

1 **Worldwide Population Structure of *Escherichia coli* Reveals Two**

2 **Major Subspecies**

3 Yu Kang^{1†}, Lina Yuan^{1†‡}, Zilong He^{1†}, Fei Chen^{1†}, Zhancheng Gao^{2†}, Shulin Liu^{3,4}, Xinmiao Jia¹, Qin
4 Ma⁵, Xinhao Jin¹, Rongrong Fu¹, Yang Yu^{6,7}, Chunxiong Luo⁸, Jiayan Wu¹, Jingfa Xiao¹, Songnian
5 Hu¹, Jun Yu^{1*}.

6 ¹ CAS Key Laboratory of Genome Sciences and Information, Beijing Institute of Genomics, Chinese Academy of
7 Sciences, 100101, Beijing, PR China.

8 ² Department of Respiratory & Critical Care Medicine, Peking University People's Hospital, Beijing, 100044, PR
9 China

10 ³ Genomics Research Center (one of The State-Province Key Laboratories of Biomedicine-Pharmaceutics of
11 China), Harbin Medical University, Harbin, 150081, PR China

12 ⁴ Department of Microbiology and Infectious Diseases, University of Calgary, Alberta, Canada, T2N 1N4.

13 ⁵ Department of Agronomy, Horticulture, and Plant Science, South Dakota State University, Brookings, SD,
14 57007, USA.

15 ⁶ School of Life Sciences, Liaoning University, Shenyang, 110036, PR China

16 ⁷ Department of Biology, University of Virginia, Charlottesville, Virginia, 22904, USA

17 ⁸ Center for Microfluidic and Nanotechnology, State Key Laboratory for Artificial Microstructures and Mesoscopic
18 Physics, School of Physics, Peking University, Beijing, 100871, PR China

19 * To whom correspondence should be addressed. CAS Key Laboratory of Genome Sciences and Information,
20 Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 100101, P.R. China. Tel: +8610-84097898;
21 Fax: +8610-84097540; Email: junyu@big.ac.cn

22 †The first 5 authors should be regarded as joint First Authors.

23 ‡Current Address: Lina Yuan, Center for Systems Medicine, Institute of Basic Medical Sciences, Chinese Academy
24 of Medical Sciences & Peking Union Medical College, Beijing100005; Suzhou Institute of Systems Medicine,
25 Suzhou 215123, P.R.China.

26

27 **Running Head:** Two Subspecies of *Escherichia coli*

28 **Word Count:**

29 Abstract 171

30 Text 5742

31

32 **ABSTRACT**

33 *Escherichia coli* is a Gram-negative bacterial species with both great biological diversity and important
34 clinical relevance. To study its population structure in both world-wide and genome-wide scales, we
35 scrutinise phylogenetically 104 high-quality complete genomes of diverse human/animal hosts,
36 among which 45 are new additions to the collection; most of them are clustered into two major clades:
37 Vig (vigorous) and Slu (sluggish). The two clades not only show distinct physiological features but
38 also genome content and sequence variation. Limited recombination and horizontal gene transfer
39 separate the two clades, as opposed to extensive intra-clade gene flow that functionally homogenizes
40 even commensal and pathogenic strains. The two clades that are genetically isolated should be
41 recognized as two subspecies both independently represent a continuum of possibilities range from
42 commensal to pathogenic phenotype. Additionally, the frequent intra-clade recombinant events, often
43 in larger fragments of over 5kb, indicates possibility of highly-efficient gene transfer mechanism
44 depending on inheritance. Underlying molecular mechanisms that constitute such recombinant
45 barrier between the subspecies deserve further exploration and investigations among broader
46 microbial taxa.

47 **IMPORTANCE**

48 The concept of bacterial species has debated over decades. The question becomes more important
49 today as human microbiomes and their health relevance are being studied extensively. The human
50 microbiomes where thousands of bacterial species co-habit need to be deciphered at minute details
51 and down to species and subspecies. In this study, we scrutinize the population genomics of *E. coli*
52 and define two subspecies that are distinct from each other concerning physiology, ecology, and
53 clinical features. As opposed to extensive genetic recombination within subspecies, limited genetic
54 flux between subspecies leads to their phenotypical distinctions and separate evolution paths. We
55 provide a key example illustrating that the divergence of a species into two subspecies depends on
56 recombination efficiency; when the recombination efficiency becomes a barrier the species appears
57 split into two. The *E. coli* scenario and its molecular mechanisms deserve further exploration in a
58 broader taxa of microbes.

59

60 INTRODUCTION

61 *Escherichia coli*, best known as a ubiquitous member of the normal intestinal bacterial microflora in
62 humans and other warm-blooded animals, is a Gram-negative bacterium of the family
63 Enterobacteriaceae. *E. coli* persists as a harmless commensal in the mucous layer of the cecum and
64 colon normally, whereas some variants have evolved to adapt pathogenic lifestyle that causes
65 different disease pathologies, including pandemic and lethal episodes (1). Depending on the site of
66 infection, pathogenic *E. coli* strains are divided into intestinal pathogenic *E. coli* (IPEC) and
67 extraintestinal pathogenic *E. coli* (ExPEC), which are able to successfully propagate both intra- and
68 extra-intestinally, respectively. Naturally, *E. coli* is a highly versatile species that survives in diversified
69 ecological habitats, such as sludgy environment of lake/river banks and tidal zones, as well as human
70 and animal intestines (2). Their environmental adaptation and lifestyle alteration, together with
71 experimental manipulation, make the species an excellent model for studying commensalism-
72 pathogenicity and genotype-phenotype relationships.

73 Phylogenetic analysis reveals that *E. coli* exhibits complex within-species sequence diversity that
74 hinders strain classification although various typing methods including ribo-typing, MLST (multi-locus
75 sequence typing), phylotyping, and whole genome phylogeny have been applied. Albeit complicated,
76 the population of *E. coli* is largely acknowledged as clonal, and four major phlotypes, A, B1, B2, and
77 D have been identified (3) which basically differ in habitat and life-style (4, 5). Phlotypes are loosely
78 associated to phenotypical characteristics, such as antibiotic-resistance and growth rate (3), and also
79 correlates with pathotype, as the ExPEC strains are normally part of B2 and D (6), whereas the IPEC
80 strains belong to A, B1, and D (7). However, recent genome-wide sequencing studies have revealed
81 dispensable/variable genes take a large part of genome plasticity contributing to biological diversity of
82 the species. Many virulence genes, including the most lethal Shiga toxins (Stx) and carbapenem
83 hydrolyase, are subject to frequent horizontal gene transfer (HGT), through which distinct pathogenic
84 and resistance phenotypes are acquired (8). Therefore, extensive HGTs interrupt the connections
85 between phenotypes and their mainstream phylogeny. Meanwhile, homologous recombination of the
86 *E. coli* core genome is found more frequent than previously realized, which also obscures phylogeny
87 and leads to either divergent or convergent characters (9). These observations on genetic flow
88 challenge the clonality of the species population, and raise the concern on the emergence of

89 disastrous “superbugs” if both lethal toxins and highly-resistant genes are recombined into a new
90 strain (10, 11).

91 To understand the population structure and genetic diversity within the species *E. coli*, we utilize
92 104 high-quality complete genome sequences from diverse geographic and host range; among them,
93 45 that have animal-host information are first released from our own sequencing efforts and the rest
94 of explicitly human hosts are from public databases. Our analysis results support the view that *E. coli*
95 is predominantly clonal, and except that a few strains of intermediate minor clades, most strains
96 cluster into two major clades. Each in successful adaptation to their own ecological niche and the two
97 clades distinct in physiological features, pathogenicity, genome content, and homologue sequences;
98 we propose that they deserve a permanent distinction as two subspecies. Further HGT analysis
99 reveals that recombination is very extensive within subspecies, which homogenizes strains into
100 continuum of genome possibilities, but rather limited when it happens in cross-subspecies manner.
101 The barrier of genetic exchanges between the two subspecies maintains clonal characters of the
102 species and drives them into separate evolutionary paths.

103

104 **RESULTS AND DISCUSSION**

105 **We contribute a unique fraction of complete genome sequences for a dataset used for**
106 **phylogenetic analysis.**

107 Our dataset for an in-depth analysis is composed of 104 complete genomes which include 59
108 human-host complete genomes from public databases and 45 newly added complete genomes of
109 animal-host isolates. The latter set is selected from a world-wide collection of 202 *E. coli* strains of
110 animal hosts, which includes strains with broad diversity in geography, climate, and host range, and
111 several ECOR (*Escherichia coli* collection of reference) strains (12). Our effort includes the
112 identification of the MLST (multi-locus sequence typing) sequences for constructing a MLST-based
113 phylogeny. The MLST-based phylogeny reveals a complex partition among animal or human hosts
114 (Figure 1A), and at the end, 45 animal-host isolates is brought on to represent diverse origin and
115 genetic heterogeneity in terms of host diet, geographic distribution (Figure 1B), as well as MLST-
116 clusters for further genome sequence finishing; their full genome sequences reveal similar

117 chromosomal organizational features to human-host isolates, such as a uniform G+C content of
118 50.58%. Their genome sizes range from 4.25 to 4.94Mb and are predicted to encode 4,728 genes in
119 average, a slightly smaller than that of pathogenic strains, but similar to commensals.

120 To construct the genome-based phylogeny, 7 draft genomes of strains isolated from wild
121 environment (13), which are phylogenetically distant from host-related strains, are included as
122 outgroups. The core genome shared by all the 111 genomes is composed of 1,095 genes and
123 collectively 1.05 Mb in length. The maximal-likelihood phylogeny of the core genome indicates
124 ancestor position of the environmental strains (Figure 2) and the host-related strains, regardless of
125 human or animal origin, are mingled together. The majority of the host-related strains clusters into two
126 major clades, each contains 56 and 31 strains respectively. The rest of 17 strains, holding closer
127 positions to the ancestral environment strains, is split into several minor clades. Overall, this
128 clustering pattern is largely in congruence with their phylotypes as previously reported (9, 14).

129 **The two major clades are physiologically distinct and adaptive to each ecological niches.**

130 The two major clades of host-related strains also exhibit distinctions in host, climate (5), and
131 pathogenicity (Table S1). According to their distinct characters, such as movement and growth, we
132 name them Vig (Vigorous) and Slu (Sluggish) clades just for convenience. The Vig, composed of
133 strains of phylotype A and B1, is featured for its strains of carnivores host and tropical geographic
134 distribution; all *E. coli* strains that have led to human pandemic infections belong to this particular
135 clade. Although some Vig strains also survive in herbivores and cold area, and many strains are
136 commensals, it appears that the Vig strains prefer warmer temperature, richer amino acid nutrients,
137 and are able to propagate rapidly under ideal conditions, albeit adapt to a wide range of ecological
138 niches. On the contrary, the Slu strains, which fall in phylotype B2, are only found in herbivores and
139 omnivores and colder climates. Besides commensals, Slu strains occasionally cause extra-intestinal
140 and antibiotic-resistant infections that are rarely seen among the Vig strains. These ecological
141 features of Slu strains suggest that they are more adaptive to lower temperatures, poor amino acid
142 nutrients, and therefore exhibit slower growth rate. This adaptation to low temperature of the Slu clade
143 (mainly B2 strains) coincides with the report of a large population investigation (15), and facilitates
144 transient period of epizooism and migration to extra-intestinal habitats together with their tolerance to

145 poor nutrition. Meanwhile, its slow growth rate increases survival rate under stresses, such as those
146 of antibiotics, and offers better opportunities to gain antibiotics-resistant genes or functional mutations.

