

1 **Rational identification of aggregation hotspots based on secondary structure and amino**
2 **acid hydrophobicity**

3

4 Daisuke Matsui^{1,2§}, Shogo Nakano^{1,2§†}, Mohammad Dadashipour^{1,2}, and Yasuhisa Asano^{1,2,*}

5

6 ¹ Biotechnology Research Center and Department of Biotechnology, Toyama Prefectural
7 University, 5180 Kurokawa, Imizu, Toyama 939-0398, Japan.

8 ² Asano Active Enzyme Molecule Project, ERATO, JST, 5180 Kurokawa, Imizu, Toyama 939-
9 0398, Japan

10

11 † Current address: Graduate School of Pharmaceutical and Nutritional Sciences, University of
12 Shizuoka, 52-1 Yada, Suruga-ku, Shizuoka 422-8526, Japan.

13 **Supplemental Information**

14 **Screening of soluble variants of 3-hydroxybutyrate dehydrogenase and tryptophan**
15 **synthase from *Thermus thermophilus* HB8**

16 The expression plasmids of 3-hydroxybutyrate dehydrogenase (*TtHBD*; GenBank accession
17 number BAD70516.1) and tryptophan synthase (*TtTST*; GenBank accession number
18 YP_144360.1) from *T. thermophilus* were from RIKEN BioResource Center (Ibaraki, Japan)¹.
19 These plasmids were transformed into *E. coli* BL21 (DE3), the transformants were grown in LB
20 medium, and the protein expression was induced under the same conditions as described in
21 Materials and Methods. *E. coli* XL-1 Red was used for the random mutagenesis according to the
22 method described in a previous report ². The soluble expressions were obtained from each of the
23 two enzymes in SDS-PAGE, and the following amino acid substitutions were identified: A25E
24 in *TtHBD*, and V136A in *TtTST*. The mutated residues were located in α -helix structures.

25

26 **Expression and characterization of *ChMOX*, *AtADC*, *DmGDH*, *DmODC*, *SuPDH*, and**
27 ***MpLUC***

28 The genes *chmox*, *supdh*, and *mpluc* were obtained from our laboratory stocks, and *atadc*,
29 *dmgdh*, and *dmodc* for the construction of expression plasmids were cloned from the *A. thaliana*
30 or *D. melanogaster* cDNA library. The sequence sizes of *ChMOX*, *AtADC*, *DmGDH*, *DmODC*,
31 *SuPDH* and *MpLUC* are 584, 702, 535, 394, 379, and 208 residues, respectively, and the
32 theoretical molecular weights are 63,800, 76,200, 59,900, 44,200, 41,300, and 22,500,
33 respectively (Fig. S1 and S2). These genes were expressed using the pET, pUC or pCold system
34 in *E. coli* BL21 (DE3) and were cultivated and induced with IPTG. Most of these genes were
35 expressed in the insoluble fractions in SDS-PAGE assays. No activity of *ChMOX*, *AtADC*,
36 *DmGDH*, *DmODC* or *MpLUC* could be detected in the crude extracts (soluble fractions), and

37 the activities of *Su*PDH were very low.

38 *Ch*MOX and *Mp*LUC have no rare codons, but there are six rare codons (Arg200, Arg229,
39 Arg268, Arg358, Arg512, and Arg630) in *At*ADC, three rare codons (Arg18, Arg472, and
40 Arg513) in *Dm*GDH, three rare codons (Arg54, Arg71, and Arg288) in *Dm*ODC, and one rare
41 codon (Arg233) in *Su*PDH. Because the enzyme was produced in the *E. coli* BL21 (DE3)
42 expression system, the gene sequences do not affect the production of the mRNAs for the
43 enzymes. The same results were obtained in the *E. coli* BL21-CodonPlus (DE3)-RIL strain
44 (Stratagene, CA, USA), which contains extra copies of the *E. coli argU*, *ileY*, and *leuW* tRNA
45 genes. It is suggested that expression speed, translational factors, chaperone recognitions, or
46 posttranslational modifications such as glycosylation affect the soluble and active expression of
47 the genes.

