1   **Ontology-Aware Deep Learning Enables Novel Antibiotic Resistance**

2   **Gene Discovery Towards Comprehensive Profiling of ARGs**

3

4   Yuguo Zha[1], Cheng Chen[2], Qihong Jiao[2], Xiaomei Zeng[1, *], Xuefeng Cui[2, *], Kang

5   Ning[1, *]

6   [1]Key Laboratory of Molecular Biophysics of the Ministry of Education, Hubei Key

7   Laboratory of Bioinformatics and Molecular-imaging, Center of AI Biology,

8   Department of Bioinformatics and Systems Biology, College of Life Science and

9   Technology, Huazhong University of Science and Technology, Wuhan 430074, Hubei,

10  China

11  [2]School of Computer Science and Technology, Shandong University, Qingdao 266237,

12  Shandong, China

13  [*]Correspondence should be addressed to K.N (Email: ningkang@hust.edu.cn), X.C

14  (Email: xfcui@email.sdu.edu.cn) and X.Z (Email: xmzeng@hust.edu.cn)

15

## Abstract

17  Antibiotic resistance genes (ARGs) have emerged in pathogens and spread faster than

18  expected, arousing a worldwide concern. Current methods are suitable mainly for the

19  discovery of close homologous ARGs and have limited utility for discovery of novel

20  ARGs, thus rendering the profiling of ARGs incomprehensive. Here, an

21  ontology-aware deep learning model, ONN4ARG (http://onn4arg.xfcui.com/), is

22  proposed for the discovery of novel ARGs based on multi-level annotations.

23  Experiments based on billions of candidate microbial genes collected from various

24  environments show the superiority of ONN4ARG in comprehensive ARG profiling.

25  Enrichment analyses show that ARGs are both environment-specific and host-specific.

26  For example, resistance genes for rifamycin, which is an important antibacterial agent

27    active against gram-positive bacteria, are enriched in Actinobacteria and in soil

28    environment. Case studies verified ONN4ARG's ability for novel ARG discovery.

29    For example, a novel streptomycin resistance gene was discovered from oral

30    microbiome samples and validated through wet-lab experiments. ONN4ARG

31    provides a complete picture of the prevalence of ARGs in microbial communities as

32    well as guidance for detection and reduction of the spread of resistance genes.

33    **Keywords:** antibiotic resistance gene, ontology-aware, deep learning, novel ARG,

34    microbiome

35

## Introduction

37    With the development of metagenomics and next-generation sequencing, many new

38    microbial taxa and genes have been discovered, but different kinds of "unknowns"

39    remain. For instance, the microbes found in the human gut microbiome involve 25

40    phyla, more than 2,000 genera, and 5,000 species[1]. However, the functional diversity

41    of microbiomes has not been fully explored, and about 40% of microbial gene

42    functions remain to be discovered[2]. A typical example is the antibiotic resistance gene

43    (ARG), which is an urgent and growing threat to public health[3]. In the past few

44    decades, problems caused by antibiotic resistance have drawn the public's attention[4].

45    Antibiotic resistance in pathogens has been an increasing threat to human health over

46    the past decade, and it is widely accepted that antibiotic resistance development and

47    spread in microbes can be largely attributed to the abuse and misuse of antibiotics. A

48    direct correlation between antimicrobial use and the extent of antimicrobial resistance

49    has been reported[4]. Antimicrobial resistance genomic data is an ever-expanding data

50    source, with many new ARG families discovered in recent years[5,6]. The discovery of

51    resistance genes in diverse environments offers possibilities for early surveillance,

52    actions to reduce transmission, gene-based diagnostics, and, ultimately, improved

53    treatment[7].

54

55    Existing annotated ARGs have been curated manually or automatically for decades.

56    Presently, there are 2,979 annotated ARGs in the reference database CARD[5,6] (v3.1.2,

57    released in April 2021), 3,159 in the ResFinder database[8] (as of May 2021), and 2,675

58    in SwissProt[9] (as of May 2021). These annotated ARGs are categorized into antibiotic

59    resistance types, which are organized in an ontology structure (see **Methods**,

60    **Supplementary Figure 1**), in which higher-level ARG types cover lower-level ARG

61    types. For example, AHE40557.1 is annotated in the CARD database as a

62    streptomycin resistance gene, which belongs to a lower-level ARG type

63    aminoglycoside and a higher-level ARG type non-beta-lactam. Current ARG

64    databases are far from complete: though no ARG database contains more than 4,000

65    well-annotated ARGs, NCBI non-redundant database searches yielded more than

66    7,000 putative genes annotated with "antibiotic resistance" as of May 2021. Therefore,

67    we deemed that there is a large gap between the genes annotated in ARG databases

68    and the possible ARGs that already exist in general databases, not to mention ARGs

69    that are not yet annotated.

70

71    Many ARG prediction tools have been proposed in the past few years[8,10-20]. These

72    tools can generally be divided into two approaches. One approach is

73    sequence-alignment, such as BLAST[21], USEARCH[22], and Diamond[23], which uses

74    homologous genes to annotate unclassified genes. A confident prediction requires a

75    homolog with sequence identity greater than 80% in many programs, such as

76    ResFinder[8,11]. The other approach is deep learning, such as DeepARG[12] and

77    HMD-ARG[16], which uses neural network models to predict and annotate ARGs. The

78    input of deep learning approach can be bit-score (for DeepARG) or one-hot encoding

79    vector of protein sequence (for HMD-ARG).

80

81    Several limitations still preclude comprehensive profiling of antibiotic resistance

82    genes. A more comprehensive set of ARGs could be roughly defined as having more

83    ARGs in type and number with less false-positive entries, regardless of the homology

84    with known ARGs, and many of these ARGs could be experimentally validated.

85    Based on this definition, existing tools fall short in comprehensive profiling of ARGs.

86  First, existing tools are limited to a few types of ARGs due to the fact that the datasets

87  used for building models are specialized and therefore cannot reconstruct the

88  comprehensive profile of ARGs across various environments. For example,

89  HMD-ARG[16] identifies only 15 types of resistance genes, and PATRIC[13] is limited to

90  identifying ARGs encoding resistance to carbapenem, methicillin, and beta-lactam

91  antibiotics. Second, existing tools fall short in discovering novel ARGs, which usually

92  lack homology to known sequences in the reference databases. For instance, the gene

93  POCOZ1 (VraR) that confers resistance to vancomycin has a sequence identity of

94  only 24% to the homolog from the CARD[12]. Recognizing such remote homologs

95  requires the ability to perceive the correlation between the internal features of genes,

96  which is challenging for existing tools. Therefore, there is an urgent need for a new

97  approach to address these limitations.

98

99  Here, we propose an ontology-aware deep learning approach, ONN4ARG, which

100  allows comprehensive identification of ARGs. ONN4ARG is an ontology-aware

101  neural network model that employs a novel ontology-aware layer and generates

102  multi-level annotations of antibiotic resistance types (**Figure 1**). Systematic

103  evaluations show that the ONN4ARG model has a profound performance

104  improvement over state-of-the-art models such as DeepARG, especially for the

105  detection of remotely homologous ARGs. The application of ONN4ARG has

106  uncovered a total of 120,726 ARGs from the microbiome, which has greatly expanded

107  the existing ARG repositories. Enrichment analyses have confirmed the enrichment

108  patterns of ARG types across multiple environments, showing that ARGs are both

109  environment-specific and host-specific. For example, resistance genes for rifamycin,

110  which is an important antibacterial agent active against gram-positive bacteria, are

111  enriched in Actinobacteria and in soil environment. Case studies have also verified

112  the ability of ONN4ARG for novel ARG discovery. For example, a recently

113  experimentally validated ARG gene GAR[7], which is not in the CARD database, could

114  not be identified by DeepARG or HMD-ARG but was predicted by ONN4ARG. A

115  novel streptomycin resistance gene was also discovered by ONN4ARG from oral

116    microbiome data and validated through wet-lab experiments.