147 To confirm some of the distinct features between the clades Vig and Slu, such as optimal
148 temperature and amino acid preference, we design a few straightforward physiological experiments
149 with resuscitated strains in our collection (23 Vig and 15 Slu strains). First, we measure growth rates
150 for the strains at various temperatures and find that the Vig strains grow significantly faster than the
151 Slu strains at 37°C and 41°C, which resemble the intestinal environment of warm-blood animals
152 (Figure 3A). However, at lower temperature of 27°C and 32°C, their differences are not significant.
153 The results well explain the higher prevalence of the Vig strains and the preference of Slu strains to
154 colder climates. Second, we compare their survival rates after heat shock at 50°C, and the faster-
155 growing Vig strains exhibit vulnerability and poorer survival rates as compared to the Slu strains after
156 20 and 30 minutes of heat stress (Figure 3B), which is in congruent with the measurement of their
157 growth rates. Third, we measure mobility of these strains in response to chemotaxis amino acid—
158 phenylalanine—in a custom-designed microfluidic device. At each time point, more cells of the Vig
159 strains reach the destination pool at the same chemoattraction, indicating their faster response to
160 amino acids and higher mobility than the Slu (Figure 3C). This ability of Vig strains ensures rapid
161 approaching toward amino acid nutrients and allocates to their carnivoral inhabitation. Finally, we
162 calculate the strain growth rate in various carbohydrate sources relative to glucose. All *E. coli* strains
163 grow much rapidly in monosaccharides and disaccharides than in polysaccharides. However, when
164 compared to the Slu strains, the Vig strains grow slightly faster in monosaccharides and
165 disaccharides, but a little slower in polysaccharides (Figure 3D). Although the difference is not
166 significant possibly due to small sample size, it seems that the Slu strains may be more adaptive to
167 herbivore host where polysaccharides are more abundant. From all the physiological experiments, it
168 is apparent that the two clades have diverged from each other in phenotypical features in adaptation
169 to distinct ecological niches as well as leading to distinct clinical features.

170 **Genomic distinctions between the two clades correspond to their phenotypes.**

171 The distinct physiological and ecological features between clades Vig and Slu lead to the
172 speculation that the two clades should be distinct in genome content and homologous sequences
173 especially in genes related to metabolism, energy, and mobility. To confirm this, we formulate a

174 parameter – pair-wised genome content distance – of all strains, finding that the Vig and Slu strains
175 are clearly separated in a neighbor-joining tree that derives from the distance matrix (Figure 4A); the
176 result indicates that the two clades have a significant number of different genes. We subsequently
177 apply two-sample Kolmogorov-Smirnov test to all dispensable genes and identify 279 and 336
178 orthologues that exhibit significant enrichment in the Vig and Slu clades, respectively (Table S2).
179 These genes represent genome content that are characteristic of the two clades; when annotated
180 according to the COG database for their functions, as expected, over half of them are in the
181 categories of metabolism. However, although different between the Vig and Slu clades, the
182 distribution (presence/absence) of these characteristic genes does not vary so much with host dietary
183 preference (Figure 4B). For metabolizing carbonates and lipids, the Vig and Slu contain a set of
184 diversely characteristic genes. For genes of mobility related, purine and amino acid metabolic, the Vig
185 strains are apparently richer than the Slu, which are in accordance with their adaptation to carnivore
186 intestines (Figure 4B).

187 Next, we search for all virulent genes within dispensable genes. For both clades, unlike metabolic
188 genes, the content of virulent genes varies greatly among strains. Pathogenic strains have much
189 more virulent genes than commensals, including all kinds of toxins and iron-uptake genes. However,
190 the two clades differ in their virulent gene content: the Vig pathogens are rich of T3SS and other
191 secretion system, whereas the Slu pathogens have more adhesion and invasion genes that facilitates
192 extra-intestinal infections (Figure 4C). We also scrutinize beta-lactamase genes and the notorious
193 lethal Shiga-like toxins (Stx) that often leads to pandemic infections when transferred to *E. coli* strains
194 (16). These genes are usually carried in plasmid, but can be inserted into chromosome by mobile
195 elements under strong selections (17). Although both clades contain strains with narrow spectrum
196 beta lactamases as TEM-1 and OXA-1 (18), the extended spectrum beta lactamases (ESBLs) are
197 rarely found among Vig strains, whereas Stx exclusively shows up in Vig strains. (Figure 4D). These
198 differences between the two clades in their genome content lead to explanation for their phenotypical
199 and pathogenic characteristics.

200 In addition to gene content, we further investigate sequence variances between homologs that are
201 shared by the two clades. First, the Vig and Slu clades are phylogenetically distinct in their core
202 genome, and from the orthologous pairs we identify 126 and 227 non-synonymous polymorphic sites

203 which reside in 97 and 168 genes and are specific for the Vig and Slu clades, respectively (Table S3).
204 These polymorphic sites exclusively found in each clade are named as lineage-associated variations
205 or LAVs. The ratios of non-synonymous to synonymous polymorphism (dN/dS) of these LAV-
206 containing genes are in average 0.02 and below 0.4, as if they are under strong purifying selection.
207 We extract all metabolic genes in carbonate and amino acid pathways from the core genes and
208 concatenate them to construct maximal-linkage tree. Clearly, genes of Vig and Slu strains cluster
209 separately (Figure 5A), likely to contribute to their metabolic characteristics. We also investigate
210 homologous sequences of the highly prevalent dispensable genes (orthologues present in over 80%
211 strains of the two clades and calculate pair-wised gene distances based on amino acid sequences
212 between orthologues. The result shows much larger distances between inter-clade pairs than intra-
213 clade pairs (Figure 5B), indicating that each clade utilizes their own preferred alleles for highly shared
214 dispensable genes. Finally, we check the GC content difference between the two clades, and
215 unexpectedly find that GC% of the Vig strains are a slightly, but significantly, higher than that of the
216 Slu strains in the core genes, dispensable genes, and whole genome (Figure 5C). We speculate that
217 the underlying reasons are adaptation to the higher temperature habitat of the Vig clade. In any case,
218 the difference in GC content definitely leads to a global effect on nucleic acid composition of all genes
219 and can be an independent factor driving genome evolution (19). Apparently, all above genetic
220 characteristics in genome content and sequence diversity between the two clades may explain in part
221 their phenotypic/physiological differences and molecular mechanisms for adaptation to their distinct
222 ecological niches.

223 **Barrier to inter-clade recombination separates the clades genetically into subspecies.**

224 *E. coli* is generally clonal (20, 21) and the Vig and Slu strains fit in such a framework, diverging
225 from each other through accumulation of independent mutations and limiting exchange of genetic
226 materials mutually. However, in previous studies, the *E. coli* recombination rate is evaluated as
227 comparable to a mutation rate with the ratio of $r/m \approx 0.9$ (9, 22, 23). Theoretically, such a
228 recombination rate is able to confound clonal framework and intermingles strains into an unstructured
229 population. The question becomes how the species keeps clonal structure and sustains a relatively
230 high recombination rate.

231 To scrutinize clonal status, we first evaluate general recombination rate based on the concept of
232 homoplasy – polymorphisms shared by two or more strains but not present in their common ancestor
233 (22, 24). Practically, homoplastic polymorphisms can only be inferred for core genome where
234 recombination are mainly introduced through homolog sequences. Our results indicate that the r/m of
235 the *E. coli* core genome is about 0.5, much lower than previously estimated (22) (Table 1). The
236 reason is the fact that dispensable genes are more subjected to HGT and lead to over-estimated r/m
237 ratio when they landed in a core genome. The large sample size of our current study results in a
238 minimal definition of the core genome that excludes almost all the dispensable genes and thus
239 narrows r/m ratio. However, the r/m ratios of the Vig and Slu clades are 0.949 and 0.745, much higher
240 than the entire species (Table 1); the deviations point out the fact that there are more within-clade
241 recombination events than that of between-clade, as it has been proposed in a previous study (9).
242 The high rate of within-clade homologous recombination is further confirmed by an analysis based on
243 a Bayesia-based method—BRATNextGen, which has been used for analysing homologous
244 recombination events between and within clades on the basis of a specified degree of sequence
245 divergence (25, 26). The result illustrates the fact that strains share more recombinant fragments
246 within clades than cross clades (Figure 6A). Both analyses demonstrate that *E. coli* strains rarely
247 exchange genetic material with distant relatives whereas closely related strains recombine more
248 frequently, and thus reconciling the controversy between clonal structure and high overall
249 recombination rate of the species.

250 The above estimations of recombination rate are applied for the core genome which only takes
251 less than one-fourth of the average *E. coli* genome, however, recombination rate varies with
252 chromosomal regions, and genes differ in their possibilities of being transferred—dispensable
253 genome contains more mobile elements and are more prone to be horizontally transferred (27) To
254 obtain complete understanding of genetic exchange in term of whole genome and all strains of entire
255 species, we try to scan recombinant events across the entire genome length; recombination plays
256 critical roles in transferring and shuffling dispensable genes that shape genome content in significant
257 ways (28). In general, recombinant events between close-related strains are not easy to identify due
258 to weak signals of recombined sequences (29, 30). Based on sequence alignment, we first identify all
259 near identical fragments for each genome pair as candidate recombinant fragments (31). Since
260 recombinant often insert genes into different location, using complete genome sequences, we also

261 compare synteny (linear order) of the candidate recombinant fragments and define non-syntenic
262 fragments as results of true recombinant events. Our result shows that strains with intra-clade
263 genome-pairs exchange more genome content (nearly 10x) than those with inter-clade pairs; the sum
264 of total recombinant fragment lengths approaches one half of the genome in some cases and never
265 less than 300kb which is even larger than the upper limit of inter-clade pairs (Figure 6B). The
266 extensive intra-clade genetic exchange in the dispensable genome strengthens the clonal structure
267 that is defined by the core genome. Furthermore, we find that intensive genetic material exchange
268 where virulent genes can be included between commensal and pathogenic strains depends on their
269 clade statuses: higher in within-clade events (Figure 6B).

270 We further correlate recombinant frequency (here we use the number of recombination fragments
271 per genome-pair) to phylogenetic distance between paired genomes. When plotted against core
272 genome distance (from which the core genome phylogeny has been inferred) for both Vig and Slu
273 strains, the regression curves of recombination frequency shows a reverse-S shape, i.e., there is a
274 rapid decline over transition between intra- and inter-clade genome-pairs (Figure 6C). Strains of the
275 same clade recombine more frequently, which often have nearly one hundred recombinant fragments
276 per genome, and show a decrease in the number of recombination fragments due to overlapping of
277 larger fragments when very closely related. In fact, very close sister strains can exchange almost half
278 of their genome. The phenomenon of chromosomal fragment transfer over 100kb or even 1Mb has
279 been experimentally validated in other species (32, 33). However, the very efficient genetic exchange
280 of large fragment and high frequency is not seen between clades, where such fragments rarely
281 exceed 5kb (Figure 6D). The result implies that closely related strains (within clade) may have a
282 unique highly-efficient molecular mechanism for recombination, whereas genetic material exchange
283 between remote relatives (cross clade) is confined to some low-efficient ways.

