).

657    4.    Schomburg, I. *et al.* The BRENDA enzyme information system–From a database to an

658          expert system. *J. Biotechnol.* **261**, 194–206 (2017).

659    5.    Wittig, U., Rey, M., Weidemann, A., Kania, R. & Müller, W. SABIO-RK: an updated

660          resource for manually curated biochemical reaction kinetics. *Nucleic Acids Res.* **46**,

661          D656–D660 (2018).

662    6.    Bar-Even, A. *et al.* The moderately efficient enzyme: evolutionary and physicochemical

663          trends shaping enzyme parameters. *Biochemistry* **50**, 4402–4410 (2011).

664    7.    Chen, Y. & Nielsen, J. Mathematical modelling of proteome constraints within

665          metabolism. *Curr. Opin. Syst. Biol.* (2021).

666   8.    Davidi, D. & Milo, R. Lessons on enzyme kinetics from quantitative proteomics. *Curr.*

667         *Opin. Biotechnol.* **46**, 81–89 (2017).

668   9.    Heckmann, D. *et al.* Machine learning applied to enzyme turnover numbers reveals

669         protein structural correlates and improves metabolic models. *Nat. Commun.* **9**, 1–10

670         (2018).

671   10.   Nilsson, A., Nielsen, J. & Palsson, B. O. Metabolic models of protein allocation call for

672         the kinetome. *Cell Syst.* **5**, 538–541 (2017).

673   11.   Kitchin, J. R. Machine learning in catalysis. *Nat. Catal.* **1**, 230–232 (2018).

674   12.   Shrivastava, A. D. & Kell, D. B. FragNet, a Contrastive Learning-Based Transformer

675         Model for Clustering, Interpreting, Visualizing, and Navigating Chemical Space.

676         *Molecules* **26**, (2021).

677   13.   Zrimec, J. *et al.* Deep learning suggests that gene expression is encoded in all parts of a

678         co-evolving interacting gene regulatory structure. *Nat. Commun.* **11**, 6141 (2020).

679   14.   Kroll, A., Heckmann, D. & Lercher, M. J. Prediction of Michaelis constants from

680         structural features using deep learning. *Preprint* at

681         https://doi.org/10.1101/2020.12.01.405928 (2020).

682   15.   Ryu, J. Y., Kim, H. U. & Lee, S. Y. Deep learning enables high-quality and high-

683         throughput prediction of enzyme commission numbers. *Proc. Natl. Acad. Sci.* 201821905

684         (2019).

685   16.   Kanehisa, M., Furumichi, M., Tanabe, M., Sato, Y. & Morishima, K. KEGG: new

686         perspectives on genomes, pathways, diseases and drugs. *Nucleic Acids Res.* **45**, D353–

687         D361 (2017).

688   17.   Yep, A., Kenyon, G. L. & McLeish, M. J. Saturation mutagenesis of putative catalytic

689     residues of benzoylformate decarboxylase provides a challenge to the accepted

690     mechanism. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 5733–5738 (2008).

691  18.  Lin, Y.-H. T., Huang, C. L. V., Ho, C., Shatsky, M. & Kirsch, J. F. A general method to

692     predict the effect of single amino acid substitutions on enzyme catalytic activity. *Preprint*

693     at https://doi.org/10.1101/236265 (2017).

694  19.  Bahdanau, D., Cho, K. & Bengio, Y. Neural machine translation by jointly learning to

695     align and translate. *Preprint* at https://arxiv.org/abs/1409.0473v7 (2014).

696  20.  Erion, M. D. *et al.* Purine nucleoside phosphorylase. 1. Structure-function studies.

697     *Biochemistry* **36**, 11725–11734 (1997).

698  21.  feiranl, hongzhonglu, Domenzain, I. & Yuan, L. SysBioChalmers/Yeast-Species-GEMs:

699     Yeast-Species-GEM. (2021). data sets. zenodo https://doi:10.5281/zenodo.4568962

700  22.  Nam, H. *et al.* Network context and selection in the evolution to enzyme specificity.

701     *Science* **337**, 1101–1104 (2012).

702  23.  Kryazhimskiy, S. & Plotkin, J. B. The population genetics of dN/dS. *PLoS Genet.* **4**,

703     e1000304 (2008).

704  24.  Ringe, D. & Petsko, G. A. Biochemistry. How enzymes work. *Science* **320**, 1428–1429

705     (2008).

706  25.  Li, G. *et al.* Bayesian genome scale modelling identifies thermal determinants of yeast

707     metabolism. *Nat. Commun.* **12**, 1–12 (2021).

