

# 1    **Ontology-aware deep learning for antibiotic resistance gene** 2    **prediction: novel function discovery and comprehensive profiling** 3    **from metagenomic data**

4

5    Yuguang Zha<sup>1</sup>, Cheng Chen<sup>2</sup>, Qihong Jiao<sup>2</sup>, Xiaomei Zeng<sup>1,\*</sup>, Xuefeng Cui<sup>2,\*</sup>, Kang  
6    Ning<sup>1,\*</sup>

7    <sup>1</sup>Key Laboratory of Molecular Biophysics of the Ministry of Education, Hubei Key  
8    Laboratory of Bioinformatics and Molecular-imaging, Center of AI Biology,  
9    Department of Bioinformatics and Systems Biology, College of Life Science and  
10    Technology, Huazhong University of Science and Technology, Wuhan 430074, Hubei,  
11    China

12    <sup>2</sup>School of Computer Science and Technology, Shandong University, Qingdao 266237,  
13    Shandong, China

14    \*Correspondence should be addressed to K.N (Email: ningkang@hust.edu.cn), X.C  
15    (Email: xfcui@email.sdu.edu.cn) and X.Z (Email: xmzeng@hust.edu.cn)

16

## 17 **Abstract**

18 Antibiotic resistance genes (ARGs) have emerged in pathogens and arousing a  
 19 worldwide concern, which is estimated to cause millions of deaths each year globally.  
 20 Accurately identifying and classifying ARGs is a formidable challenge in studying the  
 21 generation and spread of antibiotic resistance. Current methods could identify close  
 22 homologous ARGs, have limited utility for discovery of novel ARGs, thus rendering  
 23 the profiling of ARGs incomprehensive. Here, an ontology-aware neural network  
 24 (ONN) approach, ONN4ARG, is proposed for comprehensive ARG discovery.  
 25 Systematic evaluation shows ONN4ARG is advanced than previous methods such as  
 26 DeepARG in efficiency, accuracy, and comprehensiveness. Experiments using 200  
 27 million candidate microbial genes collected from 815 microbial community samples  
 28 from diverse environments or hosts have resulted in 120,726 candidate ARGs, out of  
 29 which more than 20% are not yet present in public databases. These comprehensive  
 30 set of ARGs have clarified the environment-specific and host-specific patterns. The  
 31 wet-experimental functional validation, together with structural investigation of  
 32 docking sites, have also validated a novel streptomycin resistance gene from oral  
 33 microbiome samples, confirming ONN4ARG's ability for novel ARGs identification.  
 34 In summary, ONN4ARG is superior to existing methods in efficiency, accuracy, and  
 35 comprehensiveness. It enables comprehensive ARG discovery, which is helpful  
 36 towards a grand view of ARGs worldwide. ONN4ARG is available at  
 37 <https://github.com/HUST-NingKang-Lab/ONN4ARG>, and online web service is  
 38 available at <http://onn4arg.xfcui.com/>.

39

40 **Keywords:** antibiotic resistance gene, ontology-aware neural network, deep learning,  
 41 novel ARG, microbiome

42

## Introduction

With the development of metagenomics and next-generation sequencing, many new microbial taxa and genes have been discovered, but different kinds of “unknowns” remain. For instance, the microbes found in the human gut microbiome involve 25 phyla, more than 2,000 genera, and 5,000 species [1]. However, the functional diversity of microbiomes has not been fully explored, and about 40% of microbial gene functions remain to be discovered [2]. A typical example is the antibiotic resistance gene (ARG), which is an urgent and growing threat to public health [3]. In the past few decades, problems caused by antibiotic resistance have drawn the public’s attention [4]. Antimicrobial resistance genomic data is an ever-expanding data source, with many new ARG families discovered in recent years [5, 6]. The discovery of resistance genes in diverse environments offers possibilities for early surveillance, actions to reduce transmission, gene-based diagnostics, and improved treatment [7].

Existing annotated ARGs have been curated manually or automatically for decades. Presently, there are 4,661 annotated ARGs in the reference database CARD [5, 6] (v3.2.5, released in September 2022), 3,131 in the ResFinder database [8] (as of December 2022), and 2,476 in SwissProt [9] (as of December 2022). These annotated ARGs are categorized into antibiotic resistance types, which are organized in an ontology structure (**Methods, Supplementary Figure S1**), in which higher-level ARG types cover lower-level ARG types. Current ARG databases are far from complete: though no ARG database contains more than 4,000 well-annotated ARGs, NCBI non-redundant database searches yielded more than 7,000 putative genes annotated with “antibiotic resistance” as of May 2021. Therefore, we deemed that there is a large gap between the genes annotated in ARG databases and the possible ARGs that already exist in general databases, not to mention ARGs that are not yet annotated.

72 Many ARG prediction tools have been proposed in the past few years [8, 10-20].  
 73 These tools can generally be divided into two approaches. One approach is  
 74 sequence-alignment, such as BLAST [21], USEARCH [22], and Diamond [23], which  
 75 uses homologous genes to annotate unclassified genes. A confident prediction requires  
 76 a homolog with sequence identity greater than 80% in many programs, such as  
 77 ResFinder [8, 11]. The other approach is deep learning, such as DeepARG [12] and  
 78 HMD-ARG [16], which uses neural network models to predict and annotate ARGs.

79

80 Several limitations still preclude comprehensive profiling of ARGs. A more  
 81 comprehensive set of ARGs could be roughly defined as having more ARGs in type  
 82 and number with less false-positive entries, regardless of the homology with known  
 83 ARGs, and many of these ARGs could be experimentally validated. Based on this  
 84 definition, existing tools fall short in comprehensive profiling of ARGs. First, existing  
 85 tools are limited to a few types of ARGs due to the fact that the datasets used for  
 86 building models are specialized. For example, HMD-ARG [16] identifies only 15  
 87 types of resistance genes, and PATRIC [13] is limited to identifying ARGs encoding  
 88 resistance to carbapenem, methicillin, and beta-lactam antibiotics. Second, existing  
 89 tools fall short in discovering novel ARGs, which usually lack homology to known  
 90 sequences in the reference databases. For instance, the gene POCOZ1 (VraR) that  
 91 confers resistance to vancomycin has a sequence identity of only 24% to the homolog  
 92 from the CARD [12]. Therefore, there is an urgent need for a new approach to address  
 93 these limitations.

94

95 Here, we propose an ontology-aware deep learning approach, ONN4ARG, which  
 96 allows comprehensive identification of ARGs. Systematic evaluation based on the  
 97 ONN4ARG-DB, CARD, and ResFinder datasets shows that the ONN4ARG model  
 98 outperforms state-of-the-art models such as DeepARG, especially for the detection of  
 99 remotely homologous ARGs. Experiments based on more than 200 million candidate  
 100 microbial genes collected from 815 samples in various environments have resulted in  
 101 120,726 candidate ARGs, out of which more than 20% are not yet present in public

102 databases. Our experiments confirmed that ARGs are both environment-specific and  
103 host-specific, exemplified by the rifamycin resistance genes which are enriched in  
104 Actinobacteria and in soil environment. Case study of a recently experimentally  
105 validated ARG gene GAR [7] have also verified the ability of ONN4ARG for novel  
106 ARG discovery. We also validated a novel streptomycin resistance gene from oral  
107 microbiome samples by wet-lab experiment. In summary, ONN4ARG enables  
108 comprehensive ARG discovery, which provides a relatively complete picture of the  
109 prevalence of ARGs, as well as leads a way towards a grand view of ARGs  
110 worldwide.

111

## 112 **Results**

### 113 **ONN4ARG model employs an ontology-aware neural network for ARG** 114 **identification and classification**

115 To address the large gap between the genes annotated in ARG databases and the  
116 possible ARGs that already exist in general databases along with the ARGs that are  
117 not yet annotated, we propose ONN4ARG, which is an ontology-aware neural  
118 network model (**Figure 1, Supplementary Figure S1**) that predict ARGs in a  
119 comprehensive manner. ONN4ARG takes similarities (e.g., identity, e-value, bit-score)  
120 between the query gene sequence and ARG gene sequences and profiles (i.e., PSSM)  
121 as inputs and predicts ARG annotations (**Figure 1B**). These sequence-alignment  
122 similarities and profile-alignment similarities are pre-processed by calling Diamond  
123 [23] and HHblits [24]. ONN4ARG generates hierarchical annotations of antibiotic  
124 resistance types, which are compatible with the antibiotic resistance ontology  
125 structure (**Figure 1A, C**). One advantage of ONN4ARG over state-of-the-art models  
126 is that ONN4ARG employs a novel ontology-aware layer that incorporates ancestor  
127 and descendent annotations to enhance annotation accuracies (**Methods**). To train and  
128 evaluate our ONN4ARG model and for rapid deployment of ARG discovery in  
129 multiple contexts, we also built an ARG database (**Figure 1D**), namely,  
130 ONN4ARG-DB, which comprises ARGs from CARD and UniProt (see **Methods**).

131

## 132 **Systematic evaluation and comparison**

133 Systematic evaluation based on the ONN4ARG-DB showed our model's high  
 134 efficiency, high accuracy, and comprehensiveness for ARG identification. ONN4ARG  
 135 is fast since it could complete ARG identification for all genes in the testing dataset  
 136 within four hours, which is equivalent to one second per gene identification. As  
 137 shown in **Figure 2A**, ONN4ARG was more accurate for ARG identification (overall  
 138 accuracy of 97.70%, **Table 1**) compared to sequence alignment (overall accuracy of  
 139 69.11%), and ONN4ARG has a slight advantage over DeepARG (overall accuracy of  
 140 96.39%). Moreover, ONN4ARG achieved an overall precision of 75.59% and an  
 141 overall recall of 89.93%, which were higher than DeepARG's overall precision of  
 142 68.30% and overall recall of 77.84% (**Figure 2B**, **Table 2**). It is natural that  
 143 ONN4ARG could not outperform DeepARG in all resistance types and this is  
 144 exemplified by results on pleuromutilin due to the small number of sequences for  
 145 pleuromutilin in the ONN4ARG-DB. ONN4ARG demonstrates an advantage over  
 146 other methods in identification of remotely homologous ARGs whose sequences are  
 147 not similar to existing ARG sequences (**Tables 2 and 3**). In this context, when testing  
 148 with only remotely homologs (i.e., the masking threshold of testing set equal to 0.4,  
 149 see **Methods**), ONN4ARG achieves an accuracy of 94.26%, which is largely  
 150 improved from 89.85% of DeepARG. These results validate ONN4ARG's better  
 151 generalization abilities than sequence-alignment and DeepARG, which makes  
 152 ONN4ARG especially suitable for identification of remotely homologous ARGs and  
 153 indicates ONN4ARG's ability for novel ARG discovery (**Tables 1–3**).

154

155 We have also tested ONN4ARG on a verification set built from the CARD database  
 156 version 3.1.3. Results showed our model outperformed other methods in terms of  
 157 accuracy and efficiency, i.e., high accuracy and less time usage, given that the  
 158 memory usage is acceptable for a regular laptop (**Supplementary Table S1**). We have  
 159 also evaluated ONN4ARG on the ResFinder database version 4.1, which involves  
 160 thousands of manually curated ARGs [8]. Results showed that ONN4ARG achieved

161 an accuracy higher than 90% for most types of resistance, while DeepARG was less  
162 accurate than ONN4ARG, except for the fosfomycin resistance (**Supplementary**  
163 **Table S2**).