117

118    In summary, ONN4ARG is a comprehensive deep learning method for ARG

119    discovery, which provides a complete picture of the prevalence of ARGs in microbial

120    communities as well as guidance for detection and reduction of the spread of

121    resistance genes.

122

123    **Results**

124    **ONN4ARG model employs an ontology-aware neural network for ARG**

125    **identification and classification**

126    To address the large gap between the genes annotated in ARG databases and the

127    possible ARGs that already exist in general databases along with the ARGs that are

128    not yet annotated, we propose ONN4ARG, which is an ontology-aware neural

129    network model (**Figure 1a**), that could be used to predict ARGs in a comprehensive

130    manner. ONN4ARG takes similarities (e.g., identity, e-value, bit-score) between the

131    query gene sequence and ARG gene sequences and profiles (i.e., PSSM) as inputs and

132    predicts ARG annotations for the query gene. These sequence-alignment similarities

133    and profile-alignment similarities are pre-processed by calling Diamond[23] and

134    HHsearch[24]. ONN4ARG generates multi-level annotations of antibiotic resistance

135    types, which are compatible with the antibiotic resistance ontology structure. One

136    advantage of ONN4ARG over state-of-the-art models is that ONN4ARG employs a

137    novel ontology-aware layer that incorporates ancestor and descendent annotations to

138    enhance annotation accuracies. ONN4ARG outperforms existing models, including

139    DeepARG, with higher average accuracies and better generalization ability for unseen

140    data. To train and evaluate our ONN4ARG model and for rapid deployment of ARG

141    discovery in multiple contexts, we also built an ARG database (**Figure 1b**), namely,

142    ONN4ARG-DB, which comprises ARGs from CARD and UniProt (see **Methods**).

143

144    **Systematic evaluation and comparison**

145   ONN4ARG has high efficiency, high accuracy, and comprehensiveness for ARG

146   identification based on our systematic evaluation of ONN4ARG and comparison with

147   other models. The evaluation and comparison were based on ONN4ARG-DB, with

148   28,396 positive ARGs and 17,937 negatives, out of which 75% of the dataset was

149   randomly selected for training and the remaining 25% of the dataset was selected for

150   testing (see **Methods**).

151

152   We evaluated ONN4ARG's efficiency, accuracy, and comprehensiveness. As an

153   ontology-aware deep learning model, ONN4ARG is fast: it could complete ARG

154   identification for all genes in the testing dataset within four hours, which is equivalent

155   to one second per gene identification. As shown in **Figure 2a**, ONN4ARG was more

156   accurate for ARG identification (overall accuracy of 97.70%) compared to sequence

157   alignment (overall accuracy of 69.11%) and DeepARG (overall accuracy of 96.39%).

158   Moreover, ONN4ARG achieved an overall precision of 75.59% and an overall recall

159   of 89.93%, which were higher than DeepARG's overall precision of 68.30% and

160   overall recall of 77.84% (**Figure 2b**). It is natural that ONN4ARG could not

161   outperform DeepARG in all resistance types and this is exemplified by results on

162   pleuromutilin due to the small number of sequences for pleuromutilin in the

163   ONN4ARG-DB. In addition, for most of the resistance types that have adequate

164   number of sequences, ONN4ARG's results could achieve higher precision and recall.

165   Thus, with the accumulation of annotated ARG sequences, greater advantages of both

166   ONN4ARG-DB and ONN4ARG could be expected. Furthermore, ONN4ARG was

167   more comprehensive for ARG identification: there were 4,916 ARGs in the testing set

168   (with the masking threshold of testing equal to 0.4, see **Methods**), out of which 4,913

169   were identified by the ONN4ARG model, whereas DeepARG identified 4,906

170   (**Supplementary Table 1**).

171

172   ONN4ARG demonstrates an advantage over other methods in identification of

173   remotely homologous ARGs whose sequences are not similar to existing ARG

174   sequences (**Supplementary Tables 2 and 3**). In this context, when testing with only

175   remote homologs (i.e., the masking threshold of testing set equal to 0.4), ONN4ARG

176   achieves an accuracy of 94.26%, which is significantly improved from 89.85% of

177   DeepARG. When testing with all close and remote homologs (i.e., the masking

178   threshold of testing set equal to 1.0), both ONN4ARG and DeepARG achieved high

179   accuracies. These results validate ONN4ARG's significantly better generalization

180   abilities than sequence-alignment and DeepARG, which makes ONN4ARG especially

181   suitable for identification of remotely homologous ARGs and indicates ONN4ARG's

182   ability for novel ARG discovery (**Supplementary Tables 1–3**).

183

184   In summary, ONN4ARG has high efficiency, accuracy, and comprehensiveness for

185   ARG identification, and it possesses the ability for identification of remotely

186   homologous ARGs.

187

188   **Applications of ONN4ARG on metagenomic data**

189   We collected metagenomic samples from several published studies[25,26]. These samples

190   were mainly from "marine," "soil," and "human" environments. Human-associated

191   samples consisted of two gut groups (one group from Madagascar, i.e., GutM; the

192   other group from Denmark, i.e., GutD), one oral group, and one skin group (both oral

193   and skin groups were from the HMP project). Details about these samples are

194   provided in **Supplementary Table 4**. Then, genes were obtained by calling Prodigal[27].

195   The ONN4ARG model was used to predict whether these unclassified genes were

196   ARGs and their corresponding resistance types. In total, 120,726 ARGs were

197   identified from microbiome samples, many of which are novel, which greatly expands

198   the existing ARG repositories.

199

200   **Broad-spectrum profile of predicted ARGs among diverse environments**

201   We first investigated the proportion of predicted ARGs for different sequence lengths.

202   The distribution shows that about half of the predicted ARGs have a length of

203   128–256 amino acid residues (**Figure 3a**). We found that human-associated

204   microbiome samples carry a higher abundance of ARGs, especially for the oral group,

205 in which more than one resistance gene could be observed out of a hundred genes on

206 average (**Figure 3b**, **Supplementary Table 5**).

207

208 For ARGs detected in samples from all environments, we found that about a third of

209 them (42,848 out of all 120,726 ARGs) had sequence identity of less than 40% to

210 their homologs in the ONN4ARG-DB (**Figure 3c**). We define these ARGs as novel

211 ARGs, which have low sequence identities when aligned to their homologs in the

212 reference database (i.e., ONN4ARG-DB). For example, we found 45% of predicted

213 ARGs in the marine group belonged to novel ARGs (**Figure 3c**).

214

215 In total, 31 ARG types were detected in these various environments (**Figure 3d**,

216 **Supplementary Figure 2**). The number of predicted ARG sequences for different

217 types varied greatly (**Figure 3d**), from a few (i.e., nitrofuran) to thousands (i.e.,

218 fluoroquinolone). In general, fluoroquinolone and tetracycline resistance genes were

219 more abundant than other types (**Figure 3d**). As expected, these abundant ARGs were

220 usually associated with the antibiotics used extensively in human medicine or

221 veterinary medicine, including growth promotion[28]. Novel ARG detection indicates

222 the unique ability of ONN4ARG in novel ARG discovery and ARG abundance

223 profiling in various environments, which would help researchers to better understand

224 the prevalence of antibiotic resistance genes.

