

Shedding Light on Microbial Dark Matter with A Universal Language of Life

Hoarfrost A^{1,*}, Aptekmann, A², Farfañuk, G³, Bromberg Y^{2,*}

¹ Department of Marine and Coastal Sciences, Rutgers University, 71 Dudley Road, New Brunswick, NJ 08873, USA

² Department of Biochemistry and Microbiology, 76 Lipman Dr, Rutgers University, New Brunswick, NJ 08901, USA

³ Department of Biological Chemistry, Facultad de Ciencias Exactas y Naturales, Universidad de Buenos Aires, Argentina

* Corresponding author: adrienne.hoarfrost@rutgers.edu; yana@bromberglab.org

Running Title: LookingGlass: a universal language of life

Abstract

The majority of microbial genomes have yet to be cultured, and most proteins predicted from microbial genomes or sequenced from the environment cannot be functionally annotated. As a result, current computational approaches to describe microbial systems rely on incomplete reference databases that cannot adequately capture the full functional diversity of the microbial tree of life, limiting our ability to model high-level features of biological sequences. The scientific community needs a means to capture the functionally and evolutionarily relevant features underlying biology, independent of our incomplete reference databases. Such a model can form the basis for transfer learning tasks, enabling downstream applications in environmental microbiology, medicine, and bioengineering. Here we present LookingGlass, a deep learning model capturing a “universal language of life”. LookingGlass encodes contextually-aware, functionally and evolutionarily relevant representations of short DNA reads, distinguishing reads of disparate function, homology, and environmental origin. We demonstrate the ability of LookingGlass to be fine-tuned to perform a range of diverse tasks: to identify novel oxidoreductases, to predict enzyme optimal temperature, and to recognize the reading frames of DNA sequence fragments. LookingGlass is the first contextually-aware, general purpose pre-trained “biological language” representation model for short-read DNA sequences. LookingGlass enables functionally relevant representations of otherwise unknown and unannotated sequences, shedding light on the microbial dark matter that dominates life on Earth.

33 **Availability:** The pretrained LookingGlass model and the transfer learning-derived
34 models demonstrated in this paper are available in the LookingGlass release v1.0¹.
35 The open source *fastBio* Github repository and python package provides classes and
36 functions for training and fine tuning deep learning models with biological data². Code
37 for reproducing analyses presented in this paper are available as an open source
38 Github repository³.

39 **Key words:** deep learning, bioinformatics, metagenomics, language modeling,
40 transfer learning, microbial dark matter

41

42 **Introduction**

43 The microbial world is dominated by “microbial dark matter” – the majority of
44 microbial genomes remain to be sequenced^{4,5}, while the molecular functions of many
45 genes in microbial genomes are unknown⁶. In microbial communities (microbiomes),
46 the combination of these factors compounds this limitation. While the rate of biological
47 sequencing outpaces Moore’s law⁷, our traditional experimental means of annotating
48 these sequences cannot keep pace. Scientists thus typically rely on reference
49 databases which reflect only a tiny fraction of the biological diversity on Earth.

50 Our reliance on this incomplete annotation of biological sequences propagates
51 significant observational bias toward annotated genes and cultured genomes in
52 describing microbial systems. To break out of this cycle, the scientific community
53 needs a means of representing biological sequences that captures their functional and
54 evolutionary relevance and that is independent of our limited references.

55 Deep learning is particularly good at capturing complex, high dimensional
56 systems, and is a promising tool for biology⁸. However, deep learning generally
57 requires massive amounts of data to perform well. Meanwhile, collection and
58 experimental annotation of samples is typically time consuming and expensive, and
59 the creation of massive datasets for one study is rarely feasible. The scientific
60 community needs a means of building computational models which can capture

61 biological complexity while compensating for the low sample size and high
62 dimensionality that characterize biology.

63 Transfer learning provides a solution to the high-dimensionality, low-sample-size
64 conundrum. Transfer learning^{9,10} leverages domain knowledge learned by a model in
65 one training setting and applies it to a different but related problem. This approach is
66 effective because a model trained on a massive amount of data from a particular data
67 modality of interest (e.g. biological sequences) will learn features *general* to that
68 modality in addition to the *specific* features of its learning task. This general pretrained
69 model can then be further trained, or “fine-tuned”, to predict a downstream task of
70 interest more accurately, using less task-specific data, and in shorter training time than
71 would otherwise be possible. In computer vision, for example, by starting from a
72 pretrained model trained on many images, a model of interest doesn’t relearn general
73 image features such as a curve or a corner¹¹, but instead can devote its limited dataset
74 to refining the specific parameters of the target task. In natural language processing,
75 a generic language representation model¹² has been widely applied to diverse text
76 classification tasks, including biomedical text classification^{13,14}.

77 Pretrained models lower the barrier for widespread academic and private sector
78 applications, which typically have small amounts of data and limited computational
79 resources to model relatively complex data. Natural language processing for text, and
80 language modelling in particular, is analogous to biological sequences, in that
81 nucleotides are not independent or identically distributed¹⁵ and the nucleotide *context*
82 is important for defining the functional role and evolutionary history of the whole
83 sequence.

84 In genomics and metagenomics, there is no analogous contextually-aware
85 pretrained model that can be generally applied for transfer learning on read-length
86 biological sequences. Some previous studies have obtained important results using
87 transfer learning^{16,17}, but were either limited to relatively small training sets for
88 pretraining a model on a closely related prediction task¹⁶, or relied on gene counts from
89 the relatively well-annotated human genome to compile their training data¹⁷. Previous
90 works in learning continuous representations of biological sequences^{18,19} and

91 genomes²⁰ do not account for the order in which sequences or proteins appear and
92 are thus not contextually-aware. Recent advances in full-length protein sequence
93 representation learning^{21–24} show the potential of a self-supervised learning approach
94 that accounts for sequence context, but these rely on full length protein sequences (ca.
95 1,000 amino acids or 3,000 nucleotides). Full-length protein sequences are
96 computationally difficult (and sometimes impossible) to assemble from metagenomes,
97 which can produce hundreds of millions of short-read DNA sequences (ca. 60-300
98 nucleotides) per sample. To capture the full functional diversity of the microbial world,
99 we need a contextually-relevant means to represent the functional and evolutionary
100 features of biological sequences from microbial communities, in the short, fragmented
101 form in which they are sampled from their environment.

102 A biological ‘universal language of life’ should reflect functionally and
103 evolutionarily relevant features that underly biology as a whole and facilitate diverse
104 downstream transfer learning tasks. Here, we present LookingGlass, a biological
105 language model and sequence encoder, which produces contextually relevant
106 embeddings for any biological sequence across the microbial tree of life. LookingGlass
107 is trained and optimized for read-length sequences, such as those produced by the
108 most widely used sequencing technologies²⁵. For metagenomes in particular, a read-
109 level model avoids the need for assembly, which has a high computational burden and
110 potential for error. We also focus on Bacterial and Archaeal sequences, although we
111 include a discussion of the possibility for Eukaryotic and human-specific models below.

112 We demonstrate the functional and evolutionary relevance of the embeddings
113 produced by LookingGlass, and its broad utility across multiple transfer learning tasks
114 relevant to functional metagenomics. LookingGlass produces embeddings that
115 differentiate sequences with different molecular functions; identifies homologous
116 sequences, even at low sequence similarities where traditional bioinformatics
117 approaches fail; and differentiates sequences from disparate environmental contexts.
118 Using transfer learning, we demonstrate how LookingGlass can be used to illuminate
119 the “microbial dark matter” that dominates environmental settings by developing an
120 ‘oxidoreductase classifier’ that can identify novel oxidoreductases (enzymes

121 responsible for electron transfer, and the basis of all metabolism) with very low
122 sequence similarity to those seen during training. We also demonstrate LookingGlass'
123 ability to predict enzyme optimal temperatures from short-read DNA fragments; and to
124 recognize the reading frame (and thus "true" amino acid sequence) encoded in short-
125 read DNA sequences with high accuracy.

126 The transfer learning examples shown here, aside from providing useful models
127 in and of themselves, are intended to show the broad types of questions that can be
128 addressed with transfer learning from a single pretrained model. These downstream
129 models can illuminate the functional role of "microbial dark matter" by leveraging
130 domain knowledge of the functional and evolutionary features underlying microbial
131 diversity as a whole. More generally, LookingGlass is intended to serve as the scientific
132 community's 'universal language of life' that can be used as the starting point for
133 transfer learning in biological applications, and metagenomics in particular.

134

135

Methods

I. LookingGlass design and optimization

Dataset Generation.

138 The taxonomic organization of representative Bacterial and Archaeal genomes
139 was determined from the Genome Taxonomy Database, GTDB²⁶ (release 89.0). The
140 complete genome sequences were downloaded via the NCBI Genbank ftp²⁷. This
141 resulted in 24,706 genomes, comprising 23,458 Bacterial and 1,248 Archaeal
142 genomes.

143 Each genome was split into read-length chunks. To determine the distribution of
144 realistic read lengths produced by next-generation short read sequencing machines,
145 we obtained the BioSample IDs²⁷ for each genome, where they existed, and
146 downloaded their sequencing metadata from the MetaSeek²⁸ database using the
147 MetaSeek API. We excluded samples with average read lengths less than 60 or
148 greater than 300 base pairs. This procedure resulted in 7,909 BioSample IDs. The
149 average read lengths for these sequencing samples produced the 'read-length

150 distribution' (SI Fig 1) with a mean read length of 136bp. Each genome was split into
151 read-length chunks (with zero overlap in order to maximize information density and
152 reduce data redundancy in the dataset): a sequence length was randomly selected
153 with replacement from the read-length distribution and a sequence fragment of that
154 length was subset from the genome, with a 50% chance that the reverse complement
155 was used. The next sequence fragment was chosen from the genome starting at the
156 end point of the previous read-length chunk, using a new randomly selected read
157 length, and so on. To ensure that genomes in the training, validation, and test sets had
158 low sequence similarity, the sets were split along taxonomic branches such that
159 genomes from the *Actinomycetales*, *Rhodobacterales*, *Thermoplasmata*, and
160 *Bathyarchaeia* were partitioned into the validation set; genomes from the
161 *Bacteroidales*, *Rhizobiales*, *Methanosarcinales*, and *Nitrososphaerales* were
162 partitioned into the test set; and the remaining genomes remained in the training set.
163 This resulted in 529,578,444 sequences in the training set, 57,977,217 sequences in
164 the validation set, and 66,185,518 sequences in the test set. We term this set of reads
165 the *GTDB representative* set (Table 1).

