

1 Gene-gene relationships in an *Escherichia coli*
2 accessory genome are linked to function and
3 mobility

4 Rebecca J. Hall^{1,2}, Fiona J. Whelan¹, Elizabeth A. Cummins^{1,2},
Christopher Connor², Alan McNally², James O. McInerney^{1*}

5 ¹School of Life Sciences, University of Nottingham, Nottingham, NG7 2UH, UK

6 ²Institute of Microbiology and Infection, College of Medical and Dental Sciences,
7 University of Birmingham, Birmingham, B15 2TT, UK

8 Corresponding author: james.mcinerney@nottingham.ac.uk

9 **Abstract**

10 The pangenome contains all genes encoded by a species, with the core genome
11 present in all strains and the accessory genome in only a subset. Coincident gene re-
12 lationships are expected within the accessory genome, where the presence or absence
13 of one gene is influenced by the presence or absence of another. Here, we analysed the
14 accessory genome of an *Escherichia coli* pangenome consisting of 400 genomes from
15 20 sequence types to identify genes that display significant co-occurrence or avoid-
16 ance patterns with one another. We present a complex network of genes that are
17 either found together or that avoid one another more often than would be expected
18 by chance, and show that these relationships vary by lineage. We demonstrate that
19 genes co-occur by function, and that several highly connected gene relationships are
20 linked to mobile genetic elements. We find that genes are more likely to co-occur
21 with, rather than avoid, another gene, suggesting that cooperation is more com-
22 mon than conflict in the accessory genome. This work furthers our understanding
23 of the dynamic nature of prokaryote pangenomes and implicates both function and
24 mobility as drivers of gene relationships.

25 **Data summary**

26 All Supplementary Data files and the Python scripts used in the analyses are avail-
27 able at doi.org/10.17639/nott.7103.

28 **Impact statement**

29 The pangenome of a species encompasses the core genes encoded by all genomes, as
30 well as the accessory genes found in only a subset. Much remains to be understood
31 about the relationships and interactions between accessory genes; in particular, what
32 drives pairs of genes to appear together in the same genome, or what prevents them
33 from being in the same genome together, more often than expected by chance. How
34 these co-occurrence and avoidance relationships develop, and what effect they have
35 on the dynamics and evolution of the pangenome as a whole, is largely unknown.
36 Here, we present a springboard for understanding prokaryote pangenome evolution
37 by uncovering significant gene relationships in a model *Escherichia coli* pangenome.
38 We identify mobile genetic elements and the sharing of common function as possible
39 driving forces behind the co-occurrence of accessory genes. Furthermore, this work
40 offers an extensive dataset from which gene relationships could be identified for any
41 gene of interest in this *E. coli* accessory genome, providing a rich resource for the
42 community.

43 **Introduction**

44 *Escherichia coli* is one of the most widely used and studied bacterial species in
45 microbiology. Recent efforts have resulted in a cataloguing of the essential genes
46 in this species using transposon-directed insertion site sequencing (TraDIS) (1),
47 thereby defining which genes are indispensable and which are not. It is arguable
48 that the accessory genes found in the *E. coli* pangenome are not essential for the
49 survival of the species, making it somewhat of a curiosity that such a large set of
50 accessory genes is maintained. A separate study of a dataset of 53 *E. coli* genomes
51 identified more than 3,000 metabolic innovations that all arose as a consequence of
52 the acquisition via horizontal gene transfer (HGT) of a single piece of DNA less than
53 30 kb in length, and that 10.6% of innovations were facilitated by earlier acquisitions
54 (2). This suggests therefore that HGT has the ability to bring together sets of genes
55 that, when combined, provide benefits over and above the benefits that the genes
56 could confer on their own. Indeed, it is likely that there are situations where the
57 genes on their own might be deleterious, but in combination they confer a fitness
58 advantage (3), thereby contributing to the maintenance of a large accessory genome.

59 Analyses of genome content shows that *E. coli* pangenome evolution is driven
60 by differential gain and loss of accessory genes, including plasmids, phage, and
61 pathogenicity islands (4; 5; 6). What is not yet clear is the underlying structure
62 of the pangenome and which genes are key influencers of the presence or absence
63 of other genes in a given genome. Plasmid mobility, for instance, is anticipated to
64 be an important agent of pangenome structuring (7). Plasmids engage a diverse
65 range of proteins for the purpose of plasmid partitioning and maintenance, but the
66 presence and absence of these protein encoding genes in accessory genomes has not
67 been established concretely.

68 Variation in overall genome sequence has the potential to influence which other genes
69 can or cannot be present in a genome. A gene might be beneficial to one strain of
70 a species, but deleterious in a different strain. Sets of genes encoding proteins that
71 together form an essential biosynthetic pathway, for example, are expected to co-
72 occur in the same genome, likewise those that form multi-protein complexes. The
73 presence of a gene or set of genes can also exclude, or ‘avoid’, others from the same
74 genome. Competitive exclusion is exemplified in *Salinispora* species, where only one
75 of two iron-chelating siderophore gene clusters is ever found in a given strain, despite
76 evidence of frequent HGT (8). An understanding of gene-gene dependencies and
77 influences can therefore provide a real insight into how pangenomes originate and
78 are maintained. We can address at least two questions when we look at gene-gene
79 co-occurrence and avoidance patterns. We can ask whether the pangenome has been
80 structured by random genetic drift or by natural selection, and we can ask whether
81 gene co-occurrence has been more important than avoidance in a pangenome.

82 In practical terms, the large *E. coli* accessory genome can serve as a testbed for cali-
83 brating how much gene co-occurrence and avoidance can teach us about pangenome
84 origins and evolution. A considerable amount of diversity is found across all known
85 sequence types (STs) (4; 9; 10; 11; 12), while a large body of experimental validation
86 of gene function has had the effect of reducing the number of unknown open reading
87 frames (ORFs) in the pangenome. It is anticipated that there may be collections
88 of lineage-specific coincident gene-gene relationships, but the process driving these
89 relationships is not yet clear.

90 Here, we interrogate *E. coli* pangenome dynamics in order to identify coincident
91 gene-gene relationships (13). In a large sample of the total *E. coli* pangenome,
92 comprising 400 judiciously selected genomes spanning 20 different STs, we have
93 identified significant gene-gene co-occurrences and avoidances. We find connected
94 components that are enriched in genes that share a common function, suggesting
95 that the role that genes play in a system can partially explain their significant co-
96 occurrence. We also find that MGEs extensively influence gene co-occurrence and
97 avoidance, including for genes linked to detoxification and antimicrobial resistance
98 (AMR). We identified more co-occurrence relationships than avoidance relationships,
99 which, according to our data and method of analysis, suggests that gene cooperation
100 is more common than conflict in the pangenome. Our results provide an extensive
101 set of possible empirical experiments that, together with our *in silico* predictions,
102 further our understanding of the complex ecological interactions and dynamics in
103 the *E. coli* pangenome.

104 **Methods**

105 ***E. coli* genomes and pangenome**

106 A set of 400 *E. coli* genomes were downloaded from EnteroBase using a custom
107 Python script (github.com/C-Connor/EnterobaseGenomeAssemblyDownload). The
108 set contained 20 genomes from 20 different STs; ST3, ST10, ST11, ST12, ST14,

109 ST17, ST21, ST28, ST38, ST69, ST73, ST95, ST117, ST127, ST131, ST141, ST144,
110 ST167, ST372, ST648. The STs were chosen to include multiple extraintestinal
111 pathogenic *E. coli* (ExPEC), enterohemorrhagic *E. coli* (EHEC), and commensal
112 lineages. The sequences were annotated using Prokka (14) and a gene presence-
113 absence matrix was generated with Panaroo (15) using the default settings and the
114 -a core flag to generate a core gene alignment.

115 Core gene phylogeny

116 The core gene alignment file produced by Panaroo was trimmed using trimAl (16)
117 with a -gt value of 1. Phylogenetic relationships between the strains were inferred
118 using the core gene set from the pangenome. A core gene phylogeny was constructed
119 from the trimmed alignment using the IQ-Tree software (17) with the GTR+I+G
120 substitution model. Phylogenetic tree visualisation was carried out using the Inter-
121 active Tree of Life v5 (iTOL) (18).

