1    **Protein Language Models Uncover Carbohydrate-Active Enzyme Function in**

2    **Metagenomics**

3

4    Kumar Thurimella[1,2,3,4], Ahmed M. T. Mohamed[1,2], Daniel B. Graham[1,2], Róisín M. Owens[3],

5    Sabina Leanti La Rosa[5], Damian R. Plichta[1,2,*], Sergio Bacallado[6,*], Ramnik J. Xavier[1,2,*]

6

7    [1]Broad Institute of MIT and Harvard, Cambridge, MA, USA

8    [2]Center for Computational and Integrative Biology and Department of Molecular Biology,

9    Massachusetts General Hospital, Harvard Medical School, Boston, MA, USA

10    [3]Department of Chemical Engineering and Biotechnology, University of Cambridge, Cambridge,

11    UK

12    [4]School of Medicine, University of Colorado Anschutz Medical Campus, Aurora, CO, USA

13    [5]Faculty of Chemistry, Biotechnology and Food Science, Norwegian University of Life Sciences,

14    Ås, Norway

15    [6]Department of Pure Mathematics and Mathematical Statistics, University of Cambridge,

16    Cambridge, UK

17

18    *Correspondence: damian@broadinstitute.org (D.R.P.), sb2116@cam.ac.uk (S.B.),

19    xavier@molbio.mgh.harvard.edu (R.J.X.)

20    **Abstract**

21

22    In metagenomics, the pool of uncharacterized microbial enzymes presents a challenge for

23    functional annotation. Among these, carbohydrate-active enzymes (CAZymes) stand out due to

24    their pivotal roles in various biological processes related to host health and nutrition. Here, we

25    present CAZyLingua, the first tool that harnesses protein language model embeddings to build a

26    deep learning framework that facilitates the annotation of CAZymes in metagenomic datasets.

27    Our benchmarking results showed on average a higher F1 score (reflecting an average of

28    precision and recall) on the annotated genomes of *Bacteroides thetaiotaomicron*, *Eggerthella*

29    *lenta* and *Ruminococcus gnavus* compared to the traditional sequence homology-based method

30    in dbCAN2. We applied our tool to a paired mother/infant longitudinal dataset and revealed

31    unannotated CAZymes linked to microbial development during infancy. When applied to

32    metagenomic datasets derived from patients affected by fibrosis-prone diseases such as

33    Crohn's disease and IgG4-related disease, CAZyLingua uncovered CAZymes associated with

34    disease and healthy states. In each of these metagenomic catalogs, CAZyLingua discovered

35    new annotations that were previously overlooked by traditional sequence homology tools.

36    Overall, the deep learning model CAZyLingua can be applied in combination with existing tools

37    to unravel intricate CAZyme evolutionary profiles and patterns, contributing to a more

38    comprehensive understanding of microbial metabolic dynamics.

39

40   **Introduction**

41

42   Rapid advancements in sequencing technologies have led to an abundance of genomic data,

43   outpacing the capacity to annotate and decipher the functions of these sequences[1]. A significant

44   challenge arises in contextualizing the vast number of unknown functions present in microbes[2,3]

45   and, as a consequence, a substantial fraction of microbial proteins remains unannotated[4–6]. The

46   Unified Human Gastrointestinal Protein (UHGP) catalog alone holds greater than 170 million

47   protein sequences of which 40% lack any functional annotation[2]. Elucidating the function of

48   these sequences has the potential to provide insights into microbial metabolic behaviors and

49   niches within a particular ecosystem, including the dynamics of microbial-host interactions[7–10].

50

51   In microbial genomics, accurate annotations of the biological functions of enzymes is critical, as

52   these molecules have important roles in catalyzing essential biochemical reactions with high

53   specificity and efficiency[11–14]. Carbohydrate-active enzymes (CAZymes) play fundamental roles

54   in various biological processes, including cell structure, signaling, energy storage, and nutrient

55   processing[15–17]. Metagenomic sequencing and functional 'omics have shown that CAZymes

56   support the growth of beneficial microbes in infants by catabolizing human milk oligosaccharides

57   (HMOs)[18,19]. CAZymes have also been found to play a role in the microbiomes of patients with

58   inflammatory diseases like Crohn's disease (CD)[20] and IgG4-related disease (IgG4-RD), in

59   which there is upregulation of glycan-related pathways[21].

60

61   Historically, functional annotation tools have relied on hidden Markov models (HMMs)[22,23] that

62   are built by aligning many amino acid sequences or using sequence homology tools like BLAST,

63   which employs a pairwise alignment strategy between query and target sequences[24,25]. The

64   current state-of-the-art tool for annotating CAZymes, dbCAN2, similarly relies on sequence

3

65 homology or HMMs[26]. While having achieved significant effectiveness in genomic sciences,

66 these methods are not able to assign a biological role to one-third of all bacterial proteins[27].

67 Advancements in deep learning have significantly aided the functional annotation of proteins

68 and comprehension of their diverse functions[28–35]. Protein language models (pLMs), such as

69 those used for structural prediction and other tasks, demonstrate remarkable capabilities in

70 decoding the intricate amino acid language of proteins, which facilitates their functional

71 annotation through a distinct approach compared to sequence-based alignment methods[30,36–39].

72 CAZymes are classified into distinct classes of glycoside hydrolases (GHs), polysaccharide

73 lyases (PLs), glycosyltransferases (GTs) and carbohydrate esterases (CEs). Within a class, the

74 enzymes share a conserved fold, mechanism, and catalytic residues[16]. With this fine grained

75 ontology and a set of distinct enzymatic reactions, CAZymes represent an ideal training dataset

76 for pLMs.

77

78 Here, we present CAZyLingua, the first annotation tool to harness pLMs for the accurate

79 classification of CAZymes. We applied CAZyLingua to gene catalogs derived from human

80 microbiome metagenomic datasets and identified CAZymes implicated in health and disease

81 states. Our first gene catalog was constructed from paired mother/infant metagenomes[40]

82 consisting of ~2,000,000 proteins from which we uncovered ~27,000 CAZymes previously

83 undetected by dbCAN2 or eggNOG. Early persistence of diverse microbial strains in the gut has

84 been linked with metabolic pathways utilizing CAZymes, including breakdown of HMOs and

85 dietary polysaccharides and metabolism of mucin in the colon[41]. CAZyLingua was then applied

86 to a metagenomic dataset derived from patients with inflammatory and fibrosis-prone diseases,

87 including CD and IgG4-RD. We observed that a greater percentage of genes significantly less

88 abundant in CD were predicted to be CAZymes, while in IgG4-RD, we found an expansion of

89 hundreds of CEs in particular. We demonstrate that CAZyLingua achieves high model accuracy

90 compared to standard sequence homology tools and can be used to augment the functional

91    annotation of CAZymes in metagenomic studies, providing valuable insights into the diversity

92    and functional potential of these microbial enzymes.

93

94    **Results**

95

96    *CAZyLingua Model and Performance*

97

98    The CAZyLingua pipeline consists of multiple components (Figure 1a). First, the pLM ProtT5[38] is

99    used to generate embeddings for a given query of amino acid sequences. Second, a quadratic

100   discriminant analysis (QDA) classifier[42], which takes as an input the ProtT5 embedding, is

101   applied to predict whether the query is a CAZyme or not. Finally, if the query is predicted to be a

102   CAZyme, a multiclass classifier is used to make an annotation in the CAZy database ontology,

103   returning either a family or subfamily. The multiclass classifier was built to return probabilities

104   associated with the given family or subfamily annotation and can return a top *k* number of family

105   labels for a given protein sequence.

106

107   We trained CAZyLingua on a subset of the CAZy database[16,43] (Figure 1b). CAZymes were

108   selected from every family, spanning GHs, GTs, PLs, and CEs, to create a representative

109   training dataset. To benchmark our method, we followed a procedure similar to dbCAN2, the

110   current state-of-the-art automated CAZyme annotation tool in the community[26]. We specifically

111   chose the DIAMOND+CAZy option in dbCAN2 as this was the closest representation to

112   BLASTp sequence homology. We performed a taxonomic split on the original CAZy database

113   sequences and selected 3 bacterial genomes with pre-annotated CAZymes in each genome:

114   *Bacteroides thetaiotaomicron, Eggerthella lenta,* and *Ruminococcus gnavus*. We selected these

115   bacteria based on the varying proportions of CAZymes per number of total proteins (*B.*

116   *thetaiotaomicron*: 7.6%, *E. lenta*: 1.1%, and *R. gnavus*: 3.0%) as well as biological relevance: *E.*

117    *lenta* is very prevalent and found in the gut microbiomes of 80% of humans[44], *R. gnavus* is

118    linked to patients with CD and produces a proinflammatory carbohydrate[45], and *B.*

119    *thetaiotaomicron* is one of the most prevalent members of the gut microbiota and dedicates a

120    large portion of its genome to the processing and utilization of carbohydrates[46]. We obtained

121    these exact protein sequences from the CAZy sequence database to use as the reference set

122    for dbCAN2 DIAMOND+CAZy.