166 The amount of data needed for training was also evaluated (SI Fig 2).
167 Progressively larger amounts of data were tested by selecting at random 1, 10, 100,
168 or 500 read-length chunks from each of the *GTDB representative* genomes in the
169 *GTDB representative training set*. Additionally, the performance of smaller but more
170 carefully selected datasets, representing the diversity of the microbial tree of life, were
171 tested by selecting for training one genome at random from each taxonomic class or
172 order in the *GTDB* taxonomy tree. In general, better accuracy was achieved in fewer
173 epochs with a greater amount of sequencing data (SI Fig 2); however, a much smaller
174 amount of data performed better if a representative genome was selected from each
175 *GTDB* taxonomy class.

176 The final LookingGlass model was trained on this class-level partition of the
177 microbial tree of life. We term this dataset the *GTDB class set* (Table 1). The training,
178 validation, and test sets were split such that no classes overlapped across sets: the
179 validation set included 8 genomes from each of the classes Actinobacteria,

180 Alphaproteobacteria, Thermoplasmata, and Bathyarchaeia (32 total genomes); the
181 test set included 8 genomes from each of the classes Bacteroidia, Clostridia,
182 Methanosarcinia, and Nitrososphaeria (32 total genomes); and the training set
183 included 1 genome from each of the remaining classes (32 archaeal genomes and 298
184 bacterial genomes for a total of 330 genomes). This resulted in a total of 6,641,723
185 read-length sequences in the training set, 949,511 in the validation set, and 632,388
186 in the test set (SI Table 1).

187 ***Architecture design and training.***

188 Recurrent Neural Networks (RNNs) are a type of neural network designed to take
189 advantage of the context dependence of sequential data (such as text, video, audio,
190 or biological sequences), by passing information from previous items in a sequence to
191 the current item in a sequence²⁹. Long Short Term Memory networks (LSTMs)³⁰ are
192 an extension of RNNs, which better learn long-term dependencies by handling the
193 RNN tendency to “forget” information farther away in a sequence³¹. LSTMs maintain a
194 “cell state” which contains the “memory” of the information in the previous items in the
195 sequence. LSTMs learn additional parameters which decide at each step in the
196 sequence which information in the “cell state” to “forget” or “update”.

197 LookingGlass uses a three-layer LSTM encoder model with 1,152 units in each
198 hidden layer and an embedding size of 104 based on the results of hyperparameter
199 tuning (see below). It divides the sequence into characters using a kmer size of 1 and
200 a stride of 1, i.e. is a character-level language model. LookingGlass is trained in a self-
201 supervised manner to predict a masked nucleotide, given the context of the preceding
202 nucleotides in the sequence. For each read in the training sequence, multiple training
203 inputs are considered, shifting the nucleotide that is masked along the length of the
204 sequence from the second position to the final position in the sequence. Because it is
205 a character-level model, a linear decoder predicts the next nucleotide in the sequence
206 from the possible vocabulary items ‘A’, ‘C’, ‘G’, and ‘T’, with special tokens for
207 ‘beginning of read’, ‘unknown nucleotide’ (for the case of ambiguous sequences), ‘end
208 of read’ (only ‘beginning of read’ was tokenized during LookingGlass training), and a
209 ‘padding’ token (used for classification only).

210 Regularization and optimization of LSTMs require special approaches to dropout
211 and gradient descent for best performance³². The *fastai* library³³ offers default
212 implementations of these approaches for natural language text, and so we adopt the
213 *fastai* library for all training presented in this paper. We provide the open-source *fastBio*
214 python package² which extends the *fastai* library for use with biological sequences.

215 LookingGlass was trained on a Pascal P100 GPU with 16GB memory on
216 Microsoft Azure, using a batch size of 512, a back propagation through time (bptt)
217 window of 100 base pairs, the Adam optimizer³⁴, and utilizing a Cross Entropy loss
218 function (SI Table 2). Dropout was applied at variable rates across the model (SI Table
219 2). LookingGlass was trained for a total of 12 days for 75 epochs, with progressively
220 decreasing learning rates based on the results of hyperparameter optimization (see
221 below): for 15 epochs at a learning rate of 1e-2, for 15 epochs at a learning rate of 2e-
222 3, and for 45 epochs at a learning rate of 1e-3.

223 ***Hyperparameter optimization.***

224 Hyperparameters used for the final training of LookingGlass were tuned using a
225 randomized search of hyperparameter settings. The tuned hyperparameters included
226 kmer size, stride, number of LSTM layers, number of hidden nodes per layer, dropout
227 rate, weight decay, momentum, embedding size, bptt size, learning rate, and batch
228 size. An abbreviated dataset consisting of ten randomly selected read-length chunks
229 from the *GTDB representative set* was created for testing many parameter settings
230 rapidly. A language model was trained for two epochs for each randomly selected
231 hyperparameter combination, and those conditions with the maximum performance
232 were accepted. The hyperparameter combinations tested and the selected settings are
233 described in the associated Github repository³.

234

235 **II. LookingGlass validation and analysis of embeddings**

236 **Functional relevance**

237 ***Dataset generation.***

238 In order to assess the ability of the LookingGlass embeddings to inform the
239 molecular function of sequences, metagenomic sequences from a diverse set of
240 environments were downloaded from the Sequence Read Archive (SRA)³⁵. We used
241 MetaSeek²⁸ to choose ten metagenomes at random from each of the ‘environmental
242 packages’ defined by the MIxS metadata standards³⁶: ‘built environment’, ‘host-
243 associated’, ‘human-gut’, ‘microbial mat/biofilm’, ‘miscellaneous’, ‘plant-associated’,
244 ‘sediment’, ‘soil’, ‘wastewater/sludge’, and ‘water’, for a total of 100 metagenomes. The
245 SRA IDs used are available in (SI Table 3). The raw DNA reads for these 100
246 metagenomes were downloaded from the SRA with the NCBI e-utilities. These 100
247 metagenomes were annotated with the *mi-faser* tool³⁷ with the --read-map option to
248 generate predicted functional annotation labels (to the fourth digit of the Enzyme
249 Commission (EC) number), out of 1,247 possible EC labels, for each annotatable read
250 in each metagenome. These reads were then split 80%/20% into ‘training’/‘validation
251 candidate’ sets of reads. To ensure that there was minimal overlap in sequence
252 similarity between the training and validation set, we compared the ‘validation
253 candidate’ sets of each EC annotation to the training set for that EC number with CD-
254 HIT³⁸, and filtered out any reads with >80% DNA sequence similarity to the reads of
255 that EC number in the training set (the minimum CD-HIT DNA sequence similarity
256 cutoff). In order to balance EC classes in the training set, overrepresented ECs in the
257 training set were downsampled to the mean count of read annotations (52,353 reads)
258 before filtering with CD-HIT. After CD-HIT processing, any underrepresented EC
259 numbers in the training set were oversampled to the mean count of read annotations
260 (52,353 reads). The validation set was left unbalanced to retain a distribution more
261 realistic to environmental settings. The final training set contained 61,378,672 reads,
262 while the validation set contained 2,706,869 reads. We term this set of reads and their
263 annotations the *mi-faser functional set* (Table 1).

264 As an external test set, we used a smaller number of DNA sequences from genes
265 with experimentally validated molecular functions. We linked the manually curated
266 entries of Bacterial or Archaeal proteins from the Swiss-Prot database³⁹ corresponding
267 to the 1,247 EC labels in the *mi-faser functional set* with their corresponding genes in

268 the EMBL database⁴⁰. We downloaded the DNA sequences, and selected ten read-
269 length chunks at random per coding sequence. This resulted in 1,414,342 read-length
270 sequences in the test set. We term this set of reads and their annotations the *Swiss-
271 Prot functional set* (Table 1).

272 ***Fine-tuning procedure.*** We fine-tuned the LookingGlass language model to
273 predict the functional annotation of DNA reads, to demonstrate the speed with which
274 an accurate model can be trained using our pretrained LookingGlass language model.
275 The architecture of the model retained the 3-layer LSTM encoder and the weights of
276 the LookingGlass language model encoder, but replaced the language model decoder
277 with a new multi-class classification layer with pooling (with randomly initialized
278 weights). This pooling classification layer is a sequential model consisting of the
279 following layers: a layer concatenating the output of the LookingGlass encoder with
280 min, max, and average pooling of the outputs (for a total dimension of $104 \times 3 = 312$), a
281 batch normalization⁴¹ layer with dropout, a linear layer taking the 312-dimensional
282 output of the batch norm layer and producing a 50-dimensional output, another batch
283 normalization layer with dropout, and finally a linear classification layer that outputs the
284 predicted functional annotation of a read as a probability distribution of the 1,247
285 possible mi-faser EC annotation labels. We then trained the functional classifier on the
286 *mi-faser functional set* described above. Because the >61 million reads in the training
287 set were too many to fit into memory, training was done in 13 chunks of ~5-million
288 reads each until one total epoch was completed. Hyperparameter settings for the
289 functional classifier training are seen in SI Table 2.

290 ***Encoder embeddings and MANOVA test.*** To test whether the LookingGlass
291 language model embeddings (before fine-tuning, above) are distinct across functional
292 annotations, a random subset of ten reads per functional annotation was selected from
293 each of the 100 SRA metagenomes (or the maximum number of reads present in that
294 metagenome for that annotation, whichever was greater). This also ensured that reads
295 were evenly distributed across environments. The corresponding fixed-length
296 embedding vectors for each read was produced by saving the output from the
297 LookingGlass encoder (before the embedding vector is passed to the language model

298 decoder) for the final nucleotide in the sequence. This vector represents a contextually
299 relevant embedding for the overall sequence. The statistical significance of the
300 difference between embedding vectors across all functional annotation groups was
301 tested with a MANOVA test using the R stats package⁴².

302

303 **Evolutionary relevance**

304 ***Dataset generation.***

305 The OrthoDB database⁴³ provides orthologous groups (OGs) of proteins at
306 various levels of taxonomic distance. For instance, the OrthoDB group '77at2284'
307 corresponds to proteins belonging to 'Glucan 1,3-alpha-glucosidase at the Sulfolobus
308 level', where '2284' is the NCBI taxonomy ID for the genus *Sulfolobus*.

309 We tested whether embedding similarity of homologous sequences (sequences
310 within the same OG) is higher than that of nonhomologous sequences (sequences
311 from different OGs). We tested this in OGs at multiple levels of taxonomic distance –
312 genus, family, order, class, and phylum. At each taxonomic level, ten individual taxa at
313 that level were chosen from across the prokaryotic tree of life (e.g. for the genus level,
314 *Acinetobacter*, *Enterococcus*, *Methanosarcina*, *Pseudomonas*, *Sulfolobus*, *Bacillus*,
315 *Lactobacillus*, *Mycobacterium*, *Streptomyces*, and *Thermococcus* were chosen). For
316 each taxon, 1,000 randomly selected OGs corresponding to that taxon were chosen;
317 for each of these OGs, five randomly chosen genes within this OG were chosen.