225

226 **Enrichment of predicted ARGs among diverse hosts and environments**

227 Rapid deciphering of potential antimicrobial-resistant pathogens is necessary for

228 effective public health monitoring. The host-tracking of ARGs allows for accurate

229 identification of pathogens. Therefore, we conducted Kraken2[29] analysis to track the

230 hosts of these predicted ARGs. Results showed that there are 949 genera, each genus

231 carries at least one type of ARG (**Supplementary Table 6**). The host composition and

232 distribution of all classified ARGs for the most abundant 20 genera are displayed in

233 **Supplementary Figure 3**. The host distribution shows that these predicted ARGs are

234 primarily affiliated with Proteobacteria (38.2%), including *Candidatus Pelagibacter*,

235  *Pseudomonas*, *Bradyrhizobium*, and *Escherichia* (**Supplementary Figure 3**). The

236  most abundant ARGs carried by the 20 genera were resistance types of

237  fluoroquinolone, macrolide, peptide, penam, and tetracycline, accounting for about

238  half of the total detected ARGs (**Supplementary Figure 3**). We used network

239  inference based on strong (Pearson's correlation ρ > 0.8) and significant (P-value <

240  0.01) correlations to investigate the co-occurrence patterns among ARG types and

241  microbial taxa (**Supplementary Figure 4**, **Supplementary Note**). The co-occurrence

242  network indicated the co-occurrence patterns between ARGs and microbial taxa. For

243  example, ARGs that belong to beta-lactam resistance type (e.g., cephamycin, penam,

244  penem, and monobactam) were observed to appear together in Proteobacteria.

245

246  Enrichment analyses showed that ARGs are both environment-specific and

247  host-specific (**Figure 4**). We found that the proportion of certain types of ARGs was

248  significantly higher in certain environments than in others. For example, rifamycin

249  resistance genes were found enriched in Actinobacteria (with proportion of 0.1%) and

250  enriched in the soil environment (with proportion of 4.7%) (**Figure 4**). Rifamycin is

251  an important antibacterial agent active against gram-positive bacteria, and it has a

252  wide range of applications[30,31]. The enrichment results were not surprising because

253  *Actinomycetes* is a representative genus widely distributed in various soil

254  environments, and its rifamycin resistance is compatible with its ability for rifamycin

255  production[32-35].

256

257  **Evaluation of the ability for novel ARG identification using a recently annotated**

258  **ARG**

259  We further evaluated ONN4ARG's ability for novel ARG identification based on the

260  assessment of a newly annotated aminoglycoside resistance gene, GAR[7]. GAR is a

261  recently reported aminoglycoside resistance gene (e.g., gentamicin, micronomicin)

262  that belongs to non-beta-lactam, which is not present in CARD (v3.0.3), UniProt (as

263  of May 2021), DEEPARG-DB (v1.0.2), HMD-ARG-DB (as of May 2021), and

264  ONN4ARG-DB. We searched the sequence of GAR with both DeepARG and

265    HMD-ARG models, and the results showed that both of these models indicated it as

266    non-ARG. We searched the sequence of GAR against all the sequences in

267    ONN4ARG-DB using Diamond and did not find any homologous gene as well.

268    However, the prediction by ONN4ARG identified GAR as an ARG resistant to

269    non-beta-lactam with high confidence (probability score = 100%). We should

270    emphasize that though ONN4ARG could only predict GAR as non-beta-lactam and

271    not as sub-type of aminoglycoside, it was the only method used in this study that

272    could predict GAR as an ARG gene, which again confirms ONN4ARG's better

273    generalization ability for novel ARG discovery.

274

275    **Functional verification of candidate novel resistance genes**

276    To identify promising putative novel resistance genes, we used four criteria: (i)

277    remote homologs to reference ARGs, (ii) prediction with high confidence, (iii)

278    predicted to be single-type resistance, and (iv) the host is known. Despite the large

279    number of candidate genes discovered by the ONN4ARG model (**Supplementary**

280    **Table 5**), only 4,365 ARGs fulfilled all mentioned criteria (**Supplementary Table 7**).

281

282    We selected one candidate ARG (Candi_60363_1) for further experimental validation

283    (**Supplementary Tables 8 and 9**). Candi_60363_1, detected in *Streptococcus* in the

284    oral environment, was predicted to be streptomycin (belonging to aminoglycoside)

285    resistant with high confidence by the ONN4ARG model, and the closest homolog of

286    Candi_60363_1 in ONN4ARG-DB is P12055 (sequence identity of 37.2%). One

287    positive control from CARD (AHE40557.1, streptomycin resistance gene) was used

288    in our experiments for verification of the experimental system. All these genes were

289    heterologously expressed in the *E. coli* BL21 (DE3) host by the induction of Isopropyl

290    β-D-1-thiogalactopyranoside (IPTG) and tested for minimal inhibitory concentration

291    (MIC) (**Figure 5a**). The result showed that the mRNA level of the genes increased

292    with the addition of 1 mM IPTG compared with that without IPTG (**Figure 5b**),

293    which verified the expression of the genes induced by IPTG. Furthermore, the MIC of

294    the strain containing the positive control gene AHE40557.1 was more than 1,024

295    µg/ml (**Supplementary Figure 5**), which is consistent with previous reports[36,37]. This

296    verified that our MIC measuring experimental system works well. Our results showed

297    that the MIC of the strain containing Candi_60363_1 was significantly higher than the

298    negative control containing no insert (**Figure 5c**, **Supplementary Figure 5**), which

299    demonstrated the increased resistance to streptomycin of the novel candidate gene

300    Candi_60363_1 and verified the good performance of our model.

301

302    **Phylogeny and structure of Candi_60363_1**

303    There are remote similarities between Candi_60363_1 and all known ARGs in the

304    reference database, including aminoglycoside resistance genes (closest homolog is

305    P12055, sequence identity of 37.2%). The function annotation of P12055 shows that it

306    has the catalytic activity of reaction between streptomycin and ATP, and it is required

307    for streptomycin resistance (https://www.uniprot.org/citations/3357770). Additionally,

308    the search result of Candi_60363_1 using InterPro (the Integrated Resource of Protein

309    Domains and Functional Sites) shows the protein family matching to Candi_60363_1

310    is IPR007530, which is also known as aminoglycoside 6-adenylyltransferase that

311    confers resistance to aminoglycoside antibiotics. Then, we used BLAST to search

312    homologs of Candi_60363_1 from the NCBI non-redundant protein database. The

313    BLAST result showed that there are 44 homologs with sequence identity greater than

314    80%, and they are from various organisms (**Supplementary Table 10**), such as

315    *Streptococcus oralis*, *Peptoniphilus lacrimalis DNF00528*, and *Mycobacteroides*

316    *abscessus subsp. Abscessus*. Considering that Candi_60363_1 is harbored by distantly

317    related species, it obviously has mobility. Notably, the most similar protein of

318    Candi_60363_1 from the NCBI non-redundant protein database (87.5% identity,

319    SHZ78752.1) is also annotated as aminoglycoside adenylyltransferase

320    (**Supplementary Table 10**). The result of BLAST search against the NCBI

321    non-redundant protein database and other databases showed that Candi_60363_1,

322    which is absent in all the existing ARG databases, is highly likely to be an ARG that

323    confers resistance to aminoglycoside antibiotics.

324

325      Aminoglycoside modifying enzymes are the most clinically important resistance

326      mechanism against aminoglycosides[38]. Aminoglycoside modifying enzymes are

327      divided into three enzymatic classes, namely, aminoglycoside N-acetyltransferase

328      (AAC), O-nucleotidyltransferase (ANT), and O-phosphotransferase (APH). We

329      investigated the phylogenetic relationship between Candi_60363_1 and the known

330      aminoglycoside modifying enzymes. The phylogenetic tree of Candi_60363_1 and

331      related proteins (**Figure 6a**) shows that Candi_60363_1 is clearly separated from the

332      known aminoglycoside modifying enzymes and is located among proteins mostly

333      annotated as aminoglycoside adenylyltransferase. Phylogenetic analysis indicated its

334      evolutionarily close relationships with known aminoglycoside adenylyltransferase.

335

336      Protein structure prediction results confirmed the anti-microbial functionality of

337      Candi_60363_1. The optimal Candi_60363_1-streptomycin complex structure and the

338      corresponding interaction details are described in **Figure 6b**. The optimal binding

339      affinity between the Candi_60363_1 and streptomycin is -7.7 kcal/mol

340      (**Supplementary Table 11**), which is 1.6 kcal/mol lower than the negative control. As

341      shown in **Figure 6b**, the Streptomycin ligand can fit the ARG protein structure well

342      and generate a geometric and energetic docking complex.