48

49 **CD spectra of *Ch*MOX WT and its variants**

50 The CD spectra of the *Ch*MOX, *Mp*LUC, *Dm*ODC, and *At*ADC WT and its variant V455D
51 were measured utilizing a Jasco J-715CD spectrometer (Fig. S3). The buffer contained 10 mM
52 potassium phosphate (pH 7.0) and 50 mM sodium chloride, and 0.1 mg/ml of the enzyme was
53 utilized in the measurement. Far-ultraviolet spectra were recorded from 195 to 280 nm
54 according to the method described in a previous report³.

55

56 **Soluble expression of carbonyl reductase from yeast *Ogataea polymorpha* NBRC 0799**

57 For the expression of the carbonyl reductase (*Og*CRD; GenBank accession number
58 LC176491) gene (*ogcrd*), the already constructed plasmid pET-11a-*ogcrd*, was used in this
59 study. The plasmids were transformed into *E. coli* BL21 (DE3). The transformants were grown
60 in LB medium and the protein expression was induced under the same conditions as described

61 in Materials and Methods. The enzyme activity for the reduction of acetone was assayed at 30°C
62 by measuring the oxidation of NADH to NAD⁺ at pH 6.0. One unit of enzyme activity was
63 defined as the amount of enzyme catalyzing the oxidation of one micromole of NADH per min.

64

65 **Soluble expression of human crystalline aldehyde dehydrogenase and growth hormone**

66 A cDNA of human crystalline aldehyde dehydrogenase (ALDH3A1; GenBank accession
67 number NP_000682.3) was synthesized and amplified using Tks Gflex DNA polymerase and
68 the primers P26 and P27 listed in Table S1. After digestion of pET-15a by *NdeI* and *BamHI*, the
69 amplified ALDH3A1 gene was ligated to pET-15a using an In-Fusion HD Cloning Kit. A cDNA
70 of human growth hormone (GenBank accession number KJ608193, hGH) was synthesized and
71 amplified using Tks Gflex DNA polymerase and the primers P26 and P27 listed in Table S1.

72 After digestion by *NdeI*, the amplified gene was ligated to pET15b with an In-Fusion HD
73 Cloning Kit. The plasmids were transformed into *E. coli* BL21 (DE3). The transformants were
74 grown in LB medium and the protein expression was induced under the same conditions as
75 described in Materials and Methods. LDH3A1 activity was measured by monitoring the
76 production of β-NADH at 340 nm, following the procedure described in previous reports ⁴. The
77 soluble expression levels of hGHs were determined by hGH ELISA kit (Roche, Mannheim,
78 Germany).

79 From the analysis of hGH, seven residues, Leu46, Phe55, Leu82, Leu88, Arg95, Val97, and
80 Leu114, which were on α-helices and had high or low HiSol scores (more than 1.0 or less than -
81 1.0), were selected as aggregation hot-spots (Table 1). As expected, the mutations L46K, F55H,
82 L82R, L88E, R95S, V97E, and L114K enhanced the solubility compared with WT (Fig. 4C).

83

84 **Additional information**

85 We declare that there are no competing financial interests in this work.

86

87 **References**

88

89 1 Yokoyama, S. *et al.* Structural genomics projects in Japan. *Nat Struct Biol* **7 Suppl**, 943-
90 945 (2000).

91 2 Matsui, D. & Asano, Y. Heterologous production of L-lysine ϵ -oxidase by directed
92 evolution using a fusion reporter method. *Biosci Biotechnol Biochem* **79**, 1473-1480
93 (2015).

94 3 Nakano, S. & Asano, Y. Protein evolution analysis of S-hydroxynitrile lyase by complete
95 sequence design utilizing the INTMSAlign software. *Sci. Rep.* **5** (2015).

96 4 Voulgaridou, G.-P., Mantso, T., Chlichlia, K., Panayiotidis, M. I. & Pappa, A. Efficient *E.*
97 *coli* expression strategies for production of soluble human crystallin ALDH3A1. *PloS*
98 *one* **8**, e56582 (2013).