164

#### 165 **Applications of ONN4ARG on metagenomic data**

166 We collected metagenomic samples from several published studies [25, 26]. These  
167 samples were mainly from “marine,” “soil,” and “human” environments.  
168 Human-associated samples consisted of two gut groups (one group from Madagascar,  
169 i.e., GutM; the other group from Denmark, i.e., GutD), one oral group, and one skin  
170 group (both oral and skin groups were from the HMP project). For details on these  
171 samples, see **Supplementary Table S3**. Then, genes were obtained by calling  
172 Prodigal [27] with default parameters. The ONN4ARG model was used to predict  
173 whether these unclassified genes were ARGs and their corresponding resistance types.  
174 In total, 120,726 ARGs were identified from microbiome samples, many of which are  
175 novel, which greatly expands the existing ARG repositories.

176

#### 177 **Broad-spectrum profile of predicted ARGs among diverse environments**

178 We investigated the broad-spectrum profile of these predicted ARGs among diverse  
179 environments. First, we investigated the proportion of predicted ARGs for different  
180 sequence lengths. The distribution shows that about half of the predicted ARGs have a  
181 length of 128–256 amino acid residues (**Figure 3A**). We also analyzed the protein  
182 domain of these predicted ARGs by searching the conserved domain database (CDD,  
183 last update Aug 2022) using RPS-BLAST tool version 2.9.0. Results showed that  
184 most of these predicted ARGs (over 97%) have protein domains that resemble those  
185 with known catalytic activity and/or may bind to the antimicrobials they are predicted  
186 to elicit resistance against (**Supplementary Table S4**). Second, we found that  
187 human-associated microbiome samples carry a higher abundance of ARGs, especially  
188 for the oral group, in which more than one resistance gene could be observed out of a  
189 hundred genes on average (**Figure 3B, Supplementary Table S5**). Third, we tested  
190 the novelty of these predicted ARGs. We found that about a third of them (42,848 out

191 of all 120,726 ARGs) had sequence identity of less than 40% to their homologs in the  
 192 ONN4ARG-DB (**Figure 3C**). We define these ARGs as candidate novel ARGs, which  
 193 have low sequence identities when aligned to their homologs in the reference database  
 194 (i.e., ONN4ARG-DB). For example, we found 45% of predicted ARGs in the marine  
 195 group were candidate novel ARGs (**Figure 3C**).

196

197 In total, 31 ARG types were detected in these various environments (**Figure 3D**,  
 198 **Supplementary Figure S2**). The number of predicted ARG sequences for different  
 199 types varied greatly, from a few (i.e., nitrofurantoin) to thousands (i.e., fluoroquinolone).  
 200 In general, fluoroquinolone and tetracycline resistance genes were more abundant  
 201 than other types (**Figure 3D**). As expected, these abundant ARGs were usually  
 202 associated with the antibiotics used extensively in human medicine or veterinary  
 203 medicine, including growth promotion [28].

204

## 205 **Enrichment of predicted ARGs among diverse hosts and environments**

206 Rapid deciphering of potential antimicrobial-resistant pathogens is necessary for  
 207 effective public health monitoring. The host-tracking of ARGs allows for accurate  
 208 identification of pathogens. Therefore, we conducted taxonomy analysis to track the  
 209 hosts of these predicted ARGs by using Kraken2 [29]. Results showed that there are  
 210 949 genera, each genus carries at least one type of ARG (**Supplementary Table S6**).  
 211 The host composition and distribution of all classified ARGs for the most abundant 20  
 212 genera are displayed in **Supplementary Figure S3**. The host distribution shows that  
 213 these ARGs are primarily affiliated with Proteobacteria (38.2%). The most abundant  
 214 ARGs carried by the 20 genera were resistance types of fluoroquinolone, macrolide,  
 215 peptide, penam, and tetracycline, accounting for about half of the total ARGs.  
 216 Network inference based on strong (Spearman's  $\rho > 0.8$ ) and significant (Welch's  
 217  $t$ -test,  $P$ -value  $< 0.01$ ) correlations showed the co-occurrence patterns among ARGs  
 218 and microbial taxa (**Supplementary Figure S4, Supplementary File S1**). For  
 219 example, ARGs that belong to beta-lactam resistance type (e.g., cephamycin, penam,  
 220 penem, and monobactam) were observed to appear together in Proteobacteria.



221

222 Enrichment analyses showed that ARGs are both environment-specific and  
223 host-specific (**Figure 4**). We found that the proportion of certain types of ARGs was  
224 higher in certain environments than in others. For example, rifamycin resistance genes  
225 were found enriched in the soil environment (with proportion of 0.1%) and enriched  
226 in the Actinobacteria (with proportion of 4.7%) (**Figure 4**). Rifamycin is an important  
227 antibacterial agent active against gram-positive bacteria, and it has a wide range of  
228 applications [30, 31]. The enrichment results were not surprising because  
229 *Actinomycetes* is a representative genus widely distributed in various soil  
230 environments, and its rifamycin resistance is compatible with its ability for rifamycin  
231 production [32-35].

232

## 233 **Evaluation of the ability for novel ARG identification using a recently annotated** 234 **ARG**

235 We further evaluated ONN4ARG's ability for novel ARG identification based on a  
236 newly annotated aminoglycoside resistance gene, GAR, which has been reported in a  
237 previous study by Böhm et al [7]. GAR is a recently reported aminoglycoside  
238 resistance gene, which is not present in CARD (v3.2.5), UniProt (as of December  
239 2022), DEEPARG-DB (v1.0.2), HMD-ARG-DB (as of December 2022), and  
240 ONN4ARG-DB. We searched the sequence of GAR with both DeepARG and  
241 HMD-ARG models, and the results showed that both of these models indicated it as  
242 non-ARG. We searched the sequence of GAR against all the sequences in  
243 ONN4ARG-DB using Diamond and did not find any homologous gene as well.  
244 Reassuringly, the prediction by ONN4ARG identified GAR as an ARG resistant to  
245 non-beta-lactam with high confidence (probability score = 100%). We should  
246 emphasize that though ONN4ARG predict GAR as non-beta-lactam and not as  
247 sub-type of aminoglycoside, ONN4ARG can give information about ancestors (or  
248 categories at higher levels) of the novel ARG, provide clues about novel knowledge.

249

## 250 **Functional verification of candidate novel resistance genes**

251 To identify promising putative novel resistance genes, we used four criteria: (i)  
252 remotely homologs to reference ARGs, (ii) prediction with high confidence, (iii)  
253 predicted to be single-type resistance, and (iv) the host is known. Despite the large  
254 number of candidate genes discovered by the ONN4ARG model, only 4,365 ARGs  
255 fulfilled all mentioned criteria (**Supplementary Table S7**).

256

257 To showcase the actual function of the predicted ARGs, we analyzed tens of ARGs  
258 belonging to the streptomycin resistance, and all of these ARGs have high confidence  
259 predicted by the ONN4ARG model. The experiment results showed that the  
260 Candi\_60363\_1 is one of the most promising ARG, which showed a high minimal  
261 inhibitory concentration (MIC) compared to negative control. Thus, we selected the  
262 Candi\_60363\_1 for further experimental validation (**Supplementary Table S8** and  
263 **S9**). Candi\_60363\_1, detected in *Streptococcus* in the oral environment, was predicted  
264 to confer resistance to streptomycin (belonging to aminoglycoside). One positive  
265 control from CARD (AHE40557.1, streptomycin resistance) was used for verification  
266 of the experimental system. All these genes were heterologously expressed in the *E.*  
267 *coli* BL21 (DE3) host by the induction of Isopropyl  $\beta$ -D-1-thiogalactopyranoside  
268 (IPTG) and tested for minimal inhibitory concentration (MIC) (**Figure 5A**). Results  
269 showed that the mRNA level of the genes increased with the addition of 1 mM IPTG  
270 compared with that without IPTG (**Figure 5B**), which verified the expression of the  
271 genes induced by IPTG. Furthermore, the MIC of the strain containing the positive  
272 control gene AHE40557.1 was more than 1,024  $\mu$ g/ml (**Supplementary Figure S5**),  
273 which is consistent with previous reports [36, 37]. This verified that our MIC  
274 measuring experimental system works well. Our results showed that the MIC of the  
275 strain containing Candi\_60363\_1 was significantly higher than the negative control  
276 containing no insert (Welch's t-test, one-tailed, P-value = 3.5e-3), which demonstrated  
277 the increased resistance to streptomycin of the novel candidate gene Candi\_60363\_1  
278 (**Figure 5C, Supplementary Figure S5**).

279

280 **Phylogeny and structure of Candi\_60363\_1**

281 There are remotely similarities between Candi\_60363\_1 and all known ARGs in the  
282 reference database, including aminoglycoside resistance genes. The InterPro search  
283 results showed the protein family matching to Candi\_60363\_1 is IPR007530, which is  
284 also known as aminoglycoside 6-adenylyltransferase that confers resistance to  
285 aminoglycoside antibiotics. Then, we used BLAST to search homologs of  
286 Candi\_60363\_1 from the NCBI non-redundant protein database. The BLAST result  
287 showed that there are 44 homologs with sequence identity greater than 80%, and they  
288 are from various organisms (**Supplementary Table S10**), such as *Streptococcus*  
289 *oralis*, *Peptoniphilus lacrimalis* DNF00528, and *Mycobacteroides abscessus* subsp.  
290 *Abscessus*. Considering that Candi\_60363\_1 is harbored by distantly related species,  
291 it obviously has mobility. Notably, the most similar protein of Candi\_60363\_1 from  
292 the NCBI non-redundant protein database (87.5% identity, SHZ78752.1) is also  
293 annotated as aminoglycoside adenylyltransferase (**Supplementary Table S10**). Taken  
294 together, Candi\_60363\_1 is highly likely to be an ARG that confers resistance to  
295 aminoglycoside antibiotics.

296

297 Aminoglycoside modifying enzymes are the most clinically important resistance  
298 mechanism against aminoglycosides [38]. They are divided into three enzymatic  
299 classes, namely, aminoglycoside N-acetyltransferase (AAC), O-nucleotidyltransferase  
300 (ANT), and O-phosphotransferase (APH). We investigated the phylogenetic  
301 relationship between Candi\_60363\_1 and the known aminoglycoside modifying  
302 enzymes. The phylogenetic tree of Candi\_60363\_1 and related proteins (**Figure 6A**)  
303 shows that Candi\_60363\_1 is clearly separated from the known aminoglycoside  
304 modifying enzymes and is located among proteins mostly annotated as  
305 aminoglycoside adenylyltransferase. Phylogenetic analysis indicated its evolutionarily  
306 close relationships with known aminoglycoside adenylyltransferase.

307

308 Protein structure prediction results confirmed the anti-microbial functionality of  
309 Candi\_60363\_1. The optimal Candi\_60363\_1-streptomycin complex structure and the  
310 corresponding interaction details are described in **Figure 6B**. The optimal binding

311 affinity between the Candi\_60363\_1 and streptomycin is -7.7 kcal/mol  
 312 (Supplementary Table S11), which is 1.6 kcal/mol lower than the negative control.  
 313 From wet-lab experiments, phylogenetic analysis, and protein structure docking, we  
 314 consider that Candi\_60363\_1 predicted by ONN4ARG is highly likely a real ARG  
 315 gene.

