

1 **Metagenomic screening of global microbiomes identifies**
2 **pathogen-enriched environments**

3 **Xiaofang Li^{1,2}**

4 ¹ Key Laboratory of Agricultural Water Resources, Centre for Agricultural Resources Research, Institute of
5 Genetics and Developmental Biology, Chinese Academy of Sciences, Shijiazhuang 050021, China

6 ² Environment Centres, Sustainable Minerals Institute, The University of Queensland, Brisbane, Queensland
7 4072, Australia

8 **Correspondence** to XL xfli@sjziam.ac.cn; Tel 86 311 85822874; Fax 86 31185815093.

9

10 **Abstract**

11 **Background** Human pathogens are widespread in the environment, and examination of pathogen-enriched
12 environments in a rapid and high-throughput fashion is important for development of pathogen-risk
13 precautionary measures.

14 **Methods** In this study, a Local BLASTP procedure for metagenomic screening of pathogens in the environment
15 was developed using a toxin-centered database. A total of 27 microbiomes derived from ocean water, freshwater,
16 soil, feces, and wastewater were screened using the Local BLASTP procedure. Bioinformatic analysis and
17 Canonical Correspondence Analysis were conducted to examine whether the toxins included in the database
18 were taxonomically associated.

19 **Results** The specificity of the Local BLASTP method was tested with known and unknown toxin sequences.
20 Bioinformatic analysis indicated that most toxins were phylum-specific but not genus-specific. Canonical
21 Correspondence Analysis implied that almost all of the toxins were associated with the phyla of *Proteobacteria*,
22 *Nitrospirae* and *Firmicutes*. Local BLASTP screening of the global microbiomes showed that pore-forming
23 RTX toxin and adenylate cyclase Cya were most prevalent globally in terms of relative abundance, while
24 polluted water and feces samples were the most pathogen-enriched.

25 **Conclusions** A Local BLASTP procedure was established for rapid detection of toxins in environmental
26 samples. Screening of global microbiomes in this study provided a quantitative estimate of the most prevalent
27 toxins and most pathogen-enriched environment.

28 **Keywords**

29 Metagenomics, Microbiome, Local BLASTP, Toxins, Pathogens

30

31 **Introduction**

32 Rapid identification of pathogens in a particular environment is important for pathogen-risk management.
33 Human pathogens are ubiquitous in the environment, and infections from particular environments have been
34 reported worldwide. For example, soil-related infectious diseases are common [1, 2]. *Legionella longbeachae*
35 infection has been reported in many cases, mainly due to potting mixes and composts [3]. Survival of enteric
36 viruses and bacteria has also been detected in various water environments, including aquifers and lakes [4-7].

37 Examination of pathogens from infected individuals with a particular clinical syndrome has been a major
38 achievement of modern medical microbiology [8]. Nevertheless, we still know little about the magnitude of the
39 abundance and diversity of known common pathogens in various environments, which is very important to the
40 development of appropriate precautions for individuals who work or play with certain environmental substrates.
41 This can be realized through metagenomic detection of pathogenic factors in a time-efficient and high-
42 throughput manner using next-generation sequencing methods.

43 Metagenomic detection of pathogens can be accomplished through different schemes. Li et al. examined
44 the level and diversity of bacterial pathogens in sewage treatment plants using a 16S rRNA amplicon-based
45 metagenomic procedure [9]. Quantitative PCR has also been applied for monitoring specific pathogens in
46 wastewater [10]. More studies have applied the whole-genome-assembly scheme to detect one or multiple
47 dominant pathogens, most of which were for viral detection in clinical samples [11-14]. Although metagenomic-
48 based whole-genome-assembly for bacterial pathogen detection can be conducted at the single species level [15],
49 its computational requirements are high if in a high-throughput fashion. In 2014, Baldwin et al. [16] designed
50 the PathoChip for screening pathogens in human tissues by targeting unique sequences of viral and prokaryotic
51 genomes with multiple probes in a microarray. This approach can screen virtually all pathogen-enriched samples
52 in a high-throughput manner.

53 Despite the aforementioned progress in metagenomic tools for pathogen detection, metagenomic screening
54 for bacterial pathogens in environments such as soil, where microbial diversity is tremendous, is still
55 challenging. This is mostly due to difficulty in assembling short reads generated by next-generation sequencing
56 [8]. The whole-genome-assembly approach is efficient at identifying viromes, but not at dealing with bacterial
57 communities. Amplicon-based approaches are able to detect bacterial pathogens in a high-throughput manner;
58 however, it is well known that phenotypic diversity exists widely across and within microbial species of a genus
59 because of divergent evolution [17, 18]. This also holds true for pathogenic factors [19]. Moreover, toxin factors,
60 such as the Shiga toxin (*stx*) of *Shigella*, are primarily transferable through lateral gene transfer, which leads to
61 the continuous evolution of pathogen species [20]. Therefore, it is necessary to examine the pathogen diversity
62 in environmental metagenomes using essential virulence genes as biomarkers.

63 In this study, a toxin-centered virulence factors database was established, and the well-developed Local
64 BLASTP method was applied to detect virulent factors in various environments globally. This procedure is
65 metagenome-based and can be conducted in a high-throughput fashion, which greatly simplifies development of
66 precautions for pathogen-enriched environments.

67 **Methods and Materials**

68 **Environments and their metagenomes**

69 Twenty-seven metagenomes were selected and downloaded from the MG-RAST server (Table 1). These
70 metagenomes were derived from ocean water, freshwater, wastewater, natural soil, deserts and feces,
71 representing the major environmental media found worldwide. Sequencing methods of the metagenomes include
72 the Illumina, Ion Torrent and 454 platforms, and predicted proteins in the metagenomes ranged from 33,743
73 (fresh water, ID mgm4720261) to 1,966,121 (weedy garden soil, ID mgm4679254). The gene calling results
74 were used for toxin factor screening in this study. The taxonomic composition at the genus level was also
75 retrieved from the MG-RAST server for each metagenome.