99 **Supplemental Figures**

100

101 **Figure S1. Amino acid sequences of *Ch*MOX, *At*ADC, *Dm*GDH, *Dm*ODC, *Su*PDH, and**
102 ***Mp*LUC (A, C, E, G, I, K, respectively) and their homology model structures generated**
103 **using SWISS-MODEL (B, D, F, H, J, respectively, omitting *Mp*LUC).**

104 The residues in the sequences are highlighted based on the following factors: the α -helix
105 structure in the sequence (underlined), the target α -helix structures (bold font), and the mutated
106 residues (squared). The structure of *Mp*LUC is a novel one that could not be predicted.

107

108 **Figure S2. Depiction of helical wheels of target α -helices of *Su*PDH and *Mp*LUC.** Helical
109 wheels of target α -helices: residues 218-227 (EQAIADIQKL) of *Su*PDH (A), residues 317-340
110 (PARVLAKTENIYTSLLLEV²FHQAEQ) of *Su*PDH (B), residues 76-89 (LEVLIEMEANARKA)
111 of *Mp*LUC (C), and residues 176-201 (SALLKKWLPDRCASFADKIQSEVDNI) of *Mp*LUC
112 (D)). The hydrophobic amino acids, the hydrophilic amino acids, and the target amino acids are
113 represented by black filled circles, white circles and underlined residue numbers, respectively.

114

115 **Figure S3. Comparisons of the CD spectra of WT proteins (filled circle) and variants**
116 **(open circles).**

117 CD spectra of the refolded WT of *Ch*MOX and the V455E variant are shown in A, and the
118 changes in CD at 222 nm measured after heat treatment are shown in B. The CD spectra of the
119 refolded WT of *Mp*LUC and the A177D variant and the changes in CD at 222 nm induced by
120 heat treatment were also measured (C, D). The changes in CD at 222 nm induced by heat
121 treatment were measured in *Dm*ODC WT and K117L (E) and *At*ADC WT and K441L (F).

122

123 **Figure S4. Comparisons of the thermal stability of *MpLUC* WT proteins (filled circle), the**
124 **I80K, and the A177D variant (filled triangle).**

125 The each luminescence was measured after heat treatment for 30 min.

126

127 **Supplemental Tables**

128

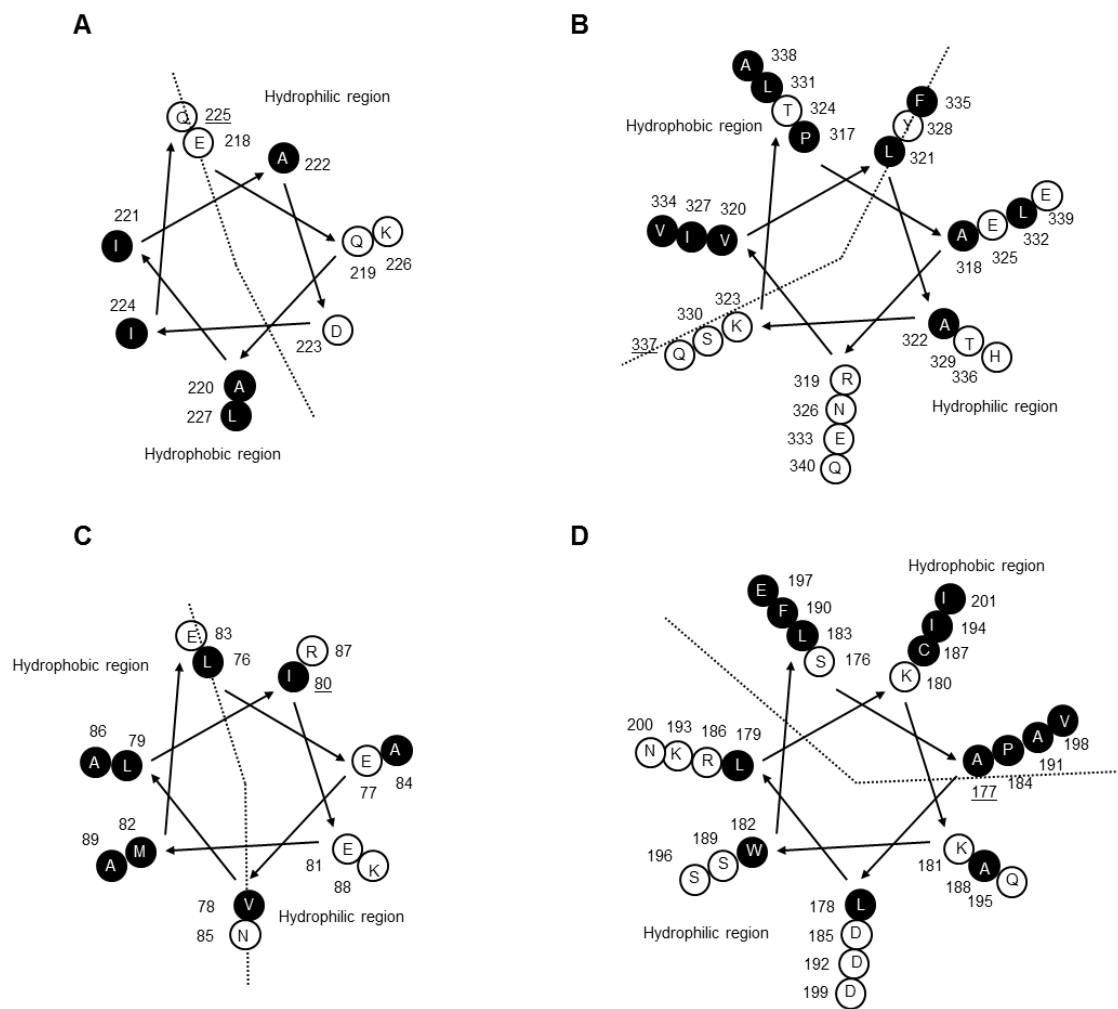
129 **Table S1.** Designed oligonucleotides used to perform random mutagenesis and site-directed
130 mutagenesis.