122 Coincident gene identification and visualisation

123 Gene co-occurrences (associations) and avoidance (dissociations) were identified us-
124 ing Coinfinder (13) using a threshold of 0.01 for low abundance filtering and em-
125 ploying the Bonferroni correction to account for multiple testing. Networks were
126 visualised using the Gephi software program (v0.9.2) with the Fructerman-Reingold
127 algorithm used for graph layout. Heatmaps were generated using seaborn (v0.10.1).
128 Hub genes are identified as genes forming a number of gene co-occurrences or avoid-
129 ances that is 1.5 times the interquartile range (IQR), as in (19). Gene names were
130 taken from the Panaroo gene cluster identifications.

131 Prediction of plasmid-encoded genes

132 Genes were determined to be chromosomal or plasmid-encoded by collecting each
133 contig containing the gene across all strains in which it was present, and then as-
134 sessing whether those contigs were predicted to be plasmid- or chromosome-derived
135 sequences using mlplasmids (20).

136 Phosphotransferase system analysis

137 A list of all phosphotransferase system (PTS) genes was obtained from the KEGG
138 database (21) and the *E. coli* gene presence-absence matrix was used to highlight
139 those genes found in the *E. coli* pangenome.

140 Data availability

141 All Python scripts used for data analysis and figure construction are freely avail-
142 able at doi.org/10.17639/nott.7103. The gene presence-absence matrix generated
143 by Panaroo, core gene phylogeny, outputs from Coinfinder, and a list of genomes
144 mapped to their ST are also available as Supplementary Data at the same location.
145 Descriptions of Coinfinder outputs can be found in the original software publication
146 (13).

147 Results

148 The *E. coli* pangenome is highly structured

149 An *E. coli* pangenome, composed of 400 genomes from 20 different STs, was con-
150 structed (Fig. S1). Panaroo was used for this analysis as preliminary investigations
151 found that the clustering produced fewer incidences of false positives in the avoid-
152 ance network than when the gene presence-absence matrix was constructed using
153 Roary (22). This pangenome sample consists of 3,191 core, 120 soft core, 2,935 shell,
154 and 11,665 cloud genes (Fig. S2), replicating previous observations that *E. coli* has
155 an open pangenome (9). Using Coinfinder (13) we identified a gene co-occurrence
156 network that consisted of 8,054 nodes joined by a total of 500,654 edges (Fig. 1a),
157 and an avoidance network consisting of 3,203 nodes joined by 203,503 edges (Fig.
158 1b). The co-occurrence network is therefore larger, both in terms of numbers of
159 nodes and also numbers of interactions, though both networks have similar average
160 numbers of connections per node. Of all gene clusters in the gene presence-absence
161 matrix, 45.0% form at least one co-occurrence pair, and 17.9% at least one avoidance.
162 Of the accessory genes analysed by Coinfinder, 77.8% form at least one co-occurring
163 or avoidant pair.

164 Co-occurring genes share function

165 The co-occurrence network contains 224 connected components, including one large
166 connected component with extensive interactions (Fig. 1a). If indeed co-occurrence
167 is shaped by functional interactions and dependencies, we could expect co-occurrence
168 analyses to pick out known subsystems. Component 150, which is made up of a total
169 of 15 genes, consists in part of twelve genes involved in the type II secretion system
170 (T2SS) (*epsC-H, epsLM, gspK, pppA* and *xcpVW*) that facilitates the translocation
171 of a wide variety of proteins from the periplasm to the exterior of the cell. Similarly,
172 component 50 includes *epsE* (T2SS), and the type IV secretion system (T4SS) genes
173 *vir1,4,8,10,11*. Component 150 has a broader distribution across the STs (n=18)
174 than component 50 (n=6), and is found at a much higher abundance. In fact, for 11
175 of the STs, all 20 genomes contain component 150 (Fig. 1c). These data illustrate
176 that indeed analyses of co-occurrence patterns, as implemented by Coinfinder, can

177 pick out functional dependencies.

178 The constituent genes in components 3, 13, 60, and 149 are outlined in Table 1
179 and are known to function in DNA replication. Of note in these components are
180 the plasmid copy control gene *repB* and the plasmid partitioning protein-encoding
181 gene *parB*, both known to be plasmid-encoded, as well as the chromosome-encoded
182 *dnaJ* that functions in plasmid replication (23; 24). The fact that these functions
183 are found across four separate co-occurrence components shows that within this set
184 of functions related to mobility, there are co-occurring communities consisting of
185 distinct groups of genes that preferentially co-occur. These components are found
186 in a variety of STs, though in low abundance (Fig. 1c).

187 **Co-occurrence gene hubs are linked to virulence and mobile** 188 **elements**

189 We found considerable variation in the number of significant co-occurrence or avoid-
190 ance relationships for any individual gene in the *E. coli* accessory genome. From the
191 degree distribution in the co-occurrence network we identified a large group of 427
192 highly connected "hub genes" (Fig. 2a), with a high of 896 co-occurrences recorded
193 for an individual gene cluster (group_2839, identified as *dnaC_4* from the non-
194 unique gene name). These hub genes either facilitate or promote the existence of a
195 large number of other genes in any given genome. The most common known func-
196 tions of co-occurrence gene hubs are attack and defence (including toxin-antitoxin
197 systems, Shiga toxin, CRISPR system subunit, and genes encoding hemolysin and
198 proteins involved in detoxification), metabolism, and DNA and RNA processes (in-
199 cluding DNA primase, helicase, and replication proteins, tyrosine recombinases, and
200 proteins involved in DNA repair) (Fig. 2b). Full lists of the number of pairs formed
201 by individual genes are provided as Supplementary information.

202 A notable subset of the hub genes is linked to DNA exchange and MGEs. Some are
203 known plasmid-encoded genes (*toxB*, *hlyACD*, and *ssb* (24; 25)), and the Shiga toxin
204 subunits *stxAB* (7) are encoded on a single phage (26). The tyrosine recombinases
205 *xerCD* function in plasmid segregation, as does *parM*, and there are several genes
206 either prophage-encoded (*recE* (27)) or related to phage functions (*intQ*, *tfaE*).
207 The transposase *tnpA* is also a hub gene, forming a co-occurring pair with 841
208 other genes. These findings indicate a process where hub genes that have a role
209 in lateral mobility of genetic material are specifically co-occurring with genes that
210 confer fitness advantages, on average, when mobilised. This suggests a process of
211 mutual benefit; the mobility enablers co-occurring with the genes most likely to
212 confer positive fitness effects.

213 **The avoidance network is characterised by root excluders**

214 The avoidance network has substantially fewer connected components (n=14) than
215 the co-occurrence network, and many of these components are characterised by the

216 presence of a single "root excluder"; a situation where one gene avoids a gene set
217 that in turn all co-occur with one another (Fig. 1b, Fig. S3). Over half of the
218 components in this network show this pattern, specifically components 5-7, 9-12,
219 and 14 (Table S1). These components are found in at least 19 genomes of all
220 STs (Fig. 1d). None of the root excluders form avoidance hubs (Fig. 2a). As an
221 example, within component 14, the root excluder *dhaR*, a transcriptional regulator of
222 the dihydroxyacetone kinase operon (28) avoids 11 genes involved in the production
223 of T2SS proteins (*gspHK*, *epsEFLM*, *pulDG*, *outC*, *xcpVW*), amongst other known
224 and hypothetical genes. The plasmid-associated nature of some T2SS genes (24)
225 suggests an active process of dissociation or avoidance as a result of their mobility.

226 Avoidance hub genes are lineage-specific

227 The avoidance network is smaller than the co-occurrence network with fewer hub
228 genes, though high degree nodes can be identified with the highest showing 749
229 avoidances for a single gene cluster (*atoA*, *atoD*, *zraR_2_spo0F*, and *zraS_2*) (Fig.
230 2a). The *ato* genes encode proteins involved in the degradation and transport of
231 short-chain fatty acids (29), and the *zra* genes encode a two-component system
232 linked to antimicrobial tolerance in *E. coli* (30). The most enriched function of the
233 avoidance hub genes is in metabolism, and when grouped by function, regulation is
234 the sole category where the number of avoidance hubs is greater than the number
235 of co-occurrence hubs (Fig. 2b).

