1 **Ultra-Accurate Classification and Discovery of Functional Protein-Coding Genes**

2 **from Microbiomes Using FunGeneTyper: An Expandable Deep Learning-Based**

3 **Framework**

4 Guoqing Zhang[#,1,2,3,4], Hui Wang[#,1,2,3], Zhiguo Zhang[1,2,3], Lu Zhang[1,2,3], Guibing Guo[5],

5 Jian Yang[6,7], Fajie Yuan[1,2,3,4]*, Feng Ju[1,2,3,4,6,7]*

6 [1] Research Center for Industries of the Future, Westlake University, Hangzhou,

7 Zhejiang 310030, China

8 [2] Center of Synthetic Biology and Integrated Bioengineering, Westlake University,

9 Hangzhou, Zhejiang 310030, China

10 [3] Key Laboratory of Coastal Environment and Resources of Zhejiang Province,

11 School of Engineering, Westlake University, Hangzhou, Zhejiang, 310030, China

12 [4] Institute of Advanced Technology, Westlake Institute for Advanced Study, 18

13 Shilongshan Road, Hangzhou, Zhejiang 310024, China

14 [5] Software College, Northeastern University, Shenyang, China

15 [6] School of Life Sciences, Westlake University, Hangzhou, Zhejiang 310030, China.

16 [7] Westlake Laboratory of Life Sciences and Biomedicine, Hangzhou, Zhejiang

17 310024, China

18 [#] The authors contributed equally to the work

19 * Correspondence to

20 Dr. Feng Ju, Email: jufeng@westlake.edu.cn

21 Dr. Fajie Yuan, Email: yuanfajie@westlake.edu.cn

22 Address: Westlake University, 600 Dunyu Road, Hangzhou 310030, China

23 Tel.: 571-87963205 (lab), 571-87380995 (office), Fax: 0571-85271986

24

25 **Abstract**

26 High-throughput DNA sequencing technologies open the gate to tremendous

27 (meta)genomic data from yet-to-be-explored microbial dark matter. However,

28 accurately assigning protein functions to new gene sequences remains challenging. To

29 this end, we developed FunGeneTyper, an expandable deep learning-based framework

30 with models, structured databases and tools for ultra-accurate (>0.99) and fine-grained

31 classification and discovery of antibiotic resistance genes (ARGs) and virulence factor

32 or toxin genes. Specifically, this new framework achieves superior performance in

33 discovering new ARGs from human gut (accuracy: 0.8512; and F1-score: 0.6948),

34 wastewater (0.7273; 0.6072), and soil (0.8269; 0.5445) samples, beating the state-of-

35 the-art bioinformatics tools and protein sequence-based (F1-score: 0.0556-0.5065)

36 and domain-based (F1-score: 0.2630-0.5224) alignment approaches. We empowered

37 the generalized application of the framework by implementing a lightweight, privacy-

38 preserving and plug-and-play neural network module shareable among global

39 developers and users. The FunGeneTyper[*] is released to promote the monitoring of

40 key functional genes and discovery of precious enzymatic resources from diverse

41 microbiomes.

42

45 * The codes and database resources are available at: https://github.com/emblab-

46 westlake/FunGeneTyper.

47 **Main**

48 High-throughput DNA sequencing and metagenomics have generated extensive

49 protein-coding gene (PCG) sequences from diverse environmental and human

50 microbiomes[1-3]. Accurate classification of genes into related protein functions is the

51 key to effective gene discovery. However, these datasets pose significant

52 computational challenges in metagenomic studies. Sequence alignment (SA),

53 implemented using NCBI's BLAST[4], usearch[5], and Diamond[6], is commonly used for

54 functional annotation of PCGs[7]. To minimize false-positives, SA-based methods are

55 routinely conducted with strict user-defined cutoffs or thresholds (alignment identity,

56 coverage, and bit scores) to retain high-confidence best hits for each query sequence

57 from validated databases. This practice is widely implemented in the development of

58 tools for categorizing genes, including antibiotic resistance genes (ARGs)[8,9] and

59 virulence factor genes (VFGs)[10]. SA-based approaches effectively predict functions

60 between genes that share high homology (>80% identity [8,9]), but exclude distantly

61 homogeneous genes that fall below arbitrarily-defined and one-size-fits-all cutoffs

62 that may represent the majority of targeted functional genes in environmental samples

63 (e.g. core ARGs in activated sludge[11] and soil[12]). Therefore, these SA approaches with

64 stringent bioinformatic cutoffs unavoidably generate numerous false-negative results

65 and heavily underestimate true novelties and diversity of functional genes in largely

66 uncultured bacteria, thus biasing research outcomes or conclusions. Therefore, it is

67 crucial to develop intelligent and accurate classification paradigm and bioinformatic

68    tools to overcome limitations of existing SA-based classification approaches.

69    Importantly, this endeavor will accelerate discovery of new genes in future

70    metagenomic-based environmental and human microbiome studies[13,14].

71       Hidden Markov models (HMM) with manual-crafted sequence alignments and

72    scoring functions are powerful tools for protein domain-based functional gene

73    annotation for detecting remote gene homologues with low sequence identity ($< 30\%$)

74    to known proteins[15,16]. However, these methods rely on token (amino acid) matching,

75    which fail to detect high-level semantic representation similarity or structure-level

76    representation similarity, leading to false-positives[17], and thus cannot distinguish

77    functions of proteins in the same family[18]. In contrast, deep learning (DL) methods

78    excel at learning rich and high-level semantic representations when sufficient training

79    data are available, and are effective at identifying proteins with structural and

80    functional similarities[19-22]. Specifically, ground-breaking big language models initially

81    developed for natural language processing tasks have been successfully applied to

82    protein function prediction tasks[23,24], often termed protein language models (PLMs).

83    The high-level semantic representations learned from PLMs establish valid

84    connections between sequences and function[25,26]. Notwithstanding the power of PLMs,

85    gene classification tasks, particularly identifying fine-grained protein function

86    subclasses, pose challenges for data-hungry deep learning paradigms because of

87    limited supervised training dataset for certain genes. Additionally, it remains unclear

88    whether advanced PLMs perform better than state-of-the-art metagenomic

89    bioinformatics tools at microbiome gene classification and discovery.

90    Here, we propose FunGeneTyper, a PLM-based deep-learning framework for

91    accurate and expandable prediction of PCG function. FunGeneTyper implements a

92    two-stage pipeline that separately handles the assignment of the main types and

93    subtypes of PCG functional classes, reducing issues associated with insufficient

94    training data during subtype-level predictions. To improve conciseness, it first

95    performs standard classification of genes of the main types and then performs fine-

96    grained retrieval by comparing similarities between learned protein subtype

97    representations. FunGeneTyper models classify ARGs with ultra-high accuracy (>0.99)

98    and outperforms the state-of-the-art SA and HMM-based methods and tools.

99    Furthermore, we also demonstrate the generalized application of FunGeneTyper

100   models in ultra-high classification of VFGs and introduce the adapter module, a

101   lightweight neural network that can be inserted into the current backbone architecture

102   to realize parameter-efficient training. The adapter-tuning-based FunGeneTyper

103   models are expandable to the classification of various categories of genes and enables

104   sharing of both task-agnostic and task-specific parameters without accessing the

105   private training dataset. Thus, FunGeneTyper offers a unified and innovative way of

106   integrating the global efforts of the microbiome and bioinformatics communities,

107   endowing the FunGeneTyper framework with the ability to conduct unlimited

108   prediction of functional gene categories beyond the ARGs and VFGs demonstrated

109   here, which is key to accelerating the global discovery of new and precious genetic

110    and enzymatic resources from microbiomes.

111    **Results**

112    **FunGeneTyper framework, structured database, and deep learning models**

113    FunGeneTyper is the first unified framework that utilizes DL models and structured

114    functional gene datasets (SFGD) to develop new DL-based classifiers for any gene

115    category via transfer learning. This novel framework achieves highly accurate PCG

116    classification from metagenomic studies and extends the models to efficiently predict

117    broad categories of gene functions from large varieties of microbiomes with

118    corresponding customizable SFGD.

119    *Structured functional gene datasets*

120    We deployed a transferable strategy to collect high-quality reference gene sequences

121    to meet FunGeneTyper's training requirements with high reliability (Fig. 1a).

122    Experimentally-confirmed reference sequences of target genes from literature and/or

123    expert-curated databases were used as the core dataset, and highly homologous

124    protein sequences (at least 80% identity and 80% coverage) were extracted from

125    Uniref100 database and used as the expanded functional genes dataset. A non-target

126    sequence dataset was selected from Swiss-Prot database (version: June 2021) by

127    excluding perfect matches to the target genes, and used as the negative training set so

128    that FunGeneTyper could learn sufficient features of non-target genes. Core and

129    expanded functional gene datasets and the non-target dataset were integrated to form

130    the SFGD, which was organized hierarchically into a secondary structure based on

131    gene ontology. The SFGD was divided into training, validation, and testing sets (ratio

132    6:2:2) and used to train the following two DL models.

