

DIVERSE GENOMIC EMBEDDING BENCHMARK FOR FUNCTIONAL EVALUATION ACROSS THE TREE OF LIFE

Anonymous authors

Paper under double-blind review

ABSTRACT

Biological foundation models hold significant promise for deciphering complex biological functions. However, evaluating their performance on functional tasks remains challenging due to the lack of standardized benchmarks encompassing diverse sequences and functions. Existing functional annotations are often scarce, biased, and susceptible to train-test leakage, hindering robust evaluation. Furthermore, biological functions manifest at multiple scales, from individual residues to large genomic segments. To address these limitations, we introduce the Diverse Genomic Embedding Benchmark (DGEB), inspired by natural language embedding benchmarks. DGEB comprises six embedding tasks across 18 expert curated datasets, spanning sequences from all domains of life and encompassing both nucleic acid and amino acid modalities. Notably, four datasets enable direct comparison between models trained on different modalities. Benchmarking protein and genomic language models (pLMs and gLMs) on DGEB reveals performance saturation with model scaling on numerous tasks, especially on those with underrepresented sequences (e.g. Archaea). This highlights the limitations of existing modeling objectives and training data distributions for capturing diverse biological functions. DGEB is available as an open-source package with a public leaderboard at [URLhiddenforanonymity](https://urlhiddenforanonymity).

1 INTRODUCTION

Biological sequences encode complex molecular, evolutionary and biophysical information that govern biological function. Deep learning models have been proposed as promising methods for extracting biologically relevant functional information from sequence data. The promise of "biological foundation models" enabling functional interpretation of sequences has resulted in many modeling efforts in protein (Rives et al., 2021; Madani et al., 2023; Elnaggar et al., 2022) and genomic (Dalla-Torre et al., 2023; Hwang et al., 2024; Nguyen et al., 2024) sequence modalities. While the field has seen major advances in AI-enabled structure prediction of protein sequences (Jumper et al., 2021; Baek et al., 2021), validated successes for AI-enabled function prediction remain limited (Li et al., 2024). Slow progress in function prediction of sequences can be attributed to the following main challenges:

1. **Unlike for structural prediction tasks, objective measurements of function do not exist.** Structure prediction tasks benefit from objective evaluation metrics based on quantifiable atomic distances (Mariani et al., 2013). However, biological function is inherently multifaceted and context-dependent, making direct quantitative assessment difficult.
2. **Functional labels are sparse, biased, and prone to leakage.** Labels are heavily biased towards model organisms (e.g. Human), therefore performance on species-specific evaluation tasks are not guaranteed to transfer to other organisms. Furthermore, functional annotations in databases are rarely standardized in format, necessitating careful curation (e.g. unification of synonymous text labels requires expert knowledge). Critically, all biological sequences are related through evolu-

047 tion. Without carefully designed parameters, train-test leakage can frequently occur, resulting in
048 unreliable evaluation results (Fang, 2023).

- 049
050 3. **Biological function takes place across diverse scales.** Single nucleotide polymorphisms can have
051 phenotypic effects, while entire segments of genomes can be coordinated to carry out singular
052 functions (e.g. biosynthetic gene clusters). These challenges innate to biological data have led to
053 the lack of diverse benchmarks, resulting in independent evaluations of models on biased sets of
054 "in-house" tasks, preventing comprehensive and objective model comparisons.

055 The Diverse Genomic Embedding Benchmark (DGEB) is inspired by text embedding evaluation bench-
056 marks that have advanced the field of natural language modeling. DGEB aims to span diverse types of
057 downstream embedding tasks, scopes of function, and taxonomic lineages. DGEB consists of 18 datasets
058 covering 117 phyla across all three domains of life (Bacteria, Archaea and Eukarya). Similar to Massive
059 Text Embedding Benchmark (MTEB) (Muennighoff et al., 2023), DGEB evaluates embeddings using six
060 different embedding tasks: Classification, BiGene mining, Evolutionary Distance Similarity (EDS), Pair
061 classification, Clustering, and Retrieval. DGEB focuses on evaluating the representations of higher-order
062 functional and evolutionary relationships of genomic elements, and is designed to complement existing
063 benchmarks that focus on residue-level representations ((Notin et al., 2023) (Marin et al., 2024)).

064 We provide DGEB as an open source software, facilitating the evaluation of custom models, and enabling
065 the addition and revision of datasets. Biological labels for function are limited and rely on careful curation
066 by domain experts. DGEB provides a much-needed infrastructure for allowing experts to contribute new
067 benchmarks and revise datasets upon acquisition of new knowledge. Community driven efforts to collect
068 and standardize diverse datasets will move the emerging interdisciplinary field of Machine Learning and
069 Biology forward.

070 2 RELATED WORKS

071 2.1 NATURAL LANGUAGE EMBEDDING BENCHMARKS

072 Embedding benchmarks (e.g. SentEval (Conneau & Kiela, 2018); BEIR (Thakur et al., 2021); MTEB
073 (Muennighoff et al., 2023)) in natural language processing (NLP) aim to evaluate how the structure of
074 word/sentence representations match the geometric structure of their semantics. For natural language, tasks
075 are typically either zero-shot or few-shot; examples of such tasks range from distance-based matching of
076 translated texts to classifying tweets based on the labeled sentiment. NLP benchmarks highlight the need
077 for holistic evaluation of models through a diverse set of tasks, as model performance can vary significantly
078 across tasks and datasets.

079 2.2 BIOLOGICAL SEQUENCE AND LANGUAGE MODELS

080 Biological sequence language models are unsupervised models trained on biological sequence data such
081 as proteins or genomic segments. Protein language models (pLMs) have been shown to encode features
082 for protein structure prediction (Lin et al., 2023), enzyme function prediction (Yu et al., 2023) and remote
083 homology search (Liu et al., 2024). More recently, genomic language models (gLMs) have been evaluated
084 on classification of various genomic motifs (e.g. regulatory elements, chromatin features, splicing) (Dalla-
085 Torre et al., 2023) and mutation fitness prediction (Nguyen et al., 2024).

086 2.3 BIOLOGICAL FUNCTION BENCHMARKS

087 Existing benchmarks rely mainly on two types of evaluation to measure biological function:

- 088
089 1. **Fitness prediction of mutations using large-scale datasets collected from deep mutational
090 scanning (DMS) data.** DMS (Fowler & Fields, 2014) uses large-scale mutagenesis and high-
091 throughput sequencing to model fitness landscapes of various mutations (e.g. substitutions and
092 indels) in a single protein. ProteinGym (Notin et al., 2023) leverages diverse DMS datasets to
093

094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140

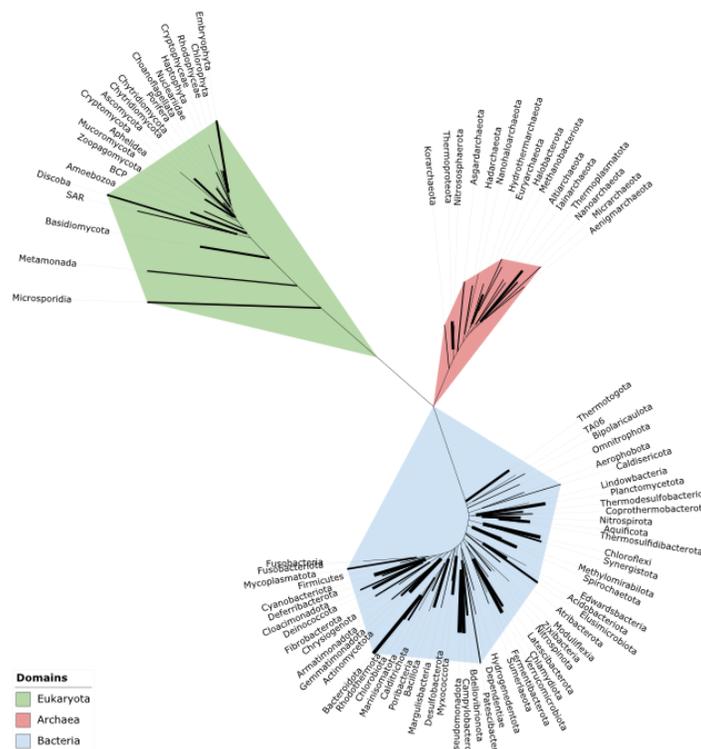


Figure 1: **Phylogenetic tree of all phyla represented in DGEb.** One representative 16S/18S sequence for each phylum represented in any DGEb dataset was obtained from SILVA (Quast et al., 2013), where available. Phylogeny was estimated using iQ-TREE 2. Widths of tree branches correspond to how well a given phylum is represented across multiple datasets.

evaluate a model’s ability to predict fitness scores of mutants in either zero-shot or supervised regimes. While fitness prediction serves as a meaningful proxy for evaluating model understanding of genotype to phenotype relationships at the residue-level for a single protein, this metric cannot be used to determine how well a model can abstract evolutionary and functional relationships between non-homologous proteins.

2. **Classification of proteins on their biophysical properties.** For example, PEER (Xu et al., 2022) benchmarks protein models on various general biophysical properties, such as fluorescence, localization and solubility. These are important properties, they are too coarse in scope to evaluate whether a model has learned more granular functional information (e.g. enzymatic function, protein interaction)

141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187

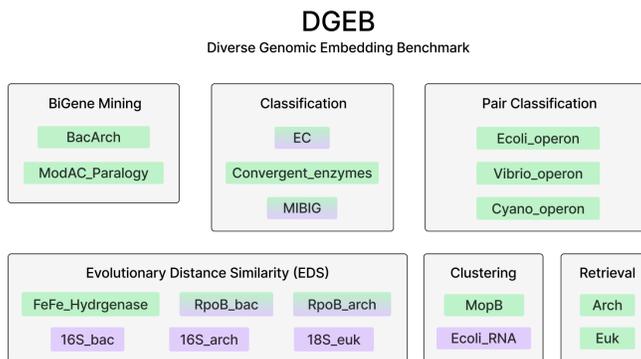


Figure 2: **Overview of tasks and datasets in DGE B.** Nucleic acid (NA) and amino acid (AA) modality specific datasets are marked in purple and green respectively, and datasets that support both modalities are marked with both colors.