131

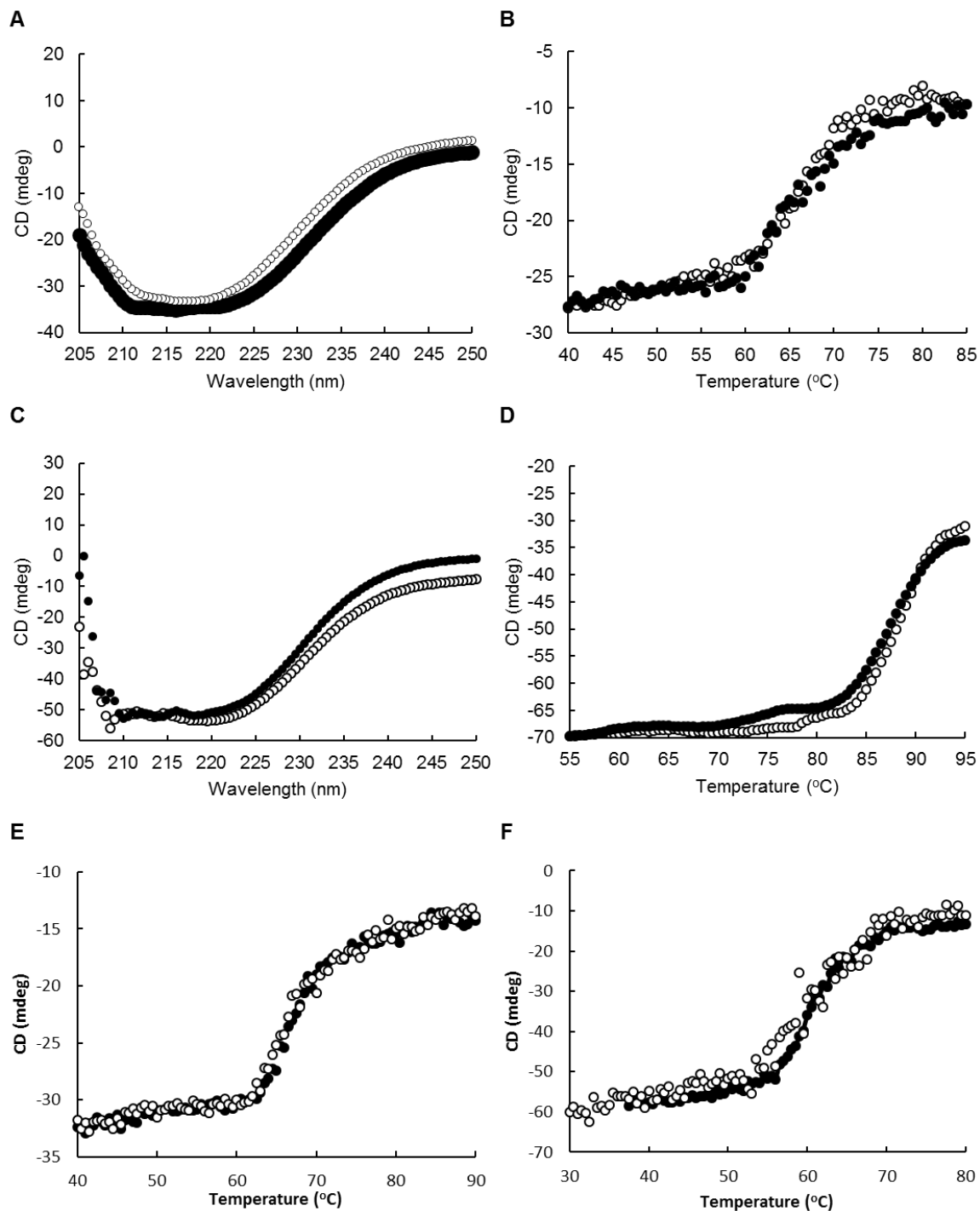
132 **Table S2.** Enzyme activity measurement of the Val444 and Val455 variants of *ChMOX* and the
133 Leu435 and Lys441 variants of *AtADC* generated by saturation mutagenesis.