133    *Deep learning models*

134    The framework has a top-down protein function prediction workflow featuring two

135    DL models (Fig. 1b), FunTrans and FunRep, which progressively classify protein

136    sequences from the upper (type) to lower (subtype) functional levels. FunGeneTyper

137    was pre-trained on ESM-1b, a large-scale pre-trained protein sequence model based

138    on the transformer architecture released by Facebook[22]. ESM-1b is composed of the

139    33-layer transformer architecture consisting of 650 million parameters trained on

140    Uniref50. It has the superior capacity to infer fundamental structural and functional

141    characteristics of proteins from sequences that can significantly increase the

142    performance metrics for sequence-function tasks. Despite sharing the 33-layer

143    transformer architecture, FunTrans and FunRep were independently constructed,

144    trained, and optimized, to complete two-level functional classification tasks that

145    successively assigned a PCG to its best matches of functional type and subtype in a

146    structured database, respectively.

147    FunTrans distinguishes protein sequences and classify them into specific functional

148    types equivalent to gene families with the same or similar functions. It is inspired by

149    the fact that proteins with similar structures and functions clustered closer together in

150    the embedding space. The main FunTrans structure is a 33-layer transformer that

151    implements initial classification of input data (Fig. 1c). Adapter modules are inserted

152    into each transformer block as trainable parameters. The adapter enables efficient

153    fine-tuning of parameters for different gene classification tasks. High parameter

154    sharing is achieved under the premise that the parameters of the original network

155    remain unchanged. Adapter modules enable flexible and parameter-efficient transfer

156    learning and prevent overfitting[27,28]. FunTrans adds a nonlinear classification layer at

157    the end of the sequence semantic representation for functional classification.

158      FunRep has a structure similar to that of FunTrans and can be used to embed

159    representation retrieval for further subtype classification of protein function (Fig. 1d).

160    It uses embedding representation retrieval to accurately predict functional subtypes of

161    the FunTrans output results for classification. FunRep also adds an adapter layer to

162    increase robustness and insight into a broader range of gene classification.

**FunGeneTyper classification performance and learning ability**

164    The spread of antibiotic resistance has raised public health concerns globally[29].

165    Reliable ARG model classification is important for surveillance and control of the

166    spread of antibiotic resistance, and achieving sufficient model sensitivity to remote

167    homologues is key to discovering new ARGs. Therefore, we first classified ARGs and

168    demonstrated the ability of the FunGeneTyper framework to achieve this goal. Before

169    building the ARGs classification models, we constructed a hierarchical structured

170    ARG database (SARD) based on antibiotic resistance ontology of the comprehensive

171    antibiotic resistance database (CARD)[7]. Based on CARD's ontological rules, ARGs

172 were assigned to class and group hierarchies based on the types of drugs to which

173 they confer resistance, and the subtypes of genes with the same resistance function,

174 respectively (Dataset S1). To test and improve the sensitivity of the model, we used

175 different identity thresholds to collect four non-target sequence sets from Swiss-Prot

176 database —excluding ARGs — as negative training datasets, for model training

177 (Supplementary Figure 1, see Methods). The addition of a negative training set allows

178 the model to learn features of non-targeted genes, which gives the model the ability to

179 directly classify targeted (e.g., ARGs) and non-targeted genes (e.g., non-ARGs) from

180 new datasets to be tested. We evaluated the impact of four identity thresholds of the

181 negative datasets on the learning features of the model. The results of five-fold cross-

182 validation revealed that the model with 0% identity as the threshold for recruiting

183 non-target sequences had the best performance metrics, including accuracy, recall,

184 precision, and F1-score (Fig. 2a). Under these optimized conditions, the positive

185 SARD set contained 61874 ARG sequences, including 2972 experimentally-

186 confirmed core sequences inherited from the CARD and 58902 homology-predicted

187 (>80% identity and >80% coverage) expanded ARG sequences from Uniref100. All

188 ARG reference sequences were hierarchically assigned to 19 classes and 2972 groups

189 (Dataset S2 and Supplementary Figure 2).

190 To demonstrate the powerful utility of FunGeneTyper, we used the reference

191 protein sequences in SARD to train two transformer models (FunTrans and FunRep)

192 and developed them as a new deep-learning ARGTyper deep-learning classifier. We

193    used the trained ARGTyper to classify the testing set to validate the performance of

194    the ARGTyper. The overall ARGTyper performance metrics prove that FunGeneTyper

195    provides an excellent and robust framework for gene classification. Specifically, the

196    optimal FunTrans model at the ARG class level reached an accuracy of 0.9979, a

197    precision of 0.9830, a recall rate of 0.9683, and an F1 score of 0.9756 (Fig. 2b).

198    Moreover, the prediction precision and recall of all 17 ARG classes exceeded 0.96

199    (Fig. 2c), apart from fusidic acid and triclosan, which showed lower precision and

200    recall because they have only 21 and 53 reference sequences, respectively, in SARD

201    (Dataset S3). More training data helps the model learn more features. Nonetheless, the

202    power of FunTrans to classify these temporarily less-represented classes of ARGs will

203    improve as more functionally-verified sequences will be available for model training.

204        The vector space generated by FunGeneTyper was semantically rich and encoded

205    structural, evolutionary, and functional information. To explain what our model

206    intuitively learns, we obtained representations of all classes of ARG and non-ARG

207    sequences in the training set. We used uniform manifold approximation and projection

208    (UMAP) to reduce data dimensions in each layer to two. Visualizations performed in

209    the four essential representative layers (1st, 5th, 32nd, and 33$^{rd}$) revealed the learning

210    process of the model (Fig. 2d). All ARG sequences were highly entangled at the first

211    level of encoding input. However, they became increasingly separated as the

212    transformer model got deeper. Each type of ARG undergoes a process from dispersion

213    to aggregation. This finding verified that FunTrans can efficiently learn the

214    representation features of sequences from raw input data with high entanglement.

215    Prediction multiclass confusion matrix was used to represent the effect of

216    FunTrans on the learning features of each ARG class. The results indicated that the

217    FunTrans model was excellent at predicting all ARG classes (Fig. 2e). We continued

218    to locate significant classification errors in the ARG classes using error detection

219    counts (Fig. 2f). Prediction error was concentrated in the multidrug class. Specifically,

220    33 non-ARG sequences were mispredicted as multidrug resistance, whereas 39

221    multidrug resistance protein sequences were mispredicted as non-ARG sequences.

222    The poor prediction performance of these proteins is mainly due to their high

223    structural differences and diverse biological functions that include roles other than

224    multidrug resistance[30], making it challenging for a DL model to effectively learn

225    sufficient discriminative features in the absence of sufficient training data. Multidrug

226    efflux pumps[30] export antibiotics and other diverse extraneous substrates, including

227    organic solvents, toxic heavy metals, and antimicrobials, and also fulfill other key

228    biological functions such as biofilm formation, quorum sensing and survival and

229    pathogenicity of bacteria[30]. Therefore, multidrug resistance proteins or efflux pumps

230    were not seriously considered as ARGs[17,31] and we recommend excluding their

231    sequences from ARG analysis unless they can be reliably or unambiguously assigned

232    to resistance functions of certain antibiotic classes.

233    Once the FunTrans model was shown to be robust and accurate in identifying

234    ARGs and classifying them into 19 classes, we trained FunRep, which conducted a

235   more detailed lower-level classification of ARGs into 2972 groups (Dataset S3).

236   FunRep achieved an overall prediction accuracy of 0.9023 for all ARG groups

237   (Dataset S4). We used UMAP to visualize FunRep model's learning process. UMAP

238   was used to visualize the characteristics of the final layer of all classes except the

239   Fusidic acid class (21 sequences, Dataset S3). UMAP showed that FunRep can cluster

240   the features of each group in the main ARG classes, including beta-lactams (5909

241   sequences), Macrolides-Lincosamides-Streptogramines (MLS, 2317 sequences),

242   aminoglycosides (3483 sequences), and glycopeptides (2037 sequences)

243   (Supplementary Figure 3).

244   In summary, we demonstrated the application of the FunGeneTyper framework to

245   develop ARGTyper as the first transformer-based ARG classifier trained from a

246   customized structured ARG database (SARD). The performance metrics of the testing

247   set show that FunTrans and FunRep can achieve highly accurate (accuracy=0.998)

248   and robust (F1-score=0.976) identification of all known types (classes) and subtypes

249   (groups) of ARGs in the authoritative CARD. Both the accuracy and robustness of

250   FunGeneTyper models outperform previously published results from DeepARG

251   (accuracy>0.97, F1-score>0.93 [9]) and HMD-ARG (accuracy=0.935, F1-score=0.893

252   [32]) on their own testing sets of ARGs.