3 THE DGE B BENCHMARK

3.1 DESIGN CHOICES

DGE B is built on the desiderata previously outlined by NLP benchmarks, in particular, MTEB. While biological sequences and their functional labels are fundamentally different from natural language, these design choices allow for a scalable and flexible framework that can be expanded and optimized as the field matures.

- Diversity:** We aim to cover sequences derived from phylogenetically diverse lineages of biology (Fig. 1). Existing functional benchmarks largely consist of human or *E.coli* K-12 sequences. Data imbalance in biology is a critical problem when training biological sequence language models (Ding & Steinhardt, 2024) and prevents the models from learning features transferable to underrepresented sequences. Benchmarks that only utilize sequences from highly overrepresented sequences in the training set perpetuate this problem of data imbalance, hindering the progress towards AI-enabled characterization, discovery and design of diverse biological sequences.
- Simplicity:** DGE B provides a simple API that can be used with any custom model that encodes biological sequences into vectors.
- Extensibility:** Given the complexity of biological function, no single dataset can fully capture its diversity, and existing functional annotations must be continuously refined and expanded. DGE B supports simple extension of tasks and datasets. New or revised datasets can be uploaded to the HuggingFace Hub and new evaluation tasks can easily be added through GitHub pull requests.
- Reproducibility:** We version both the software and the datasets and include versioning in the results, making the benchmark results fully reproducible.

3.2 TASKS AND EVALUATION

DGE B consists of 18 datasets that are evaluated using one of the six task types (Fig. 2). The tasks and their evaluation schemes are described below:

BiGene Mining BiGene Mining is inspired by Bitext Mining tasks in NLP, where the tasks typically consist of matching translated sentences between two languages using cosine similarity. For BiGene Mining, we curated functionally analogous sequences found in two phylogenetically distant taxa (e.g. Bacteria and Archaea) or interacting pairs in sets of orthologous sequences. For each gene in the first set, the best match

188 in the second set is found using the cosine similarity. F1 serves as the primary metric for BacArch BiGene
 189 Mining, while recall@50 is used as the primary metric for ModBC BiGene Mining due to the difficulty of
 190 the task; accuracy and precision are also reported.

191 **Evolutionary Distance Similarity (EDS)** This task evaluates how accurately models learn evolutionary
 192 relationships between sequences. We compute the correlation between pairwise embedding distances and
 193 their phylogenetic distances (sum of branch lengths connecting the two leaves of the calculated phylogenetic
 194 tree). Larger phylogenetic distance represents more evolutionary time since divergence. Pearson correlations
 195 are calculated and the top correlation score across three distance metrics (cosine, euclidean, and manhattan)
 196 is reported as the primary metric.

197 **Classification** Classification tasks measure the model’s ability to map from embeddings to discrete func-
 198 tional classes with few-shot supervision. For multiclass single-label classification, a logistic regression
 199 classifier is trained with up to 1000 iterations. For multiclass multi-label classification, a k-nearest neigh-
 200 bor (kNN) classifier is trained. Test performance on the test set is measured using F1 as the main metric;
 201 accuracy and average precision scores are also reported.

202 **Pair Classification** Pair classification tasks evaluate model understanding of functional relationships be-
 203 tween pairs of sequences. Inputs are pairs of sequences, where labels are binary variables denoting the
 204 existence of some particular functional relationship between the pair. Sequences are embedded and the
 205 distances between the pairs are calculated cosine similarity, dot product, euclidean distance and manhattan
 206 distance. The best binary threshold accuracy, average precision, F1, precision, and recall are calculated. The
 207 primary metric is the average precision score calculated using cosine similarity.

208 **Clustering** Clustering tasks evaluate zero-shot separability of embeddings over discrete classes. Inputs are
 209 sets of sequences with labels, and a mini-batch k-means model is trained on their embeddings. The primary
 210 metric is v-measure (Rosenberg & Hirschberg, 2007).

211 **Retrieval** Retrieval tasks evaluate how well a query embedding can retrieve functionally analogous se-
 212 quences. Dataset consists of a corpus and queries, where the objective is to rank the embeddings in the
 213 corpus by cosine similarity to each query sequence. Correct retrieval is determined by matching functional
 214 labels. An example of a retrieval task is retrieving a bacterial homolog given an archaeal query sequence.
 215 nDCG@k, MRR@k, MAP@k, precision@k and recall@k are calculated for k=5, 10, 50. MAP@5 is used
 216 as the primary metric.

217 3.3 DATASETS

218 Datasets are divided into three categories: single-element, inter-element, and multi-element, where an el-
 219 ement refers to a protein/gene or noncoding RNA. Each element can be represented in amino acid and/or
 220 nucleotide sequence modalities. Some datasets support multiple sequence modalities (AA and NA), allow-
 221 ing direct comparison between protein and genomic language models. Statistics for each dataset are found in
 222 Appendix B. All datasets are dereplicated at sequence identity thresholds of 70% using CD-hit (Huang et al.,
 223 2010), to remove sampling biases. For tasks requiring train and test splits, datasets are split with a maximum
 224 sequence identity of 10%. For tasks requiring multiple classes, we conduct class-balanced random sampling.
 225 Detailed preprocessing steps are found in Appendix A.

226 **Single Element Datasets (SE)** For SE datasets, each genomic element (protein/gene, noncoding RNA) is
 227 individually embedded, with an associated label. SE datasets in DGEB include:

- 228 • *RNA Clustering*: rRNA, sRNA, and tRNA features predicted using RFam (Kalvari et al., 2021)
 229 genomes across diverse taxa. We cluster the sequence embeddings and assess how well they match
 230 the RNA class assignments.
- 231 • *MopB Clustering*: The dimethyl sulfoxide reductase (or MopB) family is a functionally diverse set
 232 of enzymes found across Bacteria and Archaea. Sequences are sampled from Wells et al. (2023),
 233
 234

where the sequence’s catalytic functions are assigned using phylogenetic analysis. We assess how well the embeddings cluster with their catalytic function.

- *EC Classification*: Enzyme commission (EC) numbers are assigned to protein sequences. For each EC class, one sequence is randomly selected for testing, and four sequences from the corresponding class that have less than 10% sequence identity to each other and test sequence are selected for training.
- *Convergent Evolution Classification*: Examples of convergent evolution in proteins include enzymes that have different evolutionary history but have converged in the enzymatic reaction that they confer. We identify such convergent enzymes by curating a set of enzymes that have no sequence similarity to any of the other sequences in the train set with the same EC designation.
- *Archaeal Retrieval*: Given the corpus of bacterial protein sequences in SWISS-PROT (Bairoch & Apweiler, 2000), where the label is the corresponding text annotation, we query archaeal sequences with string match annotations in the bacterial corpus. We retrieve k nearest neighbors in bacterial corpus embedding space and look for matching labels to calculate the metrics@k.
- *Eukaryotic Retrieval*: Given the corpus of bacterial protein sequences in SWISS-PROT, we query eukaryotic sequences with string match annotations in the bacterial corpus. Metrics are calculated as above.

Inter-element datasets (IE) Understanding biological function relies on understanding the evolutionary and functional relationships between sequences. For IE datasets, a label is assigned for each pair of genomic embeddings. IE datasets include:

- *BacArch BiGene*: Similar to matching translated sentences between two languages, we curated functionally analogous pairs of sequences in a bacterial genome (*Escherichia coli* K-12) and an archaeal genome (*Sulfolobus acidocaldarius* DSM 639 ASM1228v1).
- *ModBC BiGene*: Identifying interacting pairs of ModB and ModC from sets of orthologs is a challenging task. ModB and ModC are interacting subunits of an ABC transporter. This dataset consists of pairs of ModB and ModC that are found to be interacting in the same genome. The goal is to correctly find the interacting ModC for each ModB given a set of orthologous ModC sequences (found in different genomes).
- *E.coli Operonic Pair Classification*: Given a pair of adjacent proteins, the label is assigned based on whether they belong to the same transcription unit in *Escherichia coli* K-12 substr. MG165.
- *Vibrio Operonic Pair Classification*: Same as *E.coli Operonic Pair Classification* except with *Vibrio cholerae* O1 biovar El Tor str. N16961.
- *Cyano Operonic Pair Classification*: Same as *Ecoli Operonic Pair Classification* except with *Synechococcus elongatus* PCC 7942.
- *FeFeHydrogenase Phylogeny*: Fe-Fe hydrogenases are complex enzymes that carry out important metabolic functions across diverse organisms. They carry out divergent and specific functions including H₂ production, H₂ sensing, H₂ uptake, and CO₂ reduction. Identifying the specific function of these hydrogenases often requires constructing a phylogenetic tree that reconstructs the evolutionary history of the catalytic, or large, subunit. This dataset includes the phylogenetic distances (sum of tree branches connecting the leaves) calculated for all pairs of Fe-Fe hydrogenase sequences.
- *RpoB Bacterial Phylogeny*: RpoB is a ribosomal protein conserved across bacteria and archaea. They are essential single-copy genes and not frequently horizontally transferred, and therefore are often used as phylogenetic marker genes. The RpoB gene is also significantly longer than the

282 Fe-Fe hydrogenase gene making this phylogeny distinctly different from the Fe-Fe hydrogenase
 283 phylogeny. We sample bacterial RpoB sequences utilized as markers in the GTDB (Parks et al.,
 284 2022) database, and calculate the tree to assign phylogenetic distances between pairs of RpoB
 285 sequences.