318 OrthoDB cross-references OGs to UniProt³⁹ IDs of the corresponding proteins.
319 We mapped these to the corresponding EMBL coding sequence (CDS) IDs⁴⁰ via the
320 UniProt database API³⁹; DNA sequences of these EMBL CDSs were downloaded via
321 the EMBL database API. For each of these sequences, we generated LookingGlass
322 embedding vectors.

323 ***Homologous and nonhomologous sequence pairs.***

324 To create a balanced dataset of homologous and nonhomologous sequence
325 pairs, we compared all homologous pairs of the five sequences in an OG (total of ten

326 homologous pairs) to an equal number of randomly-selected out-of-OG comparisons
327 for the same sequences; i.e., each of the five OG sequences was compared to 2 other
328 randomly-selected sequences from any other randomly-selected OG (total of ten
329 nonhomologous pairs). We term this set of sequences, and their corresponding
330 LookingGlass embeddings, the *OG homolog set* (Table 1).

331 ***Embedding and sequence similarity.*** For each sequence pair, the sequence
332 and embedding similarity were determined. The embedding similarity was calculated
333 as the cosine similarity between embedding vectors. The sequence similarity was
334 calculated as the Smith-Waterman alignment score using the BioPython⁴⁴ pairwise2
335 package, with a gap open penalty of -10 and a gap extension penalty of -1. The IDs of
336 chosen OGs, the cosine similarities of the embedding vectors, and sequence
337 similarities of the DNA sequences are available in the associated Github repository³.

338

339 **Environmental Relevance**

340 ***Encoder embeddings and MANOVA test .***

341 The LookingGlass embeddings and the environment of origin for each read in the
342 *mi-faser functional set* were used to test the significance of the difference between the
343 embedding vectors across environmental contexts. The statistical significance of this
344 difference was evaluated with a MANOVA test using the R stats package⁴².

345

346 **III. Oxidoreductase classifier**

347 ***Dataset generation.***

348 The manually curated, reviewed entries of the Swiss-Prot database³⁹ were
349 downloaded (June 2, 2020). Of these, 23,653 entries were oxidoreductases (EC
350 number 1.-.-) of Archaeal or Bacterial origin (988 unique ECs). We mapped their
351 UniProt IDs to both their EMBL CDS IDs and their UniRef50 IDs via the UniProt
352 database mapper API. Uniref50 IDs identify clusters of sequences with >50% amino
353 acid identity. This cross-reference identified 28,149 EMBL CDS IDs corresponding to

354 prokaryotic oxidoreductases, belonging to 5,451 unique UniRef50 clusters. We split
355 this data into training, validation, and test sets such that each UniRef50 cluster was
356 contained in only one of the sets, i.e. there was no overlap in EMBL CDS IDs
357 corresponding to the same UniRef50 cluster across sets. This ensures that the
358 oxidoreductase sequences in the validation and test sets are dissimilar to those seen
359 during training. The DNA sequences for each EMBL CDS ID were downloaded via the
360 EMBL database API. This data generation process was repeated for a random
361 selection of non-oxidoreductase UniRef50 clusters, which resulted in 28,149 non-
362 oxidoreductase EMBL CDS IDs from 13,248 unique UniRef50 clusters.

363 ~50 read-length chunks (selected from the representative read-length
364 distribution, as above) were selected from each EMBL CDS DNA sequence, with
365 randomly selected start positions on the gene and a 50% chance of selecting the
366 reverse complement, such that an even number of read-length sequences with
367 'oxidoreductase' and 'non-oxidoreductase' labels were generated for the final dataset.
368 This procedure produced a balanced dataset with 2,372,200 read-length sequences in
369 the training set, 279,200 sequences in the validation set, and 141,801 sequences in
370 the test set. We term this set of reads and their annotations the *oxidoreductase model*
371 *set* (Table 1).

372 ***Fine-tuning procedure.***

373 Since our functional annotation classifier addresses a closer classification task to
374 the oxidoreductase classifier than LookingGlass itself, the architecture of the
375 oxidoreductase classifier was fine-tuned starting from the functional annotation
376 classifier, replacing the decoder with a new pooling classification layer (as described
377 above for the functional annotation classifier) and with a final output size of 2 to predict
378 'oxidoreductase' or 'not oxidoreductase'. Fine tuning of the oxidoreductase classifier
379 layers was done successively, training later layers in isolation and then progressively
380 including earlier layers into training, using discriminative learning rates ranging from
381 $1e-2$ to $5e-4$, as previously described⁴⁵.

382

383 ***Model performance in metagenomes.***

384 16 marine metagenomes from the surface (SRF, ~5 meters) and mesopelagic
385 (MES, 175-800 meters) from eight stations sampled as part of the TARA expedition⁴⁶
386 were downloaded from the SRA³⁵ (SI Table 4, SRA accession numbers ERR598981,
387 ERR599063, ERR599115, ERR599052, ERR599020, ERR599039, ERR599076,
388 ERR598989, ERR599048, ERR599105, ERR598964, ERR598963, ERR599125,
389 ERR599176, ERR3589593, and ERR3589586). Metagenomes were chosen from a
390 latitudinal gradient spanning polar, temperate, and tropical regions and ranging from -
391 62 to 76 degrees latitude. Mesopelagic depths from four out of the eight stations were
392 sampled from oxygen minimum zones (OMZs, where oxygen <20 $\mu\text{mol/kg}$). Each
393 metagenome was rarefied to twenty million randomly selected sequences. We term
394 this set of reads the *oxidoreductase metagenome set* (Table 1, SI Table 4). Predictions
395 of “oxidoreductase” or “not oxidoreductase” were made for these sequences with the
396 oxidoreductase classifier. To compare model predictions to alternative functional
397 annotation methods, reads in the *oxidoreductase metagenome set* were annotated
398 with mi-faser³⁷ with the --read-map option, and with the MG-RAST functional
399 annotation pipeline⁴⁷ using default settings.

400

401 **IV. Reading Frame classifier**

402 ***Dataset generation.***

403 For each taxonomic order, the coding sequence (CDS) files of one of the genome
404 IDs in the *GTDB representative set* were downloaded from NCBI²⁷. These were split
405 into read-length chunks as described above. Note that because each sequence is a
406 coding sequence, the true frame of translation for each read-length chunk was known;
407 this translation frame label of (1, 2, 3, -1, -2, or -3) was recorded for each read-length
408 input³. We term this set of reads the *reading frame set* (Table 1).

409 ***Fine-tuning procedure.***

410 The translation frame classifier was adjusted with a pooling classification layer
411 with an output size of six for the six possible translation frame labels. Fine tuning was

412 performed over successive layers with discriminative learning rates ranging from 1e-3
413 to 5e-5 as described for the oxidoreductase classifier.

414

415 **V. Optimal temperature classifier**

416 ***Dataset generation.***

417 The optimal growth temperature for 19,474 microorganisms was manually
418 curated from multiple sources: BacDive⁴⁸, DSMZ⁴⁹, Pasteur Institute (PI), the National
419 Institute for Environmental Studies (NIES)⁵⁰, and a curated list from a previous work⁵¹.
420 BacDive data is available through their API, which contains calls to retrieve the species
421 list and to get all data about a specific species. For DSMZ, PI, and NIES databases we
422 used previously published⁵² data files (for DSMZ and PI) or scripts and method (NIES)
423 to query optimal growth temperature information (accessed July 2020). We finally
424 cross-referenced optimal growth temperature of these organisms to their NCBI
425 taxonomy ID⁵³.

426 Previous studies have shown a strong correlation between enzyme optimal
427 temperature and organism optimal growth temperature⁵². We assumed that core
428 housekeeping enzymes, such as those involved in transcription and translation, would
429 have the same optimal functional temperature as the organism itself. Thus, we cross-
430 referenced the 19,474 microorganisms identified above to the UniProt IDs belonging
431 to those taxa for the housekeeping genes: RNA polymerase (EC 2.7.7.6), RNA
432 helicase (EC 3.6.4.13), DNA polymerase (EC 2.7.7.7), DNA primase (EC 2.7.7.101 for
433 Bacteria, EC 2.7.7.102 for Archaea), DNA helicase (EC 3.6.4.12), DNA ligase (ECs
434 6.5.1.1, 6.5.1.2, 6.5.1.6, and 6.5.1.7), and topoisomerase (ECs 5.6.2.1 and 5.6.2.2).
435 Finally, we linked these UniProt IDs to the corresponding EMBL CDS IDs, downloaded
436 the gene sequences, and split them into read-length chunks as described above.

437 The optimal temperature label for each read was derived from the optimal growth
438 temperature from its source organism; range [4-104.5] C°. The optimal temperature
439 labels were converted to categorical labels of 'psychrophilic' for optimal temperatures
440 <15 C°, 'mesophilic' for [20-40] C°, and 'thermophilic' for >50 C°. The training,

441 validation, and test sets were split by EC number such that only sequences from EC
442 3.6.4.13 were in the validation set, only sequences from EC 6.5.1.2 were in the test
443 set, and all other EC numbers were in the training set. Finally, the inputs from each
444 label category were either downsampled or upsampled (as described above for the *mi-*
445 *faser functional set*) to a balanced number of inputs for each class. This resulted in
446 5,971,152 inputs in the training set with ~2,000,000 reads per label; 597,136 inputs in
447 the validation set with ~200,000 reads per label; and 296,346 inputs to the test set with
448 ~100,000 reads per label. We term this set of reads and their annotations the *optimal*
449 *temp set* (Table 1).

450 ***Fine-tuning procedure.***

451 The optimal temperature classifier was adjusted with a pooling classification layer
452 with an output size of three for the three possible optimal temperature labels, as
453 described above. Fine tuning was performed over successive layers with
454 discriminative learning rates ranging from 5e-2 to 5e-4 as described for the
455 oxidoreductase classifier.

456

457 **VI. Metrics**

458 Model performance metrics for accuracy (all classifiers), precision, recall, and F1
459 score (binary classifiers only) are defined as below:

$$460 \text{ **Accuracy:} \quad \frac{TP+TN}{TP+FP+TN+FN} \quad (1)**$$

461

$$462 \text{ **Precision:} \quad \frac{TP}{TP + FP} \quad (2)**$$

463

$$464 \text{ **Recall:} \quad \frac{TP}{TP + FN} \quad (3)**$$

465

466 **F1 score:**
$$2 \cdot \frac{\textit{Precision} \cdot \textit{Recall}}{\textit{Precision} + \textit{Recall}} \quad (4)$$

467

468 where TP is a true positive (correct positive label prediction), FP is a false positive
469 (incorrect prediction of the positive label), TN is a true negative (correct negative label
470 prediction), and FN is a false negative (incorrect prediction of the negative label).

