

30 **Abstract**

31 Enterohemorrhagic *Escherichia coli* O157:H7 (EHEC) is an important food-borne
32 pathogen that colonizes the colon. Transposon-insertion sequencing (TIS) was used to
33 identify genes required for EHEC and commensal *E. coli* K-12 growth in vitro and for
34 EHEC growth in vivo in the infant rabbit colon. Surprisingly, many conserved loci
35 contribute to EHEC's but not to K-12's growth in vitro, suggesting that gene acquisition
36 during EHEC evolution has heightened the pathogen's reliance on certain metabolic
37 processes that are dispensable for K-12. There was a restrictive bottleneck for EHEC
38 colonization of the rabbit colon, which complicated identification of EHEC genes
39 facilitating growth in vivo. Both a refined version of an existing analytic framework as
40 well as PCA-based analysis were used to compensate for the effects of the infection
41 bottleneck. These analyses confirmed that the EHEC LEE-encoded type III secretion
42 apparatus is required for growth in vivo and revealed that only a few effectors are critical
43 for in vivo fitness. Numerous mutants not previously associated with EHEC
44 survival/growth in vivo also appeared attenuated in vivo, and a subset of these putative
45 in vivo fitness factors were validated. Some were found to contribute to efficient type-
46 three secretion while others, including *tatABC*, *oxyR*, *envC*, *acrAB*, and *cvpA*, promote
47 EHEC resistance to host-derived stresses encountered in vivo. *cvpA*, which is also
48 required for intestinal growth of several other enteric pathogens, proved to be required
49 for EHEC, *Vibrio cholerae* and *Vibrio parahaemolyticus* resistance to the bile salt
50 deoxycholate. Collectively, our findings provide a comprehensive framework for
51 understanding EHEC growth in the intestine.

52

53 **Author Summary**

54 Enterohemorrhagic *E. coli* (EHEC) are important food-borne pathogens that infect the
55 colon. We created a highly saturated EHEC transposon library and used transposon
56 insertion sequencing to identify the genes required for EHEC growth in vitro and in vivo
57 in the infant rabbit colon. We found that there is a large infection bottleneck in the rabbit
58 model of intestinal colonization, and refined two analytic approaches to facilitate
59 rigorous identification of new EHEC genes that promote fitness in vivo. Besides the
60 known type III secretion system, more than 200 additional genes were found to
61 contribute to EHEC survival and/or growth within the intestine. The requirement for
62 some of these new in vivo fitness factors was confirmed, and their contributions to
63 infection were investigated. This set of genes should be of considerable value for future
64 studies elucidating the processes that enable the pathogen to proliferate in vivo and for
65 design of new therapeutics.

66

67 **Introduction**

68 Enterohemorrhagic *Escherichia coli* (EHEC) is an important food-borne pathogen that
69 causes gastrointestinal (GI) infections worldwide. EHEC is a non-invasive pathogen that
70 colonizes the human colon and gives rise to sporadic infections as well as large
71 outbreaks (reviewed in (1–3)). The clinical consequences of EHEC infection range from
72 mild diarrhea to hemorrhagic colitis and include the potentially lethal hemolytic uremic
73 syndrome (HUS) (4,5).

74

75 The paradigmatic EHEC O157:H7 strain, EDL933, caused the first recognized EHEC
76 outbreak in 1982 (6), and its genome shares a common 4.1 Mb DNA backbone with the
77 non-pathogenic laboratory strain of *E. coli* K-12, MG1655 (7–9). However, the EDL933
78 genome also contains 1.34Mb of chromosomal DNA that is absent from K-12, as well as
79 a 90kb virulence plasmid pO157. EDL933-specific ‘O-islands’ encode genes recognized
80 as the major EHEC virulence factors and are thought to have been acquired by
81 horizontal gene transfer; many are encoded within putative prophage elements.

82

83 Although there are a variety of EHEC serotypes and the O-island complement in
84 different EHEC isolates can differ (10), EHEC genomes all contain one or more
85 prophages encoding Shiga toxins and the Locus of Enterocyte Effacement (LEE)
86 pathogenicity island (11). These two horizontally acquired elements are critical EHEC
87 virulence determinants. Shiga toxins contribute to diarrhea and the development of HUS
88 (4,5,12). The LEE encodes a type III secretion system (T3SS) and several secreted
89 effectors. EHEC’s T3SS mediates attachment of the pathogen to colonic enterocytes,

90 effacement of the brush border microvilli, and the formation of actin-rich pedestal-like
91 structures underneath attached bacteria (reviewed in (13)). Once translocated into the
92 host cell, T3SS effectors, which are encoded both inside and outside the LEE, target
93 diverse signaling pathways and cellular processes (14,15). A functional LEE T3SS is
94 required for EHEC intestinal colonization in animal models as well as in humans
95 (12,13,16–20).

96

97 In addition to the virulence factors that prompt the key symptoms of infection, EHEC
98 also relies on bacterial factors that enable pathogen survival in and adaptation to the
99 host environment. During colonization of the human GI tract, EHEC encounters multiple
100 host barriers to infection, including, but not limited to stomach acid, bile, and other host-
101 and microbiota-derived compounds with antimicrobial properties (reviewed in (21)).
102 EHEC is known to detect intestinal cues derived from the host and the microbiota to
103 activate expression of virulence genes and to modulate gene expression both
104 temporally and spatially (22–25). However, comprehensive analyses of bacterial factors
105 that contribute to EHEC survival within the host have not been reported.

106

107 The development of transposon-insertion sequencing (TIS, aka TnSeq, InSeq, TraDis, or
108 HITS) (26–29) facilitated high-throughput and genome-scale analyses of the genetic
109 requirements for bacterial growth in different conditions, including in animal models of
110 infection (30–39). In this approach, the relative abundance of transposon-insertion
111 mutants within high-density transposon-insertion libraries provides insight into loci's
112 contributions to bacterial fitness in different environments (40,41). Potential insertion sites

113 for which corresponding insertion mutants are not recovered frequently correspond to
114 regions of the genome that are required for bacterial growth (often termed “essential
115 genes”), although the absence of a particular insertion mutant does not always reflect a
116 critical role for the targeted locus in maintaining bacterial growth (42,43). Comparative
117 analyses of the abundance of mutants in an initial (input) library and after growth in a
118 selective environment (e.g., an animal host) can be used to gauge loci’s contributions to
119 fitness in the selective condition.

120

121 Here, high-density transposon libraries were created in EHEC EDL933 and the
122 commensal *E. coli* K-12 and used to characterize their respective in vitro growth
123 requirements. The EHEC library was also passaged through an infant rabbit model to
124 identify genes required for intestinal colonization. Our data indicate that during infection
125 of the gastrointestinal tract, EHEC populations undergo a severe infection bottleneck that
126 complicates identification of genes with true in vivo fitness defects. We used two
127 complementary analytic approaches to circumvent the noise introduced by restrictive
128 bottlenecks to identify genes required for colonization of the colon. More than 240 genes
129 were found to contribute to efficient colonization of the rabbit colon. As expected, these
130 included the LEE-encoded T3SS and *tir*, a LEE-encoded effector necessary for intestinal
131 colonization (44). In addition, 2 non-LEE effectors and many additional new genes that
132 encode components of the bacterium’s metabolic pathways and stress response systems
133 were found to enable bacterial colonization of the colon. Isogenic mutants for 17 loci,
134 including *cvpA*, a gene necessary for intestinal colonization by diverse enteric pathogens
135 (40,45,46), were constructed, validated in the infant rabbit model and tested in vitro under

136 stress conditions that model host-derived challenges encountered within the
137 gastrointestinal tract. *cvpA* was found to be specifically required for resistance to the bile
138 salt deoxycholate and therefore appears to be a previously unappreciated member of the
139 bile-resistance repertoire of diverse enteric pathogens.

140

141 **Results and Discussion**

142

143 **Identification of genes required for EHEC growth in vitro**

144 The mariner-based Himar1 transposon, which inserts specifically at TA dinucleotides
145 (reviewed in (47)) was used to generate a high-density transposon-insertion library in a
146 $\Delta lacI::aacC1$ (gentamicin-resistance cassette inserted at *lacI*) derivative of EDL933. The
147 library was characterized via high-throughput sequencing of genomic DNA flanking sites
148 of transposon-insertion. To map the reads, we used the most recent EDL933 genome
149 sequence (8) and annotation (NCBI, February 2017). Since this genome, unlike the
150 initial EHEC genome (7), has not been linked to functional information (e.g., the EHEC
151 KEGG database, KEGG reviewed in (48)), we generated a correspondence table in
152 which the new genome annotations (RS locus tags) are linked to the original “Z
153 numbers” (Table S1). This correspondence table enabled us to utilize historically
154 valuable resources as well as the updated genomic sequence and should also benefit
155 the EHEC research community. 137,805 distinct insertion mutants were identified,
156 which corresponds to 52.5% of potential insertion mutants with an average of ~21 reads
157 per genotype (Fig S1A). Sensitivity analysis revealed that nearly all mutants were
158 represented within randomly selected read pools containing ~2 million reads. Increasing

159 sequencing depth to ~3 million reads had a negligible effect on library complexity,
160 suggesting that a sequencing depth of ~3 million reads is sufficient to identify virtually all
161 constituent genotypes within this EHEC library (Fig S1A).

162
163 EHEC's 6032 annotated genes were binned according to the percentage of disrupted
164 TA sites within each gene, and the number of genes corresponding to each bin was
165 plotted (Fig 1A). As expected for a high-density transposon-insertion library (41), this
166 distribution was bimodal, with a minor peak comprised of genes disrupted in few
167 potential insertion sites (Fig 1A, left), and a major peak comprised of genes that are
168 disrupted in most or all potential insertion sites (Fig 1A, right). Based on the center of
169 the major right-side peak, we estimate that ~70% of non-essential insertion sites have
170 been disrupted in this EHEC library, a degree of complexity that enabled high-resolution
171 analysis of transposon-insertion frequency.

172
173 Further analysis of insertion site distribution was performed using a hidden Markov
174 model-based analysis pipeline (EL-ARTIST, see methods and (40) for detail), that
175 classifies loci with a low frequency of transposon-insertion across the entire coding
176 sequence as 'underrepresented' (often referred to as 'essential' genes) or across a
177 portion of the coding sequence as 'regional' (Fig 1B). All other loci are classified as
178 'neutral'. Of EHEC's 6032 genes, 895 genes were classified as underrepresented, 407
179 as regional, and 4,730 as neutral (Fig 1A, Table S2). In general, neutral genes (blue)
180 were disrupted in a higher percentage of TA sites than underrepresented genes (red) or
181 regional genes (purple), which displayed low and intermediate percentage of TA site

182 disruption, respectively (Fig 1B). Neutral genes are considered to be dispensable for
183 growth in LB, whereas non-neutral genes (regional and underrepresented genes
184 combined) likely have important functions for growth in this media or are otherwise
185 refractory to transposon insertion (42,43).

186

187 We identified Z Numbers (Table S1) and the linked Clusters of Orthologous Groups
188 (COG) (49,50) and KEGG pathways associated with the 1302 genes classified as non-
189 neutral (underrepresented and regional). Although some loci have no COG assigned,
190 714 genes were assigned a COG functional category (Table S2). Each COG category
191 was plotted against its “COG Enrichment Index”, which is calculated as the percentage
192 of non-neutral genes in each COG category divided by the percent of the whole genome
193 with that COG (51). A subset of COGs, particularly cell cycle control, translation, lipid
194 and coenzyme metabolism, and cell wall biogenesis were associated with non-neutral
195 genes at a frequency significantly higher than expected based on their genomic
196 representation (Fig S2A). Collectively, the COG and a similar KEGG analysis (Table S3)
197 revealed that EHEC genes with non-neutral transposon-insertion profiles are associated
198 with pathways and processes often linked to essential genes in other organisms (52).

199

200 Non-neutral genes comprise ~22% of EHEC’s annotated genes, a proportion of the
201 genome that is substantially larger than the 8% and 9% observed in analogous TIS-
202 based characterizations of *Vibrio cholerae* and *Vibrio parahaemolyticus* (40,45). Of the
203 1,302 non-neutral EHEC genes, only 760 are homologous to a gene in *E. coli* K-12
204 MG1655 (>90% nucleotide identity or >90% amino acid identity across 90+% of gene

205 length) (Table S2, column M); thus, EDL933-specific loci comprise a high proportion
206 (~42%) of EHEC's underrepresented loci. The enrichment of underrepresented loci
207 among EHEC-specific genes, many of which were acquired by horizontal gene transfer,
208 may reflect factors that can limit transposon-insertion other than fitness costs.

209

210 Previous analyses revealed that nucleoid binding proteins such as HNS, which binds to
211 DNA with low GC content, can hinder Himar1 insertion (42). Consistent with this
212 observation, genes classified as non-neutral have a lower average GC content than
213 genes classified as neutral (Fig S2B; blue vs red distributions). The disparity in GC
214 content between neutral and non-neutral loci is particularly marked for EHEC genes that
215 do not have a homolog in K-12 (divergent; Fig 1C), although there is also a significant
216 difference between the GC content of neutral and non-neutral loci with a K-12 homolog
217 (homolog; Fig 1C). These analyses suggest that there is an association between GC
218 content and transposon-insertion frequency in EHEC, as in other organisms, and that
219 the prevalence of underrepresented loci among divergent loci may in part stem from the
220 lower average GC content of these loci (Fig S2C). Additional studies are necessary to
221 determine if the association between low GC content and reduced transposon-insertion
222 is due to HNS-binding, other nucleoid-associated proteins, or as yet unidentified fitness-
223 independent transposon insertion biases.

224

225 **TIS-based comparison of EHEC and *E. coli* in vitro growth requirements**

226 To evaluate whether the abundance of non-neutral loci was specific to EHEC or was
227 characteristic of additional *E. coli* strains, a high-density transposon-insertion library was

228 constructed in a $\Delta lac::cat$ (chloramphenicol-resistance cassette inserted at *lacI*)
229 derivative of *E. coli* K-12 MG1655. EL-ARTIST analysis of the high-density K-12 library
230 (Fig S1B) was implemented with the same parameters as those for the EHEC library
231 analysis and classified 24% of genes as underrepresented (786 underrepresented, 300
232 regional and 3397 neutral; Fig 1D, Table S4). Comparison of the gene classification of
233 homologous loci (Table S2 vs Table S4) revealed substantial concordance between the
234 sets of genes with non-neutral insertion profiles in EHEC and K-12: 83% (629/760) of
235 the non-neutral EHEC genes with homologs in K-12 were likewise classified as non-
236 neutral in the *E. coli* K-12 strain (Fig 1E). Thus, analyses of non-neutral loci suggest
237 either that the majority of ancestral loci make similar contributions to the survival and/or
238 proliferation of EHEC and K-12 strains in LB or that they are similarly resistant to
239 transposon-insertion.

240

241 We further explored the 131 underrepresented EHEC loci (Table S5) that were
242 classified as neutral (able to sustain insertions) in *E. coli* K-12. Most of these genes
243 have are linked to KEGG pathways for metabolism, particularly metabolism of
244 galactose, glycerophospholipid, and biosynthesis of secondary metabolites (Table S5,
245 Fig 1F). While this divergence could reflect the laboratory adaptation of the K-12 isolate,
246 gene acquisition during EHEC evolution may have heightened the pathogen's reliance
247 on metabolic processes that are not critical for growth of K-12. Such ancestral genes
248 may be useful targets for antimicrobial agents, as they might antagonize EHEC growth
249 without disruption of closely related commensal Enterobacteriaceae populations.

250

251 **Comparison of TIS and deletion-based gene classification**

252 The sets of genes classified as underrepresented or regional in EHEC and K-12
253 transposon libraries were compared to the 300 genes classified as essential in the K-12
254 strain BW25113 based on their absence from a comprehensive library of single gene
255 knockouts (53–55). 98% of these genes (294/300) were also classified as
256 underrepresented or regional in EDL933 and MG1655 (Table S2 and S4). The few loci
257 previously classified as essential but not found to be underrepresented or regional in
258 our analysis include several small genes, whose low number of TA sites hampers
259 confident classification. One gene in this list, *kdsC*, was found to have insertions across
260 the gene in both EDL933 and MG1655 (Fig S2D). *kdsC* knockouts have also been
261 reported previously (56), confirming that this locus is not required for K-12 growth
262 despite the absence of an associated mutant within the Keio collection. Thus,
263 underrepresented and regional loci encompass, but are not limited to, loci previously
264 classified as essential.

265
266 Several factors likely account for over-estimation of loci as underrepresented or
267 regional. First, loci can be classified as underrepresented even when viable mutants are
268 clearly present within the insertion library (Fig 1A); insertions simply need to be
269 consistently less abundant across a segment of the gene than insertions at other
270 (neutral) sites. Loci may also be classified as underrepresented due to fitness-
271 independent insertion biases, as discussed above (42,43). Additional evidence that loci
272 categorized as non-neutral by transposon-insertion studies are not necessarily essential
273 for growth was provided by a recent study of essential genes in *E. coli* K-12 (57).