- 286
- 287 • *RpoB Archaeal Phylogeny*: Same as *RpoB Bacterial Phylogeny* but with archaeal genomes in
 288 GTDB.
- 289 • *Bacterial 16S Phylogeny*: 16S rRNA genes encode ribosomal RNA and are universal across Bacteria
 290 and Archaea. 16S rRNA is often used as a taxonomic marker gene because it rarely undergoes
 291 horizontal gene transfer and has both conserved and variable regions. Bacterial 16S rRNA sequences
 292 were downloaded from the SILVA database (Quast et al., 2013) and phylogenetic distances
 293 were calculated for each pair of sequences.
- 294 • *Archaeal 16S Phylogeny*: Same as *Bacterial 16S Phylogeny* but with archaeal sequences from
 295 SILVA.
- 296 • *Eukaryotic 18S Phylogeny*: Same as *Bacterial 16S Phylogeny* but with 18S rRNA (eukaryotic
 297 homolog of 16S rRNA) from SILVA.

298
 299 **Multi-element datasets (ME)** Many biological functions are carried out by multiple genomic elements in
 300 conjunction. DGEB supports multi-element datasets, where a label is assigned to a larger genomic sequence
 301 containing more than one genes, and whereby a single embedding is calculated either by mean-pooling
 302 across genes, or segments of genome with predefined window size. DGEB currently supports one multi-
 303 element dataset:

- 304 • *MIBiG Classification*: Minimum Information about a Biosynthetic Gene cluster (MIBiG) (Terlouw
 305 et al., 2023) is a database of biosynthetic gene clusters where a genomic segment consisting of
 306 multiple genes synthesize various classes of natural products (e.g. Polyketides, NRPS, etc). A
 307 single genomic segment can synthesize molecules that belong to multiple classes, making this a
 308 multi-label, multi-class classification task. Train and test sets are split at 80/20 using stratified
 309 random sampling.

310 4 RESULTS

311 4.1 MODELS

312 We focus on evaluating self-supervised models pretrained on either amino acid (AA) or nucleic acid (NA)
 313 sequences. These are "foundation models" that are not fine-tuned for specific tasks, and we evaluate how
 314 well the pre-trained embeddings capture various aspects of biological function. For AA models, we evaluate
 315 the ESM2 (Lin et al., 2023) series, ESM3 (Hayes et al., 2024) open model, the ProGen2 (Madani et al.,
 316 2023) series, and the ProtTrans (Elnaggar et al., 2022) models. For NA models, we evaluate the DNABERT-
 317 2 (Zhou et al., 2024), Nucleotide Transformer (NT) (Dalla-Torre et al., 2023) series and the Evo (Nguyen
 318 et al., 2024) models. Notably, we include both masked language models (MLM) and causal language models
 319 (CLM) in our evaluation for both data modalities. To extract sequence-level embeddings, each model's
 320 hidden layer is mean-pool across the sequence dimension, resulting in a fixed-size representation. Model
 321 information is found in Appendix C. Additionally, we provide one-hot baselines for AA and NA sequences,
 322 where the sequence is represented as one-hot vectors per position (Appendix H).

323 4.2 ANALYSIS

324 4.2.1 LAYER PERFORMANCE

325 For all tasks, we test performances of mid- and last hidden layers in the model. For many of the tasks,
 326 the mid layer representation outperforms last layer representations (Fig. 3). This behavior has been noted
 327 in previous studies in both NLP (Rogers et al., 2020) and pLMs (Valeriani et al., 2023), where different
 328

329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375

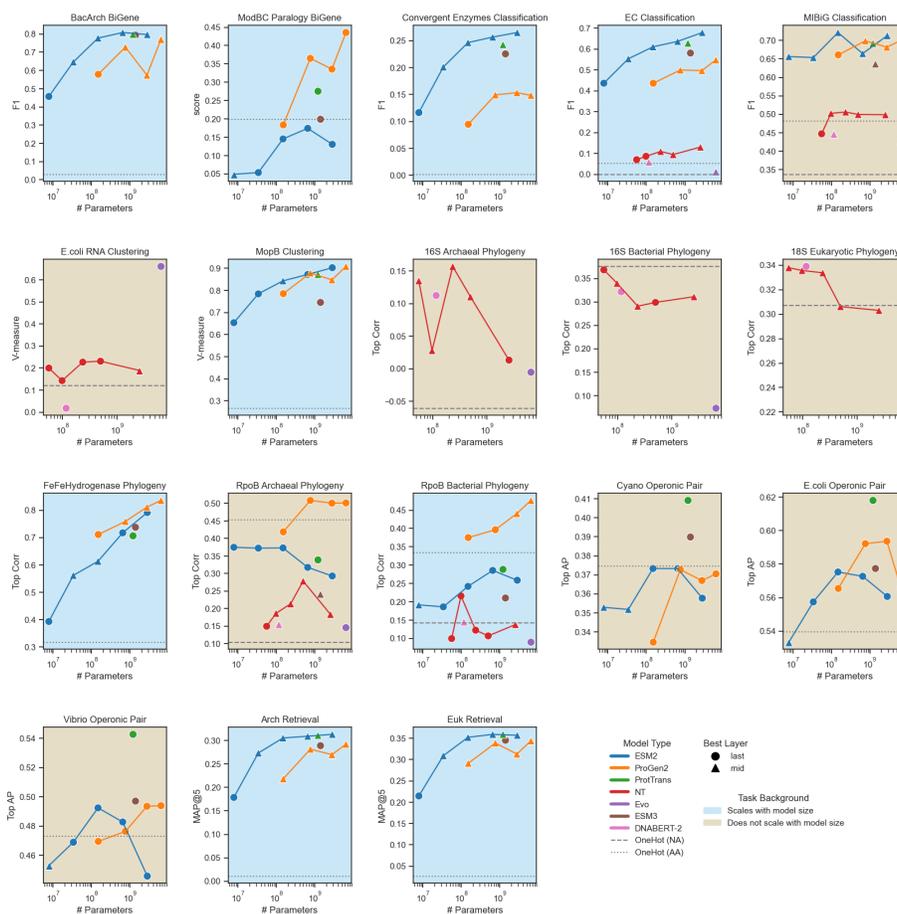


Figure 3: **Performance per task with model scaling for ESM2, ProGen2, and NT series.** Primary metric from the best scoring layer (between mid, and last) is reported for each task. Tasks where performance scales with model size for the majority of the model types are marked with a blue background. Other models plotted for reference are ProtTrans, Evo, and ESM3, and DNABERT-2.

layers specialize in learning distinct semantic information. For instance, mid-layer representations for ESM2 models perform better than last layer for enzyme function classification tasks (EC Classification, Convergent Enzyme Classification) and retrieval tasks, while phylogenetic distances are better reflected in last-layer representations (RpoB phylogenies) (Appendix D). These patterns appear specific to model type. To flexibly account for this behavior, DGEB calculates model performance for both mid and last layer and reports the best score between the two.

4.2.2 SCALING WITH MODEL SIZE

We observe scaling with model size increase for most AA tasks, except for MIBiG classification task, RpoB archaeal phylogeny, and operonic pair tasks (Fig. 3). In general, pLMs perform poorly for predicting functions of elements that span multiple genes (e.g. biosynthetic gene clusters, operons). Additionally, while we observe improved performance with model scaling for bacterial RpoB phylogeny task, we observe no

scaling in performance for archaeal RpoB phylogeny task. This may be attributed to limitations in learning due to the significant bias against archaeal sequences in training data. Interestingly, we observe little to no evidence of improvement in performance with increasing model size for NA tasks (Fig. 3 and 4). We also test performance scaling with pre-training floating point operations (FLOPs) when the information is reported or can be derived (Appendix C). We observe scaling patterns with increasing training FLOPs (Appendix E and F) similar to those observed with increasing model size. Full results can be found in Appendix G.

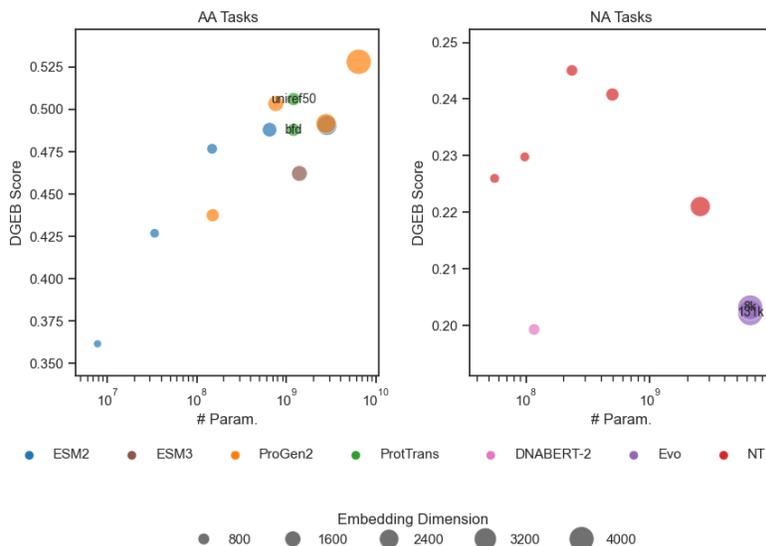


Figure 4: **Average performance across all AA and NA tasks for models benchmarked in this study.** Marker size corresponds to embedding dimension and variants of same models (e.g. evo-1-8k-base, and evo-1-131k-base) are distinguished with text labels.

4.2.3 DIRECT COMPARISON OF AMINO ACID AND NUCLEIC ACID MODELS

While both AA and NA sequences can be used to represent coding sequences, little work has been conducted on directly comparing the quality of NA-based model representations against AA-based model representations on the same task and data. DGEb includes four datasets that support both modalities as input for a given coding region of the sequence. For all such tasks, we find that NA sequence derived representations perform poorly in capturing biological function and evolutionary relationships of coding sequences (Fig. 5). This suggests that AA sequences are a more compute efficient input modality for learning functional information of coding sequences.

