# Supporting Information: Predicting Enzymatic Reactions with a Molecular Transformer

David Kreutter,<sup>a)</sup> Philippe Schwaller<sup>a), b)</sup> and Jean-Louis Reymond<sup>a)</sup>\*

<sup>*a*)</sup> Department of Chemistry and Biochemistry, University of Bern, Freiestrasse 3, 3012 Bern, Switzerland; <sup>*b*)</sup> IBM Research, Europe, Säumerstrasse 4, 8803 Rüschlikon, Switzerland

e-mail: jean-louis.reymond@dcb.unibe.ch



### Dehydrogenase frequency analysis

Figure S1. Analysis of the dehydrogenase ("DHG") diversity in the entire ENZR dataset.

# TMAP of the ENZR dataset by substrate similarity



**Figure S2**. TMAP of the ENZR dataset analyzed by substrate similarity and color-coded by "-ase" word combinations. Inset: TMAP color-coded by substrate molecular weight.

#### Cofactor importance in the prediction.



Figure S3. Examples of cofactor generator swapped or removed.

# Effect of word on the prediction.





			Database	Prediction 1	Confidence Score	Rank
(3)	Br O N	d-glucose dehydrogenase alcohol dehydrogenase ymr226c from saccharomyces cerevisiae	Br N OH	Br <u>i</u> OH	100.0%	1
(3a)	Br O N	d-glucose dehydrogenase alcohol dehydrogenase ymr226c ➤	Br i OH	Br i OH	100.0%	1
(3b)	Br O N	alcohol dehydrogenase ymr226c	Br <u><u><u></u></u> OH</u>	Br <u><u></u> <u></u> <del>U</del> OH</u> N	100.0%	1
(3c)	Br O N	d-glucose dehydrogenase ►	Br	Br <u><u><u></u></u> <u><u></u> OH</u>N</u>	92.6%	1
(3d)	Br O N	alcohol dehydrogenase	Br E OH	Br 	59.6%	1
(3e)	Br O N	dehydrogenase	Br <u><u><u></u></u> <u><u></u> <u></u> <u></u> <del><u></u></del></u> <del><u></u></del></u> <del><u></u> <del><u></u> <del><u></u></del></del></del>	Br <u><u></u> <u></u> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> <del></del> </u>	41.5%	1



		Database	Prediction 1	Confidence Score (Pred1)	Prediction 2	Confidence Score (Pred2)	Rank
(6)	omega-transaminase from arthrobacter	NH <sub>2</sub> V	NH <sub>2</sub> O	100.0%		0.0%	1
(6a)	omega-transaminase	NH <sub>2</sub>	NH <sub>2</sub>	99.0%	NH <sub>2</sub>	0.0%	2
(6b)	transaminase	NH <sub>2</sub>	NH <sub>2</sub>	58.6%	NH <sub>2</sub>	2.6%	3

F

F

99.9%

99.9%

99.8%

0.0%

0.0%

0.0%

2

2

2

imine reductase s

imine reductase

reductase

(5c)

(5d)

(5e)

			Database	Prediction	Confidence Score	Rank
(7)	ОН	ferredoxin reductase cytochrome p450 monooxygenase from rhodopseudomonas palustris cga009	о он	о он	100.0%	1
(7a)	O O O O H	ferredoxin reductase cytochrome p450 mono oxygenase from rhodopseudomonas palustris	о он	о о о о	99.8%	1
(7b)	ОН	cytochrome p450 mono oxygenase from rhodopseudomonas palustris	о он	о он	99.8%	1
(7c)	ОН	p450 monooxygenase from rhodopseudomonas palustris cga009 ➤	но он	о он	99.8%	1
(7d)	ОН	monooxygenase from rhodopseudomonas palustris ➤	но он	о он	99.9%	1
(7e)	ОН	reductase cytochrome p450 monooxygenase	- С С С С С С С С С С С С С С С С С С С		43.0%	0
(7f)	ОН	rhodopseudomonas palustris cga009 ➤	но он	ОННО	100.0%	1
(7g)	ОН	cytochrome p450	О ОН	ОН	16.4%	0
(7h)	ОН	p450 ➤	о он	HO OH	31.5%	0