76 **Toxin factor database**

77 A toxin-centered database was established for bacterial pathogen detection in metagenomes in this study.
78 Candidate toxin factors for pathogenic screening of environmental metagenomes were gathered based on well-
79 studied pathogens summarized in Wikipedia® under the entry of “pathogenic bacteria”, the Virulence Factor
80 Database [21], a soil borne pathogen report by Jeffery and van der Putten [2], and a manure pathogen report by
81 the United States Water Environment Federation [22]. Sequences of the toxin factors were then retrieved by
82 searching the UniProt database using the toxin plus pathogen names as an entry [23], while typical homologs at
83 a cutoff E value of 10^{-6} were gathered from GenBank based on BLAST results. Considering that virulence
84 process involves several essential factors including toxins, various pathogen-derived secretion proteins were
85 also included in the database, and it was tested that whether secretion proteins were as specific as toxin proteins
86 for pathogen detection. The disease relevance of all virulence factors was screened using the WikiGenes system
87 [24] and relevant publications (Table 2).

88 **Local BLASTP**

89 The Local BLASTP was applied following the procedure used in our previous study [58]. Basically, the gene
90 calling results of each metagenome were searched against the toxin factor database using BLASTP embedded in
91 BioEdit. The cutoff expectation E value was set as 10^{-6} . The results of the Local BLASTP in BioEdit were then
92 copied to an Excel worksheet, after which they were subjected to duplicate removal, quality control and
93 subtotaled according to database ID. Duplicate removal was based on the hypothesis that each sequence
94 contains one copy of a specific toxin factor, since the gene-calling results were used. For quality control of the
95 BLSAT results, a cutoff value of 40% for amino acids identity and 20 aa (1/3 of the length of the shortest toxin
96 factors (e.g., the Heat-Stable Enterotoxin C)) for query alignment length were used to filter the records. The
97 toxins abundance matrix was formed for subsequent analyses.

98 **Specificity tests of the Local BLASTP method**

99 Sequences from the toxin database established in this study, as “known sequences” to the database, were
100 selected randomly and searched against the database using the BLASTP procedure. The genome of *Clostridium*

101 *perfringens* ATCC 13124 (NC_008261), as “unknown” sequences to the database, was subject the Local
102 BLASTX procedure as well. Homologous proteins were searched exhaustively in the GenBank database using
103 BLASTP, with the representative toxin factors in the toxins database as a query. Sequences were retrieved and
104 aligned using ClustalW, and Maximum-likelihood phylogeny was conducted with MEGA 7 [59].

105 **Data analysis**

106 The toxin frequency in each metagenome was normalized to a total gene frequency of 1M to eliminate the
107 effects of gene pool size. Toxin abundance in the 27 metagenomes was visualized using Circos [60]. The genus
108 abundance of all metagenomes was calculated and sorted by genus name, followed by manual construction of a
109 genus abundance matrix for subsequent biodiversity-toxin abundance Canonical Correspondence Analysis using
110 R [61].

111 **Results and Discussion**

112 In this study, a toxin-centered database was established for bacterial pathogen screening in various microbiomes
113 globally through a Local BLASTP procedure. The specificity of the procedure was tested, the relative
114 abundance of toxins in the microbiomes was examined, and the toxin-taxonomic abundance correspondence
115 analysis was performed.

116 Like the previously established Local BLASTN method for antibiotic and metal resistance genes screening
117 [58, 62, 63], the Local BLASTP method using the toxin-centered pathogen database in this study was successful
118 at accurately identifying toxin proteins from the database. For screening of the *Clostridium perfringens* ATCC
119 13124 genome, the methods successfully detected the pore-forming genes and multiple copies of the
120 glucosyltransferase (*toxB*-like) and ADP-ribosyltransferase (*spvB*-like) genes, based on the raw data. These
121 results are consistent with the virulence genetic features of *Clostridium* sp. [21], which have not been well
122 detailed in the GenBank annotation record. Such a cross-validation positively indicated that the Local BLASTP
123 procedure established here is useful in predicting toxin genes in unknown genomes. Yet for a semi-quantitative
124 method to estimate toxin factors in metagenomes, a false positive analysis is required to examine to what level
125 mismatch is included in the Local BLASTP results. Actually, the cutoff values of identity greatly impact the
126 homolog virulence factor abundance returned. At cutoff values of 40% for identify and 20 aa for alignment
127 length, only four records for *Clostridium perfringens* ATCC 13124 genome query were returned after
128 duplication removal, one for 1-phosphatidylinositol phosphodiesterase, one for pore-forming alveolysin, one for
129 Ornithine carbamoyltransferase and one for RNA interferase NdoA. At a cutoff identity value of 35%, one more
130 record (Toxin secretion ATP binding protein) was returned. This means that the Local BLASTP procedure was
131 able to detect the virulence factors in unknown genomic dataset at least semi-quantitatively, with proper cutoff
132 values for data quality control. The accuracy of the BLASTP procedure in virulence factor detection was further
133 tested using the genomes of *Bacillus thuringiensis* serovar *konkukian* str. 97-27 (AE017355.1) and *Helicobacter*
134 *pylori* 26695 (AE000511.1) (results not shown).

135 As mentioned above, functional genes including toxin factors may partly evolve through lateral gene
136 transfer, which makes their taxonomic affiliation difficult. It is thus interesting to explore how specific toxin

137 factors are associated with the taxonomic units of pathogens. Here, I explored this issue by investigating the
138 taxonomic distribution of homologs of toxins retrieved from the GenBank database. Generally, at a lower
139 expectation value, most toxins were associated with a specific group of pathogens. For example, at a cutoff E
140 value of 10^{-6} (the default unless specified), 241 out of 242 returned records of *Mycobacterium tuberculosis*
141 RelE homologs fell within the phylum *Actinobacteria*. Moreover, 89% of these homologs were from the genus
142 *Mycobacterium*, while 99.7% of *Yersinia pestis* CdiA homologs and 92.7% of *Bordetella pertussis* cya homologs
143 belonged to *Proteobacteria*, and homologs of *Aeromonas dhakensis* repeats-in toxin (RtxA) were mostly
144 associated with the class *Gammaproteobacteria* (206 out of 242). However, no obvious genus-toxin association
145 was identified. It is worth noting that these results largely depended on the availability of toxin sequences in
146 each taxonomic unit. The lack of a genus-toxin association basically denied the possibility of detecting a
147 specific pathogen using a specific toxin as a single signature [16].