236 Some genes avoid considerably more genes than they co-occur with. For instance,
237 the putative diguanylate cyclase *ycdT*, which catalyses the production of cyclic di-
238 3',5'-guanylate and is reportedly under positive selection in uropathogenic *E. coli*
239 (12), significantly co-occurs only with *pgaA-D*. The proteins encoded by *pgaA-D*
240 are involved in biofilm formation. In contrast, *ycdT* avoids 83 other gene clusters
241 including the CRISPR system Cascade subunits *casACDE*. All 400 genomes have at
242 least one of either *ycdT* or the *casACDE* genes, but they are found in isolation in 298
243 genomes. The only STs where this *ycdT-casACDE* avoidance pattern is not observed
244 in any of the genomes are ST3, ST17, ST11, ST38, and ST128, and therefore does
245 not demonstrate a clear phylogenetic split (Fig. S1). We might speculate that there
246 is another element of genome variation that modulates whether *ycdT-casACDE* can
247 or cannot be present in the same genome.

248 The avoidance hub genes are also enriched in functions related to secretion systems
249 (Fig. 2b). We found that between 627 and 661 gene clusters avoid the T2SS-related
250 genes *gspA-M* and the probable bifunctional chitinase/lysozyme that is secreted
251 by the T2SS (31). Of these, one is particularly pertinent; *spiA* (also known as
252 *escC*), encoding a type III secretion system (T3SS) outer membrane secretin (32).
253 Incidences of *spiA* avoiding *gspA-M* gene clusters are found in all 20 STs. These data
254 provide some evidence that secretion systems may have a substantial genome-wide
255 influence.

256 Repeated co-occurrence of resistance genes and transposon 257 genes

258 We have demonstrated that genes that share functions or pathways will commonly
259 co-occur, and that certain highly connected co-occurrences and avoidances involve
260 genes associated with plasmids. We hypothesised that other agents of horizontal
261 transfer could therefore form the centre of co-occurrence hubs as a result of their
262 mobility. The transposase *tnpA*, for example, is a co-occurrence hub gene (n=841
263 gene pairs) (Fig. 2a). This is the only transposon of the 22 within the co-occurrence
264 network that is classified as a hub gene. We found that three connected components
265 (excluding the large component in the centre of Fig. 1a) contain transposons.

266 The first, component 81, is comprised of a collection of genes enriched in functions
267 related to metal detoxification. Eight genes relate to copper (*copABDR*, *pcoCE*)
268 and silver (*silPE*) resistance, and six encode proteins involved in producing a cation
269 efflux system (*cusABCFSR*) (Fig. 3a). A Tn7 transposition protein, *tnsA*, is also
270 present. These systems may be acquired and retained together in order to provide a
271 broader spectrum of metal detoxification. This connected component is only found
272 in four of the 20 STs; ST3 (number of individual genomes = 2), ST10 (n=5), ST117
273 (n=2), and ST127 (n=1) (Fig. 1c). These STs are not close phylogenetically (Fig.
274 S1), indicating that it is not simply a lineage dependent characteristic.

275 Component 53 consists predominantly of the Tn10 transposon proteins *tetCD* that
276 confer tetracycline resistance, as well as *gltS* (a sodium/glutamate symporter) and
277 five hypothetical proteins (Fig. 3b). This component is found in more than half
278 of the genomes in the multidrug resistance-associated ST648 (n=13) and ST167
279 (n=12), and in at least one genome in 16 of the 20 STs (Fig. 1c). The co-occurrence
280 of *gltS* and the hypothetical proteins with the tetracycline resistance genes could
281 suggest they may either be linked in a novel way to AMR, or, alternatively, that they
282 have been transferred with the resistance genes and their co-occurrence is purely as
283 a result of this hitchhiking effect.

284 The remaining connected component that contains a transposon is component 2.
285 This component is comprised of two quasi-cliques and consists, in part, of the T4SS
286 genes *virB1,2,4,6,8-11*, a putative transposon Tn552 DNA invertase (*bin3*), the
287 conjugal transfer protein *traG*, and the tyrosine recombinase *xerD* (Fig. 3c). The
288 presence of these genes throughout both of the quasi-cliques strongly suggests that
289 MGEs are influencing the co-occurrence relationships in this component. This com-
290 ponent is found in at least one genome in 18 different STs (Fig. 1c), providing
291 further evidence for gene mobility. The transposons could influence the transfer
292 of the genes in these components, or they could be hitchhiking between genomes
293 alongside the rest of the component.

294 Hypothetical proteins form a hidden network with a large 295 influence on the *E. coli* accessory genome

296 A theme that emerged within the connected components is the central role in struc-
297 turing the pangenome that is clearly being played by many genes with unknown
298 function. First, 67.5% of the co-occurrence hub genes (n=295) and 17.5% (n=11)
299 of the avoidant hub genes encode hypothetical proteins (Fig. 2c). This was surpris-
300 ing given the extent to which *E. coli* has been studied. While many genes in our
301 dataset have assigned functions, 11,491 ORFs in the gene presence-absence matrix
302 used as input to Coinfinder were not ascribed a function, 64.2% of the total. Second,
303 certain connected components are enriched in genes encoding hypothetical proteins,
304 with several connected components consisting solely of unknown ORFs. Examples
305 of this are co-occurrence components 198 and 89 (Fig. 3d, e), both of which are
306 found in several different STs, though no ST contains both components (Fig. 1c).
307 The entirety of component 89 consists of hypothetical proteins, while all but one of
308 component 198 is a hypothetical protein. It should also be noted that both of these
309 connected components also form a clique; every node in the connected component is
310 connected to every other node. This means that every gene in the component shows
311 a significant co-occurrence pattern with every other gene. It is therefore highly likely
312 that these groups of genes are tightly, functionally linked to one another, though
313 there is no published information related to the function of any of these genes.

314 We observe several connected components where the majority of genes in a compo-
315 nent share a similar function. This makes it tempting to imply a putative role for
316 the genes encoding the remaining hypothetical proteins in those components. Co-
317 occurrence component 90 consists of five genes with known functions; three genes
318 that are found in the dTDP-rhamnose biosynthesis pathway, a rhamnosyltransferase,
319 and a polysialic acid transporter. This leaves four genes in this connected compo-
320 nent that encode hypothetical proteins (Fig. 3f). We could speculate that these
321 hypothetical proteins may therefore function in lipopolysaccharide or capsule for-
322 mation (33). Component 2 consists of 82 gene clusters, including 53 that encode
323 hypothetical proteins and 14 that relate to secretion systems (Fig. 3c). These ob-
324 servations show that there is an important network of genes of unknown function
325 that exert a large influence on the *E. coli* accessory genome, and that our approach
326 can generate meaningful hypotheses to inform investigation of the function of those
327 genes.

328 The drivers of gene-gene relationships are multifaceted

329 We have identified co-occurrence connected components that entirely consist of, or
330 are enriched in, genes that share a common, identifiable role. We have also high-
331 lighted gene relationships that appear to be influenced by the presence of MGEs.
332 If function and mobility were the only drivers behind significant gene-gene relation-
333 ships, it might be expected that genes encoded on the same operon, or that function
334 in a discrete system, will share the same patterns of co-occurrence. To test this,
335 we focused on phosphotransferase system (PTS) genes, which are used to import

336 carbohydrates (34). These systems were chosen because several have been identified
337 in *E. coli*, they are typically multi-component (35), and the presence of one system
338 might logically be linked to the presence or absence of another.

339 We collated all genes encoding PTS components and identified those that form
340 a co-occurring pair with at least one other PTS gene. We found no correlation
341 between whether the encoded proteins are membrane-bound or cytoplasmic and the
342 likelihood that they will co-occur with another PTS gene. Of the 67 PTS genes
343 detailed in the KEGG database (21), 29 were not observed in our gene presence-
344 absence matrix, and a further ten did not manifest a significant co-occurrence with
345 any other PTS gene (Fig. 4a).

346 We found a complex pattern of co-occurrence across the systems. Certain PTS genes
347 do show a consistent pattern for the complete system. For example, *srlABE* all co-
348 occur with the same genes, and *sorABFM* only co-occur with each other (Fig. 4a,
349 b). In contrast, the *manXYZ*, *levDEFG*, *agaBCDF*, and *ulaABC* gene relationships
350 are not system-specific, varying instead by individual gene. Selection pressures on
351 these systems are complex and heterogeneous, and gene co-occurrence relationships
352 are likely driven by multiple factors.