284 **Population structure of host-related strains and relationship between phylotype and** 285 **phenotype**

286 *E. coli* is a well-studied gram-negative bacterial species due to its importance in both clinical
287 practice and biological research. However, biologically meaningful strain or population classification
288 based on genotypes and phenotypes as well as other features is still a difficult task. Our study starts
289 from phylogeny based on high-quality whole genome and, through detailed genomic analysis and

290 experimentation ends with the definition of two major clades – Vig and Slu. The two clades are distinct
291 in many aspects and their clade-centric characteristics explain many ecological and clinical features.
292 *E. coli* infections fall into two categories with different clinical outcomes and treatment strategies: 1)
293 acute intestinal infections with various severities and 2) opportunistic infections often in extraintestinal
294 loci and often resistant to common antibiotics. The two types infections appear to correlate with
295 different clinical features of the two clades and their traditional phylotypes (34, 35). Genomic analysis
296 reveals that the genetic structure underlying the physiological, ecological and clinical traits of Vig and
297 Slu strains are a pile of characteristic genes and sequence variations, especially genes involved in
298 metabolic pathways, mobility, and toxic or resistant phenotypes. And the separation of the two clades
299 is caused by the limitation of between-subspecies genetic exchange or recombination, as similar
300 scenarios found among other species of bacteria (36-38), archaea (39), and eukaryotes (40). Among
301 the characteristic genes, the extremely virulent toxin—Shige-like toxin and the most notorious
302 antibiotic resistant genes—ESBLs, are partitioned into the Vig and Slu with very little overlap.
303 Although these genes are often carried by plasmids and ready to transfer among strains, it seems that
304 strains carrying both has been reported to be rather rare (41), also supporting the between-clade
305 recombinant barrier. Therefore Vig and Slu should be regarded as two subspecies since the
306 recombinant barrier has genetically separated them and made them clearly divergent in all aspects of
307 physiology, ecology and clinical significance. In clinical, identification of subspecies for a pathogenic
308 strain will give informative clinical guidance for the treatment of the infection it caused. On the other
309 side, the genetic boundary between commensals and pathogenic strains is not very clear. Although
310 pathogenic strains bear much more toxic genes, a commensal strain can exchange genetic material
311 and acquire enough virulent genes from its close pathogenic sisters, and then become pathogenic
312 under appropriate host conditions. Intensive recombination readily alters virulent gene content and
313 thus blurs the lines between commensal and pathogen, making them genetically undistinguishable in
314 clinical practice (6).

315 Since species are “lineages evolving separately from other lineages” (42), genetic diversification
316 and geographic separation of *E coli* subspecies, represented by the Vig and Slu clades, demonstrate
317 an early process of microbial speciation, whose mechanisms and processes have been debated over
318 decades (43, 44). Until recently, technological innovation, especially the invention of next-generation
319 sequencing (NGS) technology, coupled with the emerging discipline of population genomics, has

320 been providing unprecedented tools and opportunities for the interrogation of molecular details on
321 many ongoing evolutionary processes among natural microbial populations, especially those with
322 healthcare applications. The emergence of the two *E.coli* subspecies appears initiated from distinct
323 genetic units with functional relevance, which are incubated and frequently traded among closely
324 related, structurally comparable, and geographically cohabitating strains through mutually beneficial
325 mechanisms. Our recombinant analysis reveals that the within-subspecies recombination rate is
326 much more significant than that of between-subspecies, and such a diversifying process eventually
327 drives subspecies or their populations keeping evolving into nascent species (45). In our data,
328 recombination has overall effects on both core and dispensable portions of the genome, and results in
329 hundreds of characteristic genes as well as lineage-associated variations or LAVs; these genetic and
330 functional elements form a complex background for species and its population to evolve under nature
331 selection. Certainly, physical barriers that interrupt recombination, accelerating the process of
332 speciation (43). In the case of the two *E coli* subspecies, both are widely spread and co-inhabiting,
333 such as commensals in intestines of both humans and other omnivores, and our study and
334 observations of them does not support the hypothesis of geographic isolation. Therefore, other types
335 of physical barriers such as CRISPR-Cas system (46), restriction-modification system (47), DNA
336 uptake signal sequences (48), and incompatible transfer mechanisms due to pili (49), which have
337 been reported in some species, may play roles in *E. coli* sub-speciation or speciation but remain to be
338 elucidated. Our results highlight rate difference between intra-subspecies and inter-subspecies
339 recombination as a barrier possibly due to less functional benefits. Some high-efficiency mechanisms,
340 such as distributive conjugal transfer, has been reported in a species of Mycobacteria, which is able
341 to transfer fragments over 100kb at one time but lose function when such genetic exchange
342 happened between remote relatives (50, 51). On the other hand, low-efficiency mechanisms, such as
343 phage (52) or transposon (53), usually transfer short fragments at lower frequency (54), but may still
344 work when happening across subspecies due to broader host range. These mechanisms underlying
345 cross-subspecies recombinant barrier deserve further exploration and should be investigated in a
346 broad range of microbial taxa for better understanding of microbial population structural dynamics and
347 speciation.

348

349 CONCLUSION

350 Based on an unprecedented dataset, we thoroughly studied the population structure and
351 dynamics of *E. coli*, including genetic diversity, habitat divergence, physiological features,
352 recombination rate, and gene flow. We defined two *E. coli* subpopulations among large number of
353 isolates and suggest that they appear to be distinct subspecies that have evolved to bear
354 characteristic gene content and sequence variance, which lead to distinct physiological, ecological,
355 and clinical characteristics. There is an apparent barrier of recombinant between the two subspecies,
356 which drive their genetic diversification. Although the underlying mechanism still needs further
357 demonstration, novel molecular mechanisms differentiating intra- and inter-subspecies genetic
358 material exchange may exist. The discovery of such mechanisms and confirmation in broad range of
359 microbial taxa will surely deepen our understanding of the process of bacterial speciation.

360

361 MATERIAL AND METHODS

362 **Strains and MLST typing.** A world-wide collection of 202 *Escherichia coli* strains from vertebrate
363 hosts (12) was kindly provided by Professor Shulin Liu (Genomics Research Center, Harbin Medical
364 University, Harbin, China; Department of Microbiology and Infectious Diseases, University of Calgary,
365 Calgary, Canada). Genomic DNA were extracted using a Qiagen DNeasy kit (Qiagen) for 172
366 successfully resuscitated isolates. PCR and Sanger sequencing were performed for 16S rDNA, MLST
367 genes (seven housekeeping genes of *adk*, *fumC*, *gyrB*, *icd*, *mdh*, *purA*, and *recA*, see details at
368 <http://mlst.ucc.ie/>) for these strains. Then five strains whose 16S rDNA sequence showed <97%
369 identity with *E. coli* reference or hit best to other species and 18 strains whose MLST alleles were
370 failed to be identified by this method were removed, leaving 149 animal-host strains for further
371 analysis. We further included 59 complete genomes of human-host deposited in NCBI
372 (<ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria/>) as of July 2013 and seven draft genomes published by
373 a study of environmental strains (13) from ftp://ftp.ncbi.nlm.nih.gov/genomes/Bacteria_DRAFT/, then
374 *in silico* identified their MLST loci with BLAST. Together, we aligned the seven MLST fragments of all
375 213 strains with ClustalW and concatenated them for inference of maximum-likelihood phylogeny
376 using RaxML with model GTR and 1,000 bootstraps.

377 **Genome sequencing, annotation and phylogeny inference.** Genomic DNA libraries of 45 selected
378 representative strains were prepared using a NEBNext DNA Library Prep Kit and sequenced on a
379 HiSeq 2000 for 2 × 100 bp run at the Beijing Institute of Genomics (Beijing, P.R.China). The raw reads
380 were quality filtered and trimmed using SolexaQA (-p 0.01), with an average coverage of 150× and
381 assembled with SOAPdenovo (S) (55), and then the assemblies were scaffolded into circular
382 genomes with GAAP (<http://gaap.big.ac.cn>), which is based on synteny of core genes on genomic
383 scale, and assist assembly of high quality genomes (27). For these genomes, protein-coding genes
384 were predicted by using GeneMark.hmm with a pre-trained model and annotated using BLAST
385 against the COG database. Gene sequences were mapped to metabolic pathways by using BLAST
386 against KEGG GENES for KEGG Orthology assignment using the KEGG automatic annotation server.
387 The identification of orthologous genes was performed with pan-genome analysis pipeline (PGAP) (56)
388 with identity and coverage of pairwise genes set at 0.7. Orthologues common to all 111 strains
389 (including 45 animal-host, 59 human-host and seven environment strains) are identified as core
390 genes, and others (orthologues shared by a portion of strains) as dispensable genes. The 1,095
391 single-copy core genes were concatenated into core genome and used to construct phylogenetic tree.
392 We aligned protein sequences of each core genes using ClustalW, and traced them back to nucleic
393 acid sequences. The maximum-likelihood phylogeny was inferred on the basis of SNPs in the core
394 genome based on RaxML in the model GTR+G with 1,000 bootstraps. Phylogenetic trees were
395 viewed and modified by using FigTree and EvolView (57). Phylotypes of all strains were identified *in*
396 *silicon* according to the presence of three phylotype-specific genes or fragments as previously
397 described (58).

398 **Measurement of growth rate, survival rate, and mobility.** Strains were resuscitated and cultured in
399 pre-filtered LB broth. Unless otherwise specified, each strain was normalized by cultivation at 17°C in
400 LB overnight with shaking at 200 rpm, and OD_{600nm} was measured for each culture in 96-well plate
401 (triplicates for each strain with a blank control) on an Infinite® 200 PRO Microplate Reader (Tecan)
402 every half hour.

403 *Growth rate in diverse temperature and carbonate source.* Overnight cultures were 1:50 diluted with
404 LB broth and incubated at 27, 32, 37, and 41°C or change medium supplemented with glucose, starch,

405 glycogen, heparin, cellulose, fructose, lactose, and maltose as sole carbon source. OD₆₀₀ was
406 monitored until the cultures reach 0.6.

407 *Survival rate at heat stress.* Overnight cultures of each testing strain were adjusted to OD₆₀₀=0.1.
408 Immediately, and 50µl of inoculum was diluted 1:4 with LB broth and incubated at 50°C for 0min, 10
409 min, 20min, and 30min. OD₆₀₀ was monitored until the cultures reach 0.1. The number of living cells in
410 each initial aliquot relative to a negative control was calculated as $e^{-\Delta t}$, where Δt is the time to
411 OD₆₀₀=0.1. The survival rate was defined as the ratio of live cell number after heat shock to that of
412 negative control.