123

124    We ran the protein sequences through dbCAN2 and CAZyLingua and evaluated the binary

125    classification task of detecting whether the protein is a CAZyme or not. We combined the results

126    and stratified them into three sets based on whether the protein was predicted by dbCAN2 only,

127    CAZyLingua only, or both. The precision was calculated as the number of true positives in each

128    set divided by the number of predictions made in each set, and recall was calculated as the true

129    positives in each set divided by the total number of CAZymes in each genome (Figure 2a).

130    CAZyLingua alone performed better than dbCAN2 in each measure, but the best benchmarks

131    were in the set of proteins predicted by both tools. We then calculated the F1 score as the

132    harmonic mean of the precision and recall and demonstrated that CAZyLingua outperformed

133    dbCAN2 on each test genome, notably by almost 10% for *E. lenta* (Figure 2b). We examined

134    the predictions by CAZyLingua based on CAZy classes and observed that CAZyLingua was

135    able to label all CE and GT classes in the test genomes (Figure 2c). We evaluated the

136    precision/recall and ROC curves for CAZyLingua and dbCAN2, comparing the output of the

137    decision function from the QDA and the e-value from dbCAN2. Our results showed that

138    CAZyLingua can detect up to 92% of the CAZymes while maintaining a precision of over 80%,

139    while dbCAN2 can detect approximately 82% of the CAZymes at the same precision threshold.

140    CAZyLingua has a higher true positive rate compared to dbCAN2 for this current benchmark

141    (Figure 2d).

142

143    For the CAZyme family classification step, we trained over the entire dataset more than 100

144    epochs, using RayTune[47] to select different random hyperparameter settings and the best of 20

145    different training models. The models were all trained with a cross-entropy loss, and RayTune

146    was optimized to store the model on a metric to minimize loss[48]. The best performing model

147    (lowest loss value) was saved, with the corresponding hyperparameter configuration for any

148    CAZyme family inference. The CAZyme classifier is a four-layer, feedforward neural network

149    (with two hidden layers) with an input of 1024 dimensions (fixed size from ProtT5 embeddings)

150    projected to 256 dimensions then to 512 dimensions to a final classification output layer of 574

151    corresponding to all the unique CAZyme families and subfamilies in our training dataset. We

152    used a hyperbolic tangent (Tanh) as the non-linearity between the different layers. After training,

153    the weights between the first and second layers do not correspond to any interpretable features

154    in the embedding itself (Extended Data Figure 1). When checking a micro-averaged

155    classification accuracy of all the families in the test genomes, CAZyLingua predicted 99.6% of

156    the families accurately, while dbCAN2 predicted 98.2% accurately.

157

158    *CAZyLingua Identifies Horizontally-Transferred Genes as CAZymes*

159

160    We further tested if CAZyLingua would be able to uncover CAZymes in a gene catalog of

161    microbiome samples from mother-infant pairs collected from late pregnancy to one year of

162    age[40]. We predicted CAZymes using CAZyLingua, alongside eggNOG and dbCAN2, on the

163    entirety of the gene catalog, which contained 2,327,970 genes. CAZyLingua predicted 81,498

164    CAZymes, while dbCAN2 and eggNOG predicted 77,614 and 38,862 CAZymes, respectively.

165    We stratified the dataset by number of genes per sample, then by sample month, and split the

166    observations by mother and infant. We calculated the fold change between each method and

167    CAZyLingua based on the genes per sample per month to determine how many more CAZymes

168    were predicted by CAZyLingua. CAZyLingua predicted at least 2-fold more new genes in

169     maternal and infant metagenomes compared to eggNOG and on average 1.2-fold more new

170     genes than dbCAN2 (Figure 3a). When examining the predictions made by CAZyLingua, we

171     observed 27,133 unique CAZyme predictions that were not made by dbCAN2. We distinguished

172     each unique CAZyme by CAZyme class within each sample over each sample month. We

173     observed that our model predicted many more GTs across all the samples in every month

174     (Figure 3b).

175

176     We next focused on a subset of the metagenomic data to specifically look at genes that were

177     found to be horizontally transferred between a mother/infant pair. A previous study performed a

178     sequence homology (BLASTn) analysis on DNA sequences between maternal and infant

179     metagenomes and identified 977 genes with 100% nucleotide identity that were harbored by

180     both maternal and infant species[40], a portion of which were predicted to function in carbohydrate

181     metabolism. Of the 977 genes, 12 were predicted as CAZymes by our model and either not

182     predicted or predicted as an unknown family within a CAZyme class by dbCAN2.

183

184     In order to understand the structural contributions of language models to the general predictions

185     given from ProtT5 and ultimately our pLM classifier, we searched for nearest neighbors between

186     our 12 horizontally-transferred gene embeddings in the CAZy database embeddings using

187     Euclidean distance. After identifying nearest neighbor pairs and extracting the corresponding

188     protein sequences, we computed structural predictions for those proteins using ColabFold[49]. We

189     used FoldSeek[50] to perform a structural alignment between the structures of the predicted

190     protein from CAZyLingua and the nearest protein embedding neighbor in the CAZy database.

191

192     CAZyLingua predicted four GHs, including three belonging to the families 88, 10, and 63, that

193     had a high structural homology to their nearest neighbor in the CAZy database (all with a TM

194     score > 0.50, which indicates a same fold between two proteins[51]). In contrast, when evaluating

195     sequence homology (BLASTp) between the amino acid sequences of the three proteins and the

196     nearest neighbor in the CAZy database, we found that between both sets of sequences the

197     sequence identity was lower than 35%, and for GH88 and GH63 the coverage was less than

198     30% (Figure 3c). Given these metrics, this suggests that CAZyLingua is able to predict

199     CAZymes incorporating structural homology, despite the lack of any amino acid sequence

200     homology.

201

202     The fourth GH predicted was given the annotation of GH43_18 when evaluating the ProtT5

203     nearest neighbor, while CAZyLingua classified it as a GH33 (Figure 4a). We sought to explain if

204     the classification of a GH33 was based on specific features of the unknown CAZyme. We first

205     evaluated the neighborhood of genes around the unknown CAZyme to establish if it exists in a

206     functional polysaccharide utilization locus (PUL). We found several canonical PUL features,

207     including several regulatory elements related to carbohydrate metabolism: a hybrid two-

208     component system (HTCS), TonB-dependent receptor (SusC homolog), and contiguous

209     substrate-binding lipoprotein (SusD homolog) (Figure 4b). In addition to this unknown enzyme

210     mapping to a PUL, we established the presence of a lipoprotein signal peptide in the enzyme

211     through SignalP[52]. We then explored the link between several functional sites in the GH33 and

212     the corresponding embedding generated by ProtT5. To do so, we created a sliding window of

213     10 amino acids and created more distant substitutions of the original sequence within that

214     window based on the BLOSUM62 distance. Substituting areas near the signal peptide

215     corresponded to the greatest losses in the CAZyLingua predictive value of a GH33. The first 20

216     amino acids that correspond to a signal peptide were used in a homology search, and in all

217     BLAST metrics, the signal peptide showed stronger homology to GH33: a combined percent

218     identity and coverage of 64.2% for GH33 and 55.0% for GH43_18, providing stronger evidence

219     for its classification as a GH33 (Figure 4b).

220

221     To determine if there was any structural homology between our unknown CAZyme and the

222     GH33 family, we used ColabFold[49] to fold our protein and ran a structural search with 3D crystal

223     structures found in the PDB25 database using DALI[53]. Our unknown protein had several

224     matches, with two in the top five matches being GH33-like enzymes, namely a neuraminidase

225     and a sialidase. After structurally aligning[51] our unknown structure with the neuraminidase and

226     the sialidase crystal structures, we observed that the predicted GH33 shared significant

227     structural homology (TM score > 0.5) with both. The sequence homology (BLASTp) between

228     the amino acid sequences pairwise with the unknown protein revealed sequence identities

229     <36% and coverages <31% (Figure 4c).