471

472 **VII. Software, model deployment and reproducibility**

473 The LookingGlass pretrained model, as well as the pretrained functional
474 classifier, oxidoreductase classifier, optimal temperature classifier, and reading frame
475 classifier models, are provided in the *LookingGlass* release v1.0¹. We also provide the
476 *fastBio* python package that extends the fastai³³ library for custom data loading and
477 processing functions designed for use with biological sequences². The scripts used for
478 data gathering, training of the LookingGlass model, training of models using transfer
479 learning, and analysis of the results presented in this paper are available in the
480 associated Github repository³.

481

482

Results

483 **I. LookingGlass – a “universal language of life”**

484 The LookingGlass model was trained as a 3-layer LSTM encoder chained to a
485 decoder predicting the next (masked) nucleotide in a DNA sequence fragment, on a
486 set of more than 6.6 million read-length sequences selected from microbial genomes
487 spanning each taxonomic class in the microbial tree of life (Methods).

488

489 **LookingGlass captures functionally relevant features of sequences.**

490 The LookingGlass encoder produces a fixed-length vector embedding of each
491 sequence input. In the *mi-faser functional* validation set containing metagenomic reads
492 with functional annotation labels (Methods), these sequence embeddings were distinct

493 across functional annotations (MANOVA $P < 10^{-16}$) without any additional fine-tuning.
494 Moreover, a model was fine-tuned on the *mi-faser functional set* to predict mi-faser
495 functional annotations to the 4th EC number and achieved 81.5% accuracy (Eqn 1) on
496 the validation set in only one epoch. At coarser resolution accuracy was improved: to
497 83.8% at the 3rd EC number (SI Fig 3); 84.4% at the 2nd EC number (Fig 1b); and
498 87.1% at the 1st EC number (Fig 1a). In testing on an experimentally-validated set of
499 functional annotations (*Swiss-Prot functional set*; Methods), this classifier had a lower
500 accuracy (50.8%) that was still substantially better than random (0.08%). Thus,
501 LookingGlass captures functionally relevant features of biological sequences, (1)
502 distinguishing between functional classes without being expressly trained to do so and
503 (2) enabling rapid convergence on an explicit high-dimensional functional classification
504 task at the read level.

505

506 **LookingGlass captures evolutionarily-relevant features of sequences.**

507 The embedding similarity of homologous sequence pairs in the *OG homolog set*
508 was significantly higher (unpaired t-test $P < 10^{-16}$) than that of nonhomologous pairs,
509 with no additional fine-tuning, for fine to broad levels of phylogenetic distances, i.e.
510 genus, family, order, class, and phylum (Fig 2a). LookingGlass embeddings
511 differentiate homology with ~66-79% accuracy which varied by taxonomic level (SI Fig
512 4, SI Table 5). This variation is due to variable sequence similarity across taxa, i.e.
513 sequences from species-level homologs have higher sequence similarity than
514 homologs at the phylum level. Our model attained 66.4% accuracy at the phylum level
515 (Fig 2b), 68.3% at the class level, 73.2% at the order level, 76.6% at the family level,
516 and 78.9% at the genus level. This performance is a substantial improvement over
517 random (50% accuracy), and was obtained from LookingGlass embeddings alone
518 which were not expressly trained on this task.

519 LookingGlass embeddings differentiate between homologous and
520 nonhomologous sequences independent of their sequence similarity (Smith-Waterman
521 alignments, Methods). This is particularly useful since many (e.g. 44% at the phylum

522 level, SI Table 5) homologs have very low sequence similarity (alignment score < 50;
523 Fig 2c, SI Table 5). For these, LookingGlass embedding similarity is still high, indicating
524 that our model captures evolutionary relationships between sequences, even where
525 traditional algorithmic approaches do not. In fact, embedding similarity between
526 sequences is poorly correlated with the sequence similarity alignment score (Pearson
527 $R^2=0.28-0.44$). The high accuracy with which LookingGlass identifies homologs
528 indicates that it captures high-level features reflecting evolutionary relationships
529 between sequences.

530

531 **LookingGlass differentiates sequences from disparate environmental contexts.**

532 The sequences in the *mi-faser functional set* have distinct embedding fingerprints
533 across different environments – embedding similarity between environments is
534 generally lower than embedding similarity within an environment (Fig 3, MANOVA
535 $P<10^{-16}$), even though the LookingGlass embeddings were not explicitly trained to
536 recognize environmental labels. While there is some overlap of embeddings across
537 environmental contexts, those with the most overlap are between similar environments
538 – for example, the colocalization of ‘wastewater/sludge’ with ‘human-gut’ and ‘built
539 environment’ (Fig. 3b).

540

541 **II. LookingGlass enables diverse downstream transfer learning tasks**

542

543 **Mining environmental settings for functional descriptions of “microbial dark 544 matter”.**

545 ***Using LookingGlass and transfer learning to identify novel functional groups.***

546 By using LookingGlass as a starting point, we can converge more quickly and
547 with less data on a more accurate model for assigning molecular functions at the read
548 level. Additionally, downstream models addressing similar tasks can in turn be used
549 as pretrained models for further fine-tuning. To demonstrate this, we fine-tuned the

550 LookingGlass functional classifier (described above) to predict whether a read-length
551 DNA sequence likely comes from an oxidoreductase-encoding gene (EC number 1.-.-
552 .-). Our fine-tuned model was able to correctly classify previously unseen (<50% amino
553 acid sequence-identical) oxidoreductases with 82.3% accuracy at the default
554 prediction threshold of 0.5 (Fig 4). Oxidoreductases are a deeply branched, highly
555 diverse class of enzymes, such that sequence similarity within a single functional
556 annotation (EC number) is often very low; the DNA sequence identity of
557 oxidoreductase gene sequences within a single EC number in the *oxidoreductase*
558 *model* validation set was a median of 59%, and was as low as 17%. As such,
559 oxidoreductases can be difficult to identify via sequence similarity-based homology
560 searches in environmental samples (e.g. box in Fig 2c). The oxidoreductase classifier,
561 in contrast, achieves high model performance even in such cases where sequence
562 similarity within EC annotations is low. Notably, the average model performance for a
563 given EC number was independent of the sequence similarity of genes within that EC
564 ($R^2=0.004$, SI Fig 5).

565 ***Mining novel, unannotated oxidoreductases from metagenomes along a***
566 ***latitudinal and depth gradient in the global ocean.***

567 The majority of sequencing reads from environmental metagenomes are routinely
568 unable to be functionally annotated⁵⁴. To demonstrate the advantage of the
569 oxidoreductase classifier over traditional homology-based approaches, we evaluated
570 our model on twenty million randomly-selected reads from each of 16 marine
571 metagenomes in the *oxidoreductase metagenome* set spanning broad ranges in
572 latitude (from -62 to 76 degrees), depth (from the surface, ~5 meters, to mesopelagic,
573 ~200-1,000 meters), and oxygen concentrations (including four mesopelagic samples
574 from oxygen minimum zones).

575 The percentage of reads predicted to be oxidoreductases ranged from 16.4% -
576 20.6%, and followed trends with depth and latitude (Fig 5). The relative abundance of
577 oxidoreductases was significantly higher in mesopelagic depths than in surface waters
578 (Fig 5a, ANOVA $P=0.02$), with marginally higher (albeit not statistically significant)
579 proportions of oxidoreductases in the oxygen minimum zones relative to oxygen-

580 replete mesopelagic samples ($P=0.13$). There was also a significant trend in the
581 relative abundance of oxidoreductases along latitudinal gradients in surface waters
582 (Fig 5b, $R^2=0.79$, $P=0.04$), with higher proportions of oxidoreductases in higher
583 latitudes. This latitudinal trend was reflected in a similar, but weaker, temperature-
584 driven trend ($R^2= -0.66$, $P=0.11$, SI Fig 6).

585 Two alternative functional annotation tools, mi-faser³⁷ and MG-RAST⁴⁷, were only
586 able to annotate a much smaller proportion of sequences in these metagenomes (Fig
587 5c, SI Table 6), with even smaller proportions of oxidoreductases identified. MG-RAST
588 annotated 26.7-50.3% of the reads across metagenomes, with 0.01-4.0% of reads
589 identified as oxidoreductases. Mi-faser annotated 0.17-2.9% of the reads, of which
590 0.04-0.59% were oxidoreductases. In both cases, a majority of reads remained
591 unannotated, a condition typical of homology-based functional annotation
592 approaches⁵⁴. As a result, a large proportion of enzymes in the environment are
593 unlikely to be recovered using these approaches, which may also skew the observed
594 trends across samples. Notably, the depth and latitudinal trends identified with the
595 oxidoreductase classifier were not reported by either MG-RAST or mi-faser (SI Fig 7).
596 There was no significant difference in the proportion of oxidoreductases predicted in
597 the surface vs. mesopelagic waters for either MG-RAST ($P=0.73$) or mi-faser ($P=0.60$)
598 and no significant correlation with latitude in surface waters for either mi-faser
599 ($R^2=0.58$, $P=0.17$) or MG-RAST ($R^2= -0.49$, $P=0.27$); note that MG-RAST in fact
600 observed an anticorrelation trend for the latter (although still insignificant). This
601 highlights the potential importance of unannotatable reads in driving functional patterns
602 in the environment, which can be captured by the approach and models described
603 here and would otherwise be missed using traditional approaches.

604

605 **Reference-free translation of read-length DNA sequences to peptides.**

606 While the amino acid sequence encoded in short DNA reads is difficult to infer
607 directly using traditional bioinformatic approaches, it is also a product of the non-
608 random organization of DNA sequences. We fine-tuned the LookingGlass encoder to

609 predict the translation frame start position (1, 2, 3, -1, -2, or -3) directly from read-length
610 DNA coding sequences. This reading frame classifier attained 97.8% accuracy, a
611 major improvement over random (16.7% accuracy). Note this classifier was trained
612 only on coding sequences and is currently intended only for prokaryotic sources with
613 low amounts of noncoding DNA⁵⁵.

614

615 **Prediction of enzyme optimal temperature from DNA sequence fragments**

616 The optimal temperature of an enzyme is in part dependent on DNA sequence
617 features^{56,57}, but is difficult to predict, particularly from short reads. We fine-tuned
618 LookingGlass to predict whether a read-length DNA sequence originates from an
619 enzyme with an optimal temperature that is psychrophilic (<15 C°), mesophilic (20-40
620 C°), or thermophilic (>50 C°). The optimal temperature classifier was able to predict
621 the optimal temperature category correctly with 70.1% accuracy (random accuracy
622 =33.3%).