253   **Model performance in the discovery of new genes**

254   The 'twilight zone' of protein sequence alignment is a complex, long-standing

255   problem plaguing protein function prediction[33,34], limiting the discovery of PCGs in

256    the largely uncultured microbial world. In contrast to classic SA-based tools, DL-

257    based models (FunRep and FunTrans) of the FunGeneTyper framework are designed

258    with unique features and intrinsic advantages for predicting remote homologues of

259    protein sequences with guaranteed accuracy and robustness, as previously

260    demonstrated for ARG classification.

261        To compare FunGeneTyper's ability to identify new PCGs with those of existing

262    methodologies, we evaluated its ability of its DL-based models to discover remote

263    homologues by predicting experimentally-confirmed protein sequences of new ARGs

264    discovered from three representative habitats: human gut (n = 168)[35], wastewater

265    treatment plants (n = 77)[11], and soil (n = 52)[36-39]. We computed the predictive

266    performance of FunGeneTyper classifier for ARGs (ARGTyper) and compared it with

267    that of three state-of-the-art tools: DL-based tools (HMD-ARG[32] and DeepARG[9]),

268    SA-based tools (RGI[7]), and HMM-based tools (Resfams[18]) (Table 1). FunGeneTyper

269    had higher accuracy, precision, recall, and F1-score for predicting new ARGs

270    compared with HMD-ARG[32] and DeepARG[9]. The significant improvement was

271    primarily attributed to our implementation of the protein semantic models (i.e.,

272    FunTrans and FunRep) in FunGeneTyper, which can learn more hidden features of

273    protein sequences, especially the context information[19,21], compared with the

274    traditional one-hot encoding algorithm and the convolutional neural network used by

275    HMD-ARG[32] and the multilayer perceptron used by DeepARG[9]. Moreover, the

276    overall classification performance of FunGeneTyper, as benchmarked by the F1-score

277     (0.5445 to 0.6948), was much higher than that of the classic SA-based methods

278     (0.0556 to 0.6598) and HMM-based methods (0.2630 to 0.5224) (Table 1). Although

279     RGI also achieved high accuracy (0.8830) in human intestinal data, its precision

280     (0.4545), recall (0.3968), and F1-score (0.4195) were much lower than those of the

281     FunTrans model (0.7500, 0.6642, and 0.6948, respectively) because many of the new

282     ARG sequences tested here fell below the commonly applied stringent identity cutoffs

283     (> 95% RGI). It is expected that when a strict one-size-fits-all filter cutoff is applied

284     to the alignment results, many false-negatives would result, limiting the discovery of

285     ARGs that show a more remote homology to database sequences. The superior

286     performance of the FunGeneTyper classifier over existing tools in identifying new

287     ARGs was further evident when comparative tests were performed using wastewater

288     treatment plant (WWTP) or soil samples compared with human gut samples (Table 1).

289     This indicates that FunGeneTyper has a greater capacity to predict functional genes in

290     complex environmental samples. To further resolve the superior predictive

291     performance of FunGeneTyper for remote homologues of functional genes over

292     existing tools, we divided the ARGs data into lower homology (≤50% identity) and

293     higher homology (≥50% identity) datasets (Supplementary Figure 4). FunGeneTyper

294     not only consistently achieved better classification performance of higher homology

295     ARGs in all three sample groups (WWTP, soil, and human gut), it also showed

296     outstanding performance at accurately and sensitively predicting the function of

297     remote homologous sequences (Dataset S5).

298    Taken together, our results exemplify the discovery and classification of novel

299    ARGs, especially among relatively remote homologues (<50% identity), and

300    demonstrate that FunGeneTyper is best at predicting new ARG protein sequences,

301    exhibiting unprecedented capacity to identify new genes with high accuracy,

302    sensitivity, and robustness.

303    **Evaluating the generalizability of FunGeneTyper**

304    To demonstrate the generalizability of FunTrans and FunRep in classifying other gene

305    categories, we trained the models using a calibrated and professionally expanded

306    bacterial virulence factor database, VFNet[40] and utilized them to develop a new

307    transformer-based classifier of virulence factor gene (VFG), named VFGTyper.

308    Before training the model, we built a two-level expert-curated structured database

309    based on the virulence ontology and reference sequences in the VFNet database.

310    Semantic and categorically ambiguous data were cleaned (Methods). The final

311    structured virulence factor database (SVFD) consisted of 160484 VFG sequences

312    distributed into 2837 classes in 45 families (Dataset S6).

313    We followed a strategy similar to that mentioned above to train the model,

314    collecting a non-target dataset with 551,783 sequences (excluding VFGs) from Swiss-

315    Prot, as the negative dataset (see Methods). With the merit of the proposed adapter

316    module, we only need to re-train a new adapter when building a VFGTyper. The

317    design of adapter allows us to train only a new classifier and an adapter when

318    predicting new functions. All the parameters in the backbone network can be reused.

319    Therefore, VFGTyper can be regarded as a new task branch in the FunGeneTyper

320    framework, where only the adapter and classifier differ. We verified the VFGTyper

321    using the testing set to provide evidence of its generalizability in the functional

322    genotyping process. VFGTyper achieved an accuracy of 0.9907 (Fig. 3a) in the family

323    level prediction task. The obfuscation matrix results also showed that FunTrans

324    achieved excellent classification performance for each VFG at the family level (Fig.

325    3b, Supplementary Figure 5). In addition, FunRep was 0.9499 accurate at predicting

326    different VFG classes in the second-stage prediction.

327        In conclusion, we demonstrated that FunGeneTyper can be successfully

328    generalized to develop VFGTyper as the first transformer-based VFG classifier of its

329    kind and applied the new classifier to achieve ultra-accurate classification of VFGs by

330    adding new adapters. To vividly show that FunGeneTyper can learn sufficient

331    discriminative features from different groups of functional gene datasets, we

332    visualized the learning process of FunTrans and FunRep models for VFG sequences.

333    Consistent with the learning process for ARGs (Fig. 2d), both models also achieved

334    effective feature clustering and classification of VFGs at both the family (Fig. 3c) and

335    class (Supplementary Figure 6) levels. Besides classification performance, we also

336    proved VFGTyper's full capability in the discovery of an experimentally-confirmed

337    novel VFG (NCBI accession no.: WP_034687872.1) of a toxin family in

338    *Chryseobacterium piperi* with sequence similarity to botulinum neurotoxins (BoNTs)

339    through re-analysis of published genome [41]. Specifically, of the 8 putative toxin genes

340    of *C. piperi* showing no significant (n=6) or only limited (n=2) sequence homology

341    (i.e., global identity < 10%) to known reference VFGs, 7 were effectively identified as

342    VFGs by FunGeneTyper and 4 were further classified as BoNTs (Dataset S7).

343    Compared with a conventional sequence alignment (SA)-based approach which failed

344    to predict 6 VFGs, the deep learning models of FunGeneTyper showed much greater

345    capacity for the discovery of remote homologues of known toxin genes. Therefore,

346    FunGeneTyper represents an expandable deep learning-based framework for ultra-

347    accurate classification and discovery of functional genes, as demonstrated here for

348    ARGs and VFGs.

349    **Privacy-preserving global sharing of plug-and-play adapters for functional gene**

350    **discovery**

351    To demonstrate the parameter efficiency of FunGeneTyper's adapter modules, all 650

352    million parameters of the pre-trained model are fine-tuned as a benchmark test which

353    achieved excellent prediction accuracy in ARGs class (0.9988) and VFGs family

354    (0.9930). Comparatively, with only fine-tuning of about 21 million parameters (3% of

355    all parameters) of the Adapter layer, we demonstrated that FunGeneTyper achieved

356    near-identical excellent performance of 0.9979 for ARGs class and 0.9907 for VFGs

357    family, proving that parameter-efficient lightweight plug-and-play adapter modules of

358    FunGeneTyper can be easily shared without little loss of prediction accuracy.

359        Benefiting from the parameter-efficient property, FunGeneTyper has two novel

360    merits. First, FunGeneTyper enables effective effort-sharing by the entire community

361     (Fig. 4). Specifically, a researcher who has trained our FunGeneTyper model for

362     classification or discovery of functional genes (other than ARGs and VFGs

363     demonstrated here) can submit their adapters (along with a classification layer) to the

364     adapter hub. Once the adapter has been submitted, the module can be downloaded and

365     easily inserted into the FunGeneTyper model for direct downstream user application.

366     Second, the adapter design helps solve data privacy issues. Where researchers have

367     not publicly released their own datasets, they can train FunGeneTyper with their

368     private datasets, submit only the adapter module (again along with a classification

369     layer), and provide functional descriptions of their FunGeneTyper. Thus, the private

370     datasets are protected, while the uploaded adapter models can be used without model

371     training. As the number of researchers getting involved in the development of

372     FunGeneTyper increases, the model may become a universal toolkit that can be used

373     for predicting functional genes simply by looking up related functional modules. We

374     believe that with the elegant adapter module, FunGeneTyper will facilitate adapter

375     sharing and model integration globally.