708  26.  Van Hoek, P. I. M., Van Dijken, J. P. & Pronk, J. T. Effect of specific growth rate on

709     fermentative capacity of baker's yeast. *Appl. Environ. Microbiol.* **64**, 4226–4233 (1998).

710  27.  Bekiaris, P. S. & Klamt, S. Automatic construction of metabolic models with enzyme

711     constraints. *BMC Bioinformatics* **21**, 19 (2020).

712   28.   Pfeiffer, T. & Morley, A. An evolutionary perspective on the Crabtree effect. *Front. Mol.*

713         *Biosci.* **1**, 17 (2014).

714   29.   de Alteriis, E., Cartenì, F., Parascandola, P., Serpa, J. & Mazzoleni, S. Revisiting the

715         Crabtree/Warburg effect in a dynamic perspective: a fitness advantage against sugar-

716         induced cell death. *Cell Cycle* **17**, 688–701 (2018).

717   30.   Ata, Ö. *et al.* A single Gal4-like transcription factor activates the Crabtree effect in

718         *Komagataella phaffii*. *Nat. Commun.* **9**, 1–10 (2018).

719   31.   Kamrad, S. *et al.* Pyruvate kinase variant of fission yeast tunes carbon metabolism, cell

720         regulation, growth and stress resistance. *Mol. Syst. Biol.* **16**, e9270 (2020).

721   32.   Krebs, H. A. Rate control of the tricarboxylic acid cycle. *Adv. Enzyme Regul.* **8**, 335–353

722         (1970).

723   33.   Christen, S. & Sauer, U. Intracellular characterization of aerobic glucose metabolism in

724         seven yeast species by $^{13}$C flux analysis and metabolomics. *FEMS Yeast Res.* **11**, 263–272

725         (2011).

726   34.   Blank, L. M., Lehmbeck, F. & Sauer, U. Metabolic-flux and network analysis in fourteen

727         hemiascomycetous yeasts. *FEMS Yeast Res.* **5**, 545–558 (2005).

728   35.   Chen, K. & Arnold, F. H. Engineering new catalytic activities in enzymes. *Nat. Catal.* **3**,

729         203–213 (2020).

730   36.   Markel, U. *et al.* Advances in ultrahigh-throughput screening for directed enzyme

731         evolution. *Chem. Soc. Rev.* **49**, 233–262 (2020).

732   37.   Loeb, D. D. *et al.* Complete mutagenesis of the HIV-1 protease. *Nature* **340**, 397–400

733         (1989).

734   38.   Lee, J. & Goodey, N. M. Catalytic contributions from remote regions of enzyme structure.

735          *Chem. Rev.* **111**, 7595–7624 (2011).

736   39.   Tong, H., Küken, A., Razaghi-Moghadam, Z. & Nikoloski, Z. Characterization of effects

737          of genetic variants via genome-scale metabolic modelling. *Cell. Mol. Life Sci.* **78**, 5123–

738          5138 (2021).

739   40.   Chen, Y., Li, F., Mao, J., Chen, Y. & Nielsen, J. Yeast optimizes metal utilization based

740          on metabolic network and enzyme kinetics. *Proc. Natl. Acad. Sci.* **118**, (2021).

741   41.   Kim, S. *et al.* PubChem Substance and Compound databases. *Nucleic Acids Res.* **44**,

742          D1202-13 (2016).

743   42.   Chen, F., Yuan, L., Ding, S., Tian, Y. & Hu, Q.-N. Data-driven rational biosynthesis

744          design: from molecules to cell factories. *Brief. Bioinform.* **21**, 1238–1248 (2020).

745   43.   The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids*

746          *Res.* **45**, D158–D169 (2017).

747   44.   Tsubaki, M., Tomii, K. & Sese, J. Compound-protein interaction prediction with end-to-

748          end learning of neural networks for graphs and sequences. *Bioinformatics* **35**, 309–318

749          (2019).

750   45.   LeCun, Y., Bengio, Y. & Hinton, G. Deep learning. *Nature* **521**, 436–444 (2015).

751   46.   Dong, Q.-W., Wang, X.-L. & Lin, L. Application of latent semantic analysis to protein

752          remote homology detection. *Bioinformatics* **22**, 285–290 (2006).

753   47.   Moretti, S., Tran, V. D. T., Mehl, F., Ibberson, M. & Pagni, M. MetaNetX/MNXref:

754          unified namespace for metabolites and biochemical reactions in the context of metabolic

755          models. *Nucleic Acids Res.* **49**, D570–D574 (2021).

756   48.   Shen, X.-X. *et al.* Tempo and mode of genome evolution in the budding yeast subphylum.

