

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Whole genome analysis reveals pathogenic potential of multi-drug resistant wastewater

Escherichia coli

Norhan Mahfouz^{1,*}, Serena Caucci^{2,3,*}, Eric Achatz¹, Torsten Semmler⁴, Sebastian Guenther⁴,
Thomas U. Berendonk^{2,*}, and Michael Schroeder^{1,*,#}

¹ Biotec, TU Dresden

² Institute for Hydrobiology, TU Dresden

³ United Nations University Institute for Integrated Management of Material Fluxes and of
Resources

⁴ Institute of Microbiology und Epizootics, FU Berlin

* These authors contributed equally

Correspondence: Michael Schroeder, ms@biotec.tu-dresden.de

Keywords: Antibiotic Resistance, Wastewater Treatment, Pan-Core genome, Environment

Conflict of interest statement: The authors declare no conflict of interest.

20

21 **Abstract**

22 Wastewater treatment plants play an important role in antibiotic resistance development. While it
23 has been shown that wastewater effluents contain resistant bacteria, resistance genes, and
24 antibiotics, there is little knowledge on the link between resistance genotype and phenotype.
25 Here we present the first study, which combines a culture-based phenotypic screen with the
26 analysis of whole genome sequences for the indicator species *Escherichia coli* of the inflow and
27 outflow of a sewage treatment plant. Our analysis reveals that nearly all isolates are multi-drug
28 resistant and many are potentially pathogenic. This holds in particular for the outflow of the
29 treatment plant. We devise a computational approach correlating genotypic variation and
30 resistance phenotype, which identifies known and candidate resistance genes. The identified
31 genes stem from the pan genome, which is large and thus reflects the genomic heterogeneity of a
32 treatment plant. Overall, the screen and analysis show that sewage treatment plants provide a
33 favourable environment for antibiotic resistance development and that resistant bacteria do not
34 appear to suffer from a competitive disadvantage in wastewater. These findings should find
35 consideration in future improvements of wastewater treatment.

36

37

38

39

40 **Introduction**

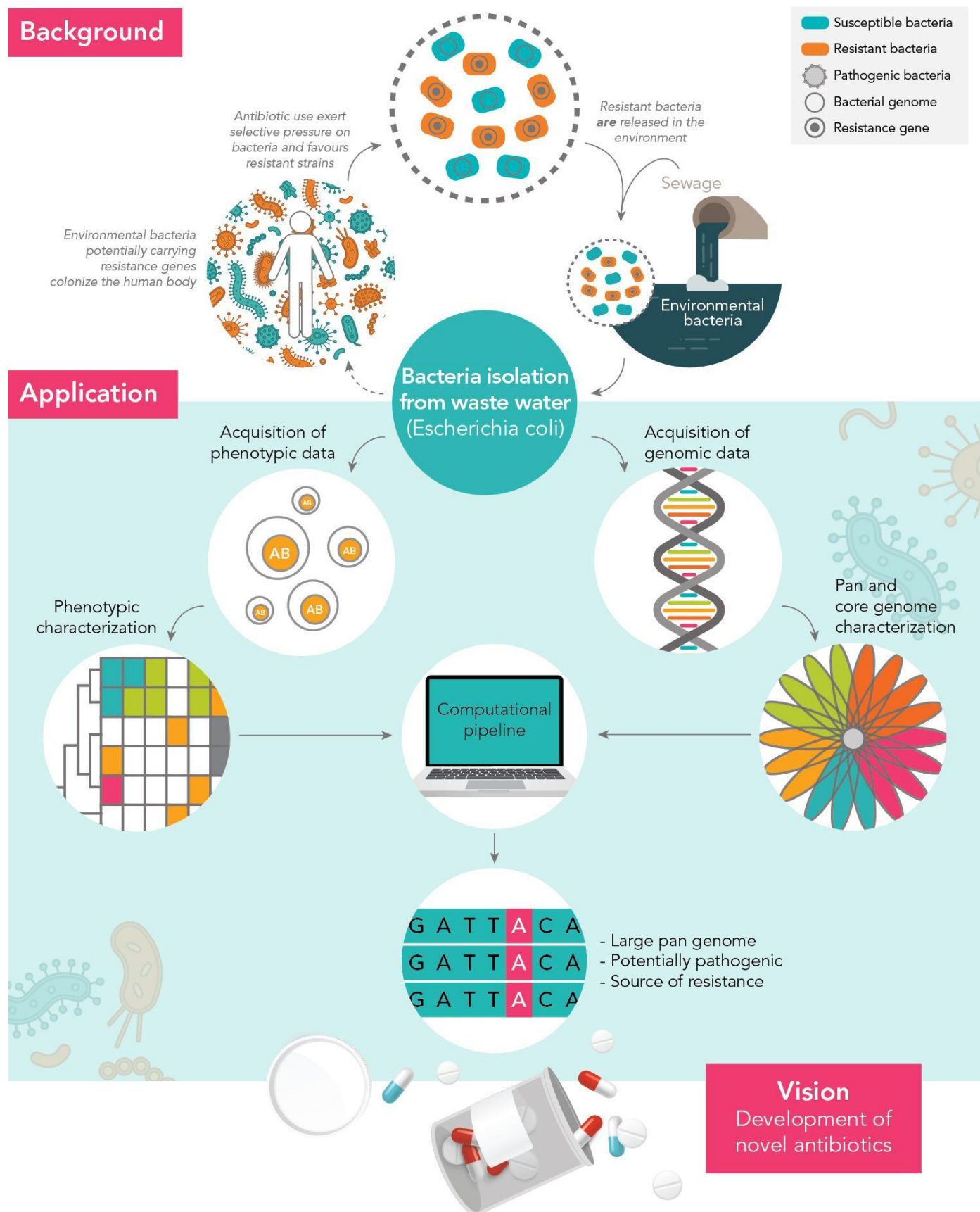
41 In 1945, Alexander Fleming, the discoverer of Penicillin, warned of antibiotic resistance. Today,
42 the WHO echoes this warning, calling antibiotic resistance a global threat to human health.

43 Humans are at the center of the modern rise of resistance. The human gut (1), clinical samples
44 (2, 3), soil (4, 5), and wastewater (6) all harbor resistant bacteria and resistance genes. At the
45 heart of modern resistance development is a human-centered network of clinics, industry, private
46 homes, farming, and wastewater. However, it is unclear, where antibiotic resistance emerges and
47 in particular, there are contradictory views on the role of wastewater treatment plants (6). On the
48 one hand, the harsh environment of a treatment plant appears unfavourable for resistant bacteria
49 (7), but on the other hand, it forms a very rich genetic reservoir, where highly diverse bacteria
50 mingle (6). Also, it is unclear, whether or not any bacteria with pathogenic potential emerge after
51 treatment.

52 To address these questions, we collected 1178 *Escherichia coli* isolates from a waste treatment
53 plant's inflow and outflow in the city of Dresden, Germany. We selected 20 antibiotics, which are
54 the most prescribed ones in the area from which the wastewater inflow originates (data provided
55 by the public health insurer AOK). We analyzed the isolates' resistance to these 20 antibiotics
56 and selected 103 isolates for whole genome sequencing. Our analysis reveals that wastewater
57 outflow harbors multi-drug resistant *Escherichia coli* with pathogenic potential and very flexible
58 genomes harboring resistance genes.

59

60



61

62 Figure 1: Wastewater plays an important role in antibiotic resistance development. Wastewater *E.*
 63 *coli* isolates are tested for antibiotic resistance and sequenced. Many isolates are multi-drug
 64 resistant and potentially pathogenic. Their large pan-genome is a source of potentially novel
 65 resistance genes.