274 However, the more expansive non-neutral classification can provide insight into loci that
275 enable optimal growth, in addition to those that are required.

276

277 **Identification of EHEC genes required for growth in vivo**

278 To identify mutants deficient in their capacity to colonize the mammalian intestine, the
279 EHEC transposon library was orogastrically inoculated into infant rabbits, an established
280 model host for infection studies (12,44,58,59). Transposon-insertion mutants were
281 recovered from the colon at 2 days post-infection, and the sites and abundance of
282 transposon-insertion mutations were determined via sequencing, as described above.

283 The relative abundance of individual transposon-insertion mutants in the library
284 inoculum was compared to samples independently recovered from the colons of 7
285 animals to identify insertion mutants that were consistently less abundant in libraries
286 recovered from the colon. Under ideal conditions, this signature is indicative of negative
287 selection of the mutant during infection, reflecting that the disrupted locus is necessary
288 for optimal growth within the intestine.

289

290 Sequencing and sensitivity analyses of the 7 passaged libraries revealed that they
291 contained substantially fewer unique insertion mutants than the library inoculum (23-
292 38% total mutants recovered, ~30,000 of 120,000) (Fig S1C-J). These data are
293 suggestive of population constrictions that could have arisen from 2 distinct but not
294 mutually exclusive causes: 1) negative selection, leading to depletion of mutants
295 deficient at in vivo survival or intestinal colonization; and/or 2) infection bottlenecks,
296 population constrictions that lead to stochastic reductions in the average number of

297 insertions per gene, independent of genotype or selective pressures. We binned genes
298 according to the percentage of TA sites disrupted within their gene sequences and
299 plotted the number of genes corresponding to each bin for both the inoculum (Fig 2A-
300 top) and a representative rabbit-passaged sample (Fig 2A-bottom). The passaged
301 sample exhibited a marked leftward shift relative to the inoculum, a signature indicative
302 of population constriction due to an infection bottleneck (41,60).

303
304 The TIS data was further analyzed using the Con-ARTIST pipeline. Con-ARTIST uses
305 iterative simulation-based normalization to compensate for experimental bottlenecks to
306 facilitate discrimination between stochastic reductions in genotype abundance and
307 reductions attributable to bona fide negative selection (mutants for which there was a
308 fitness cost in the host environment) (40). The Con-ARTIST analysis protocol and
309 subsequent gene classification is schematized in Figure 2B. For libraries recovered
310 from each rabbit, we used this workflow to classify genes as ‘conditionally depleted’
311 (red, CD), ‘queried’(blue) or ‘insufficient data’ (black) (Fig 2B) compared to the inoculum
312 library. CD genes contain sufficient insertions for analysis (see methods) and meet a
313 standard of a 4-fold reduction in read abundance that is consistent across TA sites in a
314 gene (Fig 2B). Queried genes contained sufficient insertions for analysis but failed to
315 meet this criterion. Genes classified as insufficient do not contain sufficient insertions for
316 analysis. The output of gene categorization using these thresholds is displayed for a
317 single rabbit in Fig 2C (additional animals in Fig S3A-G) and summarized for all animals
318 in Fig 2D.

319

320 The restrictive bottleneck and animal to animal variation led to differences in the
321 numbers of genes classified as CD in different rabbits (Fig 2D). Due to this variability,
322 an additional criterion that genes be classified as CD in 5 or more of the 7 animals
323 analyzed was imposed to create a consensus cutoff. In contrast to the >2000 genes
324 classified as conditionally depleted in one or more animals, only 243 genes were
325 classified as conditionally depleted across 5 or more animals (Fig S3H, Table S6).
326 These relatively stringent standards were imposed in order to identify robust candidates
327 for genes that facilitate EHEC intestinal growth, despite the limitations of the infection
328 bottleneck in this experimental model. Therefore, we do not conclude that genes
329 classified as “Queried” (3860) or “Insufficient Data” (1926) are *not* attenuated relative to
330 the wild type strain in vivo; it is likely that the list of CD loci (Table S6) is incomplete.
331
332 Using our Z correspondence table (Table S1), 89% (217/243) genes classified as
333 conditionally depleted were assigned to a COG functional category. CD genes were
334 frequently associated with amino acid and nucleotide metabolism, signal transduction,
335 and cell wall/envelope biogenesis, but only amino acid metabolism reached statistical
336 significance after correction for multiple hypothesis testing (Fig 2E). These genes are
337 also associated with KEGG metabolic pathways (particularly amino acid metabolism),
338 several two-component systems, including *qseC*, which has previously been implicated
339 in EHEC virulence gene regulation, and lipopolysaccharide biosynthesis (22) (Table S7,
340 Fig 2F). 33 of the 243 CD genes are EHEC specific, whereas the remaining 210 have
341 homologs in K-12 (Table S6), highlighting the importance of conserved metabolic
342 pathways in the pathogen’s capacity to successfully colonize its colonic niche. Similar

343 metabolic pathways were also found to be important for *V. cholerae* growth in the infant
344 rabbit small intestine (40,61), and raising the possibility of targeting metabolic pathways
345 such as those for amino acid biosynthesis with antibiotics (62–65).

346

347 The stochastic loss of individual insertion mutants in severely bottlenecked data can
348 hinder Con-ARTIST-based identification of CD genes. In particular, meeting the
349 pipeline's consistency threshold as measured by the Mann Whitney U (MWU) p-value
350 (Fig 2B) is difficult because severe bottlenecks drastically decrease the number of
351 individual transposon-insertion mutants per gene; therefore, the mutant replicates
352 needed to demonstrate consistency are not present. Queried genes often did not meet
353 the MWU p-value cut-off, even though fold change information may suggest marked
354 attenuation. To reclaim some of the mutants that were not classifiable by Con-ARTIST,
355 we also used Comparative TIS (CompTIS), a principal components analysis (PCA)-
356 based framework (66), to compare the 7 libraries recovered from rabbit colons. PCA is a
357 dimensional reduction approach used to describe the sources of variation in multivariate
358 datasets. Recently, we found that PCA is useful for identifying genes whose inactivation
359 leads to mutant growth phenotypes that are consistent across TIS replicates (66). Here,
360 we applied CompTIS as an alternative approach to identify genes with phenotypic
361 consistency (inability to colonize the rabbit colon) in all 7 rabbit replicates.

362

363 To perform CompTIS, the fold change of each gene from the seven colon libraries was
364 subjected to gene level PCA (gIPCA) (see methods and (66)), with each library
365 recovered from a rabbit colon representing a replicate. gIPC1 describes most of the

366 variation in the animals (Fig S3I) and reports a weighted average of the fold change
367 values for each gene across the 7 animals (Table S6). The signs and magnitudes of
368 PC1 were all similar (Fig S3J), indicating that each rabbit contributes approximately
369 equally to PC1, as expected for biological replicates. The distribution of gIPC1 scores is
370 continuous (Fig 2G) and describes each gene's contribution to EHEC intestinal
371 colonization. Most genes have a gIPC1 score close to zero (average PC1=0),
372 suggesting that they do not contribute to colonization. However, the distribution includes
373 a left tail beginning at PC1 scores of approximately -900 that encompassed the lowest
374 10% of scores (Fig 2G, Fig S3K), which likely correspond to genes contributing to
375 colonization. This list of 541 genes included nearly all (85%) of the genes classified as
376 CD by the more conservative Con-ARTIST analysis outlined above (Fig 2B). This
377 method allows for identification of additional candidate genes required for
378 survival/growth in vivo. For example, the PCA approach captured genes such as *ler*
379 (PC1 = -2290), a critical activator of the LEE T3SS (67), which was classified as queried
380 by Con-ARTIST due to the relative paucity of unique insertion mutants.

381

382 **Analyses of the requirement for T3SS and its associated effectors in colonization**

383 To begin to assess the accuracy of our gene classifications using the Con-ARTIST
384 consensus approach and CompTIS, we examined classifications within the LEE
385 pathogenicity island, which encodes the EHEC T3SS and plays a critical role in
386 intestinal colonization (12,16–19). The LEE is comprised of 40 genes, including genes
387 encoding the structural components of the T3SS, some of the pathogen's effectors, their
388 chaperones, and Intimin (encoded by *eae*), the adhesin that binds to the translocated

389 Tir protein. In infant rabbits, previous studies using single deletion mutants revealed that
390 *tir*, *eae*, and *escN*, the T3SS ATPase, were all required for colonization (12,44). We
391 observed a marked reduction in the abundance of insertions across nearly the entire
392 LEE in the samples from the rabbit colons relative to the simulation-normalized input
393 reads (Fig 3A). The 3 genes previously found to be required for colonization (*tir*, *eae*
394 and *escN*) were classified as CD using the Con-ARTIST consensus approach,
395 enhancing confidence in this scheme. Furthermore, 8 additional LEE-encoded genes
396 critical for T3SS activity, including translocon T3SS components (*espB*, *espD*, and
397 *espA*) and structural components (*escD*, *escQ*, *escV*, *escI*, and *escC*) were also
398 classified as CD using this scheme (Fig 3AB, Table S6) (16,68–72). Our findings
399 provide additional strong evidence that the LEE T3SS is critical for EHEC proliferation in
400 the intestine. However, many LEE-encoded genes had insufficient data to enable
401 classification via the Con-ARTIST consensus approach.

402
403 In contrast to Con-ARTIST analytic approach, with the PCA-based CompTIS analysis,
404 we were able to assess the contribution of all of the genes in the LEE and identify
405 several more genes likely to be important for in vivo colonization. With this approach,
406 most LEE genes had gIPC1 scores in the bottom 10% of the distribution (Table S6, Fig
407 3B). Notably, the genes that were not in this portion of the distribution included 4
408 effectors (*espF*, *map*, *espH*, *espG*). Previous studies in infant rabbits showed that *espG*
409 and *map* were dispensable for colonic colonization and that *espH* and *espF* mutants
410 only had modest colonization defects (44), lending credence to PCA-based
411 classification.

412
413 EHEC has a large suite of non-LEE encoded effectors (Nle), many of which reside
414 within prophage elements. Only 2 of 43 Nle genes (*nleA*, *espM1*) were classified as CD
415 by Con-ARTIST or were found within the bottom 10% of gIPC1 scores by CompTIS
416 (Fig3B, Table S6), suggesting that only a small subset of EHEC effectors are critical for
417 colonization, while other effectors likely play auxiliary roles. NleA was previously
418 reported to be important for colonic colonization by a related enteric pathogen,
419 *Citrobacter rodentium* (73), and is thought to suppress inflammasome activity (74), while
420 EspM1 is thought to modulate host actin cytoskeletal dynamics (75,76). Additional
421 studies are warranted to confirm and further explore how these 2 effectors play pivotal
422 roles promoting intestinal colonization.

423
424 **Validation of colonization defects in non-LEE encoded genes classified as**
425 **conditionally depleted**

426 We performed further studies of 17 conditionally depleted genes/operons that had not
427 previously been demonstrated to promote EHEC intestinal colonization. The Con-
428 ARTIST consensus approach and CompTIS classified all of these genes as
429 conditionally depleted except one, *hupB*, which was classified as queried by Con-
430 ARTIST but within the bottom 10% of gIPC1 scores (Fig 4A). Mutants with in-frame
431 deletions of either single loci (*agaR*, *cvpA*, *envC*, *htrA*, *hupB*, *mgtA*, *oxyR*, *prc*, *sspA*,
432 *sufI*, *tolC*, and RS09610, a hypothetical gene of unknown function) or operons with one
433 or more genes classified as conditionally depleted (*acrAB*, *clpPX*, *envZompR*, *phoPQ*,
434 *tatABC*) were generated. Then, each mutant strain was barcoded with unique sequence

435 tags integrated into a neutral locus in order to enable multiplexed analysis. The in vitro
436 growth of the barcoded mutants was indistinguishable from that of the WT strain (Fig
437 S4A), suggesting that the transposon mutants' in vivo attenuation is not explained by a
438 generalized growth deficiency.

439

440 The barcoded mutants, along with the barcoded WT EHEC, were co-inoculated into
441 infant rabbits to compare the colonization properties of the mutants and WT. The
442 relative frequencies of WT and mutant EHEC within CFU recovered from infected
443 animals was enumerated by deep sequencing of barcodes, and these frequencies were
444 used to calculate competitive indices (CI) for each mutant (i.e., relative abundance of
445 mutant/WT tags in output normalized to input). 14 of the 17 mutants tested had CI
446 values significantly lower than 1, validating the colonization defects inferred from the
447 TIS data (Fig 4). In aggregate, these observations support our experimental and
448 analytical approaches and suggest that many of the genes classified as CD by the Con-
449 ARTIST consensus approach and/or have low PC-1 scores may also contribute to
450 intestinal colonization.

451

452 **Many conditionally depleted loci exhibit reduced T3SS effector translocation**
453 **and/or increased sensitivity to extracellular stressors**

454 The many new genes implicated in EHEC colonization by the TIS data could contribute
455 to the pathogen's survival and growth in vivo by a large variety of mechanisms. Given
456 the pivotal role of EHEC's T3SS in intestinal colonization, as well as previous
457 observations that factors outside the LEE can regulate T3SS gene expression and/or

458 activity (reviewed in (67)), we assessed whether T3SS function was impaired in the 11
459 mutants with CIs <0.3 (Fig 4). Translocation of EspF (an effector protein) fused to a
460 TEM-1 beta-lactamase reporter into HeLa cells was used as an indicator of T3SS
461 functionality (77). An Δ *escN* mutant, which lacks the ATPase required for T3SS
462 function, was used as a negative control.

463

464 Deletions in three protease-encoded genes, *clpPX*, *htrA*, and *prc*, were associated with
465 reduced EspF translocation (Fig 5A). Both ClpXP and HtrA have been implicated in
466 T3SS expression/activity in previous reports (78–81). The ClpXP protease controls LEE
467 gene expression indirectly by degrading LEE-regulating proteins RpoS and GrIR (82).
468 The periplasmic protease HtrA (aka DegP) has been implicated in post-translational
469 regulation of T3SS as part of the Cpx-envelope stress response (80,81). Interestingly,
470 *prc*, which also encodes a periplasmic protease (82), also appears required for robust
471 EspF translocation. Prc has been implicated in the maintenance of cell envelope
472 integrity under low and high salt conditions in *E. coli* K-12 (83). Consistent with this
473 observation, in high osmolarity media a Δ *prc* EHEC mutant exhibited cell shape defects
474 (Fig S4B). Deficiencies in the cell envelope associated with absence of Prc may impair
475 T3SS assembly and/or function, perhaps also by triggering the Cpx-envelope stress
476 response. Together, these observations suggest that in vivo these three proteases
477 modulate T3SS expression/function, thereby promoting EHEC intestinal colonization.

478

479 We also investigated the capacity of each of the 11 mutant strains to survive challenge
480 with three stressors – low pH, bile, and high salt (osmotic challenge) – that the

481 pathogen may encounter in the gastrointestinal tract. Relative to the WT strain, all but
482 one (*sufI*) of the mutant strains exhibited reduced survival following one or more of
483 these challenges (Fig 5BC), suggesting that exposure to these host environmental
484 factors may contribute to the in vivo attenuation of these mutants. Many of the EHEC
485 mutants exhibited sensitivities to external stressors that are consistent with previously
486 described phenotypes in other organisms and experimental systems. For example, the
487 EHEC Δ *acrAB* locus, which was associated with bile sensitivity in EHEC (Fig 5B), is
488 known to contribute to a multidrug efflux system that can extrude bile salts, antibiotics,
489 and detergents (84). Our observation that mutants lacking the oxidative stress response
490 gene *oxyR* are sensitive to bile and to acid pH is also concordant with previous reports
491 linking both stimuli to oxidative stress (85–87). Furthermore, the heightened sensitivity
492 to bile, acid, and elevated osmolarity of EHEC lacking the two-component regulatory
493 system EnvZ/OmpR is consistent with previous reports that EnvZ/OmpR is a critical
494 determinant of membrane permeability, due to its regulation of outer membrane porins
495 OmpF and OmpC. Mutations that activate this signaling system (in contrast to the
496 deletions tested here) have been found to promote *E. coli* viability in vivo and to
497 enhance resistance to bile salts (88).

498

499 The EHEC Δ *tatABC* mutant exhibited a marked colonization defect and a modest
500 increase in bile sensitivity. The twin-arginine translocation (Tat) protein secretion
501 system, which transports folded protein substrates across the cytoplasmic membrane
502 (reviewed in (89,90)), has been implicated in the pathogenicity of a variety of Gram-
503 negative pathogens, including enteric pathogens such as *Salmonella enterica* serovar

504 Typhimurium (91–93), *Yersinia pseudotuberculosis* (94,95), *Campylobacter jejuni* (96),
505 and *Vibrio cholerae* (97). Attenuation of Tat mutants can reflect the combined absence
506 of a variety of secreted factors. For example, the virulence defect of *S. enterica*
507 Typhimurium *tat* mutants are likely due to cell envelope defects caused by the inability
508 to secrete the periplasmic cell division proteins AmiA, AmiC and SufI (92). Notably,
509 single knock-outs of any of these genes did not cause attenuation (92), but altogether
510 their absence renders the cell-envelope defective and more sensitive to cell-envelope
511 stressors, such as bile acids (93).