757          *Cell* **175**, 1533–1545 (2018).

758    49.    Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**,

759            1586–1591 (2007).

760    50.    Domenzain, I. *et al.* Reconstruction of a catalogue of genome-scale metabolic models with

761            enzymatic constraints using GECKO 2.0. *Preprint* at

762            https://doi.org/10.1101/2021.03.05.433259 (2021).

763    51.    Popovic, M. Thermodynamic properties of microorganisms: determination and analysis of

764            enthalpy, entropy, and Gibbs free energy of biomass, cells and colonies of 32

765            microorganism species. *Heliyon* **5**, e01950 (2019).

766    52.    Yu, R. *et al.* Nitrogen limitation reveals large reserves in metabolic and translational

767            capacities of yeast. *Nat. Commun.* **11**, 1881 (2020).

768    53.    Metzl-Raz, E. *et al.* Principles of cellular resource allocation revealed by condition-

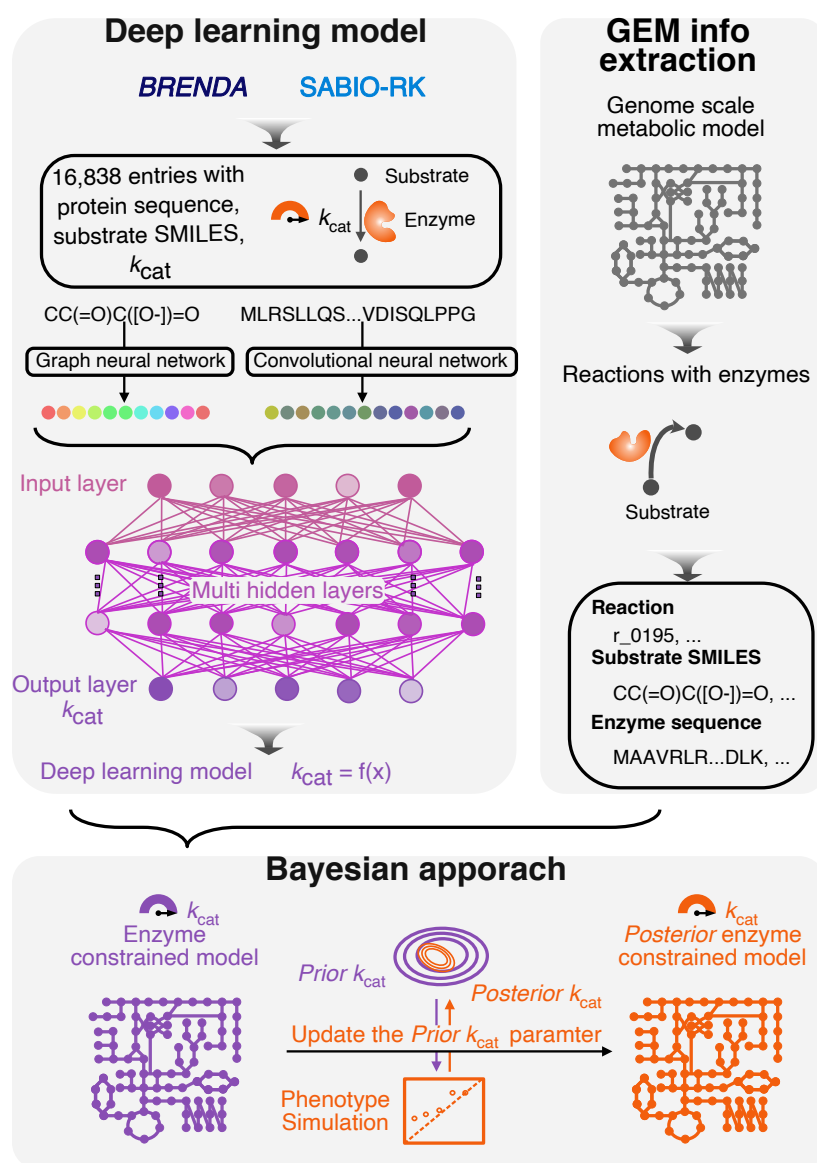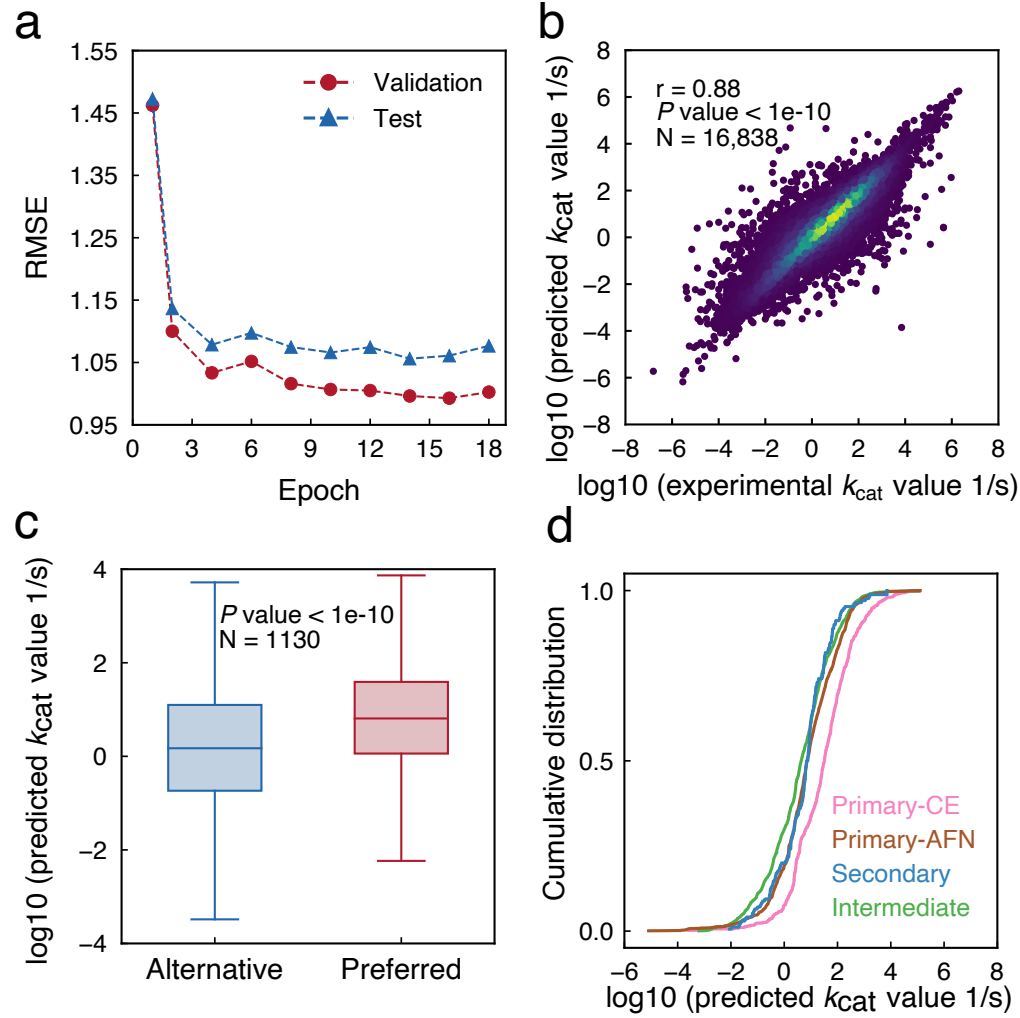769            dependent proteome profiling. *Elife* **6**, (2017).