353 Discussion

354 The open pangenome of *E. coli* (9) provides a rich testbed to help understand
355 genome and pangenome dynamics. Factors such as the immediate microenvironment
356 in which a strain lives are known to influence the accessory gene content of any given
357 genome and, consequently, the pangenome (36; 37; 38). Though there are known
358 examples of how the fitness of one gene in a genome is influenced by the presence or
359 absence of other genes, there has been no systematic, large-scale study of how genetic
360 background, in terms of the presence or absence of genes in a genome, influences
361 the fitness effect of an incoming gene (13). Recently, however, it has been shown
362 that the genetic background of a genome has a direct effect on whether or not a
363 gene is essential (39). We have little knowledge of why some lineages encode genes
364 that others do not, and the extent to which the observed encoded genes influence
365 the likelihood of successfully integrating other incoming genes. To begin to unpack
366 this, we have presented a network of coincident gene relationships in a model *E.*
367 *coli* pangenome. We found that nearly half of all gene clusters in the pangenome
368 form a significant pair, and that this relationship is more likely to be one of co-
369 occurrence than avoidance. This indicates an ecological relationship between genes
370 in that cooperation is more likely than conflict; genes co-occur because they share
371 function, whereas direct, antagonistic avoidances are less common. This shows that
372 the *E. coli* pangenome is highly structured.

373 Previously, little was known about structure in prokaryote pangenomes; whether
374 pangenomes arise as a result of drift, or whether they are maintained by selec-
375 tion (40; 41; 42). Recent genus-level analysis of a *Pseudomonas* pangenome found
376 significant co-occurrence of genes that share common function, providing support

377 for selection as a driver of pangenome evolution (19). We build upon this work
378 with the observation that many co-occurring connected components are enriched
379 in genes that share function, broadly defined. We also propose, alongside the role
380 that sharing function plays, that gene mobility could be an additional mechanism
381 behind the formation of these gene-gene relationships. MGEs are a known link to
382 gene essentiality and virulence in *E. coli* (4; 39; 43), and they have been implicated
383 in driving accessory genome differences in a *Listeria monocytogenes* pangenome
384 (44). Here, we have demonstrated that they also influence gene co-occurrences by
385 uncovering hub genes and connected components linked to or encoded on MGEs
386 (24; 25; 7; 45; 46; 47; 48; 49). This includes known phage-encoded virulence factors
387 (26; 50), transposons, and genes involved in conjugation. Together, this progresses
388 current understanding of prokaryote pangenomes by suggesting that they are struc-
389 tured and dynamic, but also further underscores the importance of HGT in driving
390 pangenome evolution (51).

391 Furthermore, we found that these gene relationships are often non-randomly dis-
392 tributed across the pangenome, with some being more frequently observed in spe-
393 cific STs. For example, co-occurrences related to resistance phenotypes are found
394 through the different STs, but are particularly evident in the MDR-associated ST167
395 and ST648 (52; 53). ST-specific differences that confer pathogenicity and resistance
396 in *E. coli* are well-studied (43; 53; 54), but this work provides a new layer of under-
397 standing into how the accessory genome may interact to this end. We suggest that
398 certain gene collections are required by specific lineages and that this may be driven
399 in part by MGEs.

400 Alongside these gene co-occurrences of known function, we have also uncovered a net-
401 work of genes with unknown function that influence the structure of the pangenome
402 through the high number of genes that they co-occur with and avoid. *E. coli* has
403 fewer genes of unknown function than most prokaryotes, although there is evidence
404 that many of the accessory genes in lineages such as ST131 are of hypothetical func-
405 tion (40). The number of coincident gene relationships formed by such genes here
406 highlights the challenges in understanding the global prokaryote pangenome. Given
407 the prevalence of MGEs in the co-occurrence network, it is tempting to conclude
408 that at least some of the high proportion of co-occurrence hub genes identified as
409 encoding hypothetical proteins may be related to mobility.

410 The data presented here support the concept that pangenomes create pangenomes.
411 The diversity of gene content in a cosmopolitan species such as *E. coli* means that the
412 fitness effect of gaining and losing individual genes is not the same for all constituent
413 genomes. We observed the presence or absence of some genes only when other
414 genes are also present or absent. We consistently observed some pairs of genes co-
415 occurring or avoiding repeatedly across the diversity of the group of genomes, and
416 root excluders consistently avoiding certain cliques. We observed unknown ORFs
417 that significantly co-occur in the same genome, even though their distribution in the
418 pangenome is patchy. There is therefore an emerging logic to the *E. coli* pangenome
419 that clearly identifies natural selection to have frequently dominated over genetic
420 drift.

421 Acknowledgements

422 RJH was supported by the BBSRC (BB/N018044/2), awarded to JOM. FJW was
423 funded by a Marie Skłodowska-Curie Individual Fellowship (GA no. 793818). EAC
424 was funded by the Wellcome AAMR DTP and CC by the Wellcome Midas DTP.
425 We thank M.R. Domingo-Sananes and S. Thorpe for their insightful feedback on the
426 work presented here.

427 Conflicts of interest

428 The authors declare that there are no conflicts of interest.

429 References

- 430 [1] Goodall, E. C. *et al.* The essential genome of *Escherichia coli* K-12. *mBio* **9**,
431 02096–17 (2018).
- 432 [2] Pang, T. Y. & Lercher, M. J. Each of 3,323 metabolic innovations in the
433 evolution of *E. coli* arose through the horizontal transfer of a single DNA seg-
434 ment. *Proceedings of the National Academy of Sciences of the United States of*
435 *America* **116**, 187–192 (2019).
- 436 [3] Domingo-Sananes, M. R. & McInerney, J. O. Mechanisms That Shape Microbial
437 Pangenomes. *Trends in Microbiology* (2021).
- 438 [4] Ogura, Y. *et al.* Comparative genomics reveal the mechanism of the parallel
439 evolution of O157 and non-O157 enterohemorrhagic *Escherichia coli*. *Proceed-*
440 *ings of the National Academy of Sciences* **106**, 17939–17944 (2009).
- 441 [5] Hentschel, U. & Hacker, J. Pathogenicity islands: The tip of the iceberg.
442 *Microbes and Infection* **3**, 545–548 (2001).
- 443 [6] Johnson, T. J. *et al.* Separate F-Type Plasmids Have Shaped the Evolution of
444 the H30 Subclone of *Escherichia coli* Sequence Type 131. *mSphere* **1**, 00121–16
445 (2016).
- 446 [7] Nakamura, K. *et al.* Differential dynamics and impacts of prophages and
447 plasmids on the pangenome and virulence factor repertoires of Shiga toxin-
448 producing *Escherichia coli* O145:H28. *Microbial Genomics* **6** (2020).
- 449 [8] Bruns, H. *et al.* Function-related replacement of bacterial siderophore pathways.
450 *ISME Journal* **12**, 320–329 (2018).
- 451 [9] Rasko, D. A. *et al.* The pangenome structure of *Escherichia coli*: Comparative
452 genomic analysis of *E. coli* commensal and pathogenic isolates. *Journal of*
453 *Bacteriology* **190**, 6881–6893 (2008).