316

## 317 Discussion

318 In this study, we proposed an ontology-aware deep learning method, ONN4ARG, for  
 319 the detection and understanding of ARGs. To complement ONN4ARG for ARG  
 320 mining applications, we have also created a custom ARG database, ONN4ARG-DB,  
 321 that contains 28,396 well-curated ARGs. The application of ONN4ARG uncovered  
 322 120,726 ARGs from microbiome samples, out of which 42,848 are novel, which  
 323 substantially expands the existing ARGs repositories.

324

325 The novelty of this work is in three contexts. First, ONN4ARG has the potential for  
 326 detection of remotely homologous ARGs and thus generates a more comprehensive  
 327 set of ARGs. The ability of ONN4ARG to identify remotely homologs allows more  
 328 accurate prediction. The antibiotic resistance ontology used in the ONN4ARG model  
 329 consists of four levels and more than 100 resistance subtypes (i.e., terms in the most  
 330 informative level on the ontology), which substantially expand the classification space  
 331 of current tools (e.g., 30 types supported for DeepARG and 15 types supported for  
 332 HMD-ARG). Therefore, ONN4ARG greatly reduces false negatives and offers a  
 333 powerful approach for accurate and comprehensive profiling of ARGs.

334

335 Second, it enabled the comprehensive enrichment analysis of ARGs, species-wise and  
 336 environment-wise. The environment-specific and host-specific enrichment of ARGs  
 337 may be caused by specific bacteria evolving to possess specific types of ARGs in  
 338 response to specific environments, and horizontal gene transfer may be one of the  
 339 mediating pathways of this process. For example, one published study has reported

340 that *Amycolatopsis* in the soil environment produces rifamycin and thus gains  
341 ecological advantages over other bacteria [32].

342

343 Third, our study demonstrates the importance and potential of complementing the  
344 computational work with wet-lab experimental validation of gene function. Functional  
345 verification of a novel streptomycin resistance gene (i.e., Candi\_60363\_1) with  
346 wet-lab experiments demonstrated the ability of the ONN4ARG model for novel ARG  
347 discovery. Moreover, phylogenetic analysis and protein structure docking further  
348 confirmed that Candi\_60363\_1 is highly likely to be an ARG that confers resistance  
349 to aminoglycoside antibiotics. Another validation of a recently annotated ARG (i.e.,  
350 GAR) also indicated the ability of the ONN4ARG model for novel ARG discovery.

351

## 352 **Conclusions**

353 We proposed an ontology-aware deep learning approach, ONN4ARG, which is  
354 superior to existing methods such as DeepARG in efficiency, accuracy, and  
355 comprehensiveness. It enables comprehensive ARG discovery. It has detected novel  
356 ARGs that are remotely homologous to existing ARGs. Whereas ONN4ARG has  
357 provided one of the most comprehensive profiles of ARGs, it could be further  
358 optimized. For more comprehensive ARG prediction, continuous improvement of  
359 curating ARG nomenclature and annotation databases is required. For novel ARG  
360 prediction, especially those belonging to entirely new ARG families, deep learning  
361 models might need to consider more information other than sequence alone, such as  
362 protein structure. We believe these efforts could lead to a holistic view about ARGs in  
363 diverse environments around the globe.

364

## 365 **Methods**

### 366 **Dataset**

367 The ARGs we used in this study for model training and testing were from the  
368 Comprehensive Antibiotic Resistance Database, CARD v3.0.3 [5, 6]. We also used

protein sequences from the UniProt (SwissProt and TrEMBL) database to expand our training dataset. First, genes with ARG annotations were collected from CARD (2,587 ARGs) and SwissProt (2,261 ARGs). Then, their close homologs (sequence identity > 90% and coverage > 98%) were collected from TrEMBL (23,728 homologous genes). These annotated and homologous ARGs made up our ARG dataset. The non-ARG dataset was made from non-ARG genes that had relatively low sequence similarities to ARG genes (sequence identity < 90% and bit-scores < alignment lengths) but not annotated as ARG genes in SwissProt (17,937 non-ARG genes). Finally, redundant genes with identical sequences were filtered out. As a result, our ARG gene dataset, namely, ONN4ARG-DB, contained 28,396 ARG genes and 17,937 non-ARG genes. The gene clustering of the 681 newly added ARGs in CARD v3.1.3 was performed using the MMseqs2 tool (version 10) with an identity of 90% and coverage of 98%. The ResFinder dataset was obtained in Jun 2022 from [https://bitbucket.org/genomicepidemiology/resfinder\\_db/src/master/](https://bitbucket.org/genomicepidemiology/resfinder_db/src/master/).

383

### 384 **Antibiotic resistance ontology**

385 The antibiotic resistance ontology was organized into an ontology structure, which  
386 contains four levels (**Figure 1A**). The root (first level) is a single node, namely, “arg”  
387 (**Supplementary Table S12**). There are 1, 2, 34, and 277 nodes from the first level to  
388 the fourth level, respectively. For instance, there are “beta-lactam” and  
389 “non-beta-lactam” in the second level, “acridine dye” and “aminocoumarin” in the  
390 third level, and “acriflavine” and “clorobiocin” in the fourth level.

391

### 392 **Framework of ONN4ARG**

#### 393 **ONN4ARG model**

394 Considering a query gene  $q$  represented by its protein sequence, as well as its  
395 potential resistance categories represented by the antibiotic resistance ontology  $O$ , to  
396 predict resistance categories  $\hat{y}_q$  of query gene  $q$ , we employed ontology-aware  
397 neural network to learn a mapping  $M$  from a set of base genes  $b \in S$  to their  
398 resistance categories  $\hat{y}_b = (y_b^1, y_b^2, y_b^3, y_b^4)$ . Here,  $S$  is the set of base genes (i.e.,

ONN4ARG-DB),  $y_b^1$  is the resistance category for base gene  $b$  in the first level of the antibiotic resistance ontology. Then, we apply  $M$  on  $q$  to determine the potential resistance categories of query gene.

$$\hat{y}_q = (y_q^i)_{1 \leq i \leq 4} = M(q)$$

## Feature encoding

The task of feature encoding is to abstract the homologous signal of a query gene. ONN4ARG takes homologous signals (e.g., identity, e-value, bit-score) between the protein sequence of query gene and protein sequences and profiles (i.e., position-specific scoring matrix) of base genes as features. The homologous signal abstraction works as following. First, a protein sequence library of base genes was made by using “makedb” function of Diamond software. Then, protein sequences of query genes and base genes were aligned by using “blastp” function of Diamond program (**Figure 1B**). Second, profile hidden Markov models (HMMs) of base genes were generated by using “HHblits” function of HH-suite3 software (version 3.2.0). Then, protein sequence of query genes and profile HMMs of base genes were aligned by using “HHblits” function of HH-suite3 software (**Figure 1B**). Third, these homologous signals were normalized (i.e., divided by alignment length) and saved as vectors. The vector sizes at the two-layers of feature embedding network are decided based on the number of sequences and profiles in the ONN4ARG-DB. The vector size of the sequence features is 25,868, and 9,564 for the profile HMMs features.

## Architecture of the ontology-aware neural network

PyTorch version 1.7.1 was used for generating the ONN model. The architecture of the ontology-aware neural network could be described in four functional layers, including feature embedding layer, residual layer, compress layer and ontology-aware layer (**Supplementary Figure S1**). Details about the four functional layers are available at **Supplementary File S1**.

## 427 **Training and testing**

428 We performed 4-fold cross-validation in the systematic evaluation of ONN4ARG  
429 model. In each fold, we divided the ONN4ARG-DB into training set and testing set,  
430 the training set contains 75% randomly selected genes from the ONN4ARG-DB,  
431 whereas the remaining 25% genes were selected as testing set. We create binary label  
432 vector for each protein sequence. If a protein sequence is annotated with a resistance  
433 type from the ontology, then we assign 1 to the type's position in the binary label  
434 vector. Otherwise, we assign 0.

435

## 436 **Masking threshold**

437 To simulate remotely homologous ARG genes in our experiments, homologous  
438 signals between the query protein and its close homologs with sequence identities  
439 greater than a threshold were masked as zeros (i.e., no signals). For instance, when the  
440 masking threshold of testing set equaled 0.4, homologous signals between the query  
441 protein (in the testing set) and its close homologs (in the training set) with sequence  
442 identities greater than 40% were masked as zeros. Occasionally, all homologs were  
443 masked for a query protein, and such query proteins were removed during testing  
444 (**Table 1**). For example, if query *X* had two homologs, *M* and *N*, and assuming the  
445 identity of *M* is 0.45 and the identity of *N* is 0.95, when the masking threshold of the  
446 testing set equaled 0.9, homologous signals between query *X* and homolog *N* were  
447 masked as zeros. When the masking threshold of the testing set equaled 0.4, query *X*  
448 was removed during testing (see **Table 1** for details).

449

## 450 **Other methods**

451 We used Diamond (version 0.9.0) [23] as the sequence-alignment tool for comparison.  
452 We used the same training and testing sets as in the ONN4ARG model to evaluate the  
453 sequence-alignment method. For queries in the testing set, we searched them against  
454 the training set. The target with the highest identity was defined as the closest  
455 homologous gene for each query. Then, we compared whether the actual annotation of  
456 the query was consistent with the annotation of its closest homologous gene to



457 evaluate the performance. DeepARG [12] is a newly developed tool that applies a  
 458 plain neural network (e.g., several fully connected layers) to predict ARGs. Here, we  
 459 reconstructed the DeepARG model with PyTorch by using the same architecture of  
 460 original DeepARG model, and used the same training and testing sets as in the  
 461 ONN4ARG model to train and test the DeepARG model. For queries in the testing set,  
 462 we used the reconstructed DeepARG model to predict their ARG annotations, and  
 463 compared whether the actual annotations were consistent with the predicted  
 464 annotations to evaluate the performance.

465

## 466 **Performance measures**

467 To assess the performance of ONN4ARG model and other methods, we used accuracy  
 468 measure with the following formula:

$$Accuracy = \frac{N_{corp}}{N_{pred}}$$

469 where  $N_{corp}$  is the number of correct predictions, and  $N_{pred}$  is the number of total  
 470 predictions. Notably, a prediction was defined to be correct if and only if all ARG  
 471 annotations (including ancestor annotations from ARG ontology) were correctly  
 472 predicted.

473

474 Furthermore, we used precision, recall, F1, AUROC, and AUPRC measures to assess  
 475 the performance of ONN4ARG model and other methods on each antibiotic resistance  
 476 type:

$$\begin{aligned} Precision(f) &= \frac{TP(f)}{TP(f) + FP(f)} \\ Recall(f) &= \frac{TP(f)}{TP(f) + FN(f)} \\ F1 &= \frac{2 \times Precision(f) \times Recall(f)}{Precision(f) + Recall(f)} \\ TPR(f) &= \frac{TP(f)}{TP(f) + FP(f)} \\ FPR(f) &= \frac{FP(f)}{FP(f) + TN(f)} \end{aligned}$$

477

where  $f$  represents one resistance type,  $TP(f)$  is the number of true positive predictions of resistance type  $f$ ,  $FP(f)$  is the number of false positive predictions of resistance type  $f$ ,  $TN(f)$  is the number of true negative predictions of resistance type  $f$ , and  $FN(f)$  is the number of false negative predictions of resistance type  $f$ . AUROC is the area under the  $TPR-FPR$  curve, and AUPRC is the area under the *Precision-Recall* curve.