230

231     *Analysis of Enriched CAZymes in Inflammatory Disease Metagenomic Gene Catalogs*

232

233     We next focused our attention on applying CAZyLingua to two metagenomic datasets derived

234     from patients with inflammatory and fibrosis-prone diseases: one from 68 CD patients and 34

235     control subjects [54] and another from 58 IgG4-RD patients and 165 healthy controls[21]. Both of

236     these disease states have unique microbial signatures potentially underlying pathologic

237     mechanisms.

238

239     To investigate disease-associated genes that may be unannotated CAZymes, we first used a

240     linear model against the CD gene catalog[55,56] (Methods) and identified 3,499 genes that were

241     significantly more abundant (two-sided *t*-test, $p < 1 \times 10^{-2}$, log fold change > 2) and 30,125 genes

242     that were significantly less abundant (two-sided *t*-test, $p < 1 \times 10^{-2}$, log fold change < -2) in CD.

243     Among these, CAZyLingua predicted 30 more abundant genes and 569 less abundant genes to

244     be CAZymes (Figure 5a, Supplementary Table 1). Given the ~10-fold difference between more

245     abundant genes in controls versus CD, we observe many more glycan-related pathways

246     associated with health compared to CD.

10

247

248    Following the same analysis procedure, we built a linear model for a differential gene

249    abundance analysis for IgG4-RD metagenomes. We stratified genes based on the same

250    criteria. Compared with the CD dataset, we noticed a higher proportion of genes were

251    significantly more abundant in IgG4-RD compared to a healthy state. We observed 9,225 genes

252    that were significantly more abundant compared to 7,284 genes that were significantly less

253    abundant in IgG4-RD. CAZyLingua predicted 65 more abundant and 87 significantly less

254    abundant CAZymes in IgG4-RD (Figure 5b, Supplementary Table 2).

255

256    We then broadened our focus to all the CAZymes in the IgG4-RD dataset, irrespective of their

257    significance to disease from the linear model. CAZyLingua predicted 437 CAZymes that

258    dbCAN2 did not. Specifically in IgG4-RD, there was a higher number of CEs that only

259    CAZyLingua predicted. CE sequences comprise only 4% of all the sequences in the CAZy

260    database; the low representation of certain sequence examples can pose a challenge for

261    sequence homology tools, which may explain the lower number of hits identified by dbCAN2. In

262    our set of genes predicted by CAZyLingua only, we observed that ~34% were CEs. Families of

263    CEs that were particularly represented included CE1, CE3, CE4, and CE12 (Figure 5c). All of

264    these families share SGNH (Ser-Gly-Asn-His) hydrolase activity, which is a conserved structural

265    feature of the enzymes in these families, suggesting that these enzymes may have low

266    sequence homology but higher structural homology within each class[57–59].

267

268    The increase in annotations by CAZyLingua for these specific CE families may be due to the

269    unique structural features of the families that otherwise would be hard to annotate by traditional

270    sequence homology methods. Given the distinct set of CAZyme families that CAZyLingua was

271    able to predict, we sought to determine the extent of overlap between CAZyLingua predictions

272    and the set of CAZymes that dbCAN2 annotated. To learn about the binary classification of

273   CAZyme/non-CAZyme given by the QDA predictions and the results from dbCAN2, we varied

274   the QDA decision boundary. We calculated the percentage of CAZymes that CAZyLingua

275   labeled as CAZyme that dbCAN2 also predicted against the percentage of the entire IgG4-RD

276   gene set that CAZyLingua labeled as CAZyme. Our QDA model was benchmarked where ~5%

277   of the dataset was labeled CAZyme by CAZyLingua and that represents ~60% of all the genes

278   that dbCAN2 also predicted as CAZyme. At ~30% of the dataset being labeled as CAZyme by

279   CAZyLingua, we captured ~80% of all the dbCAN2-predicted CAZymes. As we relaxed our

280   decision boundary and increased the number of genes in the dataset CAZyLingua labeled as

281   CAZyme, we observed a relatively linear relationship between the genes labeled as CAZyme by

282   both dbCAN2 and CAZyLingua (Figure 5d). This linear relationship describes a relative

283   discordance between the annotations from the two different tools. The divergence of

284   annotations generated by CAZyLingua compared to dbCAN2 can add to existing CAZyme

285   annotations in the analysis of large metagenomics studies.

286

287   **Discussion**

288

289   In this study, we introduced CAZyLingua, a novel approach that leverages pLMs to enhance the

290   identification and functional annotation of CAZymes in metagenomic datasets. Our method

291   mitigates the ongoing challenge of assigning functions to the vast array of unannotated

292   microbial enzymes within these datasets, shedding light on their potential roles in various

293   biological processes. The use of pLMs has emerged as a powerful tool for unraveling protein

294   functions in microbial genomics[28–30], and our results further emphasize their efficacy in this

295   context. When compared with traditional sequence homology, CAZyLingua improved the F1

296   score of classifying a protein as a CAZyme by 6.1% for each of the benchmarked test genomes

297   with gold standard annotations.

298

299    CAZyLingua's efficacy is evident in its successful identification of previously undiscovered

300    CAZymes within a longitudinal microbiome dataset of mother-infant pairs. We detected over

301    27,000 unique putative CAZymes that were missed by dbCAN2. Furthermore, our identification

302    of horizontally-transferred CAZymes between mothers and infants highlights the ability of

303    CAZyLingua to uncover potentially crucial enzymatic functions that traditional sequence

304    homology methods might overlook. When investigating GHs that were missed by dbCAN2, we

305    noticed that these GH structures shared low sequence homology (sequence identity < 40%) to

306    the most homologous protein in the embedding latent space. Our analysis of structural

307    similarities between CAZyLingua-predicted enzymes and GH structures highlights the potential

308    of CAZyLingua to predict enzyme functions based on structural conservation (TM score > 0.5),

309    thereby offering insights into their catalytic roles. We note that these findings are based on

310    structural predictions from ColabFold, not crystal structures or experimentally validated

311    enzymes. One advantage to our choice of ColabFold as a structural prediction tool is that the

312    process of generating a prediction is heavily dependent on a multiple sequence alignment

313    (MSA) between an unknown sequence and a large reference of sequences. The goal of using

314    ColabFold over popular pLM- based structural prediction tools (e.g., ESM-fold, OmegaFold) was

315    for there to be less of a bias between predictions based on embeddings in a process similar to

316    CAZyLingua and how ProtT5 may be trained versus a standard MSA.

317

318    We focus on an example of a horizontally-transferred GH33 that was not predicted by dbCAN2,

319    eggNOG, or a nearest neighbors search using ProtT5 in the CAZy database. Upon using

320    ColabFold to fold this GH33, we performed a sensitive structural search using DALI[53] against

321    experimentally-characterized crystal structures (PDB25) and found the top hits to include other

322    GH33 enzymes (a sialidase/neuraminidase), with significant structural homology (TM score >

323    0.5, Z score > 2).  A recent study examining the early colonization of microbes in a murine

324    model[60] highlights an example of vertical transmission of a GH33 sialidase (NanH) between

325    dams and pups. The NanH gene is triggered by sialylated host glycans and aids in the early

326    colonization of *Bacteroides fragilis*. The putative GH33 discovered by CAZyLingua that was

327    transmitted between a maternal *Alistipes finegoldii* strain and an infant *Alistipes putredinis* strain

328    might exhibit similar properties as NanH and could be part of a mechanism to aid in the

329    establishment of *Alistipes putredinis* in the infant gut. Again, sequence homology between our

330    putative GH33 and NanH was low (33.93% identity, 26% coverage) despite a similar predicted

331    function, indicating that existing sequence homology methods might have overlooked the

332    putative GH33 as a functional homolog. This highlights the strengths of pLMs as alternative

333    tools to augment functional protein homology discovery.

334

335    We then extended the utility of CAZyLingua to metagenomic datasets from patients with CD and

336    IgG4-RD. Both diseases share pathological features of fibrotic lesions despite having distinct

337    clinical presentations. Patients with CD have been shown to have lower microbial diversity and

338    carbohydrate utilization pathways in their gut microbiota[61–63]. Unique microbial signatures have

339    been strongly associated with IgG4-RD, and those signatures included genes linked to

340    carbohydrate metabolism[21]. Our initial analysis focused on genes that were upregulated in

341    IgG4-RD, where we found a distinct set of CAZymes using CAZyLingua. Investigating the

342    taxonomy of those genes, we found several from *Streptococcus* species that are typically found

343    in the oral cavity. In the previous study[21], many *Clostridium* and typically oral *Streptococcus*

344    species were overabundant in the disease phenotype while *Alistipes* and *Bacteroides* species

345    were depleted. Six of the top 20 (30%) putative CAZymes predicted by CAZyLingua mapped to

346    *Streptococcus mutans,* and we observed that many genes from this microbe were upregulated

347    in disease. We observed enrichment of CEs within this species and postulated that there may

348    be several CAZymes that help *Streptococcus mutans* adapt to an ecological niche in the

349    gastrointestinal tract of patients with IgG4-RD.