- 454 [10] Dobrindt, U. (Patho-)Genomics of Escherichia coli. *International Journal of*
455 *Medical Microbiology* **295**, 357–371 (2005).
- 456 [11] Decano, A. G. & Downing, T. An Escherichia coli ST131 pangenome atlas re-
457 veals population structure and evolution across 4,071 isolates. *Scientific Reports*
458 **9**, 1–13 (2019).
- 459 [12] Chen, S. L. *et al.* Identification of genes subject to positive selection in
460 uropathogenic strains of Escherichia coli: A comparative genomics approach.
461 *Proceedings of the National Academy of Sciences of the United States of Amer-*
462 *ica* **103**, 5977–5982 (2006).
- 463 [13] Whelan, F. J., Rusilowicz, M. & McInerney, J. O. Coinfinder: Detecting sig-
464 nificant associations and dissociations in pangenomes. *Microbial Genomics* **6**,
465 e000338 (2020).
- 466 [14] Seemann, T. Prokka: rapid prokaryotic genome annotation. *Bioinformatics*
467 **30**, 2068–2069 (2014).
- 468 [15] Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the
469 Panaroo pipeline. *Genome biology* **21**, 180 (2020).
- 470 [16] Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for
471 automated alignment trimming in large-scale phylogenetic analyses. *Bioinform-*
472 *atics* **25**, 1972–1973 (2009).
- 473 [17] Nguyen, L.-T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: A
474 Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood
475 Phylogenies. *Molecular Biology and Evolution* **31**, 268–274 (2015).
- 476 [18] Letunic, I. & Bork, P. Interactive Tree of Life (iTOL) v4: Recent updates and
477 new developments. *Nucleic Acids Research* **47**, W256–W259 (2019).
- 478 [19] Whelan, F. J., Hall, R. J. & Mcinerney, J. O. Evidence for selection in a
479 prokaryote pangenome. *bioRxiv* 2020.10.28.359307 (2020).
- 480 [20] Arredondo-Alonso, S. *et al.* mlplasmids: a user-friendly tool to predict plasmid-
481 and chromosome-derived sequences for single species. *Microbial Genomics* **4**,
482 e000224 (2018).
- 483 [21] Kanehisa, M. & Goto, S. KEGG: Kyoto Encyclopedia of Genes and Genomes.
484 Tech. Rep. 1 (2000).
- 485 [22] Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis.
486 *Bioinformatics* **31**, 3691–3693 (2015).
- 487 [23] Kawasaki, Y., Wada, C. & Yura, T. Roles of Escherichia coli heat shock pro-
488 teins DnaK, DnaJ and GrpE in mini-F plasmid replication. *MGG Molecular &*
489 *General Genetics* **220**, 277–282 (1990).
- 490 [24] Makino, K. *et al.* Complete nucleotide sequences of 93-kb and 3.3-kb plasmids
491 of an enterohemorrhagic Escherichia coli O157:H7 derived from Sakai outbreak.
492 *DNA Research* **5**, 1–9 (1998).

- 493 [25] Schmidt, H., Beutin, L. & Karch, H. Molecular analysis of the plasmid-encoded
494 hemolysin of Escherichia coli O157:H7 strain EDL 933. *Infection and Immunity*
495 **63** (1995).
- 496 [26] Waldor, M. K. Bacteriophage biology and bacterial virulence. *Trends in Mi-*
497 *crobiology* **6**, 295–297 (1998).
- 498 [27] Handa, N. & Kobayashi, I. Type III Restriction Is Alleviated by Bacterio-
499 phage (RecE) Homologous Recombination Function but Enhanced by Bacterial
500 (RecBCD) Function. *Journal of Bacteriology* **187**, 7362 – 7373 (2005).
- 501 [28] Bachler, C., Schneider, P., Bahler, P., Lustig, A. & Erni, B. Escherichia coli
502 dihydroxyacetone kinase controls gene expression by binding to transcription
503 factor DhaR. *The EMBO Journal* **24**, 283–293 (2005).
- 504 [29] Jenkins, L. S. & Nunn, W. D. Genetic and molecular characterization of the
505 genes involved in short-chain fatty acid degradation in Escherichia coli: The
506 ato system. *Journal of Bacteriology* **169**, 42–52 (1987).
- 507 [30] Rome, K. *et al.* The Two-Component System ZraPSR Is a Novel ESR that
508 Contributes to Intrinsic Antibiotic Tolerance in Escherichia coli. *Journal of*
509 *Molecular Biology* **430**, 4971–4985 (2018).
- 510 [31] Francetic, O., Belin, D., Badaut, C. & Pugsley, A. P. Expression of the endoge-
511 nous type II secretion pathway in Escherichia coli leads to chitinase secretion.
512 *EMBO Journal* **19**, 6697–6703 (2000).
- 513 [32] Miki, T., Okada, N., Kim, Y., Abe, A. & Danbara, H. DsbA directs efficient
514 expression of outer membrane secretin EscC of the enteropathogenic Escherichia
515 coli type III secretion apparatus. *Microbial Pathogenesis* **44**, 151–158 (2008).
- 516 [33] Silver, R. P., Aaronson, W. & Vann, W. F. The K1 capsular polysaccharide of
517 Escherichia coli. *Reviews of Infectious Diseases* **10 Suppl 2** (1988).
- 518 [34] Postma, P. W. & Lengeler, J. W. Phosphoenolpyruvate: Carbohydrate phos-
519 photransferase system of bacteria (1985).
- 520 [35] Tchieu, J. H., Norris, V., Edwards, J. S. & Saier, M. H. The complete phos-
521 photransferase system in Escherichia coli. *Journal of Molecular Microbiology*
522 *and Biotechnology* **3**, 329–346 (2001).
- 523 [36] Manges, A. R. *et al.* Global extraintestinal pathogenic Escherichia coli (ExPEC)
524 lineages. *Clinical Microbiology Reviews* **32** (2019).
- 525 [37] Sokurenko, E. V. *et al.* Pathogenic adaptation of Escherichia coli by natural
526 variation of the FimH adhesin. *Proceedings of the National Academy of Sciences*
527 *of the United States of America* **95**, 8922–8926 (1998).
- 528 [38] McNally, A. *et al.* Diversification of colonization factors in a multidrug-resistant
529 Escherichia coli lineage evolving under negative frequency- dependent selection.
530 *mBio* **10** (2019).

- 531 [39] Rousset, F. *et al.* The impact of genetic diversity on gene essentiality within
532 the *E. coli* species. *bioRxiv* 2020.05.25.114553 (2020).
- 533 [40] McInerney, J. O., McNally, A. & O’Connell, M. J. Why prokaryotes have
534 pangenomes. *Nature Microbiology* **2**, 1–5 (2017).
- 535 [41] Shapiro, B. J. The population genetics of pangenomes (2017).
- 536 [42] Andreani, N. A., Hesse, E. & Vos, M. Prokaryote genome fluidity is dependent
537 on effective population size. *ISME Journal* **11**, 1719–1721 (2017).
- 538 [43] Cusumano, C. K., Hung, C. S., Chen, S. L. & Hultgren, S. J. Virulence plas-
539 mid harbored by uropathogenic *Escherichia coli* functions in acute stages of
540 pathogenesis. *Infection and Immunity* **78**, 1457–1467 (2010).
- 541 [44] Kuenne, C. *et al.* Reassessment of the *Listeria monocytogenes* pan-genome
542 reveals dynamic integration hotspots and mobile genetic elements as major
543 components of the accessory genome. *BMC Genomics* **14**, 47 (2013).
- 544 [45] Firth, N. & Skurray, R. Characterization of the F plasmid bifunctional conju-
545 gation gene, *traG*. *MGG Molecular & General Genetics* **232**, 145–153 (1992).
- 546 [46] Wallden, K., Rivera-Calzada, A. & Waksman, G. Type IV secretion systems:
547 Versatility and diversity in function (2010).
- 548 [47] Schmidt, H., Henkel, B. & Karch, H. A gene cluster closely related to type
549 II secretion pathway operons of Gram-negative bacteria is located on the large
550 plasmid of enterohemorrhagic *Escherichia coli* O157 strains. *FEMS Microbiology*
551 *Letters* **148**, 265–272 (2006).
- 552 [48] Arciszewska, L. & Sherratt, D. Xer site-specific recombination in vitro. *The*
553 *EMBO Journal* **14**, 2112–2120 (1995).
- 554 [49] Cornet, F., Hallet, B. & Sherratt, D. J. Xer recombination in *Escherichia coli*:
555 Site-specific DNA topoisomerase activity of the XerC and XerD recombinases.
556 *Journal of Biological Chemistry* **272**, 21927–21931 (1997).
- 557 [50] Denamur, E., Clermont, O., Bonacorsi, S. & Gordon, D. The population ge-
558 netics of pathogenic *Escherichia coli*. *Nature Reviews Microbiology* **19**, 37–54
559 (2020).
- 560 [51] Brockhurst, M. A. *et al.* The Ecology and Evolution of Pangenomes (2019).
- 561 [52] Zhang, X. *et al.* First identification of coexistence of bla_{NDM-1} and bla_{CMY-}
562 42 among *Escherichia coli* ST167 clinical isolates. *BMC Microbiology* **13**, 282
563 (2013).
- 564 [53] Schaufler, K. *et al.* Genomic and functional analysis of emerging virulent and
565 multidrug-resistant *Escherichia coli* lineage sequence type 648. *Antimicrobial*
566 *Agents and Chemotherapy* **63** (2019).
- 567 [54] Hibbing, M. E., Dodson, K. W., Kalas, V., Chen, S. L. & Hultgren, S. J. Adap-
568 tation of Arginine Synthesis among Uropathogenic Branches of the *Escherichia*
569 *coli* Phylogeny Reveals Adjustment to the Urinary Tract Habitat. *mBio* **11**
570 (2020).