343

344      From wet-lab experiments, phylogenetic analysis, and protein structure docking, we

345      consider that Candi_60363_1 predicted by ONN4ARG is highly likely a real ARG

346      gene.

347

348      **Discussion**

349      In this study, we proposed an ontology-aware deep learning method, ONN4ARG, for

350      the detection and understanding of antibiotic resistance genes. The ONN4ARG model

351      is capable of accurately identifying ARGs from coarse to fine levels and discovering

352      novel ARGs that lack homology to known sequences in the reference databases. To

353      complement ONN4ARG for ARG mining applications, we have also created a custom

354  ARG database, ONN4ARG-DB, that contains 28,396 well-curated ARGs. The

355  application of ONN4ARG uncovered 120,726 ARGs from microbiome samples, out

356  of which 42,848 are novel, which substantially expands the existing ARGs

357  repositories.

358

359  The novelty of this work is in three contexts. First, ONN4ARG has the potential for

360  detection of remotely homologous ARGs and thus generates a more comprehensive

361  set of ARGs. The advantage of our ONN4ARG model over state-of-the-art models is

362  that ONN4ARG employs a novel ontology-aware layer that incorporates ancestor and

363  descendant annotations to enhance annotation accuracies. The comprehensive

364  antibiotic resistance ontology used in the ONN4ARG model consists of four levels

365  and more than 100 resistance types (**Supplementary Table 12**), which includes

366  hierarchical antibiotic resistance annotations from the most popular ARG database,

367  CARD. Thus, the classification range of the ONN4ARG model is substantially larger

368  than current tools (e.g., 30 types supported for DeepARG and 15 types supported for

369  HMD-ARG). The ability of ONN4ARG to identify remote homologs (i.e., sequence

370  identity between 30% and 40%) allows more accurate prediction for those

371  misclassified by sequence-alignment based tools as false negatives. Therefore,

372  ONN4ARG greatly reduces false negatives and offers a powerful approach for

373  comprehensive and accurate profiling of ARGs.

374

375  Second, it enabled the comprehensive enrichment analysis of ARGs, species-wise and

376  environment-wise. For the actual application of the ONN4ARG model, we

377  investigated the presence of ARGs in a variety of environments, including water, soil,

378  and the human gut, and the results showed that ARGs are environment-specific and

379  host-specific (**Figure 4**). The environment-specific and host-specific phenomenon of

380  ARGs may be caused by specific bacteria evolving to possess specific types of ARGs

381  in response to specific environments, and horizontal gene transfer may be one of the

382  mediating pathways of this process. For example, one published study has reported

383  that *Amycolatopsis* in the soil environment produces rifamycin and thus gains

384    ecological advantages over other bacteria[32].

385

386    Third, the novel ARGs predicted by ONN4ARG could be functionally validated.

387    Functional verification of a novel streptomycin resistance gene (i.e., Candi_60363_1)

388    with wet-lab experiments demonstrated the ability of the ONN4ARG model for novel

389    ARG discovery. Although the MIC test value of Candi_60363_1 was only two times

390    higher than that of the control (**Figure 5**), this increase was still sufficient to indicate

391    the presence of resistance. Moreover, phylogenetic analysis and protein structure

392    docking further confirmed that Candi_60363_1 is highly likely to be an ARG that

393    confers resistance to aminoglycoside antibiotics. Another validation of novel ARG

394    identification based on the assessment of a recently annotated ARG (i.e., GAR) also

395    indicated the ability of the ONN4ARG model for novel ARG discovery. GAR is a

396    novel ARG that is resistant to a variety of aminoglycosides (e.g., gentamicin and

397    micronomicin). We searched the sequence of GAR using other tools (i.e., DeepARG

398    and HMD-ARG), and the results showed that both of those models indicated it as

399    non-ARG. We emphasize that the ONN4ARG model only identified GAR as

400    non-beta-lactam. This shows that the multi-level annotations of ONN4ARG allow low

401    resolution recognition, which can greatly decrease the false negative rate.

402

403    In summary, ONN4ARG is a deep learning approach for ARG identification. It allows

404    in-depth gene mining on large-scale metagenomic data and helps researchers discover

405    novel ARGs. ONN4ARG provides a complete picture of the prevalence of ARGs in

406    the microbial communities and guidance for detection and reduction of the spread of

407    resistance genes in such scenarios, including clinical research, environmental

408    monitoring, and agricultural management.

409

410    ONN4ARG could be improved in a few ways. For more comprehensive ARG

411    prediction, continuous improvement of curating ARG nomenclature and annotation

412    databases is required. For novel ARG prediction, especially those belonging to

413    entirely new ARG families, deep learning models might need to consider more

414 information other than sequence alone. We believe these efforts could lead to a

415 holistic view about ARGs in diverse environments around the globe.

416

## Methods

417

### Dataset

418

419 The ARGs we used in this study for model training and testing were from the

420 Comprehensive Antibiotic Resistance Database (CARD[5,6], v3.0.3). We also used

421 protein sequences from the UniProt (SwissProt and TrEMBL) database to expand our

422 training dataset. First, genes with ARG annotations were collected from CARD (2,587

423 ARGs) and SwissProt (2,261 ARGs). Then, their close homologs (with sequence

424 identities greater than 90%) were collected from TrEMBL (23,728 homologous genes).

425 These annotated and homologous ARGs made up our positive dataset. The negative

426 dataset was made from non-ARG genes that had relatively weak sequence similarities

427 to ARG genes (with sequence identities smaller than 90% and bit-scores smaller than

428 alignment lengths) but not annotated as ARG genes in SwissProt (17,937 genes).

429 Finally, redundant genes with identical sequences were filtered out. As a result, our

430 ARG gene dataset, namely, ONN4ARG-DB, contained 28,396 positive and 17,937

431 negative genes. For evaluation and comparison of ONN4ARG, 75% of the dataset

432 was randomly selected for training, and the remaining 25% of the dataset was selected

433 for testing.

434

### Antibiotic resistance ontology

435

436 The antibiotic resistance ontology was organized into an ontology structure, which

437 contains four levels. The root (first level) is a single node, namely, "arg"

438 (**Supplementary Table 12**). There are 1, 2, 34, and 277 nodes from the first level to

439 the fourth level, respectively. For instance, there are "beta-lactam" and

440 "non-beta-lactam" in the second level, "acridine dye" and "aminocoumarin" in the

441 third level, and "acriflavine" and "clorobiocin" in the fourth level. For example,

442 AHE40557.1 is annotated in the CARD database as a streptomycin resistance gene,

443     which belongs to a lower-level ARG type aminoglycoside and a higher-level ARG

444     type non-beta-lactam (**Supplementary Figure 1**).

445

446     **Protein annotations**

447     The protein sequences for training and testing were annotated according to the

448     antibiotic resistance ontology. For example, AHE40557.1 is annotated in the CARD

449     database as a streptomycin resistance gene, which belongs to a lower-level ARG type

450     aminoglycoside and a higher-level ARG type non-beta-lactam. Accordingly, this

451     protein will be annotated as "arg" at the first level, "non-beta-lactam" at the second

452     level, "aminoglycoside" at the third level, and "streptomycin" at the fourth level.

453

454     **Sequence-alignment**

455     We used Diamond[23] as the sequence-alignment tool for comparison. For queries in the

456     testing set, we searched them against the training set. The target with the highest

457     identity was defined as the closest homologous gene for each query. Then, we

458     compared whether the actual annotation of the query was consistent with the

459     annotation of its closest homologous gene to evaluate the accuracy.

460

461     **DeepARG**

462     DeepARG[12] is a newly developed tool that applies a neural network to identify

463     antibiotic resistance genes. For queries in the testing set, we used the DeepARG[12]

464     model to predict their annotations. Then, we compared whether the actual annotation

465     of the query was consistent with the predicted annotation to evaluate the accuracy.