350

14

351 CEs themselves were sparsely populated in our training dataset for CAZyLingua and similarly in

352 the CAZy database of sequences. Due to the imbalance of this class of enzymes, we postulate

353 that sequence homology may fail to annotate these enzymes. During our training procedure, we

354 use a weighted cross entropy loss, where the weights are proportional to the number of training

355 examples for a given CAZyme family or subfamily. By allowing a more stringent penalty on

356 incorrectly annotating a rare family, we are able to predict more rare families like CEs through

357 CAZyLingua.

358

359 The implications of our findings extend beyond the specific datasets analyzed in this study.

360 CAZyLingua's demonstrated ability to accurately predict CAZymes has broader implications for

361 deciphering the functional potential of microbial communities. A similar procedure of fine-tuning

362 pLM embeddings can be broadly applied to other enzyme classes and protein domains to aid in

363 functional annotation. As an ever-growing number of metagenomic datasets become available,

364 the incorporation of deep learning tools like CAZyLingua into existing methods offers a

365 promising avenue for comprehensive and accurate functional annotation.

366

367 **Methods**

368

369 *CAZyme training dataset curation*

370

371 The CAZy database found at http://www.cazy.org/IMG/cazy_data/cazy_data.zip is cataloged by

372 the dbCAN tool maintainers and a fasta file is available at

373 https://bcb.unl.edu/dbCAN2/download/. We downloaded the CAZy database as of August 06,

374 2022 containing 2,428,817 sequences as it was the latest version that was available for when

375 we began training the model. We chose to focus on the four main classes CAZymes: 173

376 families and 177 subfamilies in glycoside hydrolases (GHs), 115 families in

377    glycosyltransferases (GTs), 20 families in carbohydrate esterases (CEs), and 42 families and 60

378    subfamilies in polysaccharide lyases (PLs). We removed everything that did not belong to one

379    of these families and any sequences that were larger than 5000 amino acids in length to prevent

380    GPU out of memory errors when generating embeddings. The entire number of remaining

381    sequences was 2,413,796: 1,221,013 in GH, 1,027,247 in GT, 122,413 in CE, and 43,123 in PL.

382

383    Using the CD-HIT software tool[64], we clustered our CAZy database at 60% sequence identity.

384    CD-HIT returns a representative sequence for a given cluster. The clusters were created such

385    that, in the resulting database (nr.CAZy.60.fasta), no two sequences had a sequence similarity

386    greater than 60%. The resulting database preserved all of the original families and subfamilies

387    while reducing the redundancy in the database. The database in nr.CAZy.60.fasta contained

388    232,736 sequences, of which 92,385 sequences were in GH, 125,240 in GT, 10,177 in CE, and

389    4,934 in PL.

390

391    Following the curation of the CAZy sequences, we used ProtT5[38] to generate embeddings for

392    each of these sequences using a V100 GPU. We stored the embeddings in h5 files, following

393    the hierarchical data format (HDF). This embedding database served as the training dataset for

394    both of the classifiers in CAZyLingua.

395

396    *Quadratic discriminant analysis training and testing*

397

398    To build the CAZyme/non-CAZyme binary classification step in the CAZyLingua pipeline we

399    modeled the embeddings from the CAZy training dataset as our positive case (CAZyme) and

400    used a combination of data from protein families database Pfam and the Kyoto Encyclopedia of

401    Genes and Genomes (KEGG) to construct our negative examples (non-CAZyme). We started

402    with the 1,296,280 Pfam seeds as a dataset from which to construct negative examples. Pfam

16

403    seeds serve as the basis for hidden Markov model (HMM) profiles and are highly curated to

404    span a diversity of domains[65]. This dataset has been previously described as building the HMMs

405    that contribute to greater than 75% of all the functional annotations of Uniprot sequences in

406    Pfam[28]. We additionally supplemented the negative examples with 3,435 enzymes from KEGG

407    that were non-CAZymes using the KEGG Enzyme database[66].

408

409    In order to create a set of negatives on which to train, we used the ultra-sensitive parameter of

410    DIAMOND[67] in the BLASTp setting between the Pfam seeds against the CAZy database and

411    then the KEGG enzymes against the CAZy database. We removed any Pfam seeds or KEGG

412    enzymes that were listed as hits from the DIAMOND output. The remaining 56,244 Pfam seeds

413    and 3,429 KEGG non-CAZyme enzymes were combined to create a non-CAZyme dataset. We

414    sampled 5,000 CAZymes from nr.CAZy.60.fasta spanning all families and subfamilies in each

415    class as our positive example.

416

417    We built our model using scikit-learn[42], importing the function QuadraticDiscriminantAnalysis

418    with the store_covariance parameter selected as true. We used the library skops to pickle and

419    save the state of our trained model. For a given set of embeddings, the QDA classifier will label

420    them as CAZyme or non-CAZyme and store the results of the CAZy embeddings in an h5 file.

421

422    To model our QDA, we model the distribution of each embedding whether it is a CAZyme or not

423    a CAZyme. These form our two classes $c$: CAZyme and non-CAZyme.

424

425    We model the prior probability of a class $c$, $P(y = c)$ by the empirical proportion of training

426    samples in that class. The conditional probability of a protein's embedding $x \in \mathbb{R}^{1024}$ given its

17

427   class $c$, $P(x \mid y = c)$ is modeled by a multivariate Gaussian distribution with probability density

428   function:

429

$$P(x \mid y = c) = \frac{1}{(2\pi)^{d/2} \left| \Sigma_c \right|^{1/2}} \exp\left( -\frac{1}{2} \left( x - \mu_c \right)^t \Sigma_c^{-1} \left( x - \mu_c \right) \right)$$

430

431

432   The parameters $\mu_c$ and $\Sigma_c$ for each class $c$ are the maximum likelihood estimators given the

433   training samples in that class. If the training samples are $(x_i, y_i), i = 1, ..., N$, the maximum

434   likelihood estimators are given by

435

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{\substack{i=1 \\ y_i=c}}^{N} x_i$$

436

437

$$\hat{\Sigma}_c = \frac{1}{N_c - 1} \sum_{\substack{i=1 \\ y_i=c}}^{N} (x_i - \mu_i)(x_i - \mu_i)^T$$

438

439

440   where $N_c$ is the number of samples class $c$.

441

442   Predictions for a protein with embedding $x^*$ are made by assigning the class $c^*$ which maximizes

443   the posterior probability, given by Bayes' rule:

444

$$P(y = c \mid x^*) = \frac{P(x^* \mid y = c)P(y = c)}{P(x)}$$

445

446

447   where only the numerator depends on $c$. A decision surface is created for the QDA based on the

448   two classes, CAZyme and non-CAZyme[42].

18

449

450    In constructing ROC curves, the decision function that we used is the logarithm of the posterior

451    probability.

452

453    *Feed forward neural network architecture*

454

455    The final stage in the CAZyLingua model is the multiclass classification for a given CAZyme

456    family based on the embeddings selected as CAZyme from the QDA. The feedforward neural

457    network architecture has three overall layers with two hidden layers. The fixed size input of 1024

458    dimensions from ProtT5 embeddings are projected to 256 dimensions then to 512 dimensions

459    to a final classification output layer of 574, which reflects the number of CAZyme families and

460    subfamilies. We implemented this model using Pytorch Lightning[68] to create a classifier that

461    included all of the training, validation, and testing steps.

462

463    The model used a Cross Entropy Loss from PyTorch[48] with the weights parameter set to

464    balance the number of sequences from the different families and subfamilies. In order to prevent

465    over training on highly represented families, the loss function penalty for a given family was

466    calculated as the inverse of the number of sequences per family. This ensures that if the model

467    is incorrectly labeling a family with very few training examples there will be a stronger penalty in

468    comparison to incorrectly labeling a family with a higher proportion of the training examples.

469

470    *Hyperparameter optimization and neural network training*

471

472    The multiclass classification neural network in the CAZyLingua pipeline was trained using

473    RayTune[47], a hyperparameter tuning library. The hyperparameters that were tested were the

474    size of layer 1, the size of layer 2, the batch size, and the learning rate. In order to find the

19

475    optimal hyperparameters to select the most accurately trained model, 20 models were tested in

476    parallel with random sampled hyperparameters selected by RayTune (Supplementary Table 3).