376     **Discussion**

377     Metagenomics has provided an opportunity for identifying microbiome diversity and

378     novel functionalities. However, the speed at which high-throughput DNA sequencing

379     technologies unravel the vast genetic novelties of uncultured microbes in nature

380     outpaces our capacity to understand their function. Previous approaches for functional

381     classification of PCGs were based on sequence alignment using tools such as BLAST[4],

382     usearch[5], and Diamond[6] or conserved motifs and domains using Hidden Markov

383     Models. Selection of uniform cutoffs and thresholds usually limit the accuracy and/or

384     sensitivity of these methods for functional gene prediction. Protein semantic

385     algorithms based on NLP methods have been developed [20,24]. However, these

386     algorithms are not optimized for classifying different categories of microbial genes,

387     and a unified thinking paradigm is required to meet the needs for accelerated

388     discovery of new genes.

389         Our study provides an expandable deep learning-based framework for efficient

390     and robust gene function prediction, which represents an emerging methodological

391     paradigm for global developers and users to tackle unprecedented challenges and

392     meet the above-mentioned urgent needs in the classification and discovery of any

393     group of functional PCGs. We propose an end-to-end FunGeneTyper framework for

394     the classification prediction of gene functions. We exemplify the framework by

395     developing two transformer-based classifiers, ARGTyper and VFGTyper, based on

396     deep learning models coupled with expert-curated structured databases (SARD and

397     SVFD) to realize robust functional classification of bacterial ARGs and VFGs, which

398     are two categories of genes key to WHO's one health approach for human, animal,

399     and environmental health protection[42].

400         A series of experimental validations, including five-fold cross-validation, testing

401     set    validation,    and    experimentally-confirmed    protein    sequence    validation,

402    demonstrate the effectiveness and robustness of FunGeneTyper. Using ARG as an

403    example, FunGeneTyper models are more effective than SA-based and DL-based

404    models in predicting protein sequences of new ARGs from the human gut, WWTP,

405    and soil microbiomes with relatively low homology (< 50% similarity) to known

406    ARGs. This shows that ARGTyper has an unmatched advantage in discovering ARGs,

407    primarily because of the powerful learning ability of protein semantic models. Since

408    experimentally-confirmed sequences of the major categories, types, and subtypes of

409    genes are not sufficient, expanding the database based on sequence homology is

410    common and necessary to obtain sufficient training sequence data. UMAP analysis

411    showed that the expanded sequences represent reliable datasets and support our model

412    to better learn discriminative protein semantic features to achieve satisfactory

413    performance in identifying functional genes, including ARGs and VFGs.

414        Accurately classifying target genes from the huge interference of non-target gene

415    data is a problem. Therefore, we purposefully introduced non-functional genetic

416    datasets as part of the training set. Although this operation increases training

417    complexity, it enables our model to accurately classify target genes from noisy data

418    when used to analyze large-scale metagenomic sequence datasets from environmental

419    or animal microbiome samples. Some machine learning methods rely on sequence

420    alignment tools to create a similarity score matrix of potential gene sequences and

421    databases[9,43]. Such practices will inevitably be affected (and limited) by the selection

422    of arbitrary thresholds for the results. The FunGeneTyper framework proposed here

423   can accurately annotate genes via classification through discriminative features

424   learned from multiple sequences. The limited number of training sequences may

425   prevent the models from learning sufficient features. This transient issue would,

426   however, be easily solved once more experimentally-confirmed reference protein

427   sequences of target genes are available for model retraining and refinement.

428   Meanwhile, the robustness of deep learning to noise labels[44] can also help our

429   framework models and classifiers outperform existing ones in discovering new genes.

430         In particular, once large amounts of (meta)genomic data are freely available, a

431   uniform and convenient understanding of the relationship between microbial gene

432   sequences and protein function becomes a perennial challenge that can be tackled to

433   create opportunities for gene discovery. There are other gene categories, apart from

434   ARGs and VFGs, including those associated with microbially-driven global

435   biogeochemical cycling (carbon, nitrogen, phosphorus, and sulfur) or microbial

436   biodegradation (bioremediation and bio-restoration) and biosynthesis (biomedicine

437   and bioresources) (Fig. 5), such as those well established by the RDP's FunGene

438   database[45]. Building a dynamic metagenomic bioinformatics community will help us

439   better understand gene function. In principle, FunGeneTyper can predict the function

440   of any gene category based on prior parameters of the pre-trained model and the

441   adapter's transfer-learning ability. The adapter module used in FunGeneTyper is a

442   lightweight plug-and-play neural network that only fine-tunes and maintains a small

443   set of parameters and is conducive for sharing and promotion. Therefore, other

444    researchers can use the framework and training parameters we provide to train their

445    own core datasets to easily develop predictive deep learning models of genes of

446    interest. Researchers can also share a trained adapter through the adapter sharing

447    community (ASC) without disclosing their private datasets. The future prosperity and

448    collaboration of the ASC under the guidance of FunGeneTyper framework provide an

449    interactive, dynamic, and continuously improving or evolving platform for functional

450    classification of various PCG sequences. More importantly, FunGeneTyper and ASC

451    are expected to contribute significantly to advances in industrial biotechnology, health

452    and medicine, food and agriculture, environmental biotechnology, and bioenergy (Fig.

453    5), as they are increasingly applied to accelerate the discovery of new genes and

454    enzymatic resources from microbiomes.

455        In conclusion, FunGeneTyper provides an innovative and unified framework with

456    deep learning models (i.e., FunTrans and FunRep), expandable classifier toolkits (e.g.,

457    ARGTyper and VFGTyper) and customizable structured databases for the ultra-

458    accurate classification and discovery of functional genes (e.g., ARGs and VFGs) that

459    have scientific and biotechnological significance. This framework will contribute to

460    the robust monitoring of function genes and discovery of novel enzymatic resources

461    from diverse microbiomes and uncultured microbes therein, which is critical to

462    understand and harness the microbiome sciences underlying environment

463    (biogeochemistry, bio-restoration, and bioremediation)[14] bioeconomy (bioenergy and

464    bioresources)[13], and human systems (food and health)[20,46].

465 **Methods**

466 **Collection and expansion of the core dataset**

467 The core dataset used for FunGeneTyper model training is a set of experimentally-

468 confirmed reference sequences of target functional genes collected from literature

469 and/or expert-curated databases. Because the core dataset does not always contain a

470 sufficient number of experimentally-confirmed sequences (no more than 10

471 sequences[40]) for every type or subtype of functional gene, it is expanded to retrieve

472 more sequence data to improve and optimize the training of deep learning models. In

473 the subsequent training method, which separates the extended categories of five or

474 more sequences into the training set, verification set, and testing set at a ratio of 6:2:2,

475 any categories that are unsuitable for inclusion in these five sequences are included in

476 the training set.

477 **Construction of structured antibiotic resistance database (SARD)**

478 *Core ARGs dataset*

479 To ensure the professionalism and accuracy of the training dataset, reference protein

480 sequences of ARGs defined by homologs in the authoritative CARD were selected as

481 core data for downstream model training. The sequences were clustered using CD-

482 HIT[47] ($v$4.8.1) at an amino acid sequence identity of 100%, and all protein sequences

483 and their ontological information were manually checked to ensure that each ARG

484 was properly classified into class (type) and group (subtype) based on their

485  ontological information. Generally, class is equivalent to CARD's ontology terms for

486  antibiotic drug types, and group is equivalent to the specific sequence category.

487  Macrolides, lincosamides, and streptogramins were combined into the MLS class.

488  Based on the above procedures, a core dataset of 2972 non-redundant sequences

489  representing 2972 groups of ARGs from 19 classes was obtained and used to build the

490  SARD, which was used in subsequent analyses.

491  *Expanded ARGs dataset*

492  To ensure sufficient training data, the core dataset was expanded by retrieving close

493  homologues of its ARGs from the Uniref100 database following strict screening

494  criteria. Briefly, Diamond[6] (version 2.0.15) was used to index the ARG sequences in

495  the core dataset and to search for homologous sequences with an amino acid identity

496  and coverage greater than or equal to 80%. The extracted candidate sequences were

497  dereplicated and used as expanded datasets.

498  *Negative dataset*

499  To ensure that the model can learn sufficient features of non-target gene function,

500  which is essential for robustly predicting target function directly from metagenomic

501  data, we used the Swiss-Prot database, an expert-validated protein database, to

502  generate a negative dataset for use as a non-ARG training set. First, protein sequences

503  associated with antibiotic resistance in the Swiss-Prot database were screened out

504  using the keywords KW-0046. The remaining sequences were aligned against the core

505  ARGs dataset using Diamond software. Sequences with an alignment coverage

506    greater than 80% were extracted and categorized into four negative datasets based on

507    their sequence alignment identity (ID): identity 0 (ID ≤0%), identity 30 (ID ≤30%),

508    identity 50 (ID ≤50%), and identity 80 (ID ≤80%).