466

467     **Evaluation and comparison**

468     In this study, the performance of ONN4ARG was evaluated and compared to

469     state-of-the-art models, including sequence-alignment and DeepARG. For these three

470     models, the training dataset was used to train the model parameters, and the testing

471     dataset was used to calculate the prediction accuracies. Both DeepARG and

472     ONN4ARG are deep learning models that use millions of parameters. Unlike deep

473    learning models, sequence-alignment (i.e., Diamond) has only one parameter (i.e., the

474    identity cutoff to distinguish ARG and non-ARG genes).

475

476    **Masking threshold**

477    To simulate remotely homologous ARG genes in our experiments, similarities

478    between the query protein and its close homologs with sequence identities greater

479    than a threshold were masked as zeros (i.e., no signals). For instance, when the

480    masking threshold of testing set equaled 0.4, similarities between the query protein (in

481    the testing set) and its close homologs (in the training set) with sequence identities

482    greater than 40% were masked as zeros. Occasionally, all homologs were masked for

483    a query protein, and such query proteins were removed during training and testing.

484    For example, if query $X$ had two homologs, $M$ and $N$, and assuming the identity of $M$

485    is 0.35 and the identity of $N$ is 0.85, when the masking threshold of the testing set

486    equaled 0.4, similarities between query $X$ and homolog $M$ were masked as zeros.

487    When the masking threshold of the testing set equaled 0.9, query $X$ was removed

488    during testing.

489

490    **Benchmark method**

491    In this study, a prediction was defined to be correct if and only if all ARG annotations

492    (including ancestor annotations from ARG ontologies) were correctly predicted. The

493    accuracy of the tested model was defined as the number of correct predictions over

494    the total number of predictions. The precision of the tested model was defined as the

495    number of true positive predictions over the total number of positive predictions, and

496    the recall was defined as the number of true positive predictions over the total number

497    of true positive plus false negative predictions.

498

499    **ARG mining on metagenomic data**

500    We collected microbiome sequencing data from several published studies

501    (**Supplementary Table 4**), including samples from soil, water, and human body. The

502    gene contigs were processed by Prodigal[27]. Protein sequences were also obtained by

503     the Prodigal program. Then, the ARG annotations of these protein sequences were

504     predicted by using ONN4ARG.

505

**Taxonomy annotation**

507     Kraken2[29] program was used to identify the host of gene contigs. Then, each ARG

508     predicted by ONN4ARG was annotated according to the host of its gene contigs.

509

**Phylogenetic tree**

511     Sequences of the 44 proteins most closely related to Candi_60363_1 were collected

512     using BLASTP with default parameters on the NCBI non-redundant protein database.

513     The retrieved proteins, Candi_60363_1 and all aminoglycoside resistance proteins

514     from ResFinder[8] (https://bitbucket.org/genomicepidemiology/resfinder_db/src/master,

515     last update March 2021), were aligned with ClustalW. The phylogenetic tree was

516     calculated by MEGA[39] (v10) using the maximum likelihood algorithm with default

517     parameters. The Interactive Tree of Life (iTOL v6) online tool[40] was used to prepare

518     the phylogenetic tree for display.

519

**Protein model and docking**

521     Rosetta[41] was utilized to predict the protein structure using ab initio protein folding

522     (http://robetta.bakerlab.org/). The top five protein pockets were generated for docking

523     calculation with Surface Topography of proteins[42] (CASTp). We used the Cambridge

524     Structure Database[43] to generate streptomycin conformers. The 3D protein-ligand

525     complexes were obtained from AutoDock Vina[44].

526

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

529     The candidate resistance gene Candi_60363_1 and a positive control resistance gene

530     AHE40557.1 were synthesized and subcloned into pUC19 vector, replacing *lacZ'*

531     gene. The recombinant plasmids were then transformed into *E. coli* BL21 (DE3). The

532     expression of resistance genes was induced by Isopropyl β-D-1-thiogalactopyranoside

533    (IPTG) and verified by quantitive Real-time PCR (qRT-PCR) assay. Briefly, bacteria

534    were grown in LB supplemented with ampicillin (100 μg/ml) to OD600 of 0.5-0.6 by

535    incubation at 37 °C with 220 rpm agitation, and the bacterial cultures were continued

536    to grow until OD600 reached to 1.0 by adding or without adding 1 mM IPTG. The

537    cells were harvested and total RNAs were purified using Bacterial RNA Extraction

538    Kit (Vazyme Biotech). RNA reverse transcription was performed by using HiScript®

539    II Q Select RT SuperMix for qPCR kit (Vazyme Biotech). qRT-PCR was performed

540    by using SYBR Green Master Mix-High ROX Premixed (Vazyme Biotech) in a

541    Stepone Plus system (Applied Biosystems). The *ldh* gene was used as internal control

542    in all reactions. The relative fold changes were determined using the $2^{-\Delta\Delta Ct}$ method, in

543    which *ldh* was used for normalization. The protein sequences of the synthesized genes

544    were presented in **Supplementary Table 8** and the primer sequences for qRT-PCR

545    were listed in **Supplementary Table 9**.

546

**MIC determination**

548    Minimal inhibitory concentrations (MICs) of the antibiotic for the strains containing

549    resistance genes were determined using E-tests. Single colonies of the strains were

550    incubated in 3 ml Mueller-Hinton (MH) medium with the addition of 100 μg/ml

551    ampicillin at 35 °C for 4 hours, and the cells equal to $1.5 \times 10^8$ cells/ml were spread on

552    MH agar plates with the addition of 100 μg/ml ampicillin and 1 mM IPTG, and

553    streptomycin MIC Test Strips (Liofilchem®) were put in the middle of the plates. The

554    plates were incubated at 35 °C for 18-24 hours, and the MICs were read. The strain

555    containing empty vector was used as a negative control.

556

**Data availability**

558    We collected metagenomic samples from several published studies[25,26], and these

559    samples are mainly from "marine", "soil" and "human" associated environments. For

560    human associated samples, including two gut groups (one group from Madagascar,

561    i.e., GutM, the other group from Denmark, i.e., GutD), one oral group and one skin

562 group (both oral and skin groups are from HMP project). Details and links about these

563 samples are shown in **Supplementary Table 4**. The ONN4ARG-DB dataset could be

564 accesses at: http://onn4arg.xfcui.com/.

565

## Code availability

567 All source codes can be accessed at: https://github.com/xfcui/onn4arg, and

568 http://onn4arg.xfcui.com/.

569

## Acknowledgments

571 We are grateful to Mingyue Cheng and Hui Chong for insightful discussions. We

572 thank LetPub (www.letpub.com) for its linguistic assistance during the preparation of

573 this manuscript. This work was partially supported by National Science Foundation of

574 China grant 81774008, 81573702, 31871334 and 31671374, and the Ministry of

575 Science and Technology's national key research and development program grant (No.

576 2018YFC0910502).

577

## Author contributions

579 KN, XC conceived and proposed the idea, and designed the study. YZ, CC, QJ, XZ,

580 XC performed the experiments and analyzed the data. YZ, CC, XZ, KN and XC

581 contributed to editing and proof-reading the manuscript. All authors read and

582 approved the final manuscript.

583

## Competing interests

585 The authors declare that they have no competing interests.

586

## Ethics approval and consent to participate

588 Not applicable

589

590

## References

1       Thomas, A. M. & Segata, N. Multiple levels of the unknown in microbiome research. *bmc biology* **17**, doi:10.1186/S12915-019-0667-Z (2019).

2       Li, J. *et al.* An integrated catalog of reference genes in the human gut microbiome. *nature biotechnology* **32**, 834-841, doi:10.1038/NBT.2942 (2014).

3       Brogan, D. M. & Mossialos, E. A critical analysis of the review on antimicrobial resistance report and the infectious disease financing facility. *globalization and health* **12**, 8-8, doi:10.1186/S12992-016-0147-Y (2016).