512

513 In EHEC, the Tat system has been implicated in Stx1 export (98), but since Stx1 was
514 not a hit in our screen and is not thought to modulate intestinal colonization (12), it is not
515 likely to explain the marked colonization defect of the EHEC Δ *tatABC* mutant. The suite
516 of EHEC Tat substrates has not been experimentally defined, although putative Tat
517 substrates can be identified by a characteristic signal sequence (89,90). A few
518 substrates, including SufI, OsmY, OppA, MglB, and H7 flagellin, have been detected
519 experimentally (98). *sufI*, interestingly, was also a validated hit in our screen and is the
520 only Con-ARTIST defined CD gene that has a predicted Tat-secretion signal. However,
521 the Δ *sufI* mutant did not display enhanced bile sensitivity, suggesting that attenuation of
522 this mutant, and perhaps of the Δ *tatABC* mutant as well, reflects deficiencies in other
523 processes. SufI is a periplasmic cell division protein that localizes to the divisome and
524 may be important for maintaining divisome assembly during stress conditions (99,100).
525 *E. coli tat* mutants have septation defects (101), presumably from loss of SufI at the
526 divisome. Interestingly, *envC*, another validated CD gene, encodes a septal murein

527 hydrolase (102) that is required for cell division, and the $\Delta envC$ mutant also displayed
528 increased bile sensitivity. Consistent with this hypothesis, in high osmolarity media, the
529 $\Delta sufl$, $\Delta envC$, and $\Delta tatABC$ mutants exhibited septation or cell shape defects (Fig S4B).
530 Collectively, these data suggest that an impaired capacity for cell division may reduce
531 EHEC's fitness for inraintestinal growth, and that at times this may reflect increased
532 susceptibility to clearance by host factors such as bile.

533

534 **CvpA promotes EHEC resistance to deoxycholate**

535 We further characterized EHEC $\Delta cvpA$ because other TIS-based studies of the
536 requirements for colonization by diverse enteric pathogens (*Vibrio cholerae*, *Vibrio*
537 *parahaemolyticus* and *Salmonella enterica* serovar Typhimurium) also classified *cvpA*
538 as important for colonization, but did not explore the reasons for the colonization
539 deficiency of the respective mutants (40,45,46).

540

541 *cvpA* encodes a putative inner membrane protein and has been linked to colicin V
542 export in *E. coli* K-12 (103) as well as curli production and biofilm formation in UPEC
543 (104). The EHEC $\Delta cvpA$ mutant did not exhibit an obvious defect in biofilm formation or
544 curli production (Fig S5AB), suggesting that *cvpA* may have a distinct role in EHEC
545 pathogenicity.

546

547 To further characterize the sensitivity of EHEC $\Delta cvpA$ mutant to bile, we exposed the
548 mutant to the two major bile salts found in the gastrointestinal tract, cholate (CHO) and
549 deoxycholate (DOC) (Fig 6AC) (85,105). Cholate is a primary bile salt, produced in the

550 liver and released into the biliary tract, while deoxycholate is a secondary bile salt that is
551 generated from cholate by intestinal bacteria in the colon. In contrast to WT EHEC,
552 which displayed equivalent sensitivity to the two bile salts in MIC assays (MIC = 2.5%
553 for both), the $\Delta cvpA$ mutant was much more sensitive to DOC than to CHO (MIC=
554 0.08% versus 1.25%). The $\Delta cvpA$ mutant's sensitivity to deoxycholate was present both
555 in liquid cultures and during growth on solid media (Fig. 6ABC).

556

557 Growth of the $\Delta cvpA$ mutant in the presence of deoxycholate was partially restored by
558 introduction of *cvpA* under the control of an inducible promoter, confirming that
559 sensitivity is linked to the absence of *cvpA* (Fig. 6A). *cvpA* lies upstream of the purine
560 biosynthesis locus *purF*, and some $\Delta cvpA$ mutant phenotypes have been attributed to
561 reduced expression on *purF* due to polar effects (103,106). The growth of the EHEC
562 $\Delta cvpA$ mutant was not impaired in the absence of exogenous purines (Fig S5C),
563 suggesting the *cvpA* deletion does not adversely modify *purF* expression.

564

565 Bile sensitivity has been associated with defects in the bacterial envelope or with
566 reduced efflux capacity (reviewed in (105)). We assessed the growth of the $\Delta cvpA$
567 mutant in the presence of a variety of agents that perturb the cell envelope to assess
568 the range of the defects associated with the absence of *cvpA*. The MICs of WT and
569 $\Delta cvpA$ EHEC were compared to those of an $\Delta acrAB$ mutant, whose lack of a broad-
570 spectrum efflux system provided a positive control for these assays. Notably, the $\Delta cvpA$
571 mutant did not exhibit enhanced sensitivity to any of the compounds tested other than
572 bile salts. In marked contrast, the $\Delta acrAB$ mutant displayed increased sensitivity to all

573 agents assayed (Fig 6C). These observations suggest that the sensitivity of the *cvpA*
574 mutant to DOC is not likely attributable to a general cell envelope defect in this strain. *V.*
575 *cholerae* and *V. parahaemolyticus* $\Delta cvpA$ mutants also exhibited sensitivity to DOC (Fig
576 S5D), implying a similar role in bile resistance in these distantly related enteric
577 pathogens.

578

579 A variety of bioinformatic algorithms (PSLPred, HHPred, Phobius, Phyre2) suggest that
580 CvpA is an inner membrane protein with 4-5 transmembrane elements similar to small
581 solute transporter proteins (Fig 6D). Phyre2 and HHPred reveal CvpA's partial similarity
582 to inner membrane transporters in the Major Facilitator Superfamily of transporters
583 (MFS) and the small-conductance mechanosensitive channels family (MscC). Additional
584 protein classification schemes group CvpA with proteins involved in solute transport. For
585 example, the PFAM database groups CvpA (PF02674) in the LysE transporter
586 superfamily (CL0292), a set of proteins known to enable solute export. In conjunction
587 with findings presented above, these predictions raise the possibility that CvpA is
588 important for the export of a limited set of substrates that includes DOC. Additional
589 studies to confirm this hypothesis and to establish how CvpA enables export are
590 warranted, particularly because this protein is widespread amongst enteric pathogens.

591

592 **Conclusions**

593 Here, we created a highly saturated transposon library in EHEC EDL933 to identify
594 the genes required for in vitro and in vivo growth of this important food-borne pathogen
595 using TIS. This approach has transformed our capacity to rapidly and fairly

596 comprehensively assess the contribution an organism's genes to growth in different
597 environments (41,107,108). However, technical and biologic issues can confound
598 interpretation of genome-scale transposon-insertion profiles. For example, we found
599 that EHEC genes with low GC content or those without homologs in K-12 were less
600 likely to contain transposon-insertions (Fig S2BC, Fig1C). Many of these genes were
601 likely acquired during EHEC evolution via lateral gene transfer; they constitute some of
602 the ~1.4MB of DNA that distinguishes EHEC from K-12 strains. Unexpectedly, more
603 than 100 of the genes conserved between EHEC and K-12 appear to promote the
604 growth of the pathogen in rich media but not that of K-12 (Table S5), suggesting that the
605 ~1.4MB of laterally acquired DNA that distinguishes EHEC and K-12 has enabled
606 divergence of the metabolic roles of ancestral *E. coli* genes in these backgrounds.

607

608 In animal models of infection, bottlenecks that result in marked stochastic loss of
609 transposon mutants can severely constrain TIS-based identification of genes required
610 for in vivo growth. Analysis of the distributions of the EHEC transposon-insertions in
611 vitro and in vivo (Fig. 2) revealed that there is a large infection bottleneck in the infant
612 rabbit model of EHEC colonization. Both Con-ARTIST, which applies conservative
613 parameters to define conditionally depleted genes (Fig 2B), and a PCA-based
614 approach, CompTIS, were used to circumvent the analytical challenges posed by the
615 severe EHEC infection bottleneck. These approaches should also be of use for similar
616 bottlenecked data that often hampers interpretation of TIS-based infection studies.
617 Validation studies, which showed that 14 of 17 genes (82%) classified as CD were
618 attenuated for colonization, suggest that these approaches are useful. Besides the LEE-

619 encoded T3SS, more than 200 additional genes were found to contribute to EHEC
620 survival and/or growth within the intestine. This set of genes should be of considerable
621 value for future studies elucidating the processes that enable the pathogen to proliferate
622 in vivo and for design of new therapeutics.

623

624 **Materials and Methods**

625

626 **Ethics statement**

627 All animal experiments were conducted in accordance with the recommendations in the
628 Guide for the Care and Use of Laboratory Animals of the National Institutes of Health
629 and the Animal Welfare Act of the United States Department of Agriculture using
630 protocols reviewed and approved by Brigham and Women's Hospital Committee on
631 Animals (Institutional Animal Care and Use Committee protocol number 2016N000334
632 and Animal Welfare Assurance of Compliance number A4752-01)

633

634 **Bacterial strains, plasmids and growth conditions**

635 Strains, plasmids and primers used in this study are listed in Supplementary Tables 8
636 and 9. Strains were cultured in LB medium or on LB agar plates at 37°C unless
637 otherwise specified. Antibiotics and supplements were used at the following
638 concentrations: 20 µg/mL chloramphenicol (Cm), 50 µg/mL kanamycin (Km), 10 µg/mL
639 gentamicin (Gent), 50 µg/mL carbenicillin (Carb), and 0.3 mM diaminopimelic acid
640 (DAP).

641

642 A gentamicin-resistant mutant of *E. coli* O157:H7 EDL933 ($\Delta lacI::aacC1$) and a
643 chloramphenicol-resistant mutant of *E. coli* K-12 MG1655 ($\Delta lacI::cat$) were used in this
644 study for all experiments, and all mutations were constructed in these strain
645 backgrounds except where specified otherwise. The $\Delta lacI::aacC1$ and $\Delta lacI::cat$
646 mutations were constructed by standard allelic exchange techniques (109) using a
647 derivative of the suicide vector pCVD442 harboring a gentamicin resistance cassette
648 amplified from strain TP997 (Addgene strain #13055) (110) or a chloramphenicol
649 resistance cassette from plasmid pKD3 (Addgene plasmid #45604) (53) flanked by the
650 5' and 3' DNA regions of the *lacI* gene. Isogenic mutants of EDL933 $\Delta lacI::aacC1$ were
651 also constructed by standard allelic exchange using derivatives of suicide vector pDM4
652 harboring DNA regions flanking the gene(s) targeted for deletion. *E. coli* MFD λ pir (111)
653 was used as the donor strain to deliver allelic exchange vectors into recipient strains by
654 conjugation. Sequencing was used to confirm mutations.

655

656 A $\Delta cvpA$ strain was also constructed using standard allelic exchange in a streptomycin-
657 resistant mutant (Sm^R) of *V. parahaemolyticus* RIMD 2210633. A *cvpA::tn* mutant was
658 used from a *Vibrio cholerae* C6706 arrayed transposon library (112).

659

660 **Transposon-insertion library construction**

661 To create transposon-insertion mutant libraries in EHEC EDL933 $\Delta lacI::aacC1$,
662 conjugation was performed to transfer the transposon-containing suicide vector pSC189
663 (113) from a donor strain (*E. coli* MFD λ pir) into the EDL933 recipient. Briefly, 100 μ L of
664 overnight cultures of donor and recipient were pelleted, washed with LB, and combined

665 in 20 μ L of LB. These conjugation mixtures were spotted onto a 0.45 μ m HA filter
666 (Millipore) on an LB agar plate and incubated at 37°C for 1 h. The filters were washed in
667 8 mL of LB and immediately spread across three 245x245 mm² (Corning) LB-agar
668 plates containing Gent and Kn. Plates were incubated at 37°C for 16 h and then
669 individually scraped to collect colonies. Colonies were resuspended in LB and stored in
670 20% glycerol (v/v) at -80°C as three separate library stocks. The three libraries were
671 pooled to perform essential genes analysis, and one library aliquot was used to as an
672 inoculum for infant rabbit infection studies.

673
674 To create TIS mutant libraries in *E. coli* K-12 MG1655 Δ *lacI::cat*, conjugation was
675 performed as above. 200 μ L of overnight culture of the donor strain (*E. coli* MFD λ pir
676 carrying pSC189) and the recipient strain (MG1655 Δ *lacI::cat*) were pelleted, washed,
677 combined and spotted on 0.45 μ m HA filters at 37°C for 5.5 hours. Cells were collected
678 from the filter, washed, plated on selective media (LB Kan, Cat), and incubated
679 overnight at 30°C. Colonies were resuspended in LB and frozen in 20% glycerol (v/v).
680 An aliquot was thawed and gDNA isolated for analysis.

681
682 **Infant rabbit infection with EHEC transposon-insertion library**
683 Mixed gender litters of 2-day-old New Zealand White infant rabbits were co-housed with
684 a lactating mother (Charles River). To prepare the EHEC transposon-insertion library for
685 infection of infant rabbits, 1 mL from one library aliquot was thawed and added to 20 mL
686 of LB. After growing the culture for 3 h at 37°C with shaking, the OD₆₀₀ was measured
687 and 40 units of culture at OD₆₀₀=1 (about 8 mL) were pelleted and resuspended in 10

688 mL PBS. Dilutions of the inoculum were plated on LB agar plates with Gent and Km for
689 precise dose determination. An aliquot of the inoculum was saved for subsequent gDNA
690 extraction and sequencing (input). Each infant rabbit was infected orogastrically with
691 500 μ l of the inoculum (1×10^9 cfu) using a size 4 French catheter. Following inoculation,
692 the infant rabbits were monitored at least 2x/day for signs of illness and euthanized 2
693 days postinfection. The entire intestinal tract was removed from euthanized rabbits,
694 and sections of the mid-colon were removed and homogenized in 1 mL of sterile PBS
695 using a minibeadbeater-16 (BioSpec Products, Inc.). 200 μ L of tissue homogenate from
696 the colon were plated on LB agar + Gm + Km to recover viable transposon-insertion
697 mutants. Plates were grown for 16 h at 37°C. The next day, colonies were scraped and
698 resuspended in PBS. A 5 mL aliquot of cells was used for genomic DNA extraction and
699 subsequent sequencing (Rabbits 1-7).

700

701 **Characterization of transposon-insertion libraries**

702 Transposon-insertion libraries were characterized as described previously. Briefly, for
703 each library, gDNA was isolated using the Wizard Genomic DNA extraction kit
704 (Promega). gDNA was then fragmented to 400-600 bp by sonication (Covaris E220)
705 and end repaired (Quick Blunting Kit, NEB). Transposon junctions were amplified from
706 gDNA by PCR. PCR products were gel purified to isolate 200-500bp fragments. To
707 estimate input and ensure equal multiplexing in downstream sequencing, purified PCR
708 products were subjected to qPCR using primers against the Illumina P5 and P7
709 hybridization sequence. Equimolar DNA fragments for each library were combined and
710 sequenced with a MiSeq.

711
712 Reads were first trimmed of transposon and adaptor sequences using CLC Genomics
713 Workbench (QIAGEN) and then mapped to *Escherichia coli* O157:H7 strain EDL933
714 (NCBI Accession Numbers: chromosome, NZ_CP008957.1; pO157 plasmid,
715 NZ_CP008958.1) using Bowtie without allowing mismatches. Reads were discarded if
716 they did not align to any TA sites, and reads that mapped to multiple TA sites were
717 randomly distributed between the multiple sites. After mapping, sensitivity analysis was
718 performed on each library to ensure adequate sequencing depth by sub-sampling reads
719 and assessing how many unique transposon mutants were detected (Fig S2). Next, the
720 data was normalized for chromosomal replication biases and differences in sequencing
721 depth using a LOESS correction of 100,000-bp and 10,000-bp windows for the
722 chromosome and plasmid, respectively. The number of reads at each TA site was
723 tallied and binned by gene and the percentage of disrupted TA sites was calculated.
724 Genes were binned by percentage of TA sites disrupted (Fig 1A, 1C).

725
726 For essential gene analysis, EL-ARTIST was used as in (45). Protein-coding genes,
727 RNA-coding genes, and pseudogenes were included in this analysis. Briefly, EL-
728 ARTIST classifies genes into one of three categories (underrepresented, regional, or
729 neutral), based on their transposon-insertion profile. Classifications are obtained using a
730 hidden Markov model (HMM) analysis following sliding window (SW) training ($p < 0.05$,
731 10 TA sites). Insertion-profiles for example genes were visualized with Artemis.