5 LIMITATIONS

DGEb includes multiple zero-shot tasks, as ground-truth labels for biological function are sparse and biased. These tasks rely on embedding geometry to evaluate model performance. The assumption that models capturing important features of biological function have geometry directly matching the given tasks is not guaranteed. Future research could explore methods for identifying and leveraging relevant subspaces within model embeddings. For the EDS task, we acknowledge the limitation of Euclidean embeddings for representing phylogenetic tree structures and the possibility that certain regions of the phylogeny may be of low confidence (due to the inherent uncertainty in reconstructing the ground-truth phylogeny). However, this task provides a useful starting point for comparing model performance, and will be important for evaluating

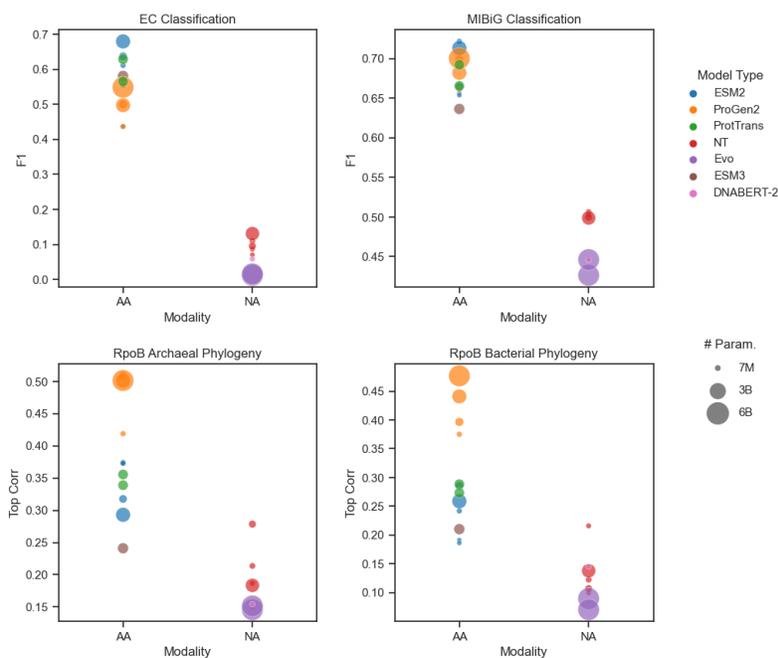


Figure 5: Comparison of AA and NA model representations on tasks that support both modalities. Marker color corresponds to the model type and the size corresponds to the model size.

novel hyperbolic architectures, baselined by Euclidean embedding model results. While DGEb is designed to support both NA and AA models, current suite is biased towards coding sequences with only four tasks targeting non-coding elements (16S_Arch, 16S_Bac, 16S_Euk, Ecoli_RNA), limiting ability to evaluate NA model representations of regulatory elements (e.g. promoters, transcription binding sites). Furthermore, DGEb’s current evaluation suite focuses on single-element, inter-element, and multi-element scales of representations, and is designed to complement existing benchmarks that focus on residue-level representations (e.g. mutational effects (Notin et al., 2023) (Marin et al., 2024)).

6 CONCLUSION

We developed DGEb to assess how well learned embeddings of biological sequences capture various aspects of biological function. Our expert-curated datasets feature diverse sequences spanning all three domains and major phyla in the tree of life. We benchmarked 20 models that are trained on either AA or NA sequences. Our results demonstrate that there is no single model that performs well across all tasks. Importantly, there are many tasks where performance does not scale with model size for existing models, particularly in tasks that feature poorly represented sequences (e.g. Archaeal genes), or tasks that assess functions that require large context lengths (e.g. biosynthetic gene cluster product class classification, operon prediction). For many tasks, there is large headroom for improvement (e.g. ModBC matching, convergent enzyme classification). DGEb also supports direct comparison of models trained on AA and NA data modalities, and our results show that NA models are yet to learn important aspects of biological function. We open-source DGEb to facilitate community-driven dataset addition and revision. We hope that DGEb and the leaderboard allow transparent comparison of biological foundation models and drive the field forward.

470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516

ETHICS STATEMENT

This study aims to advance open science by developing a open-source, reproducible benchmark for genomics. All sequences and labels are curated from public repositories. As the data originates from environmental samples, no personally identifiable information is associated with the datasets.

517 REFERENCES

- 518
519 Minkyung Baek, Frank DiMaio, Ivan Anishchenko, Justas Dauparas, Sergey Ovchinnikov, Gyu Rie Lee, Jue
520 Wang, Qian Cong, Lisa N Kinch, R Dustin Schaeffer, Claudia Millán, Hahnbeom Park, Carson Adams,
521 Caleb R Glassman, Andy DeGiovanni, Jose H Pereira, Andria V Rodrigues, Alberdina A van Dijk, Ana C
522 Ebrecht, Diederik J Opperman, Theo Sagmeister, Christoph Buhlheller, Tea Pavkov-Keller, Manoj K
523 Rathinaswamy, Udit Dalwadi, Calvin K Yip, John E Burke, K Christopher Garcia, Nick V Grishin, Paul D
524 Adams, Randy J Read, and David Baker. Accurate prediction of protein structures and interactions using
525 a three-track neural network. *Science*, August 2021.
- 526 A Bairoch and R Apweiler. The SWISS-PROT protein sequence database and its supplement TrEMBL in
527 2000. *Nucleic Acids Res.*, 28(1):45–48, January 2000.
- 528 Tyler P. Barnum, Alexander Crits-Christoph, Michael Molla, Paul Carini, Henry H. Lee, and Nili Os-
529 trov. Predicting microbial growth conditions from amino acid composition. *bioRxiv*, 2024. doi:
530 10.1101/2024.03.22.586313. URL [https://www.biorxiv.org/content/early/2024/03/](https://www.biorxiv.org/content/early/2024/03/22/2024.03.22.586313)
531 [22/2024.03.22.586313](https://www.biorxiv.org/content/early/2024/03/22/2024.03.22.586313).
- 532 Lionel Breuza, Sylvain Poux, Anne Estreicher, Maria Livia Famiglietti, Michele Magrane, Michael Tognolli,
533 Alan Bridge, Delphine Baratin, Nicole Redaschi, and UniProt Consortium. The UniProtKB guide to the
534 human proteome. *Database*, 2016, February 2016.
- 535 Salvador Capella-Gutiérrez, José M Silla-Martínez, and Toni Gabaldón. trimal: a tool for automated align-
536 ment trimming in large-scale phylogenetic analyses. *Bioinformatics*, 25(15):1972–1973, August 2009.
- 537 Bo Chen, Xingyi Cheng, Pan Li, Yangli ao Geng, Jing Gong, Shen Li, Zhilei Bei, Xu Tan, Boyan Wang, Xin
538 Zeng, Chiming Liu, Aohan Zeng, Yuxiao Dong, Jie Tang, and Le Song. xtrimopglm: Unified 100b-scale
539 pre-trained transformer for deciphering the language of protein. 2024. URL [https://arxiv.org/](https://arxiv.org/abs/2401.06199)
540 [abs/2401.06199](https://arxiv.org/abs/2401.06199).
- 541 Alexis Conneau and Douwe Kiela. SentEval: An evaluation toolkit for universal sentence representations.
542 *arXiv*, 2018.
- 543 Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza-Revilla, Nicolas Lopez Carranza, Adam Henryk
544 Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Bernardo P de Almeida, Hassan
545 Sirelkhatim, Guillaume Richard, Marcin Skwark, Karim Beguir, Marie Lopez, and Thomas Pierrot. The
546 nucleotide transformer: Building and evaluating robust foundation models for human genomics. Septem-
547 ber 2023.
- 548 Sebastian Deorowicz, Agnieszka Debudaj-Grabysz, and Adam Gudyś. FAMSA: Fast and accurate multiple
549 sequence alignment of huge protein families. *Sci. Rep.*, 6:33964, September 2016.
- 550 Frances Ding and Jacob Steinhardt. Protein language models are biased by unequal sequence sampling
551 across the tree of life. March 2024.
- 552 Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Yu Wang, Llion Jones, Tom
553 Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, Debsindhu Bhowmik, and Burkhard Rost.
554 ProtTrans: Toward understanding the language of life through Self-Supervised learning. *IEEE Trans.*
555 *Pattern Anal. Mach. Intell.*, 44(10):7112–7127, October 2022.
- 556 Jianwen Fang. The role of data imbalance bias in the prediction of protein stability change upon mutation.
557 *PLoS One*, 18(3):e0283727, March 2023.
- 558 Douglas M Fowler and Stanley Fields. Deep mutational scanning: a new style of protein science. *Nat.*
559 *Methods*, 11(8):801, August 2014.