770

771

772    **Figures**



773

774    **Figure 1** Deep learning of enzyme turnover numbers ($k_{cat}$) for genome scale metabolic model

775    (GEM) parameterization. Firstly, we developed an approach for $k_{cat}$ prediction by combining a

776    graph neural network (GNN) for substrates and a convolutional neural network (CNN) for proteins.

777    Secondly, we extracted information from GEMs as the input for the deep learning model to predict

778    $k_{cat}$ values. Thirdly, we developed a Bayesian facilitated pipeline to reconstruct enzyme-

779    constrained GEMs (ecGEMs) using the predicted $k_{cat}$ profiles from deep learning model.

780

781

782    **Figure 2** Deep learning model performance for $k_{cat}$ prediction. (a) The RMSE of $k_{cat}$ prediction

783    during the training process. (b) Performance of the final deep learning model trained by GNN and

784    CNN. The correlation between predicted $k_{cat}$ value and those present in the whole dataset was

785    evaluated. The brightness of color represents the density of data points. (c) Enzyme promiscuity

786    analysis on the whole dataset. For enzymes with multiple substrates, we divided the substrates as

787    preferred and alternative by their experimental measured $k_{cat}$, then used the predicted $k_{cat}$ values

788    for this boxplot. A two-sided Wilcoxon rank sum test was used to calculate $P$ value. (d) Cumulative

789    distribution of deep learning-based $k_{cat}$ values for enzyme-substrate pairs belonging to different

790    metabolic contexts. Abbreviations: CE, carbohydrate and energy; AFN, amino acids, fatty acids,