477    Each model was trained over 100 epochs using the Async Successive Halving (ASHA)[69]

478    scheduler that terminates a model (early stopping) optimized to minimize the training loss.

479    Metrics for the validation accuracy were collected after each epoch, and the testing accuracy

480    was collected after the model was fully trained. Each training model was visualized using

481    TensorBoard[70] (Extended Data Figure 2).

482

| Hyperparameter | Sampling Method | Sampled Values |
|---|---|---|
| Layer 1 Size | Random Choice | (256, 512, 768) |
| Layer 2 Size | Random Choice | (512, 1024, 1536) |
| Batch Size | Random Choice | (127, 256, 512) |
| Learning Rate | Log Uniform Sample | [1e-4 – 1e-2] |

483

484    **Supplementary Table 3. Hyperparameter Tuning.** Training epochs over time to pick the

485    model with the best classification accuracy. Using RayTune, we performed a random grid

486    search of different hyperparameter values and tested 20 models in parallel. We picked the

487    model with the best accuracy and used that as the model for all further inference.

488

489    *Benchmarking of CAZyme/non-CAZyme QDA classifier*

490

491    To benchmark the QDA classifier, we used different metrics to quantify the performance of

492    CAZyLingua to dbCAN2. For the F1 score, we followed a standard formula:

493

$$F1 \text{ Score } = 2 \times \frac{\text{recall} \times \text{precision}}{\text{recall} + \text{precision}}$$

494

495

496  where we define recall and precision as follows:

497

$$\text{Precision } = \frac{\text{True Positives}}{\text{True Positives} + \text{False Positives}}$$

498

$$\text{Recall } = \frac{\text{True Positives}}{\text{True Positives} + \text{False Negatives}}$$

499

500

501  The precision-recall and ROC curves were plotted using sklearn[42] using the

502  precision_recall_curve and roc_curve using the e-values from dbCAN2 and the scores from the

503  decision function of the QDA from CAZyLingua as the target scores.

504

505  We designed two metrics to benchmark the differences between CAZyLingua's predictions,

506  dbCAN2's predictions, and the predictions shared by both methods.

507

508  $\text{True CAZymes per Genome} \in \{\text{B.Theta CAZymes , E.Lenta CAZymes , R.Gnavus CAZymes}\}$

509  $\text{CAZyLingua Only} = \{\text{CAZyLingua Predictions}\}\backslash\{\text{dbCAN2 Predictions}\}$

510  $\text{dbCAN2 Only} = \{\text{dbCAN2 Predictions}\} \backslash \{\text{CAZyLingua Predictions}\}$

511  $\text{Both Predictions} = \{\text{CAZyLingua Predictions}\} \cap \{\text{dbCAN2 Predictions}\}$

512

513  With each of these different sets, we calculated the metric to find the proportion of true

514  CAZymes to all predictions in each genome predicted by each method:

515

516  $\text{Proportion of true CAZymes in genome detected in method} = \frac{\{\text{Method}\} \cap \{\text{True CAZymes per Genome}\}}{|\text{True CAZymes per Genome}|}$

21

517

518    Each method was also benchmarked to find the proportion of annotated CAZymes that were

519    correctly labeled as being CAZymes in each method:

520

521    $$\text{Proportion of predictions in method that are correct} = \frac{\{\text{Method}\} \cap \{\text{True CAZymes per Genome}\}}{|\text{Method}|}$$

522

523    where

524

525    $\text{Method} \in \{\text{CAZyLingua Only}, \text{dbCAN2 Only}, \text{Both Predictions}\}$

526

527    *Gene catalog construction*

528

529    The metagenomes for each disease type (IgG4-related disease[21] and Crohn's disease[54]) and

530    for the mother/infant cohort[40] were assembled into their respective gene catalogs following the

531    same procedure. A quality control check was performed using Trim Galore![71] to remove

532    sequencing adapters and kneadData to remove human reads and trim low quality reads (--

533    trimmomatic-options "HEADCROP:15 SLIDINGWINDOW:1:20 MINLEN:50") to keep reads that

534    were minimum 50 bp long. All the quality controlled reads were assembled using MEGAHIT[72].

535    Each contig had all of the open readings frames predicted using Prodigal[73], and we keep both

536    gene and protein sequences. A non-redundant gene catalog was built with a sequence identity

537    threshold of 95% using CD-HIT[64]. To construct a count matrix, each read was mapped using a

538    Burrows-Wheeler Aligner with at least 95% sequence identity for the length of the read. For

539    determining the taxonomy of each contig, MMseqs2[74] was used with NCBI RefSeq as the

540    taxonomic annotation database.

541

542     The IgG4-RD non-redundant (90% sequence identity) gene catalog consisted of 2,237,319

543     genes from 58 IgG4-RD samples and 165 healthy control samples[21]. The CD non-redundant

544     (90% sequence identity) gene catalog consisted of 5,929,528 genes from 68 CD samples and

545     34 non-IBD control samples[54]. The mother/infant non-redundant (95% sequence identity) gene

546     catalog consisted of 2,327,970 genes, with 74 infants, 137 mothers, and 70 mother-infant pairs.

547     Infants were sampled each month between birth (0 months) and 12 months (and additionally at

548     0.5 months), and mothers were sampled at gestational week 27 (approximately 3 months prior

549     to the birth of the child) and at 3, 6, 9, and 12 months after the birth[40]. Each of these gene

550     catalogs were constructed in each respective prior study and directly utilized in the analysis

551     presented in this paper.

552

553     *Analysis of mother/infant gene catalog*

554

555     The entire mother/infant gene catalog was run through dbCAN2 (diamond blastp -d

556     ${CAZy_reference} -q ${query_file} -o ${output_str}.matches.tsv -e 1e-102 -k 1 -p 2 -f 6) and

557     eggNOG on default parameters. Additionally, embeddings were generated for the entire

558     mother/infant dataset using ProtT5, with CAZyLingua running inference on the entire gene

559     catalog.

560

561     We took the 977 horizontally-transferred gene subset and collected all of the dbCAN2 and

562     CAZyLingua results. We took the 12 genes that only CAZyLingua predicted and performed a

563     structural prediction on each of the protein sequences. We performed a Euclidean distance

564     search between those 12 embeddings and the nr.CAZy.60.fasta database to find the closest

565     embedding and subsequently the CAZyme family. We then used ColabFold[49] to fold each of the

566     12 proteins and their nearest neighbor to generate PDBs for each horizontally-transferred gene

567     and neighbor pair. A structural alignment was computed on each of these pairs using

568    Foldseek[50], which returns the overlapped structures and a TM score for each pair. To compute

569    sequence homology metrics, we selected the "Align two or more sequences" option in the

570    BLASTp suite on the NCBI website

571    (https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastp&PAGE_TYPE=BlastSearch&LINK_

572    LOC=blasthome).

573

574    The putative GH33 and each of the GH33 and GH43_13 in nr.CAZy.60.fasta were ordinated

575    through tSNE (sklearn TSNE package)[42] and plotted using matplotlib[75]. A structural prediction of

576    the putative GH33 was produced from ColabFold[49] and the amino acid residue substitution

577    analysis was done using custom scripts. To search against known, experimentally-characterized

578    structures, the DALI option to pairwise search against PDB25[53] was used. To structurally align a

579    pairwise hit from putative GH33 to a structure from PDB25, we used US-align[51] to generate

580    aligned structures and a TM score.