623

624

624 **Discussion**

625 Microbes perform a vast diversity of functional roles in natural environments as
626 well as in industrial and biomedical settings. They play a central role in regulating
627 Earth's biogeochemical cycles⁵⁸, and have a tremendous impact on the health of their
628 human hosts⁵⁹, but the complex functional networks that drive their activities are poorly
629 understood. Microbial genomes record a partial history of the evolution of life on
630 Earth⁶⁰, but much of this information is inadequately captured by homology-based
631 inference. Microbial communities are a subject of great interest for developing natural⁶¹
632 and synthetic⁶² products for bioengineering applications, but our ability to describe,
633 model, and manipulate the systems-level functions of these microbiomes is limited.

634 The LookingGlass 'universal language of life' creates representations of DNA
635 sequences that capture their functional and evolutionary relevance, independent of
636 whether the sequence is contained in reference databases. The vast majority of
637 microbial diversity is uncultured and unannotated⁴⁻⁶. LookingGlass opens the door to

638 harnessing the potential of this “microbial dark matter” to improve our understanding
639 of, and ability to manipulate, microbial systems. It is a broadly useful, ‘universal’ model
640 for downstream transfer learning tasks, enabling a wide diversity of functional
641 predictions relevant to environmental metagenomics, bioengineering, and biomedical
642 applications.

643 We demonstrate here the ability of LookingGlass to be fine-tuned to identify novel
644 oxidoreductases, even those with low sequence similarity to currently known
645 oxidoreductases. Applying the oxidoreductase classifier to 16 marine metagenomes
646 identified patterns in the relative abundance of oxidoreductases that follow global
647 gradients in latitude and depth. These observations are in line with previous studies
648 that have identified greater overall functional and taxonomic richness^{46,63}, as well as a
649 greater diversity of oxidoreductases specifically⁶⁴, in deep marine waters relative to
650 shallow depths. Studies conflict, however, about whether taxonomic and functional
651 diversity increases^{63,65–67} or decreases^{68–70} with absolute latitude. Notably, neither the
652 latitudinal nor depth trends in oxidoreductase relative abundance observed by the
653 oxidoreductase classifier were captured by traditional homology-based functional
654 annotation tools. The inconsistent results produced by traditional annotation tools in
655 this study and others further demonstrates the importance of unannotated functional
656 diversity for cross-sample comparisons, and the potential of the approach described in
657 this study.

658 There may be multiple ecological mechanisms driving the observed latitudinal
659 and depth patterns in oxidoreductase relative abundance; for example, the
660 streamlining of genomes⁷¹ that preserves oxidoreductases relative to less essential
661 genes under resource limitation or temperature stress, or a reflection of a higher
662 abundance of anaerobic respiration genes in mesopelagic waters relative to surface
663 waters⁷². Future efforts to capture and compare the full functional diversity of
664 environmental settings using the approaches described here can further illuminate and
665 differentiate between these mechanisms.

666 The reads predicted to be from novel oxidoreductases are candidates for further
667 functional characterization, and for targeted assembly of novel oxidoreductase genes.

668 Shining light on these “dark matter” oxidoreductases can enable more complete
669 comparisons of oxidoreductase composition and diversity across environmental
670 gradients. Future efforts to fine tune LookingGlass for additional functional targets can
671 expand the classes of enzymes identified and create a fuller picture of microbial
672 functional diversity in environmental settings. By definition, poorly-studied
673 environments contain the greatest amount of unknown functional diversity, and a tool
674 such as LookingGlass provides a novel and important way to evaluate this functional
675 diversity.

676 LookingGlass was also fine-tuned to correctly identify the reading frame, and thus
677 the amino acid translation, of short-read DNA coding sequences. Translated amino
678 acid sequences are used for a variety of bioinformatics applications, most notably for
679 molecular function annotation. There are two categories of function annotation tools –
680 those that annotate from short sequencing reads directly^{37,47,73,74} and those that
681 annotate from assembled genes/contigs^{47,75}. In both cases, DNA reads must first be
682 converted to amino acid sequences. For short-read annotation tools, six-frame
683 translation of each DNA sequence produces all six possible amino acid sequences for
684 alignment to reference databases, which increases the computational burden of
685 alignment six-fold. For tools that annotate from assemblies, datasets are first
686 assembled and open reading frames (ORFs) predicted before amino acid sequences
687 can be inferred. This procedure is computationally intensive, error-prone, and throws
688 away reads that can’t be assembled or for which coding regions can’t be identified,
689 particularly for members of the rare biosphere or in highly diverse environments. Direct
690 translation from DNA reads thus could enable much more efficient computation for any
691 bioinformatics application that uses read-derived amino acid sequences. Note that the
692 reading frame classifier described here focuses on prokaryotic genomes, which
693 generally have only ~12-14% noncoding DNA⁵⁵. For eukaryotes, a classifier will need
694 to be created to distinguish between coding and noncoding DNA and predict reading
695 frames for only the coding sequences.

696 Finally, we demonstrated the ability of LookingGlass to be fine tuned to predict
697 optimal enzyme temperatures from DNA sequences. Importantly, this was possible

698 from short reads alone, although a classifier trained on assembled genes would likely
699 yield even better results. This result demonstrates that LookingGlass can be used to
700 discover environmentally relevant features, as well as evolutionary and functional
701 ones. Our optimal temperature classifier may be useful across both academic and
702 commercial applications – for instance, to compare the optimal temperatures of
703 microbial communities across environmental gradients in temperature or
704 geochemistry, or to identify candidate proteins of a particular function and optimal
705 temperature of interest for industrial applications. In addition, it may also be possible
706 to adapt the optimal temperature classifier presented here as a generative model to
707 guide protein design of a desired function and optimal temperature.

708 The LookingGlass model, and the framework for transfer learning presented
709 here, provides a framework for future efforts toward modelling of complex biological
710 systems. LookingGlass captures the complexity of biology and its interactions with the
711 environment, leveraging the full potential of the functional information contained in the
712 massive amount of sequencing data being generated by the scientific community. The
713 LookingGlass model presented here focuses on Bacterial and Archaeal DNA
714 sequences, but low hanging fruit may include a specialized Eukaryotic DNA model, a
715 model specific to the human genome, or a model specialized to a particular
716 environment such as the human gut or soil microbiome. As the scientific community
717 continues to grapple with new approaches to represent and model biological systems
718 in ways that harness the full potential of our expanding data resources, we hope that
719 LookingGlass can provide a foundation for transfer learning-based exploration of life
720 on Earth.

721

722

Acknowledgements

723 The authors would like to thank Paul Falkowski and the rest of the Rutgers
724 ENIGMA team for productive discussions of the deep transfer learning approach and
725 inspiration for downstream applications of the LookingGlass model.

726

727

Author Contributions

728

AH conceived of the project, compiled data, carried out training, validation, and application of models, and deployed open source code and software. YB provided feedback throughout the project. AA and GF curated the optimal growth temperature dataset. All authors contributed to writing of the manuscript.

729

730

731

732

733

Funding

734

This work was supported by a NASA Astrobiology Postdoctoral Fellowship (to AH) within the NAI Grant Number: 80NSSC18M0093 (to YB). YB was also supported by the NSF (National Science Foundation) CAREER award 1553289. Additional computing resources were provided by a Microsoft AI For Earth grant (to AH).

735

736

737

738

739

Competing Interest Statement

740

The authors declare no competing interests.

741 **References**

- 742 1. Hoarfrost, A. LookingGlass release v1.0.
743 <https://github.com/ahoarfst/LookingGlass/>. (2020).
744 doi:10.5281/zenodo.4382930
- 745 2. Hoarfrost, A. fastBio: deep learning for biological sequences. Github repository
746 and python package. <https://github.com/ahoarfst/fastBio/>. (2020).
747 doi:10.5281/zenodo.4383283
- 748 3. Hoarfrost, A. Github repository - LoL: learning the Language of Life.
749 <https://github.com/ahoarfst/LoL/>. doi:10.5281/zenodo.4362588
- 750 4. Lloyd, K. G., Steen, A. D., Ladau, J., Yin, J. & Crosby, L. Phylogenetically
751 Novel Uncultured Microbial Cells Dominate Earth Microbiomes. *mSystems* **3**,
752 e00055-18 (2018).
- 753 5. Steen, A. D. *et al.* High proportions of bacteria and archaea across most
754 biomes remain uncultured. *ISME J.* **13**, 3126–3130 (2019).
- 755 6. Lobb, B., Tremblay, B. J. M., Moreno-Hagelsieb, G. & Doxey, A. C. An
756 assessment of genome annotation coverage across the bacterial tree of life.
757 *Microb. Genomics* **6**, (2020).
- 758 7. Metagenomics versus Moore’s law. *Nat. Methods* **6**, 623 (2009).
- 759 8. Eraslan, G., Avsec, Ž., Gagneur, J. & Theis, F. J. Deep learning: new
760 computational modelling techniques for genomics. *Nat. Rev. Genet.* **20**, 389–
761 403 (2019).
- 762 9. Thrun, S. Is Learning The n-th Thing Any Easier Than Learning The First? *Adv.*
763 *Neural Inf. Process. Syst.* **7** (1996). doi:10.1.1.44.2898
- 764 10. Pan, S. J. & Yang, Q. A survey on transfer learning. *IEEE Trans. Knowl. Data*
765 *Eng.* **22**, 1345–1359 (2010).
- 766 11. Yosinski, J., Clune, J., Bengio, Y. & Lipson, H. How transferable are features in
767 deep neural networks? *Adv. Neural Inf. Process. Syst.* 1–9 (2014).
- 768 12. Devlin, J., Chang, M. W., Lee, K. & Toutanova, K. BERT: Pre-training of deep
769 bidirectional transformers for language understanding. *NAACL HLT 2019 -*
770 *2019 Conf. North Am. Chapter Assoc. Comput. Linguist. Hum. Lang. Technol. -*
771 *Proc. Conf.* **1**, 4171–4186 (2019).