148 It is still not clear whether virulence secretion proteins are specific for pathogen detection as signatures,
149 through they are essential for virulence process [20]. For example, the contact-dependent toxin delivery protein
150 CdiA was found to be widespread in bacteria [37]. The relative abundance of secretion proteins in the 27
151 microbiomes was determined as well as that of the toxins which are essential to virulence processes. The results
152 of the present study showed that the abundance of secretion proteins selected in the database was strongly
153 correlated with the toxin abundance ($R^2 = 0.80$, Figure 1). The most abundant secretion proteins included *L.*
154 *waltersii* toxin secretion protein (LWT1SS), *L. pneumophila* toxin secretion protein ApxIB, and *Aeromonas*
155 *hydrophila* RTX transporter (RtxB) (data not shown). Further exploration indicated that although *A. hydrophila*
156 RtxB homologs from GenBank were found in all *Proteobacteria* classes, most of the RtxB-harboring species
157 have been reported to be pathogens, including *Vibrio* spp. [64], *Pseudomonas* spp., *Neisseria meningitidis* [65],
158 *Ralstonia* spp. [66], and *Yersinia* spp. [21]. This may imply the pathogen-specific nature of secretion proteins
159 included in the database, and that toxin secretion proteins can be used as signatures for pathogen detection as
160 well.

161 Toxin-phyla CCA results showed that all phyla can be clearly separated into two groups, and that almost all
162 toxins were associated with *Proteobacteria*, *Nitrospirae* and *Firmicutes* (Figure 2). Considering the phylum-
163 specificity of the toxins stated above, these results can be biased because of the taxonomic affiliation of toxins
164 included in the Local BLASTP database. The taxonomic distribution proportion of currently available genomes
165 of identified pathogens was reflected in the toxin database, with *Proteobacteria* and *Firmicutes* accounting for
166 the majority of the genomes. However, the CCA results may also indicate, at least in part, a proportional lack of
167 pathogens in some phyla, such as *Crenarchaeota*, *Euryarchaeota*, *Verrucomicrobia* and *Bacteroidetes* [67].
168 Archaea cannot easily absorb phage particles because of their extracellular structures, which differ from bacteria
169 [68]. A recent study by Li et al. also found that the five most abundant bacterial pathogens were from either
170 *Proteobacteria* or *Firmicutes* in wastewater microbiomes [9]. Taken together, these findings could indicate that
171 *Proteobacteria* or *Firmicutes* were evolutionarily enriched with pathogens when they dominated most
172 environmental microbiomes on the planet [69, 70].

173 Interestingly, there was a strong association between the phylum *Nitrospirae* and toxins of RNase
174 inteferases (MvpA and MapC) and *Listeria monocytogenes* 1-phosphatidylinositol phosphodiesterase PLC.

175 Further searches against the UniProt database [71] revealed no homologous records of MvpA and PLC from
176 *Nitrospirae*, and only 109 out of 15,574 bacterial records for VapC were from *Nitrospirae*. These findings imply
177 that there are many more *Nitrospirae* pathogens harboring MvpA and PLC that have yet to be discovered.

178 The screening of toxins in the 27 global microbiomes revealed the most prevalent toxins and pathogen-
179 enriched environment. Specifically, the results showed that the RTX toxin RtxA and adenylate cyclase Cya were
180 most prevalent globally in terms of relative abundance. RTX toxins comprise a large family of pore-forming
181 exotoxins. Known homologs in the GenBank database of *Aeromonas dhakensis* RtxA were mainly in the genera
182 of *Aeromonas*, *Pseudomonas* (e.g., CP015992), *Vibrio* (e.g., CP002556) and *Legionella* (e.g., CP015953). These
183 genera are well known to be associated with gastroenteritis, eye and wound infections, cholera and legionellosis,
184 and RTX toxins are a key part of the virulence systems of each of these conditions [72-75]. Cya is an essential
185 unit of *Bacillus anthracis* virulence that causes anthrax and may lead to mammalian death [76]. Known
186 homologs in the GenBank database of *Bacillus anthracis* Cya were mainly from *Bacillus* spp., *Bordetella* spp.,
187 *Pseudomonas aeruginosa*, *Yersinia pseudotuberculosis*, and *Vibrio* spp. Their presence in the environment
188 should be carefully examined and precautions should be taken to prevent infection by these organisms since
189 many of them are associated with very common diseases such as whooping cough.

190 The main purpose of the Local BLASTP method established here was to screen pathogen-enriched
191 environments to enable development of precautionary measures. Our results clearly indicated that contaminated
192 lake water, feces and wastewater microbiomes were rich in pathogens (Figure 3). Although there was no
193 detailed background information regarding these environments in this study, the results presented herein may
194 provide important implications for pathogen-related risk control. Surprisingly, two lake water microbiomes from
195 Nanjing, China contained the highest toxin factors among the 27 samples. Further investigation of the location
196 and contamination status supported the sewage-nature of the lake water. In China, most polluted lakes receive
197 sewage that includes feces materials [77]. According to an official survey conducted in 2015, Nanjing has 28
198 lakes with a total area of 14 km², among which 96.7% are classified as polluted (Class V of the national
199 standard). Studies have documented that pathogens tend to be enriched in polluted waters [14]. It is not
200 surprising to find that feces samples had very high abundance of toxins. Epidemical statistics have indicated that
201 feces are the most important pathway for diarrheal diseases, which is a leading cause of childhood death
202 globally [78]. Thus, the present study provides a method for obtaining quantitative estimates of pathogen
203 enrichment of various environments, and polluted freshwater systems are found to be highly pathogen-enriched
204 relative to safer environments such as ocean water and natural soils.

205 **Conclusions**

206 A Local BLASTP procedure was established for rapid detection of toxins in environmental samples.
207 Screening of global microbiomes in this study provided a quantitative estimate of the most prevalent toxins and
208 most pathogen-enriched environments.

209 **Declarations**

210 **Acknowledgement**

211 I thank Dr. Philip L. Bond and The University of Queensland for providing training in bioinformatics. We would
212 like to thank LetPub (www.letpub.com) for providing linguistic assistance during the preparation of this
213 manuscript. I also thank the founders of the existing pathogen-relevant database, particularly the Virulence
214 Factor Database, which provided valuable reference for the buildup of the toxin database in this study.

215 *Funding*

216 This project was financially supported by the pioneer “Hundred Talents Program” of the Chinese Academy of
217 Sciences and the Hebei Science Fund for Distinguished Young Scholars (D2018503005).

218 *Competing interests*

219 The author declares no conflict of interest.

220 *Availability of data and materials*

221 The toxin database is available in the Supplementary Materials. All toxin abundance data in this study can be
222 provided by the author upon request.

