

ESM-Ezy: A deep learning strategy for the mining of novel multicopper oxidases with superior properties

Hui Qian^{1,2,#}, Yuxuan Wang^{1,2#}, Xibin Zhou^{1#}, Tao Gu^{1,2}, Hui Wang¹, Hao Lyu¹, Zhikai Li¹, Xiuxu Li^{1,2}, Huan Zhou³, Chengchen Guo¹, Fajie Yuan^{1,*} and Yajie Wang^{1,2,3,4,5*}

¹School of Engineering, Westlake University, Hangzhou 310014, Zhejiang, China

²The Center for Synthetic Biology and Integrated Bioengineering, Westlake University, Hangzhou 310014, Zhejiang, China

³Westlake Laboratory of Life Sciences and Biomedicine, Xihu District, Hangzhou 310024, Zhejiang Province, China

⁴School of life Science, Westlake University, Hangzhou 310014, Zhejiang, China

⁵Muyuan laboratory, Zhengzhou, Henan , China

#These authors contributed equally

*Corresponding authors: wangyajie@westlake.edu.cn, yuanfajie@westlake.edu.cn

Table of contents

Supplementary Tables

- Table S1 The plasmids and strains used in this study.
- Table S2 All multicopper oxidase used in this study.
- Table S3 The sequences of all multicopper oxidases mined in this study.
- Table S4 The evaluation for different positive threshold.
- Table S5 The five fold evaluation results of the MCOs model.
- Table S6 Characterization of the multicopper oxidase.
- Table S7 The enzymes with remote distances tested in this Study.
- Table S8 The sequences of enzymes with remote distances tested in this Study.
- Table S9 Refinement statistics obtained for the final Sulfur model.
- Table S10 Refinement statistics obtained for the final Bfre model.
- Table S11 Comparison of protein secondary structures and residue interactions between Sulfur and Eclac.
- Table S12 The five fold evaluation results of the L-Asparaginases model
- Table S13 Characterization of the L-Asparaginases.

Supplementary Figures

- Figure S1 Changes in accuracy during the model fine-tuning process.
- Figure S2 The embeddings of the multicopper oxidases in the ESM-1b model.
- Figure S3 The phylogenetic tree of multicopper oxidases.
- Figure S4 The sequence similarity network (SSN) of Eclac (query enzyme, QE).
- Figure S5 The cluster analysis of all multicopper oxidases characterized in this research
- Figure S6 The amino acids coordinated with metal atoms in the active center.
- Figure S7 Superimposed alignments of the Eclac AlphaFold structure with its corresponding crystal structures.
- Figure S8 Root mean square fluctuation (RMSF) of C α atoms per residue of (A) SulfurA and (B) Eclac at 37 °C
- Figure S9 Fo-Fc Omit and Electron Density Maps of the Crystal Structures.
- Figure S10 SDS-PAGE analysis of enzymes utilized in this study.

Development of a deep learning model to identify MCOs with low sequence similarities

Table S1 The plasmids and strains used in this study

Strains and plasmid	Description
<i>E. coli</i> Top10	plasmid construction
<i>E. coli</i> BL21(DE3)	the host for protein expression
pET28a	the vector for protein expression
pET28a-sumo	the vector for protein expression

Table S2 All multicopper oxidases used in this study

Name	Entry	Organism
Eclac series		
Eclac(query)	P36649	<i>Escherichia coli</i> K-12
Sulfur	A0A2K2VPI4	<i>Sulfurimonas sp</i>
Salvi	A0A2N9X7P6	<i>Snodgrassella alvi</i>
Pph	A0A238JCR0	<i>Pelagimonas phthalicica</i>
Ypes	Q8ZBK0	<i>Yersinia pestis</i> biovar <i>Microtus str. 91001</i>
Psoli	A0A328AKG1	<i>Phenylobacterium soli</i>
Faest	I0KHI7	<i>Fibrella aestuarina</i>
Mint	A0A4R6VBE7	<i>Mesocricetibacter intestinalis</i>
HR03 series		
HR03(query)	ACM46021	<i>Bacillus sp.</i> HR03
Aory	A0A5C7ESL7	<i>Azospira oryzae</i>
Talbi	A0A5B9EN11	<i>Terriglobus albidus</i>
Bfre	VEF49110	<i>Bacillus freudenreichii</i>
Tocean	D9RY72	<i>Thermosediminibacter oceanii</i>
BMS3	A0A2H6K631	bacterium BMS3Abn01
Acidi	A0A0S8GAJ7	<i>Acidithiobacillales</i> bacterium SM23_46
DSM13 series		
DSM13(query)	Q65MU7	<i>Bacillus licheniformis</i> DSM13
Bcece	M7PAR1	<i>Bhargavaea cecembensis</i> DSE10
Scla	D5SLU3	<i>Streptomyces clavuligerus</i>
Arose	A0A7W9W9Q1	<i>Armatimonas rosea</i>
Cmeth	PWB55314	<i>Candidatus Methanoperedenaceae archaeon</i>
Cohnella	A0A433RH29	<i>Cohnella sp.</i> AR92

Table S3 The sequences of all multicopper oxidases mined in this study

Sulfur
<u>ggatcc</u> TGTGAAAAGGGTAGTGATAAAAGCGTTGAACGGATAAAAAAGTGTGAAAC CCAGAAAAGCCGCCAGGAACCGCCGAAAAATATTGCCAATAATAAATTGAGCCTTCACC CAGAAACTGAAAATTCCGAAAGAAATTGATTCGAACACGTTGCAAAAGCAAATTAAATG CCCAGAAAAGCCTGAGCGCACTGTATAAAGAAAAAAACCGATATCCTGACCTTCAGG GCGATCTGCCGAATCCGACCATTCGCATTAAAATGGCGATGATTTGAACGGATTTAC AATAGCCTGGAAAACCGACCATTATTCAATTGGCATGGCCTGCTGGTGCCGAGCAATGG ATGGCCATCCGAAAGATGCCATTGCCACCCAGATGCTGAAAGAAATATCGCTATAAAGTGAAT CAGCGTGCCTGGTACCTTTGGTATCATACCCATCCGCATGGCCGCACCGGTGAAGAAATT TTATGGCCTGGCGGTCTGTATATTATTGAAGATGATAATGAGAAGGCGCTGAATCTGCCGA GCGCGAATTGAACGCCGCTGATTATTCAAGGATCGTCGTTTGATAAAGAAGGCGATCTG ATTTATAAGGAAACCCCGCAGGATAATAATGGCGTTGGGTGATGTTATGGTGAATAG TACCGTGCATCCGTATAAAAATGTGAAAATACCAAGTACCGCCTGCGCATTCTGAATGGTA GTAGCGCCCGTACCTATAAACTGGCATTGAAGGTATTGAAGATTATGCTGATCGCAGTGGCGAAC GATGGCGGTCTGCTGGAAGAACCGATTATTGTTAAAGATATTGATCGCAGTGGCGAAC GTATTGATATTATTGTGGATTAAAGGACAAGAACGGTGGTGAAAGTGTACCCCTGAAAACC CTGGGTTAAAGAACAAATAATTGTGACCAACCCGGCATATCCGGATAGCGGTGCAA AAATGGATATTATGCGTTAAAGTGACCGAACTGAGTACCCAGAATAGTCAGATTCCGAAA AAACTGAGTACCATGCCAAATGAAAGCCAGTGATGCAAGTAAAAGCCGTACCAATTACCA TGGAAATTATTGAAGGTGGCGTTGGACCCCTGAATAAAAACCGTATGATATGCATCGTGT GATGAAAAGTTAAACTGGTAGTACCGAAATTGGAAATTAAAATAGCGCACATATGG CACATCCGTTCATATGCATGGCGTCATTCAGGTTCTGGAACGTACCAAGCAGCATTGATT TTCCGACCGATAAAGGTTGGAAAGATACCGTTCTGGTATGCCCTGGAAAGTGTCTGTATT ATTGTTAAGTTACCATCCGGTCTGTTCATCATTGCCATTCTGGAACATGAAGAT CATAGTATGATGCCAATTCTGGTTGAATAA <u>aagctt</u>
GSCEKGSDKSVELDKSDETQKAAQEPPKNIANNKFEPFTQKLKIPKEIDFEHVAKAKFN AQKSLSALYKEKKTDILTFQGDLPNPTIRIKNGDDFELDFTNSLEKPTIIHWHLVPEAMDGH KDAIATQMLKEYRYKVNRAGTFWYHTHPHGRTEEIYYLAGLYIIEDDNEKALNLPSGEFE LPLIIQDRRFDKEGDLIYKETPDNNNGVLGVDMVNSTVHPYKNVKNTKYRLRILNGSSARTY KLAPEGIEDFMLIGTDGLLLEPIIVKDILIAVAERIDIIVDFDKKVGESVTLKTLGFKEANNFV TNPAYPDGAKMDIMRFKVTELSTQNSQIPKKLSTIAKMKASDASKSRTITMEIIEGGVWTLNK KPYDMHRVDEKVKLGSTEIWEIKNSAHMAHPFHMHGTVHFQVLERTSSIDFPTDKGWKDTVLV MPLESVRIIVKFTIPGLFVHHCHILEHEDHSMMANFLVE
Salvi
<u>ggatcc</u> CTGAATAACCTGAATCTGCCGGTGCTGGCATTCCGCCGCTGCTGGATAACCGTA AAGATAATCGTATTCACTGACCAATTCACTGACGAGGTACCAAGCCGTTTGGCAATTATCAGACC CAGACCTGGGTTATAATGGTGCCTGCTGGCCCGCAATTCTGATTTCGCGGTAAACC GGTCAGATTAAACTGCAGAATCAGCTGACCGAAGCCACCACCGTCATTGGCATGGCCTG GAAGTTCCGGCGAAGTGGATGGCGGTCCGCAGGCATGCATTCCGCCGGTCAGGCCGC GTTATTGAATTACCAAGCAGCCAGCCGGCAGCAACCTGCTGGTTCATCCGCATCAGCATGG TCAGACCGGCCAGGGTCAATGGGTCTGGCCGGTCTGCTGCTGGAAAGATACCGTG AGTACCGGTCTGCCATTCCGAAATTGGGGCTGGATGATATTCCGCTGATTGTCAGGA TAAAAAAATTAAATACCGCGGGTCAGATTGATTACAGCTGGATGTTAGCGCCGAATTG

GTTGGTTGGTATACCTGCTGTCAATGCCAGATTATCCGCAGCATAATGCCCGC GGCTGGCTCGCTCGCTTAATGGCTGAATGCCGCACCCCTGAATTGATGTAG CGATGCCGTCCGCTGTATGTGATTGCCAGTGATGGTGGCTTCTGCCGCAGCCGGTACCG TTAGTGGCCTGAGCGTGCTGCCGGGAACGTTGAAGTTATGATTGATACCAGTAATGGT AAAGCATTGATCTGCTGACCCGCCGGTAAACAGATGGGCATGAATCTGCCGCCCTTG ATCAGCCGTCCGATTGTGCGCGTGAGTGTCTGAGCAGCCGAATGATGCAAGCCTGCC GGATCAGCTGACCAGCAGTCTGCCGGCACTGCCGGATTAAGTCAGCTGGTCATCGTCAT TTTCAGCTGAGTATGGACCCCTGAACACTGGATAGCATGGGTATGCAGGCCCTGCAGCAGCTGG CAGGCAAAGCGCAATGCCGATGGAACATATGAGCATGCATCAGAGCCATATGAAATGAT GCATAATATGCATACCGATAAAACGCCCTGGATCTGATGCATGCAAATCGTATTAATGCCCT GACCTTGATATGAATAAACCGGCCCTTGAGCAAGCAAAGGAAATATGAATGCTGGAA ATTAGTGGCCTGGCGATATGATGCTGCATCCGTTCATATTGACGCCAGTTCTGATT CTGAATGAAAATGGTCAGCCGCCAGCCCATCGCAGCGGTTGAAAGATACCGTTCGC GTGGAAGGTGGTGTGAGTAAAGTGTGGTAAATTATCATAGCGCCCCGGCAGCTTC CGTATATGGCACATTGTCATATTGGAACATGAAGATACCGGTATGATGCTGGCTTACCG TGTAA <u>aagctt</u>
LNTLNLPVLAIPPLDTGKDNRQLTIQQQTSRGNYQTWTGYNGALLPAIRFRGKPV QIKLQNQLTEATTVWHGLEVPGEVDGGPQACIPAGQARVIEFTSSQPAATCWPHQHGQTG RQVAMGLAGLLLEDVTGLPIPKNWGLDDIPLIVQDKKFNTAGQIDYQLDVMMSAIGWFGD TLLCNGQIYPQHNAPRGWLRLRLNGCNARTLNFAACSDGRPLYVIASDGFLQPVTVSLSV LPGERFEVMIDTSNGKAFLLTPVKQMGMNLPPFDQPCPIRVSVLSSRNDAASLPDQLTSSL ALPALSQVLHRHFQLSMDELDMSGMQLQQLAGKSAMPMEHMSMHQSHMEMMHNMHTD KNALDLMHANRINGLTDFMNKPAAASKGKYECWEISGLGDMLHPFIHGTQFRILNENGQ PPAAHRSGWKDTVRVEGGVSKVLVFNHSAPASFPYMAHCHILEHEDTGMLGFTV
Pph
<u>ggatcc</u> CAGAGTCCCCCGTCCTCTGCCGCTGATTCCGCTGACCGATCTGACCGGTGGT ATTGAAGGTCGATTAGTCTGCCCTGAATAAACGCCAGCCATGATTGACCTGGCACC GGTC GAGTGAACACCTTGGTATTAAATCATGATTACCTGGTCCGGTCTGCGTGTGAAACAGGGTC AGACCCCTGCCGTTGATGTTGAAATAATATTGGCAGGAAACAGCACCCTGCATTGGCATGGT CTGCATATTCCGGCGATGTTGATGGTGGCCCGATCAGGAAATTGCACATGGCGAACCT GGAGTCGGATGTCGATTGTTGAGCATGCAAGTATGAATTGGTTCATAGTCATATGCAT GGCCGCACCGCACGCCAGACCTATAGTGGCCTGGCCGGTTCTGCTGGTGAAGATGATG CAAGTCTGAGCGCAGATCTGCCAAAACCTATGGTGTGATGATTTCACCTGATTCTGCAG GATAAAATGTTGATGGCAGTGGCAAATGACCTATGCACTGACCGCCGAAGTGGTGAAG ATGGCTTGAAAGGTGATACCCGTACCGTTAATGGTCAATTGCACTGACCGGTGCAAGAGTGT GCCGACCGGGCTGGTCCGCTGCGTATTCTGAATGCCGTAAATGCACTGACGCTTCTGGAAC GCATGCAGAGCGGCCCGATTACCGTGATTGCAAGTGTGGCGCTTCTGGCAGCCCGGT GGAAGCCGCCAGCATTCTGATGAGTCCGGCGAACGTTATGAAGTGTGGTGAATGAGC ACCGTGGATAGCAATGAACTGACCGTGAAATGGATGGTGGCGCAGCTTTGCCAGTC TGTGGTGGTAAATGCAAGCAACCACCGTGCTGACCCGTGACCGCTGCCGAAGCCGGTT TGATGGCGCAATGCCGCCGTCTGGCAAATCTGGCACC GCCAAAGCCAGTGAAGCCAG TGTTACCCGCAGTTTCAGCTGGAAATGGATGTGGCGCCGATCTGCCGCCCTGGCAGCA TCATGGGATAATTGTTGTGGTGTGAGGTGCAATGCCATTAAATGCCAGCCGATGAAAT GGATCGTATTGATGAAGTGGTGCAGTAAAGGTAATACCGAAATTGGCGCGTGAGCGTGGAT

GATCAGCTGCATCCGTTCATATTCA~~TGGCTGCAGCTTCGTATTCTGAGCCAGAATGGCGC~~
AGCACCGCCGGCATATGCCGCCGGTGGAAAGATATGGTCATGTGAAGAAGGCTGGAGC
GAAGTCTGGTCGCTTGATTATGAAGCACCAGAAGATGCACCGTATATGTATCATTGCCA
TATTCTGGAACATGAAGATTGCGGTATGATGGTCAGTTACCGTGGTTAAaagctt

QSAPRPLPLIPLTDLTGGIEGRISLALNKASHDFGTGAQSETFGINHDYLGPVLRVKQGQTL
PFDVVNNIGETSTLHWHLHIPGDVDDGPHQEIAHGATWSPDVPIVQHASMNWFHSHMHGRT
ARQTYSGLAGVLLVEDDASLADLPKTYGVDDFTLILQDKMFDSKGKMTYALTAEVFEDGFE
GDTLVNGAIAPVAQSVPTGLVRLRILNACNARFLELSMQSGPITVIASDGGFLAAPVEAASILM
SPGERYEVLVDMSTVDSNELTVNMDGGGSFFASLFGGNAATTVLTTRRAEAGFDGAMPRL
ANLAPPKASEASVTRSFQLEMDVGADLAALAASWDNFCDAGAMAINGQPMKMDRIDEVV
RKGNTEIWRVSVDDQLHPFIHGCSFRILSQNGAAPPAYAAGWKDMVHVEEGWSEVLVRFDY
EAPKDAPYMYHCHILEHEDCGMMGQFTVG