#### **Taxonomy annotation**

Kraken2 (version 2.1.2) [29] program with default parameters was used to identify the host of gene contigs. Then, each ARG predicted by ONN4ARG was annotated according to the host of its gene contigs.

#### **Phylogenetic tree**

Protein sequences of the most closely related to Candi\_60363\_1 were collected using BLASTP with default parameters on the NCBI non-redundant protein database. The retrieved proteins, Candi\_60363\_1 and all aminoglycoside resistance proteins from ResFinder [8] ([https://bitbucket.org/genomicepidemiology/resfinder\\_db/src/master](https://bitbucket.org/genomicepidemiology/resfinder_db/src/master), last update Jun 2022), were aligned with ClustalW. The phylogenetic tree was calculated by MEGA [39] (v10) using the maximum likelihood algorithm with default parameters. The Interactive Tree of Life (iTOL v6) online tool [40] was used to prepare the phylogenetic tree for display.

#### **Protein model and docking**

Rosetta [41] was utilized to predict the protein structure using ab initio protein folding (<http://rosetta.bakerlab.org/>). The top five protein pockets were generated for docking calculation with Surface Topography of proteins [42] (CASTp). We used the Cambridge Structure Database [43] to generate streptomycin conformers. The 3D protein-ligand complexes were obtained from AutoDock Vina [44].

#### **ARG candidate gene expression plasmids construction and expression**

## 508 **verification**

509 The candidate resistance gene Candi\_60363\_1 and a positive control resistance gene  
510 AHE40557.1 were synthesized and subcloned into pUC19 vector, replacing *lacZ'*  
511 gene. The recombinant plasmids were then transformed into *E. coli* BL21 (DE3). The  
512 expression of resistance genes was induced by Isopropyl  $\beta$ -D-1-thiogalactopyranoside  
513 (IPTG) and verified by quantitative Real-time PCR (qRT-PCR) assay. Briefly, bacteria  
514 were grown in LB supplemented with ampicillin (100  $\mu$ g/ml) to OD600 of 0.5-0.6 by  
515 incubation at 37 °C with 220 rpm agitation, and the bacterial cultures were continued  
516 to grow until OD600 reached to 1.0 by adding or without adding 1 mM IPTG. The  
517 cells were harvested and total RNAs were purified using Bacterial RNA Extraction  
518 Kit (Vazyme Biotech). RNA reverse transcription was performed by using HiScript®  
519 II Q Select RT SuperMix for qPCR kit (Vazyme Biotech). qRT-PCR was performed  
520 by using SYBR Green Master Mix-High ROX Premixed (Vazyme Biotech) in a  
521 Stepone Plus system (Applied Biosystems). The *ldh* gene was used as internal control  
522 in all reactions. The relative fold changes were determined using the  $2^{-\Delta\Delta C_t}$  method, in  
523 which *ldh* was used for normalization. The protein sequences of the synthesized genes  
524 and the primer sequences for qRT-PCR were listed in **Supplementary Table S8** and  
525 **S9**.

526

## 527 **MIC determination**

528 Minimal inhibitory concentrations (MICs) of the antibiotic for the strains containing  
529 resistance genes were determined using E-tests (three repeats). Single colonies of the  
530 strains were incubated in 3 ml Mueller-Hinton (MH) medium with the addition of 100  
531  $\mu$ g/ml ampicillin at 35 °C for 4 hours, and the cells equal to  $1.5 \times 10^8$  cells/ml were  
532 spread on MH agar plates with the addition of 100  $\mu$ g/ml ampicillin and 1 mM IPTG,  
533 and streptomycin MIC Test Strips (Liofilchem®) were put in the middle of the plates.  
534 The plates were incubated at 35 °C for 18-24 hours, and the MICs were read. The  
535 strain containing empty vector was used as a negative control.

536

## 537 **Statistical test**

538 According to the normality of the data distribution verified by the Shapiro–Wilk test  
539 and Levene’s test, the ARG abundance data distribution is Gaussian and unequal  
540 variance. Thus, statistical test of the enrichment analysis was performed utilizing the  
541 Welch’s *t*-test (one-tailed), at the significance level of  $\alpha=0.005$  [45]. For all the tests,  
542 when the *P* value associated is lower than the significance level, one should reject the  
543 null hypothesis  $H_0$  (ARGs are not enriched in the environment or host), and accept  
544 the alternative hypothesis  $H_a$  (ARGs are enriched in the environment or host).

545

## 546 **Key Points**

- 547 • We developed an ontology-aware deep learning approach, ONN4ARG, which is  
548 superior to existing methods such as DeepARG in efficiency, accuracy.
- 549 • ONN4ARG has the potential for detection of remotely homologous ARGs and  
550 thus generates a more comprehensive set of ARGs.
- 551 • ONN4ARG enabled the comprehensive enrichment analysis of ARGs,  
552 species-wise and environment-wise.
- 553 • Our study demonstrates the importance and potential of complementing the  
554 computational work with wet-lab experimental validation of gene function.

555

## 556 **Declarations**

### 557 **Ethics approval and consent to participate**

558 Not applicable

559

### 560 **Consent for publication**

561 Not applicable

562

### 563 **Competing interests**

564 The authors declare that they have no competing interests.

565

## 566 **Data availability**

567 We collected metagenomic samples from several published studies [25, 26], and these  
568 samples are mainly from marine, soil and human associated environments. For human  
569 associated samples, including two gut groups (one group from Madagascar, i.e., GutM,  
570 the other group from Denmark, i.e., GutD), one oral group and one skin group (both  
571 oral and skin groups are from HMP project). Details and links about these samples are  
572 shown in **Supplementary Table S3**. The ONN4ARG-DB dataset could be accesses at:  
573 <https://github.com/HUST-NingKang-Lab/ONN4ARG>.

574

## 575 **Code availability**

576 All source codes have been uploaded to the website at:  
577 <https://github.com/HUST-NingKang-Lab/ONN4ARG>, and online web service can be  
578 accessed at: <http://onn4arg.xfcui.com/>.

579

## 580 **Authors' contributions**

581 K.N, X.C conceived and proposed the idea, and designed the study. Y.Z, C.C, Q.J,  
582 X.Z, X.C performed the experiments and analyzed the data. Y.Z, C.C, X.Z, K.N and  
583 X.C contributed to editing and proof-reading the manuscript. All authors read and  
584 approved the final manuscript.

585

## 586 **Acknowledgments**

587 We are grateful to Mingyue Cheng for insightful discussions. This work was partially  
588 supported by National Natural Science Foundation of China (Grant Nos. 81774008,  
589 81573702, 31871334 and 31671374), and the National Key R&D Program (Grant No.  
590 2018YFC0910502).

591

592 **Yuguo Zha** is a PhD student at the Huazhong University of Science and Technology.  
593 His research interests are in microbiome associated data mining, including gene  
594 mining, species mining, and pattern mining.

595

596 **Cheng Chen** is a PhD student at Shandong University. His research interests are in  
597 computational biology and bioinformatics.

598

599 **Qihong Jiao** is a PhD student at Shandong University. His research interests are in  
600 computational biology and bioinformatics.

601

602 **Xiaomei Zeng** is a professor at Department of Bioinformatics and Systems Biology,  
603 College of Life Science and Technology, Huazhong University of Science and  
604 Technology.

605

606 **Xuefeng Cui** is a professor at the School of Computer Science and Technology,  
607 Shandong University. His research interests are in computational biology and  
608 bioinformatics.

609

610 **Kang Ning** is a professor at Department of Bioinformatics and Systems Biology,  
611 College of Life Science and Technology, Huazhong University of Science and  
612 Technology.

613

# References

1. Thomas AM, Segata N. Multiple levels of the unknown in microbiome research. BMC Biol. 2019;17:48.
2. Li J, Jia H, Cai X, Zhong H, Feng Q, Sunagawa S, et al. An integrated catalog of reference genes in the human gut microbiome. Nat Biotechnol. 2014; 32:834-841.
3. Brogan DM, Mossialos E. A critical analysis of the review on antimicrobial resistance report and the infectious disease financing facility. Global and Health. 2016;12:8.
4. Goossens H, Ferech M, Stichele RV, Elseviers M. Outpatient antibiotic use in Europe and association with resistance: a cross-national database study. Lancet 2005;365:579-587.
5. Jia B, Raphenya AR, Alcock B, Waglechner N, Guo P, Tsang KK, et al. CARD 2017: expansion and model-centric curation of the comprehensive antibiotic resistance database. Nucleic Acids Res. 2017;45:D566-D573.
6. Alcock BP, Raphenya AR, Lau TTY, Tsang KK, Bouchard M, Edalatmand A, et al. CARD 2020: antibiotic resistome surveillance with the comprehensive antibiotic resistance database. Nucleic Acids Res. 2020;48:D517-D525.
7. Böhm M-E, Razavi M, Marathe NP, Flach C-F, Larsson DGJ. Discovery of a novel integron-borne aminoglycoside resistance gene present in clinical pathogens by screening environmental bacterial communities. Microbiome 2020;8:1-11.
8. Bortolaia V, Kaas RS, Ruppe E, Roberts MC, Schwarz S, Cattoir V, et al. ResFinder 4.0 for predictions of phenotypes from genotypes. J Antimicrob Chemother. 2020;75:3491-3500.
9. Bateman A, Martin MJ, O'Donovan C, Magrane M, Apweiler R, Alpi E, et al. UniProt: A hub for protein information. Nucleic Acids Res. 2015;43: D204-D212.
10. Rowe W, Baker KS, Verner-Jeffreys D, Baker-Austin C, Ryan JJ, Maskell DJ,

643 et al. Search Engine for Antimicrobial Resistance: a cloud compatible pipeline  
644 and web interface for rapidly detecting antimicrobial resistance genes directly  
645 from sequence data. PLoS One. 2015;10:e0133492.

646 11. Kleinheinz KA, Joensen KG, Larsen MV. Applying the ResFinder and  
647 VirulenceFinder web-services for easy identification of acquired antibiotic  
648 resistance and E. coli virulence genes in bacteriophage and prophage  
649 nucleotide sequences. Bacteriophage. 2014;4:e27943.

650 12. Arango-Argoty G, Garner E, Pruden A, Heath LS, Vikesland PJ, Zhang L.  
651 DeepARG: a deep learning approach for predicting antibiotic resistance genes  
652 from metagenomic data. Microbiome. 2018;6:23.

653 13. Davis JJ, Boisvert S, Brettin T, Kenyon RW, Mao C, Olson R, et al.  
654 Antimicrobial resistance prediction in PATRIC and RAST. Sci Rep.  
655 2016;6:27930.

656 14. Lakin SM, Kuhnle A, Alipanahi B, Noyes NR, Dean C, Muggli M, et al.  
657 Hierarchical Hidden Markov models enable accurate and diverse detection of  
658 antimicrobial resistance sequences. Commun Biol. 2019;2:294.

659 15. Doster E, Lakin SM, Dean CJ, Wolfe C, Young JG, Boucher C, et al.  
660 MEGARes 2.0: a database for classification of antimicrobial drug, biocide and  
661 metal resistance determinants in metagenomic sequence data. Nucleic Acids  
662 Res. 2020;48:D561-D569.

663 16. Li Y, Xu Z, Han W, Cao H, Umarov R, Yan A, et al. HMD-ARG: hierarchical  
664 multi-task deep learning for annotating antibiotic resistance genes.  
665 Microbiome 2021;9:40.