732

733 For identification of mutants conditionally depleted in the rabbit colon as compared to
734 the input inoculum, Con-ARTIST was used as in (114). First, the input library was
735 normalized to simulate the severity of the bottleneck as observed in the libraries
736 recovered from rabbit colons using multinomial distribution-based random sampling
737 ($n=100$). Next, a modified version of the Mann-Whitney U (MWU) function was applied
738 to compare these 100 simulated control data sets to the libraries recovered from the
739 rabbit colon. All genes were analyzed, but classification as “conditionally depleted” was
740 restricted to genes that had sufficient data (≥ 5 informative TA sites), met our standard of
741 attenuation (mean \log_2 fold change ≤ -2), met our standard of phenotypic consistency
742 (MWU p-value of ≤ 0.05), and had a consensus classification in 5 or more of the 7
743 animals analyzed. Genes with ≥ 5 informative TA sites that fail to exceed both standards
744 of attenuation and consistency are classified as “queried” (blue), whereas genes with
745 less than 5 informative TA sites are classified as “insufficient data”.

746
747 Gene-level PCA (gIPCA) was performed using CompTIS, a principal component
748 analysis-based TIS pipeline, as described in (66). Briefly, \log_2 fold change values were
749 derived by comparing read abundance in each sample to 100 control-simulated
750 datasets as in Con-ARTIST. These fold change values were weighted to minimize noise
751 due to variability (for details, see (66)). Next, genes that did not have a fold change
752 reported for all 7 animals were discarded. The fold change values were then z-score
753 normalized. Weighted PCA was performed in Matlab (Mathworks) with the PCA
754 algorithm (`pca`).

755

756 **GC content**

757 The GC content of classified genes was compared using a Mann-Whitney U statistical
758 test and a Bonferroni correction for multiple hypothesis correction when more than one
759 comparison was made. A p-value <0.05 was considered significant for one comparison,
760 p<0.025 for two. A Fisher's exact two-tailed t-test was used to compare ratios of
761 classifications between groups, where a p-value of <0.01 was considered significant.

762

763 **In vivo competitive infection**

764 Barcodes were introduced into $\Delta lacI::aacC1$ and isogenic mutant strains as described
765 previously (45,115) (46,63) Briefly, a 991bp fragment of *cynX* (RS02015) that included
766 51bp of the intergenic region between *cynX* and *lacA* (RS02020) was amplified using
767 primers that contained a 30 bp stretch of random sequence and cloned into *SacI* and
768 *XbaI* digested pGP704. The resulting pSoA176.mix was transformed into *E. coli*
769 MFD λ pir. Individual colonies carrying unique tag sequences were isolated and used as
770 donors to deliver pSoA176 barcoded derivatives to EDL933 $\Delta lacI::aacC1$ and each
771 isogenic mutant strain. Three barcodes were independently integrated into EDL933
772 $\Delta lacI::aacC1$, and three barcodes into each isogenic mutant via homologous
773 recombination in the intergenic region between *cynX* and *lacA*, which tolerates
774 transposon-insertion in vitro and in vivo, indicating this locus is neutral for the fitness of
775 the bacteria. Correct insertion of barcodes was confirmed by PCR and sequencing.

776

777 To prepare the culture of mixed EHEC-barcoded strains for the multi-coinfection
778 experiment, 100 μ l of overnight cultures of the barcoded strains were mixed in a flask

779 and 1 mL of this mix was added to 20 mL LB. After growing the culture for 3 h at 37°C
780 with shaking, the OD₆₀₀ was measured and 40 units of culture at OD₆₀₀=1 (about 8 mL)
781 were pelleted and resuspended in 10 mL PBS. Dilutions of the inoculum were plated in
782 LB agar plates with Gent and Carb for precise cfu determination. 10 infant rabbits were
783 inoculated and monitored as described above, and colon samples collected. Tissue
784 homogenate was plated, and CFU were collected the following day. gDNA was
785 extracted and prepared for sequencing as in (115).

786

787 The quantification of sequence tags was done as described by (115). In brief, sequence
788 tags were amplified from the inoculum culture and libraries recovered from rabbit
789 colons. The relative in vivo fitness of each mutant was assessed by calculating the
790 competitive index (CI) as follows.

791

792 We compare two strains ($\Delta lacI::aacC1$ and isogenic mutant) in a population with
793 frequencies f_{wt} and $f_{mut,x}$ respectively where x is one of 17 mutant strains with a deletion
794 in gene x. For simplicity, we assume here that both expand exponentially from a time
795 point t_0 to a sampling time point t_s , their relative fitness (offspring/generation) is

796 proportional to the competitive index CI: $\ln\left(\frac{f_{mut,x,s}/f_{mut,x,0}}{f_{wt,s}/f_{wt,0}}\right) = \ln(CI)$. Here, $f_{wt,0}$ and

797 $f_{mut,x,0}$ are the frequencies of the strains in the inoculum, measured in triplicates, and $f_{wt,s}$
798 and $f_{mut,x,s}$ describe the frequencies at the sampling time point in the animal host.

799 Because the WT strain was tagged with 3 individual tags and the inoculum was

800 measured in triplicate, we have $3 \times 3 = 9$ measurements of the ratio $f_{wt,s} / f_{wt,0}$. The same
801 is true for all mutant strains, such that we have 9 measurements of the ratio
802 $f_{mut,x,s} / f_{mut,x,0}$. In total, we therefore have $3 \times 3 \times 3 \times 3 = 81$ CI measurements for each
803 mutant per animal. To determine intra-host variance in these 81 measurements, a 95%
804 confidence interval of the CI in single animal hosts was determined by bootstrapping.
805 For combining the CIs measured across all 10 animal hosts, we performed a random-
806 effects meta-analysis using the metafor package (116) in the statistical software
807 package R (version 3.0.2). The pooled rate proportions and 95% confidence intervals
808 were calculated using the estimates and the variance of CIs in each animal determined
809 by bootstrapping and corrected for multiple testing using the Benjamini-Hochberg
810 procedure.

811

812 **In vitro growth**

813 Each bacterial strain was grown at 37°C overnight. The next day, cultures were diluted
814 1:1000 into 100 uL of LB in 96-well growth curve plates in triplicate. Plates were left
815 shaking at 37°C for 10-24 hours. Absorbance readings at 600nm were normalized to a
816 blank, and the average of each triplicate was taken as the optical density.

817

818 **T3SS translocation assays**

819 T3SS functionality was assessed by translocation of the known EHEC T3SS effector
820 protein EspF into HeLa cells as described previously (77). Briefly, the plasmid encoding
821 the effector protein EspF fused to TEM-1 beta-lactamase was transformed into each of
822 the bacterial strains to be tested. Overnight cultures of each bacterial strain were diluted

823 1/50 in DMEM supplemented with HEPES (25mM), 10% FBS and L-glutamine (2mM)
824 and incubated statically at 37°C with 5% CO₂ for two hours. This media is known to
825 induce T3SS expression (117). HeLa cells were seeded at a density of 2x10⁴ cells in
826 96-well clear bottom black plates and infected for 30 minutes at an MOI of 100. After 30
827 minutes of infection IPTG was added at a final concentration of 1mM to induce the
828 plasmid-encoded T3SS effector. After an additional hour of incubation, monolayers
829 were washed in HBSS solution and loaded with fluorescent substrate CCF2/AM solution
830 (Invitrogen) as recommended by the manufacturer. After 90 minutes, fluorescence was
831 quantified in a plate fluorescence reader with excitation at 410nm and emission was
832 detected at 450nm. Translocation was expressed as the emission ratio at 450/520nm to
833 normalize beta-lactamase activity to cell loading and the number of cells presented at
834 each well, and then normalized to WT levels of translocation.

835

836 **Biofilm, curli production, and purine assays**

837 Biofilm and curli production assays were performed as described previously (104). For
838 biofilm assays, bacterial cultures were grown in yeast extract-Casamino Acids (YESCA)
839 medium until they reached an OD₆₀₀ ~ 0.5 and 1/1000 dilution of this culture was used
840 to seed 96-well PVC plates. The cultures were grown at 30°C for 48 hours and biofilm
841 production was quantitatively measured using crystal violet staining and absorbance
842 reading at 595nm. Relative biofilm production was normalized to the average of three
843 WT samples. To test curli production, bacterial cultures were grown in YESCA medium
844 until they reached an OD₆₀₀ ~ 0.5 and then were struck to single colonies onto YESCA
845 agar plates supplemented with Congo Red. Red colonies indicate curli production. To

846 test if our $\Delta cvpA$ deletion had polar effects on *purF*, the mutant and WT were struck
847 onto minimal media lacking exogenous purines.

848

849 **Acid shock assays**

850 An adaptation of the acid shock method described in (118) was performed. Briefly,
851 bacterial cultures were grown until mid-exponential phase ($OD_{600} \sim 0.6$), then diluted 20-
852 fold in LB pH 5.5 and incubated for 1 hour before preparing serial dilutions and plating
853 each culture to determine the relative percentage of survival in comparison to the wild-
854 type EDL933 strain. The pH of the LB broth was adjusted using sterilized 1mM HCl and
855 buffered with 10% MES. Values are expressed as percent survival normalized to WT.

856

857 **MIC assays**

858 MIC assays were performed using an adaptation of a standard methodology with
859 exponential-phase cultures (119). Briefly, the different compounds to be tested (see
860 Fig5B, 6B) were prepared in serial 2-fold dilutions in 50 ul of LB in broth in a 96-well
861 plate format. To each well was added 50 ul of a culture prepared by diluting an
862 overnight culture 1,000-fold into fresh LB broth, growing it for 1 h at 37°C, and again
863 diluting it 1,000-fold into fresh medium. The plates were then incubated without shaking
864 for 24 h at 37°C.

865

866 **Bile salts survival assays**

867 Bile salt sensitivity assays were adapted from (120). For plate sensitivity assays, each
868 bacterial strain was grown at 37°C until they reached mid-exponential phase of growth

869 (OD_{600nm} of 0.5) and the culture was serially diluted and spot-titered onto LB agar plates
870 supplemented with either 1% DOC or 1% CHO. Spots were air dried and plates
871 incubated at 37°C for 24 h. For complementation, strains were grown in media and on
872 plates supplemented with 0.2% arabinose. For sensitivity assays done in liquid culture,
873 each bacterial strain was grown at 37°C until it reached mid-exponential phase of
874 growth (OD_{600nm} of 0.5) and then cultures were split and supplemented with either DOC,
875 CHO or buffer (PBS) and bacterial growth was assessed by absorbance at 600nm.

876

877 **Growth in high-salt media**

878 Bacterial strains were grown in either LB or LB supplemented with 0.3M NaCl until mid-
879 exponential phase and analyzed by phase microscopy.

880

881 **Computational Analysis**

882 To enable comprehensive functional/pathway analyses in EHEC we carried out BLAST-
883 based comparisons between the old EHEC genome sequence and annotation system
884 (NCBI Accession Numbers AE005174 and AF074613) and the new sequence and
885 annotation system (NZ_CP008957.1 and NZ_CP008958.1) (Table S1). This
886 comparison links the new annotations (RS locus tags) to the original 'Z numbers' from
887 (7) and their associated function and pathway annotation.

888

889 To make the correspondence table (Table S1) between the old EHEC annotation
890 system (Z Numbers) and the new system (RS Numbers), local BLAST was used. First,
891 a reference nucleotide database was generated from the newest EHEC sequence and

892 annotation (NZ_CP008957.1 and NZ_CP008958.1). The EHEC genome sequence
893 containing Z number annotations (AE005174 and AF074613) was used as the query.
894 Best matches were taken as equivalent loci.

895

896 To find the K-12 homolog for EHEC genes (Table S2, column M), local BLAST was also
897 used. A reference nucleotide and amino acid database was generated from MG1655 K-
898 12 (NC_000913.3), and the newest EHEC genome sequence was used as the query.
899 For pseudogenes and genes coding for RNA, $\geq 90\%$ nucleotide identity across $\geq 90\%$ of
900 the gene length was considered a homolog. For protein coding genes, $\geq 90\%$ amino acid
901 identity across $\geq 90\%$ of the amino acid sequence was considered a homolog.

902

903 To find KEGG pathways and COG assignments for genes of interest, the Z
904 correspondence table was used to look up the Z number of each gene. The Z number
905 and corresponding functional information was searched on the EHEC KEGG database.

906

907 To determine if COGs were enriched in certain groups of genes (such as conditionally
908 depleted genes), a COG enrichment index was calculated as in (51). The COG
909 Enrichment Index is the percentage of the genes of a certain category (essential genes
910 or CD genes) assigned to a specific COG divided by the percentage of genes in that
911 COG in the entire genome. A two-tailed Fisher's exact test was used to determine if this
912 ratio was independent of grouping. A Bonferroni correction was applied for multiple
913 hypothesis testing. A p-value of < 0.002 was considered to be significant.

914

915 Sequencing saturation of TIS libraries was determined by randomly sampling 100,000
916 reads from each library and identifying the number of unique mutants in that pool.
917 Libraries are sequenced to saturation when no new mutants are identified as additional
918 reads are added. 2-4 million reads are sufficient to capture the depth of libraries used
919 here.

920

921 Several protein prediction programs (PSLPred, HHPred, Phobius, Phyre2) (121–123)
922 were used to analyze the CvpA amino acid sequence. Protter (124) was used to
923 compile information from several of these searches and generate a topological diagram.
924 PRED-TAT (125) was used to search for tat-secretion signals in the list of CD genes.

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950 **References**

- 951
- 952 1. World Health Organization. Shiga toxin-producing *Escherichia coli* (STEC) and
953 food: attribution, characterization, and monitoring. Microbiological Risk
954 Assessment Series. 2018.
- 955 2. Hartland EL, Leong JM. Enteropathogenic and enterohemorrhagic *E. coli*:
956 ecology, pathogenesis, and evolution. *Frontiers in Cellular and Infection*
957 *Microbiology*. 2015;3(15):1–3.
- 958 3. Davis KT, Van De Kar NC, Tarr PI. Shiga Toxin/Verocytotoxin-Producing
959 *Escherichia coli* Infections: Practical Clinical Perspectives. *Enterohemorrhagic*
960 *Escherichia coli* and Other Shiga Toxin-Producing *E. coli*. 2014;321–39.
- 961 4. Karmali MA, Petric M, Steele BT, Lim C. Sporadic Cases of Haemolytic-Uraemic
962 Syndrome Associated with Faecal Cytotoxin-producing *Escherichia coli* in stools.
963 *The Lancet*. 1983 Mar 19;321(8325):619–20.
- 964 5. Boyce TG, Swerdlow DL, Griffin PM. *Escherichia coli* O157:H7 and the
965 Hemolytic–Uremic Syndrome. *New England Journal of Medicine*.
966 1995;333(6):364–8.
- 967 6. Riley L, Remis R, Helgerson S, McGee H, Wells J, Davis B, et al. Hemorrhagic
968 colitis associated with a rare *Escherichia coli* serotype. *New England Journal of*
969 *Medicine*. 1983;308(12):681–5.
- 970 7. Perna NT, Plunkett G, Burland V, Mau B, Glasner JD, Rose DJ, et al. Genome
971 sequence of enterohaemorrhagic *Escherichia coli* O157:H7. *Nature*.
972 2001;409(6819):529–33.
- 973 8. Latif H, Li HJ, Charusanti P, Palsson BO, Aziz RK. A Gapless, Unambiguous
974 Genome Sequence of the Enterohemorrhagic *Escherichia coli* O157:H7 Strain
975 EDL933. *Genome Announcements*. 2014;2(4).
- 976 9. Blattner F, Plunkett G, Bloch CA, Perna N, Burland V, Riley M, et al. The
977 Complete Genome Sequence of *Escherichia coli* K-12. *Science*.
978 1997;277(5331):1453–62.
- 979 10. Ogura Y, Ooka T, Iguchi A, Toh H, Asadulghani M, Oshima K, et al. Comparative
980 genomics reveal the mechanism of the parallel evolution of O157 and non-O157
981 enterohemorrhagic *Escherichia coli*. *Proceedings of the National Academy of*
982 *Sciences*. 2009;106(42):17939–44.
- 983 11. Delannoy S, Beutin L, Fach P. Discrimination of Enterohemorrhagic *Escherichia*
984 *coli* (EHEC) from Non-EHEC Strains Based on Detection of Various Combinations
985 of Type III Effector Genes. *Journal of Clinical Microbiology*. 2013;51(10):3257–62.