4       Goossens, H., Ferech, M., Stichele, R. V. & Elseviers, M. Outpatient antibiotic use in Europe and association with resistance: a cross-national database study. *the lancet* **365**, 579-587, doi:10.1016/S0140-6736(05)17907-0 (2005).

5       Jia, B. *et al.* CARD 2017: expansion and model-centric curation of the comprehensive antibiotic resistance database. *nucleic acids research* **45**, doi:10.1093/NAR/GKW1004 (2017).

6       Alcock, B. P. *et al.* CARD 2020: antibiotic resistome surveillance with the comprehensive antibiotic resistance database. *nucleic acids research* **48**, doi:10.1093/NAR/GKZ935 (2019).

7       Böhm, M.-E., Razavi, M., Marathe, N. P., Flach, C.-F. & Larsson, D. G. J. Discovery of a novel integron-borne aminoglycoside resistance gene present in clinical pathogens by screening environmental bacterial communities. *microbiome* **8**, 1-11, doi:10.1186/S40168-020-00814-Z (2020).

8       Bortolaia, V. *et al.* ResFinder 4.0 for predictions of phenotypes from genotypes. *journal of antimicrobial chemotherapy* **75**, 3491-3500, doi:10.1093/JAC/DKAA345 (2020).

9       Bateman, A. *et al.* UniProt: A hub for protein information. *nucleic acids research* **43**, doi:10.1093/NAR/GKU989 (2015).

10      Rowe, W. *et al.* Search Engine for Antimicrobial Resistance: A Cloud Compatible Pipeline and Web Interface for Rapidly Detecting Antimicrobial Resistance Genes Directly from Sequence Data. *plos one* **10**, doi:10.1371/JOURNAL.PONE.0133492 (2015).

11      Kleinheinz, K. A., Joensen, K. G. & Larsen, M. V. Applying the ResFinder and VirulenceFinder web-services for easy identification of acquired antibiotic resistance and E. coli virulence genes in bacteriophage and prophage nucleotide sequences. *Bacteriophage* **4**, e27943-e27943, doi:10.4161/bact.27943 (2014).

12      Arango-Argoty, G. *et al.* DeepARG: a deep learning approach for predicting antibiotic resistance genes from metagenomic data. *microbiome* **6**, 23-23, doi:10.1186/S40168-018-0401-Z (2018).

13      Davis, J. J. *et al.* Antimicrobial Resistance Prediction in PATRIC and RAST. *scientific reports* **6**, 27930-27930, doi:10.1038/SREP27930 (2016).

14      Lakin, S. M. *et al.* Hierarchical Hidden Markov models enable accurate and diverse detection of antimicrobial resistance sequences. *Commun Biol* **2**,
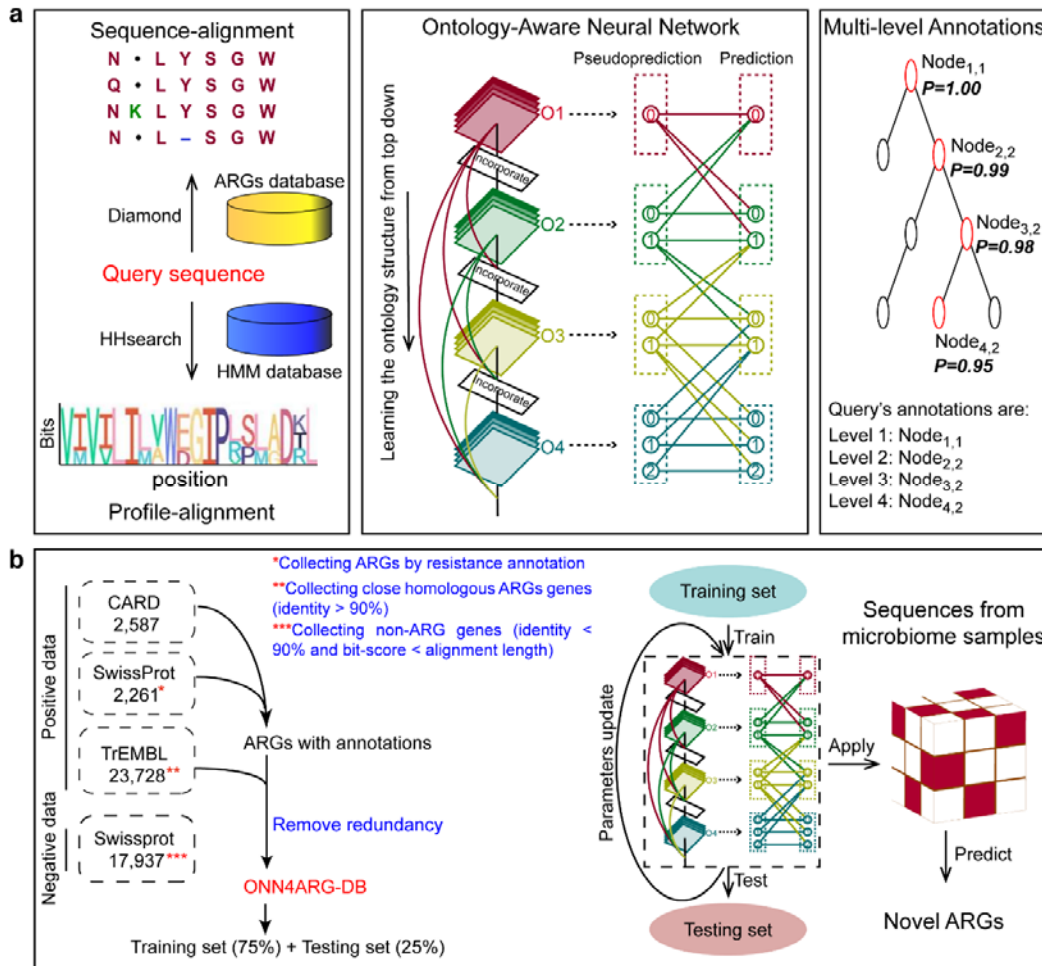
634         294-294, doi:10.1038/s42003-019-0545-9 (2019).

635   15    Doster, E. *et al.* MEGARes 2.0: a database for classification of antimicrobial
636         drug, biocide and metal resistance determinants in metagenomic sequence data.
637         *Nucleic acids research* **48**, D561-D569, doi:10.1093/nar/gkz1010 (2020).

638   16    Li, Y. *et al.* HMD-ARG: hierarchical multi-task deep learning for annotating
639         antibiotic     resistance      genes.      *microbiome*      **9**,      40-40,
640         doi:10.1186/S40168-021-01002-3 (2021).

641   17    Gupta, S. K. *et al.* ARG-ANNOT, a new bioinformatic tool to discover
642         antibiotic resistance genes in bacterial genomes. *Antimicrobial agents and*
643         *chemotherapy* **58**, 212-220, doi:10.1128/AAC.01310-13 (2014).

644   18    Feldgarden, M. *et al.* Validating the AMRFinder Tool and Resistance Gene
645         Database    by    Using    Antimicrobial    Resistance    Genotype-Phenotype
646         Correlations   in   a   Collection   of   Isolates.  *Antimicrobial   agents   and*
647         *chemotherapy* **63**, e00483-00419, doi:10.1128/AAC.00483-19 (2019).

648   19    Inouye, M. *et al.* SRST2: Rapid genomic surveillance for public health and
649         hospital    microbiology    labs.    *genome    medicine*    **6**,    90-90,
650         doi:10.1186/S13073-014-0090-6 (2014).

651   20    Rowe, W. P. M. & Winn, M. D. Indexed variation graphs for efficient and
652         accurate    resistome    profiling.    *bioinformatics*    **34**,    3601-3608,
653         doi:10.1093/BIOINFORMATICS/BTY387 (2018).

654   21    Altschul, S. F., Gish, W., Miller, W. C., Myers, E. W. & Lipman, D. J. Basic
655         Local Alignment Search Tool. *journal of molecular biology* **215**, 403-410,
656         doi:10.1016/S0022-2836(05)80360-2 (1990).