Ypes

ggatccGCAGATTAGTCCGCTGCCGATTCCGCCGCTGCTGCAGCCTGATGCAAATGGTA
AAATTAATCTGAATATCCAGACCGGTAGCGTGGTGTGGCTGCCGAGCACCGCAACCCAGAC
CTGGGGTTATAATGGTAATCTGCTGGGCCCGCAATTGCCTGCAGCGTGGCAAAGCAGTG
ACCATTGATATTACCAATGCCCTGCCGGAAGCAACCACCGTTATTGGCATGGCCTGGAAAT
TCCGGGTGAAGTTGATGGCGGTCCGCAGGCCCTGATTCA~~GAGCCGGTGCAGGT~~AAAACGTCAGGT
TACCTTGCA~~GCTGAACAGCCGGCAGCCACCTGCTGGTT~~CATCCGCATA~~ACCCATAGTAAAAA~~
CCGGCCATCAGGTTGAATGGCCTGGGTGGCCTGGTTCTGATTGATGATGATAGTGA~~TCGA~~
AACCTGCCCTGCCGAAACAGTGGGGCGTGGATGATATTCCGGTTATTCTGCAGGATAAA
CTGCTGGATAAACATGGTCAGGTTGATTATCAGCTGGATGTGATGACCGCAGCAGTGGCTG
GTTGGTGA~~CTGATGCTGACCAATGGCGTCCGTATCCGCAGCAGATTACCCCGTGGCT~~
GGGTGCGCCTCGTACTGAATGGCTGTAATGCACGTAGTCTGAATCTGGCACTGAGTGAT
GGTCGTCGATGTATGTTATTGCAAGTGTGATGGCGGTCTGCTGGCGAACCGTGGTTGTGC
GTGA~~ACTGCCGATTCTGATGGCGAACGTTGAAGTTCTGGTGGATACCGCGATGGCCA~~
GAGCCTGGATCTGGTACCCCTGCCGGT~~GACCCAGATGGCATGACCCCTGGCCCCGTTGAT~~
CAGCCGCTGCCGGTCTCGTATTAGCCGAGTCTGGCAATTGGTAGCCAGGTTCTGCCGG
AAAGCCTGGTTATTCCGA~~ACTGGCCATGTGACCCGGTTCAGGAACGCTGGTTCA~~
GCTGATGATGGACCC~~TAAACTGGATATGCTGGTATGCAGGCCCTGGTGGCACGCTATGGTA~~
TGAAAGCCATGGCAGGCATGAATATGAATCATGGT~~GATATGGTGCCATGGATCATGGTAAT~~
CGTCCGGATATGAGTCAGGGCAA~~AAATGAAAGGCATGGATCATGGCACCATGAATGGTGC~~
CGGCATTAA~~TTAGTCATGCAAATCGCATTACGGTAAAGCATTAGTATGACCGAACCG~~
GCATTGATGCCAACAGGGCAA~~ATGGACCATTAGTGGCGAACCGGATATGA~~
TGCTGCATCCGTTCATGTCATGGTACCCAGTTCTGATTCTGACCGAAA~~ATGGTAAACCG~~
CCGGCCGAACATCGTCGCCGGTGGAAAGATATTGTCGCGTTGAAGGCGCCCGAGTGA~~A~~
ATTCTGGTGC~~GT~~TTAATTATCTGGCACCGGCAAGTACCCGTATATGGCACATTGCCATCTG
CTGGAACATGAAGATACCGGATGCTGGTTTACCGTTAGCGCATAaagctt

ADFSPLPIPPLLQPDANGKINLNIQTGSVWLPSTATQTWGYNGNLLGPAIRLQRGKAVTI
DITNALPEATTVWHGLEIPGEVDGGPQALIQPGAKRQVTFAVEQPAATCWFPHTHSKTGHQ
VAMGLGGVLIDDSDSETLPLPKQWGVDDIPVILQDKLLDKHGQVDYQLDVMTAAVGWFGD
RMLTNGVPYPQQITPRGVWVRLRLNGCNARSLNLALSDGRPMYVIASDGGLAEPVVVRELPI
LMGERFEVLVDTRDGQS~~LDLVTLPVTQMGMT~~APFDQPLPVLR~~IQPSLAIGSQVLPESLV~~VI~~EL~~
ADVTGVQERWFQLMMDPKLDM~~LGMQALVARYGMKAMAGMNMNHGDMGAMDHGNRPD~~

MSQGKMKGMHDHGTMNGAPAFNFSHANRINGKAFSMTEPAFDAKQGKYEKWTISGEGDMML
HPFHVHGTQFRILTENGKPPAEHRRGWKDIVRVEGARSEILVRFNYLAPASTPYMAHCHLLEHE
DTGMLMLGFTVSA

Psoli

ggatccCAACCTGGAAGGTTACTACACCACTACCCGTTCCGCCACTGATCGATGCGCAA
GCACAGGGTAATAGCATTGTTGGCGCAAACCTGGTCGCCATGCGTTGTTGAAGGTC
GTCCGGTTGCGACCTACGGCTACAGCGGTCCGGTCTTGGTCCGACGATCCGTATGCAGTC
CGGTGCGCAAGTTGACGTTGCCATTGAGAACCGTATCGACACGGATACCATTGTACATTGG
CATGGTTCTGATTCCGTCGAGCGCGACGGTGGTCCGACGATCACATTGCTCCAGGTG
AAACCTGGCGTCGCGTGCTGCCGGTTCAGCAGCCGAAACTACTGCGTGGTATCATCCGCA
CCCGCATCGTACACCGGTCGTCAGGTTATTCGGCATGGCGGGTTGGTGGTATCGAG
GACGGTTCAGGCGCAAGATTAGACCTGCCACGTACCTATGGTGTGACGACCTGCCTCTGG
TCCTGCAAGACCGTCTGTTGATGGCAACGGTGTATCTGATCTACCCGGCCAATCCGATGACC
GTAATGCAGGGCGCGTGGCAACACCATTATCGTAACGGCGCGATGCCCGCATCGCCC
GTGTTCCGGCGGGTCTCGTGCCTTCGATTGCTGAATGGCTCGAACGCTCGTAATCTGAT
CTTGGCTCGACGATGGCCGCCCGTTGCACGTGATTGCTAGCGACGGTGGTACCTGGGCT
CCCCGGTGAGCATGACCCGCTGCTGATCGCACCGAGCGAGCGCTCGAGATCCTGGTGA
CTTAGCGATGGTAGGCCGGGGTTCTGCAGACCGGGCCGGACCCGTTCCGCCTATGATG
GGCATGATGATGGGAGAGGCCAATCTGCTGGGTGATATTATGAGCTCGCACCGGACC
GTGATCGCCCAGGTGCGAAGCGAAAATCCGAGCACCCCTCGGGACATGCCGGCGGG
CACCGAGTGTCAACTGCTGCGTCGTTGCCCTGAACGATATGGCGGAATGATGGG
CGGTATGGCGGCATGATGGTCGTGGTATGATGGCGGTATGATGGCGGCATGATG
GGTGGGGGATGATGGTGGCGGGATGATGGGTATGGGACGCGGGTCCGGTCTGGTA
TTAACGGCCAGTCTTGACATGCAACGTATCGATGTGGCTGTGGCGTGGACAGCGTGA
GGTTGGGAGGTGTCGGCGCAGAGCATGGCGCACCGTTCCACGTTACGGCGTCAATT
CGTATTAAAGCCTGGACGGTGCTGCTCCGCCCTCGCACCTGCAAGGTTGAAAGGATA
TGTGGTGAGCCGCTCTCGGGAACTGCTCATCCACTCACCCAGCCGGCGTGCAC
CCCGTTATGTTCAATTGTCACATCTGGAACATGAAGATGCAGGTATGATGGTCAGTATG
TTGCGCATAAaagctt

QPGRFTTPLPVPLIDAQAQGNSIRLVAQTGRHAFVEGRPVATYGYSGPVLGPTIRMQSGA
QVDVAIENRIDTDIVWHGFLIPSERDGGPHDHIAPGETWRRVLPVQQPETAWYHPHPHRD
TGRQVYFGMAGLVLIEDGSGARLDLPRTYGVDDLPLVLQDRLFDGNGDLIYPANPMVMQGA
RGNTIIVNGAIPIARVPAGLVRLLNGSNARNLDLGFDGRPLHVIASDGGYLGPVSMTRL
LIAPSERFEILVDFSDGRPAVLQTGPDPFPMMGMMGRGQSAAGDIMSFAPDRDRPGAKRKI
PSTLVDMMPAAAPSAQLLRRRFALNDGGMMGGMMGRGMMGGMMGGMMGGMMGGMM
GGGMMGMGRGGSGLGINGQSFDMQRIDVAVALDSREVWEVSAQSMAHPFHVHGQFRILSL
DGAAPPHLQGWKDTVLVRSÆLLIHTQPALRAHPFMFHCHILEHEDAGMMGQYVCA

Faest

ggatccTGCAATAGTCATGATATGAGTGAATGAATATGGGTGGCGAAACGCCACCGCA
GTTACCGAAGGTAGTTTACCAACCCGCTGCGCTGCTGGATACCGCAAGCCGAGTGGTC
CGCTGAGCGCAAAGCACCACCGAAGCAATTGTGGCAGGTAAAATGCCGTGTTCTGG
GTTATCGCGATGGTATGCTGGGCCGACCTTCGTGTTACCAGCGCGCCACCGTTGATCTG
CGCTTCTGAATGCCCTGACCGAAGAACCAATATTGATGGCATGGCTGCTGGTCCGGC
CAATATGGATGCCATCCGGCACAGCTGGTGGCGCCGGTCAGAGTTAATTACCTTC

GTCTGGATCAGGCAGCAACCATGGCATGGTATCATCCGCATCCGCATGGCAAAACCGCCAA
ACAGGCCTATATGGGCTGGCCGGCTGTTATTGTTGAAACCCGACCGAAAAAGCCCTG
AGCCTGCCGAGCGGTGCATATGAAC TGCCGCTGATTCTGCAGGATAAACGTCTGGATGCCA
GCGGTAGCCCGCAGTATAATCCGAGTATGAATGATGTGATGCTGGTTATATGGCGAAATT
GTTACCGTTAATGGCGTGGCAAGTCCGATGCATAGTGTGCTGCCACCCGCATGTATGCCCTGC
CCTGGTTAATGGCAGTAATGGCGCTGTATAATCTGCCCTGAGTACCGCGCACAGTTT
GGGTGATTGGTAGTGATGGTGGCCTGCTGAGTGCCCCGAAGGTGTTACCGACTGCTGCTG
GGCCCCGGCGAACCGCCTGATCTGCTGGTGGATTAGCAGCGTCCGGTGGTACCGAA
GTGTATCTGAAAAGCAATACCTTACCGGTGGTGGCGCACAGGGCAGCAGGCATTAAATC
TGCTGAAATTGTTGTGAGTAAGGCCAACCGAACCTTCGTTACCGGCCGCTGGA
AACCGTGACCCGCTGGCAGCAAGCGTTGCAACCAAACCCGCACCTTGATATTGGTATT
GCCATGCAGGCCATGCAGGGCATGAATATGACCGGCATGCATACCATAATGGCAAAGTGT
TAGCATGAATCGCATTGATGAAACCGTAAACTGGGTGATACCGAAATTGGGAAATTAAATA
ATACCCAGGGCGATGAACCGCATCCGATGCATCTGCATGGTCATTTTCAGGTGCTGAGC
CGCACCGCGGTGCAATGCACTGACCGCCAGTGAAAAGGCTGAAAGATACCGTTCTG
GTTATGCCGGCGAACGTGTTCGTATTGTTGCCGTTACCAAACCGGGCACCTTGTGTT
TCATTGTCATAATCTGGAACATGGTGATGGCATGATGGTAATTATCAGGTGGCATAAaag
ctt

CNSHDMSAMNMGETATAVTEGSFTPLRLDTASPSGPLSAKSTTEAIVAGKNARVLGY
RDGMLGPTFRVTSATVDLRFLNALTEETNIWHGLLVPANMDGHPAQLVGAGQSFNYTFRL
DQAATMAWYHPHPHGKTAKQAYMGLAGLFIVETPTEKALSLPSGAYELPLILQDKRLDASGSP
QYNPSMNDVMLGYMGEIVTVNGVASPMHSVATRMYRLRVNGSNRLYNLALSTGAQFWVI
GSDGGLLSAPEGVTSL LAPGERADLLVDFSSVPVGTEVYLKSNTFSGGGAQGQQAFNLLKFV
VSKAETETFRLPARLETVTPLAASVATKTRTDIGIAMQMAMQGMNMTGMHTINGKVFSMNRID
ETVKLGDEIWEFNNTQGDEPHPMHLHGAFFQVLSRTGGRNALTASEKGWKDTVLVMPERV
RIVVPFTKPGTFVFHCHNLEHGDDGMMGNQVA

Mint

ggatccCGTGAACGCCCTGCTGCCATTCCGCCGTGCTGGAAGCCGATGCCGCAGGT
AATATTAAATCTGAGTATT CAGGCAGGTCTGAGTGCAATTCTGGCAGGCAAACAGACCCGA
CCTGGGGCTATAATGGCGCACTGCTGGCCCGCCCTGAAAGCTAAAGCAGGCCAGAATCT
GCGCGTGTGATGATGAAATACGCCCTGGATCAGCCAGTACCCCTGCATTGGCATGGTCTGGAAA
TTGGTGGTGAAGCAGATGGCGGTCCGCAGGCAGTGATTGCACCGCAGCAGCGCGAAG
TGAAATTTCGCGTGGATCAGGCAGCAGCCACCTGTTGTTTACCGCATAACCATCATCAG
ACCGGCTATCAGGTGCAATGGCTGGTGGCCTGTTATTGAAGATGAATTAGTCG
CGGTCTGAAACTGCCGAACGTTGGGTATTGATGATGTTCCGCTGATTCTGCAGGATAAA
CGTTTGATCCGCAAGGTAATATTGATTATAACTGGATGTGATGACCGCAGCAGTGGTTG
GTTTGGCGATGTGATGCTGACCAATGGCAGCCTGTATCGAAACATATTGCACCGAAAGGTT
GGCTGCGCCTGCGCCTGCTGAATGGTGTAAATGCCGTAGCCTGCGTCTGGCAGCAAGTGA
TGGTCGTCGATTATGTTATTGCCGGCATGGCGGCTTCTGCCGAACCGGTGGCGTGC
AGGAACGTGAGTCTGCTGATGGCGAACGCTTGAAGTTCTGCTGGATTGCAGCGATGGCAA
AGCATTGATCTGCTGAGTCTGCCGGTTACAGATGGGTATGAATATGCCCGTTGATC
AGCCGCTGGCACTGCTGAGCATTGAAACCCGCTGATGCCGGCACCGCCGCTTACCGG
ATAGTCTGATTAGCCTGCCGCCGCTGCCGGTACCGATAATCTGCCGTGCGTCATCTGAAA
CTGGAAATGGATCAGAAACTGGATCATCAGGGTATGATGCCCTGATGCCCGTTATGGCCA

GCAGGCAATGCAGCAGGGTCATCATATGCATCATATGCCACCATATGCAGGGCCAGGATGCCA
TGGCCACCGAAGAAGTGCAGATTATGGGTCCAATCGTATTAATGGTCAGAGTTTAGCATG
CATGAACCGATTTGATGTTAAAGTGGCCAGTATGAAAAATGGATTATTAGCGGTACCGG
CGATATGATGCTGCATCCGTTCATATTGATGGTACCCAGTTCGCATTCTGACCAGAAAATGG
CCAGACCCCCGCTGCCGCATCGTCGTCGCAAAGATATTGTGAAAGTGGAAAGCGGTGT
GAGCGAAGTTCTGGTGCAGTTAACATAAAGCAGATAAAGCAAACGCATATATGGCACAT
TGTCACTGCTGGAACATGAAGATAACGGCATGATGATGAGTTTACCGTTAGCAAATAAaag
ctt

RERALLPIPPLLEADAAGNIINLSIQAGLSAFLAGKQTPTWGYNGALLGPALKAKAGQNLR
VMINNRLDQPSTLHWHLGEIGGEADGGPQAVIAPQQREVKFRVDQAAATCWFHPTHHQTG
YQVAMGLGLFIIEDFSRGLKLPERWGIDDVPLILQDKRFDPQGNIDYKLDVMTAAVGWFGD
VMLTNGSLYPKHIAPKGWLRLRLNGCNARSLRLAASDGRPIYVIAGDGGFLPEPVAVQELSLL
MGERFEVLLDCSDGKAFLLSLPVHQMGMNMAPFDQPLALLSIETGLMPGTGRLPDSLISLPP
LPVTDNLPVRHLKLEMDQKLDHQGMMALMARYGQQAMQQGHMHMHMQGDAMAT
EEVQIMGANRINGQSFSMHEPMFDVKVGQYEKWIISGTGDMMLHPFHIHTQFRILTENGQTP
LPHRRGRKDIVKVEGGVSEVLVQFKHKADKANAYMAHCHLLEHEDTGMMMSFTVSK