581

582    *Disease metagenomic differential abundance analysis*

583

584    In each disease gene catalog, linear modeling was used to regress different disease covariates

585    onto each gene in the catalog to find differentially abundant genes (features). An abundance

586    filter was applied to the entire count matrix to remove any genes with <10% prevalence across

587    samples. A zero-inflation was applied to any zeros in the count matrix, where the zero value

588    would be replaced by the minimum non-zero value in the given feature and divided by 2. The

589    fold change was calculated by dividing the mean of the disease group by the control group, and

590    taking the $log_2$ of the value. Each value is $log_2$ transformed and a z-score is calculated for every

591    value in a given feature using the scipy[76] library. A linear model, from the statsmodels[55] library,

592    is then applied to each feature. For IgG4-RD, the metadata covariates modeled were: age, on

593    treatment, rituximab, prednisone, other treatments, sex, and cohort. In CD the variables

24

594    modeled were: age, on antibiotics, mesalamine, and steroids. A significance threshold was

595    established for all of the analyses: we followed a multiple testing adjustment, and p-values were

596    corrected using Benjamini-Hochberg correction, with a false discovery rate (FDR)-corrected p

597    value (q-value) of 0.25. The volcano plots were labeled based on four conditional arguments for

598    the CD and IgG4-RD metagenomic catalogs. For CD, the criteria for the displayed labels were:

599        1.   logFC > 2 and p-value < $1 \times 10^{-5}$

600        2.   logFC < -2 and p-value < $1 \times 10^{-8}$

601        3.   logFC > 3  and p-value < $1 \times 10^{-2.5}$

602        4.   logFC < -4.5 and p-value < $1 \times 10^{-3}$

603

604    For IgG4-RD, the criteria for the displayed labels were:

605        1.   logFC > 2 and p-value < $1 \times 10^{-5}$

606        2.   logFC < -2 and p-value < $1 \times 10^{-3.5}$

607        3.   logFC > 3  and p-value < $1 \times 10^{-2.5}$

608        4.   logFC < -3.5 and p-value < $1 \times 10^{-2}$

609

610    **Acknowledgements**

611

615

616    **Figure Legends**

617

618    **Figure 1. CAZyLingua: a deep learning model used for the classification of proteins as**

619    **CAZymes. a)** The workflow of CAZyLingua starts with raw embeddings from ProtT5 followed by

25

620     the use of those embeddings as input through two classifiers to distinguish 1) whether the

621     embedding was a CAZyme and if so, 2) to which CAZyme family it belongs to. **b)** The training

622     strategy for CAZyLingua began with a 60% sequence identity clustering to remove redundancy

623     from the CAZy database in order to train on distinct CAZymes. The Cross Entropy loss function

624     was applied for training and the loss function that was used included a weighted balancing

625     function to proportionally sample the number of representative sequences per CAZyme

626     class/family/subfamily in the database. This strategy was employed so as not to oversample on

627     highly represented families.

628

629     **Figure 2. CAZyLingua performance relative to the BLAST-based CAZyme annotation tool**

630     **dbCAN2.** CAZyLingua was compared to the dbCAN2 DIAMOND+CAZy annotation tool option

631     (benchmarked with an e-value $< 1x10^{-102}$). A similar procedure as dbCAN2 was followed by

632     picking 3 bacterial strains with manual annotations and varying CAZyme counts per strain. **a)**

633     For predictions by CAZyLingua only, dbCAN2 only, and shared between the two methods, the

634     proportion of correct predictions made by each method (left) and the proportion of true

635     CAZymes made by each method (right) were calculated. **b)** F1 scores (harmonic means of

636     precision and recall) of all CAZyLingua predictions, all dbCAN2 predictions, and all predictions

637     combined, whether shared between the methods or not. **c)** Ground truth CAZymes were

638     stratified by class, and the percentage of accurate predictions per CAZy class from our

639     Quadratic Discriminant Analysis (QDA) binary classifier was calculated. **d)** Precision/recall (left)

640     and ROC (right) curves comparing CAZyLingua to dbCAN2. The output of the decision function

641     of the boundary that was trained for CAZyLingua and the e-value for dbCAN2 were used for

642     target scores.

643

644     **Figure 3. Application of CAZyLingua to metagenomes in paired mothers and infants. a)**

645     Comparison of CAZyLingua to eggNOG and dbCAN2 on a large metagenomics gene catalog

646     from mothers and their infants. Time of the sample is in months relative to childbirth (month 0).

647     Dotted lines represent no fold change. **b)** CAZyLingua predicted 27,133 genes that dbCAN2 did

648     not, shown by CAZy class for all infant and maternal samples at each sample month. Boxplots

649     in **a** and **b** show medians and interquartile ranges (IQRs), with whiskers showing ± 1.5 IQR. **c)**

650     Predicted structures of proteins from CAZyLingua (red) and the protein embedding nearest

651     neighbor (grey) structurally aligned with TM scores, and BLAST metrics for GH88, GH10, and

652     GH63.

653

654     **Figure 4. CAZyLingua distinguishes GH33 CAZyme from nearest neighbors of raw ProtT5**

655     **embeddings. a)** tSNE of (left) ProtT5 embeddings from the GH33 and GH43_18 families and

656     the CAZyme predicted by CAZyLingua (GH unknown) and (right) a segment of the last layer of

657     CAZyLingua. **b)** GH33 protein residues were mutated in a sliding window of ten residues over

658     the entire sequence, and ProtT5 embeddings were generated for each sliding window mutation.

659     Known features are overlaid along sections of the sequence. The probability of the CAZyLingua-

660     predicted classification being a GH33 was calculated for each sliding window mutation (top).

661     The predicted GH mapped to a PUL containing several regulatory elements consistent with a

662     CAZyme (bottom left). BLAST metrics on the predicted GH signal peptide compared with GH33

663     and GH43_18 sequences (bottom right). **c)** Overlays of the predicted GH protein structure

664     generated using ColabFold with a sialidase (top) and a neuraminidase (bottom).

665

666     **Figure 5. Application of CAZyLingua to CAZymes in metagenomes of patients with**

667     **inflammatory and fibrosis-prone diseases.** Genes enriched and depleted in the gene

668     catalogs of patients with **a)** CD and **b)** IgG4-RD selected on the fringe of the volcano plot (see

669     Methods for labeling criteria). **c)** Predicted CEs in the enriched IgG4-RD gene set, stratified to

670     analyze only the genes CAZyLingua predicted. **d)** The proportion of dbCAN2-predicted

671     CAZymes also predicted by CAZyLingua as the decision function between CAZyme/non-

27

672    CAZyme of the QDA classifier in CAZyLingua was varied. The Venn diagram shows the

673    numbers of CAZymes predicted by CAZyLingua, dbCAN2, and both on our current model

674    benchmarks of the QDA.

675

676    **Extended Data Figure Legends**

677

678    **Extended Data Figure 1. Embedding weights from first layer to next, no interpretable**

679    **chemical features.** We extracted the weights ($W$) from the CAZyLingua multiclass classifier

680    between the input layer and first hidden layer, which is a matrix of dimension 1024x256. After

681    applying a transpose to get $W^T$ we multiplied the two matrices, $W \cdot W^T$ which produced a

682    symmetric matrix, $S$ of dimensions 1024x1024. After taking the $diag(S)$ we obtained a vector of

683    size 1024, which is the size of the original embedding from ProtT5. We plotted the values in the

684    vector to visualize if there were any features or positions in specific regions of the embedding

685    that are specific to CAZymes.

686

687    **Extended Data Figure 2. Training runs for finding the best model.** RayTune ran 20 models

688    in parallel over each epoch and pruned any models that began to stagnate or have a decline in

689    training accuracy. The models were evaluated on the metric of minimizing training loss, and the

690    model with the minimal loss was stored as a checkpoint. There were 100 epochs over which

691    training occurred, and the metrics were stored and written to a TensorBoard that produced

692    these visualizations.

693

694    **References**

695

696    1.    Qin, J. *et al.* A human gut microbial gene catalogue established by metagenomic

697          sequencing. *Nature* **464**, 59–65 (2010).

698    2.   Almeida, A. *et al.* A unified catalog of 204,938 reference genomes from the human gut

699          microbiome. *Nat. Biotechnol.* **39**, 105–114 (2021).

700    3.   Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over

701          150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* **176**,

702          649-662.e20 (2019).

703    4.   Rinke, C. *et al.* Insights into the phylogeny and coding potential of microbial dark matter.

704          *Nature* **499**, 431–437 (2013).

705    5.   Perdigão, N. *et al.* Unexpected features of the dark proteome. *Proc. Natl. Acad. Sci. U. S.*

706          *A.* **112**, 15898–15903 (2015).

707    6.   Zhang, Y. *et al.* Discovery of bioactive microbial gene products in inflammatory bowel

708          disease. *Nature* **606**, 754–760 (2022).

709    7.   Cho, I. & Blaser, M. J. The human microbiome: at the interface of health and disease. *Nat.*

710          *Rev. Genet.* **13**, 260–270 (2012).

711    8.   Ma, B., Charkowski, A. O., Glasner, J. D. & Perna, N. T. Identification of host-microbe

712          interaction factors in the genomes of soft rot-associated pathogens Dickeya dadantii 3937

713          and Pectobacterium carotovorum WPP14 with supervised machine learning. *BMC*

714          *Genomics* **15**, 508 (2014).