413 *Mobility and response to chemoattractive amino acids.* Each strain was evaluated on a custom-
414 designed PDMS microfluidic device as previously described (59). Briefly, there were two pools—sink
415 and source pool—which were connected with a 600µm channel. At the end of the channel adjacent to
416 the source pool, a thin layer of precast gel obstructed cells into the source pool, in front of which a
417 400µm² observation channel was set there for cell counting. Normalized cells of single clones were
418 re-suspended in an amino acid-free basal medium with OD₆₀₀=0.1. Immediately, 30 µl of cells and 30
419 µl of chemotaxis solution (with 100 µM phenylalanine in basal medium) were injected in the sink and
420 source pool respectively. Cells reached the observation channel was immediately observed and auto-
421 tracked by using a Nikon Ti-E inverted microscope system every 5 min for 1 h. The cell count of
422 observation channel in each images, which serves as an indicator of cell mobility in response to
423 chemoattractive amino acid, was automatically measured by using ImageJ software.

424 **Genome content similarity, gene distribution, and homologue distance between Vig and Slu.**

425 The similarities between paired genomes were calculated on the basis of the Bray–Curtis dissimilarity
426 index. A dissimilarity index of d is calculated as $1 - [2 * S_{ij}/(S_i + S_j)]$, where S_{ij} is the number of
427 dispensable genes shared by strains i and j , and S_i and S_j are the numbers of dispensable genes in
428 strains i and j , respectively. Pairwise dissimilarity indices were used to construct a distance matrix,
429 which was used to construct the neighbour-joining tree, and genome similarity ($1 - d$) was used to
430 produce a heatmap. Then we applied two-sample Kolmogorov-Smirnov test for each dispensable
431 orthologue on its balanced distribution/presence in Vig and Slu. When $p < 0.01$, the orthologue was
432 significantly enriched in Vig or Slu. Antibiotic-resistance genes were identified by using BLAST
433 against CBMAR database (<http://14.139.227.92/mkumar/lactamasedb>) with E value of 10^{-5} and best hit.

434 Similarly, Shiga-like toxin genes were identified by BLAST against the protein sequences of that from
435 *E. coli* O157:H7 (|cl|AB015056.1_prot_BAA88123.1_1 for A-subunit and
436 |cl|AB015056.1_prot_BAA88124.1_2 for B-subunit) downloaded from NCBI database. To calculate
437 pairwise sequence distances for protein between homologs, individual orthologues of high-prevalence
438 dispensable genes (present in >80% of strains) were first aligned with ClustalW, and then calculated
439 pairwise distance of genes in one orthologue from various strains of one species by using Protdist in
440 the Phylip package.

441

442 **Recombination inference.** We calculated r/m statistics for the core genome of *E. coli* using a
443 computational method PHI which recognize homoplastic sites of amino acid sequences in a 1-kb
444 overlapping window. We used all the default parameters and collected homoplastic sites in windows of
445 $p < 0.01$ as polymorphism sites caused by recombination (r) and the others as those caused by
446 mutation (m). We also inferred recombination fragments of core genome with a Bayesian algorithm-
447 based method—BratNextGen (25). The strain clusters were a priori defined on the basis of a PSA
448 tree. In the default procedure, alpha was set at 3.58. One hundred iterations were performed until the
449 model parameters converged, and the significances of inferred recombinant fragments ($p < 0.05$) were
450 assessed using 100 replicate permutations of sites in the genome.

451 We utilized a programme— gmos (Genome MOsaic Structure) to compute local alignments between
452 paired query and subject genomes and reconstructed the query mosaic structure of recombinant
453 fragments over 600bp (31). These fragments, although almost identical caused by very recent
454 recombination, still held the possibility of vertically inherited under strong selection, and thus only
455 regarded as candidate recombinant fragments. We also identified fragments that changed their
456 genome locations as parsimony recombination fragments for each genome pair. To achieve this, we
457 assigned consecutive number for each candidate fragment along their order in subject genome. We
458 also recorded their relative order in the query genome as a new sequence where we identified the
459 longest increasing/decreasing subsequence as fragments that keep their original locations, and the
460 other fragments were recognized as parsimonious recombination fragments.

461 **Data Availability.** The final 45 genome sequences we contributed were deposited in GenBank under
462 the accession numbers CP012758~CP012800, CP012806, and CP012807.

463

464 **LIST OF ABBREVIATIONS**

465 dN/dS, ratios of non-synonymous to synonymous polymorphism

466 ECOR, *Escherichia coli* collection of reference

467 ESBLs, extended spectrum beta lactamases

468 ExPEC, extraintestinal pathogenic *E. coli*

469 HGT, horizontal gene transfer

470 IPEC, intestinal pathogenic *E. coli*

471 LAVs, Linage Associated Variations

472 MLST, multi-locus sequence typing

473 NGS, next-generation sequencing

474 *r/m*, ratio of polymorphisms caused by recombination to mutation

475 Slu Sluggish clade

476 Stx, Shiga toxins

477 Vig Vigorous clade

478

479 **DECLARATION**

480 **Availability of data and materials:** The final 45 complete genomes contributed by our laboratory has
481 been deposit in the GenBank under the accession numbers CP012758~CP012800, CP012806, and
482 CP012807.

483 **Competing interests:** The authors declare no competing interests.

484 **Funding:** This work was supported by the National Scientific Foundation of China [31470180,
485 31471237, and 30971610].

486 **Author contribution:** J.Y. and Y.K. conceived the project and led the writing; S.L. provided the
487 strains; L.Y., Z.H., X.J., R.F., Y.Y., and C.L. compiled the data; and C.F., Z.G., X.J., X.G., Q.M., Y.Z.,
488 J.W., J.X., and S.H. analysed the data. All authors contributed to the writing and/or intellectual
489 development of the manuscript.

490 **Acknowledgement:** Acknowledgement to W. Chen who provided manuscript feedback.

491 REFERENCES

- 492 1. Leimbach A, Hacker J, Dobrindt U. 2013. E. coli as an all-rounder: the thin line between
493 commensalism and pathogenicity. *Curr Top Microbiol Immunol* 358:3-32.
- 494 2. Anonymous. 2017. E.coli (Escherichia coli) | E.coli | CDC, on Centers for Disease Control and
495 Prevention. <https://www.cdc.gov/ecoli/>. Accessed
- 496 3. Gordon DM. 2004. The Influence of Ecological Factors on the Distribution and the Genetic
497 Structure of Escherichia coli. *EcoSal Plus* 1.
- 498 4. Carlos C, Pires MM, Stoppe NC, Hachich EM, Sato MI, Gomes TA, Amaral LA, Ottoboni LM.
499 2010. Escherichia coli phylogenetic group determination and its application in the
500 identification of the major animal source of fecal contamination. *BMC Microbiol* 10:161.
- 501 5. Gordon DM, Cowling A. 2003. The distribution and genetic structure of Escherichia coli in
502 Australian vertebrates: host and geographic effects. *Microbiology* 149:3575-86.
- 503 6. Picard B, Garcia JS, Gouriou S, Duriez P, Brahimi N, Bingen E, Elion J, Denamur E. 1999. The
504 link between phylogeny and virulence in Escherichia coli extraintestinal infection. *Infect*
505 *Immun* 67:546-53.
- 506 7. Pupo GM, Karaolis DK, Lan R, Reeves PR. 1997. Evolutionary relationships among pathogenic
507 and nonpathogenic Escherichia coli strains inferred from multilocus enzyme electrophoresis
508 and mdh sequence studies. *Infect Immun* 65:2685-92.
- 509 8. Martinez-Castillo A, Muniesa M. 2014. Implications of free Shiga toxin-converting
510 bacteriophages occurring outside bacteria for the evolution and the detection of Shiga toxin-
511 producing Escherichia coli. *Front Cell Infect Microbiol* 4:46.
- 512 9. Didelot X, Méric G, Falush D, Darling AE. 2012. Impact of homologous and non-homologous
513 recombination in the genomic evolution of Escherichia coli. *BMC Genomics* 13:256.
- 514 10. Zhang L, Levy K, Trueba G, Cevallos W, Trostle J, Foxman B, Marrs CF, Eisenberg JN. 2015.
515 Effects of selection pressure and genetic association on the relationship between antibiotic
516 resistance and virulence in Escherichia coli. *Antimicrob Agents Chemother* 59:6733-40.
- 517 11. Pitout JD. 2012. Extraintestinal Pathogenic Escherichia coli: A Combination of Virulence with
518 Antibiotic Resistance. *Front Microbiol* 3:9.
- 519 12. Souza V, Rocha M, Valera A, Eguiarte LE. 1999. Genetic structure of natural populations of
520 Escherichia coli in wild hosts on different continents. *Appl Environ Microbiol* 65:3373-85.
- 521 13. Luo C, Walk ST, Gordon DM, Feldgarden M, Tiedje JM, Konstantinidis KT. 2011. Genome
522 sequencing of environmental Escherichia coli expands understanding of the ecology and
523 speciation of the model bacterial species. *Proc Natl Acad Sci U S A* 108:7200-5.
- 524 14. Kaas RS, Friis C, Ussery DW, Aarestrup FM. 2012. Estimating variation within the genes and
525 inferring the phylogeny of 186 sequenced diverse Escherichia coli genomes. *BMC Genomics*
526 13:577.