- 564 Andrew G. Garrow, Alison Agnew, and David Robert Westhead. Tmb-hunt: An amino acid composition
565 based method to screen proteomes for beta-barrel transmembrane proteins. *BMC Bioinformatics*, 6:56 –
566 56, 2005. URL <https://api.semanticscholar.org/CorpusID:6476508>.
567
- 568 Tomas Hayes, Roshan Rao, Halil Akin, Nicholas J Sofroniew, Deniz Oktay, Zeming Lin, Robert Verkuil,
569 Vincent Q Tran, Jonathan Deaton, Marius Wiggert, Rohil Badkundri, Irhum Shafkat, Jun Gong, Alexander
570 Derry, Raúl Santiago Molina, Neil Thomas, Yousuf A Khan, Chetan Mishra, Carolyn Kim, Liam J Bartie,
571 Matthew Nemeth, Patrick D Hsu, Tom Sercu, Salvatore Candido, and Alexander Rives. Simulating 500
572 million years of evolution with a language model. July 2024.
- 573 Ying Huang, Beifang Niu, Ying Gao, Limin Fu, and Weizhong Li. CD-HIT suite: a web server for clustering
574 and comparing biological sequences. *Bioinformatics*, 26(5):680–682, March 2010.
- 575 Yunha Hwang, Andre L Cornman, Elizabeth H Kellogg, Sergey Ovchinnikov, and Peter R Girguis. Genomic
576 language model predicts protein co-regulation and function. *Nat. Commun.*, 15(1):1–13, April 2024.
577
- 578 John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn
579 Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, Alex Bridgland, Clemens Meyer, Simon
580 A A Kohl, Andrew J Ballard, Andrew Cowie, Bernardino Romera-Paredes, Stanislav Nikolov, Rishub
581 Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman, Ellen Clancy, Michal Zielinski, Martin
582 Steinegger, Michalina Pacholska, Tamas Bergthammer, Sebastian Bodenstein, David Silver, Oriol Vinyals,
583 Andrew W Senior, Koray Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. Highly accurate protein
584 structure prediction with AlphaFold. *Nature*, 596(7873):583–589, August 2021.
- 585 Ioanna Kalvari, Eric P Nawrocki, Nancy Ontiveros-Palacios, Joanna Argasinska, Kevin Lamkiewicz, Manja
586 Marz, Sam Griffiths-Jones, Claire Toffano-Nioche, Daniel Gautheret, Zasha Weinberg, Elena Rivas,
587 Sean R Eddy, Robert D Finn, Alex Bateman, and Anton I Petrov. Rfam 14: expanded coverage of
588 metagenomic, viral and microRNA families. *Nucleic Acids Res.*, 49(D1):D192–D200, January 2021.
- 589 Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray,
590 Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. 2020. URL
591 <https://arxiv.org/abs/2001.08361>.
592
- 593 Peter D Karp, Richard Billington, Ron Caspi, Carol A Fulcher, Mario Latendresse, Anamika Kothari, In-
594 grid M Keseler, Markus Krummenacker, Peter E Midford, Quang Ong, Wai Kit Ong, Suzanne M Paley,
595 and Pallavi Subhraveti. The BioCyc collection of microbial genomes and metabolic pathways. *Brief.*
596 *Bioinform.*, 20(4):1085–1093, July 2019.
- 597 Francesca-Zhoufan Li, Ava P Amini, Yisong Yue, Kevin K Yang, and Alex X Lu. Feature reuse and scaling:
598 Understanding transfer learning with protein language models. February 2024.
- 599 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert
600 Verkuil, Ori Kabeli, Yaniv Shmueli, Allan dos Santos Costa, Maryam Fazel-Zarandi, Tom Sercu, Salva-
601 tore Candido, and Alexander Rives. Evolutionary-scale prediction of atomic-level protein structure with
602 a language model. *Science*, March 2023.
- 603 Wei Liu, Ziyue Wang, Ronghui You, Chenghan Xie, Hong Wei, Yi Xiong, Jianyi Yang, and Shanfeng Zhu.
604 PLMSearch: Protein language model powers accurate and fast sequence search for remote homology.
605 *Nat. Commun.*, 15(1):1–12, March 2024.
- 606 Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton, Jose Luis
607 Olmos, Caiming Xiong, Zachary Z Sun, Richard Socher, James S Fraser, and Nikhil Naik. Large language
608 models generate functional protein sequences across diverse families. *Nat. Biotechnol.*, 41(8):1099–1106,
609 January 2023.
- 610

- 611 Valerio Mariani, Marco Biasini, Alessandro Barbato, and Torsten Schwede. IDDT: a local superposition-
612 free score for comparing protein structures and models using distance difference tests. *Bioinformatics*, 29
613 (21):2722–2728, November 2013.
- 614
615 Frederikke Isa Marin, Felix Teufel, Marc Horlacher, Dennis Madsen, Dennis Pultz, Ole Winther, and Wouter
616 Boomsma. BEND: Benchmarking DNA language models on biologically meaningful tasks. In *The*
617 *Twelfth International Conference on Learning Representations*, 2024. URL <https://openreview.net/forum?id=uKB4cFNQFg>.
- 618
619 Bui Quang Minh, Heiko A Schmidt, Olga Chernomor, Dominik Schrempf, Michael D Woodhams, Arndt
620 von Haeseler, and Robert Lanfear. IQ-TREE 2: New models and efficient methods for phylogenetic
621 inference in the genomic era. *Mol. Biol. Evol.*, 37(5):1530–1534, May 2020.
- 622
623 Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. MTEB: Massive text embedding
624 benchmark. *arXiv*, 2023.
- 625
626 Eric P Nawrocki. Annotating functional RNAs in genomes using infernal. *Methods Mol. Biol.*, 1097:163–
627 197, 2014.
- 628
629 Eric Nguyen, Michael Poli, Matthew G Durrant, Armin W Thomas, Brian Kang, Jeremy Sullivan, Made-
630 lena Y Ng, Ashley Lewis, Aman Patel, Aaron Lou, Stefano Ermon, Stephen A Baccus, Tina Hernandez-
631 Boussard, Christopher Ré, Patrick D Hsu, and Brian L Hie. Sequence modeling and design from molecular
632 to genome scale with evo. March 2024.
- 633
634 Pascal Notin, Aaron W Kollasch, Daniel Ritter, Lood van Niekerk, Steffanie Paul, Hansen Spinner, Nathan
635 Rollins, Ada Shaw, Ruben Weitzman, Jonathan Frazer, Mafalda Dias, Dinko Franceschi, Rose Orenbuch,
636 Yarin Gal, and Debora S Marks. ProteinGym: Large-Scale benchmarks for protein design and fitness
637 prediction. *bioRxiv*, December 2023.
- 638
639 Sergey Ovchinnikov, Hetunandan Kamisetty, and David Baker. Robust and accurate prediction of residue-
640 residue interactions across protein interfaces using evolutionary information. *Elife*, 3:e02030, May 2014.
- 641
642 Donovan H Parks, Maria Chuvochina, David W Waite, Christian Rinke, Adam Skarshewski, Pierre-Alain
643 Chaumeil, and Philip Hugenholtz. A standardized bacterial taxonomy based on genome phylogeny sub-
644 stantially revises the tree of life. *Nat. Biotechnol.*, 36(10):996–1004, November 2018.
- 645
646 Donovan H Parks, Maria Chuvochina, Christian Rinke, Aaron J Mussig, Pierre-Alain Chaumeil, and Philip
647 Hugenholtz. GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically
648 consistent, rank normalized and complete genome-based taxonomy. *Nucleic Acids Res.*, 50(D1):D785–
649 D794, January 2022.
- 650
651 Christian Quast, Elmar Pruesse, Pelin Yilmaz, Jan Gerken, Timmy Schweer, Pablo Yarza, Jörg Peplies, and
652 Frank Oliver Glöckner. The SILVA ribosomal RNA gene database project: improved data processing and
653 web-based tools. *Nucleic Acids Res.*, 41(Database issue):D590–6, January 2013.
- 654
655 Alexander Rives, Joshua Meier, Tom Sercu, Siddharth Goyal, Zeming Lin, Jason Liu, Demi Guo, Myle
656 Ott, C Lawrence Zitnick, Jerry Ma, and Rob Fergus. Biological structure and function emerge from
657 scaling unsupervised learning to 250 million protein sequences. *Proceedings of the National Academy of*
Sciences, 118(15):e2016239118, April 2021.
- 658
659 Anna Rogers, Olga Kovaleva, and Anna Rumshisky. A primer in BERTology: What we know about how
660 BERT works. *Trans. Assoc. Comput. Linguist.*, 8:842–866, February 2020.

- 658 Andrew Rosenberg and Julia Hirschberg. V-Measure: A conditional Entropy-Based external cluster eval-
659 uation measure. In Jason Eisner (ed.), *Proceedings of the 2007 Joint Conference on Empirical Methods*
660 *in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*, pp.
661 410–420, Prague, Czech Republic, June 2007. Association for Computational Linguistics.
- 662 Dan Søndergaard, Christian N S Pedersen, and Chris Greening. HydDB: A web tool for hydrogenase
663 classification and analysis. *Sci. Rep.*, 6:34212, September 2016.
- 664 Martin Steinegger and Johannes Söding. MMseqs2 enables sensitive protein sequence searching for the
665 analysis of massive data sets. *Nat. Biotechnol.*, 35(11):1026–1028, November 2017.
- 666 Barbara R Terlouw, Kai Blin, Jorge C Navarro-Muñoz, Nicole E Avalon, Marc G Chevrette, Susan Egbert,
667 Sanghoon Lee, David Meijer, Michael J J Recchia, Zachary L Reitz, Jeffrey A van Santen, Nelly Selem-
668 Mojica, Thomas Tørring, Liana Zaroubi, Mohammad Alanjary, Gajender Aleti, César Aguilar, Suhad A A
669 Al-Salihi, Hannah E Augustijn, J Abraham Avelar-Rivas, Luis A Avitia-Domínguez, Francisco Barona-
670 Gómez, Jordan Bernaldo-Agüero, Vincent A Bielinski, Friederike Biermann, Thomas J Booth, Victor J
671 Carrion Bravo, Raquel Castelo-Branco, Fernanda O Chagas, Pablo Cruz-Morales, Chao Du, Katherine R
672 Duncan, Athina Gavriilidou, Damien Gayraud, Karina Gutiérrez-García, Kristina Haslinger, Eric J N
673 Helfrich, Justin J J van der Hooft, Afif P Jati, Edward Kalkreuter, Nikolaos Kalyvas, Kyo Bin Kang, Satria
674 Kautsar, Wonyong Kim, Aditya M Kunjapur, Yong-Xin Li, Geng-Min Lin, Catarina Loureiro, Joris J R
675 Louwen, Nico L L Louwen, George Lund, Jonathan Parra, Benjamin Philmus, Bitá Pourmohsenin, Lotte
676 J U Pronk, Adriana Rego, Devasahayam Arokia Balaya Rex, Serina Robinson, L Rodrigo Rosas-Becerra,
677 Eve T Roxborough, Michelle A Schorn, Darren J Scobie, Kumar Saurabh Singh, Nika Sokolova, Xiaoyu
678 Tang, Daniel Udvary, Aruna Vigneshwari, Kristiina Vind, Sophie P J M Vromans, Valentin Waschulin,
679 Sam E Williams, Jaclyn M Winter, Thomas E Witte, Huali Xie, Dong Yang, Jingwei Yu, Mitja Zdouc,
680 Zheng Zhong, Jérôme Collemare, Roger G Linington, Tilmann Weber, and Marnix H Medema. MIBiG
681 3.0: a community-driven effort to annotate experimentally validated biosynthetic gene clusters. *Nucleic*
682 *Acids Res.*, 51(D1):D603–D610, January 2023.
- 683 Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. BEIR: A het-
684 erogenous benchmark for zero-shot evaluation of information retrieval models. *arXiv*, 2021.
- 685 L Valeriani, Diego Doimo, F Cuturello, A Laio, A Ansuini, and A Cazzaniga. The geometry of hidden
686 representations of large transformer models. *Adv. Neural Inf. Process. Syst.*, abs/2302.00294, February
687 2023.
- 688 Michel van Kempen, Stephanie S Kim, Charlotte Tumescheit, Milot Mirdita, Jeongjae Lee, Cameron L M
689 Gilchrist, Johannes Söding, and Martin Steinegger. Fast and accurate protein structure search with fold-
690 seek. *Nat. Biotechnol.*, 42(2):243–246, February 2024.
- 691 Michael Wells, Minjae Kim, Denise M Akob, Partha Basu, and John F Stolz. Impact of the dimethyl
692 sulfoxide reductase superfamily on the evolution of biogeochemical cycles. *Microbiol Spectr.*, 11(2):
693 e0414522, March 2023.
- 694 Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yang Zhang, Chang Ma, Runcheng Liu, and Jian
695 Tang. PEER: A comprehensive and multi-task benchmark for protein sequence undERstanding. *Adv.*
696 *Neural Inf. Process. Syst.*, abs/2206.02096, June 2022.
- 697 Tianhao Yu, Haiyang Cui, Jianan Canal Li, Yunan Luo, Guangde Jiang, and Huimin Zhao. Enzyme function
698 prediction using contrastive learning. *Science*, 379(6639):1358–1363, March 2023.
- 699 Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta, Ramana Davuluri, and Han Liu. Dnabert-2: Efficient
700 foundation model and benchmark for multi-species genome. 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2306.15006)
701 [2306.15006](https://arxiv.org/abs/2306.15006).