715    9.   Lozupone, C. A. Unraveling Interactions between the Microbiome and the Host Immune

716          System To Decipher Mechanisms of Disease. *mSystems* **3**, (2018).

717    10. Berg, G. *et al.* Microbiome definition re-visited: old concepts and new challenges.

718          *Microbiome* **8**, 103 (2020).

719    11. Yu, T. *et al.* Enzyme function prediction using contrastive learning. *Science* **379**, 1358–

720          1363 (2023).

721    12. Yeh, A. H.-W. *et al.* De novo design of luciferases using deep learning. *Nature* **614**, 774–

722          780 (2023).

723    13. Tao, Z., Dong, B., Teng, Z. & Zhao, Y. The Classification of Enzymes by Deep Learning.

724     *IEEE Access* **8**, 89802–89811 (2020).

725     14. Li, Y. *et al.* DEEPre: sequence-based enzyme EC number prediction by deep learning.

726     *Bioinformatics* **34**, 760–769 (2018).

727     15. Cantarel, B. L., Lombard, V. & Henrissat, B. Complex carbohydrate utilization by the

728     healthy human microbiome. *PLoS One* **7**, e28742 (2012).

729     16. Lombard, V., Golaconda Ramulu, H., Drula, E., Coutinho, P. M. & Henrissat, B. The

730     carbohydrate-active enzymes database (CAZy) in 2013. *Nucleic Acids Res.* **42**, D490-5

731     (2014).

732     17. Cantarel, B. L. *et al.* The Carbohydrate-Active EnZymes database (CAZy): an expert

733     resource for Glycogenomics. *Nucleic Acids Res.* **37**, D233-8 (2009).

734     18. Charbonneau, M. R. *et al.* Sialylated Milk Oligosaccharides Promote Microbiota-Dependent

735     Growth in Models of Infant Undernutrition. *Cell* **164**, 859–871 (2016).

736     19. Wardman, J. F., Bains, R. K., Rahfeld, P. & Withers, S. G. Carbohydrate-active enzymes

737     (CAZymes) in the gut microbiome. *Nat. Rev. Microbiol.* **20**, 542–556 (2022).

738     20. Porras, A. M. *et al.* Inflammatory Bowel Disease-Associated Gut Commensals Degrade

739     Components of the Extracellular Matrix. *MBio* **13**, e0220122 (2022).

740     21. Plichta, D. R. *et al.* Congruent microbiome signatures in fibrosis-prone autoimmune

741     diseases: IgG4-related disease and systemic sclerosis. *Genome Med.* **13**, 35 (2021).

742     22. Söding, J. Protein homology detection by HMM–HMM comparison. *Bioinformatics* **21**, 951–

743     960 (2004).

744     23. Finn, R. D., Clements, J. & Eddy, S. R. HMMER web server: interactive sequence similarity

745     searching. *Nucleic Acids Res.* **39**, W29-37 (2011).

746     24. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment

747     search tool. *J. Mol. Biol.* **215**, 403–410 (1990).

748     25. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein database

749     search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).

750   26. Zhang, H. *et al.* dbCAN2: a meta server for automated carbohydrate-active enzyme

751        annotation. *Nucleic Acids Res.* **46**, W95–W101 (2018).

752   27. Price, M. N. *et al.* Mutant phenotypes for thousands of bacterial genes of unknown function.

753        *Nature* **557**, 503–509 (2018).

754   28. Bileschi, M. L. *et al.* Using deep learning to annotate the protein universe. *Nat. Biotechnol.*

755        (2022) doi:10.1038/s41587-021-01179-w.

756   29. Gligorijević, V. *et al.* Structure-based protein function prediction using graph convolutional

757        networks. *Nat. Commun.* **12**, 1–14 (2021).

758   30. Kaminski, K., Ludwiczak, J., Alva, V. & Dunin-Horkawicz, S. pLM-BLAST – distant

759        homology detection based on direct comparison of sequence representations from protein

760        language models. *bioRxiv* 2022.11.24.517862 (2022) doi:10.1101/2022.11.24.517862.

761   31. Chowdhury, R. *et al.* Single-sequence protein structure prediction using a language model

762        and deep learning. *Nat. Biotechnol.* **40**, 1617–1623 (2022).

763   32. Madani, A. *et al.* Large language models generate functional protein sequences across

764        diverse families. *Nat. Biotechnol.* 1–8 (2023).

765   33. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,

766        583–589 (2021).

767   34. Koehler Leman, J. *et al.* Sequence-structure-function relationships in the microbial protein

768        universe. *Nat. Commun.* **14**, 2351 (2023).

769   35. Rives, A. *et al.* Biological structure and function emerge from scaling unsupervised learning

770        to 250 million protein sequences. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).

771   36. Heinzinger, M. *et al.* Contrastive learning on protein embeddings enlightens midnight zone.

772        *NAR Genom Bioinform* **4**, lqac043 (2022).

773   37. Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language

774        model. *Science* **379**, 1123–1130 (2023).

775   38. Elnaggar, A. *et al.* ProtTrans: Towards Cracking the Language of Lifes Code Through Self-

776      Supervised Deep Learning and High Performance Computing. *IEEE Trans. Pattern Anal.*

777      *Mach. Intell.* **PP**, (2021).

778   39.   Bepler, T. & Berger, B. Learning the protein language: Evolution, structure, and function.

779      *Cell Syst* **12**, 654-669.e3 (2021).

780   40.   Vatanen, T. *et al.* Mobile genetic elements from the maternal microbiome shape infant gut

781      microbial assembly and metabolism. *Cell* **185**, 4921-4936.e15 (2022).

782   41.   Lou, Y. C. *et al.* Infant gut strain persistence is associated with maternal origin, phylogeny,

783      and traits including surface adhesion and iron acquisition. *Cell Rep Med* **2**, 100393 (2021).

784   42.   Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *the Journal of machine*

785      *Learning research* **12**, 2825–2830 (2011).

786   43.   Huang, L. *et al.* dbCAN-seq: a database of carbohydrate-active enzyme (CAZyme)

787      sequence and annotation. *Nucleic Acids Res.* **46**, D516–D521 (2018).

788   44.   Dong, X. *et al.* Genetic manipulation of the human gut bacterium Eggerthella lenta reveals

789      a widespread family of transcriptional regulators. *Nat. Commun.* **13**, 7624 (2022).

790   45.   Henke, M. T. *et al. Ruminococcus gnavus*, a member of the human gut microbiome

791      associated with Crohn's disease, produces an inflammatory polysaccharide. *Proceedings of*

792      *the National Academy of Sciences* **116**, 12672–12677 (2019).

793   46.   Liu, H. *et al.* Functional genetics of human gut commensal Bacteroides thetaiotaomicron

794      reveals metabolic requirements for growth across environments. *Cell Rep.* **34**, 108789

795      (2021).

796   47.   Liaw, R. *et al.* Tune: A Research Platform for Distributed Model Selection and Training.

797      *arXiv [cs.LG]* (2018).

798   48.   Paszke, A. *et al.* PyTorch: An Imperative Style, High-Performance Deep Learning Library.

799      *arXiv [cs.LG]* (2019).

800   49.   Mirdita, M. *et al.* ColabFold: making protein folding accessible to all. *Nat. Methods* **19**, 679–

801      682 (2022).

802    50. van Kempen, M. *et al.* Fast and accurate protein structure search with Foldseek. *Nat.*

803        *Biotechnol.* (2023) doi:10.1038/s41587-023-01773-0.

804    51. Zhang, C., Shine, M., Pyle, A. M. & Zhang, Y. US-align: universal structure alignments of

805        proteins, nucleic acids, and macromolecular complexes. *Nat. Methods* **19**, 1109–1115

806        (2022).

807    52. Teufel, F. *et al.* SignalP 6.0 predicts all five types of signal peptides using protein language

808        models. *Nat. Biotechnol.* **40**, 1023–1025 (2022).

809    53. Holm, L. Dali server: structural unification of protein families. *Nucleic Acids Res.* **50**, W210–

810        W215 (2022).

811    54. Franzosa, E. A. *et al.* Gut microbiome structure and metabolic activity in inflammatory

812        bowel disease. *Nat Microbiol* **4**, 293–305 (2019).

813    55. Seabold, S. & Perktold, J. Statsmodels: Econometric and statistical modeling with python.

814        in *Proceedings of the 9th Python in Science Conference* (SciPy, 2010).

815        doi:10.25080/majora-92bf1922-011.

816    56. Mallick, H. *et al.* Multivariable association discovery in population-scale meta-omics studies.

817        *PLoS Comput. Biol.* **17**, e1009442 (2021).