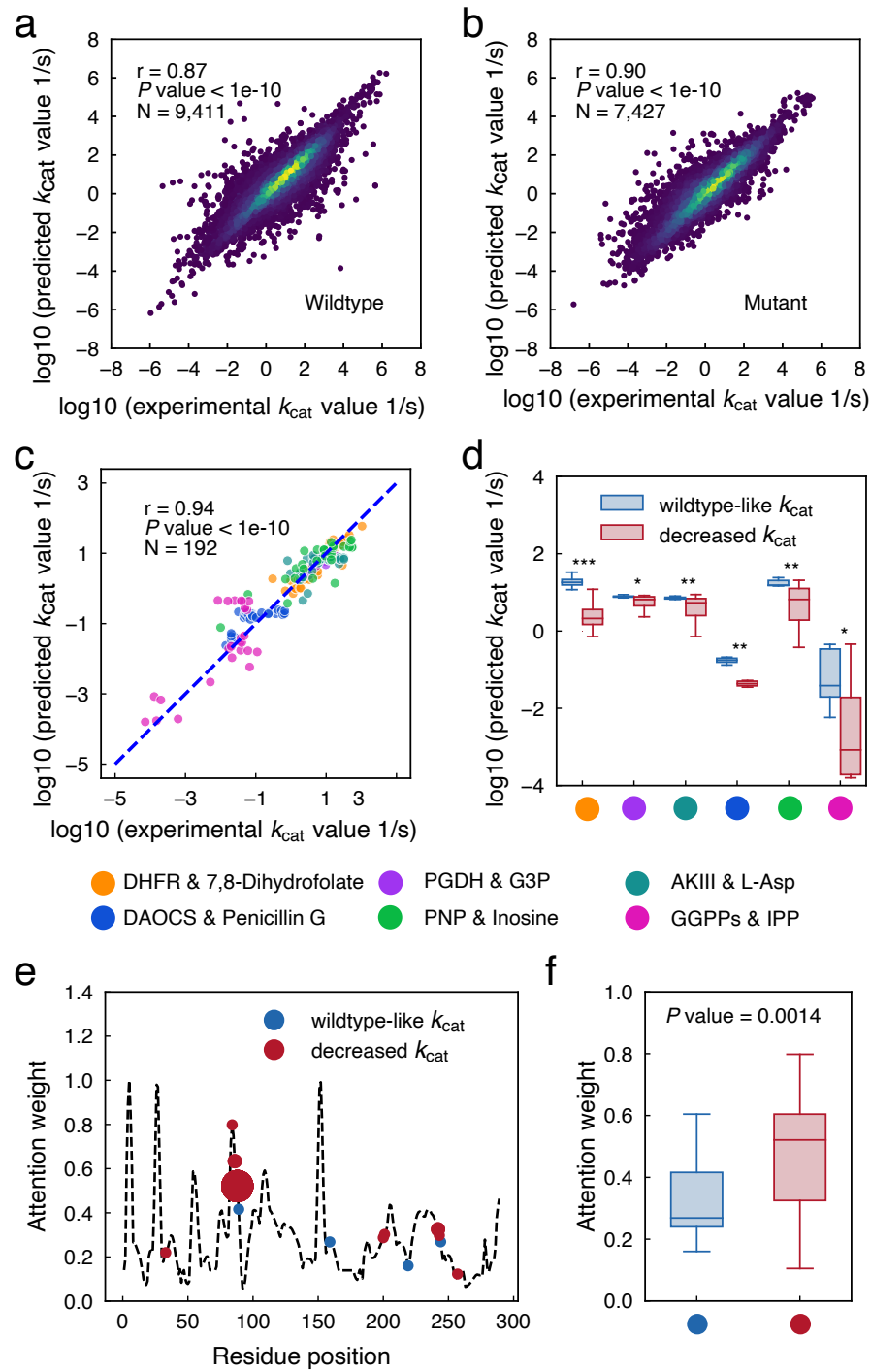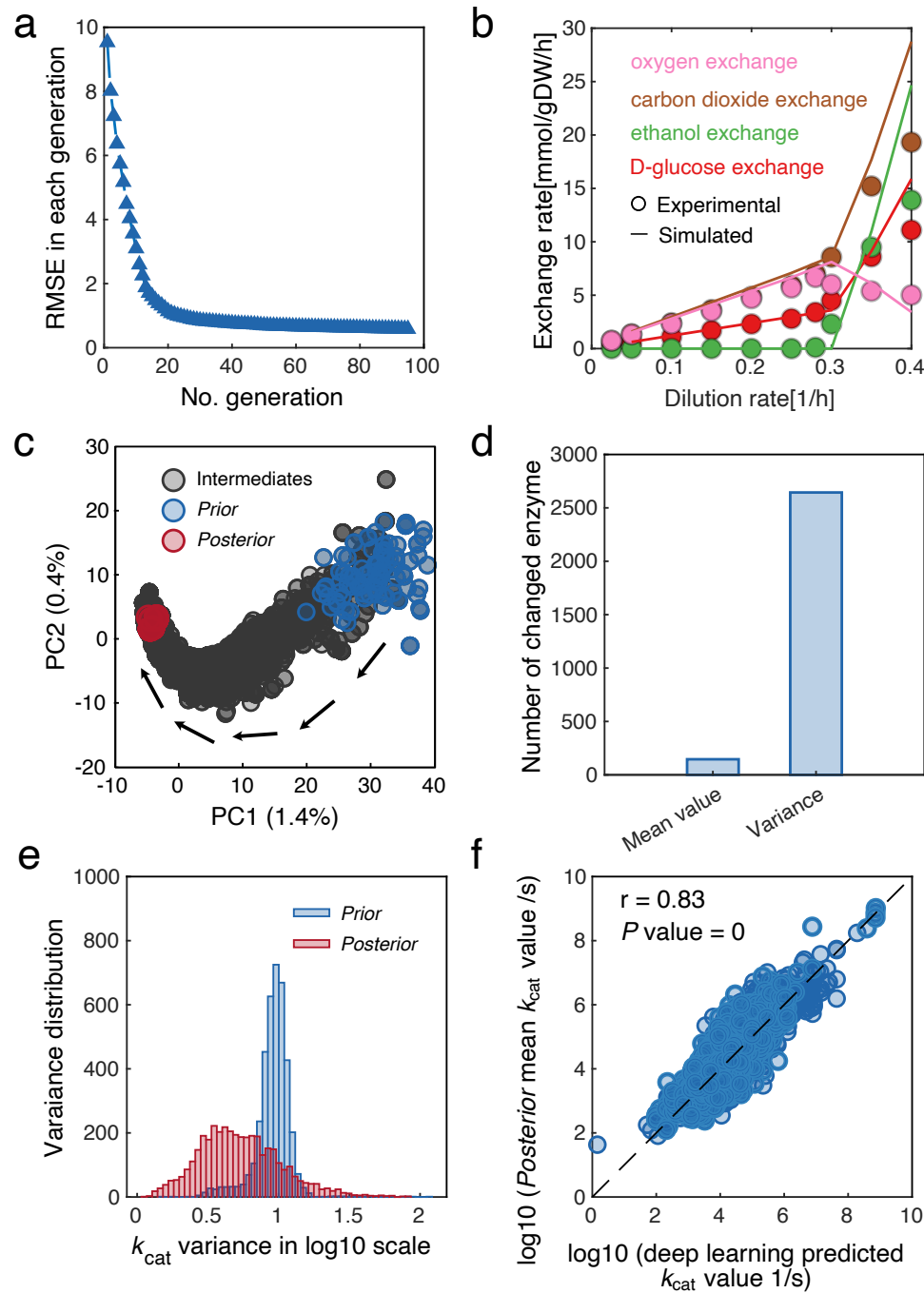791    and nucleotides.