509    **Construction of structured virulence factor database (SVFD)**

510    *Core VFGs dataset*

511    Virulence factor databases were collected from VFNet[40]. Zheng et al[40] performed a

512    detailed similarity search for known and potential VFGs in the complete bacterial

513    genome downloaded from the NCBI server using VFanalyzer[48], with Virulence Factor

514    Database (VFDB) as the core database[48]. VFNet is an expanded virulence factor

515    database that can be used directly in the training process.

516    *Negative dataset*

517    The non-VFG collection process is similar to that of the non-ARG collection process,

518    except that KW-0800 is used to filter sequences from Swiss-Prot database (version:

519    June 2021).

520    **Architecture of the FunGeneTyper model**

521    FunGeneTyper is a universal function classification framework composed of two core

522    deep learning models, FunTrans and FunRep, which share similar structures but are

523    designed to classify functional genes at the type and subtype levels, respectively. Both

524    models are modular adapter-based architectures that leverage a few extra parameters

525    to achieve efficient fine-tuning of large-scale PLMs. In detail, utilizing the state-of-

526    the-art large-scale protein PLM esm-1b as a 33-layer transformer encoder framework

527    as the foundation, we plug adapters in each transformer layer of the PLM, which are

528    individual modular units that are used as newly introduced weights to be fine-tuned

529    for specific functional tasks. Notably, ESM-1b, through self-supervised learning on

530    the UniRef50 dataset, was shown to have a superior capacity to infer fundamental

531    structural and functional characteristics of proteins from gene sequences[49].

532       The holistic architecture is depicted in Fig. 1a and consists of three main

533    components: a multi-headed self-attention, a feed-forward network, and an adapter

534    layer. Each sublayer contains layer normalization and skip connections to effectively

535    train the neural network and avoid overfitting. It is worth noting that the bottleneck-

536    shaped adapter module consists of a down-project linear $H \in \square^{d \times k} H \in \square^{d \times k}$

537    where $d$ is embedding size of the Transformer model, $k$ is the dimension of the

538    adapter and $d \gg k$, a ReLU activation followed by an up-projection $L \in \square^{k \times d} L \in$

539    $\square^{k \times d}$. The adapter layer is formulated as follows:

$$O_l = LayerNorm(T_l)$$
$$\text{Adapter}_l = L_l(\text{ReLU}(H_l(O_l)) + T_l$$

540    where $T_l$ is the hidden feature at transformer layer $l$, $d = 1280$, and $k = 256$ in the

541    actual training.

542       Following the approach of BERT[50], hidden features from the first token of the

543    sequence of the last layer are extracted. In contrast to FunTrans, which adds a

544    nonlinear layer for protein function classification after the representations of the last

545   layer, FunRep first computes the hidden features of experimentally-confirmed core

546   sequences and then annotates PCGs by finding the sequence's category with the

547   closest Euclidean distance in the representation space.

548   Here, we use a dual-tower architecture with shared parameters similar to

549   Sentence-BERT[51] for model training in order to place sequences with the same

550   category closer in the representation space. FunRep is trained by constructing

551   $< A, P, N >$ triples, where $A$ is the anchor sequence, $P$ is a positive example

552   possessing the same category as $A$, and $N$ is a negative example whose category is

553   different from $A$ and the hidden representations they obtained through FunRep are $a$, $p$,

554   and $n$, respectively. The loss function adopts Triplet Loss, which is defined as follows:

$$Loss(a, p, n) = max(D(a, p) - D(a, n)) + margin, 0)$$

555   where $D$ is the Euclidean distance between vectors, and $margin$ is an adjustable

556   threshold, set to 1.0 during model training. ARGTyper-FunRep and VFGTyper-

557   FunRep are classified at the group level with the same 21.76M learnable training

558   parameters.

559   **Training settings**

560   All datasets are divided into training, validation and testing in a 6:2:2 ratio, and the

561   five-fold cross-validation is performed. Adam optimizer with default parameters is

562   used, dropout is set to 0.2, learning rate is 1e-5, and the early stopping method is

563   adopted to prevent overfitting. The accuracy, precision, recall and F1-score are used to

564    evaluate the performance. As a result, the micro average of the F1-score also equals

565    that of precision and recall, as well as the overall accuracy. Thus, we report only the

566    overall accuracy for the micro average metrics while reporting precision, recall and

567    F1-score for the macro average metrics.

568    **Evaluation of FunGeneTyper for the discovery of new functional genes**

569    To validate the capacity of FunGeneTyper models in discovering new functional

570    genes, experimentally confirmed ARGs from functional metagenomics studies were

571    retrieved from NCBI's protein database (accession numbers in Dataset S8). After

572    removing those ARGs showing perfect sequence match to the CARD database (or

573    core dataset of ARGs), 297 experimentally confirmed ARG sequences of human gut[35]

574    (n = 168), WWTPs[11] (n = 77), and soil[36-39] (n = 52) bacteria were retained for use in

575    the downstream comparisons between FunGeneTyper and the well-established SA-

576    based (RGI[7]), HMM-based (Resfams[18]), and DL-based (DeepARG[9] and HMD-ARG[32])

577    approaches in terms of classification performance of the new ARGs. In this study,

578    four evaluation metrics including the accuracy, precision, recall and F1-score were

579    computed to assess the multi-classification results performance using the following

580    equations:

$$Accuracy = \frac{TP + TN}{TP + FP + FN + TN}$$

$$Precision = \frac{TP}{TP + FP}$$

$$Recall = \frac{TP}{TP + FN}$$

$$F1\ Score = \frac{2 * Precision * Recall}{Precision + Recall}$$

581 where TP is the number of true positives, TN is the number of true negatives, FP is the

582 number of false positives, and FN is the number of false negatives.

583     To compare the ability of FunGeneTyper for discovering new VFGs, BoNTs-like

584 sequences from the genome of *Chryseobacterium piperi* reported in a prior study [41]

585 was downloaded from NCBI's database by accession number (Dataset S7). Then,

586 VFGTyper was used to predict VFGs and their affiliated family from the BoNTs-like

587 sequences, and the output results were compared with those by a conventional

588 sequence alignment-based approach with Diamond[6] (version 2.0.15) search of the

589 BoNTs-like sequences against SVFD.

590 **Abbreviations**

591 DL: Deep Learning

592 SFGD: Structured Functional Gene Dataset

593 PCG: Protein-Coding Gene

594 ARG: Antibiotic Resistance Gene

595 VFG: Virulence Factor Gene

596 MLS: Macrolides, Lincosamides and Streptogramines

597 **Author Contributions**

598 F. Ju conceived the FunGeneTyper framework idea, obtained funding, and supervised

599 the project. F. Yuan designed the Adapter sharing mechanism. G. Zhang and H. Wang

600   (visiting student from Northeastern University) performed the model construction,

601   data analysis and visualization. F. Ju and F. Yuan co-supervised G. Zhang and H.

602   Wang on the deep-learning model construction with additional support from G. Guo.

603   G. Zhang built the structured databases and accomplished data presentation with the

604   assistance from J. Yang, Z. Zhang, L. Zhang. F. Ju and G. Zhang co-wrote the

605   manuscript with assistance from F. Yuan and H. Wang. J. All authors approved the

606   final version of the manuscript.

## Reference

1 Ju, F. *et al.* Wastewater treatment plant resistomes are shaped by bacterial composition, genetic exchange, and upregulated expression in the effluent microbiomes. *ISME J* **13**, 346-360 (2019). https://doi.org/10.1038/s41396-018-0277-8

2 Pasolli, E. *et al.* Extensive Unexplored Human Microbiome Diversity Revealed by Over 150,000 Genomes from Metagenomes Spanning Age, Geography, and Lifestyle. *Cell* **176**, 649-662 e620 (2019). https://doi.org/10.1016/j.cell.2019.01.001

3 Nayfach, S. *et al.* A genomic catalog of Earth's microbiomes. *Nat Biotechnol* **39**, 499-509 (2021). https://doi.org/10.1038/s41587-020-0718-6

4 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *Journal of Molecular Biology* **215**, 403-410 (1990). https://doi.org/10.1016/s0022-2836(05)80360-2

5 Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460-2461 (2010).