66

67

68 **Results**

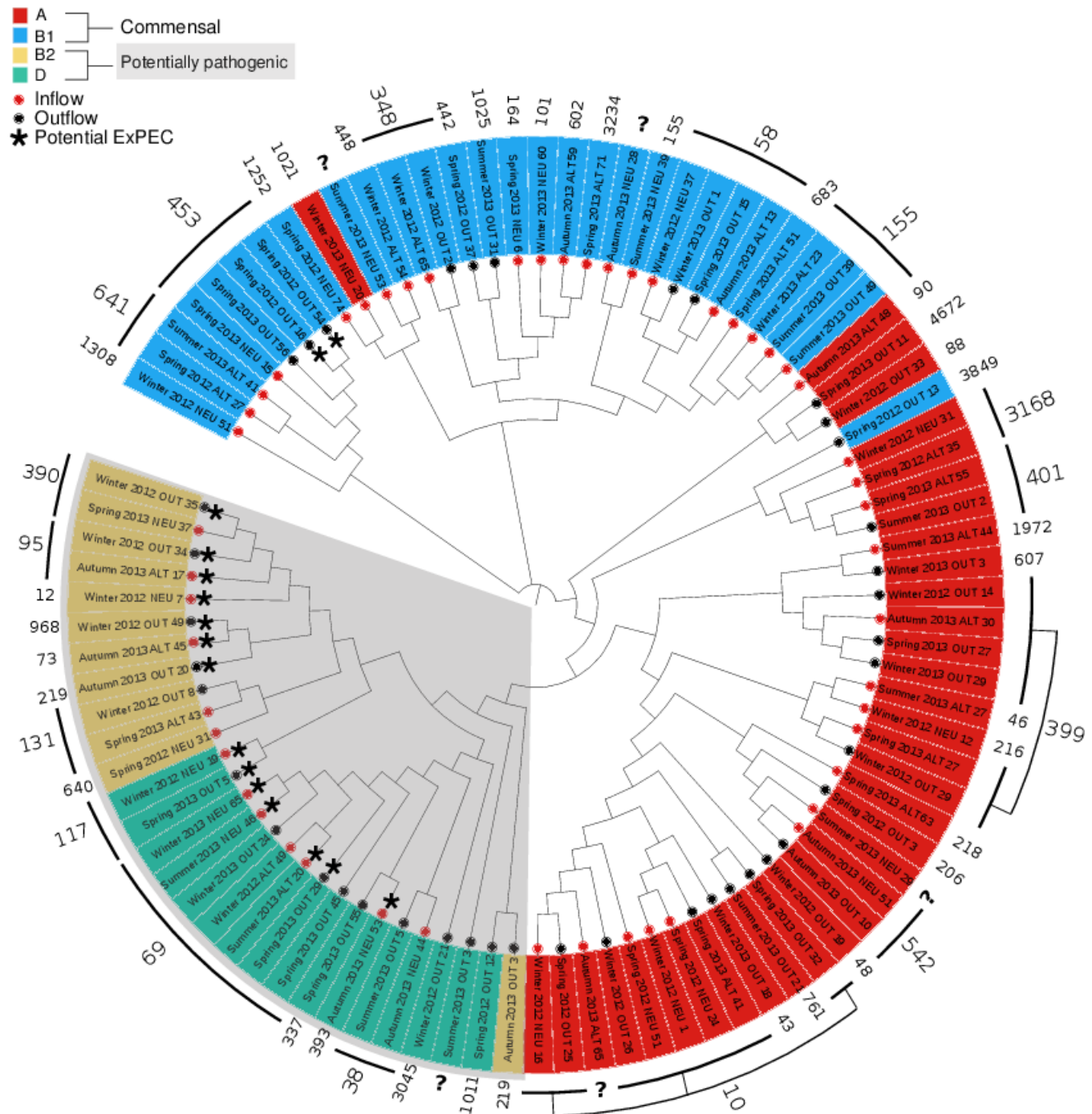
69 **Pathogenic potential.** *Escherichia coli* strains exhibit great variation. Many exist as harmless
70 commensals in the human gut, but some are intra- (InPEC) or extra-intestinal pathogenic
71 *Escherichia coli* (ExPEC). To assess the pathogenic potential without in vivo testing, clinical
72 research has developed databases of virulence factors and genotyping schemes. The sequenced
73 isolates contain some 700 of the 2000 *Escherichia coli* virulence factors in the virulence factor
74 database (8), averaging to 153 and to 155 virulence factors per isolate for inflow and outflow,
75 respectively. Hence, there is no significant difference (Welch test, CI 95%) between inflow and
76 outflow. In particular, we found combinations of virulence factors for 16 isolates (see methods),
77 which are indicative of ExPEC. Eight of these 16 isolates were obtained from the outflow of the
78 treatment plant (see Fig. 2).

79

80 Besides the presence of known virulence factors, the pathogenic potential can be assessed using
81 genotyping with multi-locus sequence types (9) and phylogroups (10). Broadly, *Escherichia coli*
82 has, among others, four phylogroups, A, B1, B2 and D. Commensal *Escherichia coli* fall mostly
83 into groups A and B1 and ExPEC into B2 and D (10). Fig. 2 shows a phylogenetic tree of the
84 sequenced wastewater *Escherichia coli* isolates along with the commensal phylogroups A (red)
85 and B1 (blue) and the pathogenicity-associated groups B2 (yellow) and D (green), as well as the
86 finer-grained multi-locus sequence types. The tree is based on genomic variations compared to
87 the reference genome of *Escherichia coli* K12 MG1655. Fig. 2 reveals that nearly one third of
88 isolates belong to group B2 and D, in which ExPEC are usually found. In particular, B2 and D
89 include 14 of the 16 potential ExPEC isolates. Remarkably, half of the B2 and D isolates are from
90 the wastewater treatment plant's outflow.

91

92



93

94 Figure 2: Phylogeny and pathogenic potential of wastewater *Escherichia coli*. Phylogenetic tree,
 95 multi-locus sequence types, and phylogroups of 92 sequenced wastewater *Escherichia coli*
 96 isolates reveal 16 potential ExPEC isolates (marked with a black star) in phylogroups B2 (yellow)
 97 and D (green), which are associated with pathogenicity. Half of the potentially pathogenic isolates
 98 stem from the outflow of the treatment plant.

99

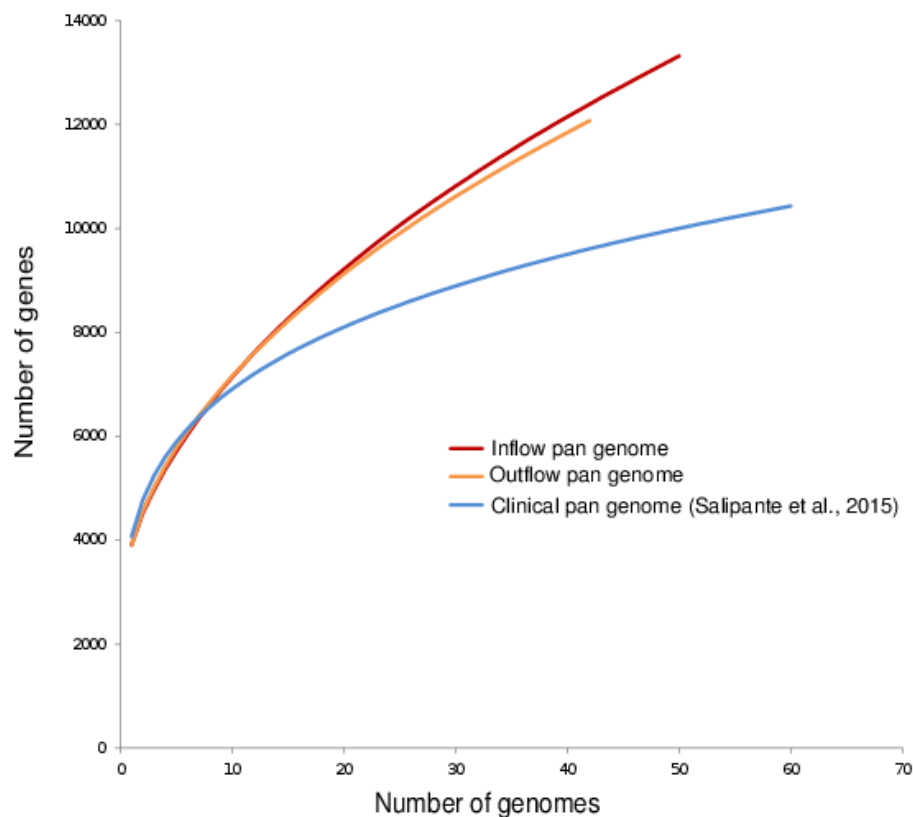
100
101 **The wastewater pan-genome.** The concept of evolution implies that genomes of organisms of
102 the same species differ. Differences range from small single nucleotide polymorphisms to large
103 genome rearrangements. As a consequence, *Escherichia coli* possesses a core of genes present
104 in all genomes, as well as genes only present in some genomes, or even just in one. The union of
105 all of these genes is called the pan-genome. It is believed, that the *Escherichia coli* core genome
106 comprises around 1400-1500 genes, while the pan-genome may be of infinite size (11).

107
108 To assess the degree of genomic flexibility of the wastewater isolates, we relate the wastewater
109 pan-genome and the wastewater core genome. At 16582 genes, the wastewater pan-genome is
110 nearly six times larger than the wastewater core genome of 2783 genes, a reservoir of some
111 14000 genes. Despite this large reservoir, the size difference of nearly 1000 genes between the
112 wastewater *Escherichia coli* core genome and the whole species core genome suggests that the
113 full diversity of *Escherichia coli* is still not covered in our wastewater sample.

114
115 The balance between maintaining the core genome and spending energy on acquisition of new
116 genetic material can be captured by the ratio of the core genome size and the average genome
117 size, which is 4700 genes in our sample. This means that only $1400/4700 = 30\%$ of genes in our
118 wastewater *Escherichia coli* are core genes. Most of the non-core genes are very unique and
119 appear only in one or two isolates each. More precisely, 50% of the pan-genome genes appear in
120 only one or two isolates each. This implies that the wastewater *Escherichia coli* studied are highly
121 individual.