HR03

ggatccATGACCCCTGGAAAAATTCTGTTGATGCCCTGCCGATTCCGGATACCCCTGAAACCG
GTGCAGCAGAGTAAAGAAAAACCTATTACGAGGTGACCATGGAAGAAATGCACCCATCAG
CTGCATCGTGATCTGCCGCCGACCAGACTGTGGGGTATAATGGTCTGTTCCGGGTCCGAC
CATTGAAGTTAACCGCAATGAAATGTGTACGTGAAGTGGATGAACAATCTGCCGAGTACC
CATTTCCTGCCGATTGATCATAACCATTCATCATAGTGATAGCCAGCATGAAGAAAGTGAAGT
GAAAACCGTGGTGCATCTGCATGGCGCGTTACCCCTGATGATAGTGATGGCTATCCGAA
GCATGGTTAGTAAAGATTCAACAGACCGGTCCGTATTTAACGTGAAGTGTATCATT
CCCTAACCGAGCAGCGCGGTGCAATTCTGTGGTATCATGATCATGCCATGGCACTGACCCGTC
TGAATGTTATGCCGGCCTGGTGGCGCTATATTATTATGATCCGAAAGAAAAGCGTCTG
AAACTGCCGAGTGATCGCTATGATGTGCCGCTGCTGATTACCGATCGCACCATTAAATGAAGA
TGGCAGTCTGTTACCCGAGTGACCGGAAAATCCGAGTCCGAGTCTGCCAATCCGAGC
ATTGTTCCGGCTTTGTGGTAAACCATTCTGGTTAATGGTAAAGTGTGGCGTATCTGGA
AGTGGAACCGCGAAATATCGCTTCGTGATTAATGCCAGCAATACCCGACCTATAATC
TGAGTCTGGATAATGGCGCGAATTCACTCAGGTGGTAGTGATGGCGGCCTGCTGCCTAG
AAAGTGTAAACTGAATAGCTTAGTCTGGCACCGGCGAACGCTATGATATTATTGATTT
CACCGCCTACGAGGGTGAAGCATTATTCTGGCCAATAGCGCAGGTTGTGGCGGTGACGTT
AATCCGGAAACCGATGCCAATATTATGCAGTTCGCTTACCAAACCGCTGGCCCAGAAAG
ATGAAAGCCGAAACCGGAATATCTGGCAAGTTATCCGAGCCTGCTGAGCATGAACGCATTCA
GAATATTCTGACCCCTGAAACTGGCGGTACCCAGGATGAATATGGCCGCTTCTGCTGC
TGAATAATAAGCGTTGGCATGATCCGGTGACCGGCGACCTAAAGTTGGTACCGACCGAAAT
TTGGAGCATTATTAATCCGACCCGTGGCACCCATCCGATTCTGCATCTGGTAGTTTCCG
CGTGCCTGGATCGCGCCCGTTGATATTGCACGCTATCAGGAAAGCGCGAACCTGAGTTATA
CCGGTCCGGCAGTGCCGCCCTCCCTAGCGCAGGAGGGTTGGAAAGATACCAATTCAAGGCAC
ATGCCGGCGAAGTGCTGCGCATTGCAGCTACATTGGTCCGTATAGTGGCGCTATGTTGG
CATTGCCATATTCTGGAACATGAAGATTACGATATGATGCGCCCGATGGATATTACCGATCGT
CATAAATAAaagctt

MTLEKFVDALPIPDTLKPVQQSKEKTYEVTMEECHQLHRDLPPTRLWGYNGLFPGPTI

EVKRNEVYVKWMNNLPSTHFLPIDHTIHSDSQHEESEVKTVVHLHGVTDDSDGYPEAW
FSKDFEQTGPYFKREVYHYPNQQRGAILWYHDHAMALTRLNVYAGLVGAYIIHDPKEKRLKL
PSDRYDVPLLTDRTRINEDGSLFYPASPENPSPSLPNPSIVPAFCGETILVNGKVWPYLEVEPRKY
RFRVINASNRTYNLSLDNGGEFIQVGSDGGLLPRSVKLNSFLAPAERYDIIDFTAYEGESIILA
NSAGCGGDVNPETDANIMQFRVTKPLAQKDESRKPEYLASYPVQHERIQNIRTLLAGQTDE
YGRPVLLNNKRWHDPVTGAPKVGTTIEWSIINPTRGTHPIHLHLSFRVLDRRPFDIARYQES
GELSYTGPAVPPPSEEGWKDTIQAHAGEVLRIAATFGPYSGRYVWHCHILEHEDYDMMRPM
DITDRHK

Aory

ggatccATGGCCCAGCCGCCCTAAATGGACCCCGCCGGCTCATGCCCGCTGGACCCCTAGA
AGTCAGAAACAGTTGTGAATCCGCTGCCGATTCCGCCGGTGGCACGCCCTGCTGTTGCAG
ATGGCACCCATTATGAAATTCCGATTGGCGAATTTCGCCAGCATCTGGCATTCCGCGATCCG
GAAACCGGCGCCCCGCTGATGACCACCGTGTGGGGTTATGCCGGTCTTATCCGGGCCCCGA
CCATTGAAGCACGTCGCGGTCGTCCGATTACCGTTCGCTGGGTGAATGCACTGTATGGTGCC
AGCGGTCCGCTGCCGCATCGCATGCCCTGTTGATACCAGTGTTCATATGGCCCGCTGCGTAA
TTGGCCGGCCAGCGGTGTTCCGACCGTTACCCATCTGCATGGTGGTCATAACGAATGGCAG
AGTGTGGTGATCCGGATGCATGGTATAACCCGGCTATGCACAGAAAGGCCCCGCTGTTG
CAAAGAAGTGTATAACCTATGATAACGATCAGGAAGCCGACTGCTGTGGTATCATGATCAT
GCACTGGGTTTACCCGCCCTGAATGTTATACCGGCCTGGCCGGTCTGTATGTGATTGCTGA
TGATTGGGAACGCAGCCTGGGTCTGCCGAGTGGTCGTTATGAACTGCCGCTGCTGATTCA
GATCGTAGTTTATGCCAATGGTGAACGTATTATCATGCCCTGGCCCAAGGAACCGAGTC
TCCGAGTCCGAGTGTCTGCCGGAACAGTTGGTATGTGATTCTGGTGAATGGTATGCCCT
GGCCGACCGCCGAAGTTGAACCGCGTAAATATGCCCTCGCTGCTGAATGGCAGTGATAG
TCGTTTTATGCCCTGGGTTAGTAGTCGCCCTGAAATTATTAGTGGCAGTGATGGTGG
TCTGCTGGATGCCCGGTTGAACTGACCGAACTGCTGCTGGCACCGGCCAACCGGCCGAT
GTGATTGTGGATTTCGCCCGCTCGCGGGTAAAACCGTTCTGCTGCTGAATGATGCCCGC
CCCGTTCCGGATGGTACCTGTTGACCGTACCAATGGTCGTTATGCCCTAAAG
TGGCAGTACCAACCAAAAGATGAAAGTCGTATTCCGATGCCCTCGCTGCAACCGAT
TCCGTTCTGGTCCGACCGCACCACCGAAACGTGCCCTGCTGTGGGAAGGTCAAGGATGCA
CATGCCGCCAGAAAAGCCTGCTGGCACCGTGCAGGATGGTCACTGCATTGAAATAGTC
CGATTACCGAAAATCCGCCCTGGGTCAACCGAAGTGTGGAAATTATAATACCA
GGATACCCATCCGGTGCATATTGTTCTGGTCAATTGAGTCGTCAGAAATTAA
GGCCGATCTGGATGAAAAAACCGGCCCTGAGCAATATTGTTACCGGCCGCT
CCGCCGCCCTCGAGGAGCGTGGTGGAAAGATACCGTCGCAATCCGGCGAAGTTA
CCCGTGTGATTGCAACCTTGATGCCCGGCCGTTATGTTGGCATTGTCATATTGAGCC
ATGAAGATCATGAAATGATGCGTCCGTATCAGGTGGCCCGATTGGTGAACCCCGGTG
TGGTGCACTGTAAagctt

MAQPPKWTPPAHPLDPRSQQKFVNPLPIPPVARPAVADGTHYEIPIGEFRQHLGIRDPE
APLMTTVWGYAGRPGPTIEARRGRPITVRWWNALYGASGPLPHRMPVDTSVHMAPLRNWPA
SGVPTVTHLHGHHTEWQSDGDPDAWYTPGYAQKGPLFAKEVYTYDNDQEALLWYHDHAL
GFTRLNVYTGLAGLYVIRDDWERSLGLPSGRYELPLLIQDRSFYANGELYYHALAQEP
SHPSVLPEQFGDVILVNGMAWPTAEVEPRKYRLRLLNGSDSRFYRLGFSSRLKFIQIG
SDGGLLDAPV
ELTELLLAPAERADVIVDFAPLRGKTVLLNDARAPFPDGDPVDPRTNGRVMAFKVG
STTKD
ESRIPDRLRREPIPFLVPTAPKRALLWEGQDAHGRQKSLLGTVQDGALHWNSP
ITENPALGAT

EVWEFYNTTPDTHPVHILVFRIMSRQKFKADLDEKTGRSLNIRFTGPPLPPPPEERGWKDTVRANPGEVTRVIATFDRPGRYVWHCHILSHEDHEMMRPYQVGPIGETPGRGAL

Talbi

ggatccATGACAGAATTACTAGGCCCTACACTCACAGTTGCCGCTGCACGTCTGTGGGGTTACGAGGGTCAGTATCCGGTCCGATTATCGACACCCGCCGTAAATAAACCGATCGCGGTTCACTGGGAGAACCGCCTGCCGACCGTCATATATTACCGATCGACCCTCATATTACGGTCTATGCCGCCGACCCCGGAAGTTCGCACTGTGCCGACCTGCACGGAGCCAATGTTCCGTCACCGTCCGACGGCTTGCCGGAAAAGTGGTCACCCGGGTCAAGCGCACGTTATGAATATCCGAATCGCCAACGTGCCGCAACCCGTGGTATCATGATCATGCAATCGGCATCACCGTCTTAACGTTATGCCGGTTGAGCAGCTTTATCTGCTGAGGGACGACGAGGAGCTGGGTATGCACTTGCCTGAGTACGAGATCCCATTGGTTCTGCAGGATCGTACCTTAGACGAGCACGGCCAACCTCTGTACGCCCGAACCCCTGGACGACGCCGTACCGCTGCCGAAAGGTGTTGGGGTCCGATGTTTCGGTGAATTGCCGGTCTGCAATGGGGCGATTATCCATACGTGGAA GTGAGACCGCAGTTACCGTATTGCTCAACAGTTCAACTCTCGCTCCTGAACCTGTACCTGAACCTGGGCGAAATCCCCGACGGATATCCCGCAATTGATCAAATTCCACCAAGATTGCGACGGACGGCGGTCTGATGCCCGTCCGGTGGCGCTGGAGAAGCTGGAGCTGGCGCCTGGCAAAGAGCGGATCTGATTGTTGATTCTCCAGCTTAGCTGGTAAGACCGTGACGCTGAGCAACGATGCTCCGGCTCCGTACCCAGGCTGGAATGTGATGAATGCAACTTGGCCGATCCTGCTGGAAATTATGCAGTTCTGTGTTACCCCTCCCTGCTGAACCGTCGCGATCATTCAAGCTGCAACCGTACCCGCACTTCATCTTGAACCGAAGAGATGGATGCGCAAGGTCGTAGCGTTGGCTTGACATTAACGGCAAGGGTTACCATGACCCGGTGACCGAAACCGTGACCCCTGGCACCCCTCGAGAAGTGGCGTTCATCACACCAGCGACGATGCGCATCCGATGCATCTCACCTGGTGCAATTCCAGATTCTGCACCGTCAAGGTTAACCTGGCACCTACCGCATGAACGGCAAGATGCAACCGGTGGGTTGACGCGTCAGCCGCGCGAATGAACAGGGTTGGAAAGATACCGCAACGGTAAACCCGGGTGACA TTTGACCATTCTGGCGCGCTCGAAGGCTACACCGTAAATACGTTTCACTGCCATATGCTGGAGCAGGAGATAACGATATGATGCGTCCATTGTTGTCGTGGCTCCGGCCAAGGCTAAagatt

MTEFTRPLHSQQLPPARLWGYEGQYPGPIIDTRRNKPIAVQWENRLPTRHILPIDPHIHGAMPPTPEVRTVPHLHGAVPSPSDGLPEKWFTPGRHSARYEYNRQRAATLWYHDHAIGITRLNVYAGLSSFYLLRDDELMHLPSEYEIPVLQDRTLDEHGQLLYAPTLDDAVPLPKGVWGPMPFFELPVVNGAIYPYVEVRPQLYRIRLLNSSNSRVNLNLAKSPTDIPQLIKFHQIGSDGGLMPRPVALEKLELAPGERADLIVDFSSLAGKTVTLSNDAPAPYPGWNVMNATWPILLEFMQFRVTLPAPRSRSFLPNVTFAKLDEGEAIRTRDFILTEEMDAQGRSVGLHINGKGYHDPVTETVTLTLEKWRFINTSDDAHPMHLHVFQILHRQGFNLGTYRMNGKIEPVGLTRQPRANEQGWKDTATVNPGDILTILARFEGYTGYVFHCHMLEHEDNDMMRPFVVVAPGQG

Bfre

ggatccATGAAGATCAAGAACAGCGTCGTATGAAAGTGGACCCTAGCAATCCGGAAACCATCCGAAATATGGATGAACTGCCGATTCCGCCGGTTGCACGTCCGCTGGCAGAAATTAAAGGTAGCCGTATTATGAAATTGCCATGCGTCAGGTGCCGATCGCTTCATCGCTGTTCCGCCGACCCGTTGGGCTATGATGGTATGCTGCCGGTCCGACCATTAAAGTCAGAAAGATGAAAAAAATCTACGTGCGTTGGAAAAAAACTGCCGAAAAACATCTGCTGCCATTGATCGCACCCCTGCATGAAACCGCAGGTCCGCCGGATGTGCGCACCGTGGTGCATCTGCATGGTCAAATGTGGCATGGGATAGTGATGGTCATCCGGAAGCCTGGTTAGTCGCATTTGC

CAAAACCGGTGCCACCTTCGTCGCAAAGTTATGAATATACCAATAAACAGATGGGTGCA
ACCTGTGGTATCATGATCATGCCATTGGCATTACCCGCCTGAATGTGTAGTGGTCTGAGC
GGCTTTATCTGATTGAAGATCCGGTGGAAAAACATTAAAAGATGGTATGAA
TATTCCGCTGATGATTCAAGGATCGCAGCTTCGTTAGTGATGGTAGTCTGAGTTATCCGGAAA
ATACCAATCCGCCGGCCCCGGTGAATCCGAGCGTCAGCCGTTTATTGGTAATACCATT
GCAGTTAACGGTAAAATTGGCGAAACTGACCGTGAACCGCGTAAATATCGCTTCGTAT
TCTGAATGCAAGCAATACCAATGCCATACCCCTGCGCCTGGGTGATGGTCGCAAATTATC
AGATTAGCACCGATGGTGGCCTGCTGACCGAACCGGGTGAACGTGACCAACCCCTGCCGCTGG
AACCGGCAGAACGCAGTGTGATTATTGATTAGTCAGCATAAGGGTAAGAAACTGAT
TCTGCAGAATACCAATCGGAAGGCAATATGGCATTATTATCGCCTTGATGTTCTGCAGC
CGCTGCGCGTCGCGATACCACTGAAATTCCGGAAAAGTATTAGCAGAACAGGGTCT
GTATGAACATCATGCAGATAAAACCCGCTGCTGAAACTGGATGCAATTCAAGGATGAATATA
ATCGCCCGTTCTGCTGCTGGATGATCGCATGTGGATGATCCGGTACCGAAAAACCGGT
TATTGGTGTACCGAAGTGTGGAAACTGATTAATGTTACCAATTGCCCATTGATTCAAT
TCATCTGATTCACTTAAATCCTGCATCGCACCCGTTGATCTGGAACGCTTCAGCAGG
ATGGCTATATTGATTATACCGTCCGCCATTGAAACCGGCCATTGAAACGTGGTGGAAA
GATACCGTGAAGCCAACCGGGTATGGTTACCAAGCGTGATTATGAAATTACCGAAAATCC
GGGTGAATATGTTGGCATTGCCATTCTGGAACATGAAGATTATGATATGATGCCCGAT
GCGCGTGGTGGAAAAAGAAAAATAAAagctt