705 APPENDIX A METHODS

706 **ModBC BiGene Mining** ModA and ModC sequence pairs were identified in Ovchinnikov et al.
 707 (2014) and downloaded from [https://gremlin.bakerlab.org/cplx.php?uni_a=2ONK_A&](https://gremlin.bakerlab.org/cplx.php?uni_a=2ONK_A&uni_b=2ONK_C)
 708 [uni_b=2ONK_C](https://gremlin.bakerlab.org/cplx.php?uni_a=2ONK_A&uni_b=2ONK_C). Original sequences queried using the UniProt IDs were used for the dataset. Sequences
 709 were dereplicated at 70% sequence identity using CD-HIT (Huang et al., 2010) and only pairs where both
 710 sequences were in the dereplicated set were included in the dataset.

711 **BacArch BiGene Dataset** RefSeq annotations were obtained for *E. coli* str. K-12 substr. MG1655
 712 (GCF_000005845.2) and *Sulfolobus acidocaldarius* (GCF_000012285.1). Orthologous genes were identi-
 713 fied as follows: 1) Genes with exactly matching annotations were identified first and added to the dataset;
 714 numerous genes with nearly identical, but not exactly matching, annotations, were also added to the dataset.
 715 2) Genes without highly similar annotations and with matching function indicated through other databases
 716 such as UniRef, were marked as orthologs. 3) Unannotated genes in *Sulfolobus* were identified as orthologs
 717 to *E. coli* sequences through a combination of genome context information, matching HMM domains, and
 718 high structural similarity identified through a Foldseek (van Kempen et al., 2024) search between the pre-
 719 dicted structure of *Sulfolobus* sequences against the structures available for *E. coli* MG1655; 4) For *Sul-*
 720 *folobus* genes with ambiguous RefSeq annotations, at least two such clues (matching UniRef annotations,
 721 genome context clues, matching HMM annotations, and Foldseek structural similarity) were required to as-
 722 sign an orthologous pair. 5) Genes where multiple homologs existed in both *E. coli* and *Sulfolobus* genomes
 723 were deliberately excluded.

724 **EC Classification Datasets** Sequences with an assigned EC number were downloaded from UniProtKB
 725 (Breuza et al., 2016) on May 16th 2024. Only "reviewed" sequences, non-fragments and sequences with
 726 a single EC designation were included. Sequences were first dereplicated at 70% sequence identity using
 727 CD-HIT and further clustered at 10% sequence identity (`--min-seq-id 0.1`) using `mmseqs cluster`
 728 (Steinegger & Söding, 2017) with coverage threshold of 30% (`-c 0.3`) and minimum alignment length of
 729 50 bp (`--min-aln-len 50`). Only EC classes with greater than five sequences after dereplication and
 730 clustering were kept. Five sequences were chosen randomly for each EC class, where one sequence was
 731 added to the test set and the remaining four were added to the train set.

732 **Convergent Enzymes Classification Dataset** Raw sequences and EC labels were downloaded from
 733 UniProtKB and dereplicated at 70% sequence identity as described above in section "EC Classification."
 734 Sequences were BLASTed against every other sequence with the same EC number designation in the derepli-
 735 cated set. Only one example per EC class with at least five examples in the same EC class without a signifi-
 736 cant BLASTP match (alignment length <10 and percent identity <0.1) were kept for testing. Five sequences
 737 in the corresponding EC class that have no significant BLASTP match to the test sequence were randomly
 738 chosen for training.

739 **MIBiG Classification Dataset** Sequences and labels (secondary metabolite classes) were downloaded
 740 from the MIBiG server version 3.1 (<https://mibig.secondarymetabolites.org/>). Secondary
 741 metabolite class "Other" was removed from the dataset. For the AA dataset, protein sequences were
 742 extracted from the MIBiG genbank files and embedded in chunks of maximum sequence length set by
 743 `--max_seq_len` (determined by the model, e.g., 1024 for ESM2) and subsequently mean-pooled across
 744 the example. For the NA dataset, DNA sequences were extracted from the MIBiG genbank files, embed-
 745 ded in chunks of sequence length set by `--max_seq_len` (e.g. 8,192 for `evo-1-8k-base`, 65,536 for
 746 `evo-1-131k-base` as sequence length 131,072 did not fit into a single 80GB GPU with batch size 1) and
 747 subsequently mean-pooled to yield a single embedding per example. Examples were split into train and test
 748 sets using 80/20 ratio random sampling with stratification on the first class label.

749 **MopB Clustering Dataset** Labeled MopB family sequences, displayed in the phylogenetic tree of Fig-
 750 ure 1 in Wells et al. (2023), were obtained from their provided Supplementary Materials ([https://](https://itol.embl.de/tree/249112161424681659917609)
 751 itol.embl.de/tree/249112161424681659917609). Wells et al. (2023) conducted one of

752 the most comprehensive and up-to-date classification of MopB family enzymes to date. Sequences were
 753 first dereplicated at 70% identity using CD-HIT. Functional groups with fewer than 60 representatives were
 754 excluded from the dataset, and functional groups with greater than 100 representatives were randomly down-
 755 sampled to only include 100 representatives. Selected sequences were aligned with FAMSA (Deorowicz
 756 et al., 2016). Alignments were trimmed with trimAL (Capella-Gutiérrez et al., 2009) v1.4.rev15 with the
 757 parameter `-gt 0.1` to remove columns consisting of $\geq 90\%$ gaps. Phylogenetic trees were estimated using
 758 iQ-TREE 2 (Minh et al., 2020) with the following parameters: `-bb 1000 -m GTR+G4+F`.

759 **E.coli RNA Clustering** RNA sequences in the *E. coli* str. K-12 substr. MG1655 genome (GenBank ID
 760 GCF_000005845.2) were identified by running the RFAM (Kalvari et al., 2021) family of models using the
 761 Infernal (Nawrocki, 2014) software suite. RNA groups with more than one identified representative included
 762 sRNAs, tRNAs, and rRNAs, and each sequence was classified using these three labels. In order to remove
 763 length bias in each RNA class (e.g. rRNAs are significantly longer than sRNAs), each sequence longer than
 764 100bp was replaced by a random subsequence of length 100bp.

765 **RpoB Phylogenies** RpoB sequences were obtained from the GTDB database (release 09-RS220). Bacte-
 766 rial RpoB sequences were identified using the TIGRFAM model TIGR02013 (rpoB_bac); Archaeal RpoB
 767 sequences were identified using the TIGRFAM model TIGR03670 (rpoB_arc) using methods described pre-
 768 viously (Parks et al., 2018). Sequences were then dereplicated at 70% identity using CD-HIT. Sequences
 769 from phyla with fewer than 10 representatives in the GTDB were excluded. For all other phyla, 10 rep-
 770 resentative sequences were chosen; where 10 or more classes were present in each phylum, one sequence
 771 each was chosen for each of 10 random classes within the phylum in order to diversify sampled sequences,
 772 otherwise the 10 representative sequences for that phylum were chosen randomly. Nucleotide coding se-
 773 quences for each chosen protein sequence were then obtained and used to construct separate phylogenies.
 774 Four phylogenies were constructed in total: Bacterial amino acid, Archaeal amino acid, Bacterial nucleotide,
 775 and Archaeal nucleotide. All alignments were performed using FAMSA. All alignments were trimmed us-
 776 ing trimAL with parameters described above. Amino acid phylogenies were estimated using iQ-TREE 2
 777 with the following parameters: `-bb 1000 -m LG+G4+F`. Nucleotide phylogenies were estimated using
 778 iQ-TREE 2 with the following parameters: `-bb 1000 -m GTR+G4+F`.

779 **FeFeHydrogenase Phylogeny** FeFe hydrogenase catalytic subunit sequences were obtained from HydDB
 780 (Søndergaard et al., 2016) and dereplicated at 70% ID using CD-HIT. The remaining sequences were then
 781 aligned using FAMSA. Alignments were trimmed using trimAL with parameters as described above. Amino
 782 acid phylogenies were then estimated using iQ-TREE 2 with the following parameters: `-bb 1000 -m`
 783 `LF+G4+F`.