666 17. Gupta SK, Padmanabhan BR, Diene SM, Lopez-Rojas R, Kempf M, Landraud  
667 L, et al. ARG-ANNOT, a new bioinformatic tool to discover antibiotic  
668 resistance genes in bacterial genomes. Antimicrob Agents Chemother.  
669 2014;58:212-220.

670 18. Feldgarden M, Brover V, Haft DH, Prasad AB, Slotta DJ, Tolstoy I, et al.  
671 Validating the AMRFinder tool and resistance gene database by using  
672 antimicrobial resistance genotype-phenotype correlations in a collection of



673 isolates. Antimicrob Agents Chemother. 2019;63:e00483.

674 19. Inouye M, Dashnow H, Raven LA, Schultz MB, Pope BJ, Tomita T, et al.

675 SRST2: rapid genomic surveillance for public health and hospital

676 microbiology labs. Genome Med. 2014;6:90.

677 20. Rowe WPM, Winn MD. Indexed variation graphs for efficient and accurate

678 resistome profiling. Bioinformatics. 2018;34:3601-3608.

679 21. Altschul SF, Gish W, Miller WC, Myers EW, Lipman DJ. Basic Local

680 Alignment Search Tool. J Mol Biol. 1990;215:403-410.

681 22. Edgar RC. Search and clustering orders of magnitude faster than BLAST.

682 Bioinformatics. 2010;26:2460-2461.

683 23. Buchfink B, Xie C, Huson DH. Fast and sensitive protein alignment using

684 DIAMOND. Nat Methods. 2015;12:59-60.

685 24. Steinegger M, Meier M, Mirdita M, Vöhringer H, Haunsberger SJ, Söding J.

686 HH-suite3 for fast remote homology detection and deep protein annotation.

687 BMC Bioinform. 2019;20:1-15.

688 25. Sunagawa S, Coelho LP, Chaffron S, Kultima JR, Labadie K, Salazar G, et al.

689 Structure and function of the global ocean microbiome. Science.

690 2015;348:1261359.

691 26. Mitchell AL, Scheremetjew M, Denise H, Potter S, Tarkowska A, Qureshi M,

692 et al. EBI Metagenomics in 2017: enriching the analysis of microbial

693 communities, from sequence reads to assemblies. Nucleic Acids Res.

694 2018;46:D726-D735.

695 27. Hyatt D, Chen GL, LoCascio PF, Land ML, Larimer FW, Hauser LJ. Prodigal:

696 prokaryotic gene recognition and translation initiation site identification. BMC

697 Bioinform. 2010;11:119.

698 28. Li B, Yang Y, Ma L, Ju F, Guo F, Tiedje JM, et al. Metagenomic and network

699 analysis reveal wide distribution and co-occurrence of environmental

700 antibiotic resistance genes. ISME J. 2015;9:2490-2502.

701 29. Wood DE, Lu J, Langmead B. Improved metagenomic analysis with Kraken 2.

702 Genome Biol. 2019;20:1-13.

- 703 30. Qi F, Lei C, Li F, Zhang X, Wang J, Zhang W, et al. Deciphering the late steps  
704 of rifamycin biosynthesis. *Nat Commun.* 2018;9:2342.
- 705 31. Floss HG, Yu T-W. Rifamycin-mode of action, resistance, and biosynthesis.  
706 *Chem Rev.* 2005;105:621-632.
- 707 32. Yao Y, Zhang W, Jiao R, Zhao G, Jiang W. Efficient isolation of total RNA  
708 from antibiotic-producing bacterium *Amycolatopsis mediterranei*. *J Microbiol*  
709 *Methods.* 2002;51:191-195.
- 710 33. Wilson MC, Gulder TAM, Mahmud T, Moore BS. Shared biosynthesis of the  
711 saliniketals and rifamycins in *Salinispora arenicola* is controlled by the  
712 sare1259-encoded cytochrome P450. *J Am Chem Soc.*  
713 2010;132:12757-12765.
- 714 34. Saxena A, Kumari R, Mukherjee U, Singh P, Lal R. Draft genome sequence of  
715 the rifamycin producer *amycolatopsis rifamycinica* DSM 46095. *Genome*  
716 *Announc.* 2014;2:e00662.
- 717 35. Huang H, Lv J, Hu Y, Fang Z, Zhang K, Bao S. *Micromonospora rifamycinica*  
718 sp. nov., a novel actinomycete from mangrove sediment. *Int J Syst Evol*  
719 *Microbiol.* 2008;58:17-20.
- 720 36. Pinto-Alphandary H, Mabilat C, Courvalin P. Emergence of aminoglycoside  
721 resistance genes *aadA* and *aadE* in the genus *Campylobacter*. *Antimicrob*  
722 *Agents Chemother.* 1990;34:1294-1296.
- 723 37. Holden MTG, Hauser H, Sanders M, Ngo TH, Cherevach I, Cronin A, et al.  
724 Rapid evolution of virulence and drug resistance in the emerging zoonotic  
725 pathogen *Streptococcus suis*. *PLoS One.* 2009;4:e6072.
- 726 38. Ramirez MS, Nikolaidis N, Tolmasky M. Rise and dissemination of  
727 aminoglycoside resistance: the *aac(6')-Ib* paradigm. *Front Microbiol.*  
728 2013;4:121.
- 729 39. Kumar S, Stecher G, Li M, Knyaz C, Tamura K. MEGA X: molecular  
730 evolutionary genetics analysis across computing platforms. *Mol Biol Evol.*  
731 2018;35:1547-1549.
- 732 40. Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new

733 developments. *Nucleic Acids Res.* 2019;47:W256-W259.

734 41. Rohl CA, Strauss CEM, Misura KMS, Baker D. Protein structure prediction  
735 using Rosetta. *Methods Enzymol.* 2004;383:66-93.

736 42. Tian W, Chen C, Lei X, Zhao J, Liang J. CASTp 3.0: Computed Atlas of  
737 Surface Topography of Proteins. *Nucleic Acids Res.* 2018;46:W363–W367.

738 43. Cole JC, Korb O, McCabe P, Read MG, Taylor R. Knowledge-based  
739 conformer generation using the cambridge structural database. *J Chem Inf*  
740 *Model.* 2018;58:615-629.

741 44. Trott O, Olson AJ. AutoDock Vina: improving the speed and accuracy of  
742 docking with a new scoring function, efficient optimization, and  
743 multithreading. *J Comput Chem.* 2010;31:455-461.

744 45. Benjamin DJ, Berger JO, Johannesson M, Nosek BA, Wagenmakers EJ, Berk  
745 R, et al. Redefine statistical significance. *Nat Hum Behav.* 2018;2:6-10.

746

## 747 **Figure Legends**

### 748 **Figure 1. Overview of the ONN4ARG model and its use for novel ARG discovery.**

749 (A) The antibiotic resistance gene ontology contains four levels. The root (first level)  
750 is a single node, namely, “arg”. There are 1, 2, 34, and 277 nodes from the first level  
751 to the fourth level, respectively. (B) The feature encoding procedure of ONN4ARG  
752 model. The sequence alignment features and profile HMMs features are encoded by  
753 calling Diamond and HHblits. (C) The architecture of the ontology-aware neural  
754 network could be described in four functional layers, including feature embedding  
755 layer, residual layer, compress layer and ontology-aware layer. The ontology-aware  
756 layer is a partially connected layer which encourage annotation predictions satisfying  
757 the ontology rules (i.e., the ontology tree structure). Specially, weight between nodes  
758 with relationship (e.g., parent and child) satisfying the ontology rules would be saved  
759 in the partially connected layer, and weights between irrelevant nodes would be  
760 masked. (D) Building the dataset for training and testing, and applying ONN4ARG  
761 model on metagenomic samples to discover candidate novel ARGs.

762

### 763 **Figure 2. Systematic evaluation and comparison between sequence-alignment,**

764 **DeepARG, and ONN4ARG.** (A) The accuracy of three models on ARG  
765 classification was assessed using a box plot. Diamond was used for  
766 sequence-alignment; significance test was based on the *t*-test. (B) The precision and  
767 recall of DeepARG and ONN4ARG on ARG classification for each antibiotic  
768 resistance type. The masking threshold of testing set equaled 0.4 (details of masking  
769 threshold are provided in **Methods**).

770

### 771 **Figure 3. Broad-spectrum profile of predicted ARGs among diverse**

772 **environments.** (A) The proportion of predicted ARGs for different protein sequence  
773 lengths. (B) The abundance ratio of predicted ARGs among diverse environments.  
774 Abundance ratio was defined as the number of ARGs divided by the number of total  
775 genes. (C) The proportion of predicted ARGs for different sequence identities among

776 diverse environments. **(D)** Number of genes in ONN4ARG-DB (left), predicted  
 777 homologous ARGs (middle), and predicted novel ARGs (right) for various resistance  
 778 types. The horizontal axis indicates the logarithmic number of genes, and the vertical  
 779 axis indicates different antibiotic resistance types. We collected metagenomic samples  
 780 from several published studies; these samples were mainly from “marine,” “soil,” and  
 781 “human” environments. Human-associated samples consisted of two gut groups (one  
 782 group from Madagascar, i.e., GutM; the other group from Denmark, i.e., GutD), one  
 783 oral group, and one skin group (both oral and skin groups were from the HMP  
 784 project).

785

786 **Figure 4. Enrichment of predicted ARGs among diverse environments and hosts.**

787 **(A)** Relative abundance and enrichment of ARGs among diverse environments.  
 788 Abundance ratio was defined as the number of ARGs divided by the number of total  
 789 genes. **(B)** Proportion and enrichment of ARGs among diverse hosts. Colors indicate  
 790 the proportion of ARGs for each phylum and resistance type. Results for the most  
 791 abundant five phyla that carry ARGs are shown. “+”: P-value < 0.005 (Welch’s *t*-test,  
 792 one-tailed).

793

794 **Figure 5. Functional validation of a predicted candidate novel ARG. (A)**

795 diagram showing the procedure of heterologous expression and functional analysis of  
 796 the predicted candidate ARG in the *E. coli* BL21 (DE3) host. **(B)** Gene expression  
 797 validation of the predicted candidate ARG. The vertical axis indicates the relative  
 798 mRNA level. **(C)** The MIC of the predicted candidate ARG and negative control. The  
 799 vertical axis indicates the MIC value. The MIC of the predicted candidate novel ARG  
 800 is significantly higher than the negative control (Welch’s *t*-test, one-tailed, P-value =  
 801 3.5e-3).

802

803 **Figure 6. Phylogenetic analysis and structure investigation of Candi\_60363\_1. (A)**

804 Phylogenetic tree of aminoglycoside resistance enzymes, Candi\_60363\_1, and its  
 805 homologs from the NCBI non-redundant protein database. ANT:

806 O-nucleotidyltransferase, AAC: N-acetyltransferase, APH: O-phosphotransferase,  
 807 AADT: aminoglycoside adenylyltransferase. (B) The optimal  
 808 Candi\_60363\_1-streptomycin complex structure (left), and the local interactions  
 809 between ligand and neighboring residues (right). The docking experiment indicates  
 810 there are six neighboring residues whose distances are less than three angstroms.