122
123 But do *Escherichia coli* maintain such a rich genome after wastewater treatment? Fig. 3 shows
124 that they do. The 42 *Escherichia coli* genomes of the plant's outflow comprise nearly 12000
125 genes and the pan-genome growth curves between in- and outflow are nearly identical, which
126 means that wastewater treatment does not affect the genetic diversity of *Escherichia coli*. Fig. 3
127 also shows a clinical dataset of ExPEC and these clinical *Escherichia coli* are more
128 homogeneous and hence their pan-genome is smaller. In contrast, the diversity of the wastewater
129 *Escherichia coli* match other datasets comprising mixtures of commensal and pathogenic

130 *Escherichia coli*, as well as *Shigella* genomes (see Table 1). This underlines the great diversity of
 131 genomes before and after wastewater treatment and leads to the question whether these diverse
 132 genomes harbour antibiotic resistance genes?



133
 134 Figure 3: The pan-genome at the outflow has the same size as at the inflow, suggesting that
 135 highly flexible *Escherichia coli* emerge from a treatment plant. The wastewater pan-genome is
 136 larger than a clinical pan-genome one and of similar size to (see Table 1) highly diverse samples
 137 comprising pathogenic, commensal, and lab *Escherichia coli*, as well as *Shigella*.
 138

Ref	Pan	Core	Strains	Path.	Comm.	Lab	Shig.
This study	16582	2783	92	28	62	0	0
Kaas et al., 2012 ¹⁵	16373	1702	186	171			15
Vieira et al., 2011 ¹³	14986	1957	29	21	8	0	6
Gordienko et al., 2013 ¹²	12000	2000	32	16	6	3	7
Lukjancenko et al., 2010 ¹⁶	13000	1472	53	35	11	7	0
Rasko et al., 2008 ¹⁷	13000	2344	17	14	1	2	0
Touchon et al., 2009 ¹⁴	11432	1976	20	10	3	0	7

139
 140 Table 1: Highly diverse samples comprising pathogenic, commensal, and lab *Escherichia coli*, as
 141 well as *Shigella*.

142
 143

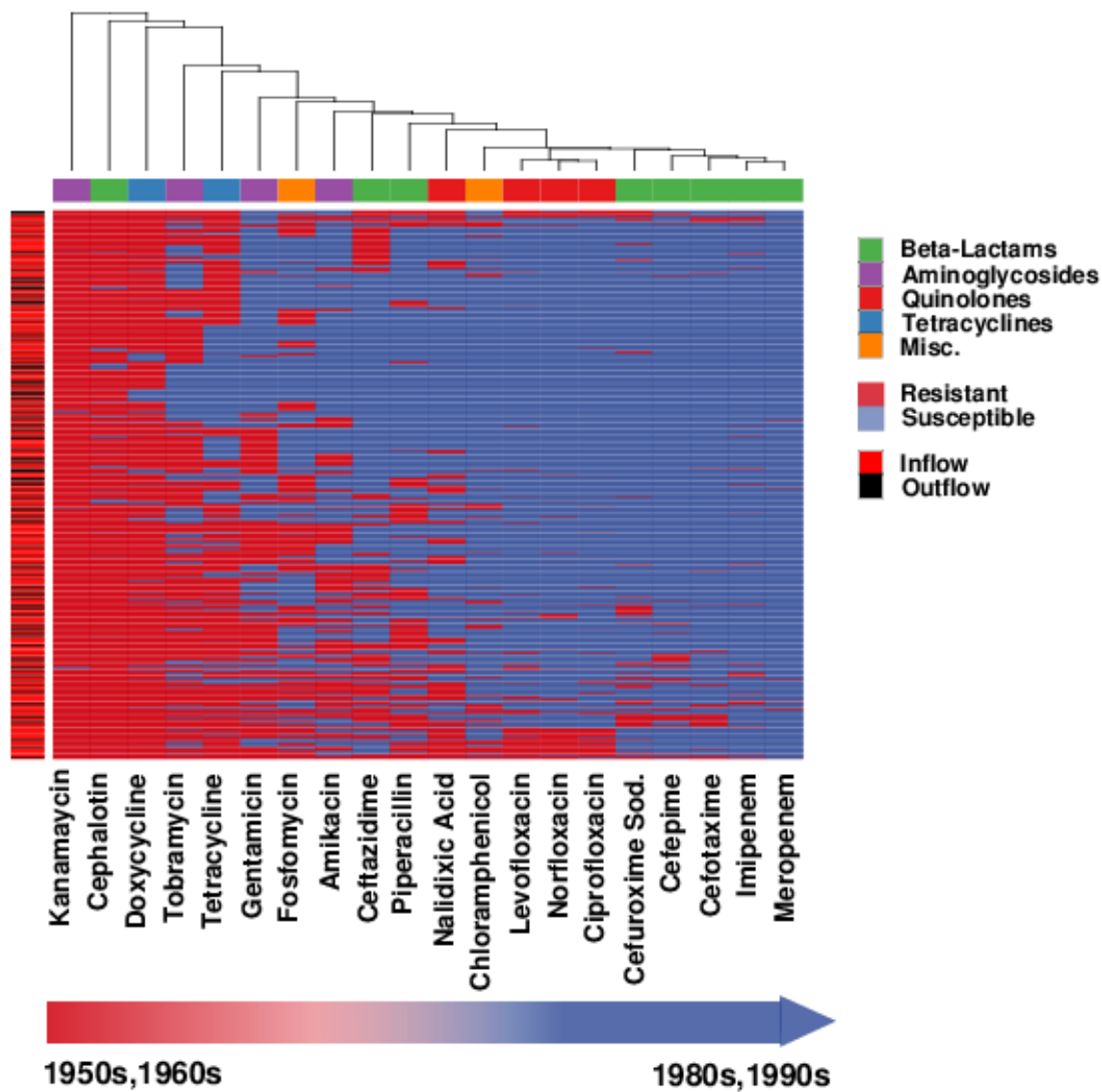
144

145

146 **Resistance genes in the wastewater pan-genome.** Wastewater *Escherichia coli* are known to
147 host antibiotic resistance genes. While there are many known resistance genes (see e.g. CARD
148 (12)), they fall mostly into a few groups, such as beta-lactamases. Here, we seek to confirm and
149 expand the space for candidate resistance genes. Firstly, we measured antibiotic resistance in all
150 1178 isolates to the 20 antibiotics. Fig. 4 reveals a high degree of resistance and huge
151 differences between different antibiotics, including a general trend indicating greater resistance to
152 antibiotics that have been available for longer. Concretely, antibiotics from the 50s and 60s have
153 a significantly different number of resistances than the more recent antibiotics (Welch test, p-
154 value < 0.0025). However, there is no significant difference in the number of resistances between
155 isolates from the inflow and the outflow (p-value 0.0001), suggesting that wastewater treatment is
156 not affecting resistance at all.

157

158



159

160 Figure 4: 1178 Wastewater *Escherichia coli* isolates are tested for antibiotic resistance to 20
 161 antibiotics. Nearly all isolates are multi-drug resistant. Generally, isolates are more susceptible to
 162 betalactams and fluoroquinolones than to tetracyclins and aminoglycosides. Surprisingly, the
 163 outflow isolates show similar resistance as inflow (p-value 0.0001), suggesting that wastewater
 164 treatment is not reducing resistance development.

165

166

167
168 Next, we correlated the presence of each gene in the sequenced isolates with their phenotypic
169 antibiotic resistance profiles. We excluded meropenem and imipenem, since nearly all isolates
170 are susceptible. For each of the 18 remaining antibiotics, we list the top ten candidate resistance
171 genes in Table 2. These 180 genes comprise 88 unique confirmed genes, including many well-
172 known resistance genes, such as efflux pumps (MT1297 and *emrE*), membrane and transport
173 proteins (*aida-I*, *yiaV*, *yijK*, *pitA*, *icsA*, and *pagN*), tetracycline (*tetA*, *tetR*, and *tetC*),
174 chloramphenicol (*cat*), and piperacillin (the beta lactamase *bla2*) resistance genes. However, the
175 180 genes also comprise a large number of open reading frames encoding hypothetical proteins
176 (41) and genes not yet linked to antibiotic resistance (116). These genes have to be studied
177 further to determine whether they are novel resistance genes or just correlating (e.g. because
178 they are on the same genetic element with a resistance gene). All but three of the known and the
179 potentially new resistance genes are present in some isolates obtained from the treatment plant's
180 outflow.