223 *Authors' contributions*

224 XL initiated the study, analyzed the data and wrote the manuscript.

225 *Ethics approval*

226 Not applicable.

227 *Consent for publication*

228 Not applicable.

229 **References**

230 [1] Baumgardner DJ. Soil-related bacterial and fungal infections. *J Am Board Fam Med*, 2012, 25: 734-
231 744

232 [2] Jeffery S, van der Putten WH, Ispra, Italy: Joint Research Centre, Institute for Environment and
233 Sustainability 2011

234 [3] Whiley H, Bentham R. *Legionella longbeachae* and legionellosis. *Emerging Infectious Diseases*, 2011,
235 17: 579-583

236 [4] Wait DA, Sobsey MD. Comparative survival of enteric viruses and bacteria in Atlanticocean seawater.
237 *Water Sci Technol*, 2001, 43: 139-142

238 [5] Dan TBB, Wynne D, Manor Y. Survival of enteric bacteria and viruses in lake kinneret, Israel. *Water*
239 *Res*, 1997, 31: 2755-2760

240 [6] Keswick BH, Gerba CP, Secor SL, et al. Survival of enteric viruses and indicator bacteria in

- 241 groundwater. *J Environ Sci Heal A*, 1982, 17: 903-912
- 242 [7] Cooper RC, Golueke CG. Survival of enteric bacteria and viruses in compost and its leachate. *Compost*
243 *Sci Land Ut*, 1979, 20: 29-35
- 244 [8] Miller RR, Montoya V, Gardy JL, et al. Metagenomics for pathogen detection in public health. *Genome*
245 *Med*, 2013, 5:
- 246 [9] Li B, Ju F, Cai L, et al. Profile and fate of bacterial pathogens in sewage treatment plants revealed by
247 high-throughput metagenomic approach. *Environ Sci Technol*, 2015, 49: 10492-10502
- 248 [10] Amha YM, Anwar MZ, Kumaraswamy R, et al. *Mycobacteria* in municipal wastewater treatment and
249 reuse: Microbial diversity for screening the occurrence of clinically and environmentally relevant species in arid
250 regions. *Environ Sci Technol*, 2017, 51: 3048-3056
- 251 [11] Yang J, Yang F, Ren LL, et al. Unbiased parallel detection of viral pathogens in clinical samples by use
252 of a metagenomic approach. *J Clin Microbiol*, 2011, 49: 3463-3469
- 253 [12] Granberg F, Vicente-Rubiano M, Rubio-Guerri C, et al. Metagenomic detection of viral pathogens in
254 spanish honeybees: Co- infection by aphid lethal paralysis, Israel acute paralysis and lake sinai viruses. *Plos*
255 *One*, 2013, 8:
- 256 [13] Nakamura S, Yang CS, Sakon N, et al. Direct metagenomic detection of viral pathogens in nasal and
257 fecal specimens using an unbiased high-throughput sequencing approach. *Plos One*, 2009, 4: e4219
- 258 [14] Bibby K. Metagenomic identification of viral pathogens. *Trends Biotechnol*, 2013, 31: 11-15
- 259 [15] Fukui Y, Aoki K, Okuma S, et al. Metagenomic analysis for detecting pathogens in culture-negative
260 infective endocarditis. *J Infect Chemother*, 2015, 21: 882-884
- 261 [16] Baldwin DA, Feldman M, Alwine JC, et al. Metagenomic assay for identification of microbial
262 pathogens in tumor tissues. *Mbio*, 2014, 5: 01714-14
- 263 [17] Achtman M, Wagner M. Microbial diversity and the genetic nature of microbial species. *Nat Rev*
264 *Micro*, 2008, 6: 431-440
- 265 [18] Liao PY, Lee KH. From snps to functional polymorphism: The insight into biotechnology applications.
266 *Biochem Eng J*, 2010, 49: 149-158
- 267 [19] Sokurenko EV, Hasty DL, Dykhuizen DE. Pathoadaptive mutations: Gene loss and variation in
268 bacterial pathogens. *Trends Microbiol*, 1999, 7: 191-195
- 269 [20] Strauch E, Lurz R, Beutin L. Characterization of a shiga toxin-encoding temperate bacteriophage of
270 *Shigella sonnei*. *Infect Immun*, 2001, 69: 7588-7595
- 271 [21] Chen LH, Zheng DD, Liu B, et al. Vfdb 2016: Hierarchical and refined dataset for big data analysis-10
272 years on. *Nucleic Acids Res*, 2016, 44: D694-D697
- 273 [22] US Water Environment Federation. Manure pathogens: Manure management, regulations, and water
274 quality protection. Water Environment Federation Alexandria, Virginia: WEF Press, 2009
- 275 [23] Bateman A, Martin MJ, O'Donovan C, et al. Uniprot: The universal protein knowledgebase. *Nucleic*
276 *Acids Res*, 2017, 45: D158-D169
- 277 [24] Hoffmann R. A wiki for the life sciences where authorship matters. *Nature Genetics*, 2008, 40: 1047-
278 1051
- 279 [25] Lesnick ML, Reiner NE, Fierer J, et al. The *Salmonella* *spvB* virulence gene encodes an enzyme that
280 ADP-ribosylates actin and destabilizes the cytoskeleton of eukaryotic cells. *Mol Microbiol*, 2001, 39: 1464-