811

812 **Table 1. Accuracy comparison of sequence-alignment, DeepARG and ONN4ARG**  
 813 **based on different masking threshold of testing set.**

814

815 **Table 2. Evaluation results of ONN4ARG for ARGs identification at different**  
 816 **masking threshold of testing set.**

817

818 **Table 3. Evaluation results of DeepARG for ARGs identification at different**  
 819 **masking threshold of testing set.**

820

## 821 **Supplementary Materials**

### 822 **Supplementary Figure S1. The architecture of the ontology-aware neural**

823 **network.** (A) The architecture of the ontology-aware neural network could be  
824 described in four functional layers, including feature embedding layer, residual layer,  
825 compress layer and ontology-aware layer. The ontology-aware layer is a partially  
826 connected layer which encourage annotation predictions satisfying the ontology rules  
827 (i.e., the ontology tree structure). Specially, weight between nodes with relationship  
828 (e.g., parent and child) satisfying the ontology rules would be saved in the partially  
829 connected layer, and weights between irrelevant nodes would be masked. (B) The  
830 weight matrix derived from the antibiotic resistance ontology and the ontology-aware  
831 layer.

832

### 833 **Supplementary Figure S2. The number of pan and core ARG types among**

834 **various environments, and gene mobility analysis for predicted ARGs.** (A) The

835 number of pan and core ARG types change as more groups are included. For core/pan  
836 counts, we only counted ARG types with the relative abundance ratio greater than  
837  $1e-4$ . The pan ARGs refer to the ARG types that are included in any environments.

838 The core ARGs refer to the ARG types that are included in all environments. (B) The

839 venn diagram shows the ARG types relationship among marine, soil and gut groups.

840 (C) The venn diagram shows the ARG types relationship among gut, oral and skin

841 groups. (D) The distribution of acquired and intrinsic ARGs in various environments.

842 (E) The line regression analysis indicates no significant correlation ( $P > 0.05$ )

843 between the abundances of MGEs and ARGs. The horizontal axis indicates the

844 abundance ratio of predicted ARGs and the vertical axis indicates the abundance ratio

845 of MGEs. Each point represents a group.

846

### 847 **Supplementary Figure S3. The host range of all classified ARGs and the**

848 **resistance composition of the most abundant 20 genera.** (A) The Sankey diagram

849 shows the host composition and distribution of all classified ARGs (the most abundant

20 genera carrying ARGs were used for display). **(B)** The bar chart indicates the diversity and relative abundance of ARGs for the most abundant 20 genera carrying ARGs.

**Supplementary Figure S4. The network analysis revealing the co-occurrence patterns among ARG types and microbial taxa, the nodes were represented by pie charts which shows the taxonomic compositions of ARG types.** A connection represents a strong (Spearman's  $\rho > 0.8$ ) and significant ( $P$ -value  $< 0.01$ ) correlation. The size of each node is proportional to the number of connections, that is, the degree.

**Supplementary Figure S5. The MIC experiment for predicted candidate ARG (top), negative control (middle) and positive control (bottom).** The MIC values are tested for three repeats.

**Supplementary Table S1. Comparison of ONN4ARG and other methods for ARG identification on the verification set.**

**Supplementary Table S2. Evaluation of ONN4ARG and DeepARG on the ResFinder dataset.**

**Supplementary Table S3. Metagenomic samples using for resistance gene mining are collected from published studies.**

**Supplementary Table S4. The number of predicted ARGs by ONN4ARG that have protein domains with known catalytic activity and/or may bind to the antimicrobials they are predicted to elicit resistance against.**

**Supplementary Table S5. Data distribution during the pipeline of ARGs prediction.**



880 **Supplementary Table S6. The hosts of predicted ARGs at different taxonomic**  
 881 **level.**

882

883 **Supplementary Table S7. The predicted ARGs which fulfilling all mentioned**  
 884 **criteria.**

885

886 **Supplementary Table S8. Protein sequences of the synthesized genes.**

887

888 **Supplementary Table S9. Real-time PCR primer sequences.**

889

890 **Supplementary Table S10. The BLAST result of Candi\_60363\_1 when search**  
 891 **against the NCBI non-redundant protein database.**

892

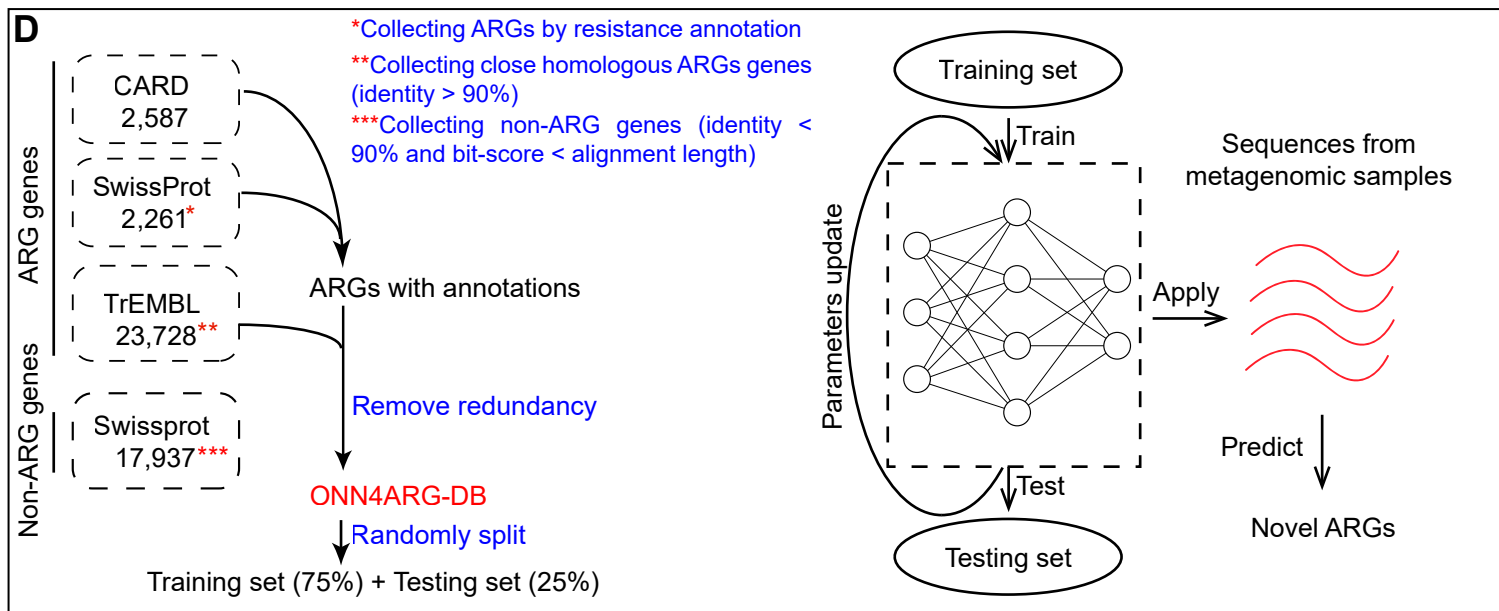
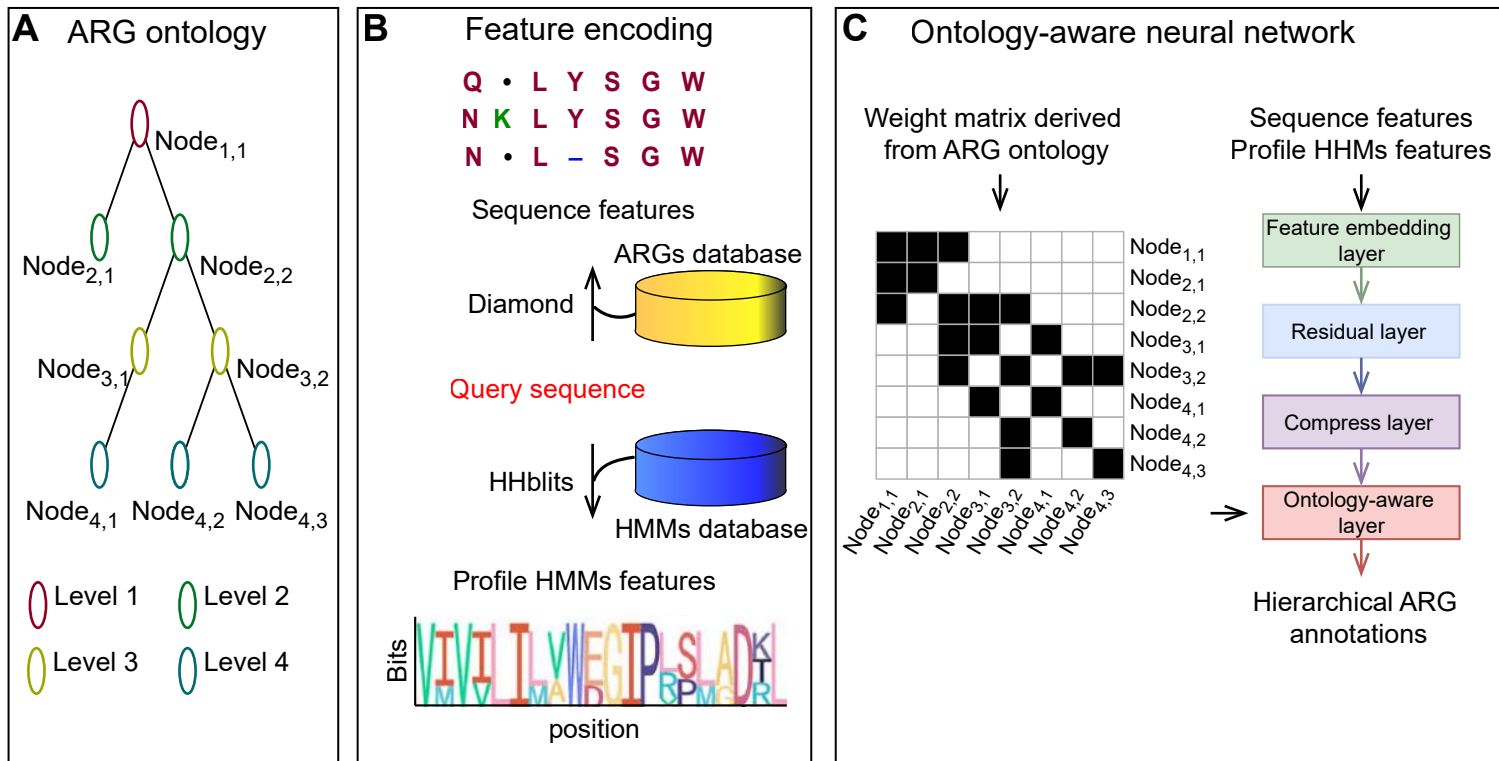
893 **Supplementary Table S11. The binding affinity of protein–ligand complexes**  
 894 **using the top five pockets.**