134

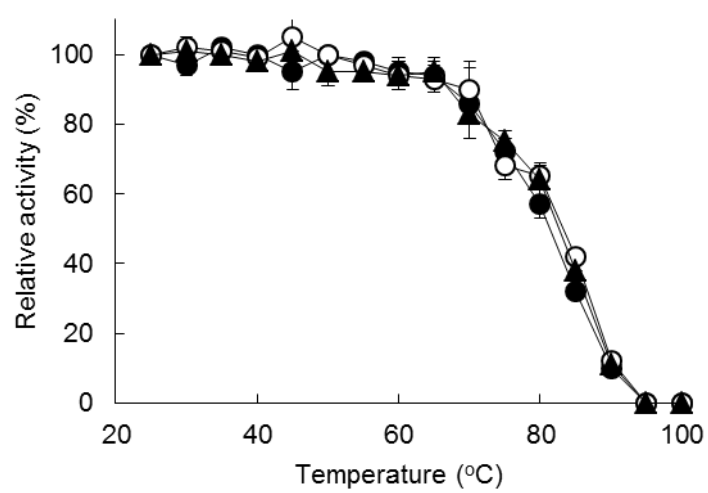
135 **Table S3.** Run time parameters of INTMSAlign for *ChMOX*, *AtADC*, *DmGDH*, *DmODC*,
136 *SuPDH* and *MpLUC*.