- 772 13. Liu, H., Perl, Y. & Geller, J. Transfer Learning from BERT to Support Insertion
773 of New Concepts into SNOMED CT. *AMIA ... Annu. Symp. proceedings. AMIA*
774 *Symp.* **2019**, 1129–1138 (2019).
- 775 14. Peng, Y., Yan, S. & Lu, Z. Transfer Learning in Biomedical Natural Language
776 Processing: An Evaluation of BERT and ELMo on Ten Benchmarking
777 Datasets. 58–65 (2019). doi:10.18653/v1/w19-5006
- 778 15. Fofanov, Y. *et al.* How independent are the appearances of n-mers in different
779 genomes? *Bioinformatics* **20**, 2421–2428 (2004).
- 780 16. Kelley, D. R., Snoek, J. & Rinn, J. L. Basset : learning the regulatory code of
781 the accessible genome with deep convolutional neural networks. 990–999
782 (2016). doi:10.1101/gr.200535.115.Freely
- 783 17. Taroni, J. N. *et al.* MultiPLIER : A Transfer Learning Framework for
784 Transcriptomics Reveals Systemic Features of Rare Article MultiPLIER : A
785 Transfer Learning Framework for Transcriptomics Reveals Systemic Features
786 of Rare Disease. *Cell Syst.* **8**, 380-394.e4 (2019).
- 787 18. Menegaux, R. & Vert, J. P. Continuous Embeddings of DNA Sequencing
788 Reads and Application to Metagenomics. *J. Comput. Biol.* **26**, 509–518 (2019).
- 789 19. EIAbd, H. *et al.* Amino acid encoding for deep learning applications. *BMC*
790 *Bioinformatics* **21**, 235 (2020).
- 791 20. Viehweger, A., Krautwurst, S., Parks, D. H., König, B. & Marz, M. An encoding
792 of genome content for machine learning. *bioRxiv* 524280 (2019).
793 doi:10.1101/524280
- 794 21. Heinzinger, M. *et al.* Modeling aspects of the language of life through transfer-
795 learning protein sequences. *BMC Bioinformatics* **20**, 1–17 (2019).
- 796 22. Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified
797 rational protein engineering with sequence-based deep representation
798 learning. *Nat. Methods* **16**, 1315–1322 (2019).
- 799 23. Rao, R. *et al.* Evaluating Protein Transfer Learning with TAPE. *33rd Conf.*
800 *Neural Inf. Process. Syst. (NeurIPS 2019)* (2019). doi:10.1101/676825
- 801 24. Rives, A. *et al.* Biological Structure and Function Emerge from Scaling
802 Unsupervised Learning to 250 Million Protein Sequences. *bioRxiv* 1–31 (2019).
803 doi:<https://doi.org/10.1101/622803>

- 804 25. Bennett, S. Solexa Ltd. *Pharmacogenomics* **5**, 433–8 (2004).
- 805 26. Parks, D. H. *et al.* A standardized bacterial taxonomy based on genome
806 phylogeny substantially revises the tree of life. *Nat. Biotechnol.* **36**, 996 (2018).
- 807 27. Agarwala, R. *et al.* Database resources of the National Center for
808 Biotechnology Information. *Nucleic Acids Res.* **46**, D8–D13 (2018).
- 809 28. Hoarfrost, A., Brown, N., Brown, C. T. & Arnosti, C. Sequencing data discovery
810 with MetaSeek. *Bioinformatics* **35**, 4857–4859 (2019).
- 811 29. Jordan, M. I. Attractor dynamics and parallelism in a connectionist sequential
812 machine. *Proc. Cogn. Sci. Soc.* 531–546 (1986).
- 813 30. Hochreiter, S. & Schmidhuber, J. Long Short-Term Memory. *Neural Comput.* **9**,
814 1735–1780 (1997).
- 815 31. Yoshua Bengio, Patrice Simard & Paolo Frasconi. Learning Long-term
816 Dependencies with Gradient Descent is Difficult. *IEEE Trans. Neural Netw.* **5**,
817 157 (2014).
- 818 32. Merity, S., Keskar, N. S. & Socher, R. Regularizing and Optimizing LSTM
819 Language Models. (2015).
- 820 33. Howard, J. & Gugger, S. Fastai: A layered api for deep learning. *arXiv* (2020).
821 doi:10.3390/info11020108
- 822 34. Kingma, D. P. & Ba, J. L. Adam: A Method for Stochastic Optimization. 1–15
823 (2015).
- 824 35. Leinonen, R., Sugawara, H. & Shumway, M. The sequence read archive.
825 *Nucleic Acids Res.* **39**, 2010–2012 (2011).
- 826 36. Yilmaz, P. *et al.* Minimum information about a marker gene sequence
827 (MIMARKS) and minimum information about any (x) sequence (MIxS)
828 specifications. *Nat. Biotechnol.* **29**, 415–420 (2011).
- 829 37. Zhu, C. *et al.* Functional sequencing read annotation for high precision
830 microbiome analysis. *Nucleic Acids Res.* **46**, (2018).
- 831 38. Li, W. & Godzik, A. Cd-hit: A fast program for clustering and comparing large
832 sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658–1659 (2006).
- 833 39. Consortium, T. U. UniProt: A worldwide hub of protein knowledge. *Nucleic*

- 834 *Acids Res.* **47**, D506–D515 (2019).
- 835 40. Kanz, C. *et al.* The EMBL nucleotide sequence database. *Nucleic Acids Res.*
836 **33**, 29–33 (2005).
- 837 41. Ioffe, S. & Szegedy, C. Batch Normalization : Accelerating Deep Network
838 Training by Reducing Internal Covariate Shift. (2015).
- 839 42. Team, R. C. R: A language and environment for statistical computing. (2017).
- 840 43. Kriventseva, E. V *et al.* OrthoDB v8: Update of the hierarchical catalog of
841 orthologs and the underlying free software. *Nucleic Acids Res.* **43**, D250–D256
842 (2015).
- 843 44. Cock, P. J. A. *et al.* Biopython: Freely available Python tools for computational
844 molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 845 45. Howard, J. & Ruder, S. Universal Language Model Fine-tuning for Text
846 Classification. *arXiv* (2018). doi:arXiv:1801.06146v3
- 847 46. Sunagawa, S. *et al.* Structure and function of the global ocean microbiome.
848 *Science* (80-.). **348**, 1–10 (2015).
- 849 47. Meyer, F. *et al.* The metagenomics RAST server - A public resource for the
850 automatic phylogenetic and functional analysis of metagenomes. *BMC*
851 *Bioinformatics* **9**, 1–8 (2008).
- 852 48. Reimer, L. C. *et al.* BacDive in 2019: Bacterial phenotypic data for High-
853 throughput biodiversity analysis. *Nucleic Acids Res.* **47**, D631–D636 (2019).
- 854 49. Parte, A. C., Carbasse, J. S., Meier-Kolthoff, J. P., Reimer, L. C. & Göker, M.
855 List of prokaryotic names with standing in nomenclature (LPSN) moves to the
856 DSMZ. *Int. J. Syst. Evol. Microbiol.* **70**, 5607–5612 (2020).
- 857 50. Kawachi, M. & Noël, M. H. Microbial Culture Collection at the National Institute
858 for Environmental Studies, Tsukuba, Japan. *PICES Press* **22**, 43 (2014).
- 859 51. Aptekmann, A. A. & Nadra, A. D. Core promoter information content correlates
860 with optimal growth temperature. *Sci. Rep.* **8**, 1–7 (2018).
- 861 52. Engqvist, M. K. M. Correlating enzyme annotations with a large set of microbial
862 growth temperatures reveals metabolic adaptations to growth at diverse
863 temperatures 06 Biological Sciences 0605 Microbiology 06 Biological Sciences
864 0601 Biochemistry and Cell Biology. *BMC Microbiol.* **18**, 1–14 (2018).

- 865 53. Wheeler, D. L. *et al.* Database Resources of the National Center for
866 Biotechnology Information. *Nucleic Acids Res.* **33**, D39–D45 (2016).
- 867 54. Tamames, J., Cobo-Simón, M. & Puente-Sánchez, F. Assessing the
868 performance of different approaches for functional and taxonomic annotation of
869 metagenomes. *BMC Genomics* **20**, 1–16 (2019).
- 870 55. Konstantinidis, K. T. & Tiedje, J. M. Trends between gene content and genome
871 size in prokaryotic species with larger genomes. *Proc. Natl. Acad. Sci. U. S. A.*
872 **101**, 3160–3165 (2004).
- 873 56. Sheridan, P. P., Panasik, N., Coombs, J. M. & Brenchley, J. E. Approaches for
874 deciphering the structural basis of low temperature enzyme activity. *Biochim.*
875 *Biophys. Acta - Protein Struct. Mol. Enzymol.* **1543**, 417–433 (2000).
- 876 57. Li, W. F., Zhou, X. X. & Lu, P. Structural features of thermozyms. *Biotechnol.*
877 *Adv.* **23**, 271–281 (2005).
- 878 58. Falkowski, P. G., Fenchel, T. & Delong, E. F. The microbial engines that drive
879 Earth's biogeochemical cycles. *Science* **320**, 1034–9 (2008).
- 880 59. Clemente, J. C., Ursell, L. K., Parfrey, L. W. & Knight, R. The impact of the gut
881 microbiota on human health: An integrative view. *Cell* **148**, 1258–1270 (2012).
- 882 60. Hug, L. *et al.* A new view of the tree of life. *Nat Microbiol* **1**, 16048 (2016).
- 883 61. Pham, J. V. *et al.* A review of the microbial production of bioactive natural
884 products and biologics. *Front. Microbiol.* **10**, (2019).
- 885 62. Song, H., Ding, M. Z., Jia, X. Q., Ma, Q. & Yuan, Y. J. Synthetic microbial
886 consortia: From systematic analysis to construction and applications. *Chem.*
887 *Soc. Rev.* **43**, 6954–6981 (2014).
- 888 63. Salazar, G. *et al.* Gene Expression Changes and Community Turnover
889 Differentially Shape the Global Ocean Metatranscriptome. *Cell* **179**, 1068-
890 1083.e21 (2019).
- 891 64. Ramírez-Flandes, S., González, B. & Ulloa, O. Redox traits characterize the
892 organization of global microbial communities. *Proc. Natl. Acad. Sci. U. S. A.*
893 **116**, 3630–3635 (2019).
- 894 65. Fuhrman, J. A. *et al.* A latitudinal diversity gradient in planktonic marine
895 bacteria. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 7774–8 (2008).

- 896 66. Ibarbalz, F. M. *et al.* Global Trends in Marine Plankton Diversity across
897 Kingdoms of Life AR OCEANS EXPEDITION Article Global Trends in Marine
898 Plankton Diversity across Kingdoms of Life. 1084–1097 (2019).
899 doi:10.1016/j.cell.2019.10.008
- 900 67. Sul, W. J., Oliver, T. A., Ducklow, H. W., Amaral-Zettlera, L. A. & Sogin, M. L.
901 Marine bacteria exhibit a bipolar distribution. *Proc. Natl. Acad. Sci. U. S. A.*
902 **110**, 2342–2347 (2013).
- 903 68. Ghiglione, J.-F. *et al.* Pole-to-pole biogeography of surface and deep marine
904 bacterial communities. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 17633–8 (2012).
- 905 69. Ladau, J. *et al.* Global marine bacterial diversity peaks at high latitudes in
906 winter. *ISME J.* **7**, 1669–77 (2013).
- 907 70. Raes, E. J. *et al.* Oceanographic boundaries constrain microbial diversity
908 gradients in the south pacific ocean. *Proc. Natl. Acad. Sci. U. S. A.* **115**,
909 E8266–E8275 (2018).
- 910 71. Giovannoni, S. J., Cameron Thrash, J. & Temperton, B. Implications of
911 streamlining theory for microbial ecology. *ISME J.* **8**, 1553–1565 (2014).
- 912 72. Ulloa, O., Canfield, D. E., DeLong, E. F., Letelier, R. M. & Stewart, F. J.
913 Microbial oceanography of anoxic oxygen minimum zones. *Proc. Natl. Acad.*
914 *Sci. U. S. A.* **109**, 15996–16003 (2012).
- 915 73. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using
916 DIAMOND. *Nat. Methods* **12**, 59–60 (2014).
- 917 74. Nazeen, S., Yu, Y. W. & Berger, B. Carnelian uncovers hidden functional
918 patterns across diverse study populations from whole metagenome sequencing
919 reads. *Genome Biol.* **21**, 1–18 (2020).
- 920 75. Seemann, T. Prokka: Rapid prokaryotic genome annotation. *Bioinformatics* **30**,
921 2068–2069 (2014).