6 Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND. *Nat Methods* **12**, 59-60 (2015). https://doi.org/10.1038/nmeth.3176

7 Alcock, B. P. *et al.* CARD 2020: antibiotic resistome surveillance with the comprehensive antibiotic resistance database. *Nucleic Acids Res* **48**, D517-D525 (2020). https://doi.org/10.1093/nar/gkz935

8 Yang, Y. *et al.* ARGs-OAP: online analysis pipeline for antibiotic resistance genes detection from metagenomic data using an integrated structured ARG-database. *Bioinformatics* **32**, 2346-2351 (2016). https://doi.org/10.1093/bioinformatics/btw136

9 Arango-Argoty, G. *et al.* DeepARG: a deep learning approach for predicting antibiotic resistance genes from metagenomic data. *Microbiome* **6**, 23 (2018). https://doi.org/10.1186/s40168-018-0401-z

10 de Nies, L. *et al.* PathoFact: a pipeline for the prediction of virulence factors and antimicrobial resistance genes in metagenomic data. *Microbiome* **9**, 49 (2021). https://doi.org/10.1186/s40168-020-00993-9

11 Munck, C. *et al.* Limited dissemination of the wastewater treatment plant core resistome. *Nat Commun* **6**, 8452 (2015). https://doi.org/10.1038/ncomms9452

12 Forsberg, K. J. *et al.* Bacterial phylogeny structures soil resistomes across habitats. *Nature* **509**, 612-616 (2014). https://doi.org/10.1038/nature13377

13 Díaz Rodríguez, C. A. *et al.* Novel bacterial taxa in a minimal lignocellulolytic consortium and their potential for lignin and plastics transformation. *ISME Communications* **2** (2022). https://doi.org/10.1038/s43705-022-00176-7

14 Royo-Llonch, M. *et al.* Compendium of 530 metagenome-assembled bacterial and archaeal genomes from the polar Arctic Ocean. *Nat Microbiol* **6**, 1561-

660       1574 (2021). https://doi.org/10.1038/s41564-021-00979-9

661   15   Feldgarden, M. *et al.* Validating the AMRFinder Tool and Resistance Gene
662       Database by Using Antimicrobial Resistance Genotype-Phenotype
663       Correlations in a Collection of Isolates. *Antimicrob Agents Chemother* **63**
664       (2019). https://doi.org/10.1128/AAC.00483-19

665   16   Xie, G. & Fair, J. M. Hidden Markov Model: a shortest unique representative
666       approach to detect the protein toxins, virulence factors and antibiotic
667       resistance genes. *BMC Res Notes* **14**, 122 (2021).
668       https://doi.org/10.1186/s13104-021-05531-w

669   17   Boolchandani, M., D'Souza, A. W. & Dantas, G. Sequencing-based methods
670       and resources to study antimicrobial resistance. *Nat Rev Genet* **20**, 356-370
671       (2019). https://doi.org/10.1038/s41576-019-0108-4

672   18   Gibson, M. K., Forsberg, K. J. & Dantas, G. Improved annotation of antibiotic
673       resistance determinants reveals microbial resistomes cluster by ecology. *ISME*
674       *J* **9**, 207-216 (2015). https://doi.org/10.1038/ismej.2014.106

675   19   Biswas, S., Khimulya, G., Alley, E. C., Esvelt, K. M. & Church, G. M. Low-N
676       protein engineering with data-efficient deep learning. *Nat Methods* **18**, 389-
677       396 (2021). https://doi.org/10.1038/s41592-021-01100-y

678   20   Ma, Y. *et al.* Identification of antimicrobial peptides from the human gut
679       microbiome using deep learning. *Nat Biotechnol* (2022).
680       https://doi.org/10.1038/s41587-022-01226-0

681   21   Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M.
682       Unified rational protein engineering with sequence-based deep representation
683       learning. *Nat Methods* **16**, 1315-1322 (2019). https://doi.org/10.1038/s41592-
684       019-0598-1

685   22   Rives, A. *et al.* Biological structure and function emerge from scaling
686       unsupervised learning to 250 million protein sequences. *Proc Natl Acad Sci U*
687       *S A* **118** (2021). https://doi.org/10.1073/pnas.2016239118

688   23   Ofer, D., Brandes, N. & Linial, M. The language of proteins: NLP, machine
689       learning & protein sequences. *Comput Struct Biotechnol J* **19**, 1750-1758
690       (2021). https://doi.org/10.1016/j.csbj.2021.03.022

691   24   Unsal, S. *et al.* Learning functional properties of proteins with language
692       models. *Nature Machine Intelligence* **4**, 227-245 (2022).
693       https://doi.org/10.1038/s42256-022-00457-9

694   25   Bileschi, M. L. *et al.* Using deep learning to annotate the protein universe. *Nat*
695       *Biotechnol* (2022). https://doi.org/10.1038/s41587-021-01179-w

696   26   Dohan, D., Gane, A., Bileschi, M. L., Belanger, D. & Colwell, L. Improving
697       Protein Function Annotation via Unsupervised Pre-training: Robustness,
698       Efficiency, and Insights.*Proceedings of the 27th ACM SIGKDD Conference on*
699       *Knowledge Discovery & Data Mining.* 2782–2791 (Association for
700       Computing Machinery).

701   27   Yuan, F., He, X., Karatzoglou, A. & Zhang, L. Parameter-Efficient Transfer

from Sequential Behaviors for User Modeling and Recommendation.*Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval.* 1469-1478.

28  Houlsby, N. *et al.* Parameter-Efficient Transfer Learning for NLP.*Proceedings of the 36th International Conference on Machine Learning.* (eds Chaudhuri Kamalika & Salakhutdinov Ruslan) 2790--2799 (PMLR).

29  Murray, C. J. L. *et al.* Global burden of bacterial antimicrobial resistance in 2019: a systematic analysis. *The Lancet* **399**, 629-655 (2022). https://doi.org/10.1016/s0140-6736(21)02724-0

30  Du, D. *et al.* Multidrug efflux pumps: structure, function and regulation. *Nat Rev Microbiol* **16**, 523-539 (2018). https://doi.org/10.1038/s41579-018-0048-6

31  Piddock, L. J. Clinically relevant chromosomally encoded multidrug resistance efflux pumps in bacteria. *Clin Microbiol Rev* **19**, 382-402 (2006). https://doi.org/10.1128/CMR.19.2.382-402.2006

32  Li, Y. *et al.* HMD-ARG: hierarchical multi-task deep learning for annotating antibiotic resistance genes. *Microbiome* **9**, 40 (2021). https://doi.org/10.1186/s40168-021-01002-3

33  Rost, B. Twilight zone of protein sequence alignments. *Protein Eng* **12**, 85-94 (1999). https://doi.org/10.1093/protein/12.2.85

34  Bepler, T. & Berger, B. Learning the protein language: Evolution, structure, and function. *Cell Syst* **12**, 654-669 e653 (2021). https://doi.org/10.1016/j.cels.2021.05.017

35  Sommer, M. O. A., Dantas, G. & Church, G. M. Functional characterization of the antibiotic resistance reservoir in the human microflora. *Science* **325**, 1128-1131 (2009). https://doi.org/10.1126/science.1176950

36  Willms, I. M. *et al.* Novel Soil-Derived Beta-Lactam, Chloramphenicol, Fosfomycin and Trimethoprim Resistance Genes Revealed by Functional Metagenomics. *Antibiotics (Basel)* **10** (2021). https://doi.org/10.3390/antibiotics10040378
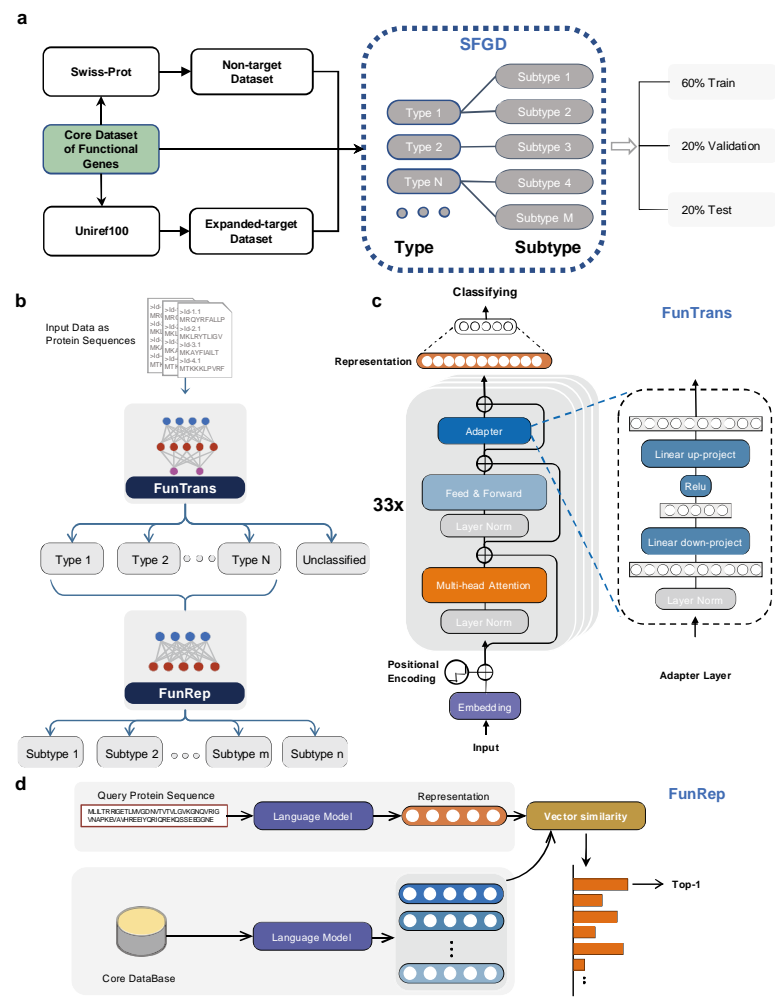
37  Wang, S. *et al.* Tetracycline Resistance Genes Identified from Distinct Soil Environments in China by Functional Metagenomics. *Front Microbiol* **8**, 1406 (2017). https://doi.org/10.3389/fmicb.2017.01406