- 281 1470
- 282 [26] Skopova K, Tomalova B, Kanchev I, et al. Camp-elevating capacity of the adenylate cyclase toxin-
283 hemolysin is sufficient for lung infection but not for full virulence of *Bordetella pertussis*. *Infect Immun*, 2017,
- 284 [27] Henner DJ, Yang M, Chen E, et al. Sequence of the *Bacillus thuringiensis* phosphatidylinositol specific
285 phospholipase-c. *Nucleic Acids Res*, 1988, 16: 10383-10383
- 286 [28] Schwarzenbacher R, Stenner-Liewen F, Liewen H, et al. Structure of the chlamydia protein Cdd
287 reveals a redox enzyme that modulates host cell apoptosis. *J Biol Chem*, 2004, 279: 29320-29324
- 288 [29] Hamon MA, Batsche E, Regnault B, et al. Histone modifications induced by a family of bacterial
289 toxins. *P Natl Acad Sci USA*, 2007, 104: 13467-13472
- 290 [30] Cossart P. The listeriolysin o-gene - a chromosomal locus crucial for the virulence of *Listeria*
291 *monocytogenes*. *Infection*, 1988, 16: S157-S159
- 292 [31] Geoffroy C, Mengaud J, Alouf JE, et al. Alveolysin, the thiol-activated toxin of *Bacillus alvei*, is
293 homologous to listeriolysin o, perfringolysin o, pneumolysin, and streptolysin o and contains a single cysteine. *J*
294 *Bacteriol*, 1990, 172: 7301-7305
- 295 [32] Rossjohn J, Polekhina G, Feil SC, et al. Structures of perfringolysin o suggest a pathway for activation
296 of cholesterol-dependent cytolysins. *J Mol Biol*, 2007, 367: 1227-1236
- 297 [33] Lyras D, O'Connor JR, Howarth PM, et al. Toxin b is essential for virulence of *Clostridium difficile*.
298 *Nature*, 2009, 458: 1176-1181
- 299 [34] Schmidt H, Scheef J, JanetzkiMittmann C, et al. An ileX trnA gene is located close to the shiga toxin ii
300 operon in enterohemorrhagic *Escherichia coli* O157 and non-O157 strains. *FEMS Microbiol Lett*, 1997, 149: 39-
301 44
- 302 [35] Labandeira-Rey M, Couzon F, Boisset S, et al. *Staphylococcus aureus* panton-valentine leukocidin
303 causes necrotizing pneumonia. *Science*, 2007, 315: 1130-1133
- 304 [36] Bukowski M, Wladyka B, Dubin G. Exfoliative toxins of *Staphylococcus aureus*. *Toxins*, 2010, 2:
305 1148-1165
- 306 [37] Aoki SK, Diner EJ, de Roodenbeke CT, et al. A widespread family of polymorphic contact-dependent
307 toxin delivery systems in bacteria. *Nature*, 2010, 468: 439-442
- 308 [38] Tian QB, Ohnishi M, Tabuchi A, et al. A new plasmid-encoded proteic killer gene system: Cloning,
309 sequencing, and analysing *hig* locus of plasmid rts1. *Biochem Biophys Res Commun*, 1996, 220: 280-284
- 310 [39] Hurley JM, Woychik NA. Bacterial toxin HigB associates with ribosomes and mediates translation-
311 dependent mRNA cleavage at a-rich sites. *J Biol Chem*, 2009, 284: 18605-18613
- 312 [40] Pullinger GD, Lax AJ. A *Salmonella*-dublin virulence plasmid locus that affects bacterial-growth under
313 nutrient-limited conditions. *Mol Microbiol*, 1992, 6: 1631-1643
- 314 [41] Korch SB, Contreras H, Clark-Curtiss JE. Three *Mycobacterium tuberculosis* Rel toxin-antitoxin
315 modules inhibit mycobacterial growth and are expressed in infected human macrophages. *J Bacteriol*, 2009, 191:
316 1618-1630
- 317 [42] Pellegrini O, Mathy N, Gogos A, et al. The *Bacillus subtilis* ydcde operon encodes an endoribonuclease
318 of the *mazF/pemK* family and its inhibitor. *Molecular Microbiology*, 2005, 56: 1139-1148
- 319 [43] Yamaguchi Y, Inouye M. Regulation of growth and death in *Escherichia coli* by toxin-antitoxin
320 systems. *Nature Rev Microbiol*, 2011, 9: 779-790

- 321 [44] Hatziloukas E, Panopoulos NJ. Origin, structure, and regulation of *argk*, encoding the phaseolotoxin-
322 resistant ornithine carbamoyltransferase in *Pseudomonas syringae* pv. *Phaseolicola*, and functional expression of
323 *argk* in transgenic tobacco. *J Bacteriol*, 1992, 174: 5895-5909
- 324 [45] Yates SP, Merrill AR. Elucidation of eukaryotic elongation factor-2 contact sites within the catalytic
325 domain of *Pseudomonas aeruginosa* exotoxin a. *Biochem J*, 2004, 379: 563-572
- 326 [46] Songer JG. Bacterial phospholipases and their role in virulence. *Trends Microbiol*, 5: 156-161
- 327 [47] Krueger KM, Barbieri JT. The family of bacterial ADP-ribosylating exotoxins. *Clin Microbiol Rev*,
328 1995, 8: 34-47
- 329 [48] Phillips RM, Six DA, Dennis EA, et al. In vivo phospholipase activity of the *Pseudomonas aeruginosa*
330 cytotoxin *exoU* and protection of mammalian cells with phospholipase a2 inhibitors. *J Biol Chem*, 2003, 278:
331 41326-41332
- 332 [49] Weiss AA, Johnson FD, Burns DL. Molecular characterization of an operon required for pertussis toxin
333 secretion. *P Natl Acad Sci USA*, 1993, 90: 2970-2974
- 334 [50] Masuzawa T, Sawaki K, Nagaoka H, et al. Relationship between pathogenicity of *Coxiella burnetii*
335 isolates and gene sequences of the macrophage infectivity potentiator (*cbmip*) and sensor-like protein (*qrsa*).
336 *FEMS Microbiol Lett*, 1997, 154: 201-205
- 337 [51] D'Auria G, Jimenez N, Peris-Bondia F, et al. Virulence factor *rtx* in *Legionella pneumophila*, evidence
338 suggesting it is a modular multifunctional protein. *Bmc Genomics*, 2008, 9:
- 339 [52] Rasmussen-Ivey CR, Figueras MJ, McGarey D, et al. Virulence factors of *Aeromonas hydrophila*: In
340 the wake of reclassification. *Frontiers in Microbiology*, 2016, 7:
- 341 [53] Howard SP, Garland WJ, Green MJ, et al. Nucleotide sequence of the gene for the hole-forming toxin
342 aerolysin of *Aeromonas hydrophila*. *J Bacteriol*, 1987, 169: 2869-2871
- 343 [54] Sandkvist M, Michel LO, Hough LP, et al. General secretion pathway (*eps*) genes required for toxin
344 secretion and outer membrane biogenesis in *Vibrio cholerae*. *J Bacteriol*, 1997, 179: 6994-7003
- 345 [55] Lioy VS, Machon C, Tabone M, et al. The ζ toxin induces a set of protective responses and dormancy.
346 *Plos One*, 2012, 7: e30282
- 347 [56] Söderberg MA, Rossier O, Cianciotto NP. The type II protein secretion system of *Legionella*
348 *pneumophila* promotes growth at low temperatures. *J Bacteriol*, 2004, 186: 3712-3720
- 349 [57] Danilchanka O, Pires D, Anes E, et al. The *Mycobacterium tuberculosis* outer membrane channel
350 protein *Cpnt* confers susceptibility to toxic molecules. *Antimicrob Agents Ch*, 2015, 59: 2328-2336
- 351 [58] Li X, Zhu YG, Shaban B, et al. Assessing the genetic diversity of *cu* resistance in mine tailings through
352 high-throughput recovery of full-length *copA* genes. *Sci Rep*, 2015, 5: 13258
- 353 [59] Kumar S, Stecher G, Tamura K. *Mega7*: Molecular evolutionary genetics analysis version 7.0 for
354 bigger datasets. *Mol Biol Evol*, 2016, 33: 1870-1874
- 355 [60] Krzywinski M, Schein J, Birol I, et al. *Circos*: An information aesthetic for comparative genomics.
356 *Genome Res*, 2009, 19: 1639-1645
- 357 [61] R Core Team. R: A language and environment for statistical computing. 2016, URL [http://www.R-](http://www.R-project.org/)
358 [project.org/](http://www.R-project.org/).
- 359 [62] Li XF, Bond PL, Huang LB. Diversity of as metabolism functional genes in Pb-Zn mine tailings.
360 *Pedosphere*, 2017, 27: 630-637