181
182

	Amikacin	Gentamicin	Kanamycin	Tobramycin	Doxycycline	Tetracycline	Cefepime	Cefotaxime	Ceftazidime	Cefuroxime Sod.	Cephalotin	Piperacillin	Ciprofloxacin	Levofloxacin	Nalidixic Acid	Norfloxacin	Chloramphenicol	Fosfomicin
1	Hypothetical Protein	4-hydroxyacetophenone monooxygenase hapE	Transposase IS200 like protein	Autotransporter precursor aida-I	Tetracycline resistance protein, class B tetA	Oxygen-dependent choline dehydrogenase betA	Ash protein family protein	Hypothetical Protein	cell division protein	Type-1 restriction enzyme R protein hsdR	GTPase era	Beta-lactamase TEM precursor bla	Virulence regulon transcriptional activator virB	Transposon Tn10 protein tetD	Mercuric resistance operon regulatory protein merR	Transposon Tn10 protein tetD	Chloramphenicol acetyltransferase cat	Invasin 184
2	Caudovirales tail fiber assembly protein	Phosphoadenosine phosphosulfate reductases	putative multidrug-efflux transporter/M T1297	putative protease yhbU precursor	Tetracycline repressor protein class B tetR	NAD/NADP-dependent betaine aldehyde dehydrogenase betB	Fibronectin type III protein	Hypothetical Protein	Plasmid stability protein	Type I restriction enzyme EcoKI M protein hsdM	Prophage CP4-57 regulatory protein alpA	Transposon Tn3 resolvase tnpR	Sporulation initiation inhibitor protein Soj	Tetracycline resistance protein, class B tetA_1	Mercuric resistance protein merC	Tetracycline resistance protein, class B tetA_1	Streptomycin 3'-adenylyltransferase ant1	Putative DNA-invertase Rac pinR
3	Swarming motility protein ybiA	putative multidrug-efflux transporter/M T1297	Phosphotransferase enzyme family protein	Chaperone protein dnaK	Transposon Tn10 TetC protein tetC	HTH-type transcriptional regulator betI	Transcriptional activator perC	Transcriptional activator perC	HTH-type transcriptional regulator cmtR	mrr restriction system protein	Hypothetical Protein	Tyrosine recombinase xerD	putative HTH-type transcriptional regulator	Tetracycline repressor protein class B from transposon Tn10 tetR	mercuric transport protein merT	Tetracycline repressor protein class B from transposon Tn10 tetR	Chromosome-partitioning ATPase soj	Transcriptional repressor dicA
4	Phospholipase ytpA	Phosphotransferase enzyme family protein	Hypothetical Protein	putative ABC transporter ATP-binding protein yjjK	HTH-type transcriptional regulator cmtR	Tetracycline resistance protein, class B tetA	Hypothetical Protein	Hypothetical Protein	Phage-related minor tail protein	Outer membrane protein lcsA precursor	Hypothetical Protein	Acetyltransferase (GNAT) family protein	DNA-binding transcriptional regulator dicC	Transposon Tn10 protein tetC	Mercuric transport protein periplasmic component precursor merP	Transposon Tn10 protein tetC	parG	Hypothetical Protein
5	Carbonic anhydrase 1 cynT	Hypothetical Protein	Streptomycin 3'-adenylyltransferase ant1	cell envelope integrity inner membrane protein tolA	Tetracycline resistance protein, class C tetA	Tetracycline repressor protein class B tetR	Hypothetical Protein	Hypothetical Protein	Phage tail protein E	Hypothetical Protein	Hypothetical Protein	Virulence regulon transcriptional activator virB	Hypothetical protein	putative HTH-type transcriptional regulator	Anti-adaptor protein iraM	CAAX amino terminal protease self-immunity	Hypothetical Protein	Hypothetical Protein
6	Hypothetical Protein	Hypothetical Protein	Hypothetical Protein	Inner membrane protein yiaV precursor	putative inner membrane transporter yedA	Transposon Tn10 TetC protein tetC	Chromosome partition protein smc	Hypothetical Protein	Hypothetical Protein	Fibronectin type III protein	Transposon Tn10 tetD protein	Transposase	Lysine-tRNA ligase lysS	DNA-binding transcriptional regulator dicC	Hypothetical protein	mRNA interferase pemK	Hypothetical Protein	Hypothetical Protein
7	Xanthine dehydrogenase molybdenum-binding subunit xdhA	Hypothetical Protein	Zinc-responsive transcriptional regulator	Entericidin B membrane lipoprotein	Tetracycline repressor protein class A from transposon 1721 tetR	High-affinity choline transport protein betT	Hypothetical Protein	Invasin	Hypothetical Protein	Hypothetical Protein	putative multidrug-efflux transporter/M T1297	Tetracycline resistance protein, class B tetA	Transposon Tn10 protein tetD	Hypothetical protein	Mercuric reductase merA_1	Antitoxin pemI	Acetyltransferase (GNAT) family protein	Molybdenum cofactor biosynthesis protein A
8	Nicotinate dehydrogenase FAD-subunit ndhF	Hypothetical Protein	merE protein	Low-affinity inorganic phosphate transporter 1 pitA	Hypothetical Protein	Formate dehydrogenase H fdhF	Aldehyde-alcohol dehydrogenase adhE	Hypothetical Protein	Tyrosine recombinase xerC	Hypothetical Protein	Phosphotransferase enzyme family protein	Tetracycline repressor protein class B tetR	Tetracycline resistance protein, class B tetA_1	CAAX amino terminal protease self-immunity	Hypothetical protein	putative HTH-type transcriptional regulator	putative multidrug-efflux transporter/M T1297	ATP-dependent zinc metalloprotease ftsH4
9	Nicotinate dehydrogenase small FeS subunit ndhS	Phage polarity suppression protein psu	Phosphoadenosine phosphosulfate reductases	Methyl-accepting chemotaxis protein II tar	Transposon Tn10 tetD protein	S-fimbrial protein subunit sfah	Aldehyde-alcohol dehydrogenase adhE	Hypothetical Protein	Hypothetical Protein	Hypothetical Protein	Outer membrane protein pagN precursor	Transposon Tn10 tetC protein	Tetracycline repressor protein class B transposon Tn10 tetR	mRNA interferase pemK	zinc-responsive transcriptional regulator	DNA-binding transcriptional regulator dicC	Phosphotransferase enzyme family protein	Molybdenum cofactor biosynthesis protein A
10	putative fimbrial-like protein EIfG precursor elfG	DNA primase traC	Caudovirales tail fiber assembly protein	Leucine-specific-binding protein precursor livK	putative multidrug-efflux transporter/M T1297	Beta-lactamase TEM precursor bla	Cob(II)yrinic acid a,c-diamide adenosyltransferase yvqK	Type-1 restriction enzyme R protein hsdR	Hypothetical Protein	Hypothetical Protein	Tetracycline resistance protein, class B tetA	Multidrug transporter emrE	Transposon Tn10 protein TetC	Antitoxin PemI	MerE protein	Caudovirales tail fiber assembly protein	Leucine-specific-binding protein precursor livK	Hypothetical Protein

Table 2: Known and candidate resistance genes from correlation of genomes to resistance phenotype. Top 10 genes for 18 antibiotics.

185 **Discussion**

186 **The cost of resistance.** There is a debate on whether the evolution of resistance is a
187 competitive disadvantage. Some evidence indicates that resistant bacteria may be outcompeted
188 by susceptible bacteria (13) and that they may be collaterally sensitive (14), i.e. resistant bacteria
189 may become susceptible through an appropriate co-treatment. In contrast, it appears that our
190 isolates do not suffer from evolutionary disadvantages despite their resistance and the harsh
191 environment of wastewater treatment. The latter includes the reduction of the bacterial population
192 as indicated by a reduction in biochemical oxygen demand, which is part of secondary
193 wastewater treatment. Despite this bacterial reduction, wastewater outflow has similar resistance
194 levels as the inflow. These findings support evidence (15) that bacteria can compensate for the
195 cost of resistance.

196 **Pathogenic potential and resistance.** Ultimate proof for pathogenicity can only be obtained
197 from in vivo studies. However, the pathogenic potential can be assessed from an analysis of a
198 genome for virulence markers. Here we chose to consider three independent approaches:
199 classification by phylogenetic groups, by multi-locus sequence tags, and by identification of
200 specific virulence factors (see methods). While the three approaches showed consistent results,
201 they are by no means proof for pathogenicity, since there can be exceptions to these
202 classification rules. As an example, consider the strain ed1a, which belongs to the phylogenetic
203 group B2, but it is not pathogenic in mice (16). Similarly, pathogenicity may not only arise from the
204 acquisition of genes, but also from the loss (17). Regarding resistance there are similar
205 confounding factors. *Escherichia coli* is inherently resistant to kanamycin and cephalotin, which is
206 also clearly shown in Fig. 4. More generally, antibiotic resistance is ancient (18) and naturally
207 occurring in the environment. Nonetheless, there are pronounced differences between pristine
208 and human environments (19). This is also supported by Fig. 4, which shows that antibiotics
209 introduced in the 50s and 60s have more resistances than those introduced later (p-value <
210 0.0025), which suggests, that the naturally occurring resistances do not play a major role in the
211 emergence of observed resistances.