784 **16S/18S rRNA phylogenies** 16S/18S sequences were obtained from SILVA release 138.2 and derepli-
 785 cated at 70% identity using CD-HIT. Sequences were then aligned with FAMSA and trimmed using trimAL
 786 with parameters as described above. Phylogenies were estimated using iQ-TREE with the following param-
 787 eters: `-m GTR+G4+F+I -bb 1000` with the addition of the +I model parameter to accommodate the
 788 presence of invariant sites in the alignment. The phylogeny in Fig. 1 was obtained by sampling one 16S or
 789 18S rRNA sequence from each phylum designated and constructed using the procedure described above.

790 **Operonic Pairs** For transcription units information and the corresponding protein sequences were ex-
 791 tracted from the BioCyc server (<https://biocyc.org/>) (Karp et al., 2019) for genomes *Escherichia*
 792 *coli* K-12 substr. MG165 *Vibrio cholerae* O1 biovar El Tor str. N16961, *Synechococcus elongatus* PCC
 793 7942. For a given consecutive gene pair, a label was assigned (1 or 0) depending on whether or not they are
 794 found in the same transcription unit.

795 **Retrieval** Protein sequences and protein name annotations were downloaded from UniProtKB on June 16
 796 2024. Only reviewed sequences and non fragments were kept for further processing. First, the sequences
 797 were partitioned into three domain (bacterial, archaeal or eukaryotic) sets using the UniProt taxonomic
 798 designation. Second, all proteins with "UPF" or "Uncharacterized protein" in the text labels were removed.

799 Third, the sequence were dereplicated at 50% sequence identity with CD-HIT with additional parameters `-c`
800 `0.5`, `-n 2`. Finally, overlapping text annotations between bacterial and archaeal, or bacterial and eukaryotic
801 sequence sets were identified, and only sequences that map to the overlapping text annotations were kept.
802 For the Arch_retrieval dataset, bacterial sequences were used as corpus with archaeal sequences as query.
803 For the Euk_retrieval dataset, bacterial sequences were used as corpus with eukaryotic sequences as query.
804 Relevance scores for each corpus-query sequence pair were calculated using fuzzy string matching (<https://github.com/seatgeek/thefuzz>): for fuzz ratio >90 between two text annotations relevance
805 score of 1 was assigned, otherwise, score of 0 was assigned.
806

807 A.1 MODEL INFERENCE

808 For all tasks except MIBiG classification task, sequences were truncated to the model's maximum sequence
809 length (predetermined by the model) using the flag `--max_seq_len`. For the MIBiG classification task,
810 sequences were chunked by the model's maximum sequence length as described above.
811

812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845

APPENDIX B DATASET STATISTICS

Overview of DGEB dataset statistics. For datasets that support both modalities (amino acids (AA) and nucleic acids (NA)), the values in parenthesis refer to the statistics for NA datasets.

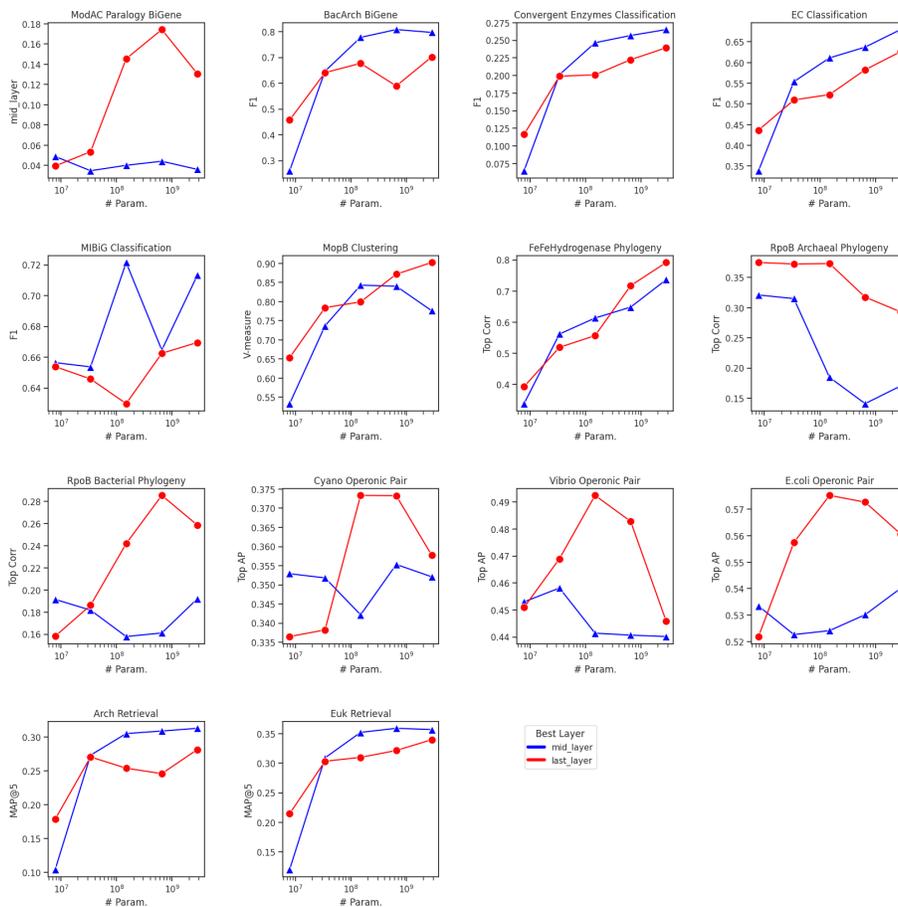
Dataset	Type	Categ.	# Phyla	# Label classes	# Train	Avg. train seq length	# Test	Avg. test seq length	Modalities
BacArch	BiGene Mining	IE	2	2	-	-	265	663	AA
ModBC	BiGene Mining	IE	36	2	-	-	1492	707	AA
FeFe Hydrogenase	EDS	IE	26	-	-	-	429	569	AA
RpoB Bac	EDS	IE	56	-	-	-	360 (360)	1305 (3927)	AA, NA
RpoB Arch	EDS	IE	13	-	-	-	170 (170)	831 (2491)	AA, NA
16S Bac	EDS	IE	31	-	-	-	545	1686	NA
16S Arch	EDS	IE	10	-	-	-	96	1423	NA
18S Euk	EDS	IE	20	-	-	-	751	2117	NA
Ecoli Operon	Pair Classification	IE	1	2	-	-	4315	310	AA
Vibrio Operon	Pair Classification	IE	1	2	-	-	2574	335	AA
Cyano Operon	Pair Classification	IE	1	2	-	-	2611	305	AA
EC	Classification	SE	38	128	512 (512)	541 (1901)	128 (128)	640 (1622)	AA, NA
Convergent Enzymes	Classification	SE	51	400	2000	415	400	433	AA
MIBIG	Classification	ME	15	6	29992 (1763)	647 (41178)	7213 (441)	638 (38206)	AA, NA
MopB	Clustering	SE	46	13	-	-	1300	817	AA
Ecoli RNA	Clustering	SE	1	3	-	-	161	83	NA
Arch	Retrieval	SE	52	-	9229	344	2343	332	AA
Euk	Retrieval	SE	44	-	3202	353	311	367	AA

APPENDIX C MODEL INFORMATION AND STATISTICS

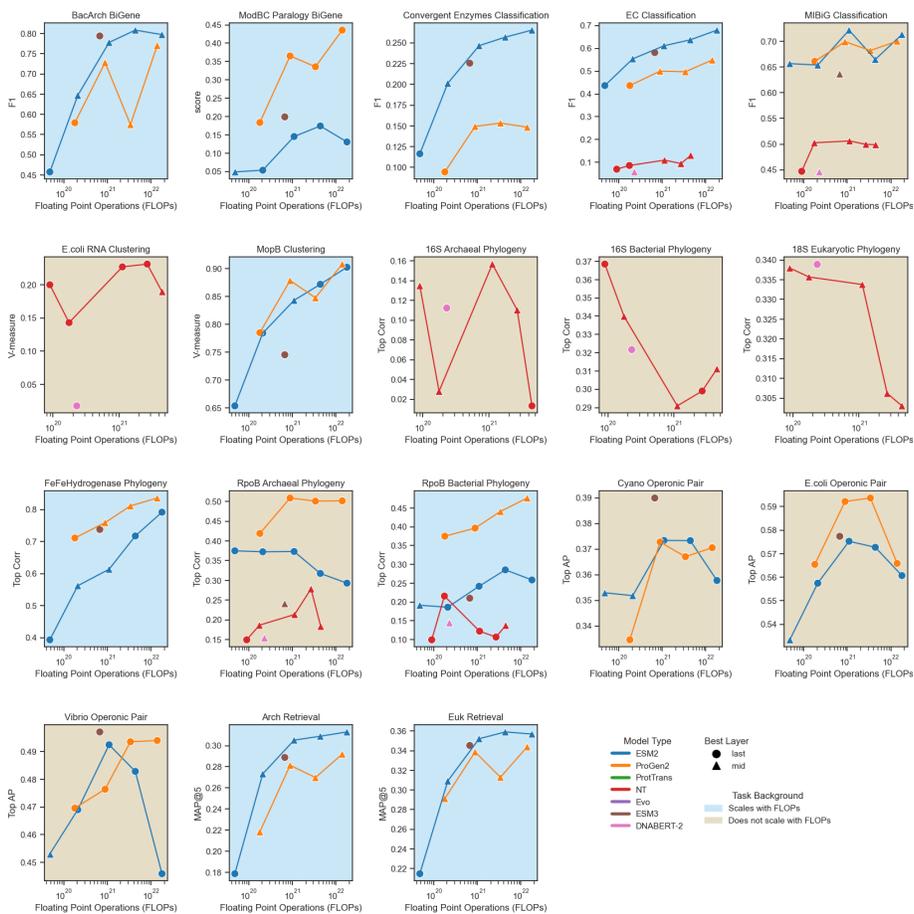
Models evaluated with DGEb are detailed with the number of parameters, number of hidden layers, and embedding dimensions. Pretraining FLOPs are estimated in (Chen et al., 2024) or from the model’s original papers when available. FLOP values with an asterisk are estimated using the formula $C = 6ND$ from (Kaplan et al., 2020), where C is the total pretraining flops, N is the model size, and D is the number of pretraining tokens.