- 527 15. Escobar-Paramo P, Grenet K, Le Menac'h A, Rode L, Salgado E, Amorin C, Gouriou S, Picard B,
528 Rahimy MC, Andremont A, Denamur E, Ruimy R. 2004. Large-scale population structure of
529 human commensal *Escherichia coli* isolates. *Appl Environ Microbiol* 70:5698-700.
- 530 16. Tarr PI, Gordon CA, Chandler WL. 2005. Shiga-toxin-producing *Escherichia coli* and
531 haemolytic uraemic syndrome. *Lancet* 365:1073-86.
- 532 17. Beyrouthy R, Robin F, Delmas J, Gibold L, Dalmaso G, Dabboussi F, Hamze M, Bonnet R.
533 2014. IS1R-mediated plasticity of IncL/M plasmids leads to the insertion of bla OXA-48 into
534 the *Escherichia coli* Chromosome. *Antimicrob Agents Chemother* 58:3785-90.
- 535 18. Poirel L, Naas T, Nordmann P. 2010. Diversity, epidemiology, and genetics of class D beta-
536 lactamases. *Antimicrob Agents Chemother* 54:24-38.
- 537 19. Wu H, Fang Y, Yu J, Zhang Z. 2014. The quest for a unified view of bacterial land colonization.
538 *ISME J* doi:10.1038/ismej.2013.247.
- 539 20. Tenaillon O, Skurnik D, Picard B, Denamur E. 2010. The population genetics of commensal
540 *Escherichia coli*. *Nat Rev Microbiol* 8:207-17.
- 541 21. Dobrindt U. 2005. (Patho-)Genomics of *Escherichia coli*. *Int J Med Microbiol* 295:357-71.
- 542 22. Bobay LM, Traverse CC, Ochman H. 2015. Impermanence of bacterial clones. *Proc Natl Acad Sci U S A* 112:8893-900.
- 544 23. Touchon M, Hoede C, Tenaillon O, Barbe V, Baeriswyl S, Bidet P, Bingen E, Bonacorsi S,
545 Bouchier C, Bouvet O, Calteau A, Chiapello H, Clermont O, Cruveiller S, Danchin A, Diard M,
546 Dossat C, Karoui ME, Frapy E, Garry L, Ghigo JM, Gilles AM, Johnson J, Le Bouguenec C,
547 Lescat M, Mangenot S, Martinez-Jehanne V, Matic I, Nassif X, Oztas S, Petit MA, Pichon C,
548 Rouy Z, Ruf CS, Schneider D, Turret J, Vacherie B, Vallenet D, Medigue C, Rocha EP,
549 Denamur E. 2009. Organised genome dynamics in the *Escherichia coli* species results in
550 highly diverse adaptive paths. *PLoS Genet* 5:e1000344.
- 551 24. Everitt RG, Didelot X, Batty EM, Miller RR, Knox K, Young BC, Bowden R, Auton A, Votintseva
552 A, Larner-Svensson H, Charlesworth J, Golubchik T, Ip CL, Godwin H, Fung R, Peto TE, Walker
553 AS, Crook DW, Wilson DJ. 2014. Mobile elements drive recombination hotspots in the core
554 genome of *Staphylococcus aureus*. *Nat Commun* 5:3956.
- 555 25. Castillo-Ramirez S, Corander J, Marttinen P, Aldeljawi M, Hanage WP, Westh H, Boye K,
556 Gulay Z, Bentley SD, Parkhill J, Holden MT, Feil EJ. 2012. Phylogeographic variation in
557 recombination rates within a global clone of methicillin-resistant *Staphylococcus aureus*.
558 *Genome Biol* 13:R126.
- 559 26. McNally A, Cheng L, Harris SR, Corander J. 2013. The evolutionary path to extraintestinal
560 pathogenic, drug-resistant *Escherichia coli* is marked by drastic reduction in detectable
561 recombination within the core genome. *Genome Biol Evol* 5:699-710.
- 562 27. Kang Y, Gu C, Yuan L, Wang Y, Zhu Y, Li X, Luo Q, Xiao J, Jiang D, Qian M, Ahmed Khan A,
563 Chen F, Zhang Z, Yu J. 2014. Flexibility and symmetry of prokaryotic genome rearrangement
564 reveal lineage-associated core-gene-defined genome organizational frameworks. *MBio*
565 5:e01867.
- 566 28. Soucy SM, Huang J, Gogarten JP. 2015. Horizontal gene transfer: building the web of life. *Nat*
567 *Rev Genet* 16:472-82.
- 568 29. Ravenhall M, Skunca N, Lassalle F, Dessimoz C. 2015. Inferring horizontal gene transfer. *PLoS*
569 *Comput Biol* 11:e1004095.
- 570 30. Nielsen KM, Bohn T, Townsend JP. 2014. Detecting rare gene transfer events in bacterial
571 populations. *Front Microbiol* 4:415.
- 572 31. Domazet-Loso M, Domazet-Loso T. 2016. gmos: Rapid Detection of Genome Mosaicism over
573 Short Evolutionary Distances. *PLoS One* 11:e0166602.
- 574 32. Chen L, Mathema B, Pitout JD, DeLeo FR, Kreiswirth BN. 2014. Epidemic *Klebsiella*
575 *pneumoniae* ST258 is a hybrid strain. *MBio* 5:e01355-14.

- 576 33. Boritsch EC, Khanna V, Pawlik A, Honore N, Navas VH, Ma L, Bouchier C, Seemann T, Supply P,
577 Stinear TP, Brosch R. 2016. Key experimental evidence of chromosomal DNA transfer among
578 selected tuberculosis-causing mycobacteria. *Proc Natl Acad Sci U S A* 113:9876-81.
- 579 34. Bukh AS, Schonheyder HC, Emmersen JM, Sogaard M, Bastholm S, Roslev P. 2009.
580 *Escherichia coli* phylogenetic groups are associated with site of infection and level of
581 antibiotic resistance in community-acquired bacteraemia: a 10 year population-based study
582 in Denmark. *J Antimicrob Chemother* 64:163-8.
- 583 35. Vading M, Kabir MH, Kalin M, Iversen A, Wiklund S, Naucler P, Giske CG. 2016. Frequent
584 acquisition of low-virulence strains of ESBL-producing *Escherichia coli* in travellers. *J*
585 *Antimicrob Chemother* 71:3548-3555.
- 586 36. Huang CL, Pu PH, Huang HJ, Sung HM, Liaw HJ, Chen YM, Chen CM, Huang MB, Osada N,
587 Gojobori T, Pai TW, Chen YT, Hwang CC, Chiang TY. 2015. Ecological genomics in
588 *Xanthomonas*: the nature of genetic adaptation with homologous recombination and host
589 shifts. *BMC Genomics* 16:188.
- 590 37. Shapiro BJ, Friedman J, Cordero OX, Preheim SP, Timberlake SC, Szabo G, Polz MF, Alm EJ.
591 2012. Population genomics of early events in the ecological differentiation of bacteria.
592 *Science* 336:48-51.
- 593 38. Zwick ME, Joseph SJ, Didelot X, Chen PE, Bishop-Lilly KA, Stewart AC, Willner K, Nolan N,
594 Lentz S, Thomason MK, Sozhamannan S, Mateczun AJ, Du L, Read TD. 2012. Genomic
595 characterization of the *Bacillus cereus* sensu lato species: backdrop to the evolution of
596 *Bacillus anthracis*. *Genome Res* 22:1512-24.
- 597 39. Cadillo-Quiroz H, Didelot X, Held NL, Herrera A, Darling A, Reno ML, Krause DJ, Whitaker RJ.
598 2012. Patterns of gene flow define species of thermophilic Archaea. *PLoS Biol* 10:e1001265.
- 599 40. Ellison CE, Hall C, Kowbel D, Welch J, Brem RB, Glass NL, Taylor JW. 2011. Population
600 genomics and local adaptation in wild isolates of a model microbial eukaryote. *Proc Natl*
601 *Acad Sci U S A* 108:2831-6.
- 602 41. Day M, Doumith M, Jenkins C, Dallman TJ, Hopkins KL, Elson R, Godbole G, Woodford N.
603 2017. Antimicrobial resistance in Shiga toxin-producing *Escherichia coli* serogroups O157 and
604 O26 isolated from human cases of diarrhoeal disease in England, 2015. *J Antimicrob*
605 *Chemother* 72:145-152.
- 606 42. De Queiroz K. 2007. Species concepts and species delimitation. *Syst Biol* 56:879-86.
- 607 43. Krause DJ, Whitaker RJ. 2015. Inferring Speciation Processes from Patterns of Natural
608 Variation in Microbial Genomes. *Syst Biol* 64:926-35.
- 609 44. Doolittle WF, Papke RT. 2006. Genomics and the bacterial species problem. *Genome Biol*
610 7:116.
- 611 45. Nosil P, Feder JL. 2012. Genomic divergence during speciation: causes and consequences.
612 *Philos Trans R Soc Lond B Biol Sci* 367:332-42.
- 613 46. Samson JE, Magadan AH, Moineau S. 2015. The CRISPR-Cas Immune System and Genetic
614 Transfers: Reaching an Equilibrium. *Microbiol Spectr* 3:Plas-0034-2014.
- 615 47. Pleska M, Qian L, Okura R, Bergmiller T, Wakamoto Y, Kussell E, Guet CC. 2016. Bacterial
616 Autoimmunity Due to a Restriction-Modification System. *Curr Biol* 26:404-9.
- 617 48. Frye SA, Nilsen M, Tonjum T, Ambur OH. 2013. Dialects of the DNA uptake sequence in
618 *Neisseriaceae*. *PLoS Genet* 9:e1003458.
- 619 49. Cehovin A, Simpson PJ, McDowell MA, Brown DR, Noschese R, Pallett M, Brady J, Baldwin GS,
620 Lea SM, Matthews SJ, Pelicic V. 2013. Specific DNA recognition mediated by a type IV pilin.
621 *Proc Natl Acad Sci U S A* 110:3065-70.
- 622 50. Gray TA, Krywy JA, Harold J, Palumbo MJ, Derbyshire KM. 2013. Distributive conjugal
623 transfer in mycobacteria generates progeny with meiotic-like genome-wide mosaicism,
624 allowing mapping of a mating identity locus. *PLoS Biol* 11:e1001602.
- 625 51. Wang J, Parsons LM, Derbyshire KM. 2003. Unconventional conjugal DNA transfer in
626 mycobacteria. *Nat Genet* 34:80-4.

- 627 52. Krupovic M, Prangishvili D, Hendrix RW, Bamford DH. 2011. Genomics of bacterial and
628 archaeal viruses: dynamics within the prokaryotic virosphere. *Microbiol Mol Biol Rev* 75:610-
629 35.
- 630 53. Griffiths AJF MJ, Suzuki DT, et al. 2000. Prokaryotic transposons, *An Introduction to Genetic*
631 *Analysis* 7th edition. :W. H. Freeman, New York.
- 632 54. Mortimer TD, Pepperell CS. 2014. Genomic signatures of distributive conjugal transfer
633 among mycobacteria. *Genome Biol Evol* 6:2489-500.
- 634 55. Li R, Yu C, Li Y, Lam TW, Yiu SM, Kristiansen K, Wang J. 2009. SOAP2: an improved ultrafast
635 tool for short read alignment. *Bioinformatics* 25:1966-7.
- 636 56. Zhao Y, Wu J, Yang J, Sun S, Xiao J, Yu J. 2012. PGAP: pan-genomes analysis pipeline.
637 *Bioinformatics* 28:416-8.
- 638 57. Zhang H, Gao S, Lercher MJ, Hu S, Chen WH. 2012. EvoView, an online tool for visualizing,
639 annotating and managing phylogenetic trees. *Nucleic Acids Res* 40:W569-72.
- 640 58. Clermont O, Bonacorsi S, Bingen E. 2000. Rapid and simple determination of the *Escherichia*
641 *coli* phylogenetic group. *Appl Environ Microbiol* 66:4555-8.
- 642 59. Si G, Yang W, Bi S, Luo C, Ouyang Q. 2012. A parallel diffusion-based microfluidic device for
643 bacterial chemotaxis analysis. *Lab Chip* 12:1389-94.

644

645 TABLES AND FIGURES LEGENDS

646 **Table 1. Recombination rate inferred based on homoplasy**

	No. str.	Core genome (Mb)	No. of genes	Total polymorphism sites	Homoplasy sites	<i>r/m</i>	Source
<i>E. coli</i>	111	1.05	1,095	12,185	4,076	0.503	This study
Vig	56	1.05	1,095	3,293	1,603	0.949	This study
Slu	31	1.05	1,095	3,578	1,528	0.745	This study
<i>E. coli</i>	19	3.20	3,022	242,688	116,288	0.92	(22)

647

648 **Figure 1. Phylogenetic and geographic diversity of strains used for this study.**

649 (A) Maximal-linkage tree based on concatenated sequences of seven MLST loci of strains from
650 animal/human host and environment. Strains labelled with blue circle are selected for genome
651 sequencing. (B) Geographic and host distributions of all host-related strains.