MKIKNRRMKVDPSPNPETIPKYMDELPIPPVARPLAEIKGSPYYEIAMRQVPHRFHRLFPP
TTVWGYDGMLPGPTIKVQKDEKIYVRWKNKLPEKHLPIRTLHETAGPPDVRTVVLHGAN
VAWDSDGHPEAWFSRDFAKTGATFRRKVYEYTNKQMGATLWYHDHAIGTRLNVYSGLSFY
LIEDPVEKHLKLPKDGYDIPLMIQDRSFRSDGSLSYPENTNPPAPVNPSVQPFIGNTIAVNGKIW
PKLTVEPRKYRFRILNASNTNAYTLRLGDRKFYQISTDGLLTEPVELTLPLEPAERSDVIIDF
SQHKGKKLILQNTNAEGNMGIIMRFDVLQPLRGRTSEIPAKLISEEQVLYEHADKTRLLKLD
AIQDEYNRPVLLDDRMWHDPVTEKPVIDTEVWKLINVTNFAPIHIHLIQFKILHRTPDFDE
RFQQDGYIDYTGPPIEPAVHERGWKDTVKAEPGMVTSVIMKFTENPGEYVWHCHILEHEDYD
MMRPMRVVEKEK

Tocean

ggatccATGCATACCATGAATTGTAAGGTGAATCAGAATTATCCGTTTATAAGGAGAT
GATCAAACGGATCGCAATTCTGTATCTGTATGGTGGCCGTGGCTATTATGGCCCGAA
ACTGGAAAAATTAAAGATCCGCTGCCGGTCCGAATTATTGCCCGGTGGCACCAAA
AATGGCATTCCGTATTATGAACCTGAAATGACCGAATTACCAAGATCTGCATAGTGTATG
CCGGATACCGCGTGTGGGCTATGAAAATAGCTATCCGGGCCGACCATTGAAGCATATAA
AGGCCAGCCGATTATGTGCGTTGGATTAAATGCCCTGCCGGAAAACATTCTGCCGATTG
ATAAAACCATTATGGTGCCAAAGATAATCCGAAGTTCGCACCGTGGTCATCTGCATGGT
CTGAATGTTCGTCCGGATAGCGATGGCTTCCGGATGATTGGTTACCCGGTAAAGTGC
ACTGTATTATTATCCGAATAAGCAGCAGGCCGAACCCCTGTGGTATCATGATCATGCAGTTG
GTATTACCGTCTGAATGTGTATGCAGGTCTGGTGGTCTGTATCTGATTGCGATGAACGT
GAAGAACGCTCTGAATCTGCCGAGCGCGAATATGAAATTCCGCTGATTATTCAAGGATAAAG
ATTTAATGACGACGGTAGCCTGAGTTATCCGGCACCGAATAATAATGGTGTGGAACCGAGT
ATTGTTCCGGGCTTTGGTAATTATGCACTGGTGAATGGCAAAGTTGGCGTATCTGGTT
GTGAAACCGCGTAAATATCGCTTCGCATTCTGAATGGCAGTAATAGCCGTAGCTATAAACT
GCGCTTAGTGGTGAAGCCGGTGTATTCTGCCGGCTGGTATCAGATTGGCACCGATGGC

GGTTTCTGGAACGTCCGGTGAAACTGGATAGCCTGAGCCTGCAGCCGCCAACCGCGCT
GATGTTATTGTGGATTTACCGGTCTGGAAAATCAGACCTTACCGTATTAAATGAAGAAATT
CCGCCGTTGGCCGCCGCTGCCGGAAATTATGCAGTTCGCGTAGTGGCACCCGGTGAA
AAGATAATAGTAGCCTGCCGCTGATTCTGAATCGTATTAAACCGCTGCCGGTGAAAAAAAGC
CAAACAGCGCAATATTGAAATTGTGGTTGGTCAGGATGAACGGCCGTTATGTTATGC
TGGAAAATAAAAGTGGACCGATCCGGTACCATTAAAACCAAACGGCAATGTTGAAGT
GTGGAATATTATTAATACCGCCGAGCCGCCATCCGATTGATGTCATCTGGTCAGTTCA
GATTCTGAATGCCAGCCGTTGATGTCAGGAATACCGCAACCGGACCCCTGAAATTAA
TTGGTCCGCCGGTCTGCCGGATGAAAATGAAAAAGGCTGGAAAGATACCGTGCCTGCCG
AACCGGGCATGTGACCCGATTATTGCCGTTGGCATTACCGCATTATCCGTTTC
ATTGCCATATTCTGGAACATGAAGATCATGAAATGATGCGCCGTTGAAGTGTAAACAT
AAAAGCGGTGGTAAAGATAAGAAAGATAGTGAAGAAATCCACACCGTTAACCGGATACC
AAAGCAGATAAAGAAGAACCGCCAATAAaagctt

GSMHTMNCKVNQNFIPIFYKEMIKLDRAILYLYGGRGYYMAPKLEKFKDPLPVPKFIRPG
TKNGIPYYELEMTEFYQNLHSDMPDTRVWGYENSYPGPTIEAYKGQPIYVRWINRLPEKHFLPI
DKTIHGAKDNPEVRTVVHLHGLNVRPDSDGFPDDWFTPGKSALYYYPNKQQATLWYHDHA
VGITRLNVYAGLVGLYLIRDEREERLNLPSEYEIPLIIQDKDFNDDGSLSYAPNNNGVEPSIVP
GFFGNYALVNGKVWPYLVVKPRKYRFRILNGNSRSYKLRFSGEAGRILPAWYQIGTDGGFLE
RPVKLDSSLQPAERADVIVDFTGLENTFTLINEEIPPPGPPPLPEIMQFRVSGTPVKDNSSLPLIL
NRIKPLPVKKAKQRNIEIVVGQDELGRFMFMLENKKWTDPVTIKTKLGNEVWNINTAAAAA
HPIHVHLVQFQILNRQPFDVQEYTATGTLKFIGPPVLPDENEKWKDTVRAEPGHVTRIIARFG
DFTGIYPFHCHILEHEDHEMMRPFEVFHKSGGKDKKDSEEIHTVKPDTKADKEEPAK

BMS3

ggatccTGCGGTGATAATCAGCAGACCGCCGCAACCACCCGCCGGTGATCCTGTGACC
ACCAGTAAAGTTATCGTACCCCAGCACGTATTGATCATTATTGATCCGCTGCCGGTCCG
CCGACCCCTGAGTCCTGATACCGAACGCTATCCGGATACCGATTATTGAAATTGTATGAGT
CAGTCGGTCAGGAACACTGCATAGCCAGCTGCCGGAAACCACCGTTGGGCTATGAAAAAA
AGTTTCCGGGCCCCGACCATTGAAGCCCGTCGTGGCCGTACCCGCGTTAAATGGATTAA
ATGATGATCTGCCGCCGCCCATTGCTGCAGGCCTCAATTGATCCGACCATTATCAGGGC
CGTGAATTCCGGATGTGCGACCATTGTGCATCTGCATGGTAGCTTGTCCGTCCGGATT
TGATGCCAGCCGGATGCCCTGGAGCAGCCGGGTGCAGGTAATACCGGGGGTCATT
GGTAGTGAATTACCTATCGAATGATCAGCCGGCACCATGCTGTGGTATCATGATCATGC
AATGCATAGCACCCTGAATGTTATGCCGGTCTGCCGGTCTGTATTATTGCGATGA
AGAAGAAGAAAAGCTGAGTCTGCCGGTGGCAATCATGAAATTCCGCTGGTTATTCAAGGAT
CGCACCCCTGGCGATGATGCAAGTCTGAGTTATCCGACCACCGCATTACCCGGTCATCC
GGTGGGTGATGCGCTTCTGGGTGATATGCCGTGATTAATGGTAGCCCTATCGTATCT
GGATGCCGAACCGCGCCATTGCCCTGCCGTGATTAATGCCAGTAATACCGCACCTGG
AATCTGTGGTTGATGCCGGTGGTGGCCCGTTCCGTTATGTGATTGGTAGTGTGGTGG
TTTCTGCCGGCCCCGGCCCGCTGAAAGTCTGCGTCTGCCGGCAGAACCGCGCAGAT
GTGATTCTGGATCTGAGCGCAGCCGATCCGGCACCGTTTACCCCTGCGTAATGATGCC
GGCACCGTATCCGAAGGTGGCGATCAGCCGCTGGATAATCTGATGCAAGATTGCTGAGC
CGTGAATGGATGGTAAGATCGCAGTACCCGGCAGATCGTCTGAGTCTGCCGTGCAACCA
GCACCGCAACCCGACCAAGTGGCGTCCGAAACGCATTATCTGGATGAACGTGAGCAA
TGATCAGGGCAGCGCAGTGGAACTGCTGATTAATGATGCCGTTTATGATCCGGTTGAA

GAACAGCCGGCAGCAGGCACCACCGAAATTGGGAAATTGTGAATCTGACCCCTGGATGCC
CATCCGTTCATATTCATCTGGTCAGTTCAGGTGATGAATGCCAGCCCTGGACCAGGC
ACAGTATTGGAAAGATAAAGATGCATATGTTGCCGGCACCGCACCAGTCCGGACCCCTATA
CCTATCTGACCGGCGCAGCAGTGCCGCCACCTGCAGAAACCGGTTGGAAAGATAACCG
CCCTGAGCATGCCGGCGAAGTCTGCGTCTGGCGGTGCCGTTAGCCTGCCGGCGGTGT
GAGTGGCCCGGAAGATATATTATGCCATATTCTGGAGCATGAAGATAATGATATGAT
GCGCCGTTGATGTTAAaagctt

CGDNQQTAATTRPGDPVTTSKVIRTPARIDHFIDPLPVPPTLSPTERYPDTDYYEIRMSQF
GQELHSQLPETTVWGYEKSFPGPTIEARRGRTRVKWINDDLPAHLLQASIDPTIYQGREFPD
VRTIVHLHGSFVRPEFDGQPDAWSSPGAGNTGPGHFGSEFTYPNDQPATMLWYHDHAMHSTR
LNVYAGLAGLYFIRDEEEEKLSPVGNHEIPLVIQDRTLADDASLSYPTGITPVHPVVVMRFL
GDMPLINGSAYPYLDAEPRRYRLLLNASNRTWNLFADAGGGPFPHVIGSDGGFLPAPARV
ESLRLAPAERADVLDLSAADPGTVFTLRNDAPAPYPEGGDQPLDNLMQIRLSRDLDGEDRSTP
ADRSLSPATSTATPTSGVPKRIFYLDELSNDQSAVELLINDRRFHDPVEEQPAAGTTEIWEIVNL
TLDAHPFHIHLVQFQVMNRQPLDQAQYWKDKDAYVAGTAPIPDPTYLTGAAVPPAPAETGW
KDTALSMPEAVLRLAVPFSLPAGVSGPARYIYHCHILEHEDNDMMRPFDVV

Acdi

ggatccATGGATAACGTTCTGGGCCGGATACCACCACTATCCGGGCACCGATTATTATC
ATGTTCGCATGGTTCAAGGTTCAAGCAGGATCTGGCTCTGATTGATCCGGAAACCCGTAAACC
GCTCGTACCAACCGTTGGCATATGGCGCCCAATCGTACCGAACCTATCCGGTCCG
ACCATTGTTGCCCATAGTACCTGGCAACCAATGCCAGCCGGAAACCGGTGAAAGTGC
GCTGGAAAATGCACTGCCGGATAAACATCTGCTGCCGGTTGATACCACCGTTATTGTGG
CCCGGATGCACGCCAGCAGCATAGCATTGTCGTCGTTGTGCGCACCGTTACCCATCTGC
ATGGCGGTATGTGCCGGATCATAGTGTGTTATCCGAAGCATGGTTAGCCGGGTTTT
CGTAAAAAAGGTCCGCTGTGGAGTCGTGAAGTGTATGATTATCCAATGATCAGGAAGCCG
CCACCCCTGTGGTATCATGATCATGCCATGGTATTACCGTCTGAATGTGTATGCAGGCCTGG
CAGGTTTTATCTGGCGGTGATGATAATGAACAGCGTCTGCAGCGTAGTTTCGTCGTGATGGCAGTCT
GCCCTGCTTATGATGTGCCGCTGGTGATTAGGATCGTAGTTTCGTCGTGATGGCAGTCT
GTTTATACCAGCGTAAAGCGAATTGAAACAGCGTCCGGAAGAAGCCAAAAACCGC
CGCAGCACCGGGCGAACGCCTGCCAGAGATCCGATTACCGTCAGCTGAGTAGCAGTATT
GAACCGGAATTTTGGTGATACCATTCTGGTAATGGTAAAGCCTGGCCGGTTCTGGAAG
TGGAACCGCGCAAATATGCCCTGCGCGTTCTGAATGGTAGCAATAGCGTTTATCGTCTG
ACCCTGAGCAGCGGTAGGAATTTCAGCAGATTGGTAGCGATGGCGGTTCTGAATGCC
CGGTTGCCCGCGTGAACCTGCTGGCCCCGGCAGAACCGCCGATGTGATTGTGGATT
TGCAGATCCGAAACTGGCGGTGACGACATTGTGGTCTGCAATGATGCCCGACCCCGTT
CCGAATGGCCATGCAAGTTGATCCGGCAACCAACCGGTAGGTTATGGCATTGCGTTAGCA
AACCGATGAGTGCACCCCGATGCAACCCCTGCCGCCAGTCTGCTGCTGGCGGAAATTAAAGATGAATT
AACTGCCGGGCTGCGTGCCTGCCCGTTCGTCAGCTGCTGCTGGCGGAAATTAAAGATGAATT
TGGCCGCATTAAAACCATGCTGGGCCACCGTTGAACATGGCGCACTGGCTGGATGCCCG
ATTAGCGAAACCCCGCGCCGTAATGATGTGGAAATTGGAGCGTTGTGAATGCCACCCCG
ATGCCCATCCGATGCATCTGCATATGGTTTTTCAGGTTCTGGATGCCAGAAATATGATG
CCGAAAAATTGAAGCAGGCCAACCGGCCACCCCTGCCCTGACCGTACCGCAATGGCAC
CGCCGGCCGGTGAACGCGGTTGGAAAGATACCGTTATTATGCGTCCGGGTGAAGTGACCCG
TGTATTGCCCGCTTGTATCTGCCGGTTATATGTGTGGCATTGCCATATTCTGGAACATGA

AGATCATGAAATGATGCCCGTATCGTGTCTGCCGTAAagctt

MDNVLVPDTTYPGTDYYHVRMVQVQQDLGLIDPETRKPLRTTVWAYGAANRTATYPG
PTIVAHSTLATNRQPGKPVKVRWENALPDKHLLPVDTTVHCGPDARQQHSHCRPFVRTVTHL
HGGHVPDHSDGYPEAWFSPGFREKGPLWSREVDYDYNPNDQEATLWYHDHAMGITRLNVYAG
LAGFYLVRRDDNEQRQLQRDGRLPAPAYDVPLVIQDRSFRDGSLFYTSGKGEFEQRPEEAEKTA
AAPGERLPRDPITGQLSSSIEPEFFGDTILVNGKAWPVLEVEPRKYRLRVLNGNSNSRFYRLTLSS
GQEQQIQIGSDGGFLNAPVARRELLAPAERADVIVDFADPKLAGQTIVVRNDARTPFPNGHAV
DPATTGQVMAFRVSKPMSATADATLPASLRAPLAKLPGLRARVRQLLLAEIKDEFGRIKMLGT
VEHGALGWDAPISETPRRNDVEIWSVVNATPDAHPMHLHMVFQVLDRQKYDAEKFEAGKP
ATLRLTGTAMAPPAGERGWKDTVIMRPGEVTRVIARFDLPGLYVWHCHILEHEDHEMMR PYR
VLP