657   22    Edgar, R. C. Search and clustering orders of magnitude faster than BLAST.
658         *Bioinformatics* **26**, 2460-2461, doi:10.1093/bioinformatics/btq461 (2010).

659   23    Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment
660         using   DIAMOND.   *nature   methods*   **12**,   59-60,   doi:10.1038/NMETH.3176
661         (2015).

662   24    Steinegger, M. *et al.* HH-suite3 for fast remote homology detection and deep
663         protein    annotation.    *bmc    bioinformatics*    **20**,    1-15,
664         doi:10.1186/S12859-019-3019-7 (2019).

665   25    Sunagawa, S. *et al.* Structure and function of the global ocean microbiome.
666         *science* **348**, 1261359-1261359, doi:10.1126/SCIENCE.1261359 (2015).

667   26    Mitchell, A. L. *et al.* EBI Metagenomics in 2017: enriching the analysis of
668         microbial communities, from sequence reads to assemblies. *Nucleic acids*
669         *research* **46**, D726-D735, doi:10.1093/nar/gkx967 (2018).

670   27    Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation
671         initiation    site    identification.    *bmc    bioinformatics*    **11**,    119-119,
672         doi:10.1186/1471-2105-11-119 (2010).

673   28    Li, B. *et al.* Metagenomic and network analysis reveal wide distribution and
674         co-occurrence of environmental antibiotic resistance genes. *ISME J* **9**,
675         2490-2502, doi:10.1038/ismej.2015.59 (2015).

676   29    Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with
677         Kraken 2. *genome biology* **20**, 1-13, doi:10.1186/S13059-019-1891-0 (2019).

678   30   Qi, F. *et al.* Deciphering the late steps of rifamycin biosynthesis. *Nature*
679        *Communications* **9**, 2342, doi:10.1038/s41467-018-04772-x (2018).
680   31   Floss, H. G. & Yu, T.-W. RifamycinMode of Action, Resistance, and
681        Biosynthesis. *Chemical Reviews* **105**, 621-632, doi:10.1021/cr030112j (2005).
682   32   Yao, Y., Zhang, W., Jiao, R., Zhao, G. & Jiang, W. Efficient isolation of total
683        RNA from antibiotic-producing bacterium Amycolatopsis mediterranei.
684        *Journal     of     Microbiological     Methods*    **51**,     191-195,
685        doi:doi.org/10.1016/S0167-7012(02)00078-7 (2002).
686   33   Wilson, M. C., Gulder, T. A. M., Mahmud, T. & Moore, B. S. Shared
687        biosynthesis of the saliniketals and rifamycins in Salinispora arenicola is
688        controlled by the sare1259-encoded cytochrome P450. *J Am Chem Soc* **132**,
689        12757-12765, doi:10.1021/ja105891a (2010).
690   34   Saxena, A., Kumari, R., Mukherjee, U., Singh, P. & Lal, R. Draft Genome
691        Sequence of the Rifamycin Producer Amycolatopsis rifamycinica DSM 46095.
692        *Genome Announc* **2**, e00662-00614, doi:10.1128/genomeA.00662-14 (2014).
693   35   Huang, H. *et al.* Micromonospora rifamycinica sp. nov., a novel actinomycete
694        from mangrove sediment.    **58**, 17-20, doi:doi.org/10.1099/ijs.0.64484-0
695        (2008).
696   36   Pinto-Alphandary, H., Mabilat, C. & Courvalin, P. Emergence of
697        aminoglycoside resistance genes aadA and aadE in the genus Campylobacter.
698        *antimicrobial     agents     and     chemotherapy*    **34**,     1294-1296,
699        doi:10.1128/AAC.34.6.1294 (1990).
700   37   Holden, M. T. G. *et al.* Rapid evolution of virulence and drug resistance in the
701        emerging    zoonotic    pathogen    Streptococcus    suis.    *plos    one*    **4**,
702        doi:10.1371/JOURNAL.PONE.0006072 (2009).
703   38   Ramirez, M. S., Nikolaidis, N. & Tolmasky, M. Rise and dissemination of
704        aminoglycoside resistance: the aac(6')-Ib paradigm. *frontiers in microbiology*
705        **4**, 121-121, doi:10.3389/FMICB.2013.00121 (2013).
706   39   Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. MEGA X: Molecular
707        Evolutionary Genetics Analysis across Computing Platforms. *Mol Biol Evol*
708        **35**, 1547-1549, doi:10.1093/molbev/msy096 (2018).
709   40   Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and
710        new    developments.    *Nucleic    Acids    Research*    **47**,    W256-W259,
711        doi:10.1093/nar/gkz239 (2019).
712   41   Rohl, C. A., Strauss, C. E. M., Misura, K. M. S. & Baker, D. Protein structure
713        prediction    using    Rosetta.    *methods    in    enzymology*    **383**,    66-93,
714        doi:10.1016/S0076-6879(04)83004-0 (2004).
715   42   Tian, W., Chen, C. & Liang, J. CASTp 3.0: Computed Atlas of Surface
716        Topography    of    Proteins    and    Beyond.    *biophysical    journal*    **114**,
717        doi:10.1016/J.BPJ.2017.11.325 (2018).
718   43   Cole, J. C., Korb, O., McCabe, P., Read, M. G. & Taylor, R. Knowledge-Based
719        Conformer Generation Using the Cambridge Structural Database. *journal of*
720        *chemical    information    and    modeling*    **58**,    615-629,
721        doi:10.1021/ACS.JCIM.7B00697 (2018).

722   44    Trott, O. & Olson, A. J. AutoDock Vina: improving the speed and accuracy of
723         docking with a new scoring function, efficient optimization, and
724         multithreading. *J Comput Chem* **31**, 455-461, doi:10.1002/jcc.21334 (2010).

725

726 **Figure 1**



727

728 **Figure 1. Overview of the ONN4ARG model and its use for novel ARG discovery. a.** The

729 input (left), architecture (middle), and output (right) of the ONN4ARG model. ONN4ARG takes

730 similarities between the query gene sequence and ARG gene sequences and profiles as inputs.

731 Then, ontology-aware layers (i.e., O1, O2, O3, and O4) are employed to incorporate ancestor and

732 descendant annotations to enhance annotation accuracy. ONN4ARG outputs multi-level

733 annotations of antibiotic resistance types, which are compatible with the antibiotic resistance

734 ontology structure. **b.** Building the dataset for training and testing and applying it on microbiome

735 sequencing data to discover novel ARGs.

736

737 **Figure 2**



738

739 **Figure 2. Systematic evaluation and comparison between sequence-alignment, DeepARG,**

740 **and ONN4ARG. a.** The accuracy of three models on ARG classification was assessed using a box

741 plot. Diamond was used for sequence-alignment; significance test was based on the *t*-test. **b.** The

742 precision and recall of DeepARG and ONN4ARG on ARG classification for each antibiotic

743 resistance type. The masking threshold of testing set equaled 0.4 (details of masking threshold are

744 provided in **Methods**).

745

## Figure 3



**Figure 3. Broad-spectrum profile of predicted ARGs among diverse environments. a.** The proportion of predicted ARGs for different protein sequence lengths. **b.** The abundance ratio of predicted ARGs among diverse environments. Abundance ratio was defined as the number of ARGs divided by the number of total genes. **c.** The proportion of predicted ARGs for different sequence identities among diverse environments. **d.** Number of genes in ONN4ARG-DB (left), predicted homologous ARGs (middle), and predicted novel ARGs (right) for various resistance types. The horizontal axis indicates the logarithmic number of genes, and the vertical axis indicates different antibiotic resistance types. We collected metagenomic samples from several

756    published studies; these samples were mainly from "marine," "soil," and "human" environments.