212

213 **From clinic to river.** We have shown that there are *Escherichia coli* at the wastewater outflow,
214 which are multi-drug resistant and have pathogenic potential. But are they abundant enough to
215 have an impact in the aquatic system they are released into? They do. The percentage of
216 possibly pathogenic *Escherichia coli* in the outflow is considerable and may correspond to a large
217 absolute amount. If an average of 100 *Escherichia coli* colony forming units (CFU) are released
218 per ml, then 10^{13} CFUs per day are released (assuming a release of 10^5 m³ per day). This is in
219 accordance with Manaia *et al.*, who showed that 10^{10} - 10^{14} CFU of ciprofloxacin-resistant bacteria
220 are released by a mid-sized wastewater treatment plant (20). Furthermore, a study in a Japanese
221 river shows the presence of pathogenic *Escherichia coli*. Gomi *et al.* (21) sequenced over 500
222 samples from the Yamato river and most of their prevalent multi-drug resistant and clinical strains
223 are also present in our samples. In a related study, Czekalski *et al.* found that particle-associated
224 wastewater bacteria are the responsible source for antibiotic resistance genes in the sediments of
225 lake Geneva in Switzerland (22). Assuming that the river Elbe is comparable to these aquatic
226 systems, it suggests, that clinic and river are connected with wastewater treatment plants in
227 between.

228 **Composition of phylogroups.** It is interesting to compare the breakdown into phylogenetic
229 groups of wastewater *Escherichia coli* to compare samples from human and animal
230 environments. It is, e.g., known that the phylogenetic group B2 is more abundant among
231 commensal *Escherichia coli* from human faeces (43%) than from farm animals (11%) (23).
232 Therefore, the composition of wastewater *Escherichia coli* as shown in Fig. 2 resembles
233 commensal *Escherichia coli* from farm animals more closely. Similarly, Tenailon find that groups
234 A and B1 make up one third in human faeces (23), whereas we find two thirds. This suggests that
235 animal waste plays an important role for resistance of waste water bacteria.

236 **Pan and core genome.** We identified known and candidate resistance genes by correlating their
237 presence in the genome against their resistance phenotype across isolates. By virtue of this
238 correlation, the identified genes are not to be found in the core genome, but in the pan genome.
239 As many authors have pointed out, *Escherichia coli* has a large and flexible pan genome.
240 Lapierre *et al.* argue that *Escherichia coli* appears to have unlimited ability to absorb genetic
241 material and hence its pan genome is open (11). In a recent study comprising over 2000

242 genomes Land *et al.* put this into numbers and arrive at a pan genome of 60000-89000 gene
243 families for over 2000 sequenced *E. coli* genomes (24). This upper limit shows that the
244 wastewater pan genome of 16582 genes is still not the top. Nonetheless, it is considerably larger
245 than a clinical pan genome. These differences indicate the heterogeneity of genomes. Clinical
246 *Escherichia coli* genomes are not as diverse as the ones in a wastewater pool, which comprises
247 besides human faeces also animal waste.

248 **Random sampling and hypothesis-free analysis.** The initial 1178 isolates were sampled
249 randomly over different times of the year, from two different inflows and the outflow. In contrast,
250 the 103 sequenced isolates were chosen in such way that all of the phenotypes encountered
251 were represented (see methods). Within a phenotype group isolates were chosen randomly. This
252 random, but representative choice and the subsequent link from genotype to phenotype is an
253 example of high-throughput hypothesis-free analysis. And although, there was no pre-defined
254 resistance mechanism, which we aimed to hit, many of the well-known resistance genes were
255 ranked high. This supports the hope that high-throughput, hypothesis-free methods such as deep
256 sequencing will help to uncover novel resistance mechanisms and in particular that some of the
257 candidate resistance genes will prove to have a causal link to resistance.

258

259 **Conclusion**

260 Overall, we have shown for the first time that *Escherichia coli* isolates from a wastewater outflow
261 have pathogenic potential and large pan-genomes, which harbor known and novel candidate
262 resistance genes. Together with the estimates on absolute *Escherichia coli* abundance, this
263 means that despite treatment, there is a considerable pathogenic potential at the outflow of a
264 wastewater treatment plant. These results underline the need to include wastewater treatment
265 plants in the combat against antibiotic resistance.

266

267

268 **Methods**

269 **Collection.** 1178 samples were collected from the municipal wastewater treatment plant
270 Dresden, Germany. Samples were collected on 11/4/2012 (Spring 2012), 30/7/2012 (Summer
271 2012), 21/1/2013 (Winter 2012), 27/3/2013 (Spring 2013), 6/8/2013 (Summer 2013), 14/10/2013
272 (Autumn 2013), and 17/12/2013 (Winter 2013). Samples were collected either at the outflow
273 (OUT) or at one of two inflow locations (Altstadt ALT and Neustadt NEU), representing the area
274 south and north of the river Elbe).

275 **Isolation.** *Escherichia coli* and total coliforms bacteria were enumerated via serial fold dilution
276 plating of the original wastewater (triplicate samples). Wastewaters were diluted in double distilled
277 water, until the enumeration of bacterial colonies was possible. *Escherichia coli* and coliform
278 counts were always performed in triplicates. The *Escherichia coli* colonies were selected and
279 picked after overnight growth at 37°C on a selective chromogenic media (OXOID Brilliance
280 *Escherichia coli*/Coliform Selective Agar, Basingstoke, England). To minimize the risk of colony
281 contamination, picked colonies were spiked a second time on the same selective media and pure
282 single colonies were grown overnight on LB media at 37°C and stored on glycerol stock at -80° C.

283 **Resistance phenotyping.** Antibiotic resistance phenotypes were determined by the agar
284 diffusion method using 20 antibiotic discs (OXOID, England) according to EUCAST (or CLSI
285 when EUCAST was not available) (13, 18). The selected drugs belong to the most commonly
286 prescribed antibiotics for diseases caused by bacteria according to the German health insurance
287 AOK Plus: piperacillin (100µg), nalidixic acid (30µg), chloramphenicol (30µg), imipenem (10µg),
288 cefotaxime (30µg), cephalotin (30µg), kanamycin (30µg), tetracycline (30µg), gentamicin (10µg),
289 amikacin (30µg), ciprofloxacin (5µg), fosfomycin (50µg), doxycycline (30µg), cefepime (30µg),
290 ceftazidime (10µg), levofloxacin (5µg), meropenem (10µg), norfloxacin (10µg), cefuroxime sod.
291 (30µg), tobramycin (10µg) (25). After 24 hours of incubation at 37°C, the resistance diameters
292 were measured. Clustering of antibiotics and of isolates was performed using the R function
293 heatmap.2 from the R library (26) Heatplus and hierarchical clustering of matrices based on
294 Euclidean distances between isolates and between antibiotics.

295

296 **Sequencing.** To select isolates representative of phenotype, we clustered isolates according to
297 the diameters of inhibition zone against the 20 antibiotics using k-means clustering based on
298 Euclidean distances between isolates (vectors of 20 inhibition zone diameters). The analysis and
299 graphs were produced using R version 3.2.4 (26). As clusters may be highly skewed in number of
300 cluster members, we tested all cluster numbers from 1 to 100 and plotted within class sum of
301 squares against k . At $k = 47$, the sum of squares tails off and there is a steep local decrease, so
302 that $k = 47$ was fixed as k-means parameter. We obtained 103 isolates, which were subsequently
303 used for sequencing and further analysis. To further validate the choice, we plotted the average
304 number of resistances against number of isolates and antibiotics vs. number of isolates for the
305 total 1178 and the selected 103 isolates (see Supp Fig. 1) and concluded that both distributions
306 are roughly similar. 3000ng DNA were extracted from each of the 103 selected isolates using
307 MasterPure extraction kit (Epicentre) according to the manufacturer's instructions. Sequencing
308 was performed using Illumina Flex GL.

309
310 **Assembly.** Genomes were assembled with Abyss (version 1.5.2) (27). In order to optimize k for
311 the best assembly, k-mer values had to be empirically selected from the range of 20-48 (see
312 Supp. Fig. 2) on a per sample basis to maximize contiguity (3). To determine the k-mer length
313 that achieved highest contiguity, the 28 assemblies per draft genome/isolate were compared
314 based on $N50$ values. 11 assemblies with an $N50$ statistic of less than 5×10^4 bp were excluded
315 (28).

316

317

318 **Genes.** Reference gene clusters were computed from 58 complete *Escherichia coli* genomes
319 (see Table 2) available in June 2015 from NCBI. Genes were identified in wastewater and
320 reference genomes using Prokka (version 1.11) (29). Genes were clustered at 80% using CD-
321 HIT (30) (version 4.6.3, arguments -n 4 -c 0.8 -G 1 -aL 0.8 -aS 0.8 -B 1). Genes with over 90%
322 sequence identity, but only 30% coverage, as well as genes with 80% or greater identity and
323 covered to phage and virus sequences (31) were discarded. A gene cluster is defined to be
324 present in an isolate if there is a Prokka gene in the genome, which is longer than 100 amino
325 acids and has over 80% sequence identity and coverage against the gene cluster representative.

326

327 **Pan- and core-genome.** To generate the pan- and core-genome size graph we followed the
328 procedure in (3, 16). We had 92 genomes available. We varied i from one to 92. At each subset
329 size i , we randomly selected i genomes and computed the sizes of the union (pan) and
330 intersection (core) of gene clusters. This random selection was carried out 2000 times in each
331 step.