- 986 12. Ritchie JM, Thorpe CM, Rogers AB, Waldor MK. Critical Roles for *stx2*, *eae*, and
987 *tir* in Enterohemorrhagic *Escherichia coli*-Induced Diarrhea and Intestinal
988 Inflammation in Infant Rabbits. *Infection and Immunity*. 2003;71(12):7129–39.
- 989 13. Nguyen Y, Sperandio V. Enterohemorrhagic *E. coli* (EHEC) pathogenesis.
990 *Frontiers in Cellular and Infection Microbiology*. 2012;2(90):1–7.
- 991 14. Wong ARC, Pearson JS, Bright MD, Munera D, Robinson KS, Lee SF, et al.
992 Enteropathogenic and enterohaemorrhagic *Escherichia coli*: even more
993 subversive elements. *Molecular Microbiology*. 2011;80(6):1420–38.
- 994 15. Connolly JPR, Finlay BB, Roe AJ. From ingestion to colonization: the influence of
995 the host environment on regulation of the LEE encoded type III secretion system
996 in enterohaemorrhagic *Escherichia coli*. *Frontiers in Microbiology*. 2015;6:568.
- 997 16. Donnenberg MS, Tzipori S, McKee ML, O'Brien AD, Alroy J, Kaper JB. The role of
998 the *eae* gene of enterohemorrhagic *Escherichia coli* in intimate attachment in vitro
999 and in a porcine model. *Journal of Clinical Investigation*. 1993;92(3):1418–24.
- 1000 17. Mckee ML, O'Brien AD. Investigation of Enterohemorrhagic *Escherichia coli*
1001 O157:H7 Adherence Characteristics and Invasion Potential Reveals a New
1002 Attachment Pattern Shared by Intestinal *E. coli*. *Infection and Immunity*.
1003 1995;63(5):2070–4.
- 1004 18. Tzipori S, Gunzer F, Donnenberg MS, Montigny LD, Kaper JB, Donohue-Rolfe A.
1005 The Role of the *eaeA* Gene in Diarrhea and Neurological Complications in a
1006 Gnotobiotic Piglet Model of Enterohemorrhagic *Escherichia coli* Infection. *Infection*
1007 and *Immunity*. 1995;63:7.
- 1008 19. Dean-Nystrom EA, Bosworth BT, Moon HW, O'Brien AD. *Escherichia coli*
1009 O157:H7 Requires Intimin for Enteropathogenicity in Calves. *Infection and*
1010 *Immunity*. 1998;66(0):4560–3.
- 1011 20. Cornick N, Booher S, Moon H. Intimin Facilitates Colonization by *Escherichia coli*
1012 O157:H7 in Adult Ruminants. *Infection and Immunity*. 2002;70(5):2704–7.
- 1013 21. Barnett Foster D. Modulation of the enterohemorrhagic *E. coli* virulence program
1014 through the human gastrointestinal tract. *Virulence*. 2013;4(4):315–23.
- 1015 22. Hughes DT, Clarke MB, Yamamoto K, Rasko DA, Sperandio V. The QseC
1016 Adrenergic Signaling Cascade in Enterohemorrhagic *E. coli* (EHEC). *PLoS*
1017 *Pathogens*. 2009;5(8):e1000553.
- 1018 23. Sperandio V, Torres AG, Jarvis B, Nataro JP, Kaper JB. Bacteria-host
1019 communication: The language of hormones. *Proceedings of the National*
1020 *Academy of Sciences*. 2003;100(15):8951–6.

- 1021 24. Hirakawa H, Kodama T, Takumi-Kobayashi A, Honda T, Yamaguchi A. Secreted
1022 indole serves as a signal for expression of type III secretion system translocators
1023 in enterohaemorrhagic *Escherichia coli* O157:H7. *Microbiology*. 2009;155(2):541–
1024 50.
- 1025 25. Kendall M. Interkingdom Chemical Signaling in Enterohemorrhagic *Escherichia*
1026 *coli* O157:H7. In: Lyte M, editor. *Microbial Endocrinology: Interkingdom Signaling*
1027 *in Infectious Disease and Health*. 2nd edition. Springer International Publishing.
1028 2016. p.p. 201–213.
- 1029 26. van Opijnen T, Bodi KL, Camilli A. Tn-seq: high-throughput parallel sequencing for
1030 fitness and genetic interaction studies in microorganisms. *Nature Methods*.
1031 2009;6(10):767–72.
- 1032 27. Goodman AL, McNulty NP, Zhao Y, Leip D, Mitra RD, Lozupone CA, et al.
1033 Identifying Genetic Determinants Needed to Establish a Human Gut Symbiont in
1034 Its Habitat. *Cell Host & Microbe*. 2009;6(3):279–89.
- 1035 28. Gawronski JD, Wong SMS, Giannoukos G, Ward DV, Akerley BJ. Tracking
1036 insertion mutants within libraries by deep sequencing and a genome-wide screen
1037 for *Haemophilus* genes required in the lung. *Proceedings of the National*
1038 *Academy of Sciences*. 2009;106(38):16422–7.
- 1039 29. Langridge GC, Phan M-D, Turner DJ, Perkins TT, Parts L, Haase J, et al.
1040 Simultaneous assay of every *Salmonella Typhi* gene using one million transposon
1041 mutants. *Genome Research*. 2009;19(12):2308–16.
- 1042 30. Gao B, Vorwerk H, Huber C, Lara-Tejero M, Mohr J, Goodman AL, et al.
1043 Metabolic and fitness determinants for in vitro growth and intestinal colonization of
1044 the bacterial pathogen *Campylobacter jejuni*. *PLOS Biology*.
1045 2017;15(5):e2001390.
- 1046 31. Cole BJ, Feltcher ME, Waters RJ, Wetmore KM, Mucyn TS, Ryan EM, et al.
1047 Genome-wide identification of bacterial plant colonization genes. *PLOS Biology*.
1048 2017;15(9):e2002860.
- 1049 32. McCarthy AJ, Stabler RA, Taylor PW. Genome-Wide Identification by Transposon
1050 Insertion Sequencing of *Escherichia coli* K1 Genes Essential for *In Vitro* Growth,
1051 Gastrointestinal Colonizing Capacity, and Survival in Serum. *Journal of*
1052 *Bacteriology*. 2018;200(7).
- 1053 33. de Moraes MH, Desai P, Porwollik S, Canals R, Perez DR, Chu W, et al.
1054 *Salmonella* Persistence in Tomatoes Requires a Distinct Set of Metabolic
1055 Functions Identified by Transposon Insertion Sequencing. *Applied and*
1056 *Environmental Microbiology* 2017;83(5).
- 1057 34. Armbruster CE, Forsyth-DeOrnellas V, Johnson AO, Smith SN, Zhao L, Wu W, et
1058 al. Genome-wide transposon mutagenesis of *Proteus mirabilis*: Essential genes,

- 1059 fitness factors for catheter-associated urinary tract infection, and the impact of
1060 polymicrobial infection on fitness requirements. PLOS Pathogens.
1061 2017;13(6):e1006434.
- 1062 35. Capel E, Barnier J-P, Zomer AL, Bole-Feysot C, Nussbaumer T, Jamet A, et al.
1063 Peripheral blood vessels are a niche for blood-borne meningococci. Virulence.
1064 2017;8(8):1808–19.
- 1065 36. Zhu L, Charbonneau ARL, Waller AS, Olsen RJ, Beres SB, Musser JM. Novel
1066 Genes Required for the Fitness of *Streptococcus pyogenes* in Human Saliva.
1067 2017;2(6):13.
- 1068 37. George AS, Cox CE, Desai P, Porwolik S, Chu W, de Moraes MH, et al.
1069 Interactions of *Salmonella enterica* Serovar Typhimurium and *Pectobacterium*
1070 *carotovorum* within a Tomato Soft Rot. Applied and Environmental Microbiology.
1071 2017;84(5).
- 1072 38. White KM, Matthews MK, Hughes R, Sommer AJ, Griffiths JS, Newell PD, et al. A
1073 Metagenome-Wide Association Study and Arrayed Mutant Library Confirm
1074 *Acetobacter* Lipopolysaccharide Genes Are Necessary for Association with
1075 *Drosophila melanogaster*. G3: Genes, Genomes, Genetics.
1076 2018;g3.300530.2017.
- 1077 39. Price MN, Wetmore KM, Waters RJ, Callaghan M, Ray J, Liu H, et al. Mutant
1078 phenotypes for thousands of bacterial genes of unknown function. Nature.
1079 2018;557(7706):503–9.
- 1080 40. Pritchard JR, Chao MC, Abel S, Davis BM, Baranowski C, Zhang YJ, et al.
1081 ARTIST: high-resolution genome-wide assessment of fitness using transposon-
1082 insertion sequencing. PLoS Genetics. 2014;10(11):e1004782.
- 1083 41. Chao MC, Abel S, Davis BM, Waldor MK. The design and analysis of transposon
1084 insertion sequencing experiments. Nature Reviews Microbiology. 2016;14(2):119–
1085 28.
- 1086 42. Kimura S, Hubbard TP, Davis BM, Waldor MK. The Nucleoid Binding Protein H-
1087 NS Biases Genome-Wide Transposon Insertion Landscapes. mBio. 2016;7(4).
- 1088 43. DeJesus MA, Gerrick ER, Xu W, Park SW, Long JE, Boutte CC, et al.
1089 Comprehensive Essentiality Analysis of the *Mycobacterium tuberculosis* Genome
1090 via Saturating Transposon Mutagenesis. mBio. 2017;8(1).
- 1091 44. Ritchie JM, Waldor MK. The Locus of Enterocyte Effacement-Encoded Effector
1092 Proteins All Promote Enterohemorrhagic *Escherichia coli* Pathogenicity in Infant
1093 Rabbits. Infection and Immunity. 2005;73(3):1466–74.

- 1094 45. Hubbard TP, Chao MC, Abel S, Blondel CJ, Abel zur Wiesch P, Zhou X, et al.
1095 Genetic analysis of *Vibrio parahaemolyticus* intestinal colonization. Proceedings
1096 of the National Academy of Sciences. 2016;113(22):6283–8.
- 1097 46. Chaudhuri RR, Morgan E, Peters SE, Pleasance SJ, Hudson DL, Davies HM, et
1098 al. Comprehensive Assignment of Roles for Salmonella Typhimurium Genes in
1099 Intestinal Colonization of Food-Producing Animals. PLoS Genetics.
1100 2013;9(4):e1003456.
- 1101 47. Munoz-Lopez M, Garcia-Perez J. DNA Transposons: Nature and Applications in
1102 Genomics. Current Genomics. 2010;11(2):115–28.
- 1103 48. Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M. KEGG as a
1104 reference resource for gene and protein annotation. Nucleic Acids Research.
1105 2016;44(D1):D457–62.
- 1106 49. Tatusov RL. The COG database: a tool for genome-scale analysis of protein
1107 functions and evolution. Nucleic Acids Research. 2000;28(1):33–6.
- 1108 50. Galperin MY, Makarova KS, Wolf YI, Koonin EV. Expanded microbial genome
1109 coverage and improved protein family annotation in the COG database. Nucleic
1110 Acids Research. 2015;43(D1):D261–9.
- 1111 51. Shields RC, Zeng L, Culp DJ, Burne RA. Genomewide Identification of Essential
1112 Genes and Fitness Determinants of *Streptococcus mutans* UA159. mSphere.
1113 2018;3(1):e00031-18.
- 1114 52. Luo H, Lin Y, Gao F, Zhang C-T, Zhang R. DEG 10, an update of the database of
1115 essential genes that includes both protein-coding genes and noncoding genomic
1116 elements. Nucleic Acids Research. 2014;42(D1):D574–80.
- 1117 53. Baba T, Ara T, Hasegawa M, Takai Y, Okumura Y, Baba M, et al. Construction of
1118 *Escherichia coli* K-12 in-frame, single-gene knockout mutants: the Keio collection.
1119 Molecular Systems Biology. 2006;2:2006.0008.
- 1120 54. Kato J, Hashimoto M. Construction of consecutive deletions of the *Escherichia coli*
1121 chromosome. Molecular Systems Biology. 2007;3:132.
- 1122 55. Yamamoto N, Nakahigashi K, Nakamichi T, Yoshino M, Takai Y, Touda Y, et al.
1123 Update on the Keio collection of *Escherichia coli* single-gene deletion mutants.
1124 Molecular Systems Biology. 2009;5:335.
- 1125 56. Sperandio P, Pozzi C, Dehò G, Polissi A. Non-essential KDO biosynthesis and
1126 new essential cell envelope biogenesis genes in the *Escherichia coli* yrbG–yhbG
1127 locus. Research in Microbiology. 2006;157(6):547–58.
- 1128 57. Goodall ECA, Robinson A, Johnston IG, Jabbari S, Turner KA, Cunningham AF,
1129 et al. The Essential Genome of *Escherichia coli* K-12. 2018;9(1):18.

- 1130 58. Abel S, Waldor MK. Infant Rabbit Model for Diarrheal Diseases: Infant Rabbit
1131 Model for Diarrheal Diseases. *Current Protocols in Microbiology*. 2015;6A.6.1-
1132 6A.6.15.
- 1133 59. Ritchie JM, Rui H, Bronson RT, Waldor MK. Back to the Future: Studying Cholera
1134 Pathogenesis Using Infant Rabbits. *mBio*. 2010;1(1).
- 1135 60. Abel S, Abel zur Wiesch P, Davis BM, Waldor MK. Analysis of Bottlenecks in
1136 Experimental Models of Infection. *PLOS Pathogens*. 2015;11(6):e1004823.
- 1137 61. Fu Y, Waldor MK, Mekalanos JJ. Tn-Seq Analysis of *Vibrio cholerae* Intestinal
1138 Colonization Reveals a Role for T6SS-Mediated Antibacterial Activity in the Host.
1139 *Cell Host & Microbe*. 2013;14(6):652–63.
- 1140 62. Abu Kwaik Y, Bumann D. Host Delivery of Favorite Meals for Intracellular
1141 Pathogens. *PLOS Pathogens*. 2015;11(6):e1004866.
- 1142 63. Meylan S, Porter CBM, Yang JH, Belenky P, Gutierrez A, Lobritz MA, et al.
1143 Carbon Sources Tune Antibiotic Susceptibility in *Pseudomonas aeruginosa* via
1144 Tricarboxylic Acid Cycle Control. *Cell Chemical Biology*. 2017;24(2):195–206.
- 1145 64. Murima P, McKinney JD, Pethe K. Targeting Bacterial Central Metabolism for
1146 Drug Development. *Chemistry & Biology*. 2014;21(11):1423–32.
- 1147 65. Ren W, Rajendran R, Zhao Y, Tan B, Wu G, Bazer FW, et al. Amino Acids As
1148 Mediators of Metabolic Cross Talk between Host and Pathogen. *Frontiers in*
1149 *Immunology*. 2018;9:319.
- 1150 66. Hubbard TP, D’Gama JD, Billings G, Davis BM, Waldor MK. Unsupervised
1151 Learning Approach for Comparing Multiple Transposon Insertion Sequencing
1152 Studies. 2019;4(1):13.
- 1153 67. Franzin FM, Sircili MP. Locus of Enterocyte Effacement: A Pathogenicity Island
1154 Involved in the Virulence of Enteropathogenic and Enterohemorrhagic *Escherichia*
1155 *coli* Subjected to a Complex Network of Gene Regulation. *BioMed Research*
1156 *International*. 2015;2015:1–10.
- 1157 68. Knutton S. A novel EspA-associated surface organelle of enteropathogenic
1158 *Escherichia coli* involved in protein translocation into epithelial cells. *The EMBO*
1159 *Journal*. 1998;17(8):2166–76.
- 1160 69. Kresse AU, Schulze K, Deibel C, Ebel F, Rohde M, Chakraborty T. Pas, a Novel
1161 Protein Required for Protein Secretion and Attaching and Effacing Activities of
1162 Enterohemorrhagic *Escherichia coli*. *Journal of Bacteriology*. 1998;180:10.
- 1163 70. Kresse AU, Rohde M. The EspD Protein of Enterohemorrhagic *Escherichia coli* Is
1164 Required for the Formation of Bacterial Surface Appendages and Is Incorporated

- 1165 in the Cytoplasmic Membranes of Target Cells. *Infection and Immunity*.
1166 1999;67:9.
- 1167 71. Tatsuno I, Kimura H, Okutani A, Kanamaru K, Abe H, Nagai S, et al. Isolation and
1168 Characterization of Mini-Tn5Km2 Insertion Mutants of Enterohemorrhagic
1169 *Escherichia coli* O157:H7 Deficient in Adherence to Caco-2 Cells. *Infection and*
1170 *Immunity*. 2000;68(10):5943–52.
- 1171 72. Sal-Man N, Setiaputra D, Scholz R, Deng W, Yu ACY, Strynadka NCJ, et al. EscE
1172 and EscG Are Cochaperones for the Type III Needle Protein EscF of
1173 Enteropathogenic *Escherichia coli*. *Journal of Bacteriology*. 2013;195(11):2481–9.
- 1174 73. Gruenheid S, Sekirov I, Thomas NA, Deng W, O'Donnell P, Goode D, et al.
1175 Identification and characterization of NleA, a non-LEE-encoded type III
1176 translocated virulence factor of enterohaemorrhagic *Escherichia coli* O157:H7:
1177 Identification and characterization of NleA. *Molecular Microbiology*.
1178 2004;51(5):1233–49.
- 1179 74. Yen H, Sugimoto N, Tobe T. Enteropathogenic *Escherichia coli* Uses NleA to
1180 Inhibit NLRP3 Inflammasome Activation. *PLOS Pathogens*. 2015;11(9):e1005121.
- 1181 75. Arbeloa A, Bulgin RR, MacKenzie G, Shaw RK, Pallen MJ, Crepin VF, et al.
1182 Subversion of actin dynamics by EspM effectors of attaching and effacing
1183 bacterial pathogens. *Cellular Microbiology*. 2008;10(7):1429–41.
- 1184 76. Simovitch M, Sason H, Cohen S, Zahavi EE, Melamed-Book N, Weiss A, et al.
1185 EspM inhibits pedestal formation by enterohaemorrhagic *Escherichia coli* and
1186 enteropathogenic *E. coli* and disrupts the architecture of a polarized epithelial
1187 monolayer. *Cellular Microbiology*. 2010;12(4):489–505.
- 1188 77. Munera D, Crepin VF, Marches O, Frankel G. N-Terminal Type III Secretion
1189 Signal of Enteropathogenic *Escherichia coli* Translocator Proteins. *Journal of*
1190 *Bacteriology*. 2010;192(13):3534–9.
- 1191 78. Iyoda S, Watanabe H. ClpXP Protease Controls Expression of the Type III Protein
1192 Secretion System through Regulation of RpoS and GrlR Levels in
1193 Enterohemorrhagic *Escherichia coli*. *Journal of Bacteriology*. 2005;187(12):4086–
1194 94.
- 1195 79. Tomoyasu T, Takaya A, Handa Y, Karata K, Yamamoto T. ClpXP controls the
1196 expression of LEE genes in enterohaemorrhagic *Escherichia coli*. *FEMS*
1197 *Microbiology Letters*. 2005;253(1):59–66.
- 1198 80. MacRitchie DM, Ward JD, Nevesinjac AZ, Raivio TL. Activation of the Cpx
1199 Envelope Stress Response Down-Regulates Expression of Several Locus of
1200 Enterocyte Effacement-Encoded Genes in Enteropathogenic *Escherichia coli*.
1201 *Infection and Immunity*. 2008;76(4):1465–75.