- 361 [63] Gupta SK, Padmanabhan BR, Diene SM, et al. Arg-annot, a new bioinformatic tool to discover
362 antibiotic resistance genes in bacterial genomes. *Antimicrob Agents Ch*, 2014, 58: 212-220
- 363 [64] Austin B, Zhang XH. *Vibrio harveyi*: A significant pathogen of marine vertebrates and invertebrates.
364 *Lett Appl Microbiol*, 2006, 43: 119-124
- 365 [65] Rouphael NG, Stephens DS. *Neisseria meningitidis*: Biology, microbiology, and epidemiology.
366 *Methods in molecular biology* (Clifton, NJ), 2012, 799: 1-20
- 367 [66] Xu J, Zheng HJ, Liu L, et al. Complete genome sequence of the plant pathogen *Ralstonia solanacearum*
368 strain po82. *J Bacteriol*, 2011, 193: 4261-4262
- 369 [67] Ecker DJ, Sampath R, Willett P, et al. The microbial rosetta stone database: A compilation of global and
370 emerging infectious microorganisms and bioterrorist threat agents. *Bmc Microbiol*, 2005, 5:
- 371 [68] Gill EE, Brinkman FSL. The proportional lack of archaeal pathogens: Do viruses/phages hold the key?
372 *Bioessays*, 2011, 33: 248-254
- 373 [69] Roesch LF, Fulthorpe RR, Riva A, et al. Pyrosequencing enumerates and contrasts soil microbial
374 diversity. *Isme J*, 2007, 1: 283-290
- 375 [70] Fierer N, Bradford MA, Jackson RB. Toward an ecological classification of soil bacteria. *Ecology*, 2007,
376 88: 1354-1364
- 377 [71] Apweiler R, Bairoch A, Wu CH, et al. Uniprot: The universal protein knowledgebase. *Nucleic Acids*
378 *Res*, 2004, 32: D115-D119
- 379 [72] Lin W, Fullner KJ, Clayton R, et al. Identification of a *Vibrio cholerae* rtx toxin gene cluster that is
380 tightly linked to the cholera toxin prophage. *P Natl Acad Sci USA*, 1999, 96: 1071-1076
- 381 [73] Cirillo SLG, Bermudez LE, El-Etr SH, et al. *Legionella pneumophila* entry gene *rtxa* is involved in
382 virulence. *Infect Immun*, 2001, 69: 508-517
- 383 [74] Suarez G, Khajanchi BK, Sierra JC, et al. Actin cross-linking domain of *Aeromonas hydrophila* repeat
384 in toxin a (*rtxa*) induces host cell rounding and apoptosis. *Gene*, 2012, 506: 369-376
- 385 [75] Terada LS, Johansen KA, Nowbar S, et al. *Pseudomonas aeruginosa* hemolytic phospholipase c
386 suppresses neutrophil respiratory burst activity. *Infect Immun*, 1999, 67: 2371-2376
- 387 [76] Leppla SH. Anthrax toxin edema factor - a bacterial adenylate-cyclase that increases cyclic-amp
388 concentrations in eukaryotic cells. *P Natl Acad Sci-Biol*, 1982, 79: 3162-3166
- 389 [77] Qiu Z. Current pollution status of china's lakes. *The 5th Forum for China Lakes*, 2015, 5.
- 390 [78] Liu L, Johnson HL, Cousens S, et al. Global, regional, and national causes of child mortality: An
391 updated systematic analysis for 2010 with time trends since 2000. *Lancet*, 2012, 379: 2151-2161

392

393 **Figure captions**

394 Figure 1 Correlation between relative abundance of toxins and secretion proteins in the global microbiomes.

395 Figure 2 Canonical correspondence analysis of the associations between phyla and toxins.

396 Figure 3 Circular visualization of the toxin abundance in the microbiomes selected from locations worldwide.

397 The designated environment was prefixed with the first letters of the environment names and suffixed with the

398 last three numbers of their MG-RAST ID in Table 1.

Table 1 General information regarding the metagenomes retrieved from the MG-RAST server