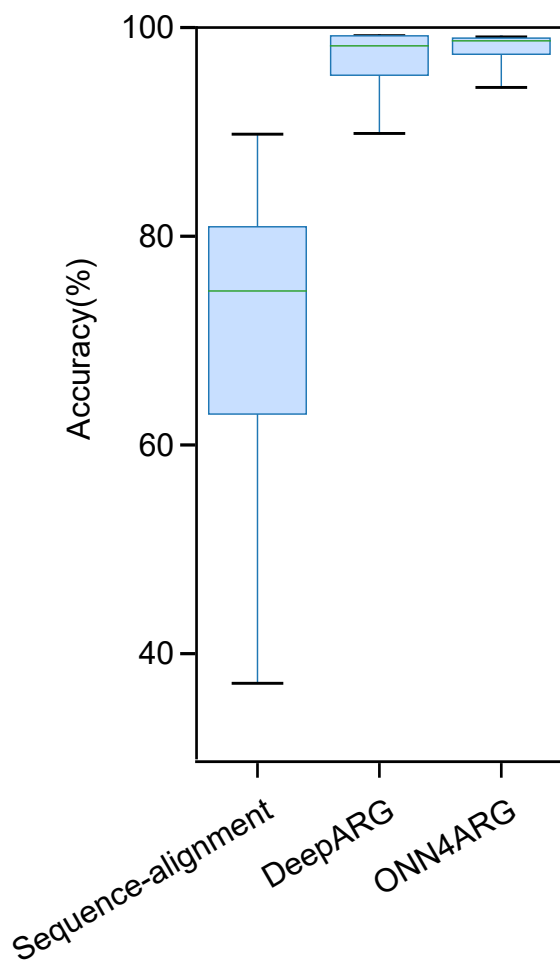
895

896 **Supplementary Table S12. The antibiotic resistance ontology used in the**  
 897 **ONN4ARG model.**