Bceee

ggatccATGGAAAATGTGAAAGCCATCAGATCCAACTCTATCCCTAAATTCA^GGA^C
CGCTACCGAAGCCGAGCGTTGCAGTCCC^GTT^CAGCACGATGACTACC^CGGATGGCAGCTA
TTACGAGCTGGAGATGAAAGAAGGTTGCACCGCTACCATCAAAACTTCCGGAAACCAA
GATTGGGGTTACGATGGTCTGATACCTGGTCCGACC^AT^CGAAGCGCTAAAGACAAAACC
ACCTATGTTAAATATTAAATAATCTGCCGGATGAGCACTTC^CTC^TGC^CCT^TGGATCGCAC^CTTG
CACAGCTCCATCGATACCGCAGATGTGCGCACCGT^CGT^CACCTGCACGGCGCTAAAGTTG
ATTGGGAAAGCGACGGT^CATCCGGAAAGCATGGTACACCAAAACTACGAGATGACCGGTC
CGACCTTC^CGT^CCGT^CAGGTTACCGT^ACC^AT^CAC^ACCAGCCGGT^CCCAC^CCTGTGGTAT
CATGACCACGCTATGAGCTTGACCCGT^TTAACGT^TT^ACTGGCCTGGCAGGCTTTACCT
GCTCGTGATGCCCTGGAAAGATCGTTGAGACTGCCGAGTGGTAAGTACGAGATTCCGATG
ATGATTCAAGATCGCAGCTTAATGATGATGGCAGCTGTTCTACCCGGACACCCCGCCGTT
CCCGGTTACGGTCAATCCGAGCATCACCCCGGCGTTAGGTGATACGATTGCGGTCAAC
GGCAAGCTGTGGCGTATCTGAACGTGGAGCCCGT^AAGTATCGTTTC^TGT^TCTCAACG
CGAGCAATCGT^CGT^TGGTTACTCCCTGTCT^CTGAGCAACGGTGGTACGATCCATCAGATTGGC
ACGGACGGCGGTTGTGGAGGCTCCGGCGAGTTGACCTCGT^CGATCTGCTGCCGGCTG
AGCGCACAGACATTATTATCGACTTTCTACGATGCAAGGT^CAGACCATCACCTGTTGAAC
AACGATCCAGAATTCCCGCTGAATGAGCACACCTCCGTGGTTATGCAGTT^CAGCGTTGTCT
GCCGTTAGATGGTGAAGATGACAGCCAGATTCCGGAGACCTTATATCCGGAAATGGCGCTG
CACACCGAACACCGC^AT^CGT^CGA^AT^TCGT^CAC^ACC^ACCCCGAT^CCA^ACC^AGGATT^CAC^C
GTCGTCCGATGCTGATGTTGAACGACCATATGTATCATGACCCGGTAGCGAACGCCGT^C
CTCGACAGCATCGAGATCTGAATT^TCATCAACACCACCCGAT^CCA^ACC^AGGATT^CAC^C
GCACCTGATTCAATTAAAATC^TGGAGCGCCGT^CCGTTGATGTTGATGTGTTACCGAAA
CCGGTGAATTGTTTACCGGTCCAGCCGAGCC^CCACGTGAATACGAGCGTGGATGGAA
GGATGTTGTGCGCGCTGATATCGGATGGTACCTCCATTGCTATGCATTGGAAAGGACTTCA
CCGGAAACTATATTGGCATTGTCATTCTGGAACACGAAGATCACGACATGATGCGTCCG
ATTACCATCATCGAAAATACG^CATCCAGTT^CAGCTGCC^CATGCTGACGAGGCCCGCG
AGGTTCCAGGCCAGCGGAACCGCAACCGACACCGTGACACCGGAGGGACCCACAGGTG
ACCGACGAAACAGCGGAGGAAACGGTTAACGAAGTTATTGAAGGTACGGATGATACCACA
GCGACCACCGACACGGAAAGGT^ATGGT^AAGCGACGCCGGACGGTGGCAGAGGGGACT
GACGATTCTAACGCAGGGACTCCCGCGAGGGCGCCGATGATTCCGACGCCGGGACCGCG
GTGGAAGGAGCCGACGACTCCGGTGC^GGGCACAACGGTGGAAAGGC^CCGATGACTCCGG
GGCGGCAACGACCGCGGAGGGCACAGACGACTCCGACACC^CGGTACCA^CGGTTGAGGGTG

CAGATGACAGCGCGCAGGCACGCCGGCGGAAGGTGCGGTGGATTCAAGATGCCGGCACG
ACCGTGGAGGGTGCACGACTCTGGCGCCGCACCCCGGCCAGGGTGCAGTCGATTCC
GACGCCGGTACCACTGTGGAGGGCACCAATGACACTGACGCTGGACCACGCCGGAGGGC
ACTGACGACTCTGGCGCTGGTACCAACCGTGGAAAGGTACGGACGACTCCGATGCAGGC
ACCGCAGAAGGTACCGACGATTCCGGCGCGGGTACCAACCGTCGAGGGCACGGACGTGAGC
GACGCCGGTACGGCTGTGGTAGCACCAGATTGAAAGAGCACCATGGCGATCCAGGGCTAC
GAAGCGATTGAAAAGCAGAAAACGCGACTGACCGCAGAGACCCGCATGGCACCCGCTGA
GGGCACTGCGCCCACGGCACCGACGACACCGGTGCAGGCACCACCGTAGAGGGTACCG
ATGTAAGCGACGCTGGTACTGCGGTGGTTCAACCGACTGGAAGTCGACCATGGCCATCCA
AGGTTACGAGGCAATCGAGACCCAAAACAGCCGGTAGAGCAGGTACTGGTATGGCTCC
GACCTACAAACAGAAACCGTCAACTGCCGAAGTGCCAACCGCACCGAAGAAAAAGA
AAGGTAATAAaagctt

MGKCVKPSDPNSIPKMDPLPKPSAVPVQHDDYPDGSYYELEMKEGLHRYHQNFPEK
WGYDGLIPGPTIEARKDKTTYVKYLNNLPDEHFLPLDRTLHSSIDTADVRTVVLHGAKVDW
ESDGHPEAWYTKNYEMTGPTRRQVHAYTNQPGATLWYHDHAMSLTRLNVYSGLAGFYLL
RDALEDRRLPSGKYEIPMMIQDRSFNDDGSLFYPDTPFPVTVNPSITPGVLGDTIAVNGKLW
PYLNVEPRKYRFRVLNASNRGYSLSLNSNGTIHQIGTDGLLEAPAELTSFDLLPAERTDIIIDF
STMQQQTITLLNNDPEFPLNEHTSVVMQFSVCLPLDGEDDSQIPETLYPEMALHTEAHIVRN
PLTATTDEYGRPMLMLNDHMYHDPASERPSLDSIEIWNFINNTPIQHPIHLHLIQFKILERRPF
DVFTETGEIVFTGPAEPPREYERGWKDVVRADIGMVTISIAMHWKDFGTGNYIWHCHFLEHEDH
DMMRPITIENTHPVQLPHADEAPAEVPAPAEPPTDTVPEDPQVTDETAETVNEVIEGTDDTT
ATTDTATEGNGDSDAGTVAEGTDDSNAGTPAEGADDSDAGTAVEGADDSGAGTTVEGADDSGA
ATTAEGTDDSDTGTVEGADDSGAGTPAEGAVDSDAGTTVEGADDSGAGTPAEGAVDSDAGT
TVEGTNDTDAGTTAEGTDDSGAGTTVEGTDSDAGTTAEGTDDSGAGTTVEGTDVSDAGTAV
GSTDWKSTMAIQGYEAIETQNATDRAETGMAPAEGTAPDGTDDTGAGTTVEGTDVSDAGTAV
GSTDWKSTMAIQGYEAIETQNSPGRAGTGMAPTYKQKPCNCPKCRYARKKKKGK

Scla

ggatccGCTCCACCCGGAGGGTACCTGCAGTAAGGGACCCCTCGTTACCTGGCTCCGATT
GCGGATGCGAGAGCGGTGCCCGTTCTGTGACGCCGTGCCGCCAACCGTATTGATC
TCTCGGATGGCGGCTCGCGCACCTGGTATGGCGCCGGTACCCAAGACGTGGTGGGTGC
AGGCCTCGGTTGCGCACTCCGGTGGGGTTCCGGCGGTGGTCCGGTCCGGTGC
TAGGCCGAGCTGGCCTGGCCGACGCTGGTGCCTCCGGTCCGGTCCGGTGC
TGGCGTAATGCGCTGCCGTCCGTATCTACTACCGCTGGACACCAACTGTTCATGGCGTT
CAGCGGAACGGCCGTAGCATCGAGCACGATGGTGGTCCGACCATGTCATCTGCACGGT
GGACACAGCGACGCTGGTCCGATGGTACCCGGACGCCTGGTACACCCGGTGGTGC
AGAGGTCCGCCTTCACCGGACCCGTCCGACCTATGACAACAGCCAGGAGGCGCGACG
TTATGGTATCATGATCACACCTTAGGGCTAACTCGTCTGAACGTGTATGCGGGTTGGCTGG
CTTCTACCTGTTGCGTGACGACCCGCAACTGTCTCTGATCGAACGTGGCGCCCTCCGTCT
GGTCCGTACGAGGTGGAGCTGGCGATCCAGGATCGAGACCTGCGTCCGGACGGCAGTCTG
GCGTTCCGGACGCTCCGGCGGCAGCTCCAGATTGGCCGGTGGCCCCGTCTGTACTGCCGG
AGTTCTTGGCAGCGTTATTGTTGCAACGGTGCATGGCCGGTGTGGACGTAGAGCC
GCGTCGCTACCGTTGCGTCTGTAATGCGAGCAATTACGTTCTATCGTCTAGCGCAC
ATGGTGTTCGTCCGGTGCCTACCCAGATTGGTGTGACGGCGGTTCTGGACCGCCC
GGCAGCCCTGGACCGTCCGGTCTAGCCCCGGCGAACGTGCGGATCTGATCGTTGAT

TTTCGCGGCCACGAGGGTGCACTTATCGATCTGACGAACGATGCGCCAGTCCCCTTCCGG
ACGGCGATCCGGTCGCCCCCTCCGGCTGACCGTGTGCGCTCCCGTAAGCCTGCCCTA
TGACCGTCGTGTTCCCGGAAGCAGTTCCCGCGCACCCCTGAGAGAAACCCGTTCCGTGCC
GGTGGCGCCCCAGCGCGACCCGCCGTTGGTGTAAACCGAGCACATGGATAGATACGACC
GTCCGAAGCCGATGTTGGGCACGGTCAGCGCGCATGCTGGAGTGGATGGAACCGGTTA
CTGAAACCCCGCGCTGCACAGCACCGAGATCTGGATTACAACACCAAGATGACG
CCCACCCGGTCCATCTCCACTGGTCCAATTCACTGGTCCGACATGCCACCGGATCGGGCGCCCGCAG
GTTCAAGATCCGGTACCGGTGCGTTGTCCGACATGCCACCGGATCGGGCGCCCGCAG
GCGCTGATGAAAACGGTCCGAAAGATAACGTTAGAGCGCTGCCTGGTGAAGTTACCCGTCT
GCGCGGGTGTGATAAACCGGGTGTATTGTGTGGCATTGCCATTCTGGACCACGAAG
ATCACGAGATGATGCGTGAATATGAGGTGGTCGACGACGGTTAaagctt

APPGGLPAVRDPRYLAPIADARAVPRVTPLPRPSRIDSDGSRVLMAPVTQDVVGAGL
GLRTPVWFGAVGGPAGARPSWPGLTVALSRGRPVRWRNLPFRHLLPLDTTVHWAFSGT
GRSIEHDGVPTIVHLHGGHSDAGSDGHPDAWYTPGGARGPRFTGTRPTYDNSQEATLWYHD
HTLGLTRLNVYAGLAGFYLLRDDRSLIERGALPSGPYEVELAIQDRDLRPDGLSLAFPDAPAA
APDWPGGPSVLPEFFGSVIVVNGAAWPVLDEPRRYRLRLLNASNSRFYRLSAHGVRPVPVTQ
IGADGGFLDRPAALDRPLLLAPGERADLIVDFRGHEGALIDLTNDAPVPFPDGDPVAPPADRL
RFRVSLPYDRRVPEAVPPRTLRETPFRAAGGAPARTRLVLTEHMDRYDRPKPMGLGTVERGMLE
WMEPVETPALHSTEWEFYNTTDDAHPVHLHLVQFQILDRAPFTAQDPVTGALSIDIATGSRR
PPGADENGPKDTVRALPGEVTRLRAVFDKPGVFVWHCHILDHEMMREYEVVDDG

Arose

ggatccATGCAGCTGACCCGTATAAAGATAGTCTGACCGTGCGAAAAAACTGATTATTA
AACCGGGTAAACTGAAAATTGCCATGGTTCGACCCGCCTGCGCCTGCATAGTAACGCGC
GAATACCACCGCCCTGTGGACCTATGAAGGCAGTGCCCCGGGCGATTATTGAAGTGCCTG
CGCGGCCAGAAACTGACCGTGGATGGAAAAATCAGCTGACCGGCACCATCCGCTGACC
GCCGCACGCAGCAAACAGGTTCCGGCGATGAAGATCTGTGGCATAGCAATGTGGCAGGT
ATTGCCCTGAGCAGTGATGCAGATGAAGTGTGATTACCCGCATATTGATGATCTGCATGC
ACATACCGTTGTCATCTGCATGGTGGAAAACCAGCCGGATAGTGTGGTTGGACCGATA
ATGTTCATCTGGGTCAACCCAGGCATGCGTTATGAAAATGATATGCGTAGTACCTG
CTGTGGTATCATGATCATGCCATGGCGTTACCGCTTAATGTTTAGTGGCTGGCCGGC
TTTATATTATTCTGTGATGCAGATGATGACAAGTTATGAGTGCAC TGCGTCAAGTCATGGT
CTGACCAAAAAGCACCGAGCCGAAACCCGCTGGAACTGCCGCTGGTTATTCAAGGAT
CGTAATTTCTGACCAAAGCAAATGGCGATCTGACCGCGAACGCTGCATATTGTCAGG
ATCTGCCGCTGAATGCAGCCGGCGATAAAGGTCCGATGGAATTGGCCGTATACCTG
GTGAATGGTACCGTGTGCCGCATGTTGATGTGCGCCGGCCGTATCGTCTGCGCTTACT
GAATGGCAGCAATGCCGTACCTATGCACTGCCCTGGTGGATCAGAATAATGCCGGTT
CGAATGTTACCATGCAGCAGATTGGTACCGATGGCGCCTGCTGGCGCCCTAGACCTGT
TGAACGTATTGTGATTGCACCGCAGAACCGCTGAGCCGTTGGCAATCAGTATGT
CCGCAGGTACCAAACCGCCTTCTGAATACCGCCCTGAGCCGTTGGCAATCAGTATGT
GCCGAAACGGGTGATACCGCAGCGAGCGATCTGCTGGTGTGCCGGAAGTTATGGAATT
CGTGTGCAGGGTAAAAATGCCAGCGCCTGGATCTGCCGGCCCGTTAAGCAAATTAAAG
CCCTGACCCATGGTATGCTGACCTTCTGGAACGGAACTGGAAAAAGTGGCCGGCCAGACCC
CGACCCCTGACGGTATGCTGAATCTGCCGCTGGATATTCACTCCGGATGATCCGGCCAATCCG

CTGCCGAGTGTGGATTATCGTATTACGCCCGCATGTTCTGATGGTAATCCGTGGATGGT
GGCGAAAATGCAACCGAAGTGTGGAATTATTAAATCTGACCGGTATAACCATCCGTTTC
ATATTCACTGGTGCAGTTTAGATGCTCGCACCAGTTGATCCGAGAATTGAA
TATAATGACGGTAAATCGGAATCAGGGTCCGATTATTAGTAAAGGCACCAGTCAGCGCT
GGAAATTGGTGGAAAGATAACCGTGCCTGGACCCCTGGTAAGTGGTAGCATTGCAGTT
CCGTTGAAGGCTTAGCGGCCAGTATATGTATCATTGTCATATTCTGAAACACGAAGATCAT
GAAATGATGCGTCCGTTGTGATGAGCGAAGAAACCATGCCGTTATGCCGAGTCATG
GCCATAGTAGCGGCCATGATCATAATCATTAAaagctt

MQLPYKDSLTVPKLIIKRGKLKIAMVRTRLRLHSELPNTALWTYEGSAPGPIIEVRG
QKLTVEWKNQLTGTIPLTAARSKQVPGDEDLWHSNVAGIALSSDADEVMITPHIDDLHAHTVV
HLHGGKTSPDSGDWTNDNVHHLGQTQACVYENDMRSTLLWYHDHAMGVTRFNFSLAGFY
IIRDADDDKVMSALRASHGLTKKARSPKTPLEPLVIQDRNFLTAKANGDLTGEELLHIVQDPLN
AAGDKGPMEFFGPYTLVNGTVWPHDVPRGPYRLRLLNGSNARTYALALVDQNNSPVPNVT
MQQIGTDGGLLGAPRPVERIVIAPAERVDLIIDFTSVAAGTKLRFLNTALSPFGNQYVPKPGDTA
PSDLLVLPEVMEFRVQGKKCQRLLDPAPLSKFKALTHGDLPHHHNRLVALVEDNTFGMLTF
LELEKVGPAQTPTPDGMLNLPLDIHPDDPANPLPSVDYRITARMFRDGNPWMVAENATEVWNI
INLTGDTHPFHIHLVQFQMLRRTGFPQNFEYNDGKIANQGPIISKGTSQLPLEIGWKDTVRVDPG
ELVSIAVPFEFGSGQYMYHCHILEHEDHEMMRPFVVMSEETMPFMPSHGHSSGHDHNH