652 **Figure 2. Maximum-likelihood phylogeny of *E. coli* strains.**

653 The ML phylogeny based on core genome polymorphisms which structures host-related strains into
654 two major clades and a few minor clades. The constitution of the two major clades is generally
655 congruent with traditional phylotype classification with a few exceptions. Each *E. coli* strain is labelled

656 as phylotype (coloured branch), origin (coloured inner stripe), pathotype (coloured outer stripe) and
657 clade (Slu, the sluggish clade; Vig, the vigorous clade).

658 **Figure 3. Physiological features of the Vig and Slu clades.**

659 Physiological experiments for the Vig and Slu strains show their distinct features that explain their
660 ecology distinctions. Asterisks indicate significant difference (*, $p < 0.05$; **, $p < 0.01$) between clades
661 on the basis of *t*-test. (A) A growth rate test under variable temperatures around the normal range. (B)
662 Survival rate test after heat stress (50°C, 20 and 30 min). (C) Mobility test in response to chemo-
663 attraction of amino acids. (D) Relative growth rate in medium of various carbonates to glucose.

664 **Figure 4. Genome content of strains in Vig and Slu.**

665 (A) Genome similarity among strains. The heatmap shows pairwise genome similarity of all strains.
666 Colours indicate scale of genome similarity, as shown in the left legend. Left to the heatmap shows
667 the neighbour-joining tree of the strains according to their genome similarities. The Vig and Slu strains
668 are indicated with coloured branches. (B) Average number of characteristic genes in each category of
669 metabolism with diverse host diet, which are generally determined by clade affiliation rather than host
670 diet. (C) Average number of virulent genes in commensal or pathogenic strains, which is influenced
671 by both clade affiliation and pathogenicity. (D) Average number of antibiotic-resistant and Shiga-like
672 toxin genes shows that the two types of virulent genes are clearly separated into different clades.

673 **Figure 5. Sequence variation in shared genes of the Vig and Slu clades.**

674 (A) Maximal-likelihood tree of concatenated genes in metabolism of amino acid (left) and
675 carbohydrates (right) shared by all Vig and Slu strains. (B) Distribution of paired protein sequence
676 distance between homologs for each orthologue. The red, blue, and purple lines indicate Vig within-
677 clade pairs (Vig–Vig), Slu within-clade pairs (Slu–Slu), and between-clade pairs (Vig–Slu),
678 respectively. (C) GC content of whole genome (left), core genes (central), and dispensable genes
679 (right) of the two clades. Asterisks indicate significant difference (*, $p < 0.05$; **, $p < 0.01$) between
680 clades on the basis of *t*-test.

681 **Figure 6. Recombination in the Vig and Slu clades.**

682 (A) Homologous recombination in the core genome of 111 *E. coli* strains inferred based on
683 BratNextGen. Seven clusters are a priori that is defined according to the PSA topology. Color bars
684 indicate recombined regions with various PSA cluster origins. Thick colored regions indicate
685 recombination hotspots. Gray bars indicate gaps in alignment. (B) Distribution of total length of
686 recombination between pairs of strains. The red, blue, and purple lines indicate Vig within-clade pairs
687 (Vig–Vig), Slu within-clade pairs (Slu–Slu), and between-clade pairs (Vig–Slu), respectively. Insets
688 are distributions of recombination length between commensal/pathogenic pairs (dotted line) compared
689 with intra-clade pairs (solid line) for Vig (upper) and Slu (lower). (C) Recombination frequency against
690 genome distance for Vig (left) and Slu (right) strains. Red, blue, and green dots stand for pairing with
691 Vig, Slu, minor clades, respectively, and the purple dots are Vig–Slu pairs. (D) Distributions of
692 recombinant fragment length. Recombinant fragments when cross-clade never exceed 5kb, when
693 within-clade recombinant fragments exhibiting right-shift peak overriding the 5kb-limit. The red, blue,
694 and purple lines indicate fragments of within-clade and inter-clade fragments, respectively.

695

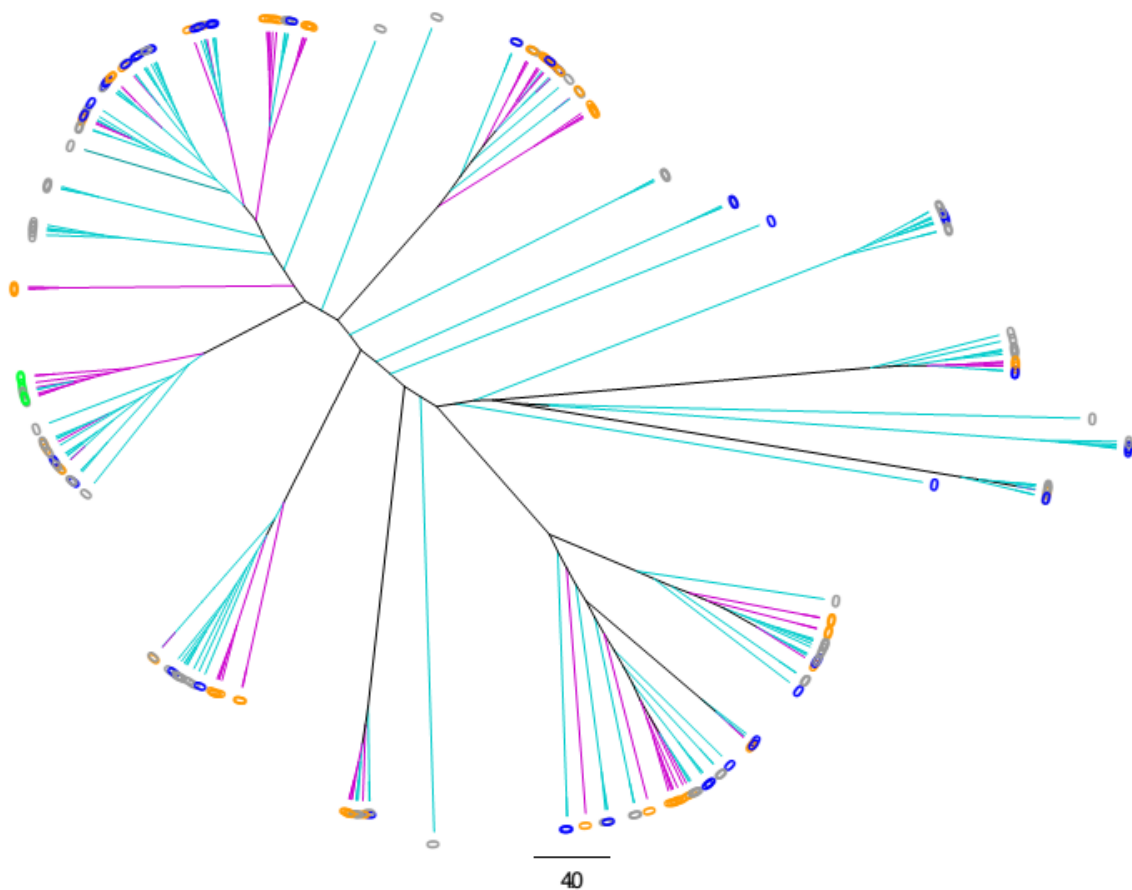
696 **ADDITIONAL FILES**

697 **Table S1. *E.coli* strains used in this study**

698 **Table S2. Characteristic genes enriched in clade Vig and Slu**

699 **Table S3. List of LAV-containing genes of clade Vig and Slu**

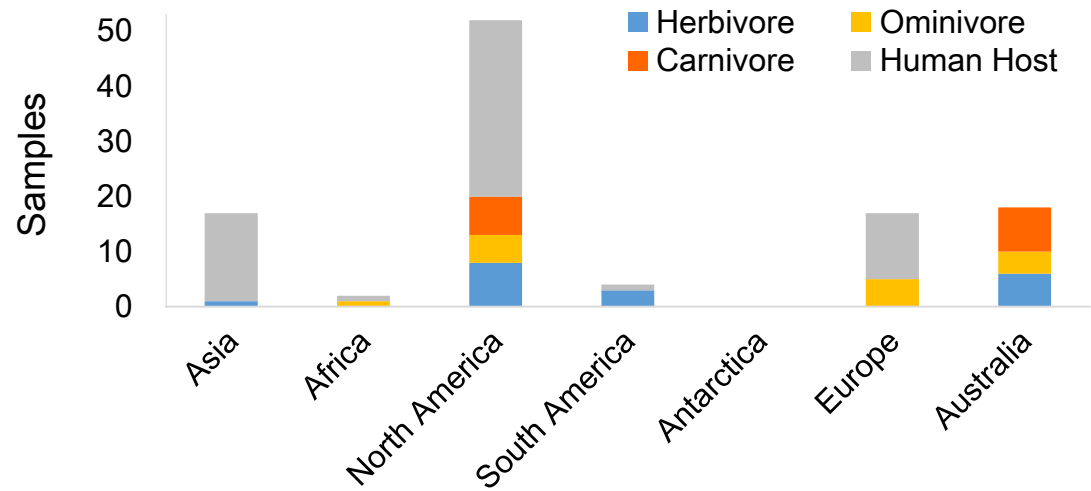
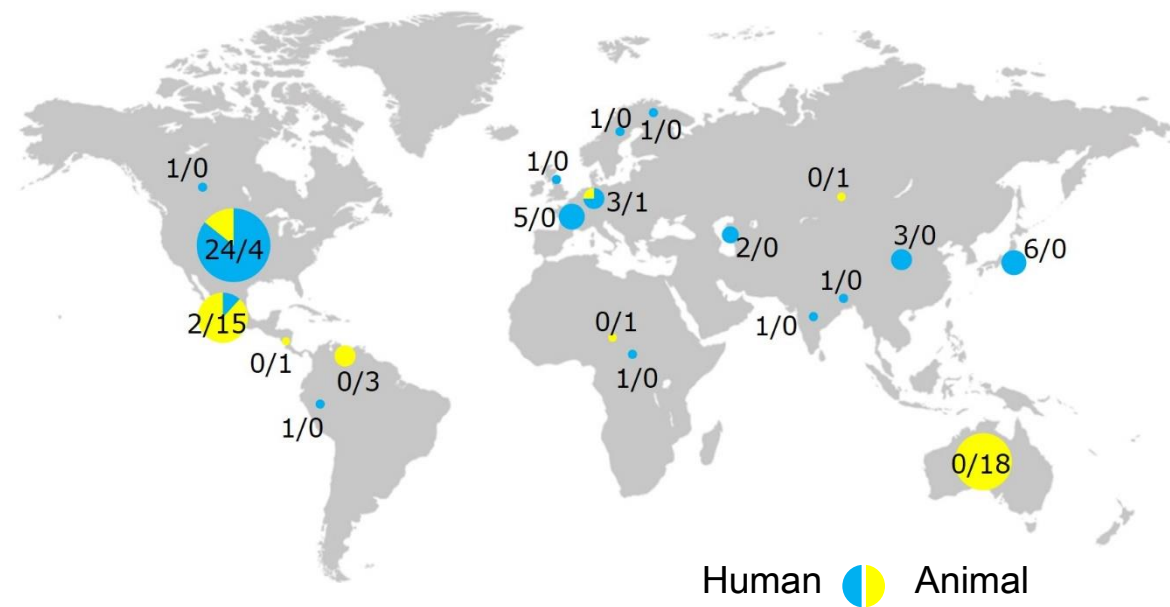
A