332

333 **Gene clusters to rank genes by correlation to phenotype.** Prokka genes were identified in all
334 isolate genomes and then clustered with CD-HIT at 60% sequence identity and 50% coverage
335 (arguments -n 4 -c 0.6 -G 1 -aL 0.8 -aS 0.5 -B 1). A 80% identity cutoff was also tried but
336 dismissed, because the 60% threshold yielded 25% less clusters while adequately clustering
337 homologous gene sequences with lower sequence similarity. This threshold value is also
338 supported by the widespread default use of the BLOSUM62 matrix, the basis of which is
339 sequences clustered by 62% sequence identity.

340

341 **Tree.** The phylogenetic tree of 92 isolates was built following the procedure of (32, 33) using
342 FastTree version 2.1 (34). Sequence reads were aligned to *Escherichia coli* K12 MG 1665 and
343 single nucleotide variant calling was carried out using GATK (35). Quality control for variant
344 calling was performed; variants supported by more than ten reads or likelihood score greater than
345 200 were always in the range of 84 – 99% of variants called per isolate with the exception of 2

346 isolates where only 59% and 60% of the variants were above the threshold for quality and
347 supporting reads. FastTree 2.1 (34) was then used to build the maximum likelihood tree based on
348 the sequences derived from variant calling.

349 **Phylogrouping.** For phylogrouping, the classification system established by Clermont *et al.* (10)
350 based on the genes *chuA* and *yjaA* and the DNA fragment TspE4.C2 was used. Blast was
351 performed to check each genome assembly for presence or absence of the aforementioned
352 elements with an identity cutoff $\geq 90\%$.

353
354 **MLST.** Concerning epidemiology and Multi-Locus Sequence Typing, we used the webserver at
355 <https://cge.cbs.dtu.dk/services/MLST/> that follows the MLST scheme in (36) for predicting MLSTs
356 from whole genome sequence data (37). 92 Draft genome assemblies were submitted and
357 results were obtained; 5 isolates were unidentified demonstrating novel sequence types.

358
359 **Virulence factors.** Virulence factors protein sequences were downloaded from VFDB: Virulence
360 Factors database (8, 38). 2000 sequences which were *Escherichia coli* related were chosen.
361 Sequences were then clustered at 80% sequence identity using CD-HIT (version 4.6.3,
362 arguments -n 4 -c 0.8 -G 1 -aL 0.8 -aS 0.8 -B 1). A virulence factor was considered present in an
363 isolate's genome if there is a Prokka gene in the genome that has over 80% sequence identity
364 and coverage against the virulence factor cluster representative.

365

366

367 **ExPEC classification.** There are intra- and extra-intestinal pathogenic *E.coli*, which can be
 368 classified from the presence of virulence factors (39-42). InPEC are characterised by the
 369 virulence factors stx1, stx2, escV, and bfpB. They are ExPEC if they contain over 20 of 58
 370 virulence factors afa/draBC, bmaE, gafD, iha cds, mat, papEF, papGII, III, sfa/foc, etsB, etsC, sitD
 371 ep, sitD ch, cvaC MPIII, colV MPIX, eitA, eitC, iss, neuC, kpsMTII, ompA, ompT, traT, hlyF, GimB,
 372 malX, puvA, yqi, stx1, stx2, escV, bfp, feob, aatA, csgA, fimC, focG, nfaE, papAH, papC, sfaS,
 373 tsh, chuA, fyuA, ireA, iroN, irp2, iucD, iutA, sitA, astA, cnf1, sat, vat, hlyA, hlyC, ibeA, tia, and pic.

374

375 **Data availability statement**

376 Genome assemblies of the analyzed isolates that support the findings of the study will be made
 377 available on the NCBI upon paper publication (see Table 3).

378 **Accession numbers of novel whole genome assemblies of analyzed isolates.**

Bioproject	Biosample	Accession	strain
PRJNA380388	SAMN06641941	NBBP00000000	Escherichia coli Win2013_WWKa_OUT_3
PRJNA380388	SAMN06641940	NBBQ00000000	Escherichia coli Win2013_WWKa_OUT_29
PRJNA380388	SAMN06641933	NBBR00000000	Escherichia coli Win2013_WWKa_OUT_18
PRJNA380388	SAMN06641932	NBBS00000000	Escherichia coli Win2013_WWKa_OUT_24
PRJNA380388	SAMN06641931	NBBT00000000	Escherichia coli Win2013_WWKa_OUT_1
PRJNA380388	SAMN06641928	NBBU00000000	Escherichia coli Win2013_WWKa_NEU_65
PRJNA380388	SAMN06641927	NBBV00000000	Escherichia coli Win2013_WWKa_NEU_20
PRJNA380388	SAMN06641926	NBBW00000000	Escherichia coli Win2013_WWKa_NEU_60
PRJNA380388	SAMN06641901	NBBX00000000	Escherichia coli Win2013_WWKa_ALT_23
PRJNA380388	SAMN06641884	NBBY00000000	Escherichia coli Win2012_WWKa_OUT_49
PRJNA380388	SAMN06641883	NBBZ00000000	Escherichia coli Win2012_WWKa_OUT_8
PRJNA380388	SAMN06641882	NBCA00000000	Escherichia coli Win2012_WWKa_OUT_34
PRJNA380388	SAMN06641881	NBCB00000000	Escherichia coli Win2012_WWKa_OUT_35
PRJNA380388	SAMN06641880	NBCC00000000	Escherichia coli Win2012_WWKa_OUT_29
PRJNA380388	SAMN06641879	NBCD00000000	Escherichia coli Win2012_WWKa_OUT_26
PRJNA380388	SAMN06641878	NBCE00000000	Escherichia coli Win2012_WWKa_OUT_33
PRJNA380388	SAMN06641877	NBCF00000000	Escherichia coli Win2012_WWKa_OUT_21
PRJNA380388	SAMN06641876	NBCG00000000	Escherichia coli Win2012_WWKa_OUT_2
PRJNA380388	SAMN06641875	NBCH00000000	Escherichia coli Win2012_WWKa_NEU_7
PRJNA380388	SAMN06641874	NBCI00000000	Escherichia coli Win2012_WWKa_OUT_14
PRJNA380388	SAMN06641873	NBCJ00000000	Escherichia coli Win2012_WWKa_NEU_51
PRJNA380388	SAMN06641872	NBCK00000000	Escherichia coli Win2012_WWKa_NEU_31
PRJNA380388	SAMN06641871	NBCQ00000000	Escherichia coli Win2012_WWKa_NEU_37
PRJNA380388	SAMN06641870	NBCR00000000	Escherichia coli Win2012_WWKa_NEU_16
PRJNA380388	SAMN06641869	NBCS00000000	Escherichia coli Win2012_WWKa_NEU_19
PRJNA380388	SAMN06641868	NBCT00000000	Escherichia coli Win2012_WWKa_NEU_12
PRJNA380388	SAMN06641867	NBCU00000000	Escherichia coli Win2012_WWKa_ALT_65
PRJNA380388	SAMN06641866	NBCV00000000	Escherichia coli Win2012_WWKa_NEU_1
PRJNA380388	SAMN06641865	NBCW00000000	Escherichia coli Win2012_WWKa_ALT_49
PRJNA380388	SAMN06641864	NBCX00000000	Escherichia coli Win2012_WWKa_ALT_54
PRJNA380388	SAMN06641863	NBCY00000000	Escherichia coli Sum2013_WWKa_OUT_5
PRJNA380388	SAMN06641862	NBCZ00000000	Escherichia coli Sum2013_WWKa_OUT_39
PRJNA380388	SAMN06641861	NBDA00000000	Escherichia coli Sum2013_WWKa_OUT_49
PRJNA380388	SAMN06641860	NBDB00000000	Escherichia coli Sum2013_WWKa_OUT_3
PRJNA380388	SAMN06641859	NBDC00000000	Escherichia coli Sum2013_WWKa_OUT_31
PRJNA380388	SAMN06641858	NBDD00000000	Escherichia coli Sum2013_WWKa_OUT_2
PRJNA380388	SAMN06641857	NBDE00000000	Escherichia coli Sum2013_WWKa_OUT_21
PRJNA380388	SAMN06641856	NBDF00000000	Escherichia coli Sum2013_WWKa_NEU_53
PRJNA380388	SAMN06641855	NBDG00000000	Escherichia coli Sum2013_WWKa_NEU_46
PRJNA380388	SAMN06641854	NBDH00000000	Escherichia coli Sum2013_WWKa_NEU_39
PRJNA380388	SAMN06641853	NBDI00000000	Escherichia coli Sum2013_WWKa_ALT_44
PRJNA380388	SAMN06641852	NBDJ00000000	Escherichia coli Sum2013_WWKa_NEU_29
PRJNA380388	SAMN06641851	NBDK00000000	Escherichia coli Spr2013_WWKa_OUT_27