140



141



143 **Table S1. Oligonucleotides used in this study.**

Primer	Comments ^a
P1	5'-GCGCATATGAATCCGACCGAAAACAAAGATG-3', Amplification of <i>MpLUC</i> gene
P2	5'-CGCCTCGAGTTAACGATCGCCTGCCAGGCCTTT-3', Amplification of <i>MpLUC</i> gene
P3	5'-GGAGATATACATATGCCTGCTCTAGCTTTTGTGGA-3', Amplification of <i>AtADC</i> gene
P4	5'-TTAGCAGCCGGATCCTCAACCGAAATAAGACCAAT-3', Amplification of <i>AtADC</i> gene
P5	5'-AGAAGGAGATATACATATGCTTCGTTATACGGCACGGATT-3', Amplification of <i>DmGDH</i> gene
P6	5'-TTTGTTAGCAGCCGGATCCCTACTGTTGGGAAATTCCGGAGC-3', Amplification of <i>DmGDH</i> gene
P7	5'-AGAAGGAGATATACATATGGCGGCCGCTACCCCTGAAAT-3', Amplification of <i>DmODC</i> gene
P8	5'-TTTGTTAGCAGCCGGATCCTCATATAGCTTGGAAGTACAGGG-3', Amplification of <i>DmODC</i> gene
P9	5'-ACCGACCCTGGAAGANNSGACATTGATACGAT-3', Saturation site-directed mutagenesis at Val444 of <i>ChMOX</i>
P10	5'-ATCGTATCAATGTCSNNTCTTCCAGGGTCGGT-3', Saturation site-directed mutagenesis at Val444 of <i>ChMOX</i>
P11	5'-GTTTCGAGGCGTACACNNSGCTCTTAACCTTTGGA-3', Saturation site-directed mutagenesis at Val455 of <i>ChMOX</i>
P12	5'-TCCAAAGTTAAGAGCSNNGTGTACGCCTCGAAC-3', Saturation site-directed mutagenesis at Val455 of <i>ChMOX</i>
P13	5'-CGTGAAAGCTGCTTGNNSTATGTTGATCAGCTG-3', Saturation site-directed mutagenesis at Leu435 of <i>AtADC</i>
P14	5'-CAGCTGATCAACATASNNCAAGCAGCTTTCACG-3', Saturation site-directed mutagenesis at Leu435 of <i>AtADC</i>
P15	5'-ATGTTGATCAGCTGNNSCAGAGATGTGTTGAAG-3', Saturation site-directed mutagenesis at Lys441 of <i>AtADC</i>
P16	5'-CTTCAACACATCTCTGSNNCAGCTGATCAACAT-3', Saturation site-directed mutagenesis at Lys441 of <i>AtADC</i>
P17	5'-GTGGGCVTCCCGGTCGAATATGGTGGCGGT-3', Amino acid substitution at Lys148 to Ile, Val or Leu of <i>SuPDH</i>
P18	5'-GATATTGTGAAGCTCGGTGGAAGCGCTGTC-3', Amino acid substitution at Gln225 to Val of <i>SuPDH</i>
P19	5'-AGTCAGGYAGCAGATGTTTTTGTTCCTTGT-3', Amino acid substitution at Gln243 to Val or Ala of <i>SuPDH</i>
P20	5'-TTCCATATCGCAGAACAGGATCATATGACA-3', Amino acid substitution at Gln337 to Ile of <i>SuPDH</i>
P21	5'-AACCGACCGRTATGGAATTTTCATCAGTAA-3', Amino acid substitution at Lys374 to Ile or Val of <i>SuPDH</i>

- P22 5'-CCGCTGGAAGTTCTGAAAGAAATGGAAGCAAAT-3', Amino acid substitution at Ile80 to Lys of *MpLUC*
- P23 5'-ATTTGCTTCCATTTCTTTCAGAACTTCCAGCGG-3', Amino acid substitution at Ile80 to Lys of *MpLUC*
- P24 5'-AATGTTAAATGTAGCGATCTGCTGAAAAAATGG-3', Amino acid substitution at Ala177 to Asp of *MpLUC*
- P25 5'-CCATTTTTTCAGCAGATCGCTACATTTAACATT-3', Amino acid substitution at Ala177 to Asp of *MpLUC*
- P26 5'-GCTAATTTTGCTCATATGGCTGCCGCGCGGCACCA-3', Amplification of ALDH3A gene
- P27 5'-ATGACCCAGCATTAAAGGATCCGGCTGCTAACAAAG-3', Amplification of ALDH3A gene
- P28 5'-CGCGGCAGCCATATGTTTCCGACCATTCCGCTGAGCC-3', Amplification of hGH gene
- P29 5'-TTAGCAGCCGGATCCTTAAAAACACAGCTACCTTCAAC-3', Amplification of hGH gene
-

144 ^a Mutation sites are underlined.