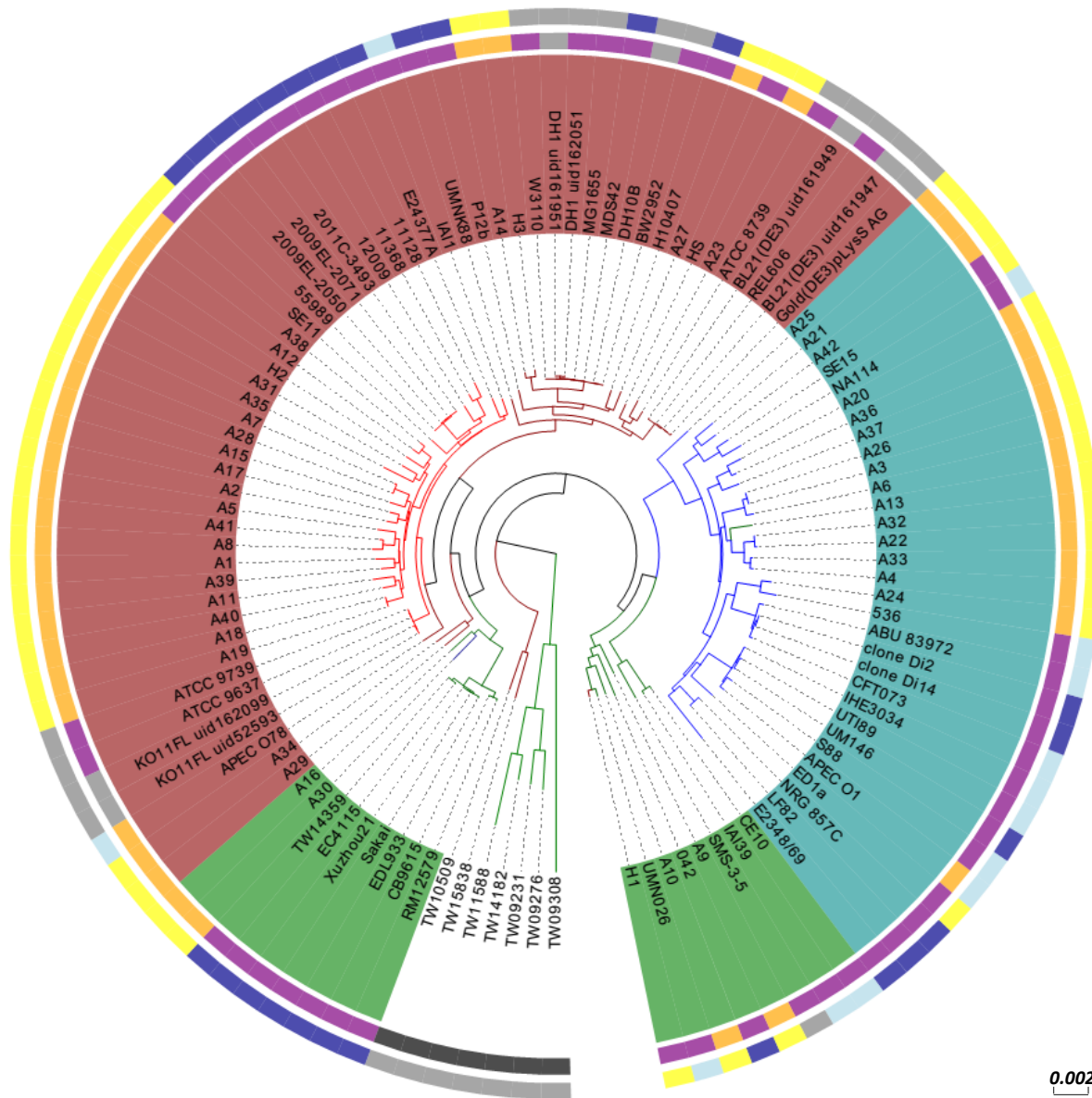


Samples

- Our collection (149 str.)
 - Self-sequenced (45 str.)
- Public collection (66 str.)
 - NCBI deposited (59 str.)
 - Drafts of environmental strains (7 str.)

B





Clade(Leaf)

- Vig
- Slu
- Minor Clades

Phylotype (Branch)

- A
- B1
- B2
- D

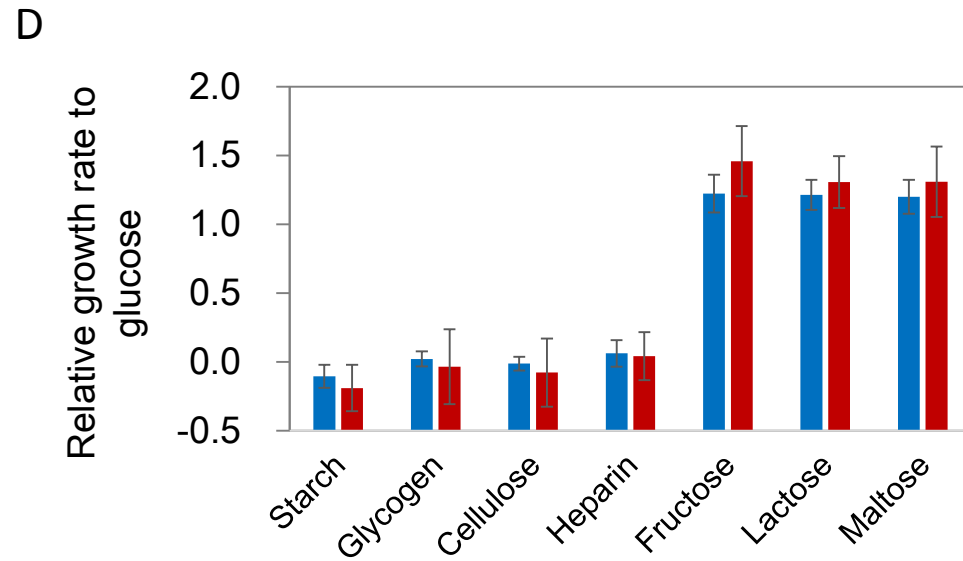
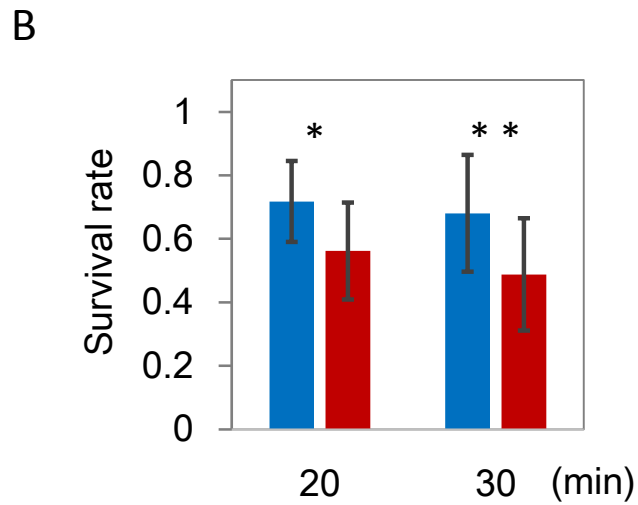
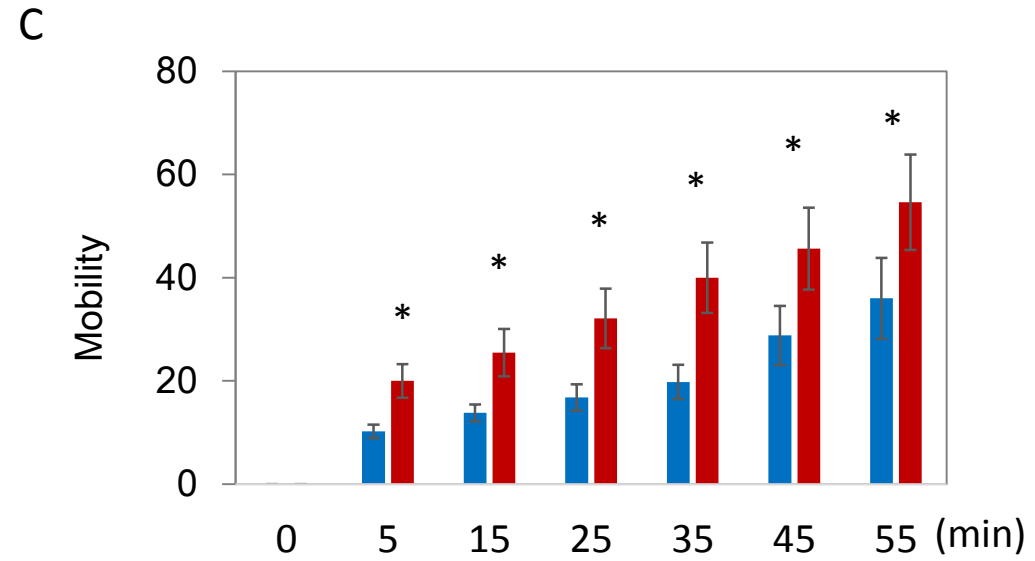
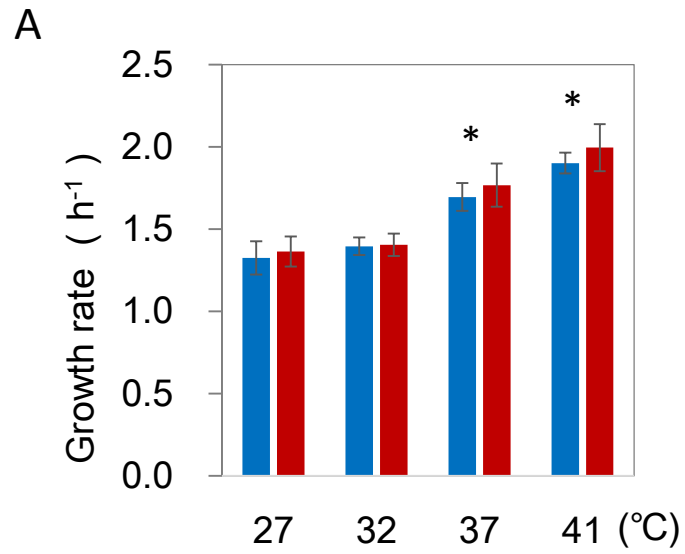
Source/Host (Inner Strip)