- 1202 81. MacRitchie DM, Acosta N, Raivio TL. DegP Is Involved in Cpx-Mediated
1203 Posttranscriptional Regulation of the Type III Secretion Apparatus in
1204 Enteropathogenic *Escherichia coli*. *Infection and Immunity*. 2012;80(5):1766–72.
- 1205 82. Hara H, Yamamoto Y, Higashitani A, Suzuki H, Nishimura Y. Cloning, mapping,
1206 and characterization of the *Escherichia coli* *prc* gene, which is involved in C-
1207 terminal processing of penicillin-binding protein 3. *Journal of Bacteriology*.
1208 1991;173(15):4799–813.
- 1209 83. Kerr CH, Culham DE, Marom D, Wood JM. Salinity-Dependent Impacts of ProQ,
1210 Prc, and Spr Deficiencies on *Escherichia coli* Cell Structure. *Journal of*
1211 *Bacteriology*. 2014;196(6):1286–96.
- 1212 84. Pos KM. Drug transport mechanism of the AcrB efflux pump. *Biochimica et*
1213 *Biophysica Acta (BBA) - Proteins and Proteomics*. 2009;1794(5):782–93.
- 1214 85. Begley M, Gahan CGM, Hill C. The interaction between bacteria and bile. *FEMS*
1215 *Microbiology Reviews*. 2005;29(4):625–51.
- 1216 86. Christman MF, Storz G, Ames BN. OxyR, a positive regulator of hydrogen
1217 peroxide-inducible genes in *Escherichia coli* and *Salmonella typhimurium*, is
1218 homologous to a family of bacterial regulatory proteins. *Proceedings of the*
1219 *National Academy of Sciences*. 1989;86(10):3484–8.
- 1220 87. Daugherty A, Suvarnapunya AE, Runyen-Janecky L. The role of OxyR and
1221 SoxRS in oxidative stress survival in *Shigella flexneri*. *Microbiological Research*.
1222 2012;167(4):238–45.
- 1223 88. De Paepe M, Gaboriau-Routhiau V, Rainteau D, Rakotobe S, Taddei F, Cerf-
1224 Bensussan N. Trade-Off between Bile Resistance and Nutritional Competence
1225 Drives *Escherichia coli* Diversification in the Mouse Gut. *PLoS Genetics*.
1226 2011;7(6):e1002107.
- 1227 89. Berks BC. The Twin-Arginine Protein Translocation Pathway. *Annual Review of*
1228 *Biochemistry*. 2015;84(1):843–64.
- 1229 90. Green ER, Meccas J. Bacterial Secretion Systems: An Overview. *Microbiology*
1230 *Spectrum*. 2016;4(1).
- 1231 91. Mickael CS, Lam P-KS, Berberov EM, Allan B, Potter AA, Koster W. *Salmonella*
1232 *enterica* Serovar Enteritidis *tatB* and *tatC* Mutants Are Impaired in Caco-2 Cell
1233 Invasion In Vitro and Show Reduced Systemic Spread in Chickens. *Infection and*
1234 *Immunity*. 2010;78(8):3493–505.
- 1235 92. Craig M, Sadik AY, Golubeva YA, Tidhar A, Slauch JM. Twin-arginine
1236 translocation system (*tat*) mutants of *Salmonella* are attenuated due to envelope
1237 defects, not respiratory defects: Role of Tat in *Salmonella* virulence. *Molecular*
1238 *Microbiology*. 2013;89(5):887–902.

- 1239 93. Fujimoto M, Goto R, Hirota R, Ito M, Haneda T, Okada N, et al. Tat-exported
1240 peptidoglycan amidase-dependent cell division contributes to *Salmonella*
1241 *Typhimurium* fitness in the inflamed gut. *PLOS Pathogens*.
1242 2018;14(10):e1007391.
- 1243 94. Avican U, Doruk T, Östberg Y, Fahlgren A, Forsberg Å. The Tat Substrate SufI Is
1244 Critical for the Ability of *Yersinia pseudotuberculosis* To Cause Systemic Infection.
1245 *Infection and Immunity*. 2017;85(4).
- 1246 95. Rajashekara G, Drozd M, Gangaiyah D, Jeon B, Liu Z, Zhang Q. Functional
1247 Characterization of the Twin-Arginine Translocation System in *Campylobacter*
1248 *jejuni*. *Foodborne Pathogens and Disease*. 2009;6(8):935–45.
- 1249 96. Lavander M, Ericsson SK, Broms JE, Forsberg A. The Twin Arginine
1250 Translocation System Is Essential for Virulence of *Yersinia pseudotuberculosis*.
1251 *Infection and Immunity*. 2006;74(3):1768–76.
- 1252 97. Zhang L, Zhu Z, Jing H, Zhang J, Xiong Y, Yan M, et al. Pleiotropic effects of the
1253 twin-arginine translocation system on biofilm formation, colonization, and
1254 virulence in *Vibrio cholerae*. *BMC Microbiology*. 2009;9(1):114.
- 1255 98. Pradel N, Ye C, Livrelli V, Xu J, Joly B, Wu L-F. Contribution of the Twin Arginine
1256 Translocation System to the Virulence of Enterohemorrhagic *Escherichia coli*
1257 O157:H7. *Infection and Immunity*. 2003;71(9):4908–16.
- 1258 99. Samaluru H, SaiSree L, Reddy M. Role of SufI (FtsP) in Cell Division of
1259 *Escherichia coli*: Evidence for Its Involvement in Stabilizing the Assembly of the
1260 Divisome. *Journal of Bacteriology*. 2007;189(22):8044–52.
- 1261 100. Tarry M, Arends SJR, Roversi P, Piette E, Sargent F, Berks BC, et al. The
1262 *Escherichia coli* Cell Division Protein and Model Tat Substrate SufI (FtsP)
1263 Localizes to the Septal Ring and Has a Multicopper Oxidase-Like Structure.
1264 *Journal of Molecular Biology*. 2009;386(2):504–19.
- 1265 101. Stanley NR, Findlay K, Berks BC, Palmer T. *Escherichia coli* Strains Blocked in
1266 Tat-Dependent Protein Export Exhibit Pleiotropic Defects in the Cell Envelope.
1267 *Journal of Bacteriology*. 2001;183(1):139–44.
- 1268 102. Bernhardt TG, De Boer PAJ. Screening for synthetic lethal mutants in *Escherichia*
1269 *coli* and identification of EnvC (YibP) as a periplasmic septal ring factor with
1270 murein hydrolase activity: *E. coli* synthetic lethal screen. *Molecular Microbiology*.
1271 2004;52(5):1255–69.
- 1272 103. Fath MJ, Mahanty HK, Kolter R. Characterization of a *purF* operon mutation which
1273 affects colicin V production. *Journal of Bacteriology*. 1989;171(6):3158–61.

- 1274 104. Hadjifrangiskou M, Gu AP, Pinkner JS, Kostakioti M, Zhang EW, Greene SE, et al.
1275 Transposon Mutagenesis Identifies Uropathogenic *Escherichia coli* Biofilm
1276 Factors. *Journal of Bacteriology*. 2012;194(22):6195–205.
- 1277 105. Urdaneta V, Casadesús J. Interactions between Bacteria and Bile Salts in the
1278 Gastrointestinal and Hepatobiliary Tracts. *Frontiers in Medicine*. 2017;4(163).
- 1279 106. Shaffer CL, Zhang EW, Dudley AG, Dixon BREA, Guckes KR, Breland EJ, et al.
1280 Purine Biosynthesis Metabolically Constrains Intracellular Survival of
1281 Uropathogenic *Escherichia coli*. *Infection and Immunity*. 2017;85(1).
- 1282 107. Kwon YM, Ricke SC, Mandal RK. Transposon sequencing: methods and
1283 expanding applications. *Applied Microbiology and Biotechnology*. 2016;100(1):31–
1284 43.
- 1285 108. Barquist L, Boinett CJ, Cain AK. Approaches to querying bacterial genomes with
1286 transposon-insertion sequencing. *RNA Biology*. 2013;10(7):1161–9.
- 1287 109. Donnenberg MS, Kaper JB. Construction of an *eae* Deletion Mutant of
1288 Enteropathogenic *Escherichia coli* by Using a Positive-Selection Suicide Vector.
1289 *Infection and Immunity*. 1991;59:8.
- 1290 110. Poteete AR, Rosadini C, St. Pierre C. Gentamicin and other cassettes for
1291 chromosomal gene replacement in *Escherichia coli*. *BioTechniques*.
1292 2006;41(3):261–4.
- 1293 111. Ferrieres L, Hemery G, Nham T, Guerout A-M, Mazel D, Beloin C, et al. Silent
1294 Mischief: Bacteriophage Mu Insertions Contaminate Products of *Escherichia coli*
1295 Random Mutagenesis Performed Using Suicidal Transposon Delivery Plasmids
1296 Mobilized by Broad-Host-Range RP4 Conjugative Machinery. *Journal of*
1297 *Bacteriology*. 2010;192(24):6418–27.
- 1298 112. Cameron DE, Urbach JM, Mekalanos JJ. A defined transposon mutant library and
1299 its use in identifying motility genes in *Vibrio cholerae*. *Proceedings of the National*
1300 *Academy of Sciences*. 2008;105(25):8736–41.
- 1301 113. Chiang SL, Rubin EJ. Construction of a mariner -based transposon for epitope-
1302 tagging and genomic targeting. *Gene*. 2002;296(1–2):179–85.
- 1303 114. Hubbard TP, Billings G, Dörr T, Sit B, Warr AR, Kuehl CJ, et al. A live vaccine
1304 rapidly protects against cholera in an infant rabbit model. *Science Translational*
1305 *Medicine*. 2018;10(445):eaap8423.
- 1306 115. Abel S, Abel zur Wiesch P, Chang H-H, Davis BM, Lipsitch M, Waldor MK.
1307 Sequence tag-based analysis of microbial population dynamics. *Nature Methods*.
1308 2015;12(3):223–6.

- 1309 116. Viechtbauer W. Conducting Meta-Analyses in R with the metafor Package.
1310 Journal of Statistical Software. 2010;36(3).
- 1311 117. Collington GK, Booth IW, Knutton S. Rapid modulation of electrolyte transport in
1312 Caco-2 cell monolayers by enteropathogenic Escherichia coli (EPEC) infection.
1313 Gut. 1998;42(2):200–7.
- 1314 118. Stincone A, Daudi N, Rahman AS, Antczak P, Henderson I, Cole J, et al. A
1315 systems biology approach sheds new light on Escherichia coli acid resistance.
1316 Nucleic Acids Research. 2011;39(17):7512–28.
- 1317 119. Dörr T, Möll A, Chao MC, Cava F, Lam H, Davis BM, et al. Differential
1318 Requirement for PBP1a and PBP1b in In Vivo and In Vitro Fitness of Vibrio
1319 cholerae. Infection and Immunity. 2014;82(5):2115–24.
- 1320 120. Cremers CM, Knoefler D, Vitvitsky V, Banerjee R, Jakob U. Bile salts act as
1321 effective protein-unfolding agents and instigators of disulfide stress in vivo.
1322 Proceedings of the National Academy of Sciences. 2014;111(16):E1610–9.
- 1323 121. Käll L, Krogh A, Sonnhammer EL. A Combined Transmembrane Topology and
1324 Signal Peptide Prediction Method. Journal of Molecular Biology.
1325 2004;338(5):1027–36.
- 1326 122. Zimmermann L, Stephens A, Nam S-Z, Rau D, Kübler J, Lozajic M, et al. A
1327 Completely Reimplemented MPI Bioinformatics Toolkit with a New HHpred Server
1328 at its Core. Journal of Molecular Biology. 2018;430(15):2237–43.
- 1329 123. Kelley LA, Mezulis S, Yates CM, Wass MN, Sternberg MJE. The Phyre2 web
1330 portal for protein modeling, prediction and analysis. Nature Protocols.
1331 2015;10(6):845–58.
- 1332 124. Omasits U, Ahrens CH, Müller S, Wollscheid B. Protter: interactive protein feature
1333 visualization and integration with experimental proteomic data. Bioinformatics.
1334 2014;30(6):884–6.
- 1335 125. Bagos PG, Nikolaou EP, Liakopoulos TD, Tsirigos KD. Combined prediction of Tat
1336 and Sec signal peptides with hidden Markov models. Bioinformatics.
1337 2010;26(22):2811–7.

1338
1339
1340
1341
1342
1343
1344
1345
1346

1347 **Figure Legends**

1348

1349 **Figure 1: Analysis of essential genes in EHEC EDL933 and comparison to K-12**
1350 **MG1655.**

1351 A) Distribution of percentage TA site disruption for all genes in EHEC EDL933. Genes
1352 are classified by EL-ARTIST as either underrepresented (red), regional (purple), or
1353 neutral (blue).

1354

1355 B) Transposon-insertion profiles of representative underrepresented, regional, and
1356 neutral genes.

1357

1358 C) GC content (%) of EDL933 genes with and without homologs in K-12 MG1655,
1359 classified by TIS classification (neutral or non-neutral). Neutral and non-neutral genes
1360 within each gene type (divergent or homolog) are compared using a Mann-Whitney U
1361 test with a Bonferroni correction. Ratios of neutral to non-neutral genes for each gene
1362 type are compared using a Fisher's Exact Test; (**) indicates a p-value of <0.01 and
1363 (***) indicate a p-value of <0.001.

1364

1365 D) Distribution of percentage TA site disruption for all genes in K-12 MG1655. Genes
1366 are classified by EL-ARTIST as underrepresented (red), regional (purple), or neutral
1367 (blue) using the same parameters as for the EDL933 library.

1368

1369 E) Genes classified as non-neutral in the EDL933 TIS library were compared to the
1370 MG1655 TIS library and categorized as either lacking a homolog (green), having the
1371 same classification in both libraries (red), or being non-neutral in EDL933 and neutral in
1372 MG1655 (blue).

1373

1374 F) KEGG pathway information for genes that are non-neutral in EDL933 and neutral in
1375 MG1655.

1376

1377 **Figure 2: Identification of EHEC genes required for intestinal colonization.**

1378 A) Distribution of percentage TA site disruption for all genes in EDL933 in the library
1379 used to inoculate infant rabbits (top) and in a representative library recovered from a
1380 rabbit colon two days after infection (bottom).

1381

1382 B) Schematic of Con-ARTIST classification scheme. Con-ARTIST utilizes iterative
1383 resampling of the inoculum data set to generate 100 simulated control data sets and
1384 compares relative abundance of mutants in these simulated control data sets relative to
1385 the passaged library. Genes with sufficient data (≥ 5 TA sites disrupted) are then
1386 classified based on a dual standard of attenuation (mean \log_2 fold change ≤ -2) and
1387 consistency (Mann Whitney U p-value <0.05) as either queried (blue) or conditionally
1388 depleted (red).

1389

1390 C) Distribution of percentage TA site disruption for all genes in the inoculum library (top)
1391 and a representative library recovered from the rabbit colon (bottom) overlaid with the

1392 classifications described in panel B. Genes with insufficient data are removed from the
1393 bottom panel.

1394
1395 D) Distribution of Con-ARTIST gene classifications (insufficient data (ID, black), queried
1396 (Q, blue), conditionally depleted (CD, red)) in each library recovered from seven infant
1397 rabbit colons two days post infection as compared to the inoculum.