922

923

924

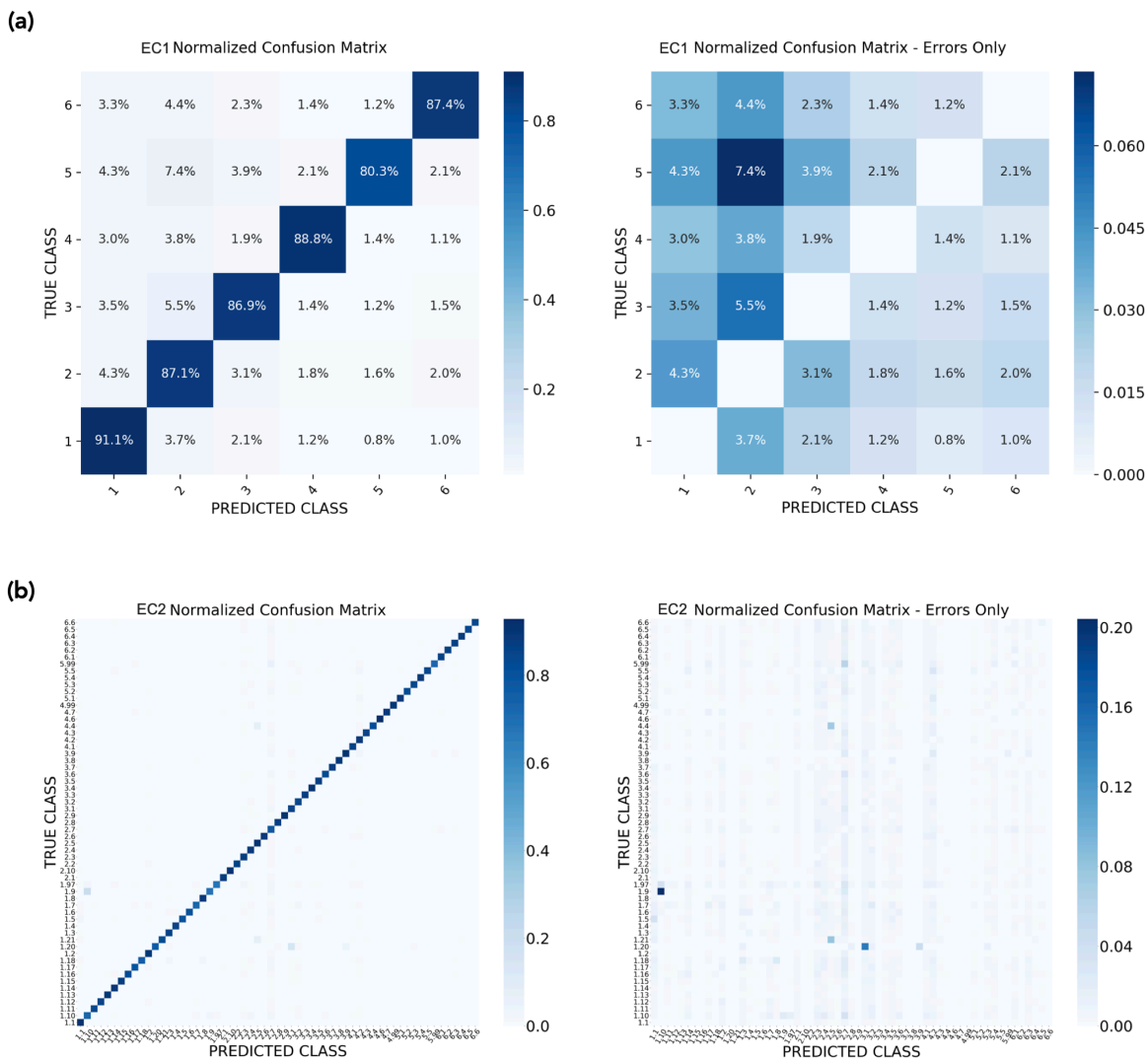
925

926

927

Figures

928



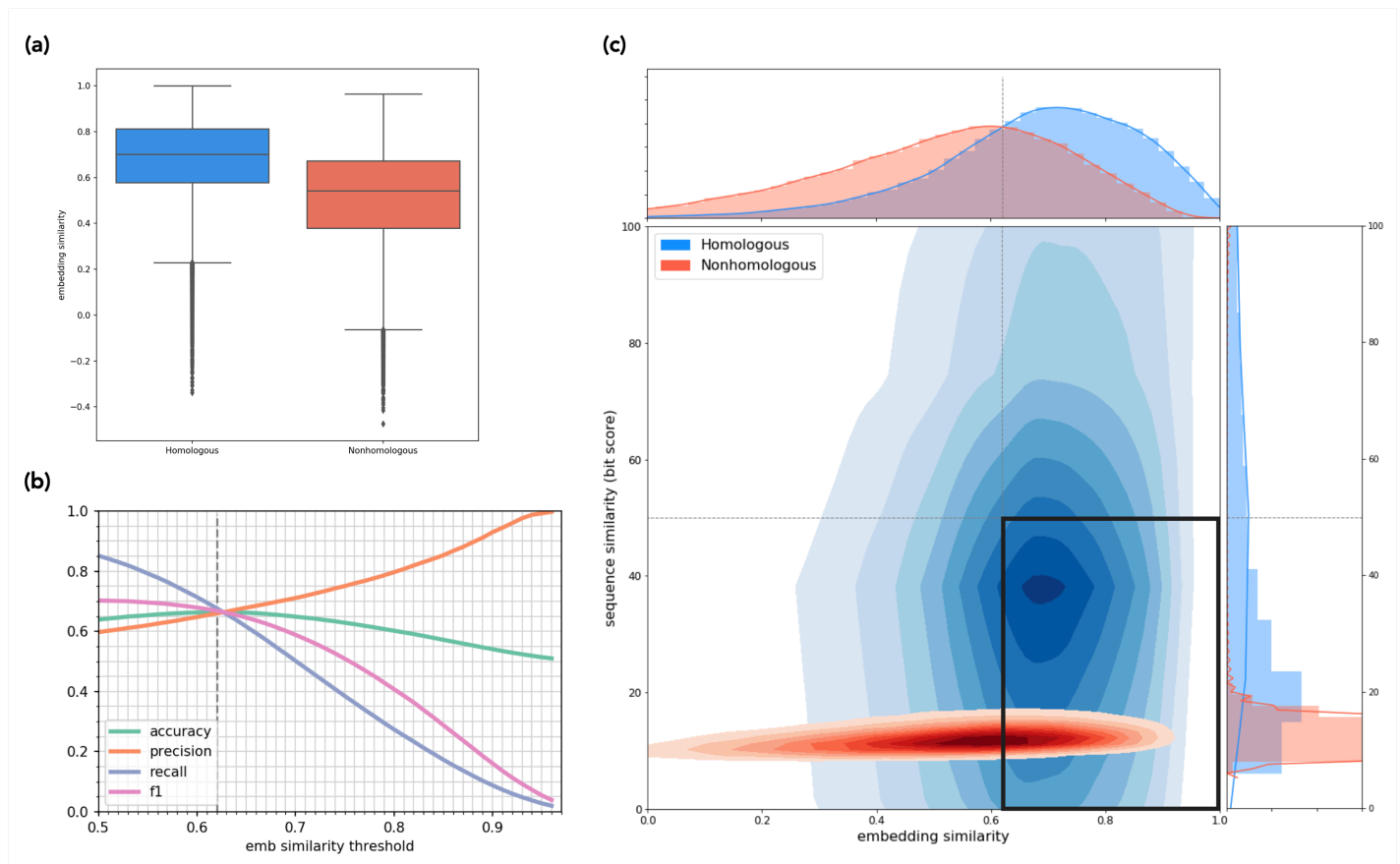
929

Fig. 1: Functional annotation prediction multiclass confusion matrix. Confusion between true (y axis) and predicted (x axis) functional annotations, shown as normalized percentages of predictions for each label including correct predictions (left) and showing errors only (right), for (a) predictions to the 1st EC number and (b) predictions to the 2nd EC number.