MG-RAST ID	bp Count	Sequence Count	Material	Location	Country	Coordinates	Sequencing Method	Protein Predicted
mgm4713202	89,726,117	442,552	Water	Moorea	Pacific Ocean	17.538, -149.829	Illumina	254,139
mgm4713205	106,474,596	476,363	Water	Moorea	Pacific Ocean	17.538, -149.829	Illumina	235,777
mgm4713197	60,417,678	272,918	water	Moorea	Pacific Ocean	17.538, -149.829	Illumina	140,195
mgm4719940	360,335,259	1,425,556	Water	Irish Sea	Atlantic Ocean	53.225, -4.159	Illumina	1,023,445
mgm4533707	134,833,790	668,257	Brackish Water	Columbia River margin, OR & WA coastal	USA	46.265, -123.999	Illumina	508,217
mgm4720261	35,487,527	6,896	Fresh Water	Galway	Ireland	53.276, -9.060	Illumina	33,743
mgm4695622	114,430,648	111,889	Fresh Water	Nanjing, Jiangsu	China	32.600, 118.160	Illumina	148,833
mgm4695626	86,732,360	78,621	Fresh Water	Nanjing, Jiangsu	China	32.400, 118.140	Illumina	111,489
mgm4589537	337,068,782	2,099,471	Surface Water	West Virginia	USA	38.094, -81.959	Illumina	1,842,975
mgm4679248	603,919,746	3,365,512	Soil	Seoul	South Korea	37.460, 126.948	Illumina	1,361,948
mgm4679254	689,019,062	3,688,750	Soil	Seoul	South Korea	37.459, 126.948	Illumina	1,966,121
mgm4514299	322,114,449	242,0832	Saline Desert Soil	Gujarat of India	India	23.7925, 71.008	Ion Torrent	1,323,378
mgm4543019	282,578,916	2,016,127	Saline Desert Soil	Gujarat of India	India	23.908, 70.538	Ion Torrent	842,475
mgm4697397	143,214,978	397,067	Organic Soil	Beijing	China	32.054, 118.763	Illumina	299,940
N/A	507,124,889	1,552,234	Shrub Land Soil	Mt Isa	Australia	20.440, 139.300	Illumina	1,413,889
N/A	532,850,584	1,632,914	Red Soil	Mt Isa	Australia	20.440, 139.300	Illumina	1,408,943
N/A	433,386,397	1,338,665	Red Soil Polluted	Mt Isa	Australia	20.440, 139.300	Illumina	1,081,822
mgm4507016	163,648,718	227,551	Feces	Bologna	Italy	44.495, 11.343	Illumina	250347
mgm4718752	329,518,322	1,312,822	Feces	Upstate NY	USA	42.668, -76.528	Illumina	950489
mgm4568577	10,065,266	50,137	Mine Water	Guangdong	China	24.503, 113.710	454	34,287
mgm4568580	12,911,442	62,018	Mine Water	Guangdong	China	22.940, 112.050	454	36,461
mgm4620491	52,759,415	244,855	Biosolides	Guelph ON Canada	Canada	43.545, -80.248	Illumina	238630
mgm4546371	84,424,005	907,785	Wastewater	Universiti Teknologi Malaysia	Malaysia	2.558, 104.642	Illumina	803,682
mgm4620487	147,523,219	696,132	Wastewater	Guelph ON Canada	Canada	43.5448, -80.248	Illumina	640,283
mgm4620488	115,131,556	578,337	Wastewater	Guelph ON Canada	Canada	43.545, -80.248	Illumina	537,267
mgm4560423	22734940	73479	Feces	Lake Eyasi, Tanzania	Tanzania	-3.635, 35.083	Illumina	76,569

mgm4440281	35439683	334386	Mine drainage	Soudan Mine	United States of America	47.819, -92.243	454	227,038
------------	----------	--------	---------------	-------------	-----------------------------	-----------------	-----	---------

Table 2 Typical virulence factors investigated in this study and their disease–relevance

Toxin factor	ID in the database	Typical pathogens and disease	Role of the toxin	Reference
Mono(ADP-ribosyl)transferase	spvB	<i>Salmonella dublin</i> (gastroenteritis)	ADP-ribosylating, destabilizing cytoskeleton	[25]
Adenylate cyclase	cyaA	<i>Sa. choleraesuis</i> (typhoid fever) <i>Bordetella pertussis</i> (whooping cough) <i>Bacillus anthracis</i> (anthrax)	Pore-forming with cAMP-elevating activity	[26]
1-phosphatidylinositol phosphodiesterase	PLC	<i>Listeria monocytogenes</i> (listeriosis)	Lysis of the phagolysosomal membrane	[27]
Chlamydia protein associating with death domains	CADD	<i>Chlamydia trachomatis</i> (trachoma, urethritis, etc.)	Inducing cell apoptosis	[28]
Listeriolysin O	hly	<i>Li. monocytogenes</i> (listeriosis)	Pore forming, hemolysin	[29, 30]
Alveolysin	alo	<i>Ba. anthracis</i>	Pore-forming	[31]
Perfringolysin O	pfo	<i>Clostridium perfringens</i> (food poisoning)	Pore-forming	[32]
Glucosyltransferase toxin B	toxB	<i>Cl. sordellii</i> <i>Cl. difficile</i> (diarrhea)	Cytopathic effects	[33]
Shiga toxin 1	stx1	<i>Escherichia coli</i> (diarrhea) <i>Shigella dysenteriae</i> (Shigellosis)	Haemolytic uraemic syndrome	[20]
Shiga-like toxin 2	stx2	<i>Enterobacteria phage 933W</i> <i>E. coli</i>	Haemolytic uraemic syndrome	[34]
Leucotoxin	luk	<i>Staphylococcus aureus</i> (sinusitis, skin abscess)	Lysis of leukocytes	[35]
Exfoliative toxin	ET	<i>St. aureus</i>	Proteolytic activity	[36]
Toxin CdiA	cdiA	<i>E. coli</i> <i>Yersinia pestis</i> (plaque)	Decreasing aerobic respiration	[37]

RNA interferase	mazF/pemK/ndoA/relE/relK/relG/yoeB/higB/mvpA	<i>Proteus vulgaris</i> (wound infections) <i>Mycobacterium tuberculosis</i> (tuberculosis) <i>E. coli</i>	and ATP levels Cleavage of cellular mRNAs, inhibiting growth tRNase activities	[38-43]
Endonuclease VapC	vapC	<i>Shigella flexneri</i> (diarrhea) <i>Sa. Dublin</i> <i>My. Tuberculosis</i> <i>Coxiella burnetii</i> (Q fever) <i>Pseudomonas savastanoi</i>		[40]
Ornithine carbamoyltransferase	argK		Promoting survival and pathogenicity	[44]
Exotoxin A	ETA	<i>Ps. aeruginosa</i> (eye and wound infections)	ADP-ribosylating eukaryotic elongation factor 2	[45]
Hemolytic phospholipase C	plcH	<i>Ps. Aeruginosa</i> <i>Clostridium perfringens</i> (food poisoning)	Membrane-damaging	[46]
ADP-ribosyltransferase toxin	exoS	<i>P. aeruginosa</i>	Inhibition of phagocytosis	[47]
Exoenzyme U	exoU	<i>P. aeruginosa</i>	Membrane-lytic and cytotoxic	[48]
Dermonecrotic toxin	dnt	<i>Bo. pertussis</i>	Stimulating the assembly of actin stress fibers and focal adhesions	[49, 50]
Pertussis toxin subunit 1	ptxA	<i>Bo. pertussis</i>	Causing disruption of host cellular regulation	[49]