792

795    **Figure 3** Deep learning model for the prediction and interpretation of $k_{cat}$ of mutated enzymes. (a)

796    Prediction performance of $k_{cat}$ values for all of the wildtype enzymes via deep learning model. The

797    brightness of color represents the density of data points. (b) Prediction performance of $k_{cat}$ values

798    for all of the mutated enzymes via deep learning model. The brightness of color represents the

799    density of data points. (c) Comparison between predicted and measured $k_{cat}$ values for several well-

800    studied enzyme-substrate pairs with rich experimental mutagenesis data. Enzyme abbreviations:

801    DHFR, dihydrofolate reductase; PGDH, D-3-phosphoglycerate dehydrogenase; AKIII,

802    aspartokinase III; DAOCS, deacetoxycephalosporin C synthase; PNP, purine nucleoside

803    phosphorylase; GGPPs, geranylgeranyl pyrophosphate synthase. Substrate abbreviations: G3P,

804    glycerate 3-phosphate; L-Asp, L-Aspartate; IPP, isopentenyl diphosphate. (d) Comparison of

805    predicted $k_{cat}$ values on several mutated enzyme-substrate pairs between 'wildtype-like $k_{cat}$' and

806    enzymes with 'decreased $k_{cat}$'. $P$ value < 0.05 (*), $P$ value < 0.01 (**) and $P$ value < 0.001 (***).