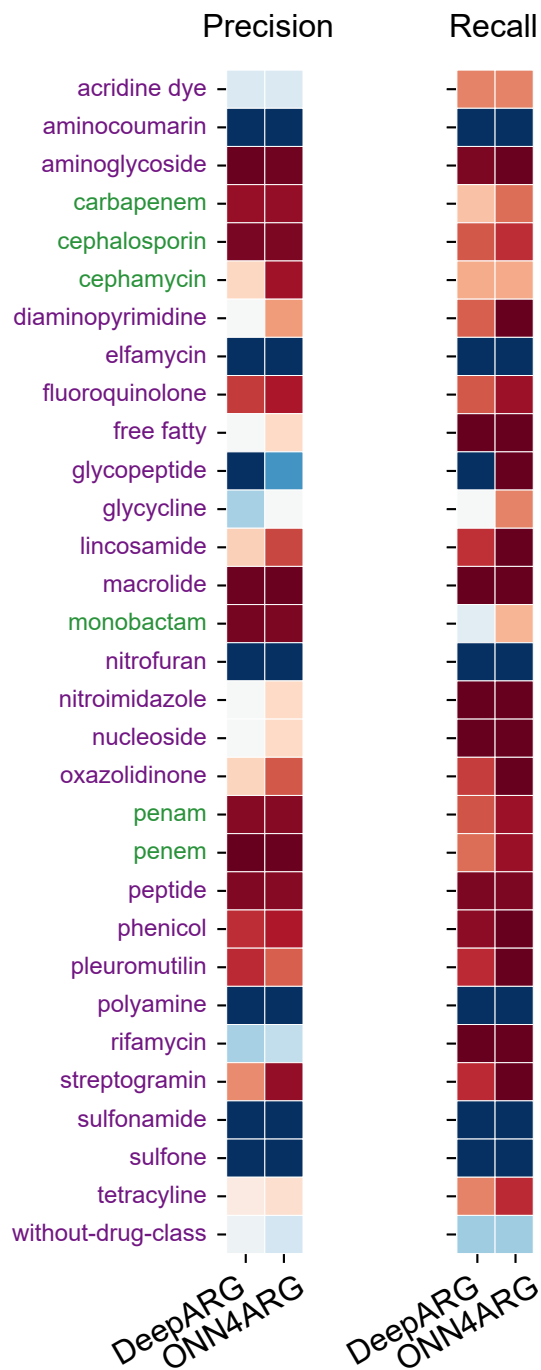
898

899 **Supplementary File S1. Supplemental information about experiments.**

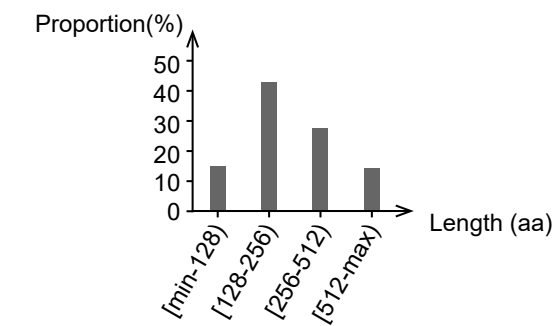


**A**

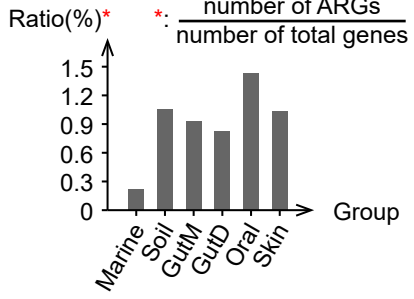
non-beta-lactam ●  
beta-lactam ●

**B**

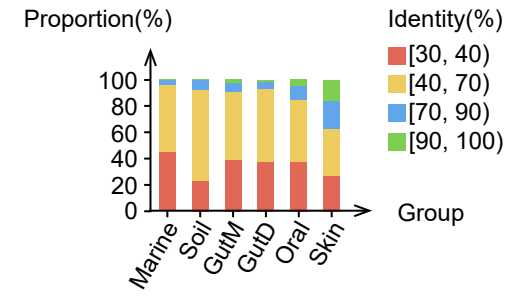
A



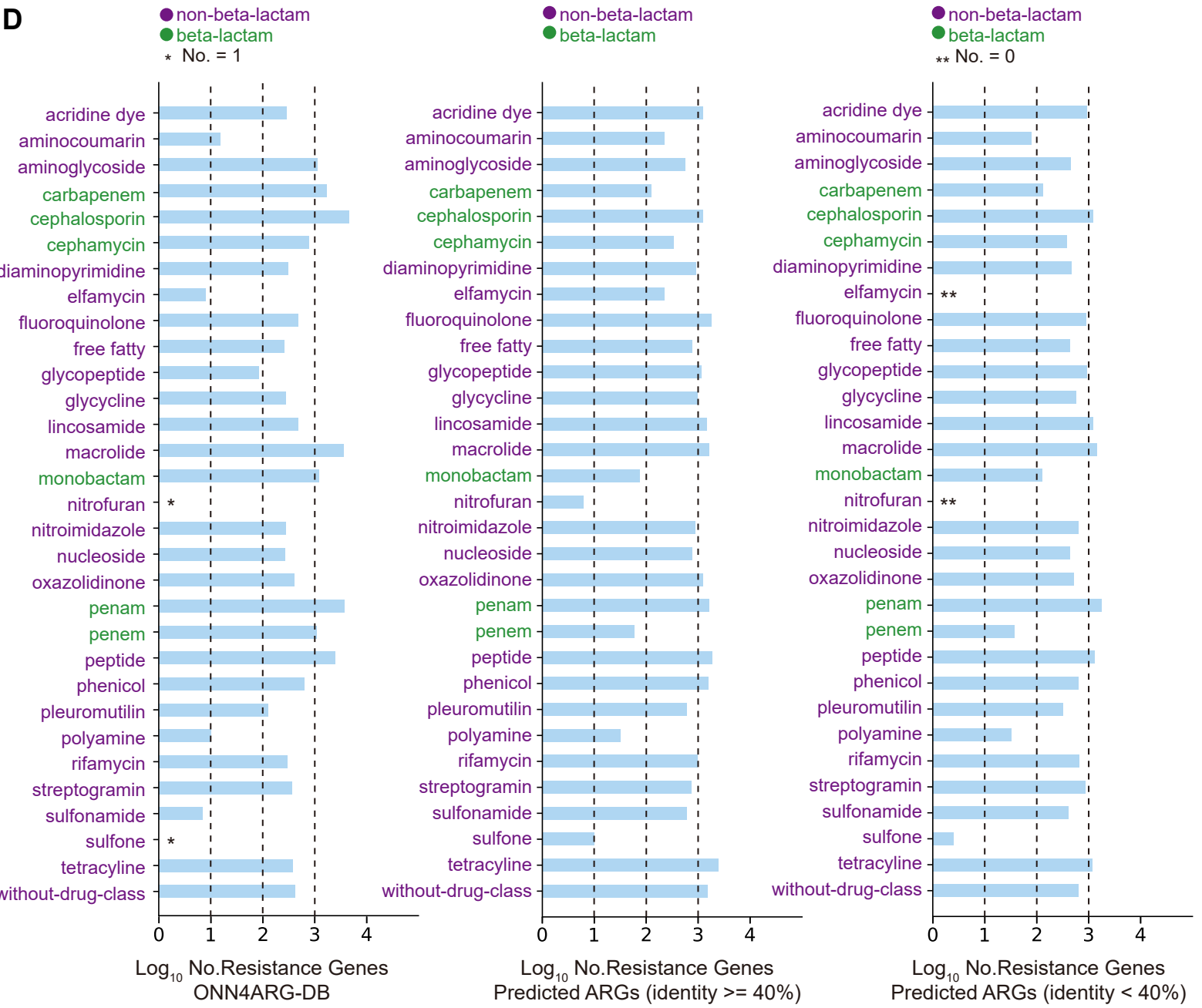
B

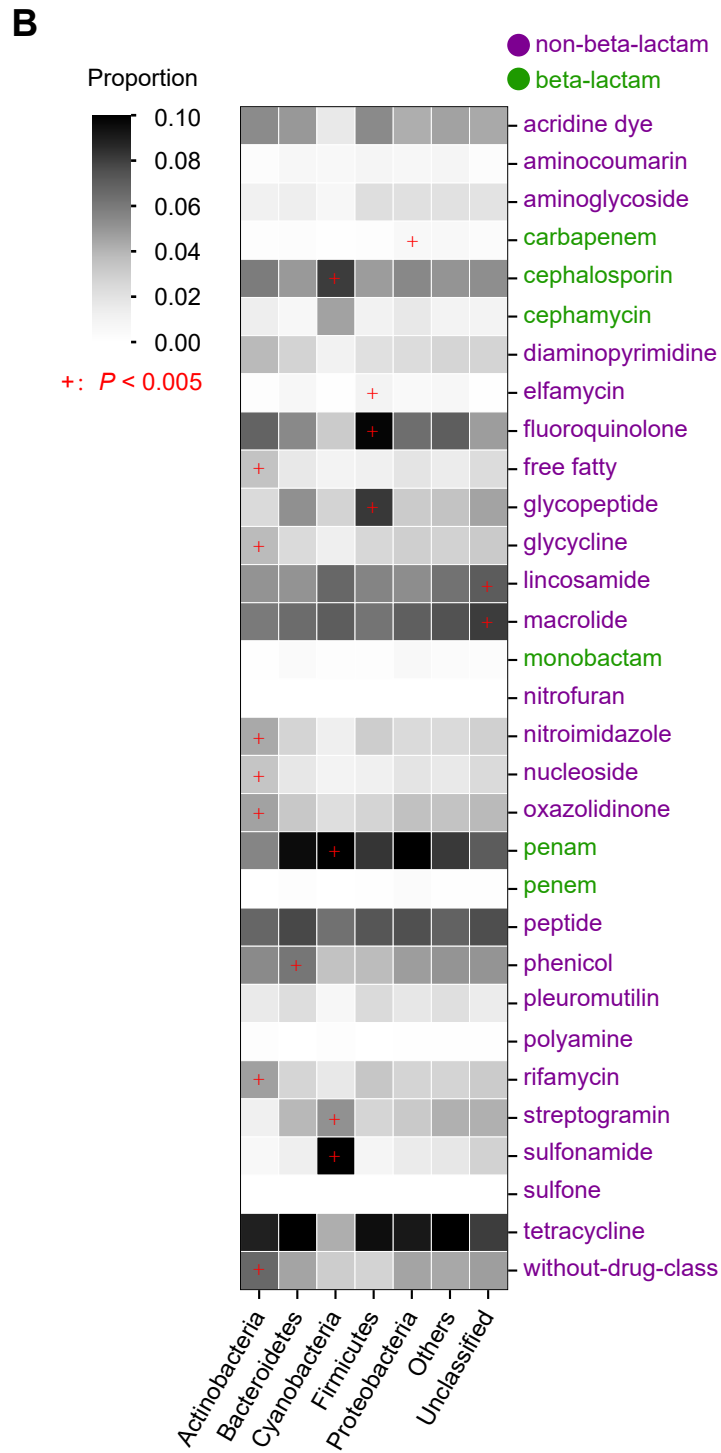
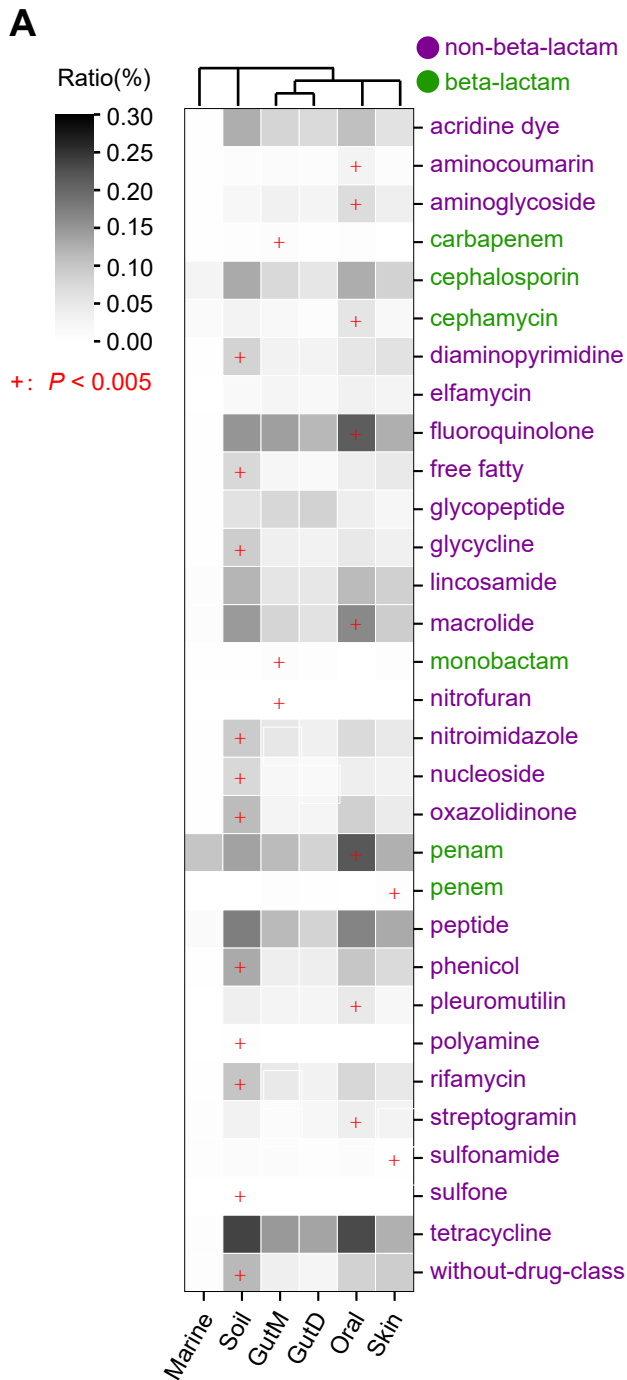


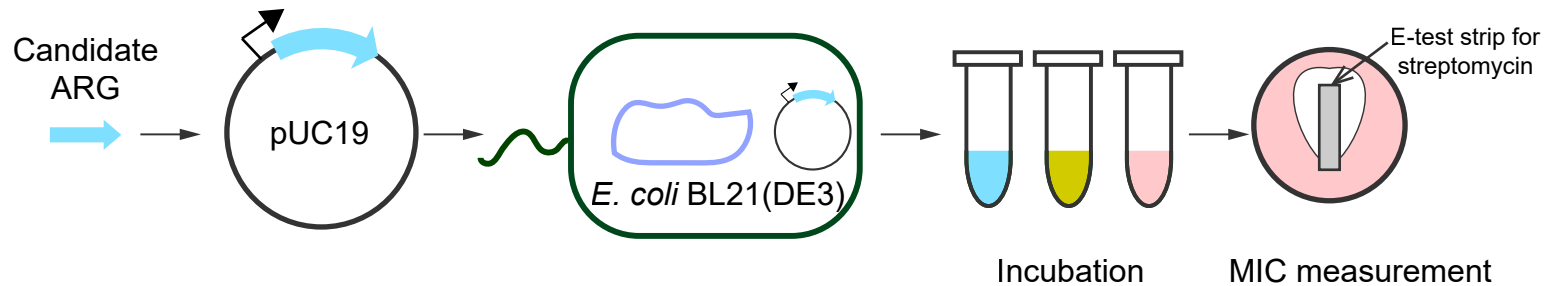
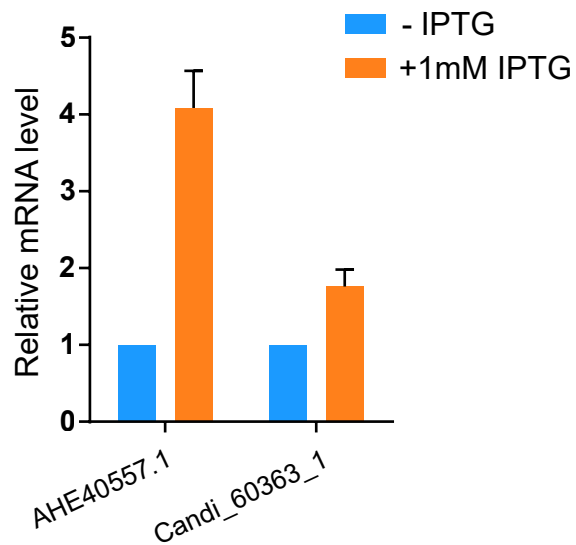
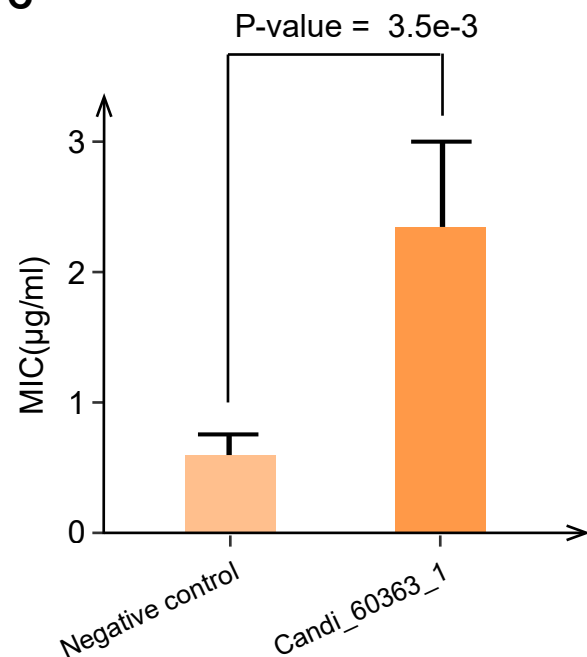
C



D





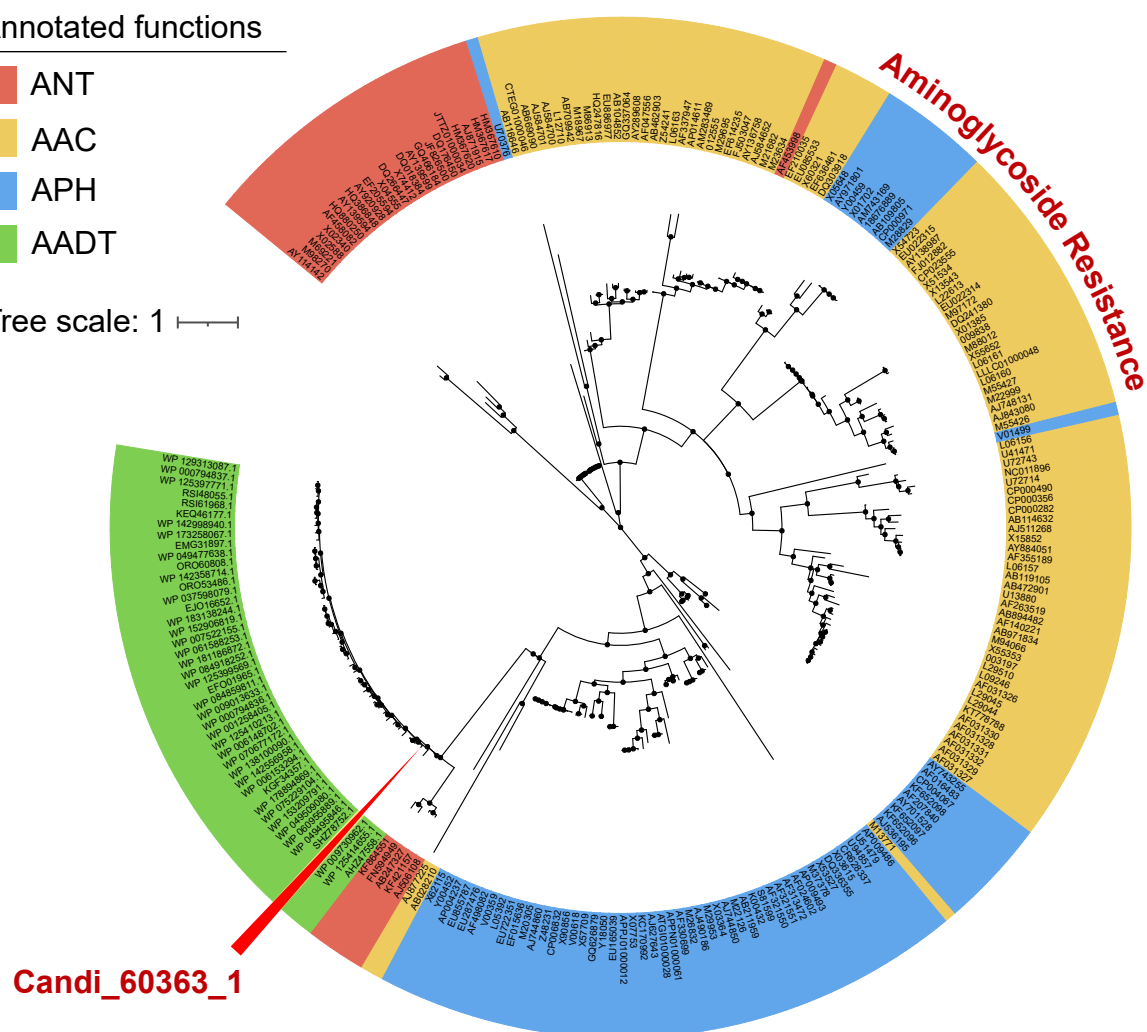
**A****B****C**

A

Annotated functions

- ANT
- AAC
- APH
- AADT

Tree scale: 1



B