Repeats-in toxin	rtxA	<i>Legionella pneumophila</i> (Legionnaires' disease) <i>Aeromonas dhakensis</i> (gastroenteritis, septicemia) <i>Aeromonas</i> spp.	Adherence and pore forming	[51, 52]
Aerolysin	aerA	<i>Le. pneumophila</i>	Cytolytic pore-forming	[53]
Cholera toxin secretion protein EpsF	epsF	<i>Le. pneumophila</i>	Toxin secretion	[54]
Zeta toxin family protein	ZETA	<i>Coxiella</i> sp. DG_40	Inhibiting cell wall biosynthesis	[55]
Toxin secretion ATP binding protein	LwT1SS	<i>Le. waltersii</i>	Toxin secretion	[56]
Outer membrane channel protein CpnT	cpnT	<i>My. tuberculosis</i>	Nutrient uptake	[57]
Type IV secretion system protein Ptl	ptlCH	<i>Bo. pertussis</i>	Secretion of pertussis toxin	[49]

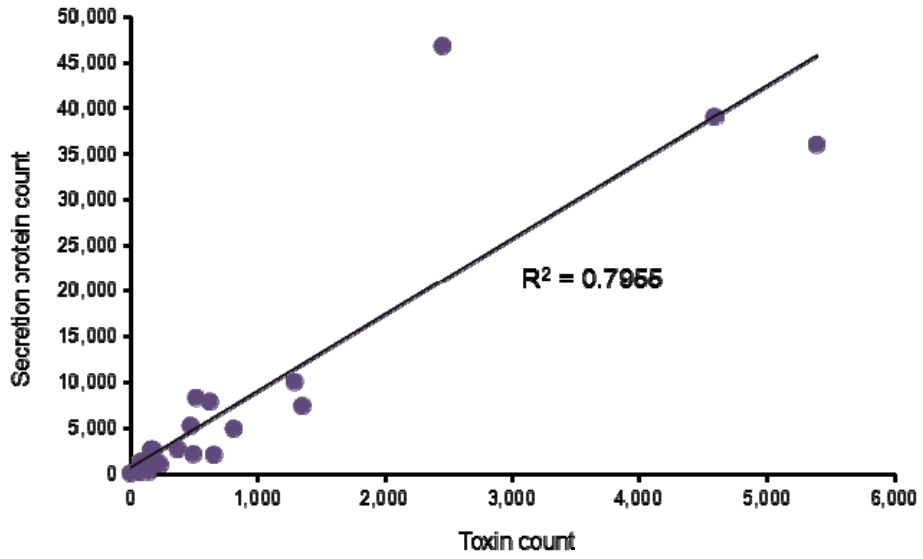


Fig 1

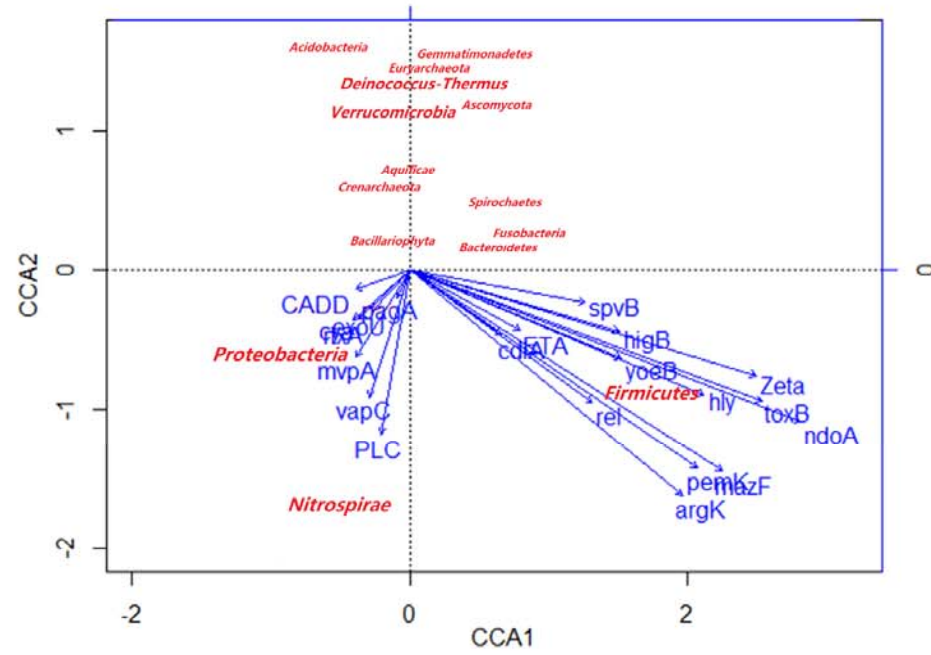


Fig 2

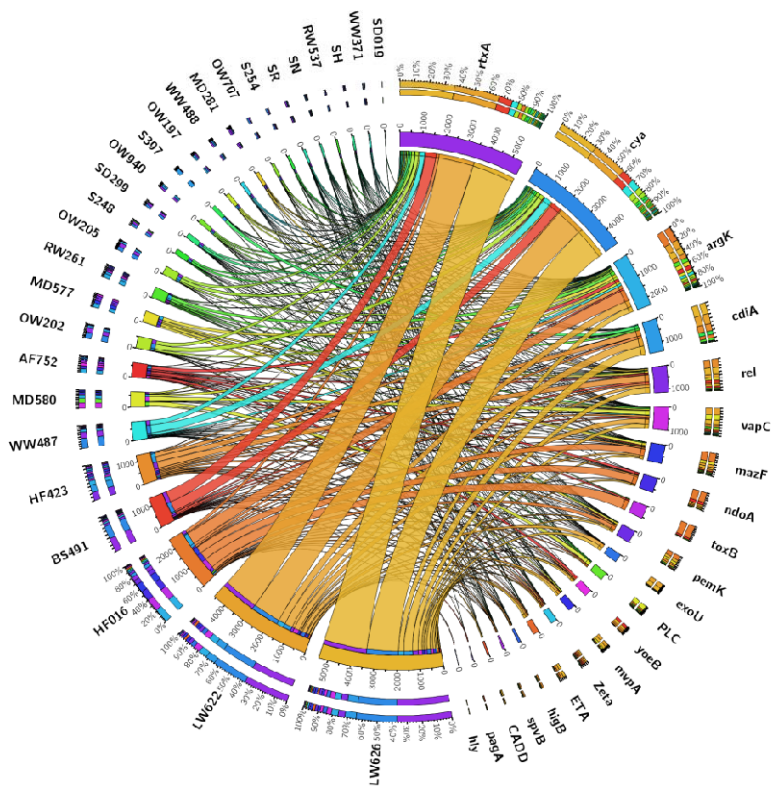


Fig 3