

1 **Full Title:**

2 **Combining *Shigella* Tn-seq data with Gold-standard *E. coli* Gene Deletion Data**
3 **Suggests Rare Transitions between Essential and Non-essential Gene Functionality**

4 **Short Title:**

5 **Transitions between Essential and Non-essential Functionality**

6

7 Nikki E. Freed^{1,2}, Dirk Bumann², and Olin K. Silander^{1,3,*}

8

9 1 Institute of Natural and Mathematical Sciences, Massey University, Auckland, New
10 Zealand

11 2 Infection Biology, Biozentrum, University of Basel, Basel, Switzerland

12 3 Computational and Systems Biology, Biozentrum, University of Basel, Basel, Switzerland

13 * Corresponding author: olinsilander@gmail.com (OKS)

14 **Abstract**

15 Gene essentiality - whether or not a gene is necessary for cell growth - is a fundamental
16 component of gene function. It is not well established how quickly gene essentiality can
17 change, as few studies have compared empirical measures of essentiality between closely
18 related organisms. Here we present the results of a Tn-seq experiment designed to detect
19 essential protein coding genes in the bacterial pathogen *Shigella flexneri* 2a 2457T on a
20 genome-wide scale. Superficial analysis of this data suggested that 451 protein-coding genes
21 in this *Shigella* strain are critical for robust cellular growth on rich media. Comparison of this
22 set of genes with a gold-standard data set of essential genes in the closely related *Escherichia*
23 *coli* K12 BW25113 suggested that an excessive number of genes appeared essential in
24 *Shigella* but non-essential in *E. coli*. Importantly, and in converse to this comparison, we
25 found no genes that were essential in *E. coli* and non-essential in *Shigella*, suggesting that
26 many genes were artefactually inferred as essential in *Shigella*. Controlling for such artefacts
27 resulted in a much smaller set of discrepant genes. Among these, we identified three sets of
28 functionally related genes; two of which have previously been implicated as critical for
29 *Shigella* growth, but which are dispensable for *E. coli* growth. The data presented here
30 highlight the small number of protein coding genes for which we have strong evidence that
31 their essentiality status differs between the closely related bacterial taxa *E. coli* and *Shigella*.
32 A set of genes involved in acetate utilization provides a canonical example. These results
33 leave open the possibility of developing strain-specific antibiotic treatments targeting such
34 differentially essential genes, but suggest that such opportunities may be rare in closely
35 related bacteria.

36 **Author Summary**

37 Essential genes are those that encode proteins required for growth and survival in a particular
38 environment. We performed experiments using transposons, genetic elements that disrupt
39 gene function, to determine the set of essential genes in the pathogenic bacteria *Shigella*
40 *flexneri*. We then compared our results to the well-characterized set of essential genes in the
41 closely related, yet non-pathogenic, bacteria *Escherichia coli*. We found only a small number
42 of genes that are important for growth in *Shigella flexneri*, yet not in *Escherichia coli*. We
43 believe these findings are interesting for several reasons; they help us better understand how
44 quickly the functions of proteins change over time; they suggest possible targets for
45 developing strain-specific antibiotic treatments; and they expand our basic understanding of
46 this pathogen's metabolic processes.

47 **Introduction**

48 One general functional characteristic of a gene is essentiality - whether that gene is required
49 for cellular viability and growth. In haploid (e.g. bacterial) genomes, this characteristic can be
50 assessed by attempting to delete a specific gene from a genome. When such a deletion is not
51 possible, this gene is frequently termed “essential” [1], implying that the gene is necessary for
52 cell growth and viability. Gene disruption, although less precise, is more commonly used to
53 infer essentiality using a similar criterion. For example, genes that cannot be disrupted by
54 transposon insertion have been inferred as being essential (e.g. [2]).