38  Allen, H. K., Moe, L. A., Rodbumrer, J., Gaarder, A. & Handelsman, J. Functional metagenomics reveals diverse beta-lactamases in a remote Alaskan soil. *ISME J* **3**, 243-251 (2009). https://doi.org/10.1038/ismej.2008.86

39  Donato, J. J. *et al.* Metagenomic analysis of apple orchard soil reveals antibiotic resistance genes encoding predicted bifunctional proteins. *Appl Environ Microbiol* **76**, 4396-4401 (2010). https://doi.org/10.1128/AEM.01763-09

40  Zheng, D., Pang, G., Liu, B., Chen, L. & Yang, J. Learning transferable deep convolutional neural networks for the classification of bacterial virulence
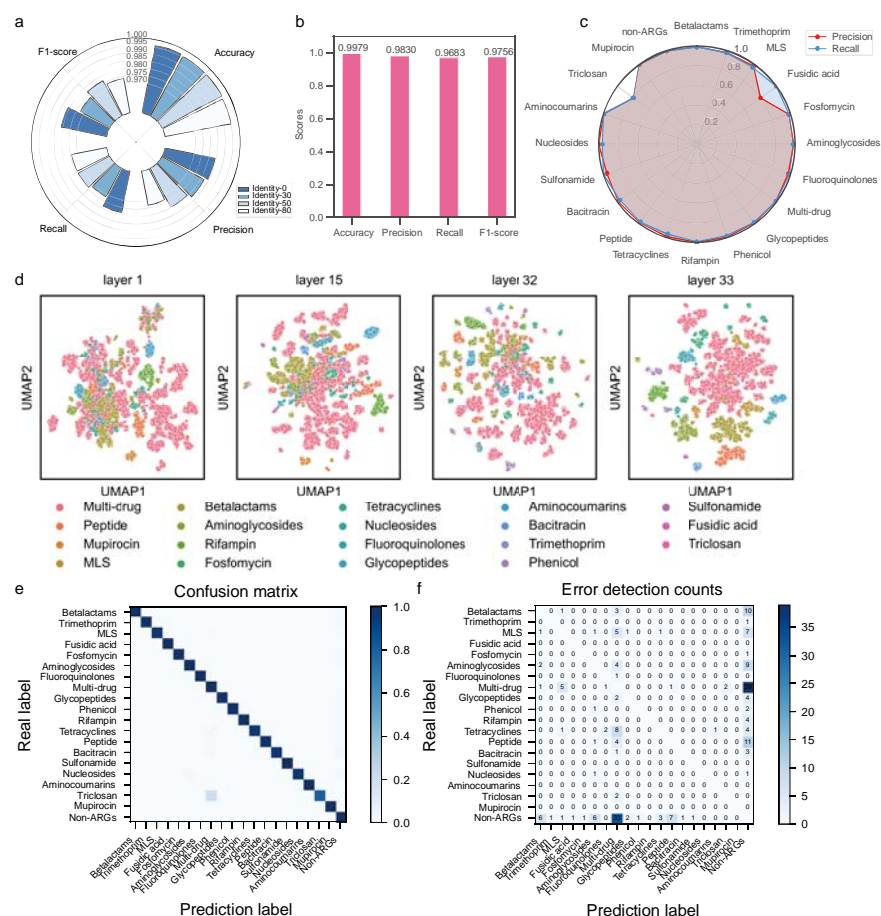
744    factors.    *Bioinformatics*    **36**,    3693-3702    (2020).
745    https://doi.org/10.1093/bioinformatics/btaa230

746  41  Mansfield, M. J. *et al.* Bioinformatic discovery of a toxin family in
747    Chryseobacterium piperi with sequence similarity to botulinum neurotoxins.
748    *Sci Rep* **9**, 1634 (2019). https://doi.org/10.1038/s41598-018-37647-8

749  42  WHO, O. One health. *World Health Organization* (2017).

750  43  Wang, Z. *et al.* ARG-SHINE: improve antibiotic resistance class prediction by
751    integrating sequence homology, functional information and deep convolutional
752    neural    network.    *NAR    Genom    Bioinform*    **3**,    lqab066    (2021).
753    https://doi.org/10.1093/nargab/lqab066

754  44  Chen, P., Ye, J., Chen, G., Zhao, J. & Heng, P.-A. Robustness of accuracy
755    metric and its inspirations in learning with noisy labels.*Proceedings of the*
756    *AAAI Conference on Artificial Intelligence.*  11451-11461.

757  45  Fish, J. A. *et al.* FunGene: the functional gene pipeline and repository. *Front*
758    *Microbiol* **4**, 291 (2013). https://doi.org/10.3389/fmicb.2013.00291

759  46  Lee, E. D., Aurand, E. R., Friedman, D. C. & Engineering Biology Research
760    Consortium    Microbiomes    Roadmapping    Working,    G.    Engineering
761    Microbiomes-Looking    Ahead.    *ACS    Synth    Biol*    **9**,    3181-3183    (2020).
762    https://doi.org/10.1021/acssynbio.0c00558

763  47  Li, W. & Godzik, A. Cd-hit: a fast program for clustering and comparing large
764    sets of protein or nucleotide sequences. *Bioinformatics* **22**, 1658-1659 (2006).
765    https://doi.org/10.1093/bioinformatics/btl158

766  48  Liu, B., Zheng, D., Jin, Q., Chen, L. & Yang, J. VFDB 2019: a comparative
767    pathogenomic platform with an interactive web interface. *Nucleic Acids Res*
768    **47**, D687-D692 (2019). https://doi.org/10.1093/nar/gky1080

769  49  Rao, R., Meier, J., Sercu, T., Ovchinnikov, S. & Rives, A. Transformer protein
770    language    models    are    unsupervised    structure    learners.    *bioRxiv*,
771    2020.2012.2015.422761 (2020). https://doi.org/10.1101/2020.12.15.422761

772  50  Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of
773    Deep Bidirectional Transformers for Language Understanding. 4171-4186
774    (2019). https://doi.org/10.18653/v1/N19-1423

775  51  Reimers, N. & Gurevych, I. Sentence-BERT: Sentence Embeddings using
776    Siamese BERT-Networks.*Proceedings of the 2019 Conference on Empirical*
777    *Methods in Natural Language Processing and the 9th International Joint*
778    *Conference on Natural Language Processing (EMNLP-IJCNLP).*  3982-3992.

779

780

**Fig. 1** FunGeneTyper model design and database construction workflows. **a**, Process

of preparing a structured functional gene dataset (SFGD). The data set is divided into

the training set, validation set and testing set in a 6:2:2 ratio. **b**, Two-level hierarchical

structure of FunGeneTyper. **c**, Schematic representation of FunTrans model. **d**,

Schematic representation of FunRep model.