Cmeth

ggatccGCTGCAGGTACCCGAATATTACCAAATTGTGGACCCCTGCTGAGTAGTATT
CGAAACTGATTCCGGTGGCCAATGATGCATATCCGGGTGCGATTATTATGAAATTGAAATG
AAACCGGGCCGCTTCAGTTAGTAGCCAGCTGCCGGTACCGCCCTGGCAGGCATTAATA
CCACCTGGGTTATGGTGGTCGTGATCTGGAAATAAGGTCCGCTGGTTCTGGGTTAT
CTGGGTCGACCATTGAAGCACAGAAAGGAAAACCGTTGTGGTAAATTATTAATAACC
TGCCGACCACCCATCCGCTGCAGGCCAGCATTGATCCGACCGTGCCGGACCCCTGCAATGTA
TAGCGATCTGTATGATACCAGCACCAATCCGCCGACCCGATTCATATTGCCGCACCAACCC
CGCATCTGCATGGCGTTTACCGCCCCGAGTTGATGCCATCCGATAGTTGGTTTAC
CCGATTAGTGTGGCGTAAAGGTAGTCATTATGCCACCCGTAGCGGTGCAGCAGCCAATG
AAGCCATTGTTAGCAATCAGCAGCCGGAACCAATCTGTGGTATCATGATCATGCC
ATGGGCATTACCGCCTGAATGTATGCCGGTCTGGCGGTCTGATTTGTTCGCGATCA
GTATGATACCGGACAGCAATCCGCCGACACCGGGCTGAATCTGCCGGCAGGTAATTAT
GAAGTTCCGCTGGTTCTGCAGGATAAACAGTTAATGCAGATGGTACCCGTTATCCGAC
CGTTGGCATTACCGCAGTTCATCCGATTGGATGCCGGAAATTTTGGTATACCCGGTGG
TGAATGGTACCGCCTATCCGTTCTGAGTGTGAACCGCGCCGTTATGTTCTGATTCTGA
ATGGTAGTCAGGCCGTTTATAATCTGTGGTTGATAATGGTCTGGCACCAGTCGTT
ATCTGATTGGCATGGAACAGAGCCTGGTCCGACCCGGTTGCAATGACAAACTGCTGAT
TGCACCGGGTGAACCGCCGATATTATTGTGGATTTGCCGGTATGCAGAATCAGATGCTGA
CCCTGAAAAATAATGCCAACAGCCCCGTATCCGGCGGAAAGGTGGCTTGGTCAGATTAT
GCAGGTGCGTGTGAATCTGCCGATGCCGGTACCGATACCACCAACCCGGCAGCCAATCTG
GTGCTGCCGCAGGTTCCCGTCTGACCGGGCCTACCAATTGCGTGAATGGTATGAAAG
AAGATGTTGATCCGGTTACCGGTCTGCCGATTGATATGAAACTGAATGGCAAATGGTTAAT
GATCTGCCGATTGACGAAACCCCGACCCAGAATGCCACCGAAGTGGCAGTTATTAAATCT
GACCGTGGATGCACATCCGATGCATATGCATCTGGTAAATTCAAGGTTAGCGATGCCAGC
CGTTGATGCCAAAGCCTATACCACCGCATGGCTGAAATATCAGGCAGGTCTGGCAGCAA

ACCGATTCTGGCAAATTTCTGACCAAAGGCCGCCGTGCTGCCGCCTCTGAAGAAATGGCTGGAAAGATAACGCCAAAAGCTATCCGGCGAAATTCTGCGCGTTATTGCAAAATTGAACTGCCGGATGCAGATAGCAATATTCCGGTAGCGGTACCCAGCTGCCGGCAGAATATGTTTATCATTGTCATATTCTGGAGCATGAAGAAAATGAAATGATGCGCCCGTTACCGTGGGTTA
Aaagctt

GSAAGTPNITKFVDPLSSIPKLIPVANDAYPGADYYEIEMKPGRFQFSQLPVTALAGINTTWGYGGRDLANKGPLVFLGLGPTIEAQKGKTVVKFINNLPTTHPLQASIDPTVPDPAMYSDLYDTSTNPPTPIHIGRTTPHLHGGFTAPQFDGHPHSWFTPISDGGKGSHYATLSGAAANEAFVF SNQQPATNLWYHDHAMGITRLNVYAGLAGLYFVRDQYDTGSNPPTPGLNLPAGNYEVPLVL QDKQFNADGTLFYPTVGITAVHPIWMPEFFGDTPVVNGTAYPFLSVEPRRYRFRILNGSQARFY NLWFDNGLAPIFHILIGMEQSLVPTVAMTKLIAPGERADIIVDFAGMQNQMLTLKNNAKAP YPGGKGGFGQIMQVRVNLPMTGTDTTTPAANLVLQPVRLTGATIMREIVMKEDVDPVTGLPIDMKLNGKWFNDLPIDETPTQNATEVWQFINLTVDAPHMHMLVKFQVSDRQPFDAKAYTTAWLKYQAGLGSKPILANFLTKGPALPPPEEMWKDTAKSYPGEILRVIAKFELPDADSNIPGSG TQLPAEYVYHCHILEHEENEMMRPTVG

Cohnella

ggatccATGACCGTGAAATCCGGCCGATCCGGAAACCATTCCGAAATTGTTGATCCGCTGC CGATTCCGCCGATTCTGAAACCGGTTCGCAAAATTAAATGGCAAACCGTATTATGAAGTGCAGC ATGCGTCAGACCCAGCAGCAGGTTCATCGCGATTCCGGTTACCCCTGTGGGGTTATAA TGGCATGTGGCCGGGTCCGACCATTGCCGACGTCGTAATTGCCGATTAAAGTCATTGG ATTAATGAACTGCCGACCACCAGTCTGCTGACCGTGGATCATAGCATTGGTGCAGAACG ACTGTTAGTGGCAATGGTACCCAGAATTTCGCCCGATGTCGAATGTGGTTCATCTGC ATGGTATTAAATGCCCGGAGAATATGATGGTACCCGGATCAGTGGTATAACCCGGGTCTG ACCCAGATTGGTCCGCTGTTAAAACCGATGTGTTGAATATCGAACGCCAGCCGGCAGC ACAGCTGCATTATCATGATCATGCAATTGGCATTACCGCTGAATGTGTATAGCGGCCTGCA GGGCTTTATTTCGCGATGCAGTTGAAGAAAAACTGCCGTGCCAAAATTCCGTATG AAATTCCGATTCTGATTACCGATCGTCAGTTAACCGATGGTACCCGGATCAGTGGTCTGTTTATCCGAAC ATTGGGAACCGCAGTTCTGGCAAATACCAATTGCAGTTAACCGATGGTAGCAATAGCCGTTTATCA GCTCGTCTGGACCCCTGTGCTGACCTTCATGTTATTGGTACCGATGATGGTCTGATGGAAC GTCCGGCAGCAGTGGAAAGAATTCTGCTGGCCCCGGCGAACGTGTGGATGTTATTGTTGA TTTTACCGGTACCGCGCAAAACCTTACCATGACCAATAGCGCACTGGCCCCGTTGAT GTTGGCCTGCCGCCAATCGAACCCGATGGCCTGATTATGCAGTTCTCGGTGAATCAGCC GATGAGCGGCAAAGATCATAGTCGATTCCGAAACGCCCTGGCACAGCTGAAACGCCCTGAAT CCGAATAAAAGTGGCAAAATTCTGTGATGTTACCGATGGCCTGATTATGCAGTTCTCGGTGAATCAGCC AAGAAGAACCGCGCAATTGTGCTGGGTACCCGTAATCTGACCGGTGTCGTAACCGCGTA TCTGTGGGGTGATAGCACCACCGAAAAACCGGTTCTGGTAGCACCAGAAATTGGCGTCTG ATTAATACCAGCCGGGACCCATCCGATTCTGCATCTGGTGCCTTCTGATTCTGGAT CGTCAGCCGTATGATGTGAAACATTGAAAAACCGGTGAACCTGGTTTACCGGCCGC GTGTGCCGCCGGCAACCTATGAACGTGGTATAAAGATGTGGTGCCTGCCAAATCCGGGTGA AGTTACCCGCATTATTGCACGTTGATTATGCCGGTCTGTATGTGTGGCATTGTCATATT CTGGAACATGAAGAAAATGAAATGATGCGTCGTATGGAAGTGGTGCCTGCCGGCTGTCGA CCAGCGAACCGGCAGCAAGCCGAGTTGGCCGCAGTTCCCGTGCAGAATCGCCATG TTAATCGTCGCCCGGTGTCATGAAATTCCGACAGCAGCATTAAACGCCGTTGCCGTGTCG

CTAA <u>aagctt</u>
MTVNPADPETIPKFVDPLPIPKPVRKINGKPYYEVMRQQTQQVHRDFPVTLWGYN GMWPGPTIAARRNCPIKVHWINELPTSSLTVDHSIFGAEALFSGNGTQNFRPDVRNVVHLHGI NAPAEYDGTPDWYTPGLTQIGPLFKTDVFYEYPNRQPAAQLHYHDHAIGITRLNVYSGLQGFY FIRDAVEEKPLPKIPIYEIPIILITDRQFNPDGSLFYPEHWEPAFLANTIAVNGKLWPYLDVEPRRY RFRILNGSNSRFYQLRLDPVLTFHVIGTDDGLMERPAAVEFLLAPGERVDVIVDFTGQRGKTF TMTNSALAPFDVGLPPNPNTDGLIMQFRVNQPMMSGKDHSRIPKRLAQLKRLNPNKVAKIRDVT LNDGTEREGEEERAIVLLGTRNLTGRTQPYLWGDSSTEKPVLGSTEIWRLINTSPGTHPIHLHLV RFLILDQPYDVEHFEKTGELVFTGPRVPPATYERGYKDVVVRANPGEVTRIIARFVDYAGLYVW HCHILEHEENEMMRRMEVVRPGCRTSEPAASRSWAAPRRQNRHVNRGGVHEIRSSIKRRCR RR

*All sequences were codon optimized and the signal peptides were removed.

Table S4 The evaluation for different positive threshold

Positive threshold	Accuracy	Precision	Recall	F1-score
50%	1.0000	1.0000	1.0000	1.0000
90%	1.0000	1.0000	1.0000	1.0000
99%	0.9990	1.0000	0.9667	0.9831

Table S5 The five-fold evaluation results of the MCOs model

Fold	Accuracy	Precision	Recall	F1-score	ROC-AUC
1	1.0000	1.0000	1.0000	1.0000	1.0000
2	0.9981	1.0000	0.9310	0.9643	0.9861
3	0.9981	1.0000	0.9310	0.9643	0.9492
4	0.9990	1.0000	0.9655	0.9825	0.9977
5	0.9990	1.0000	0.9655	0.9825	0.9858
Avg.	0.9988	1.0000	0.9584	0.9787	0.9838

Table S6 Characterization of the multicopper oxidases

enzyme	Specific activity (U/mg)	K_m (mM)	k_{cat} (s ⁻¹)	k_{cat}/K_m (s ⁻¹ /mM)	Halt-life $t_{1/2} 80^\circ\text{C}$ (min)
Eclac (query)	0.41±0.10	6.91±1.32	1.50±0.28	0.22	23.1±1.2
Sulfur	13.51±0.74	0.94±0.11	17.48±1.88	18.60	156.9±9.0
Salvi	0.61±0.13	6.08±0.77	2.03±0.13	0.34	14.6±1.5
Ypes	0.47±0.06	12.16±3.74	4.76±0.97	0.37	9.6±0.8
Psoli	0.20±0.01	9.78±2.03	0.78±0.11	0.08	27.2±1.4
Faest	0.04±0.01	32.53±5.85	0.76±0.12	0.02	16.0±1.1
Mint	0.04±0.01	12.13±2.51	0.69±0.10	0.06	6.0±0.3
Pph	0.70±0.10	0.54±0.18	0.86±0.20	1.66	44.2±7.0
DSM13(query)	3.12±0.37	0.58±0.05	2.39±0.12	4.12	23.2±1.7
Arose	6.67±0.78	8.81±1.31	13.27±1.76	1.51	10.6±1.2
Scla	13.38±1.23	6.76±1.65	7.30±1.19	1.08	62.2±5.2
Bcece	0.55±0.04	0.33±0.02	1.10±0.02	3.33	68.8±8.3
Cmeth	0.05±0.01	ND	ND	ND	ND
Cohnella	ND	ND	ND	ND	ND
HR03 (query)	1.96±0.09	12.45±1.96	6.22±1.07	0.50	66.3±2.6
Aory	0.37±0.04	1.29±0.35	2.16±0.26	1.95	43.3±4.1
Bfre	2.07±0.04	0.10±0.03	4.46±0.40	47.63	61.1±6.4
Tocean	1.61±0.04	0.71±0.16	2.01±0.46	2.88	65.6±5.1
BMS3	0.001	ND	ND	ND	ND
Talbi	0.006	ND	ND	ND	ND
Acid	ND	ND	ND	ND	ND

ND means not detected.

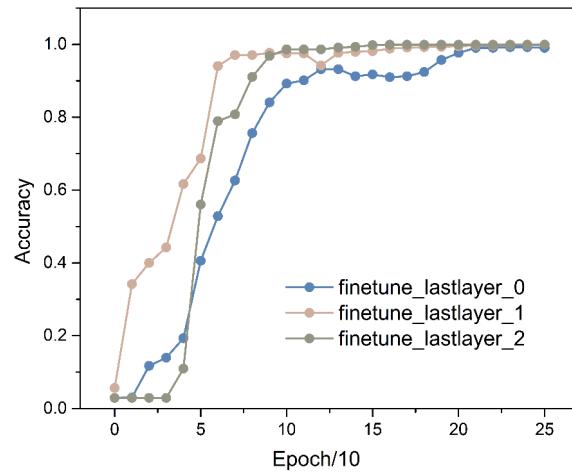


Figure S1 Changes in accuracy during the model fine-tuning process.

After fine-tuning the classification layer of the model, the classification accuracy can approach 1.0. lastlayer_0: last layer; lastlayer_1: last two layers; lastlayer_2: last three layers. Source data are provided as a Source Data file.

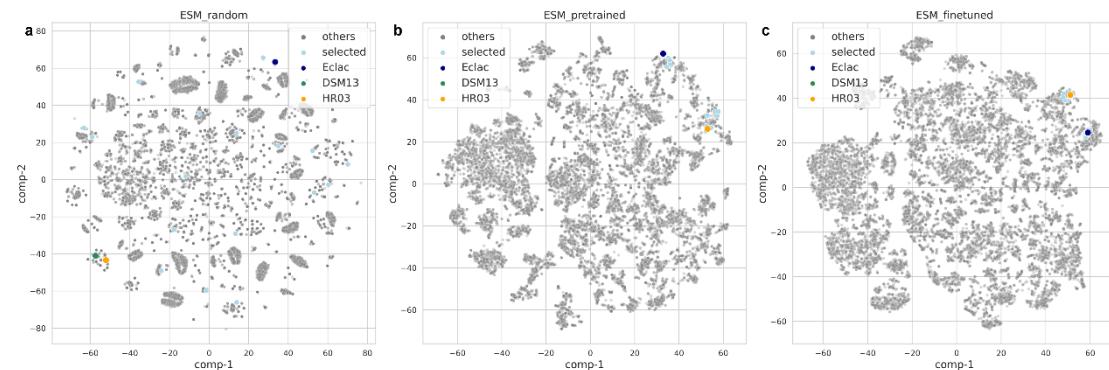


Figure S2 The embeddings of the multicopper oxidases in the ESM-1b model.

a In the random ESM-1b model, all parameters were randomly initialized, resulting in dots being randomly distributed in space. **b** In the pretrained ESM-1b model, the dots of the queried sequences and selected sequences began to cluster. **c** After the model was fine-tuned, the selected sequences became neighbors to the queries. The high-dimensional embeddings were reduced using the t-SNE algorithm. Source data are provided as a Source Data file.

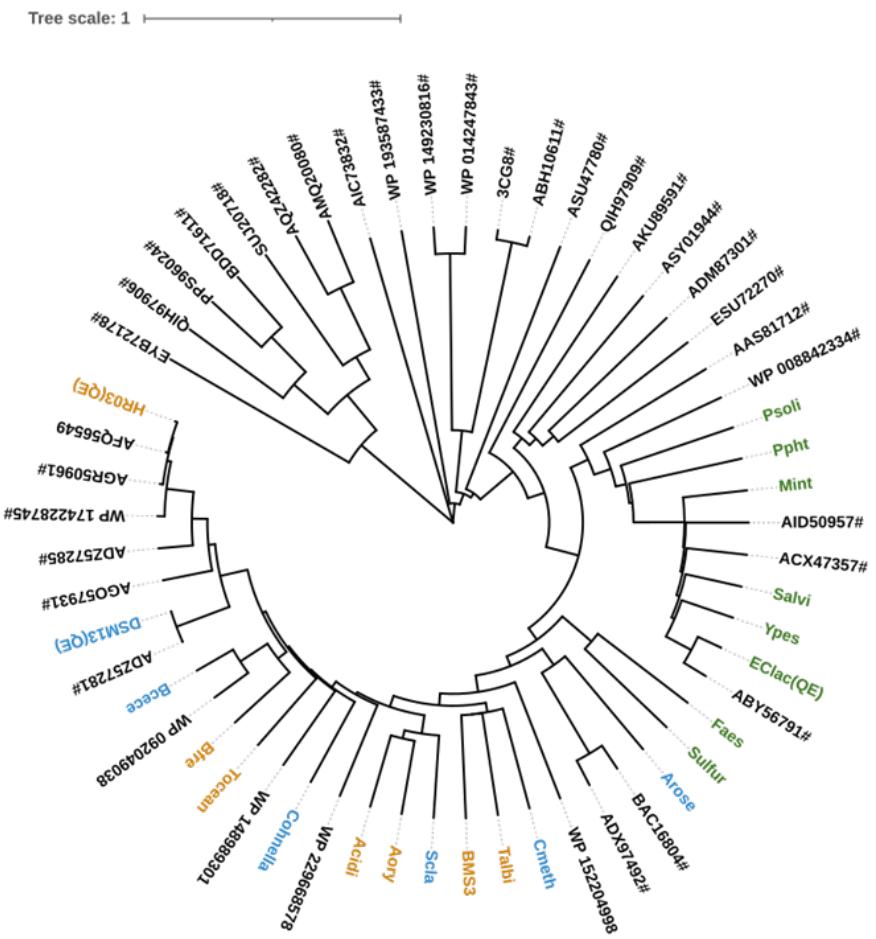


Figure S3 The phylogenetic tree of multicopper oxidases.