818    57. Nakamura, A. M., Nascimento, A. S. & Polikarpov, I. Structural diversity of carbohydrate

819        esterases. *Biotechnology Research and Innovation* **1**, 35–51 (2017).

820    58. Anderson, A. C., Stangherlin, S., Pimentel, K. N., Weadge, J. T. & Clarke, A. J. The SGNH

821        hydrolase family: a template for carbohydrate diversity. *Glycobiology* **32**, 826–848 (2022).

822    59. Prates, J. A. *et al.* The structure of the feruloyl esterase module of xylanase 10B from

823        Clostridium thermocellum provides insights into substrate recognition. *Structure* **9**, 1183–

824        1190 (2001).

825    60. Buzun, E. *et al.* A bacterial sialidase mediates early life colonization by a pioneering gut

826        commensal. *bioRxiv* 2023.08.08.552477 (2023) doi:10.1101/2023.08.08.552477.

827    61. Manichanh, C. *et al.* Reduced diversity of faecal microbiota in Crohn's disease revealed by

828     a metagenomic approach. *Gut* **55**, 205–211 (2006).

829  62. Gevers, D. *et al.* The treatment-naive microbiome in new-onset Crohn's disease. *Cell Host*

830     *Microbe* **15**, 382–392 (2014).

831  63. Lloyd-Price, J. *et al.* Multi-omics of the gut microbial ecosystem in inflammatory bowel

832     diseases. *Nature* **569**, 655–662 (2019).

833  64. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-

834     generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).

835  65. Sonnhammer, E. L. L., Eddy, S. R., Birney, E., Bateman, A. & Durbin, R. Pfam: Multiple

836     sequence alignments and HMM-profiles of protein domains. *Nucleic Acids Res.* **26**, 320–

837     322 (1998).

838  66. Kanehisa, M., Sato, Y., Kawashima, M., Furumichi, M. & Tanabe, M. KEGG as a reference

839     resource for gene and protein annotation. *Nucleic Acids Res.* **44**, D457-62 (2016).

840  67. Buchfink, B., Reuter, K. & Drost, H.-G. Sensitive protein alignments at tree-of-life scale

841     using DIAMOND. *Nat. Methods* **18**, 366–368 (2021).

842  68. *lightning: Deep learning framework to train, deploy, and ship AI products Lightning fast*.

843     (Github).

844  69. Li, L. *et al.* Massively Parallel Hyperparameter Tuning. (2018).

845  70. Abadi, M. *et al.* TensorFlow: Large-Scale Machine Learning on Heterogeneous Distributed

846     Systems. *arXiv [cs.DC]* (2016).
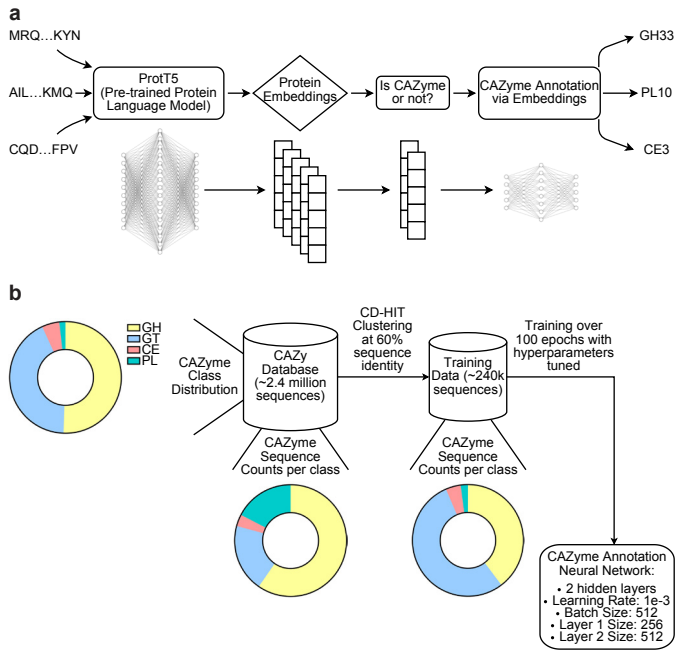
847  71. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.

848     *EMBnet.journal* **17**, 10–12 (2011).

849  72. Li, D., Liu, C.-M., Luo, R., Sadakane, K. & Lam, T.-W. MEGAHIT: an ultra-fast single-node

850     solution for large and complex metagenomics assembly via succinct de Bruijn graph.

851     *Bioinformatics* **31**, 1674–1676 (2015).
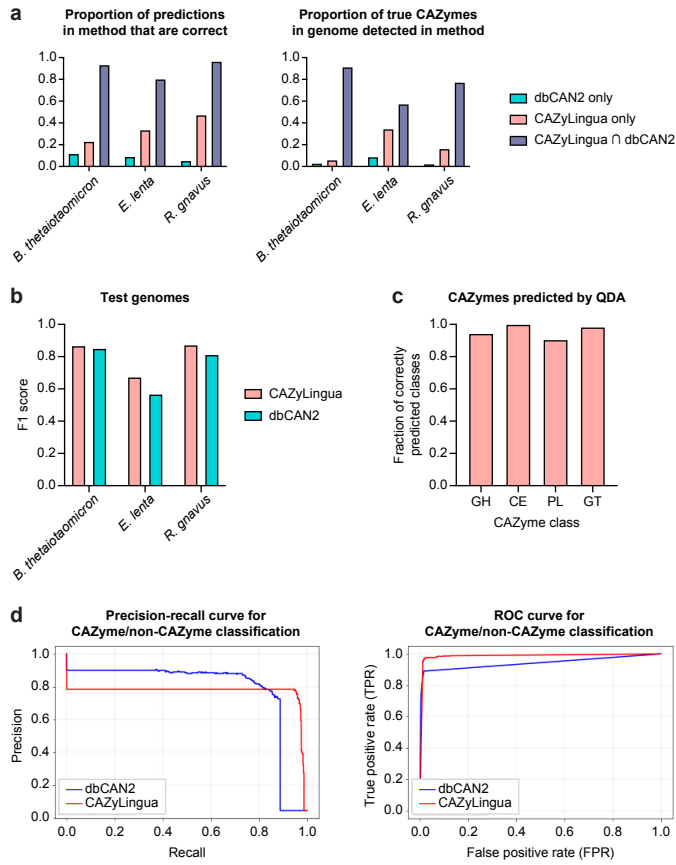
852  73. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site

853     identification. *BMC Bioinformatics* **11**, 119 (2010).

854    74. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the

855        analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).

856    75. Hunter. Matplotlib: A 2D Graphics Environment. **9**, 90–95 (2007).

857    76. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat.*

858        *Methods* **17**, 261–272 (2020).

**Figure 1**

**a**



**b**

**Figure 2**

**Figure 3**

**a**



Fold change of CAZyme genes

Infant — Mother

Legend: CAZyLingua vs eggNOG, CAZyLingua vs dbCAN2

**b**



CAZymes predicted by CAZyLingua

Infant — Mother

Legend: CE, GH, GT, PL

**c**



GH88
TM score: 0.61844
Sequence identity: 20.49%
Coverage: 29%

GH10
TM score: 0.90763
Sequence identity: 33.16%
Coverage: 95%

GH63
TM score: 0.53431
Sequence identity: 30.51%
Coverage: 12%

# Figure 4



**a** GH33 and GH43_18 embeddings (ProtT5)   GH33 and GH43_18 embeddings (CAZyLingua)

**b** Probability of GH33 after substituting amino acids in specific positions

**c**

TM score: 0.547
Sequence identity: 35.00%
Coverage: 31%

Sialidase (PDB: 1KIT) (GH33)
Unknown CAZyme (Predicted GH33)

TM score: 0.556
Sequence identity: 36.16%
Coverage: 16%

Neuraminidase (PDB: 3H6J) (GH33)
Unknown CAZyme (Predicted GH33)

Polysaccharide utilization locus

HTCS   hyp   SusC   SusD   GH33

BLAST metrics from unknown gene signal peptide

☐ GH33
▲ GH43_18

**Figure 5**



**a** Crohn's disease gene enrichment

**b** IgG4-related disease gene enrichment

**c** IgG4-related disease CE family predictions proportional to count

**d** Proportion of dbCAN2-predicted CAZymes detected by CAZyLingua as decision boundary is tuned

**Extended Data Figure 1**

## Extended Data Figure 2

Using the balanced accuracy score (micro averaging) over correct labels:
• Test accuracy: 99.6%
• dbCAN: 98.2%