571 **Figures and tables**

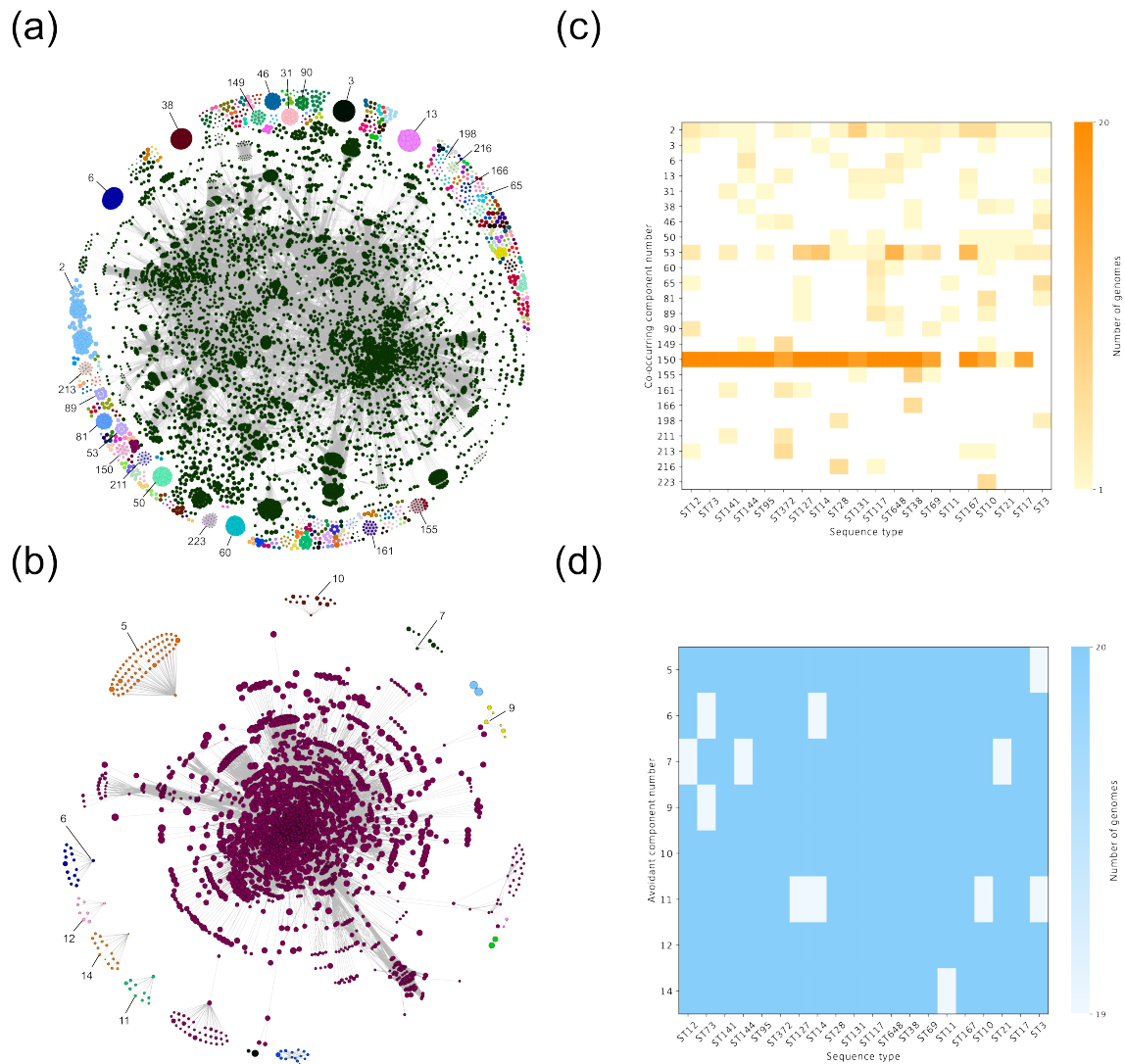


Fig. 1 Gene relationships in the *E. coli* pangenome. (a) An overview of the co-occurrence and (b) avoidance networks, coloured by connected component. Selected connected components are numbered as in the Supplementary Data Files. (c) The number of genomes in which part or all of a co-occurrence and (d) avoidance connected component appears in each ST. An absence of a co-occurrence component in a ST is depicted as colourless.

Table 1 Certain co-occurrence components function in fundamental cellular processes. The known gene clusters relating to DNA and RNA replication, regulation and repair found in the co-occurrence connected components 3, 13, 60, and 149. Genes are named as per the Panaroo gene clusters in the gene presence-absence matrix. Where the gene is identified by ‘group_’, a gene identifier, taken from the Panaroo non-unique gene name, is given in parentheses.

| Component 3 | Component 13 | Component 60 | Component 149 |
|----------------------|----------------------------|----------------------------|----------------------------|
| dnaB_2_dnaB_1_dnaB_3 | polA_2 | dnaB_3_dnaB_2 | group_4304 (<i>recQ</i>) |
| dnaJ_3_dnaJ_1_dnaJ_2 | dnaE_1_dnaE_2_dnaE1 | ssb_1_ssb_5_ssb_4 | srmB_2 |
| ssb_3_ssb_2 | dnaG_1 | group_3992 (<i>topB</i>) | group_296 (<i>rapA</i>) |
| smc__smc_1 | dnaQ_1_dnaQ_2 | | |
| | group_7368 (<i>parB</i>) | | |
| | repB_2 | | |
| | lig | | |
| | smc_2_smc | | |