PRJNA380388	SAMN06641844	NBDL00000000	Escherichia coli Sum2013_WWKa_ALT_41
PRJNA380388	SAMN06641843	NBDM00000000	Escherichia coli Sum2013_WWKa_ALT_27
PRJNA380388	SAMN06641842	NBDN00000000	Escherichia coli Spr2013_WWKa_OUT_56
PRJNA380388	SAMN06641841	NBDO00000000	Escherichia coli Sum2013_WWKa_ALT_20
PRJNA380388	SAMN06641840	NBJM00000000	Escherichia coli Spr2013_WWKa_OUT_5
PRJNA380388	SAMN06641839	NBJN00000000	Escherichia coli Spr2013_WWKa_OUT_55
PRJNA380388	SAMN06641838	NBJO00000000	Escherichia coli Spr2013_WWKa_OUT_32
PRJNA380388	SAMN06641837	NBJP00000000	Escherichia coli Spr2013_WWKa_OUT_45
PRJNA380388	SAMN06641822	NBJQ00000000	Escherichia coli Spr2013_WWKa_OUT_15
PRJNA380388	SAMN06641821	NBJR00000000	Escherichia coli Spr2013_WWKa_OUT_29
PRJNA380388	SAMN06641820	NBJS00000000	Escherichia coli Spr2013_WWKa_NEU_6
PRJNA380388	SAMN06641819	NBJT00000000	Escherichia coli Spr2013_WWKa_OUT_11
PRJNA380388	SAMN06641818	NBJU00000000	Escherichia coli Spr2013_WWKa_NEU_15
PRJNA380388	SAMN06641817	NBJV00000000	Escherichia coli Spr2013_WWKa_NEU_37
PRJNA380388	SAMN06641816	NBJW00000000	Escherichia coli Spr2013_WWKa_ALT_63
PRJNA380388	SAMN06641815	NBJX00000000	Escherichia coli Spr2013_WWKa_ALT_71
PRJNA380388	SAMN06641814	NBJY00000000	Escherichia coli Spr2013_WWKa_ALT_51
PRJNA380388	SAMN06641813	NBJZ00000000	Escherichia coli Spr2013_WWKa_ALT_55
PRJNA380388	SAMN06641812	NBKA00000000	Escherichia coli Spr2013_WWKa_ALT_43
PRJNA380388	SAMN06641811	NBKB00000000	Escherichia coli Spr2013_WWKa_ALT_27
PRJNA380388	SAMN06641810	NBKC00000000	Escherichia coli Spr2013_WWKa_ALT_41
PRJNA380388	SAMN06641809	NBKD00000000	Escherichia coli Spr2012_WWKa_OUT_37
PRJNA380388	SAMN06641808	NBKE00000000	Escherichia coli Spr2012_WWKa_OUT_54
PRJNA380388	SAMN06641807	NBKF00000000	Escherichia coli Spr2012_WWKa_OUT_25
PRJNA380388	SAMN06641806	NBKG00000000	Escherichia coli Spr2012_WWKa_OUT_3
PRJNA380388	SAMN06641805	NBKH00000000	Escherichia coli Spr2012_WWKa_OUT_16
PRJNA380388	SAMN06641804	NBKI00000000	Escherichia coli Spr2012_WWKa_OUT_13
PRJNA380388	SAMN06641803	NBKJ00000000	Escherichia coli Spr2012_WWKa_NEU_74
PRJNA380388	SAMN06641802	NBKK00000000	Escherichia coli Spr2012_WWKa_OUT_12
PRJNA380388	SAMN06641801	NBKL00000000	Escherichia coli Spr2012_WWKa_NEU_31
PRJNA380388	SAMN06641800	NBKM00000000	Escherichia coli Spr2012_WWKa_NEU_51
PRJNA380388	SAMN06641799	NBKN00000000	Escherichia coli Spr2012_WWKa_NEU_24
PRJNA380388	SAMN06641798	NBKO00000000	Escherichia coli Spr2012_WWKa_ALT_27
PRJNA380388	SAMN06641797	NBKP00000000	Escherichia coli Spr2012_WWKa_ALT_35
PRJNA380388	SAMN06641796	NBKQ00000000	Escherichia coli Aut2013_WWKa_OUT_3
PRJNA380388	SAMN06641793	NBKR00000000	Escherichia coli Aut2013_WWKa_OUT_10
PRJNA380388	SAMN06641792	NBKS00000000	Escherichia coli Aut2013_WWKa_OUT_20
PRJNA380388	SAMN06641791	NBKT00000000	Escherichia coli Aut2013_WWKa_NEU_51
PRJNA380388	SAMN06641789	NBKU00000000	Escherichia coli Aut2013_WWKa_NEU_53
PRJNA380388	SAMN06641788	NBKV00000000	Escherichia coli Aut2013_WWKa_NEU_44
PRJNA380388	SAMN06641786	NBKW00000000	Escherichia coli Aut2013_WWKa_ALT_65
PRJNA380388	SAMN06641785	NBKX00000000	Escherichia coli Aut2013_WWKa_NEU_28
PRJNA380388	SAMN06641784	NBKY00000000	Escherichia coli Aut2013_WWKa_ALT_59
PRJNA380388	SAMN06641782	NBKZ00000000	Escherichia coli Aut2013_WWKa_ALT_48
PRJNA380388	SAMN06641780	NBLA00000000	Escherichia coli Aut2013_WWKa_ALT_45
PRJNA380388	SAMN06641779	NBLB00000000	Escherichia coli Aut2013_WWKa_ALT_30
PRJNA380388	SAMN06641778	NBLC00000000	Escherichia coli Aut2013_WWKa_ALT_17
PRJNA380388	SAMN06641777	NBLD00000000	Escherichia coli Aut2013_WWKa_ALT_13
PRJNA380388	SAMN06670745	NBNO00000000	Escherichia coli Win2012_WWKa_OUT_19

379

380 Table 3: Accession numbers of 92 de novo assembled *E. coli* genomes.

381

382

383

384

385

386
387

References

- 388 1. **Hu Y, Yang X, Qin J, Lu N, Cheng G, Wu N, Pan Y, Li J, Zhu L, Wang X, Meng Z, Zhao F, Liu D,**
389 **Ma J, Qin N, Xiang C, Xiao Y, Li L, Yang H, Wang J, Yang R, Gao GF, Zhu B.** 2013.
390 Metagenome-wide analysis of antibiotic resistance genes in a large cohort of human gut
391 microbiota. *Nat Commun* **4**:2151.
- 392 2. **Sommer MO, Dantas G, Church GM.** 2009. Functional characterization of the antibiotic
393 resistance reservoir in the human microflora. *Science* **325**:1128-1131.
- 394 3. **Salipante SJ, Roach DJ, Kitzman JO, Snyder MW, Stackhouse B, Butler-Wu SM, Lee C,**
395 **Cookson BT, Shendure J.** 2015. Large-scale genomic sequencing of extraintestinal pathogenic
396 *Escherichia coli* strains. *Genome Res* **25**:119-128.
- 397 4. **Forsberg KJ, Reyes A, Wang B, Selleck EM, Sommer MO, Dantas G.** 2012. The shared
398 antibiotic resistome of soil bacteria and human pathogens. *Science* **337**:1107-1111.
- 399 5. **Riesenfeld CS, Goodman RM, Handelsman J.** 2004. Uncultured soil bacteria are a reservoir
400 of new antibiotic resistance genes. *Environ Microbiol* **6**:981-989.
- 401 6. **Rizzo L, Manaia C, Merlin C, Schwartz T, Dagot C, Ploy MC, Michael I, Fatta-Kassinos D.**
402 2013. Urban wastewater treatment plants as hotspots for antibiotic resistant bacteria and
403 genes spread into the environment: a review. *Sci Total Environ* **447**:345-360.
- 404 7. **Schaufler K, Semmler T, Pickard DJ, de Toro M, de la Cruz F, Wieler LH, Ewers C,**
405 **Guenther S.** 2016. Carriage of Extended-Spectrum Beta-Lactamase-Plasmids Does Not Reduce
406 Fitness but Enhances Virulence in Some Strains of Pandemic *E. coli* Lineages. *Front Microbiol*
407 **7**:336.
- 408 8. **Yang J, Chen L, Sun L, Yu J, Jin Q.** 2008. VFDB 2008 release: an enhanced web-based resource
409 for comparative pathogenomics. *Nucleic Acids Res* **36**:D539-542.
- 410 9. **Jaureguy F, Landraud L, Passet V, Diancourt L, Frapy E, Guigon G, Carbonnelle E,**
411 **Lortholary O, Clermont O, Denamur E, Picard B, Nassif X, Brisse S.** 2008. Phylogenetic and
412 genomic diversity of human bacteremic *Escherichia coli* strains. *BMC Genomics* **9**:560.
- 413 10. **Clermont O, Bonacorsi S, Bingen E.** 2000. Rapid and simple determination of the *Escherichia*
414 *coli* phylogenetic group. *Appl Environ Microbiol* **66**:4555-4558.
- 415 11. **Lapierre P, Gogarten JP.** 2009. Estimating the size of the bacterial pan-genome. *Trends Genet*
416 **25**:107-110.
- 417 12. **McArthur AG, Waglechner N, Nizam F, Yan A, Azad MA, Baylay AJ, Bhullar K, Canova MJ,**
418 **De Pascale G, Ejim L, Kalan L, King AM, Koteva K, Morar M, Mulvey MR, O'Brien JS,**
419 **Pawlowski AC, Piddock LJ, Spanogiannopoulos P, Sutherland AD, Tang I, Taylor PL,**
420 **Thaker M, Wang W, Yan M, Yu T, Wright GD.** 2013. The comprehensive antibiotic resistance
421 database. *Antimicrob Agents Chemother* **57**:3348-3357.
- 422 13. **Andersson DI, Hughes D.** 2010. Antibiotic resistance and its cost: is it possible to reverse
423 resistance? *Nat Rev Microbiol* **8**:260-271.
- 424 14. **Pal C, Papp B, Lazar V.** 2015. Collateral sensitivity of antibiotic-resistant microbes. *Trends*
425 *Microbiol* **23**:401-407.
- 426 15. **Schrag SJ, Perrot V, Levin BR.** 1997. Adaptation to the fitness costs of antibiotic resistance in
427 *Escherichia coli*. *Proc Biol Sci* **264**:1287-1291.
- 428 16. **Touchon M, Hoede C, Tenaillon O, Barbe V, Baeriswyl S, Bidet P, Bingen E, Bonacorsi S,**
429 **Bouchier C, Bouvet O, Calteau A, Chiapello H, Clermont O, Cruveiller S, Danchin A, Diard**
430 **M, Dossat C, Karoui ME, Frapy E, Garry L, Ghigo JM, Gilles AM, Johnson J, Le Bouguenec C,**
431 **Lescat M, Mangenot S, Martinez-Jehanne V, Matic I, Nassif X, Oztas S, Petit MA, Pichon C,**
432 **Rouy Z, Ruf CS, Schneider D, Turret J, Vacherie B, Vallenet D, Medigue C, Rocha EP,**
433 **Denamur E.** 2009. Organised genome dynamics in the *Escherichia coli* species results in
434 highly diverse adaptive paths. *PLoS Genet* **5**:e1000344.
- 435 17. **Maurelli AT, Fernandez RE, Bloch CA, Rode CK, Fasano A.** 1998. "Black holes" and bacterial
436 pathogenicity: a large genomic deletion that enhances the virulence of *Shigella* spp. and
437 enteroinvasive *Escherichia coli*. *Proc Natl Acad Sci U S A* **95**:3943-3948.
- 438 18. **D'Costa VM, King CE, Kalan L, Morar M, Sung WW, Schwarz C, Froese D, Zazula G, Calmels**
439 **F, Debruyne R, Golding GB, Poinar HN, Wright GD.** 2011. Antibiotic resistance is ancient.
440 *Nature* **477**:457-461.