Multicopper oxidases retrieved and the enzymes reported (labeled '#') were shown with a UPGMA phylogenetic tree. The multicopper oxidases candidates (18 total, labeling color) are notably diverse compared to the QEs. The protein sequences used for the analysis are provided in Supplementary Table 3. Accession numbers for additional proteins included in the tree are indicated directly on the branches.

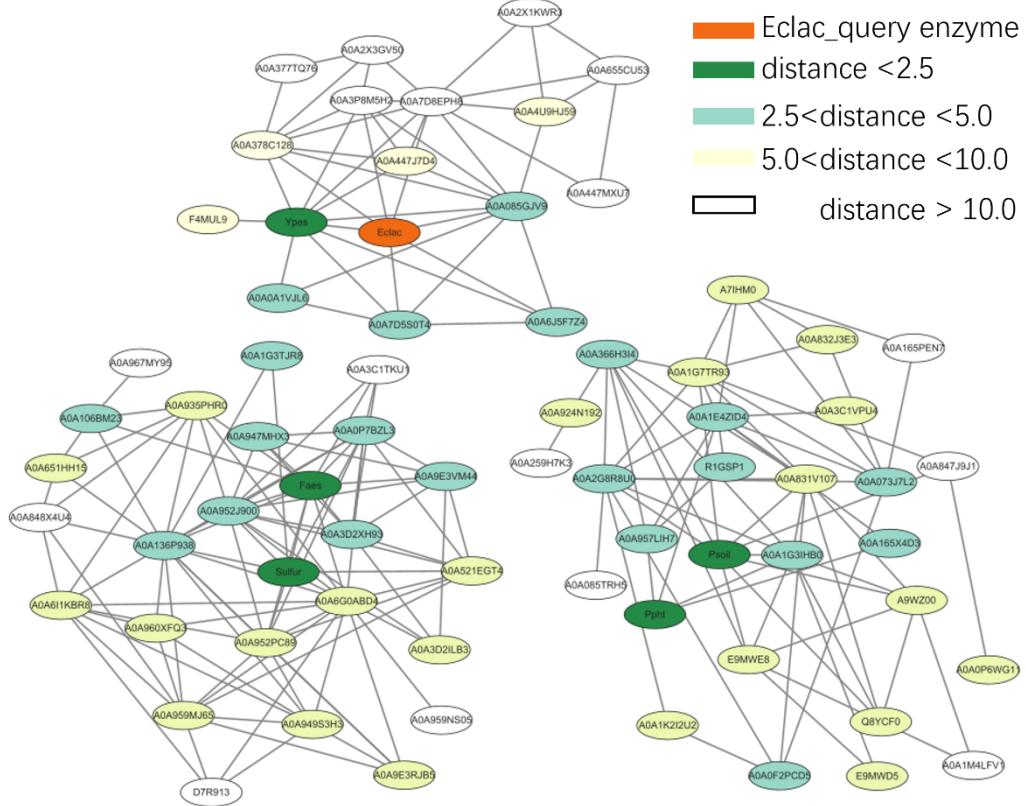


Figure S4 The sequence similarity network (SSN) of Eclac (query enzyme, QE).

The SSN visualizes the relationships between various sequences, emphasizing their Euclidean distances. The network reveals multiple distinct clusters, each representing groups of sequences with closer similarities. Notably, significant variations in Euclidean distances are observed within the same cluster, indicating a spectrum of similarity levels among the sequences, as measured by the Euclidean distance metric. Source data are provided as a Source Data file.

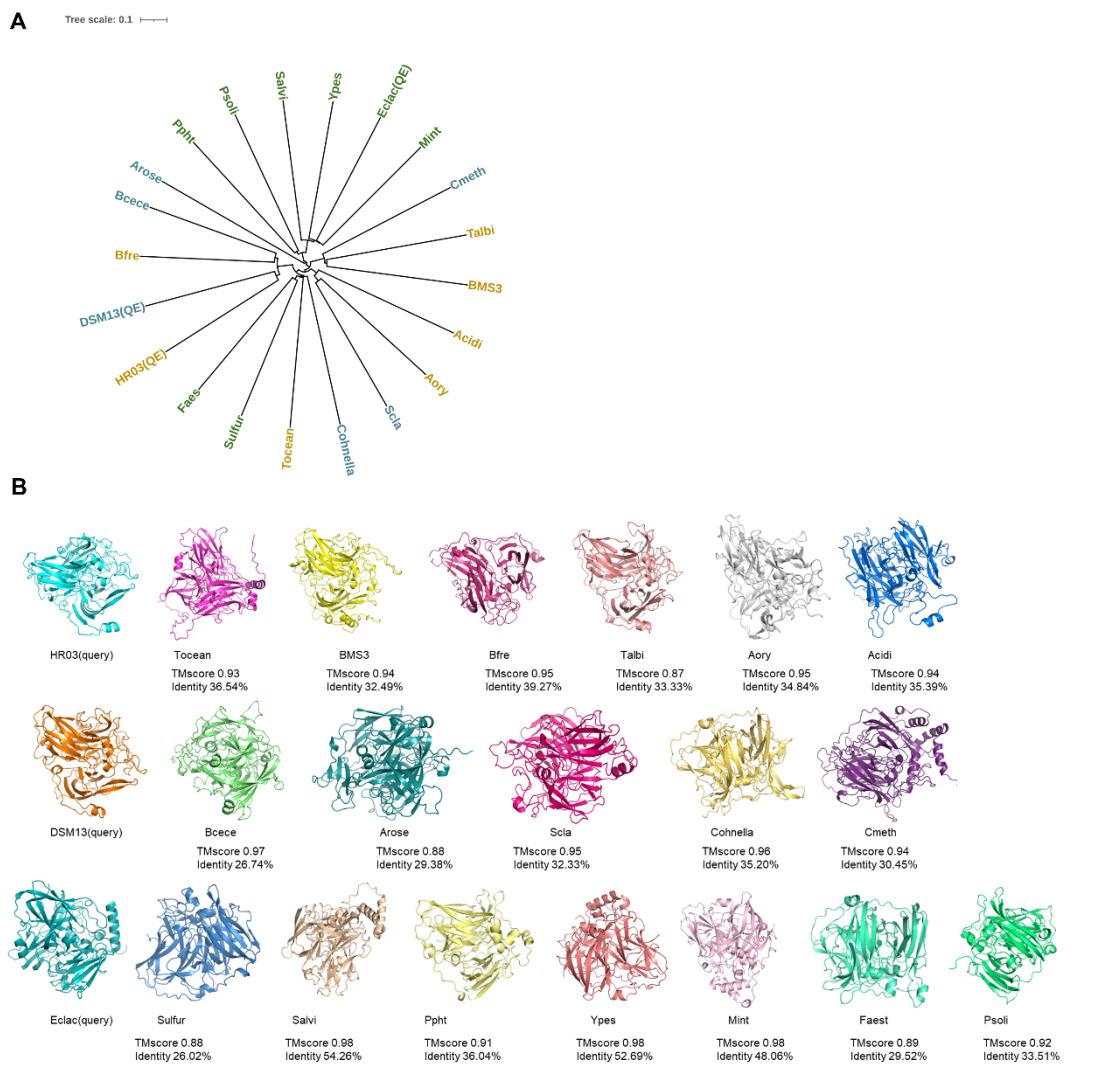


Figure S5 The cluster analysis of all multicopper oxidases characterized in this research.
a Cluster analysis of the protein structure. Each series labeled using different color modes. **b** The AlphaFold2 predicted structures and crystal structures are listed.

Table S7 The enzymes with remote distances tested in this study

Enzyme	Organism	Activity
Eclac(query)	<i>Escherichia coli</i> K-12	
A0A447MXU7	<i>Salmonella enterica</i> I	ND
A0A377TQ76	<i>Klebsiella pneumoniae</i>	ND
A0A2X3GV50	<i>Klebsiella pneumoniae</i>	ND
A0A2X1KWR3	<i>Escherichia coli</i>	ND
A0A1M4LFV1	<i>Brucella</i> sp. 10RB9215	ND

Table S8 The sequences of enzymes with remote distances tested in this Study

A0A447MXU7
ggatccGCTGAACGTCCGGCGCTGCCGATCCCGATCTGCTGACCGCGGATGCGAGCAACCGT ATGCAGCTGATCGTAAAGCAGGGTCAGAGCACCTTCGGCGGGCAAAACCGCACCACCTGG GGCTACAACGGTAACCTGCTGGTCCGGCGGTTCAAGCTGCACAAAGGTAAAAGCGTGACC GTTGATATCCACAACCAGCTGGCGGAGGATACGACCCCTGCACTGGCACGGCCTGGAAATCC CAGGGATCGTGGACGGTGGCCCGCAGGGCATCATTCCGGCGGGTGGCACCCGTACCGTTAC CTTCACCCCCGCAAGCAGCGTGCAGCAACCTGCTGGATTCAACCGCACAAGCACGGCAAAAC CGGCCGTAGGTTGCGATGGGTCTGGCGGGCTGGTTCTGATCGAAGATGATGAAATCCGT AAACTCGTCGGCGGAAACCGTTGGCCACCGTCGTTGCACCGGAGATCACAGCGGTACAG ACCCCTGCTGCGTCGTTGGCCGGATTAA <u>gaattc</u>
AERPALPIPDLTADASNRMQLIVKAGQSTFAGKNATTWGYNGNLLGPAVQLHKGKSVTVDIH NQLAEDETLHWHLGEIPGIVDGGPQGIIPAGGTRTVTFPQQRAATCWIHPHKHGKTGRQVAM GLAGLVIEDDEIRKLRLPAETVGHRRCTGDHSQTLRRWPD*
A0A2X1KWR3
ggatccGCTGAACGCCCGACCCCTGCCGATCCCGACCTGCTGACCAACCGATGCTCGTAACCGT ATCCAGCTGACCATCGGCCGGGGCCAGAGCACCTTCGGCGGAAACCGCACCACCTGG GGCTACAACGGCAACCTGCTGGGCCGGTTAAACTGCAGCGTGGTAAAGCGGTGACC GTTGACATCTACAACCAGCTGACCGAAGAAACCACCCCTGCATTGGCACGGCCTGGAAAGTT CCGGGTGAAGTTGACGGCGCCCGCAGGGGATCATCCGCCGGCGGTAAACGTAGCGTT ACCCCTGAATGTTGATCAGCCGGCGCTACCTGCTGGTTCACCCGACCAGCACGGAAAA CCGGCCGTAGGTGGCGATGGGCCTGGCCGGCTGGTTATCGAGGATGATGAAATCCT GAAGCTGATGCTGCCGAAGCAGTGGGCATAGATGATGTCCGGTTATCGTTAGGATAAA AAATTAGCGCGGATGGCCAGATCGATTACCAAGCTGGACGTTATGACCGCGGCGGTGGCT GGTCGCGATCCGTTGTTAA <u>gaattc</u>
AERPTLPIPDLTTDARNRIQLTIGAGQSTFGGKTATTWGYNGNLLGPAVKLQRGKAVTVDIYN QLTEETTLHWHLGEVPGEVDGGPQGIIPPGKRSVTLNVDQPAATCWFPHQHGKTGRQVAM GLAGLVIEDDEILKMLPKQWGIDDVPVIVQDKKFSADGQIDYQLDVMTAAVGWFAIRC*
A0A2X3GV50
ggatccATGGGTCTGGCGGGCCTGGTTCTGATTGAAGATGAAGAACATGGTCGCTGCTGCTG CCGAAACAGTGGGCATTGATGATGTTCCGGTTATTGTTCAAGATAAAAATCACGGCTGC TGGTGAATTGATTACCAAGCTGGATGTTAGAGCGCAGCGGTTGGTTGGTCGGTACACC TGCTGACCAACGGCGCGCTGTATCCGAACACGCGGCCGAGGGCTGGCTGCGTCTGC

GTCTGCTGAACGGCTGCAACGCGCGTAGCCTGAACTCGCACCAGCGATAAACCGTCCGCT
GTACGTTGCGCTCTGATGGCGGCCGCTGGCGAACCGTTAAAGTTGATGAAC TGCG
GTTCTGATGGCGAACGTTCGAAGTTCTGGTTGATACCTCCGATGGCAAACCGTTGATCT
GGTTACCCCTGCCGGTTAGCCAGATGGCATGGCGATCGCGCCGTTGATAAACCGCAGCCG
GTTCTGCGTGTTCA GCGCTGGTTATCCCTGCTAGCGGTAACACTGCTGGATACCCTGGCTGC
GTTGCCGGCTTACCATCTTAACCGGCCTGACTCAGCGTCAGCTGCAGCTGTCTATGGACC
CGATGCTGGATCGTATGGCTGCCGCTTAAgatcc

MGLAGLVLIEDEESGRLLLKPQWGIDDVPVIVQDKKFTAAGEIDYQLDVMMSAVGWFGDTLL
TNGALYPEHAAPRGWLRLNNGCNARSLNFATSDKRPLYVVASDGGLAEPVKVDELPVLM
GERFEVLVDTSDGKPFDLVTPVSQMGMIAIPFDKQPVLVQPLVIPASGKLLDTLAALPALP
SLTGLTQRQLQLSMDPMLDRMACRR*

A0A377TQ76

ggatccATGGGCCTGGCGGCCGCTGGTTCTGATCGAAGATGAAGAAAGCGGTCGTCTGCTGCTG
CCGAAACAGTGGGCATCGATGATGTTCCGGTGATCGTTAGGATAAAAAATTACCGCTG
CGGGCGAAATTGACTACCAGCTGGATGTATGTCGCGCGGTTGGTTGGTCGGTGATAC
CCTGCTGACCAACGGCGACTGTACCCGGAACACCGCGCCGCGTGGCTGGCTGCCT
GCGCCTGCTGAACGGCTGCAACGCGCTCTGAACCTCGTACCA CGCATAAACGTCCG
CTGTACGTTGTGGCTAGCGATGGCGCTGCTGGCTGAACCGGTGAAAGTTGATGAAC TG
CGGTTCTGATGGCGAACGTTCGAAGTTCTGGTTGATA CCTCTGATGGCAAACCGTTGAT
CTGGTTACTCTGCCGGTTAGCCAGATGGGTATGGCGATCGCACCGTTGATAAACCGCAGC
CGGTTCTGCGTGTTCAGCCGCTGGTTATCCCGCTAGCGCAAAC TGCTGGATACCCGG
CCGTAGTCCGGCGCTGCCGAGCCTGACCGGTCTACCCAGCGGCAGCTGCAGCTGAGCAT
GGGCTCTGATGCGCGTCCGGATGGTCACGCCGGCGCTGATGGTAAGTTGGCGTCCGG
GATGGCCGTAACGGCAGTCGTACGACGGTGCTGGCGCACGACGTATGGCAATATG
CTTCTCGTCGTCA TGAACACGAACCGCGTATTGTCATTGGCTATGGCTTGCGTCAAGCA
CGTGCACTGATTAGCATTACCCCGACCGCTAGCACCGCTAAACCGAGCACCTAgaattc

MGLAGLVLIEDEESGRLLLKPQWGIDDVPVIVQDKKFTAAGEIDYQLDVMMSAVGWFGDTLL
TNGALYPEHAAPRGWLRLNNGCNARSLNFATSDKRPLYVVASDGGLAEPVKVDELPVLM
GERFEVLVDTSDGKPFDLVTPVSQMGMIAIPFDKQPVLVQPLVIPASGKLLDTLGRSPALPS
LTGLTQRQLQLSMGSDARPDGHAGADGEVWRPGDGRNGSRHDGAWRHERHGEYASRRHEH
EPRIRHWSMACRQARALISITPTASTAKPST*

A0A1M4LFV1

ggatccGTGCTGGTTGATTTAGCAACGGCGAAGCTGTTGATCTGGTTACCTATGGTGATAATG
GTAGCGGTGATGGCCTGCATCTGATGCGTTCGCGGTTGATCCGGCATTAGAAGGTCGTGTT
GCTAAACAGCCTGAATCTTAGATGGTCCGGCGCTCCGGATGAAAAACTGTCTGTTCAGC
GTCGTAGCTTTCTTGATGAACCGCATGGCAGAAAACATGAAACTGATGATGCGTCAGCC
GTCCTCTAACCCGACGCATCTGGTGATGATGGATCACATGGAAATGGCTATGGCGG
GTATGGATCATGATATGCATGGTAGCCGTAGCGCGGCTGATGCTGGTCCAGCACTGGATGCA
CTGACCTCTGGTTCAGATGGCGATTGCGATAAACCGTTGACATGGAACGTATCGATGT
TGAAGCTAAATTGGGTTCTGGAAATCTGGGAACTGACTTCCC GTGAAATGGCTCATCCG
TTCCATATCCATGGTGC GTCTTCCGTATTCTGTCTATGAACGGTAAAAAACGCCGGCGCA
TCAGGCTGGCTGGAAAGATAACCGCACTGATCGATGGTAAAGCAGAAATCCTGGTTCACTTC
GATCGTGAAGCGCGCGTCTCACCCGTTCATGTTCCACTGCCATCTGCTGGAACATGAAG
ATGTTGGTATGATGGCGCAGTCGTTACCGTTAAgatcc

VLVDFSNGEAVDLVTYGDNGSGDGLHLMRAVDPALEGRVAKQPESLDGPAAPDEKLSVQRRS
FFFDERMAENMKLMMRQPSSNPHASGDDMDHMEMGSMAGMDHDMHGSRSAADAGPALD
ALTSGVQMAIADKPFDMERIDVEAKLGSWEIWELTSREMAHPFHIHGASFRLSMNGKKPPAH
QAGWKDTALIDGKAEILVHFDRREAARSHPFMFHCHLLEHEDVGMMMAQFVTV

*All sequences were codon optimized and the signal peptides were removed.