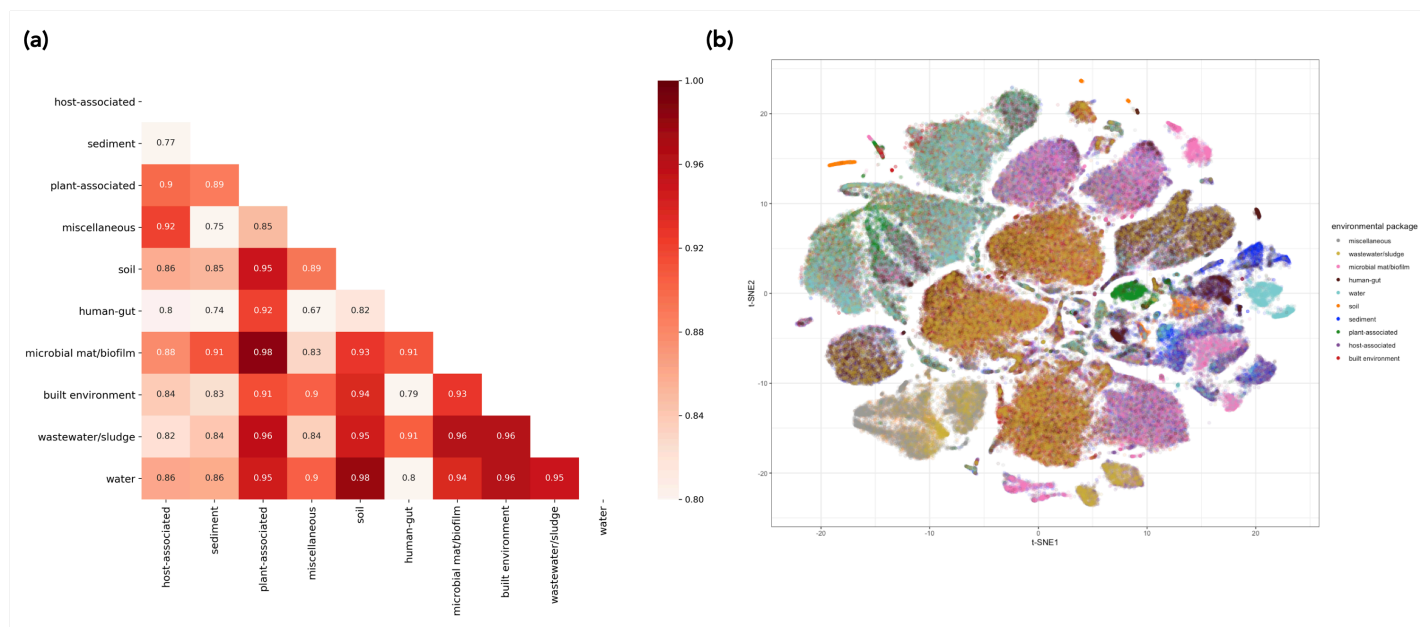
930

931

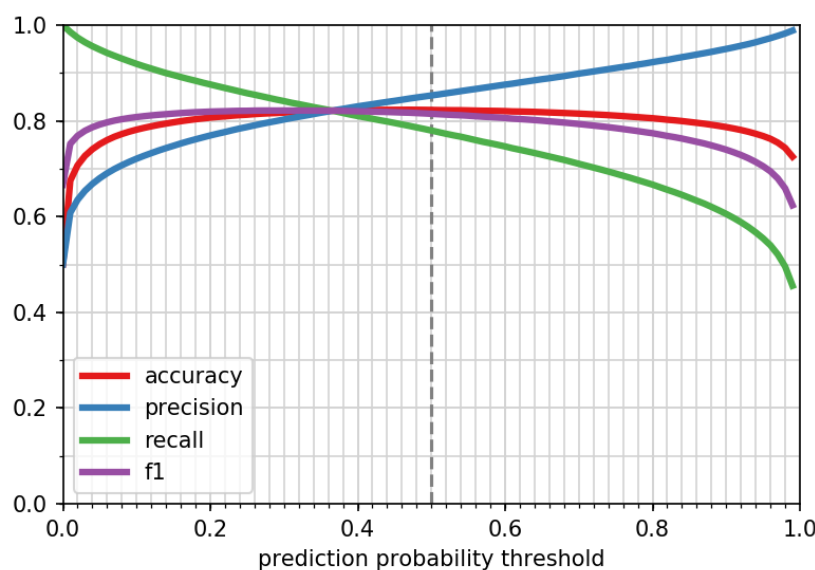
932



933 **Fig. 2: LookingGlass identifies homologous sequence pairs at the phylum level.** (a) Distribution
934 of embedding similarities for homologous (blue) and nonhomologous (red) sequence pairs are
935 significantly different ($P < 10^{-16}$). (b) Accuracy, precision, recall, and F1 metrics (Eqns 1-4) for
936 homologous/ nonhomologous predictions across embedding similarity thresholds. Default threshold of
937 maximum accuracy (0.62) shown in vertical dashed line. (c) Distribution of embedding and
938 sequencing similarities for homologous (blue) and nonhomologous (red) sequence pairs. 44% of
939 homologous sequence pairs have sequence similarity alignment scores below the threshold of 50
940 (horizontal line). Embedding similarity threshold (0.62, vertical line) separates homologous and
941 nonhomologous sequence pairs with maximum accuracy. Bold black box in the lower right indicates
942 homologous sequences correctly identified by LookingGlass that are missed using alignments.

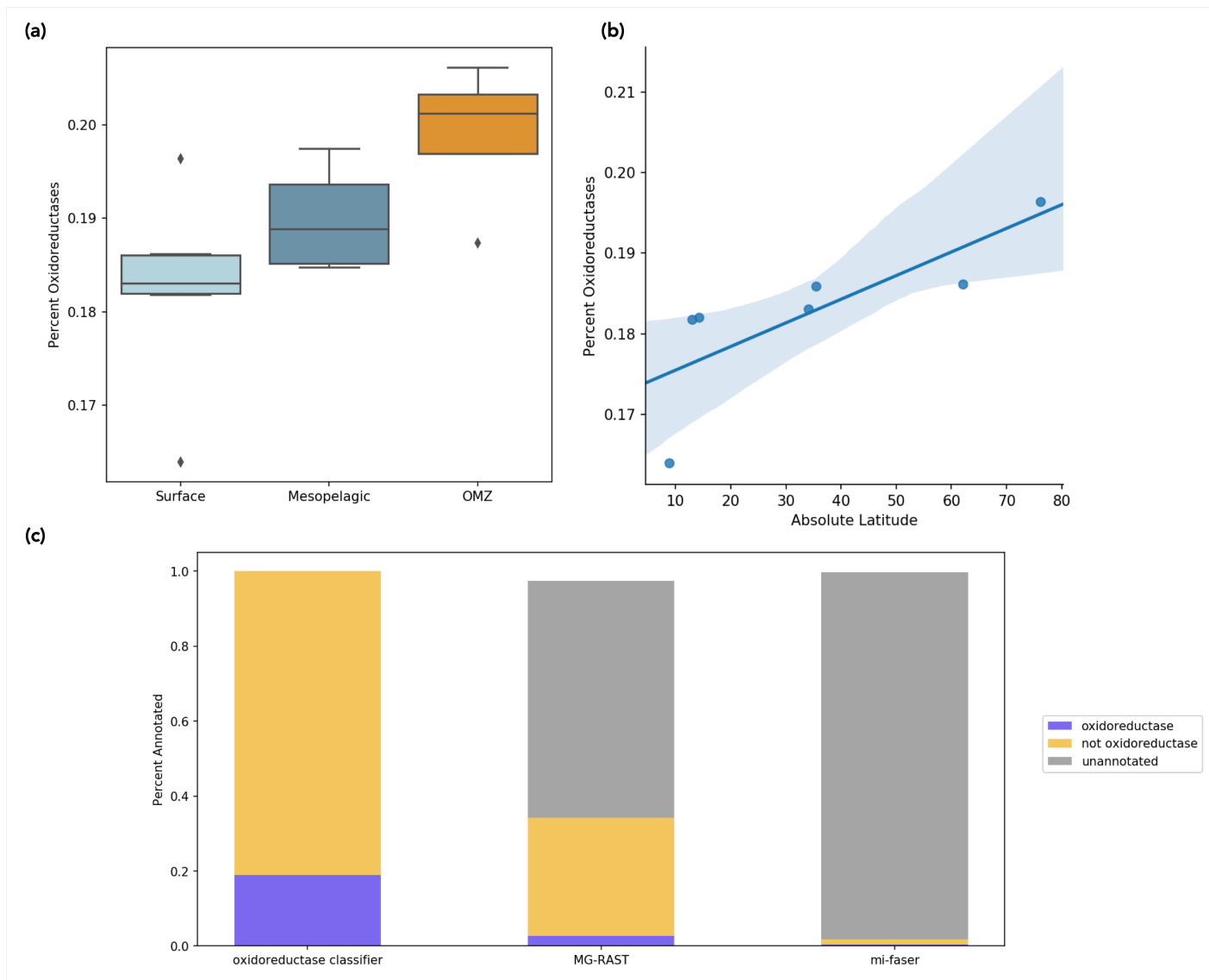


943 **Fig. 3: Distributions of LookingGlass embeddings across environmental packages.** (a) Pairwise
 944 cosine similarity among the average embeddings of 20,000 randomly selected sequences from each
 945 environmental package. (b) t-SNE visualization of the embedding space for 20,000 randomly selected
 946 sequences from each of ten distinct environmental contexts in the *mi-faser functional* validation set.
 947 Sequences from the same environmental context generally cluster together. Colors indicate
 948 environmental package. Embeddings are significantly differentiated by environmental package ($P <$
 949 10^{-16}).



950

951 **Fig. 4: Performance of the oxidoreductase classifier.** Accuracy, precision, recall, and F1 score
 952 metrics (Eqns 1-4) of the oxidoreductase classifier across prediction probability thresholds. Default
 953 threshold of 0.5 shown in vertical dashed line.



954 **Fig. 5: Oxidoreductase identification in marine metagenomes.** (a) Proportion of oxidoreductase
955 sequences (y axis) predicted by the oxidoreductase classifier in surface, mesopelagic, and oxygen
956 minimum zone (OMZ) depths. (b) Correlation between the proportion of oxidoreductases and absolute
957 degrees latitude in surface metagenomes of the *oxidoreductase metagenome set* ($R^2=0.79$, $P=0.04$).
958 (c) Proportion of sequences predicted as oxidoreductases, not oxidoreductases, or left unannotated
959 across the oxidoreductase classifier, MG-RAST, and mi-faser tools.

960

961

962

963

964

965

Tables

966

Dataset Name	Dataset Description
<i>GTDB representative set</i>	Read-length DNA sequences from each of the 24,706 Bacterial and Archaeal representative genomes in the GTDB ²⁶
<i>GTDB class set</i>	Reduced set of read-length sequences from a representative genome of each class in the GTDB ²⁶ taxonomy
<i>mi-faser functional set</i>	Functionally annotated reads from 100 metagenomes from evenly distributed environmental packages
<i>Swiss-Prot functional set</i>	DNA read-length sequences of genes with experimentally validated functions from the Swiss-Prot database
<i>OG homolog set</i>	Homologous and nonhomologous sequence pairs of gene sequences from 1,000 orthologous groups from the OrthoDB database defined at multiple taxonomic levels: genus, family, order, class, and phylum
<i>Oxidoreductase model set</i>	Read-length DNA sequences from genes corresponding to Bacterial and Archaeal oxidoreductases from the manually reviewed entries of the Swiss-Prot database
<i>Oxidoreductase metagenome set</i>	Sequencing reads from 16 marine metagenomes, rarefied to 20 million sequences each, from latitudes spanning -62 to 76 degrees and two depths – surface and mesopelagic. Mesopelagic depths at 4 stations corresponded to an oxygen minimum zone (OMZ)
<i>Reading frame set</i>	Read-length sequences, and labels corresponding to their true frame of translation, for gene coding sequences from one genome selected from each order in the GTDB taxonomy
<i>Optimal temp set</i>	Read-length sequences from core genes associated with transcription and translation, and labels corresponding to their optimal enzyme temperature, inferred from the manually curated optimal growth temperature of 19,474 genomes.

967

Table 1 – Summary table of datasets used.