-					
		Database	Prediction	Confidence Score	Rank
F	glucose dehydrogenase catalase from micrococcus lysodeikticus cytochrome p450 bm3 mono oxygenase from bacillus megaterium f87a, I188c double mutant	F	F	99.8%	1
	cytochrome p450 bm3 mono oxygenase from bacillus megaterium f87a, I188c double mutant	F ОН	F ОН	99.6%	1
F C	glucose dehydrogenase catalase cytochrome p450 mono oxygenase ────	F ОН	F ОН	100.0%	1
	cytochrome p450 bm3 mono oxygenase	F	F	37.8%	1
F	glucose dehydrogenase catalase	F ОН	F ОН	97.0%	1
	glucose dehydrogenase	F OH	F OH	61.2%	1
F	cytochrome p450	F	F ОН	99.4%	1

(8)	F
(8a)	F
(8b)	F
(8c)	F
(8d)	

(8e)

(8f)

S6



**Figure S4**. Examples of predictions from success examples from figure 4 with a variety of truncated sentences. "Rank" represent the top position prediction containing the correct product.

#### All P450 reactions from the test set.





**Figure S5**. Every reaction from the test set containing "p450" in the sentence correctly predicted by the full sentence model. Reactions sorted by decreasing confidence score.





















fusion protein catalase



**Figure S6**. Every reaction from the test set containing "p450" in the sentence incorrectly predicted by the full sentence model. "rank" showing the rank of the correct prediction assigned by the model, "0" meaning that the model did not predict the correct product within the 5 first predictions. Reactions are sorted by decreasing confidence score.



#### Oxidase wild type (WT) and mutant (M).

M = horse-radish peroxidase choline oxidase from arthrobacter cholorphenolicus, mutant s101a d250g f253r v355t f357r m359r WT = horse-radish peroxidase choline oxidase from arthrobacter cholorphenolicus

**Figure S7**. Reactions using the choline oxidase wild type (WT) and mutant (M) from Heath *et al.*<sup>1</sup> that were assigned to the training set. The numbers in parenthesis correspond to the specific activity of either the mutant or the wild type enzyme express in mU.mg<sup>-1</sup>. (n.t. = not tested).

M = horse-radish peroxidase choline oxidase from arthrobacter cholorphenolicus, mutant s101a d250g f253r v355t f357r m359r WT = horse-radish peroxidase choline oxidase from arthrobacter cholorphenolicus



**Figure S8**. Reactions using the choline oxidase wild type (WT) and mutant (M) from Heath *et al.*<sup>1</sup> that were assigned to the validation set. The numbers in parenthesis correspond to the specific activity of either the mutant or the wild type enzyme express in mU.mg<sup>-1</sup>.

M = horse-radish peroxidase choline oxidase from arthrobacter cholorphenolicus, mutant s101a d250g f253r v355t f357r m359r WT = horse-radish peroxidase choline oxidase from arthrobacter cholorphenolicus



**Figure S9**. Reactions using the choline oxidase wild type (WT) and mutant (M) from Heath *et al.*<sup>1</sup> that were assigned to the test set. All reactions were predicted correctly. The numbers in parenthesis correspond to the specific activity of either the mutant or the wild type enzyme express in mU.mg<sup>-1</sup>.



# Screening of various substrates for the same sentences.

E = d-glucose dehydrogenase alcohol dehydrogenase ymr226c from saccharomyces cerevisiae AD = alcohol dehydrogenase



E = d-glucose dehydrogenase alcohol dehydrogenase ymr226c from saccharomyces cerevisiae AD = alcohol dehydrogenase

**Figure S10**. Various substrates tested on two sentences, a simple "alcohol dehydrogenase" (AD) and the "d-glucose dehydrogenase alcohol dehydrogenase ymr226c from *Saccharomyces cerevisiae*" (E). All substrates were derivatives from **D1** and **D2** which were present in the test set<sup>2</sup> and predicted correctly. Even though products from substrates **D16** and **D19** using enzyme "E" are not chiral, the model gave those chiral centers in the output SMILES ("CC[C@H](O)CC" for **D16**, "O[C@H]1CCCCC1" for **D19**).

# Token frequencies analysis.



Figure S11. Top 40 most frequent tokens from the entire ENZR dataset.



**Figure S12**. Power law distribution of the occurrence frequencies of all tokens in the ENZR sentences sorted by frequency (total of 6,139 tokens).

# References

- R. S. Heath, W. R. Birmingham, M. P. Thompson, A. Taglieber, L. Daviet and N. J. Turner, *ChemBioChem*, 2019, 20, 276–281.
- H. Ankati, D. Zhu, Y. Yang, E. R. Biehl and L. Hua, J. Org. Chem., 2009, 74, 1658– 1662.