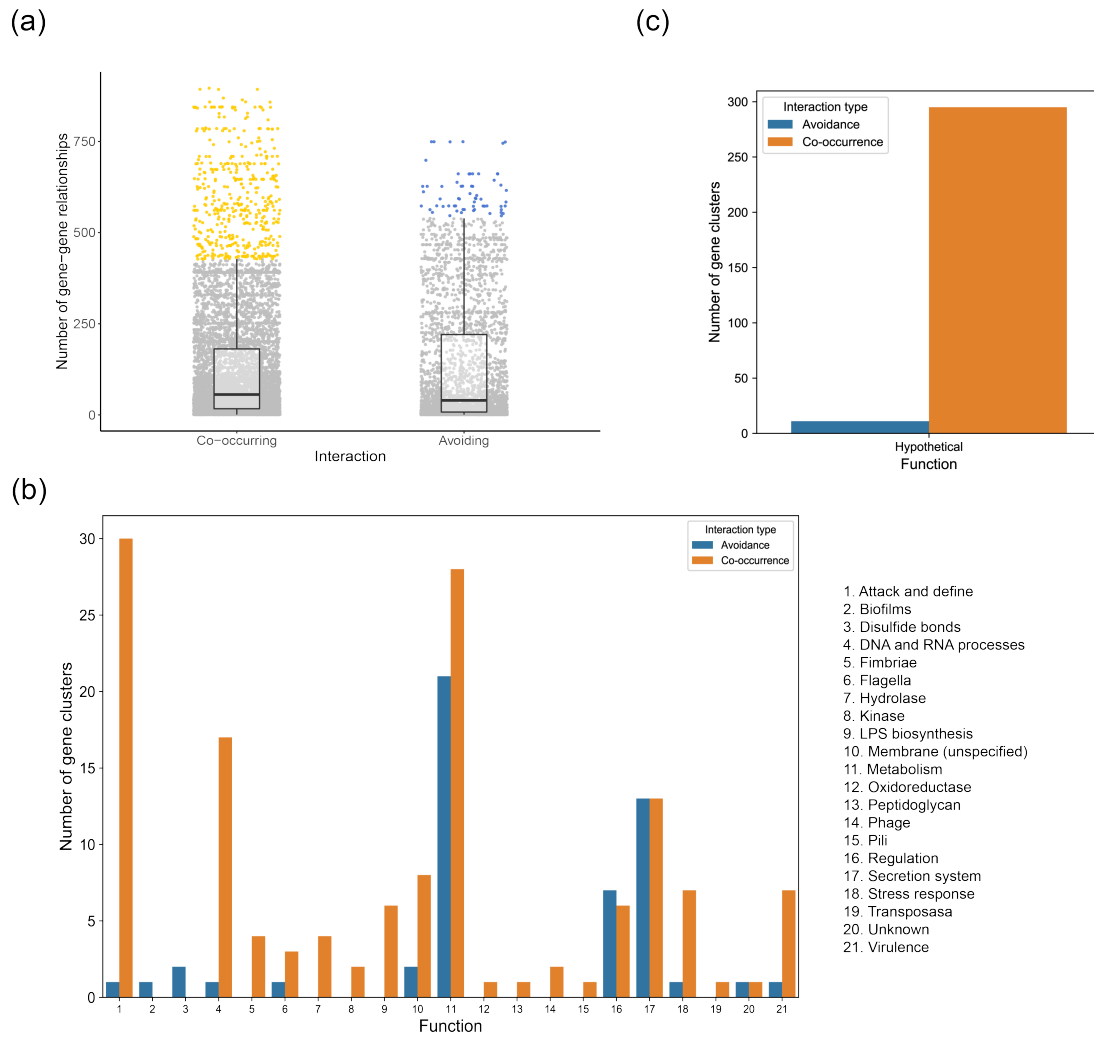


Fig. 2 Highly connected hub genes in the accessory genome. (a) The number of gene-gene relationships formed by individual genes. Hub genes (1.5 times the upper IQR) are coloured orange (co-occurrence) and blue (avoidance). (b) Co-occurrence (orange) and avoidance (blue) hub genes categorised by loose biological function. Attack and defence includes CRISPR system subunits, toxin-antitoxin systems, toxin production, detoxification. DNA and RNA processes includes replication, repair, recombination, and plasmid segregation. (c) The number of hub genes encoding hypothetical proteins.

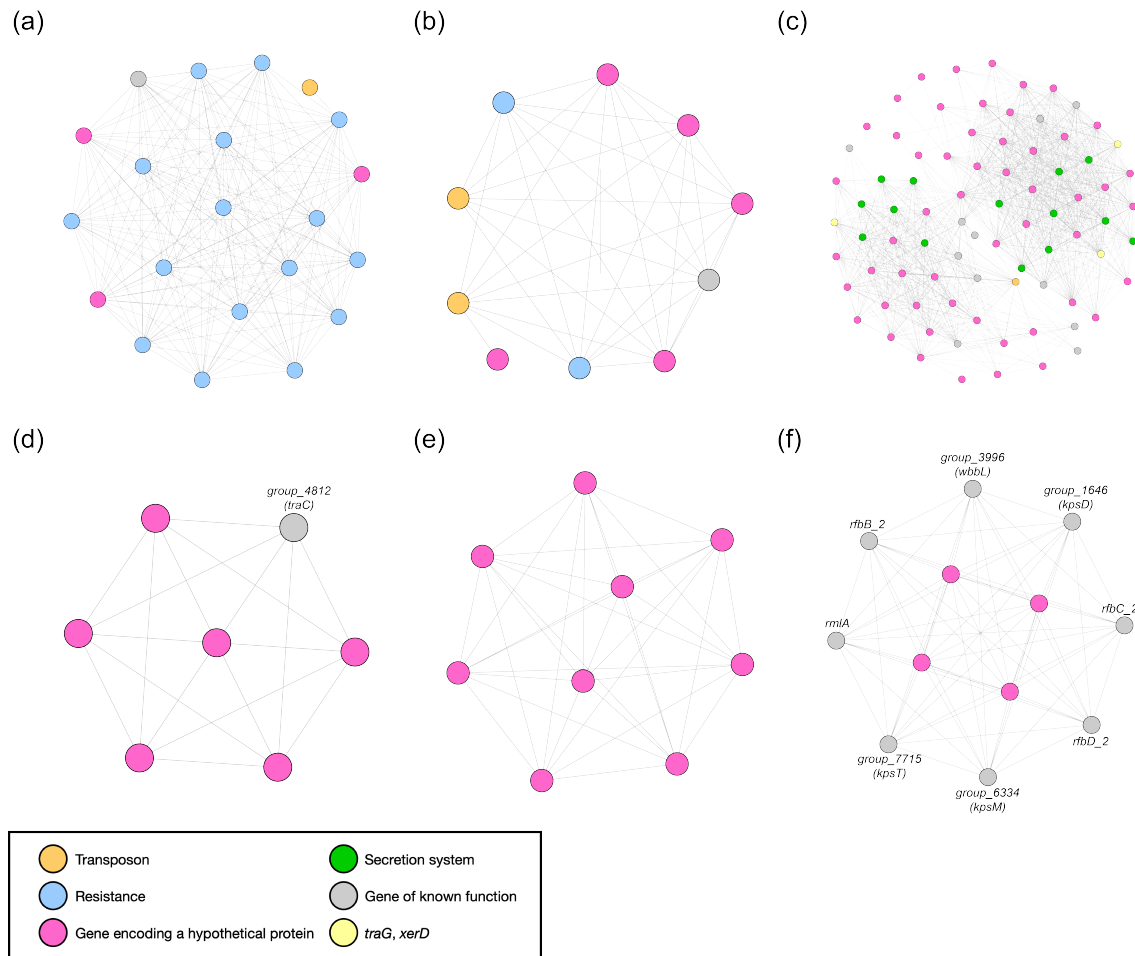


Fig. 3 Select co-occurrence connected components are linked to transposons and ones enriched in genes encoding hypothetical proteins. (a) Co-occurrence component 81, enriched in genes related to metal detoxification. (b) Co-occurrence component 53, linked to tetracycline resistance. (c) Co-occurrence component 2, comprised in part of genes encoding secretion systems. (d) Co-occurrence components 198 and (e) 89 are enriched in genes encoding hypothetical proteins. (f) Co-occurrence component 90 contains genes relating to capsule formation. Select gene clusters are labelled.

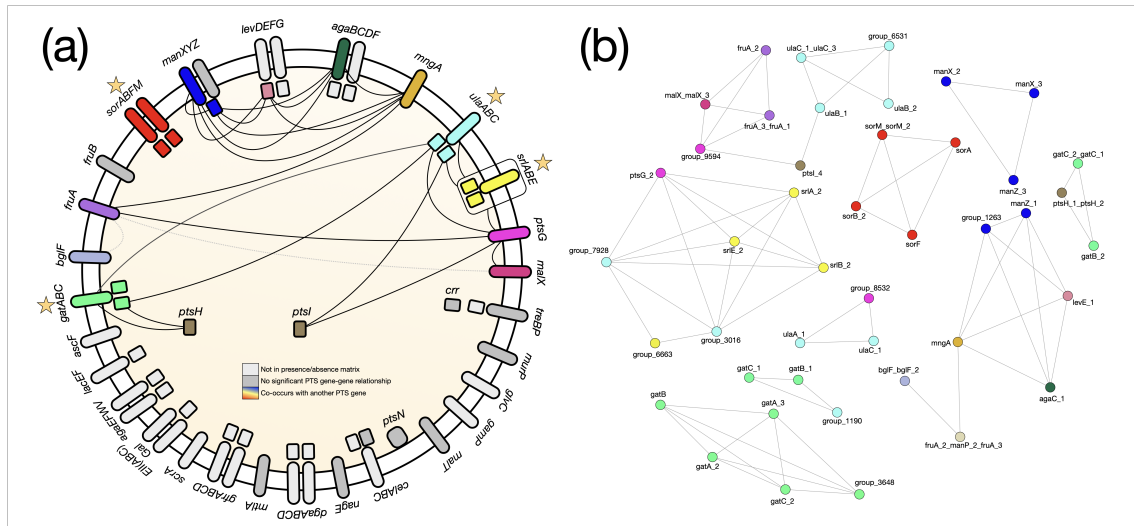


Fig. 4 PTS gene relationships are not universal. (a) A schematic overview of the co-occurrence PTS gene relationships. PTS gene systems found in the KEGG database are coloured by those not present in this *E. coli* pangenome (white), those which do not form a coincident PTS-PTS pair in the accessory genome (grey), and the clusters that form a co-occurring pair with another PTS gene (coloured by system). A significant relationship is indicated by a solid black line. Where a relationship is observed with all genes in a PTS, the line connects to a box around the system. Transporters that consist of genes that all significantly co-occur with one another are marked with a yellow star. The grey dashed line connecting *fruA* with *bglF* and *malX* indicates a significant relationship in both the co-occurrence and avoidance datasets for different clusters of the same gene identification. (b) The specific co-occurrences by gene cluster. Nodes are coloured by system as in (a).