786

787

**Fig. 2. Performance evaluation of deep-learning FunGeneTyper models with structured Antibiotic Resistance Gene Database (SARD) for classification of ARGs. a**, Evaluation of the influence of identity threshold used for selecting the negative dataset on model performance in the classification of ARGs. **b**, Performance metrics of ARGTyper developed based on FunGeneTyper models and SARD. **c**, Performance of all 19 c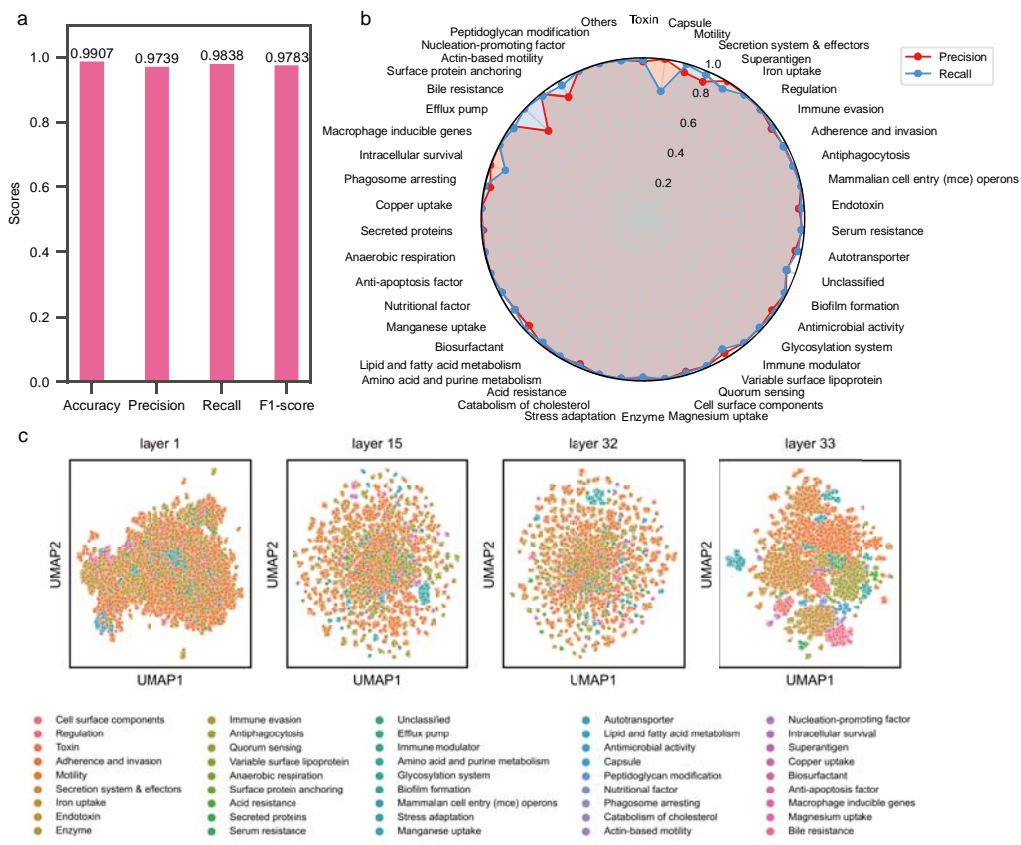lasses as indicated by precision and recall of ARGs and non-ARG classes. **d**, Visualization of feature learning at different layers during the ARGTyper training process. **e**, Confusion matrix for ARG class classification, confusion between true (y-axis) and predicted (x-axis) ARGs. **f**, Number of ARG protein sequences annotated incorrectly. MLS: Macrolides, Lincosamides and Streptogramines.
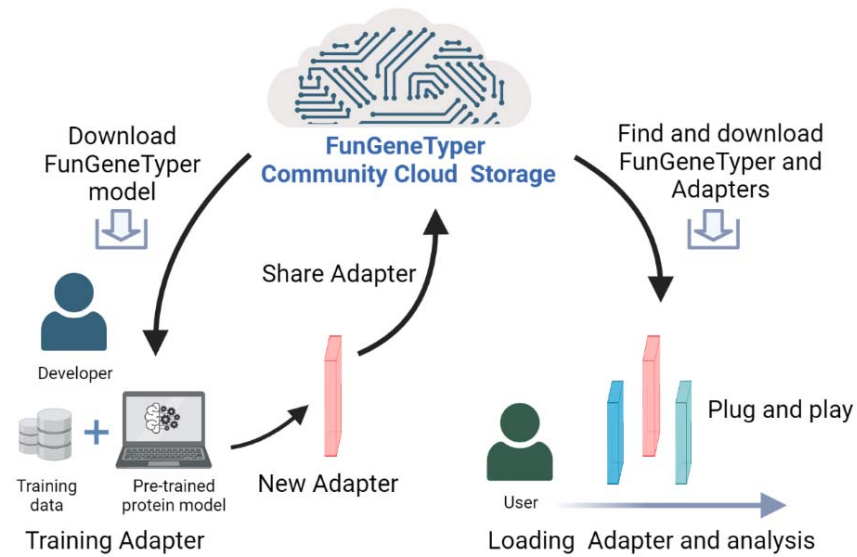
**Fig. 3. Transfer learning of FunGeneTyper models on Structured Virulence Factor Gene Database (VFGD) and performance evaluation for VFG classification.** **a**, Performance metrics of VFGTyper developed based on FunGeneTyper models and VFGD. **b**, Precision and Recall of VFGs family and non-VFGs category. **c**, Visualization of feature learning at different layers in VFGs FunTrans training. VFGs: virulence factor genes.

807

**Fig. 4. Schematic of the Adapter Sharing Community (ASC) in the framework of FunGeneTyper.** The community developers are cyber de-centralized to train customizable structure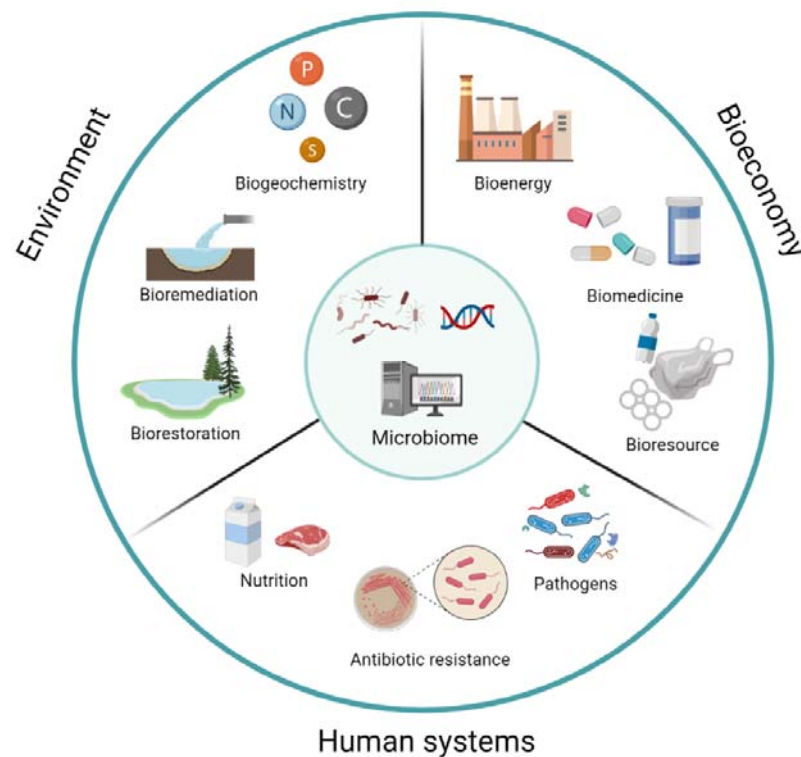d databases and develop deep learning classifiers of various categories of functional genes, while users utilize the classifiers of interest to accelerate the discovery of genes which, in turn, provide new experimentally-confirmed sequences to expand the structured databases and improve deep-learning models.

816

**Fig. 5. Potential applications of FunGeneTyper to the discovery of microbiome resources for enhancing our environment, bioeconomy, and human systems.** Metagenomic discovery of precious genetic and enzymatic resources facilitated by the Adapter Sharing Community of FunGeneTyper can contribute to follow-up microbiome, genetic and protein engineering researches for enhancing human health and eco-environment systems.

823 **Table 1 Performance comparison between FunGeneTyper and other alternative bioinformatics tools for the discovery of experimentally**

824 **confirmed new ARGs.** In total, 297 experimentally confirmed ARGs sequences of human gut[35] (n = 168), WWTPs[11] (n = 77), and soil[36-39] (n =

825 52) bacteria were included in the comparative analysis which was performed under the default settings of each deep-learning (DL)-based,

826 sequence alignment (SA)-based or Hidden Markov Model (HMM)-based tool recommended by the developers.

| Tools | Human gut (n=168) | | | | WWTP (n=77) | | | | Soil (n=52) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Accuracy | Precision | Recall | F1-score | Accuracy | Precision | Recall | F1-score | Accuracy | Precision | Recall | F1-score |
| *DL-based tools* | | | | | | | | | | | | |
| FunGeneTyper | 0.8512 | **0.7500** | **0.6642** | **0.6948** | **0.7273** | **0.7500** | **0.5403** | **0.6072** | **0.8269** | 0.5926 | **0.5529** | **0.5445** |
| HMD-ARG | 0.8452 | 0.6000 | 0.5230 | 0.5486 | 0.5714 | 0.7161 | 0.3877 | 0.4589 | 0.8077 | **0.6000** | 0.4560 | 0.5119 |
| DeepARG | 0.3512 | 0.6250 | 0.4720 | 0.5149 | 0.1688 | 0.5714 | 0.1682 | 0.2591 | 0.2885 | 0.3750 | 0.1057 | 0.1607 |
| *SA-based tools* | | | | | | | | | | | | |
| RGI | 0.3452 | 0.6250 | 0.4596 | 0.5065 | 0.0390 | 0.3750 | 0.0349 | 0.0632 | 0.1538 | 0.1250 | 0.0357 | 0.0556 |
| *HMM-based tools* | | | | | | | | | | | | |
| Resfams | **0.8830** | 0.4545 | 0.3968 | 0.4195 | 0.6234 | 0.6250 | 0.4736 | 0.5224 | 0.8088 | 0.2727 | 0.2545 | 0.2630 |

827

828