The structural characterization of the top-performing MCOs reveals unique features

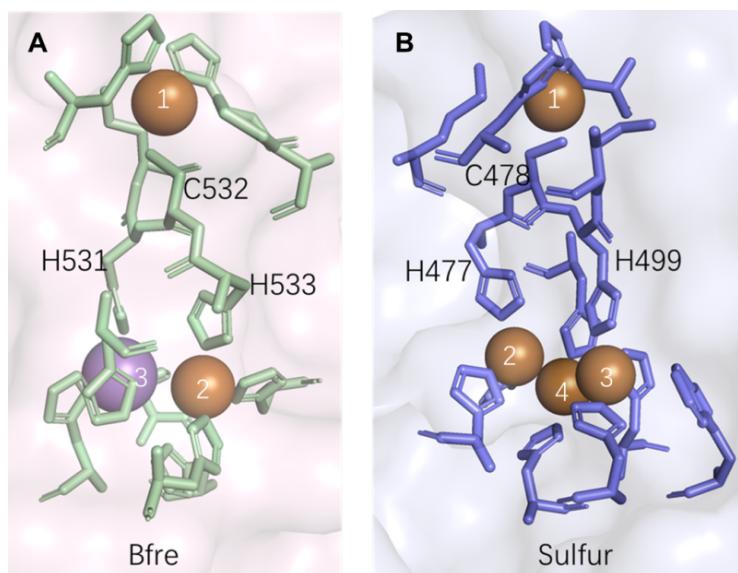


Figure S6 The amino acids coordinated with metal atoms in the active center.

A. The coordinated amino acids of Bfre are listed as follows:

- Cu1: H459, C532, H537, M542
- Cu2: H531, H202, H464
- Mn3: H533, H154, H200

B. The coordinated amino acids of Sulfur are listed as follows:

- Cu1: H424, M488, C478, H483
- Cu2: His477, His166, His429
- Cu3: His164, His126, His479
- Cu4: His124, His427

Table S9 Refinement statistics obtained for the final Bfre model

Data set	Bfre
PDB code	8Z5B
Data collection	
Wavelength (Å)	0.979
Resolution (Å)	29.02 - 1.4 (1.45 - 1.4)
Space group	P 4 ₁ 2 ₁ 2
Cell dimensions	74.59 74.59 230.946 90 90 90
Unique reflections	124241 (10796)
Completeness (%)	96.30 (85.04)
Rmeas (%)	3.0 (38.0)
CC1/2 (%)	99.8 (65.8)
Redundancy	14.3 (10.2)
$I/\sigma(I)$	24.2(1.2)
Refinement	
Rwork (%)	15.8 (2996)
Rfree (%)	17.3 (35.1)
R.m.s.d.	
Bond length (Å)	0.008
Bond angles (°)	1.04
No. of atoms	4834
protein	4143
water	688
Average B factors (Å²)	23.14
protein	21.19
water	34.92
Ramachandran plot (%)	
Favored region	96.83
Allowed region	3.17
Outliers	0

Table S10 Refinement statistics obtained for the final Sulfur model

Data set	Sulfur
PDB code	8Z59
Data collection	
Wavelength (Å)	0.979
Resolution (Å)	50-2.58 (2.62-2.58)
Space group	P 1
Cell dimensions	86.6 103.4 131.9 104.9 104.0 98.4
Unique reflections	126901 (5355)
Completeness (%)	96.5 (81.1)
Rmeas (%)	14.9 (65.6)
CC1/2 (%)	98.4 (76.6)
Redundancy	3.4 (3.0)
$I/\sigma(I)$	7.4 (1.1)
Refinement	
Rwork (%)	19.7 (29.6)
Rfree (%)	23.7 (34.0)
R.m.s.d.	
Bond length (Å)	0.012
Bond angles (°)	1.47
No. of atoms	
protein	28636
water	445
Average B factors (Å²)	
protein	50.75
water	50.86
	43.84
Ramachandran plot(%)	
Favored region	95.89
Allowed region	4.11
Outliers	0

Table S11 Comparison of protein secondary structures and residue interactions between Sulfur and Eclac

Sulfur and Eclac

	Sulfur	Eclac
Salt bridge	72	44
Hydrogen Bond	368	368
Sheet	194	174
Helix	38	56
Coil	115	135

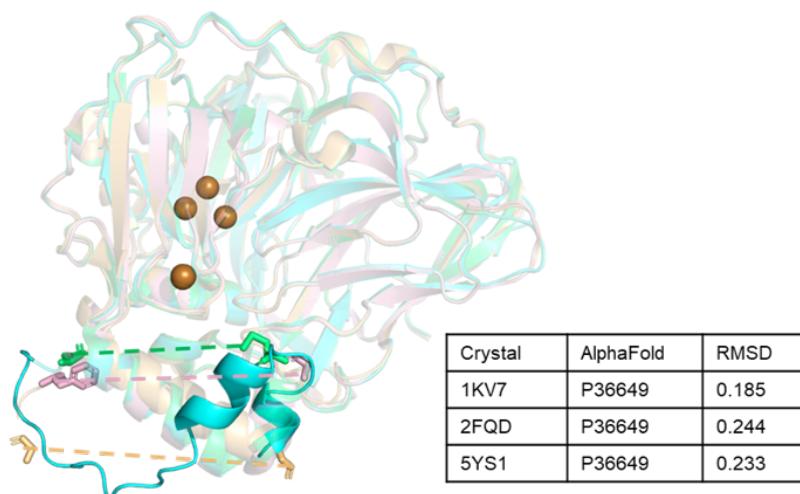


Figure S7 Superimposed alignments of the Eclac AlphaFold structure with its corresponding crystal structures. The structures are colored as follows: 1KV7 in green, 2FQD in pink, and 5YS1 in cyan. The dashed line indicates the regions that are missing in the structures. The relatively low RMSD values of the structural alignments indicate that the AlphaFold predictions are highly accurate.

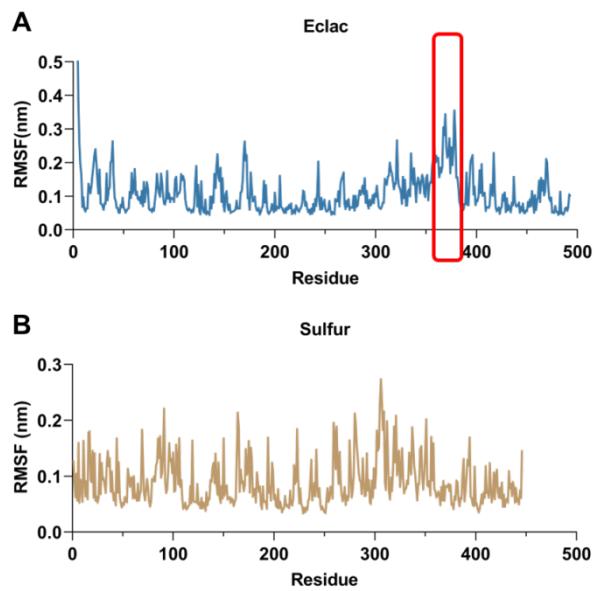


Figure S8 Root mean square fluctuation (RMSF) of C α atoms per residue of (A) SulfurA and (B) Eclac at 37 °C. Source data are provided as a Source Data file.

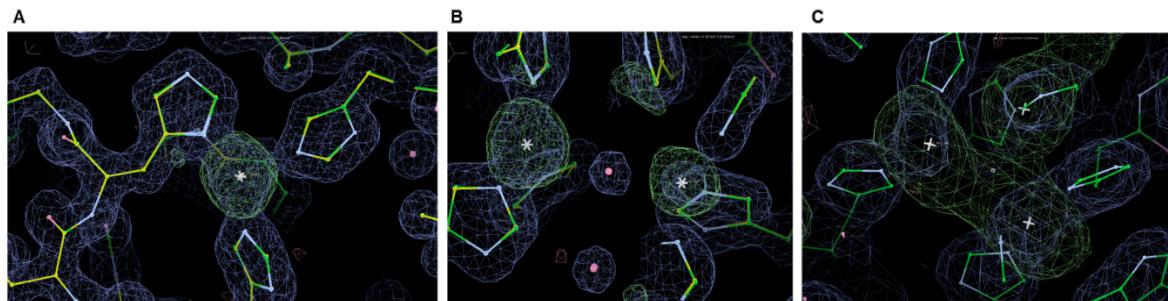


Figure S9 Fo-Fc omit and electron density maps of the crystal structures. A Mn atom in the Bfre crystal structure. **B** Cu atoms in the Bfre crystal structure. **C** Cu atoms in the Sulfur crystal structure.

Characterization of the L-Asparaginases

Development of a deep learning model to identify L-Asparaginases

To further verify the generalizability of the ESM-Ezy method, we manually collected a positive dataset of L-Asparaginases and sourced a negative dataset from Swiss-Prot, adhering to the same quality requirements as the MCO model. During the fine-tuning process of the L-Asparaginases model, we conducted five-fold cross-validation and employed the fine-tuned binary classification model to identify positive L-Asparaginases. Subsequently, we utilized the representations from the fine-tuned model to calculate the Euclidean distance between QEs and positive L-Asparaginase sequences, selecting the top-ranking candidates.

The expression and purification of the recombinant L-Asparaginases

The genes of the L-Asparaginases were synthesized by GeneScript (Nanjing, China) and cloned into plasmid pET28a. Then the strain *E.coli* BL21(DE3) was used for the expression of enzymes. Recombinant strains expressing the L-asparaginase were inoculated in 500 mL LB (kanamycin 50 mg/mL) broth and shaken at 37 °C. The 0.1 mM IPTG was added into the culture when OD600 reached 0.6. The culture was shaken at 16 °C for another overnight for protein expression. Harvested cells were first centrifuged at 8000 rpm for 10 min and re-suspended in 50 mM phosphate buffered saline (PBS, pH 7.4) supplementing with 300 mM NaCl. The cells were disrupted using ultrasonication. After centrifuging at 12,000 rpm for 0.5 h, the supernatant was used as a crude enzyme mixture, which was further purified by nickel-affinity chromatography. Before sample loading, the 2-mL nickel-charged Ni-NTA agarose column (Qiagen, Shanghai, China) was equilibrated with binding buffer (50 mM PBS, 300 mM NaCl, pH 7.4), and elution buffer (50 mM PBS, 300 mM NaCl and 500 mM imidazole, pH 7.4) was used to elute enzyme preparations by gradient elution. Imidazole was removed by ultrafiltration. All of the enzyme purifications were carried out at low temperature.

Enzymatic activity assays

The Nessler's reagent (TMRM Co., Ltd., Beijing, China) method was adopted to determine L-asparaginase activity with modifications. The enzymatic reaction solution consisted of 0.1 -1 mM enzyme, 900 µL of 20 mM substrate L-asparagine of 50 mM Tris-HCl buffer (pH 7.4) and was incubated at 37 °C for 5 min. Add 100 µL of 25% (v/v) trichloroacetic acid (Macklin Co., Ltd., Shanghai, China) solution to terminate the enzymatic reaction. Then 10 µL of the above reaction supernatant coupled with 100 µL of Nessler's reagent were mixed with 890 µL of deionized water. After 10 min, ammonia release quantity was determined by OD₄₃₆ at room temperature. One Unit of L-asparaginase activity was defined as the amount of enzyme releasing 1 µmol of ammonia per minute at pH 7.4 at 37 °C.

Substrate specificity and kinetics parameters

The substrate specificity of the enzymes was determined through measuring the relative enzymatic activity toward L-asparagine , L-glutamine. The Km and Vmax of purified variants were determined by changing the concentration of L-asparagine over a range encompassing 0.5-5 mM at pH 7.4 at 37 °C. Lineweaver-Burk plot was used to calculate the kinetics parameters K_m, k_{cat}

Table S12. The five-fold evaluation results of the L-Asparaginases model

Fold	Accuracy	Precision	Recall	F1-score	ROC-AUC
1	0.9971	0.9688	0.9394	0.9538	0.9883
2	0.9971	0.9412	0.9697	0.9552	0.9999
3	0.9990	0.9706	1.0000	0.9851	1.0000
4	0.9952	0.9118	0.9394	0.9254	0.9769
5	0.9961	1.0000	0.8788	0.9355	0.9890
Avg.	0.9969	0.9585	0.9455	0.9510	0.9908

Table S13 Characterization of the L-Asparaginases

Uniprot ID	Specific activity (U/mg)	Halt-life t _{1/2} 50°C (min)	K _m (mM)	k _{cat} (s ⁻¹)	k _{cat} /K _m (s ⁻¹ /mM)
O34482(query)	239.45±24.32	65.16±4.22	0.39±0.06	57.27±1.93	149.75±23.97
A0A3N5F6J4	469.96±58.86	23.06±1.45	1.09±0.19	112.91±8.49	105.752±12.29
A5EVZ9	262.71±16.46	42.91±4.31	0.63±0.08	82.01±3.31	131.86±18.64
H1D2G7	989.52±125.41	53.13±4.22	1.33±0.24	155.08±16.06	121.30±26.04
A0A2N0LBY7	30.06±3.09	34.06±3.54	ND	ND	ND
A0A085ZTY1	100.12±9.93	7.15±1.04	ND	ND	ND

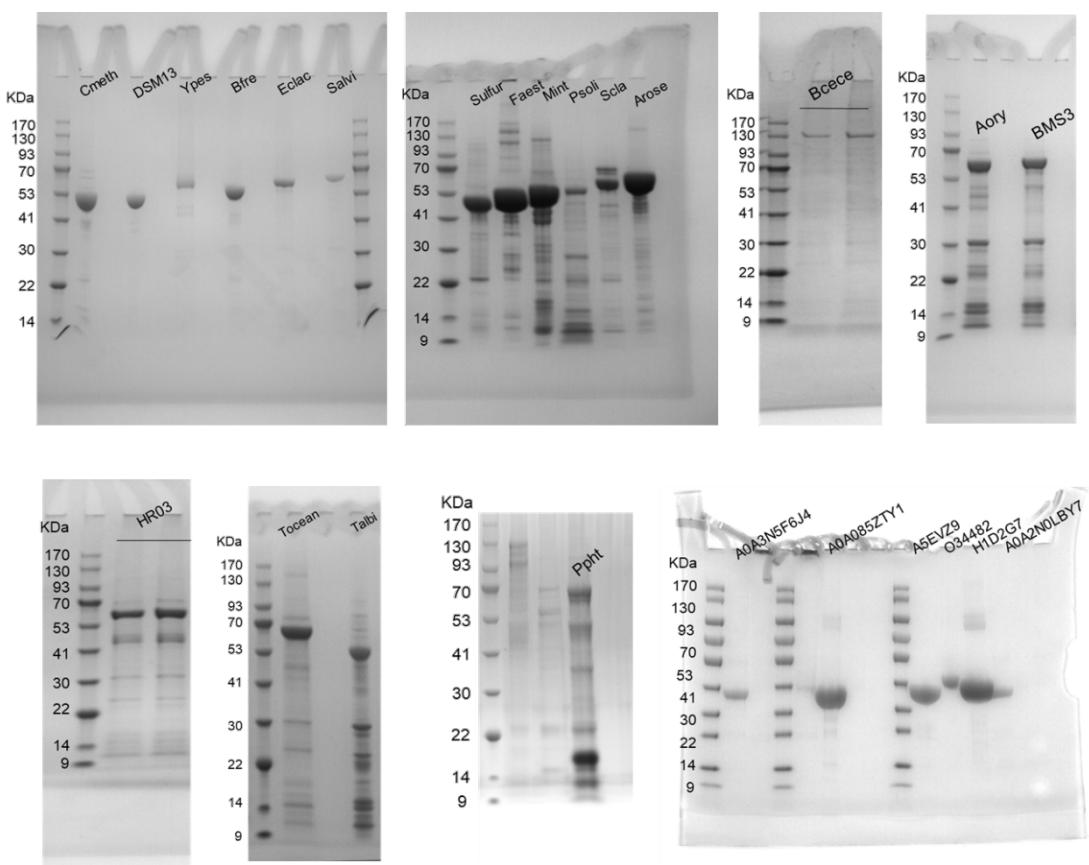


Figure S10 SDS-PAGE analysis of enzymes utilized in this study. The names of purified proteins are labeled on the gel.