572 Supplementary

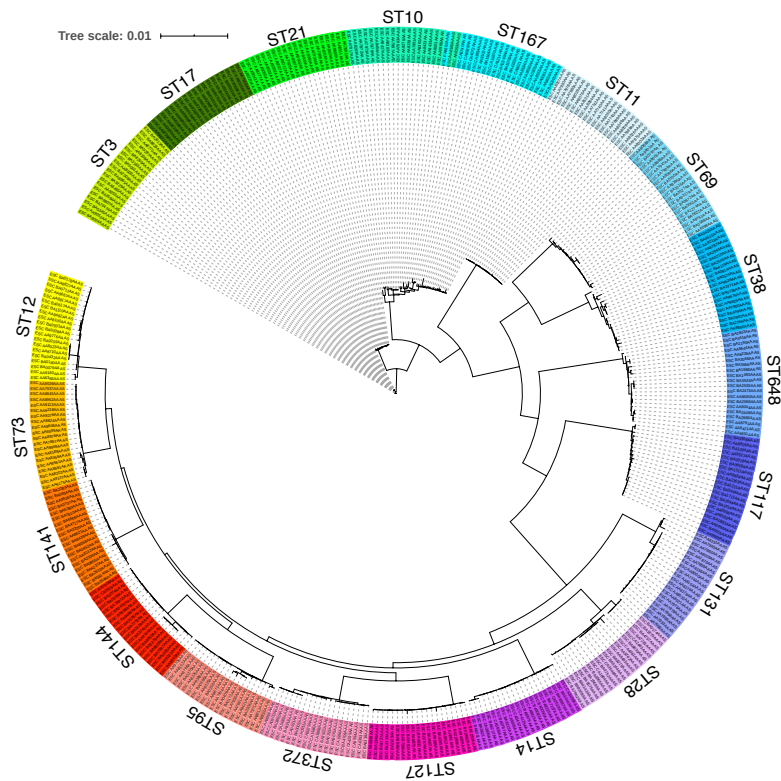


Fig. S1 An *E. coli* core gene phylogeny. Phylogeny of the 400 genomes used in this study, constructed from the trimmed core gene alignment using IQ-Tree based on the GTR+I+G substitution model. Labels are coloured by ST.

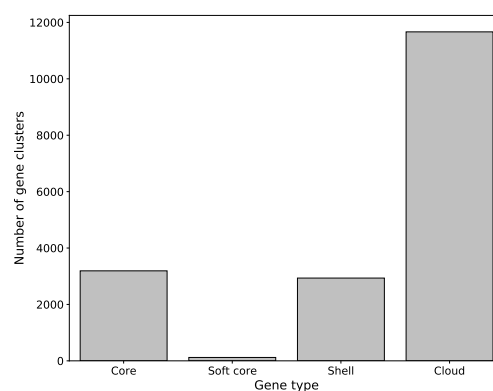


Fig. S2 Distribution of core, soft core, shell, and cloud genes in this *E. coli* pangenome.

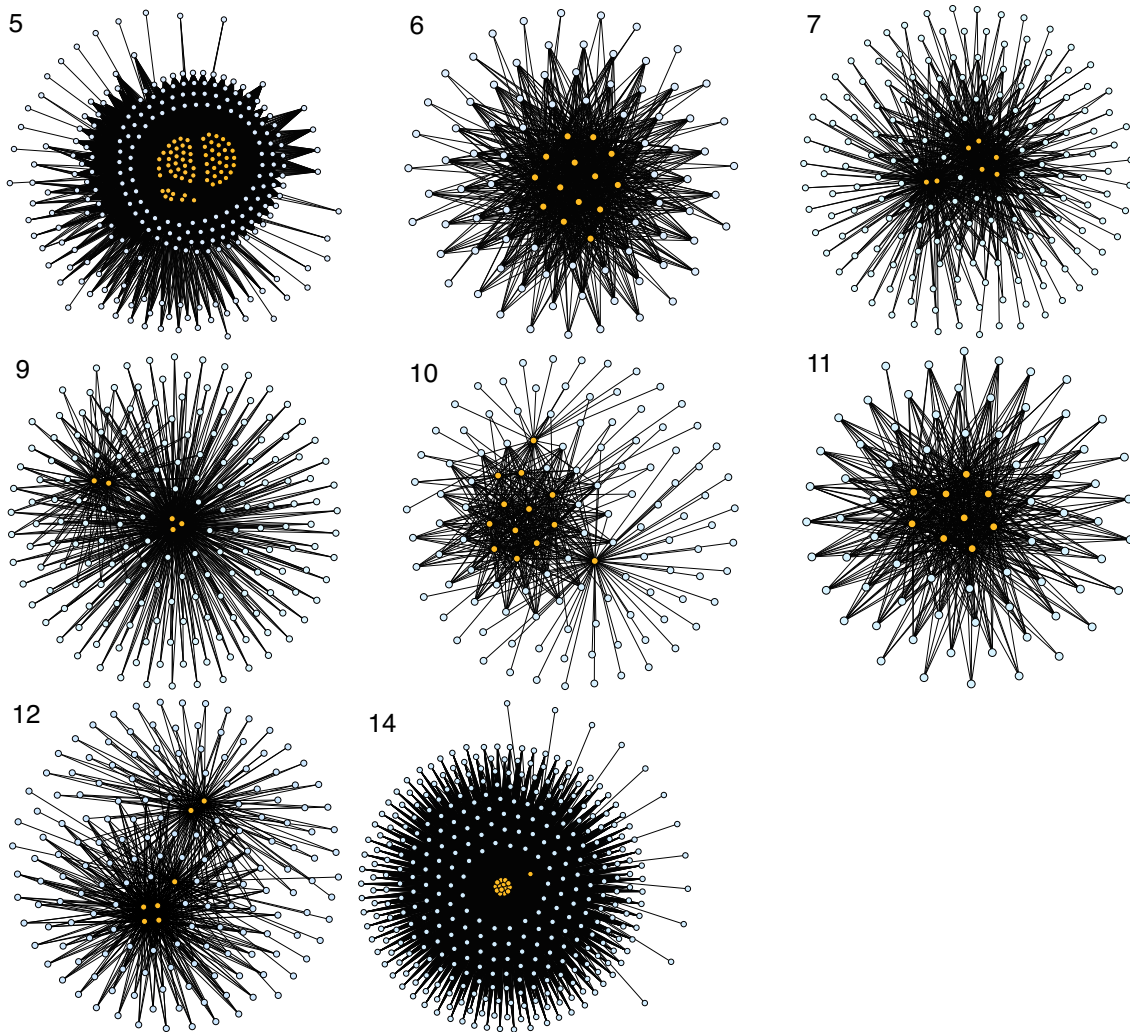


Fig. S3 Gene relationships in avoidance connected components that contain a root excluder. Genes avoided by the root excluders (orange) co-occur with one another and with other genes in the accessory genome (blue). Networks numbered according to avoidance connected component in Fig. 1b and the Supplementary Data files.

Table S1 All root excluder genes and their corresponding avoidance connected component. Component number refers to that given in the Coinfinder output provided as Supplementary Data. The gene cluster IDs are given as in the Panaroo gene presence-absence matrix.

| Component | Gene name | Gene cluster ID | Protein function |
|-----------|-------------------|-----------------------------|-------------------------------------|
| 5 | <i>rtcB</i> | rtcB_rtcB_1_rtcB_2 | RNA ligase |
| 6 | <i>ybdO</i> | ybdO_2_ybdO_1_leuO_2_ybdO_3 | Transcription regulator |
| 7 | <i>tam</i> | tam_tam_2 | Trans-aconitate 2-methyltransferase |
| 9 | <i>evgS</i> | evgS_evgS_2 | Two-component system sensor protein |
| 10 | <i>group_9436</i> | group_9436 | Hypothetical protein |
| 11 | <i>nlpA</i> | nlpA | Lipoprotein 28 |
| 12 | <i>yfkM</i> | yfkM_yfkM_1 | General stress protein |
| 14 | <i>dhaR</i> | dhaR | Transcription factor |