807    (e) Attention weight of sequence position in the wildtype PNP enzyme using inosine as the

808    substrate. The mutated enzymes (enzymes with 'wildtype-like $k_{cat}$' and enzymes with 'decreased

809    $k_{cat}$') were marked on the curve according to their mutated position. The dot size indicates the

810    number of mutated enzymes occurring in that mutated position. (f) Comparisons of the overall

811    attention weight for the PNP – Inosine pair between enzymes with 'wildtype-like $k_{cat}$' and enzymes

812    with 'decreased $k_{cat}$'. For two group comparisons in subfigure d and f, a two-sided Wilcoxon rank

813    sum test was used to calculate $P$ value.

814

815

816

817    **Figure 4** Bayesian modeling training performance for *S. cerevisiae* ecGEM. (a) RMSE for

818    phenotype measurement and prediction during Bayesian training process. (b) Simulated exchange

819    rates by *Posterior*-mean-ecGEM (line) compared with experimental data (dot). $k_{cat}$ values in the

820    *Posterior*-mean-ecGEMs here is mean values from 100 sampled *Posterior* datasets after the

821    Bayesian training process. (c) Principal component analysis (PCA) for $k_{cat}$ datasets sampled in the

822    Bayesian training approach. Each parameter in the set was standardized by subtracting the mean

823    and then divided by the standard deviation before PCA. Sampled 100 *Prior* datasets are

824    highlighted in blue, while sampled 100 *Posterior* datasets are highlighted in red. All other datasets

825    were termed as "intermediate" and marked in gray. (d) The number of enzymes with a significantly

826    changed mean values (Šidák adjusted Welch's t test $P$ value $< 0.01$, two-sided) and variance

827    (Šidák adjusted one-tailed F-test $P$ value $< 0.01$) between sampled *Prior* and *Posterior* $k_{cat}$

828    datasets. (e) Variance distribution comparison for *Prior* and *Posterior* distribution. (f) Correlation

829    between deep learning predicted $k_{cat}$ and *Posterior* mean $k_{cat}$.
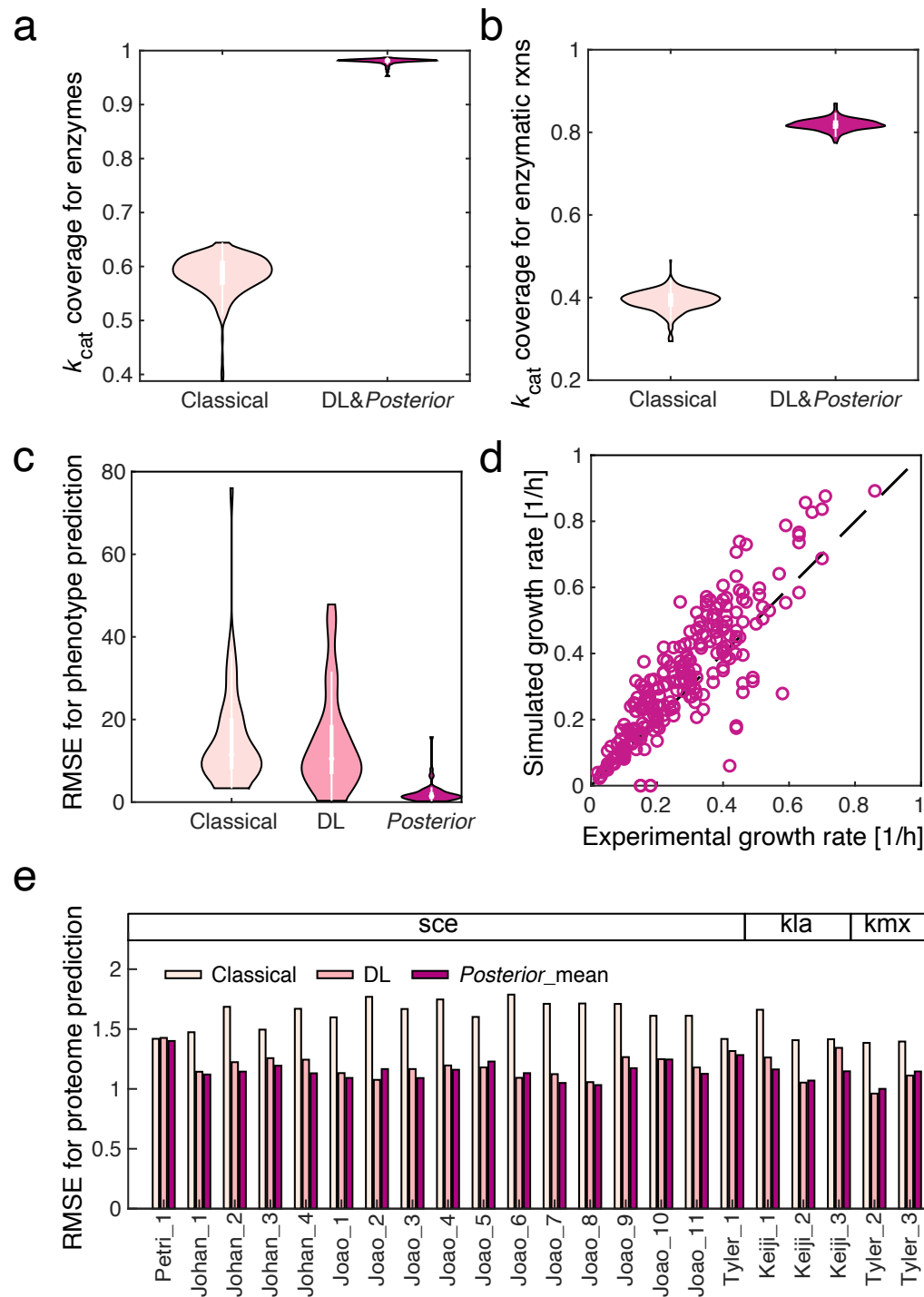
830