145 **Table S2. Saturation mutagenesis at Val444 and at Val455 of *ChMOX* and at Leu435 and at Lys441 of *AtADC***

	<i>ChMOX</i> Val444			<i>ChMOX</i> Val455			<i>AtADC</i> Leu435			<i>AtADC</i> Lys441		
	Total activity (U/ml)	Soluble protein (mg/ml) ^a	Total specific activity (U/mg)	Total activity (U/ml)	Soluble protein (mg/ml)	Total specific activity (U/mg)	Total activity (U/ml)	Soluble protein (mg/ml)	Total specific activity (U/mg)	Total activity (U/ml)	Soluble protein (mg/ml)	Total specific activity (U/mg)
Ile	ND	1.4	ND	ND	1.7	ND	ND	1.2	ND	ND	1.5	ND
Val	ND	1.7	ND	ND	1.6	ND	ND	1.2	ND	0.005	1.4	0.0036
Leu	0.028	1.2	0.024	0.022	1.4	0.016	ND	1.4	ND	0.043	1.3	0.0331
Phe	ND	0.6	ND	ND	1.3	ND	ND	1.8	ND	ND	1.3	ND
Cys	ND	1.7	ND	ND	1.3	ND	ND	1.3	ND	ND	1.1	ND
Met	0.010	1.8	0.006	0.010	1.2	0.008	ND	1.1	ND	ND	1.4	ND
Ala	0.005	0.8	0.006	0.005	1.8	0.003	ND	1.2	ND	0.009	1.8	0.0050
Gly	ND	1.4	ND	ND	1.5	ND	ND	1.3	ND	ND	1.2	ND
Thr	0.005	1.6	0.003	ND	1.9	ND	ND	1.2	ND	ND	1.2	ND
Ser	0.020	1.4	0.014	ND	1.8	ND	ND	1.8	ND	ND	1.5	ND
Trp	ND	1.8	ND	ND	1.2	ND	ND	1.2	ND	ND	1.3	ND
Tyr	0.037	1.3	0.028	ND	1.3	ND	ND	1.4	ND	0.005	1.2	0.0042
Pro	ND	1.5	ND	ND	1.8	ND	ND	1.1	ND	ND	1.8	ND
His	0.060	0.9	0.067	ND	1.5	ND	0.036	1.6	0.02	ND	1.2	ND
Glu	0.055	1.2	0.046	0.040	1.9	0.021	0.005	1.5	0.00	ND	1.3	ND
Gln	0.050	1.4	0.036	0.060	1.8	0.033	0.012	1.2	0.01	ND	1.2	ND
Asp	0.050	1.3	0.038	0.024	1.4	0.017	0.002	1.4	0.00	ND	1.2	ND

Asn	0.050	1.2	0.042	0.022	1.5	0.015	0.012	1.2	0.01	ND	1.3	ND
Lys	0.012	1.1	0.011	0.042	1.3	0.032	0.001	1.4	0.00	ND	1.1	ND
Arg	0.055	1.5	0.037	0.040	1.5	0.027	ND	1.1	ND	ND	1.5	ND

U/ml of cell-free extract prepared from a 3-ml LB culture in triplicate; ND, not determined

^aConcentration of soluble proteins means concentration of soluble fraction of crude enzyme solution.

147 **Table S3. INTMSAlign parameters for *ChMOX*, *AtADC*, *DmGDH*, *DmODC*, *SuPDH* and *MpLUC***

	<i>ChMOX</i>	<i>AtADC</i>	<i>DmGDH</i>	<i>DmODC</i>	<i>MeHNL</i>	<i>SuPDH</i>	<i>MpLUC</i>
Types of Blast	Blastp	Blastp	Blastp	Blastp	Blastp	Blastp	Blastp
Database	Non redundant	Non redundant	Non redundant	Non redundant	Non redundant	Non redundant	Non redundant
Sequence of target protein (STP)	<i>ChMOX</i> (GenBank ID; not registered)	<i>AtADC</i> (GenBank ID:15227223)	<i>DmGDH</i> (GenBank ID:24649283)	<i>DmODC</i> (GenBank ID:24586472)	<i>MeHNL</i> (GenBank ID: 55469815)	<i>SuPDH</i> (GenBank ID: 1842144)	<i>MpLUC</i> (GenBank ID; not registered)
Total number of sequences in the library	5,000	5,000	5,000	5,000	825	158	37
N_{pick}	8	8	8	8	8	8	8
N_{trial}	500	500	500	500	1000	1,000	500

148