1398
1399 E) Conditionally depleted genes (defined by Con-ARTIST consensus approach) by
1400 Clusters of Orthologous Groups (COG) classification. COG enrichment index (displayed
1401 as log₂ enrichment) is calculated as the percentage of the CD genes assigned to a
1402 specific COG divided by the percentage of genes in that COG in the entire genome. A
1403 two-tailed Fisher's exact test with a Bonferroni correction was used to test the null
1404 hypothesis that enrichment is independent of TIS classification. (***) p-value <0.0001.

1405
1406 F) KEGG pathways of EHEC genes classified as conditionally depleted by Con-ARTIST
1407 consensus approach.

1408
1409 G) Distribution of PC1 scores across all EHEC genes. Red bins fall within the lowest
1410 10% of PC1 scores.

1411
1412 **Figure 3: Con-ARTIST and CompTIS-based classification of LEE genes and T3SS**
1413 **effectors.**

1414 A) Artemis plots of reads in the LEE pathogenicity island in the control-simulated
1415 inoculum library (top) and a representative library recovered from the rabbit colon
1416 (middle). The genes in the LEE are displayed at the bottom. The color of the gene
1417 corresponds to its classification. Maroon genes were categorized as conditionally
1418 depleted by Con-ARTIST consensus and fell in the bottom 10% of glPC1 scores by
1419 CompTIS; red genes had a glPC1 score in the bottom 10% of the distribution, but not
1420 classified as CD by Con-ARTIST. Gray genes did not meet the glPC1 cutoff and were
1421 not classified by Con-ARTIST.

1422
1423 B) Schematic showing classification of the LEE genes and non-LEE-encoded effectors.
1424 Color symbols as above. Orange indicates the gene was identified as CD by Con-
1425 ARTIST and had a glPC1 score above 10%. LEE-encoded effectors, non-LEE encoded
1426 effectors, chaperones/substrate selection proteins, and regulators are indicated in
1427 boxes.

1428
1429
1430 **Figure 4: Validation of colonization defects in selected mutants.**

1431 A) Competitive indices of indicated mutants vs wild type EHEC. Bar-coded mutants
1432 were co-inoculated with bar-coded wild-type EHEC into infant rabbits and recovered two
1433 days later from the colons of infected rabbits. Relative abundance of each mutant was
1434 determined by sequencing the barcodes. (**) p-values < 0.01 and (***) p-value <0.001.
1435 *ΔhupB*, which had a glPC1 score in the bottom 10%, but was not classified as CD by
1436 Con-ARTIST, is highlighted in blue.

1437

1438

1439 **Figure 5: Effector translocation and survival in response to various gastro-**
1440 **intestinal stressors by mutants defective in colonization.**

1441 A) Normalized effector translocation of mutants compared to WT. $\Delta escN$, a mutant that
1442 abrogates T3SS activity, was used as a control. Mutants were tested for their ability to
1443 translocate EspF-TEM1 into HeLa cells, as measured by a shift in emission spectra
1444 from 520 to 450 nm. Fluorescence was normalized to WT levels. Geometric means and
1445 geometric standard deviations are plotted.

1446

1447 B) MIC for NaCl (osmotic stress) and crude bile for the indicated mutants. Bold text
1448 highlights values differing from the wild-type.

1449

1450 C) Normalized acid resistance. Mutants were tested for their ability to survive low acid
1451 shock. Survival is shown as a percentage of the acid resistance of the WT. Geometric
1452 mean and geometric standard deviation are plotted.

1453

1454 **Figure 6: CvpA promotes EHEC resistance to deoxycholate.**

1455

1456 A) Dilution series of WT, $\Delta cvpA$ mutant, and $\Delta cvpA$ mutant with arabinose-inducible
1457 *cvpA* complementation plasmid plated on LB, LB 1% deoxycholate (DOC), LB 1%
1458 cholate (CHO), or LB 1% DOC + 0.2% arabinose.

1459

1460 B) Optical density of WT and $\Delta cvpA$ grown in LB and two concentrations of DOC, added
1461 at the indicated arrow. The average of three readings is plotted with errors bars indicate
1462 standard deviation.

1463

1464 C) MIC of antimicrobial compounds for WT and $\Delta cvpA$ and $\Delta acrAB$ mutants. Units are
1465 mg/mL unless specified otherwise. Bolded values are those different than the wild-type.

1466

1467 D) Predicted CvpA topology diagram.

1468

1469

1470 **Supplementary Figure Legends**

1471

1472 **S1: Sequencing saturation of TIS libraries**

1473 Reads were randomly sampled from each library and the percentage of TA sites
1474 disrupted in each randomly selected pool were plotted for the EDL933 library (A),
1475 MG1655 library (B), the inoculum library used to infect infant rabbits (C), and the
1476 libraries recovered from 7 rabbit colons (D-J).

1477

1478 **S2: Assessment of non-neutral EHEC EDL933 genes.**

1479 A) Non-neutral genes (defined by EL-ARTIST as either regional or underrepresented)
1480 by Clusters of Orthologous Groups (COG) classification. COG enrichment index
1481 (displayed as \log_2 enrichment) is calculated as defined in (51) as the percentage of the
1482 CD genes assigned to a specific COG divided by the percentage of genes in that COG
1483 in the entire genome. A two-tailed Fisher's exact test with a Bonferroni correction was

1484 used to test the null hypothesis that enrichment is independent of TIS status. p-values
1485 considered to be significant if <0.002. Single asterisks (*) indicates p-value <0.002,
1486 double asterisks (**) indicate p-value <0.001, and triple asterisks (***) indicate p-value
1487 <0.0001.

1488
1489 B) GC content (%) of EDL933 genes classified as either neutral (blue) or non-neutral
1490 (regional + underrepresented; red) by TIS. Distributions are compared using a Mann-
1491 Whitney U non-parametric test; (***) p-value of <0.0001.

1492
1493 C) GC content (%) of EDL933 genes classified as either having homologs in MG1655
1494 (homolog) or lacking homologs (divergent). Distributions are compared using a Mann-
1495 Whitney U test; (***)p-value of <0.0001.

1496
1497 D) TA insertions across *kdsC* in EDL933 (left) and MG1655 (right).

1498
1499 **S3: Con-ARTIST and CompTIS classification of genes important for colonization**

1500 A-G) Distribution of percentage TA site disruption in libraries recovered from 7 rabbit
1501 colons. These distributions are overlaid with Con-ARTIST classification (queried, blue;
1502 CD (conditionally depleted), red) as described in Figure 2B.

1503 H) CD genes were grouped by consensus across animals. Many genes are CD in only
1504 one rabbit; fewer are classified as CD across all 7 animals. A standard of consensus of
1505 5 or more animals was chosen to determine the list of CD genes, indicated with an
1506 asterisk (*).

1507
1508 I) Variance explained by each gene-level (gl) principal component for gIPCA performed
1509 across the 7 rabbit screens.

1510
1511 J) Gene-level principal component 1 (gIPC1) coefficients for each rabbit dataset.

1512
1513 H) Heatmap of the log₂ fold change for each gene with a gIPC1 score that falls within
1514 the bottom 10% of the distribution. Each column represents genes from a separate
1515 rabbit replicate. Genes are ordered by PC1 score, lowest at the top of the heatmap and
1516 highest at the bottom.

1517
1518 **S4: In vitro growth and morphology of mutants.**

1519 A) 17 mutant strains plus the wild-type were grown in LB and turbidity measured by
1520 optical density. The average of three readings with the standard deviation is plotted.

1521
1522 B) Cell-shape defects of $\Delta sufl$, $\Delta envC$, $\Delta tatABC$, and Δprc mutants in high osmolality
1523 media. Morphology in LB (top) or LB supplemented with 0.3M NaCl (bottom) is shown.

1524
1525 **S5: Characterization of $\Delta cvpA$.**

1526 A) Biofilm production in WT and $\Delta cvpA$ using crystal violet staining and absorption.
1527 Levels were normalized to a percent of the WT value; three samples were analyzed and
1528 the geometric means and geometric standard deviation are plotted. The differences
1529 between the two groups were not significant (n.s.) by Mann-Whitney U.

1530

1531 B) WT and $\Delta cvpA$ struck to single colonies on an agar plate made with YESCA media
1532 supplemented with Congo Red to detect curli fibers.

1533

1534 C) WT and $\Delta cvpA$ struck to single colonies on an agar plate containing minimal media
1535 with no exogenous purines.

1536

1537 D) Dilution series of *Vibrio cholerae* C6706 WT and *cvpA::tn* plated on LB and LB 1%
1538 deoxycholate (DOC).

1539

1540 E) Dilution series of *Vibrio parahaemolyticus* WT and $\Delta cvpA$ plated on LB and LB 1%
1541 deoxycholate (DOC).