41

833 **Figure 5** Evaluation of three ecGEM modelling frameworks including Classical-ecGEM, DL-

834 ecGEMs and *Posterior*-mean-ecGEMs. Enzymatic constraint coverage comparison for (a)

835 enzymes and (b) enzymatic reactions. (c) RMSE for the phenotype prediction. (d) Growth

836 prediction for *Posterior*-mean-ecGEMs. (e) Performance of three types of ecGEMs in predicting

837 quantitative proteome data, Classical-ecGEM, DL-ecGEM and *Posterior*-mean-ecGEM are shown.

838 RMSE is shown on log10 scale. Classical-ecGEM is constructed following the pipeline to extract

839 $k_{cat}$ profiles from BRENDA and SABIORK, DL-ecGEMs are constructed using the $k_{cat}$ profiles

840 predicted from the deep learning model. *Posterior*-mean-ecGEMs here were parameterized by the

841 $k_{cat}$ profiles of the mean values from 100 *Posterior* datasets after the Bayesian training process.

842 Detailed conditions for those proteome datasets can be found in the Supplementary Table 6 and

843 collected proteome dataset are available in GitHub repository.

844

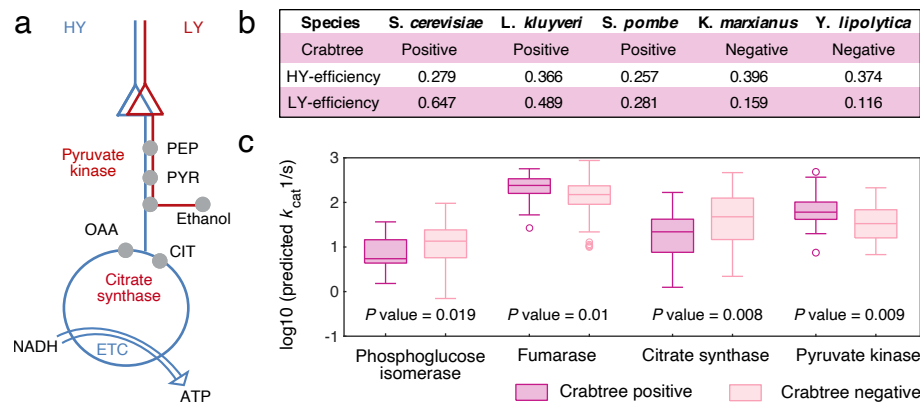| Species | S. cerevisiae | L. kluyveri | S. pombe | K. marxianus | Y. lipolytica |
|---|---|---|---|---|---|
| Crabtree | Positive | Positive | Positive | Negative | Negative |
| HY-efficiency | 0.279 | 0.366 | 0.257 | 0.396 | 0.374 |
| LY-efficiency | 0.647 | 0.489 | 0.281 | 0.159 | 0.116 |

**Figure 6** Explanation of the Crabtree effect by energy metabolism. (a) High-yield (HY) and low-yield (LY) pathway definition. (b) Model-inferred protein efficiency of energy metabolism in several common yeast species. Protein efficiency: ATP produced per protein mass per time (Unit: mmolATP/gProtein/h). (c) Enzymes with significantly different $k_{cat}$ values between Crabtree positive and negative species. A two-sided Wilcoxon rank sum test was used to calculate $P$ value.

44