- 441 19. **Durso LM, Miller DN, Wienhold BJ.** 2012. Distribution and quantification of antibiotic
442 resistant genes and bacteria across agricultural and non-agricultural metagenomes. *PLoS One*
443 **7**:e48325.
- 444 20. **Manaia CM, Novo A, Coelho B, Nunes OC.** 2010. Ciprofloxacin Resistance in Domestic
445 Wastewater Treatment Plants. *Water Air and Soil Pollution* **208**:335-343.
- 446 21. **Gomi R, Matsuda T, Matsumura Y, Yamamoto M, Tanaka M, Ichiyama S, Yoneda M.** 2017.
447 Whole-Genome Analysis of Antimicrobial-Resistant and Extraintestinal Pathogenic
448 *Escherichia coli* in River Water. *Appl Environ Microbiol* **83**.
- 449 22. **Czekalski N, Berthold T, Caucci S, Egli A, Burgmann H.** 2012. Increased levels of
450 multiresistant bacteria and resistance genes after wastewater treatment and their
451 dissemination into lake Geneva, Switzerland. *Front Microbiol* **3**:106.
- 452 23. **Tenaillon O, Skurnik D, Picard B, Denamur E.** 2010. The population genetics of commensal
453 *Escherichia coli*. *Nat Rev Microbiol* **8**:207-217.
- 454 24. **Land M, Hauser L, Jun SR, Nookaew I, Leuze MR, Ahn TH, Karpinets T, Lund O, Kora G,
455 Wassenaar T, Poudel S, Ussery DW.** 2015. Insights from 20 years of bacterial genome
456 sequencing. *Funct Integr Genomics* **15**:141-161.
- 457 25. **Caucci S, Karkman A, Cacace D, Rybicki M, Timpel P, Voolaid V, Gurke R, Virta M,
458 Berendonk TU.** 2016. Seasonality of antibiotic prescriptions for outpatients and resistance
459 genes in sewers and wastewater treatment plant outflow. *FEMS Microbiol Ecol* **92**:fiw060.
- 460 26. **R Development Core Team.** 2010. R: A language and environment for statistical computing,
461 R Foundation for Statistical Computing, Vienna, Austria. <http://www.R-project.org>.
- 462 27. **Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ, Birol I.** 2009. ABySS: a parallel
463 assembler for short read sequence data. *Genome Res* **19**:1117-1123.
- 464 28. **Hashimoto M, Ichimura T, Mizoguchi H, Tanaka K, Fujimitsu K, Keyamura K, Ote T,
465 Yamakawa T, Yamazaki Y, Mori H, Katayama T, Kato J.** 2005. Cell size and nucleoid
466 organization of engineered *Escherichia coli* cells with a reduced genome. *Mol Microbiol*
467 **55**:137-149.
- 468 29. **Seemann T.** 2014. Prokka: rapid prokaryotic genome annotation. *Bioinformatics* **30**:2068-
469 2069.
- 470 30. **Li W, Godzik A.** 2006. Cd-hit: a fast program for clustering and comparing large sets of
471 protein or nucleotide sequences. *Bioinformatics* **22**:1658-1659.
- 472 31. **Zhou Y, Liang Y, Lynch KH, Dennis JJ, Wishart DS.** 2011. PHAST: a fast phage search tool.
473 *Nucleic Acids Res* **39**:W347-352.
- 474 32. **Delsuc F, Brinkmann H, Philippe H.** 2005. Phylogenomics and the reconstruction of the tree
475 of life. *Nat Rev Genet* **6**:361-375.
- 476 33. **Kumar S, Filipski AJ, Battistuzzi FU, Kosakovsky Pond SL, Tamura K.** 2012. Statistics and
477 truth in phylogenomics. *Mol Biol Evol* **29**:457-472.
- 478 34. **Price MN, Dehal PS, Arkin AP.** 2010. FastTree 2--approximately maximum-likelihood trees
479 for large alignments. *PLoS One* **5**:e9490.
- 480 35. **McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,
481 Altshuler D, Gabriel S, Daly M, DePristo MA.** 2010. The Genome Analysis Toolkit: a
482 MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res*
483 **20**:1297-1303.
- 484 36. **Wirth T, Falush D, Lan R, Colles F, Mensa P, Wieler LH, Karch H, Reeves PR, Maiden MC,
485 Ochman H, Achtman M.** 2006. Sex and virulence in *Escherichia coli*: an evolutionary
486 perspective. *Mol Microbiol* **60**:1136-1151.
- 487 37. **Larsen MV, Cosentino S, Rasmussen S, Friis C, Hasman H, Marvig RL, Jelsbak L, Sicheritz-
488 Ponten T, Ussery DW, Aarestrup FM, Lund O.** 2012. Multilocus sequence typing of total-
489 genome-sequenced bacteria. *J Clin Microbiol* **50**:1355-1361.
- 490 38. **Chen L, Yang J, Yu J, Yao Z, Sun L, Shen Y, Jin Q.** 2005. VFDB: a reference database for
491 bacterial virulence factors. *Nucleic Acids Res* **33**:D325-328.
- 492 39. **Antikainen J, Tarkka E, Haukka K, Siitonen A, Vaara M, Kirveskari J.** 2009. New 16-plex
493 PCR method for rapid detection of diarrheagenic *Escherichia coli* directly from stool samples.
494 *Eur J Clin Microbiol Infect Dis* **28**:899-908.
- 495 40. **Johnson JR, Russo TA.** 2005. Molecular epidemiology of extraintestinal pathogenic
496 (uropathogenic) *Escherichia coli*. *Int J Med Microbiol* **295**:383-404.

- 497 41. **Johnson JR, Stell AL.** 2000. Extended virulence genotypes of *Escherichia coli* strains from
498 patients with urosepsis in relation to phylogeny and host compromise. *J Infect Dis* **181**:261-
499 272.
- 500 42. **Pitout JD.** 2012. Extraintestinal Pathogenic *Escherichia coli*: A Combination of Virulence with
501 Antibiotic Resistance. *Front Microbiol* **3**:9.
502
- 503