1542

1543 **Supplementary Tables Captions:**

1544 S1) RS to Z Annotation

1545 S2) EHEC EL-ARTIST

1546 S3) EHEC Non-Neutral KEGG

1547 S4) K12 EL-ARTIST

1548 S5) EHEC Unique Non-Neutral

1549 S6) Con-ARTIST and CompTIS

1550 S7) EHEC CD Genes KEGG

1551 S8) Strains

1552 S9) Oligos

1553

1554

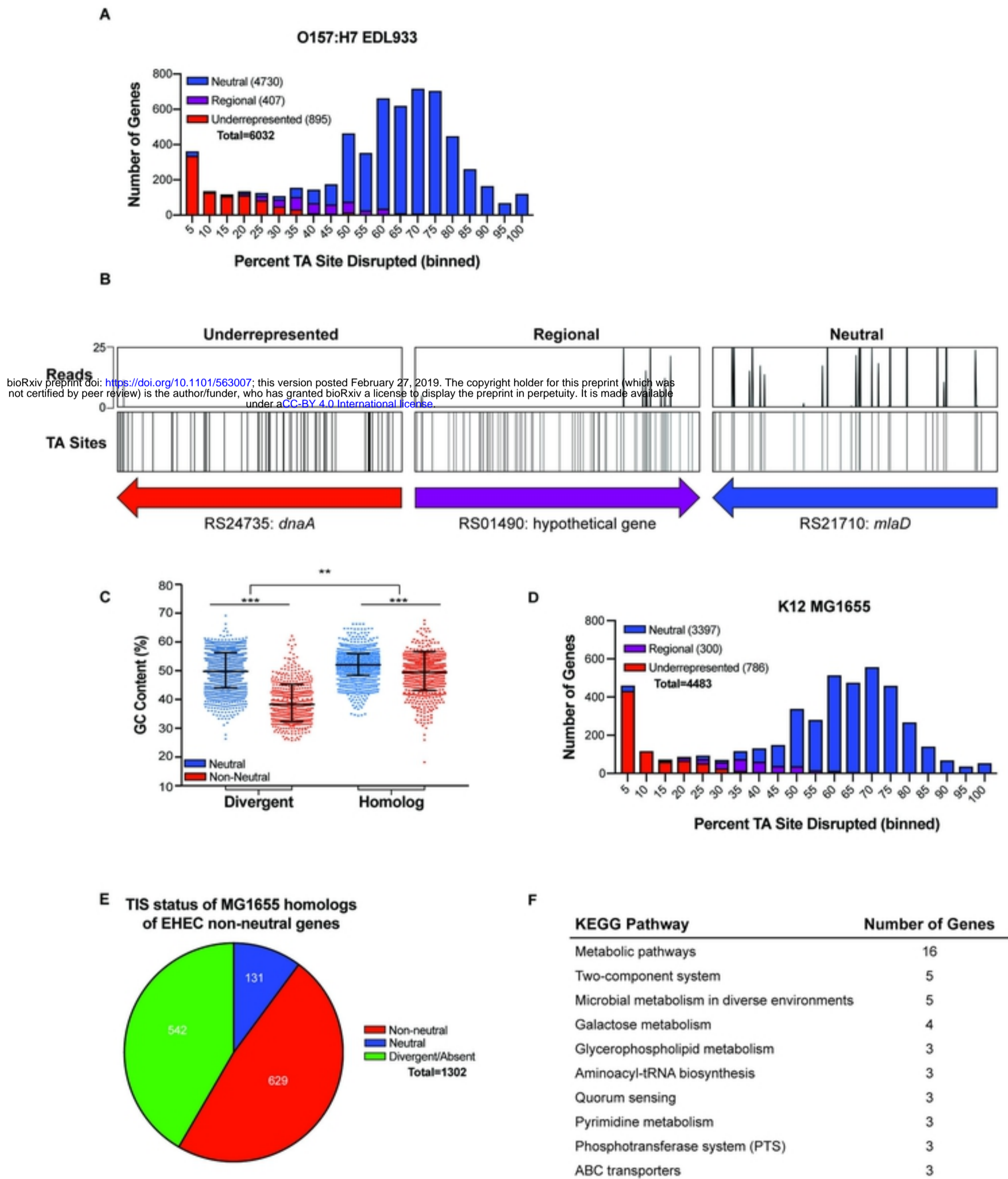


Figure 1

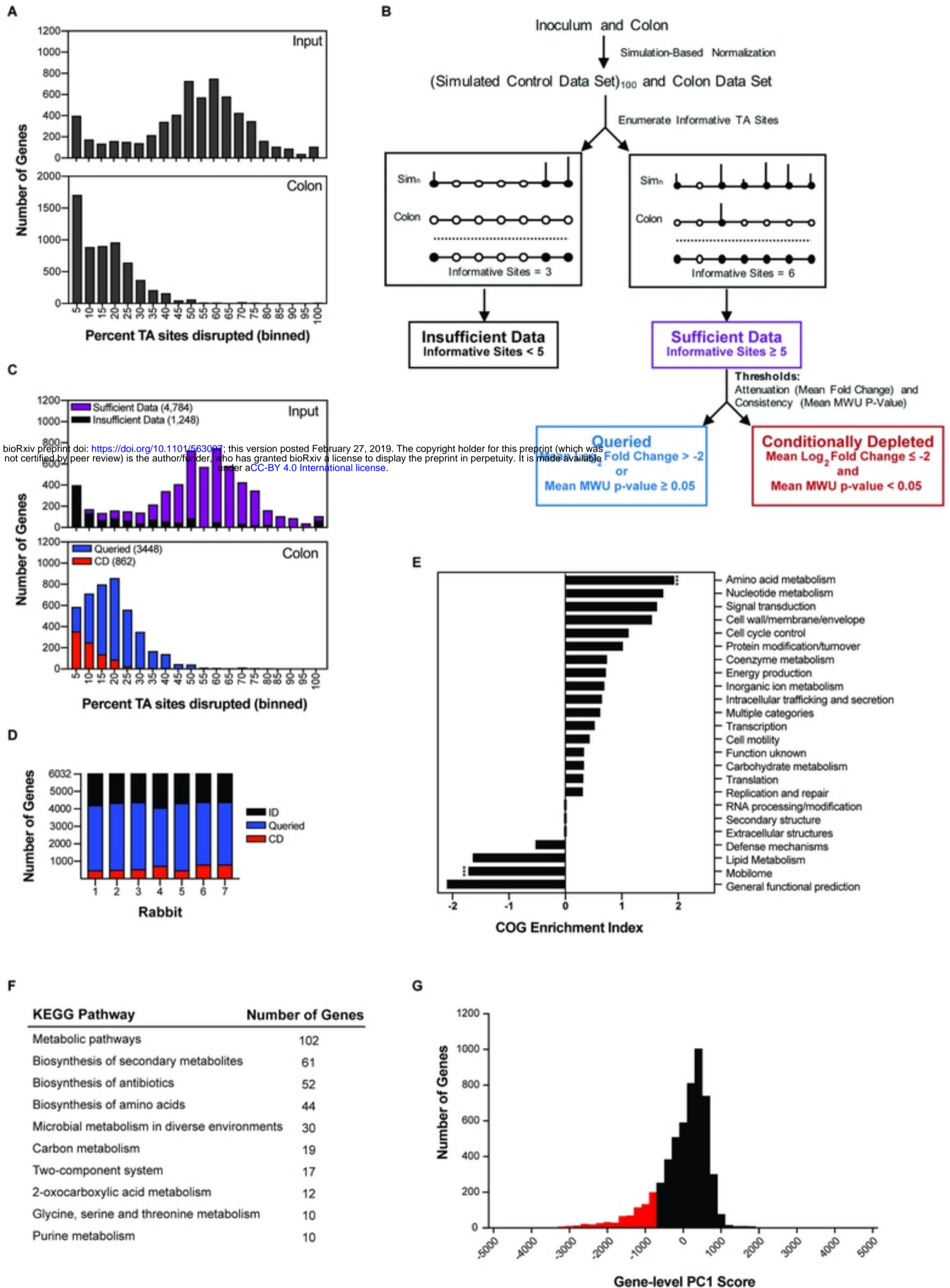


Figure 2

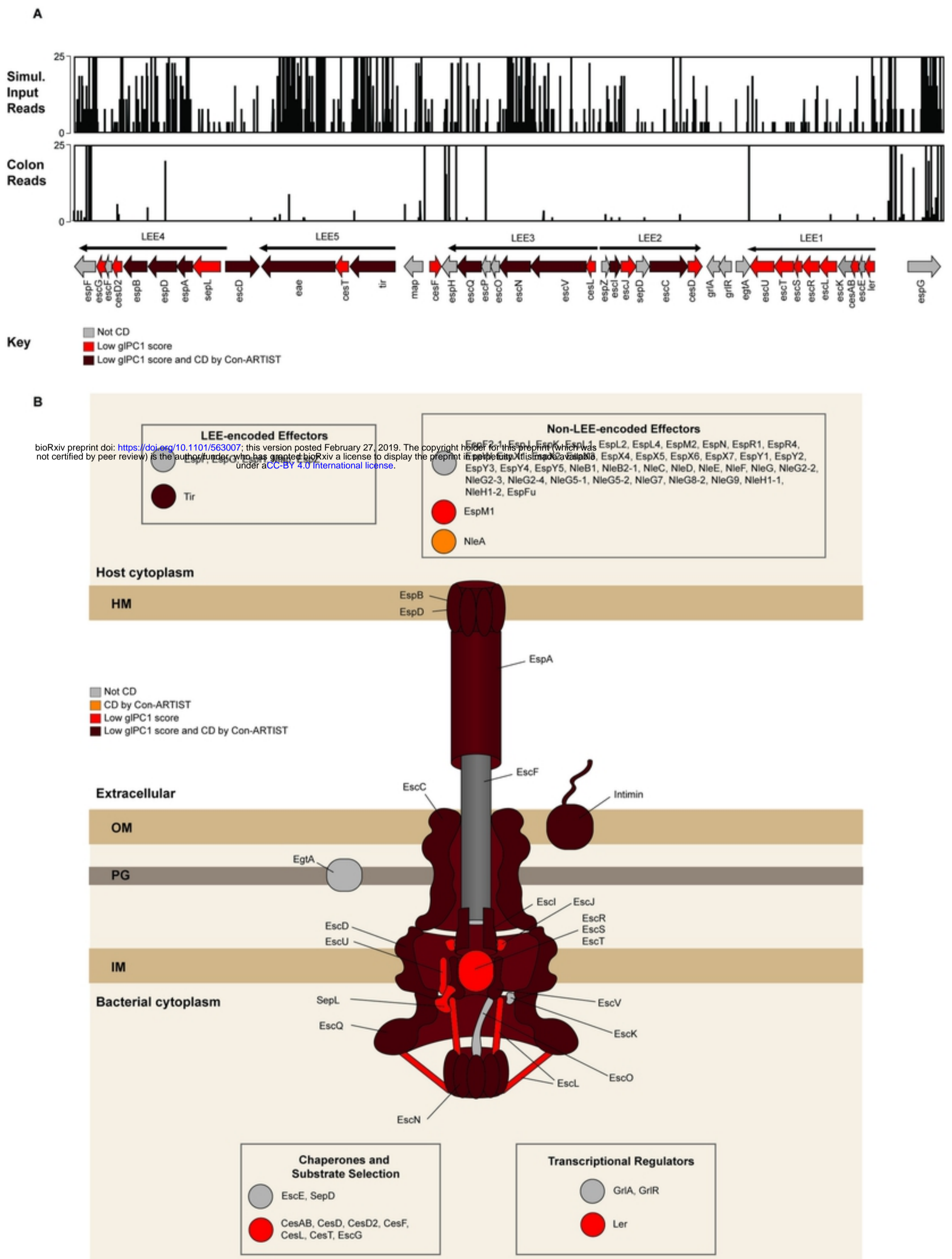
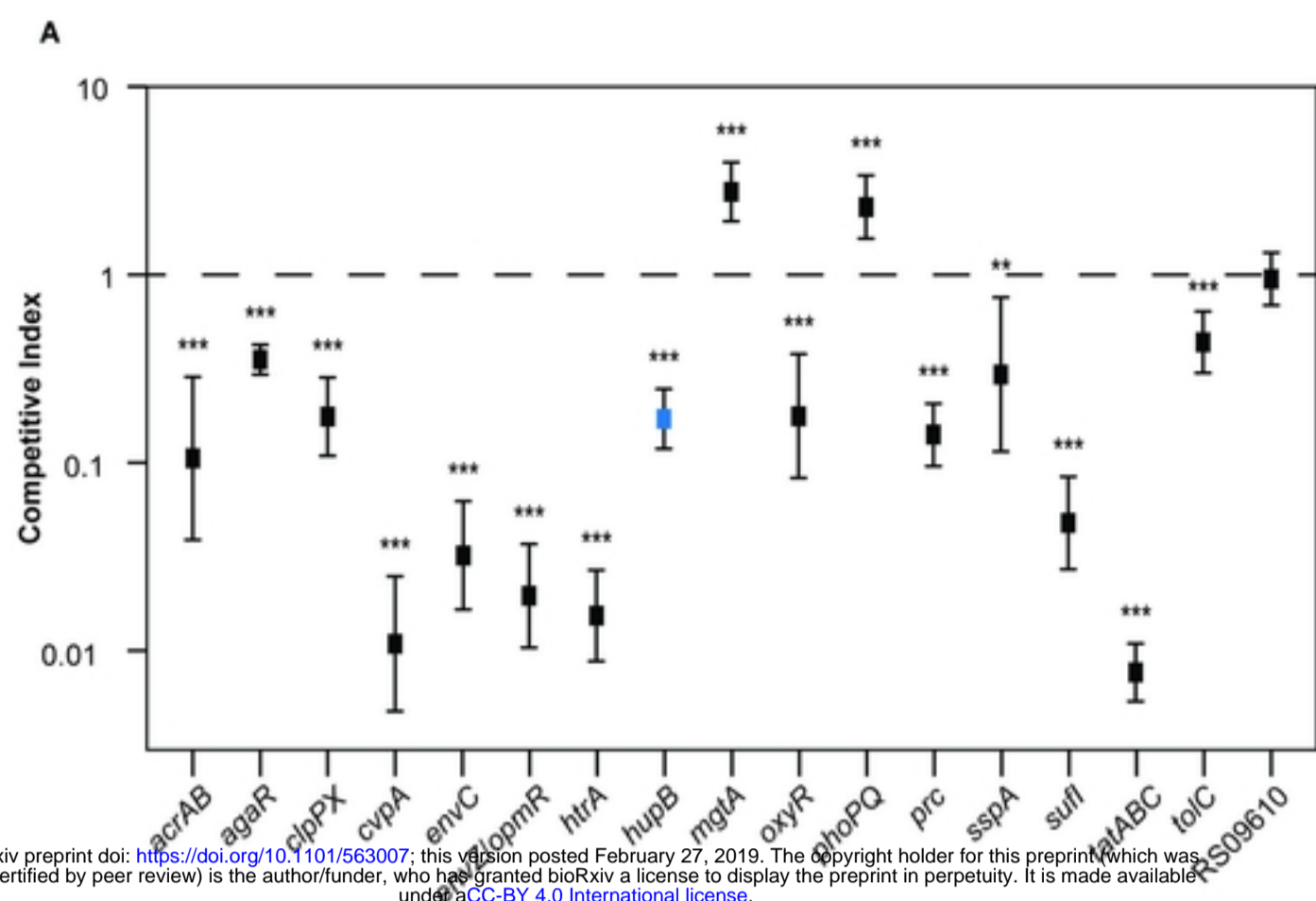
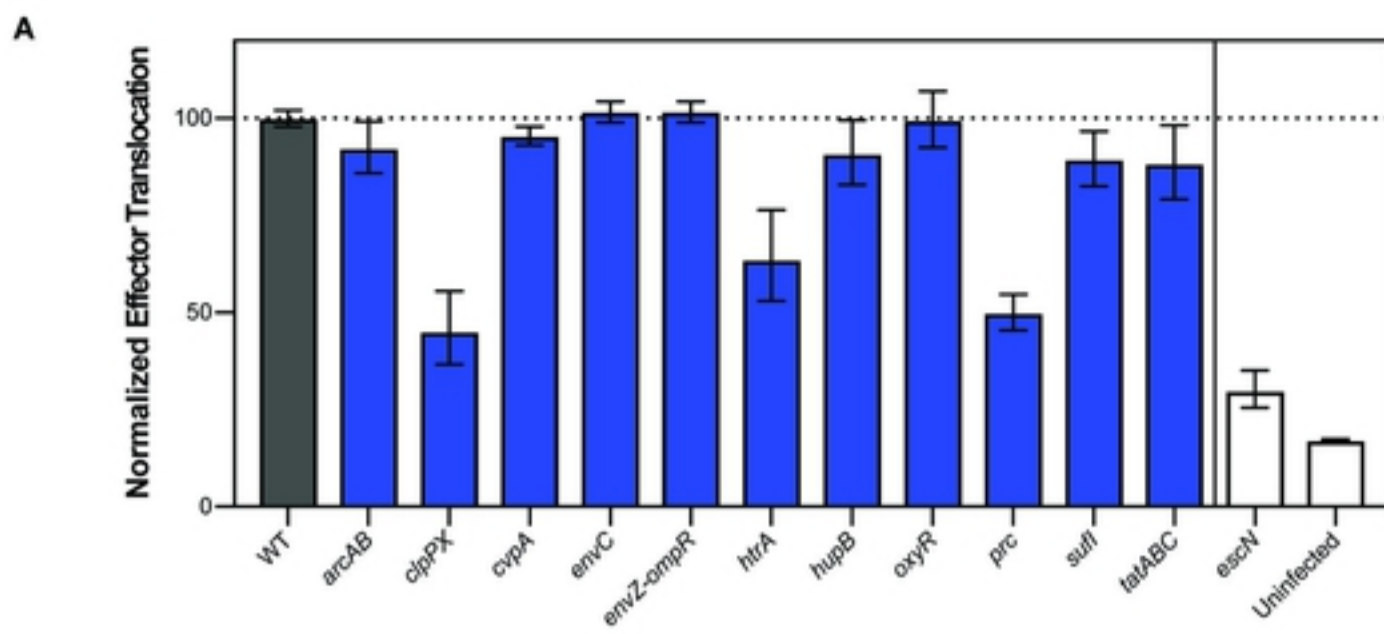


Figure 3



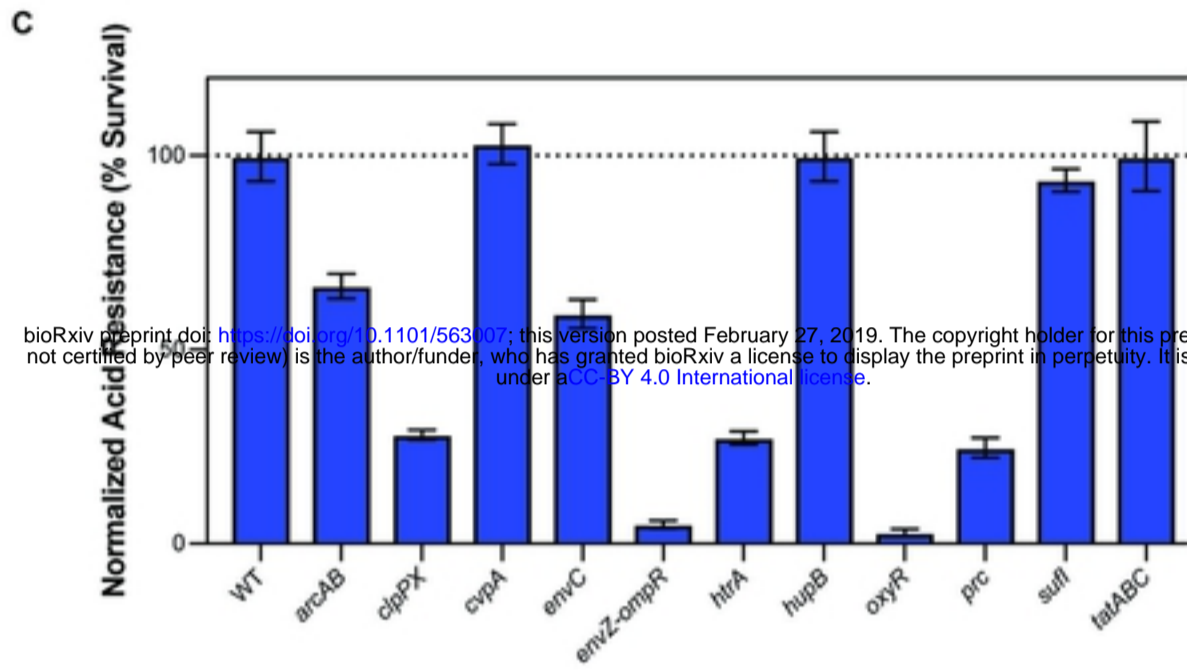
bioRxiv preprint doi: <https://doi.org/10.1101/563007>; this version posted February 27, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

Figure 4

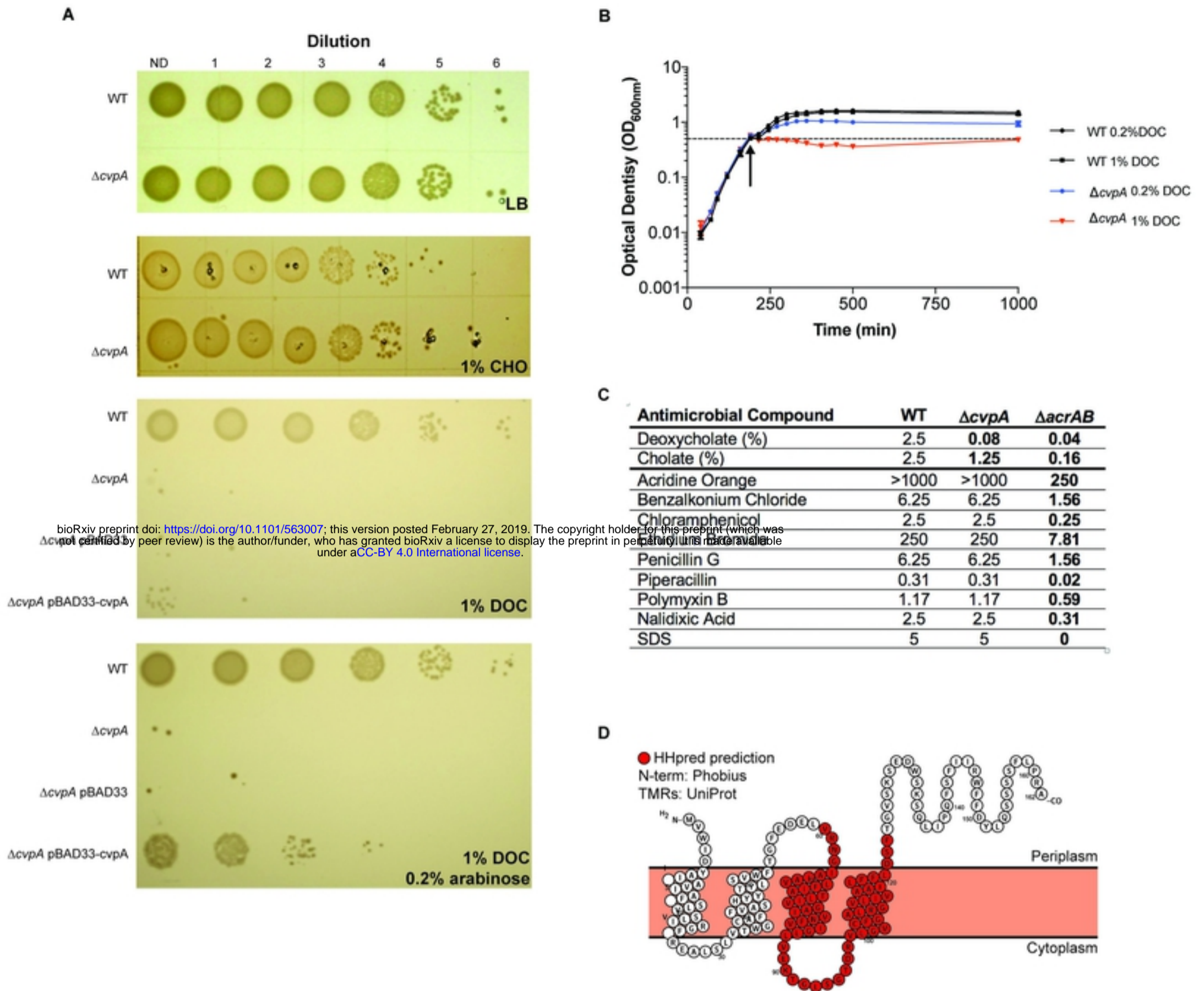


B

Strain	NaCl (mM)	Bile (%)
WT	1500	12
Δ acrAB	1500	0.05
Δ cdpPX	1500	6
Δ cypA	1500	3
Δ envC	750	0.05
Δ envZ-ompR	750	1.5
Δ htrA	1500	12
Δ hupB	1500	12
Δ oxyR	1500	0.75
Δ prc	750	6
Δ sufl	1500	12
Δ tatABC	750	6



bioRxiv preprint doi: <https://doi.org/10.1101/563007>; this version posted February 27, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



bioRxiv preprint doi: <https://doi.org/10.1101/563007>; this version posted February 27, 2019. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

Figure 6