- Human
- Animal
- Environment
- Engineered strains

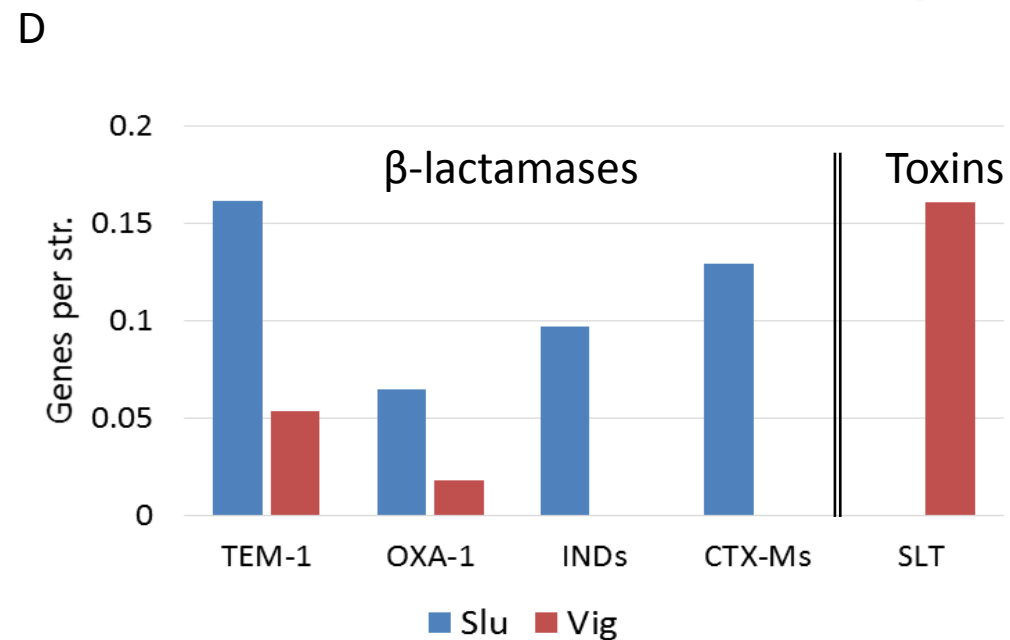
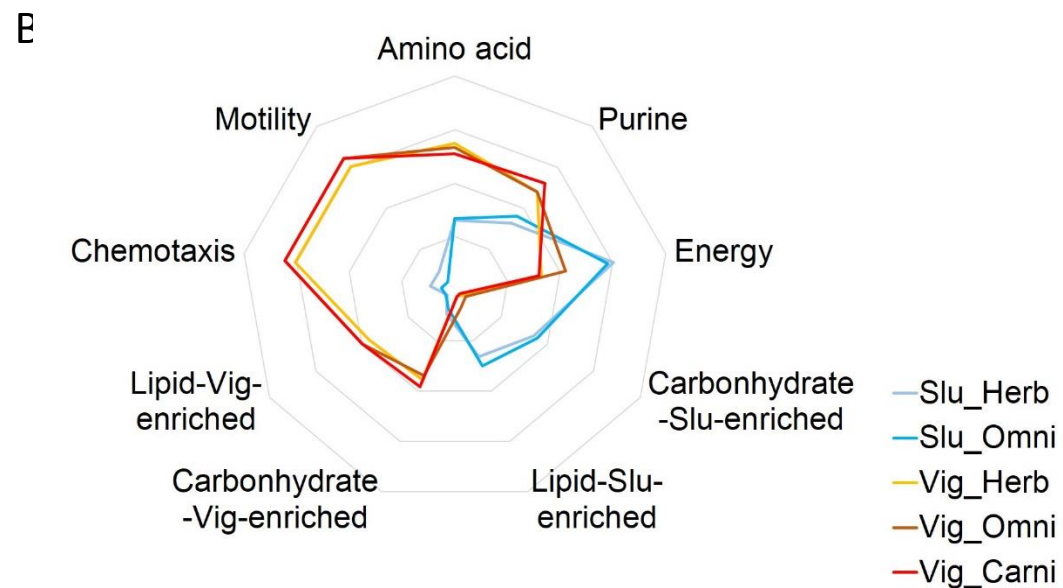
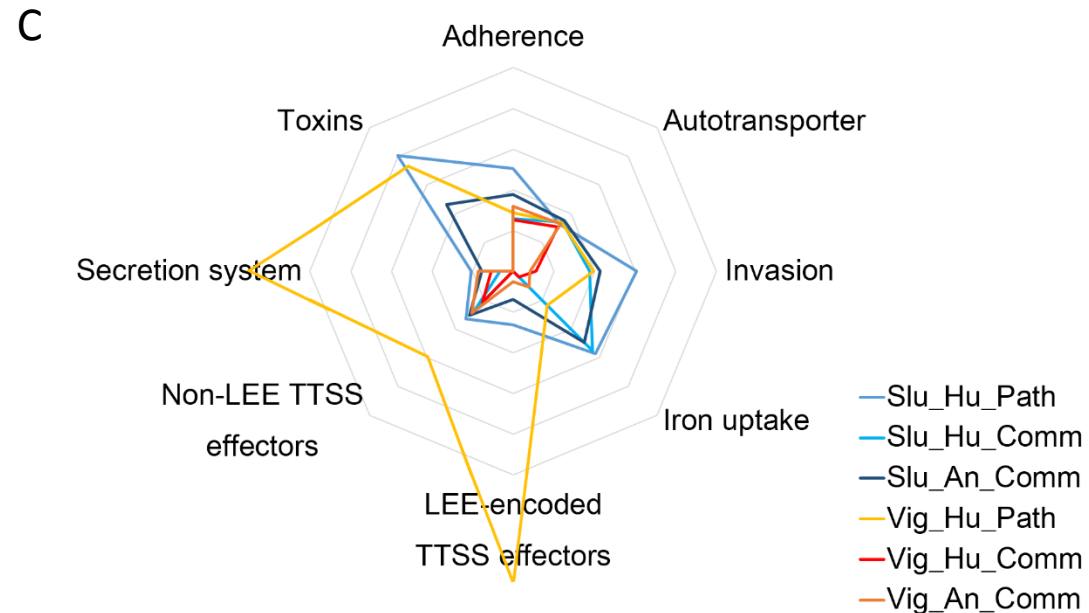
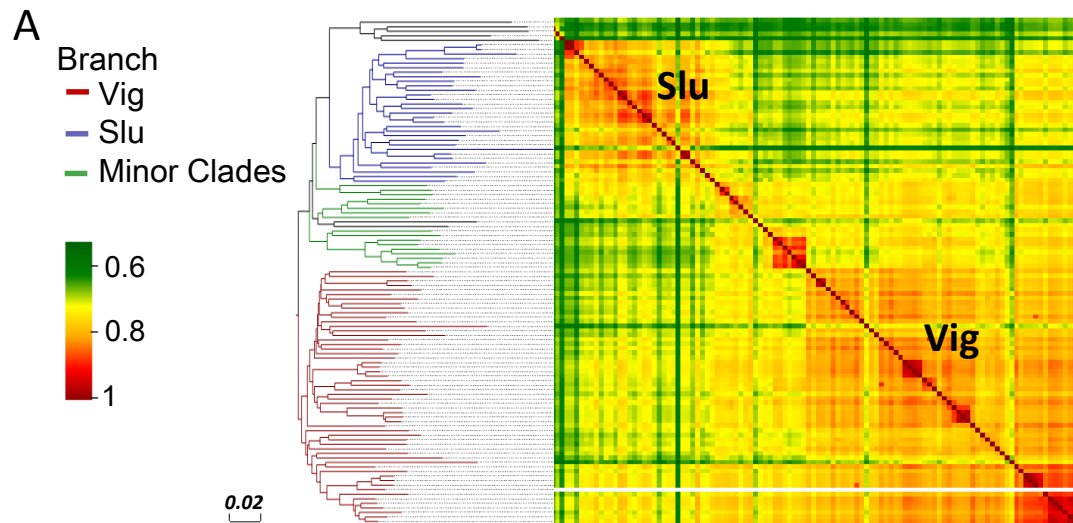
Pathogenicity (Outer Strip)

- Commensal
- IPEC
- ExPEC
- N/A

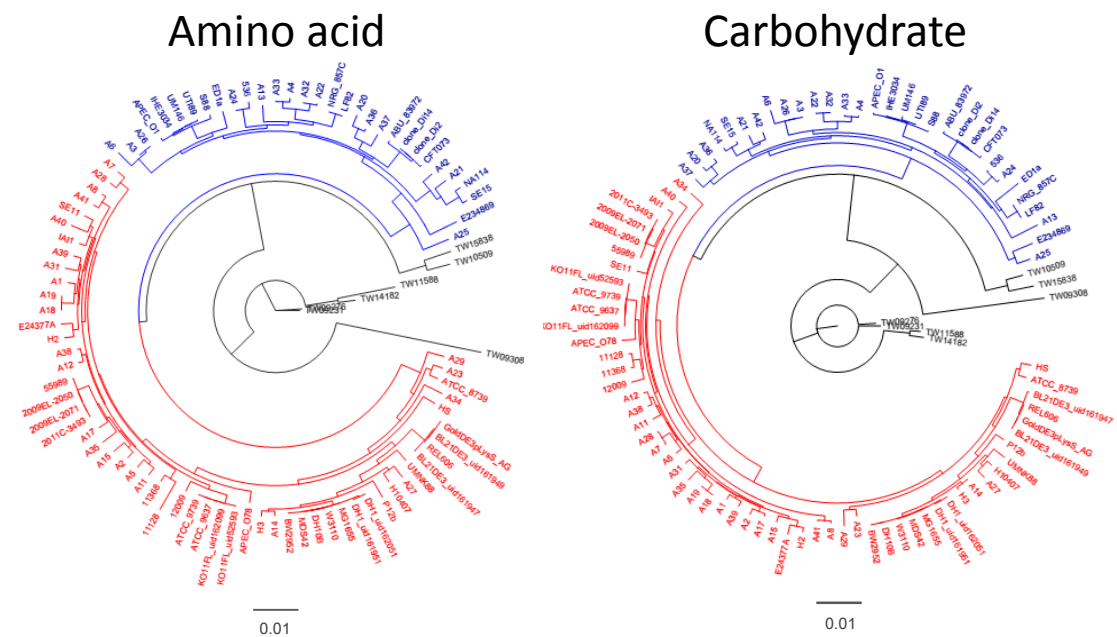
0.002



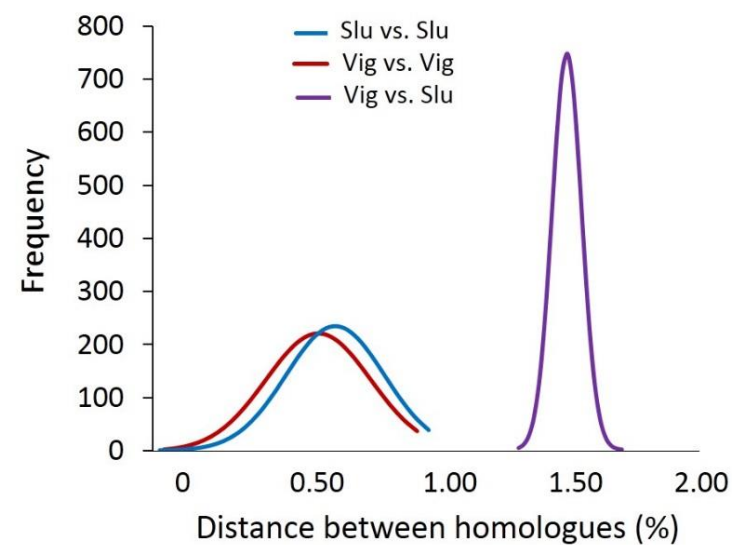
■ Slu ■ Vig



A



B



C