757    Human-associated samples consisted of two gut groups (one group from Madagascar, i.e., GutM;

758    the other group from Denmark, i.e., GutD), one oral group, and one skin group (both oral and skin

759    groups were from the HMP project).

760

761    **Figure 4**



762

763    **Figure 4. Enrichment of predicted ARGs among diverse environments and hosts. a.** Relative

764    abundance and enrichment of ARGs among diverse environments. Abundance ratio was defined as

765    the number of ARGs divided by the number of total genes. **b.** Proportion and enrichment of ARGs

766    among diverse hosts. Colors indicate the proportion of ARGs for each phylum and resistance type.

767    Results for the most abundant five phyla that carry ARGs are shown. "+": P-value < 0.05 (*t*-test);

768    "++": P-value < 0.01 (*t*-test).

769

770    **Figure 5**



772    **Figure 5. Functional validation of a predicted candidate novel ARG. a.** A diagram showing

773    the procedure of heterologous expression and functional analysis of the predicted candidate ARG

774    in the *E. coli* BL21 (DE3) host. **b.** Gene expression validation of the predicted candidate ARG.

775    The vertical axis indicates the relative mRNA level. **c.** The MIC of the predicted candidate ARG

776    and negative control. The vertical axis indicates the MIC value. The MIC of the predicted

777    candidate novel ARG is significantly higher than the negative control (*t*-test, P-value = 3.5e-3).

778

779 **Figure 6**



780

781 **Figure 6. Phylogenetic analysis and structure investigation of Candi_60363_1. a.**

782 Phylogenetic tree of aminoglycoside resistance enzymes, Candi_60363_1, and its homologs from

783 the NCBI non-redundant protein database. ANT: O-nucleotidyltransferase, AAC:

784 N-acetyltransferase, APH: O-phosphotransferase, AADT: aminoglycoside adenylyltransferase. **b.**

785 The optimal Candi_60363_1-streptomycin complex structure (left), and the local interactions

786    between ligand and neighboring residues (right). The docking experiment indicates there are six

787    neighboring residues whose distances are less than three angstroms.

788

**a**

**Sequence-alignment**

N · L Y S G W
Q · L Y S G W
N K L Y S G W
N · L – S G W

Diamond → ARGs database

**Query sequence**

HHsearch → HMM database

Bits / position

**Profile-alignment**

**Ontology-Aware Neural Network**

Learning the ontology structure from top down

Pseudoprediction    Prediction

Incorporate

O1

Incorporate

O2

Incorporate

O3

O4

**Multi-level Annotations**

Node$_{1,1}$
**P=1.00**

Node$_{2,2}$
**P=0.99**

Node$_{3,2}$
**P=0.98**

Node$_{4,2}$
**P=0.95**

Query's annotations are:
Level 1: Node$_{1,1}$
Level 2: Node$_{2,2}$
Level 3: Node$_{3,2}$
Level 4: Node$_{4,2}$

**b**

*Collecting ARGs by resistance annotation
**Collecting close homologous ARGs genes (identity > 90%)
***Collecting non-ARG genes (identity < 90% and bit-score < alignment length)

Positive data

CARD 2,587

SwissProt 2,261*

TrEMBL 23,728**

Negative data

Swissprot 17,937***

ARGs with annotations

**Remove redundancy**

**ONN4ARG-DB**

Training set (75%) + Testing set (25%)

**Training set**

Train

Parameters update

O1
O2
O3
O4

Test

**Testing set**

Apply →

Sequences from microbiome samples

Predict ↓

**Novel ARGs**

**a**

Accuracy(%)

*: $P < 0.05$

non-beta-lactam ●

beta-lactam ●

**b**

Precision    Recall

acridine dye
aminocoumarin
aminoglycoside
carbapenem
cephalosporin
cephamycin
diaminopyrimidine
elfamycin
fluoroquinolone
free fatty
glycopeptide
glycycline
lincosamide
macrolide
monobactam
nitrofuran
nitroimidazole
nucleoside
oxazolidinone
penam
penem
peptide
phenicol
pleuromutilin
polyamine
rifamycin
streptogramin
sulfonamide
sulfone
tetracyline
without-drug-class

DeepARG  ONN4ARG     DeepARG  ONN4ARG

**a**

Proportion(%)

Length (aa)
[min-128), [128-256), [256-512), [512-max)

**b**

Ratio(%)*

*: number of ARGs / number of total genes

Group
Marine, Soil, GutM, GutD, Oral, Skin

**c**

Proportion(%)

Identity(%)
[30, 40)
[40, 70)
[70, 90)
[90, 100)

Group
Marine, Soil, GutM, GutD, Oral, Skin

**d**

● non-beta-lactam
● beta-lactam
* No. = 1

acridine dye
aminocoumarin
aminoglycoside
carbapenem
cephalosporin
cephamycin
diaminopyrimidine
elfamycin
fluoroquinolone
free fatty
glycopeptide
glycycline
lincosamide
macrolide
monobactam
nitrofuran  *
nitroimidazole
nucleoside
oxazolidinone
penam
penem
peptide
phenicol
pleuromutilin
polyamine
rifamycin
streptogramin
sulfonamide
sulfone  *
tetracyline
without-drug-class

Log$_{10}$ No.Resistance Genes
ONN4ARG-DB

● non-beta-lactam
● beta-lactam

acridine dye
aminocoumarin
aminoglycoside
carbapenem
cephalosporin
cephamycin
diaminopyrimidine
elfamycin
fluoroquinolone
free fatty
glycopeptide
glycycline
lincosamide
macrolide
monobactam
nitrofuran
nitroimidazole
nucleoside
oxazolidinone
penam
penem
peptide
phenicol
pleuromutilin
polyamine
rifamycin
streptogramin
sulfonamide
sulfone
tetracyline
without-drug-class

Log$_{10}$ No.Resistance Genes
Predicted ARGs (identity >= 40%)

● non-beta-lactam
● beta-lactam
** No. = 0

acridine dye
aminocoumarin
aminoglycoside
carbapenem
cephalosporin
cephamycin
diaminopyrimidine
elfamycin  **
fluoroquinolone
free fatty
glycopeptide
glycycline
lincosamide
macrolide
monobactam
nitrofuran  **
nitroimidazole
nucleoside
oxazolidinone
penam
penem
peptide
phenicol
pleuromutilin
polyamine
rifamycin
streptogramin
sulfonamide
sulfone
tetracyline
without-drug-class

Log$_{10}$ No.Resistance Genes
Predicted ARGs (identity < 40%)

**a**

Ratio(%)
- 0.30
- 0.25
- 0.20
- 0.15
- 0.10
- 0.05
- 0.00

● non-beta-lactam
● beta-lactam

+: $P < 0.05$
++: $P < 0.01$

Rows (top to bottom): acridine dye, aminocoumarin, aminoglycoside, carbapenem, cephalosporin, cephamycin, diaminopyrimidine, elfamycin, fluoroquinolone, free fatty, glycopeptide, glycycline, lincosamide, macrolide, monobactam, nitrofuran, nitroimidazole, nucleoside, oxazolidinone, penam, penem, peptide, phenicol, pleuromutilin, polyamine, rifamycin, streptogramin, sulfonamide, sulfone, tetracycline, without-drug-class

Columns: Marine, Soil, GutM, GutD, Oral, Skin

**b**

Proportion
- 0.10
- 0.08
- 0.06
- 0.04
- 0.02
- 0.00

● non-beta-lactam
● beta-lactam

+: $P < 0.05$
++: $P < 0.01$

Columns: Actinobacteria, Bacteroidetes, Cyanobacteria, Firmicutes, Proteobacteria, Others, Unclassified

**a**

Candidate ARG → pUC19 → *E. coli* BL21(DE3) → Incubation → MIC measurement

E-test strip for streptomycin

**b**

Relative mRNA level

- IPTG
+1mM IPTG

AHE40557.1    Candi_60363_1

**c**

P-value = 3.5e-3

MIC(μg/ml)

Negative control    Candi_60363_1