Model type	Model Name	Modeling Objective	Training Data	Num Params	Num Layers	Emb. Dim.	Modality	Pretrain FLOPs
ESM2	esm2_t6_8M_UR50D	MLM	UniRef50/D	8M	6	320	AA	4.8E+19*
ESM2	esm2_t12_35M_UR50D	MLM	UniRef50/D	35M	12	480	AA	2.1E+20*
ESM2	esm2_t30_150M_UR50D	MLM	UniRef50/D	150M	30	640	AA	1.1E+21
ESM2	esm2_t33_650M_UR50D	MLM	UniRef50/D	650M	33	1280	AA	4.4E+21
ESM2	esm2_t36_3B_UR50D	MLM	UniRef50/D	3B	36	2560	AA	1.9E+22
ESM3	esm3_sm_open_v1	MLM	UniRef, MGnify; JGI (Hayes et al. 2024)	1.4B	48	1536	AA	6.72E+20
ProGen	progen2-small	CLM	UniProtKB	150M	12	1024	AA	1.8E+20
ProGen	progen2-medium	CLM	UniProtKB	765M	27	1536	AA	8.9E+20
ProGen	progen2-large	CLM	UniProtKB	2.7B	32	2560	AA	3.4E+21
ProGen	progen2-xlarge	CLM	UniProtKB	6.4B	32	4096	AA	1.4E+22
ProTrans	prot_t5_xl_uniref50	MLM	UniRef50	1.2B	24	1024	AA	-
ProTrans	prot_t5_xl_bfd	MLM	BFD (Steinegger and Söding 2018)	1.2B	24	1024	AA	1.7E+22
NT	nt-v2-50m-multi-species	MLM	Multispecies (NCBI) (Dalla-Torre et al. 2023)	55M	12	512	NA	9.0E+19*
NT	nt-v2-100m-multi-species	MLM	Multispecies (NCBI)	98M	22	512	NA	1.76E+20*
NT	nt-v2-250m-multi-species	MLM	Multispecies (NCBI)	235M	24	768	NA	1.13E+21*
NT	nt-v2-500m-multi-species	MLM	Multispecies (NCBI)	498M	29	1024	NA	2.69E+21*
NT	nt-2.5b-multi-species	MLM	Multispecies (NCBI)	2.5B	32	2560	NA	4.5E+21*
Evo	evo-1-8k-base	CLM	OpenGenome (Nguyen et al. 2024)	6.5B	32	4096	NA	-
Evo	evo-1-131k-base	CLM	OpenGenome	6.5B	32	4096	NA	2E+22
DNABERT	DNABERT2	MLM	Multispecies	117M	12	768	NA	2.3E+20

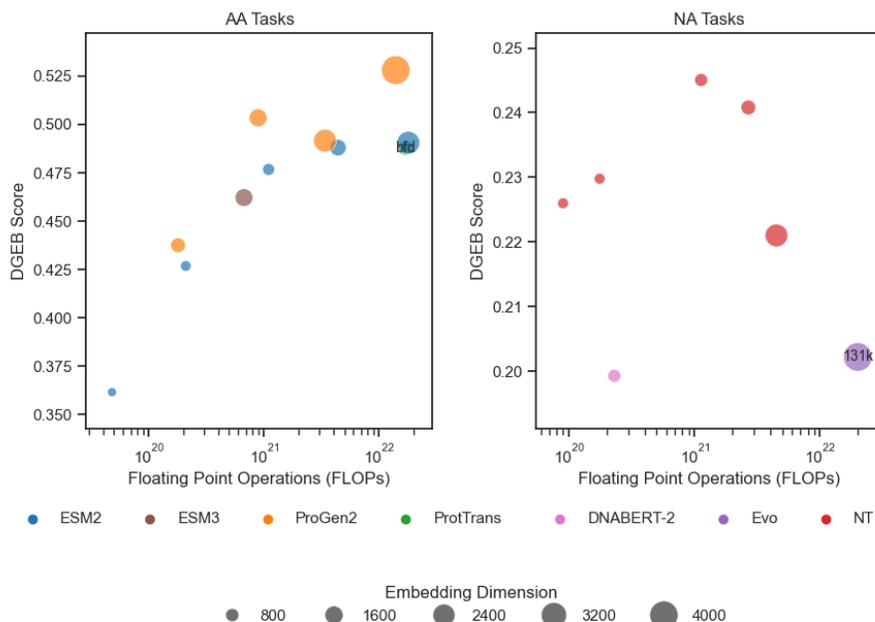
APPENDIX D COMPARISON OF MID LAYER AND LAST LAYER PERFORMANCE FOR ESM2 SERIES MODELS.



APPENDIX E PER-TASK PERFORMANCE SCALING WITH PRE-TRAINING FLOPS.



APPENDIX F AGGREGATED DGEGB SCORE RELATIVE TO PRE-TRAINING FLOPS.



APPENDIX G MODEL PERFORMANCE PER TASK

Task Type	Task	AA models											NA models									
		ESM2					ESM3	Progen				ProfTrans		Nucleotide Transformer					Evo		DNABERT	
		esm2 16.8M UR50D	esm2 112.35M UR50D	esm2 130.150M UR50D	esm2 133.650M UR50D	esm2 136.3B UR50D	esm3 3B UR50D	progen2- small	progen2- medium	progen2- large	progen2- xlarge	proL15.xl uniref50	proL15.xl bfd	NT v2-50m Multispecies	NT v2-100m Multispecies	NT v2-250m Multispecies	NT v2-500m Multispecies	NT 2.5b Multispecies	evo-1 8k-base	evo-1 131k-base	DNABERT2 117M	
BiGene Mining	ModBC BiGene BacArch BiGene	0.049 0.457	0.054 0.647	0.145 0.778	0.174 0.808	0.131 0.797	0.199 0.794	0.184 0.579	0.365 0.728	0.336 0.575	0.436 0.770	0.275 0.799	0.273 0.782	n/a n/a	n/a n/a	n/a n/a	n/a n/a	n/a n/a	n/a n/a	n/a n/a	n/a n/a	
Classification	EC Classification MIBIG Classification Convergent Enzymes Classification	0.437 0.656 0.116	0.554 0.654 0.201	0.611 0.722 0.246	0.637 0.665 0.257	0.680 0.713 0.265	0.581 0.636 0.225	0.437 0.661 0.095	0.500 0.699 0.149	0.497 0.682 0.153	0.549 0.700 0.148	0.629 0.692 0.243	0.565 0.665 0.227	0.070 0.447 n/a	0.086 0.503 n/a	0.110 0.506 n/a	0.095 0.500 n/a	0.131 0.499 n/a	0.012 0.426 n/a	0.016 0.446 n/a	0.086 0.446 n/a	
Clustering	MopB Clustering E. coli RNA Clustering	0.654 n/a	0.784 n/a	0.843 n/a	0.872 n/a	0.902 n/a	0.745 n/a	0.785 n/a	0.879 n/a	0.848 n/a	0.908 n/a	0.872 n/a	0.828 n/a	n/a 0.200	n/a 0.143	n/a 0.227	n/a 0.231	n/a 0.190	n/a 0.660	n/a 0.681	n/a 0.018	
EDS	FeFeHydrogenase Phylogeny	0.393	0.562	0.614	0.717	0.792	0.738	0.711	0.759	0.811	0.839	0.707	0.624	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	
	16S Bacterial Phylogeny	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	0.368	0.340	0.291	0.299	0.311	0.073	0.073	0.322	
	16S Archaeal Phylogeny	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	0.135	0.028	0.157	0.110	0.013	-0.005	0.019	0.112	
	18S Eukaryotic Phylogeny	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	n/a	0.338	0.336	0.334	0.306	0.303	0.223	0.161	0.339	
	RpoB Archaeal Phylogeny RpoB Bacterial Phylogeny	0.375 0.191	0.372 0.186	0.373 0.242	0.318 0.286	0.293 0.259	0.241 0.210	0.419 0.375	0.509 0.397	0.501 0.441	0.477 n/a	0.339 0.288	0.356 0.273	0.150 0.100	0.187 0.216	0.214 0.123	0.279 0.107	0.184 0.138	0.146 0.090	0.152 0.070	0.154 0.145	
Pair Classification	E. coli Operonic Pair Cyano Operonic Pair Vibrio Operonic Pair	0.533 0.353 0.453	0.557 0.352 0.469	0.575 0.373 0.492	0.573 0.373 0.483	0.561 0.358 0.446	0.577 0.390 0.497	0.565 0.335 0.470	0.592 0.373 0.476	0.594 0.367 0.494	0.566 0.371 0.409	0.618 0.407 0.543	0.626 0.541 0.541	n/a n/a n/a	n/a n/a n/a	n/a n/a n/a	n/a n/a n/a	n/a n/a n/a	n/a n/a n/a	n/a n/a n/a		
Retrieval	Euk Retrieval Arch Retrieval	0.215 0.179	0.309 0.273	0.352 0.305	0.359 0.313	0.357 0.289	0.345 0.218	0.291 0.218	0.339 0.281	0.313 0.270	0.344 0.292	0.359 0.311	0.355 0.306	n/a n/a	n/a n/a	n/a n/a	n/a n/a	n/a n/a	n/a n/a	n/a n/a		

1131 APPENDIX H ONE-HOT BASELINE DETAILS

1132 We introduce a one-hot vector representation of biological sequences as a baseline method to compare
1133 model performance. This baseline represents each residue or nucleic acid as a one-hot vector. The one-
1134 hot representation is mean-pooled across the sequence dimension, like all evaluated models (Section 4.1).
1135

1136 After mean-pooling, the one-hot baseline results in a representation equivalent to amino-acid composition for
1137 AA tasks and nucleic-acid composition for NA tasks. Previous work has shown that amino-acid composition
1138 is highly predictive of biological tasks such as transmembrane β -barrel protein identification (Garrow et al.,
1139 2005) and microbial growth conditions (Barnum et al., 2024). The baseline representation is evaluated in
1140 the same way as model embeddings, such as logistic regression for single-label classification tasks.
1141

1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177