55 One important question is how quickly essential functions change over evolutionary time. If
56 orthologous protein coding genes in two bacterial strains differ in their essentiality
57 classification, this suggests that either the biochemical nature of the protein has changed, or
58 that the cellular context in which the protein acts has changed [3]. It has been experimentally
59 established that such transitions can occur [4-6]. Here we examine how frequently proteins go
60 from being essential to non-essential and vice versa in nature.

61 A recent study quantified changes in the essentiality classifications of protein coding genes
62 between three alpha-proteobacteria: *Caulobacter crescentus*, *Brevundimonas subvibrioides*,
63 and *Agrobacterium tumefaciens* [3]. The analysis showed that although orthologous cell
64 components are well conserved, the essentiality of such components (e.g. those involved in
65 the cell cycle) had changed considerably, with only 106 orthologous genes being essential in
66 all three organisms, despite their relatively close evolutionary relationship (89%-93% identity
67 in 16S RNA genes).

68 In this study we combine dense transposon mutagenesis with high-throughput sequencing
69 (Tn-seq [7]) to quantify gene essentiality in *Shigella flexneri* 2a 2457T (hereafter referred to
70 as *Shigella*). We compare the essentiality classifications of protein coding genes in *Shigella*
71 with a gold-standard assessment of essentiality in the closely related strain *Escherichia coli*
72 K12 BW25113 (hereafter referred to as *E. coli*) [1]. These two stains are 99.5% identical in
73 their 16S RNA genes and share approximately 70% of their genomic content.

74 This proximity in evolutionary distance, and the use of a gold-standard data set, brings two
75 unique advantages that have not been available in other studies that have used Tn-Seq or
76 similar methods to quantify gene essentiality [3, 7-19]. First, by relying on the null hypothesis
77 that protein coding genes do maintain their essentiality characteristics, we can objectively
78 assess which quantitative features in the *Shigella* Tn-seq data best predict essentiality or non-
79 essentiality of their orthologous counterparts in *E. coli*; such a comparison to a gold-standard
80 has not yet been used to assess the quality and sensitivity of Tn-Seq data [20], although
81 several studies have validated a small number of Tn-seq-inferred growth defects using clean

82 deletion methods (e.g. [10]). Second, the use of very closely related taxa allows us to quantify
83 on a much shorter time scale the fraction of the essential gene complement that has changed,
84 providing a fine scale window into the rate with which orthologues change in their essential
85 functions.

86 The data presented here suggest that the essential gene complement of *Shigella* and *E. coli*
87 overlap considerably. Indeed, we find no strong evidence that there are any protein-coding
88 genes that are essential in *E. coli* but not *Shigella*. Conversely, we do find a small number of
89 genes that play critical roles for *Shigella* growth, but which have dispensable roles in *E. coli*,
90 or which are absent entirely from *E. coli*. This implies that the functional correspondence, in
91 terms of essentiality, has changed for only a small number of protein-coding genes.

92 However, our analysis also suggests that some protein-coding genes that we observe as
93 undisrupted by transposon insertions are in fact not essential for cell growth. Instead, they are
94 either essential for transposon insertion to occur successfully, or their disruption (but not
95 clean deletion) is detrimental to cell growth. This result emphasizes that in high throughput
96 transposon mutagenesis studies, false positive inferences of essentiality may be common, and
97 that simply increasing the resolution or precision of a dataset cannot necessarily solve this
98 problem.

99 Taken together, our data suggests that the essential gene complement is relatively static over
100 short time scales. However, when protein-coding genes do change from being non-essential to
101 being essential, this appears more likely to occur in pathogenic organisms, perhaps because
102 host environments absolve the organism from manufacturing its own nutrients, or because
103 such organisms have smaller population sizes and are prone to the accumulation of
104 deleterious mutations. It would be interesting to see if this pattern is observed when
105 comparing other pathogens to their free-living sister taxa. If antibiotics can be directed against
106 the function or expression of such differentially essential genes, this may allow targeting such
107 treatments toward specific bacterial strains.

108 **Results and Discussion**

109 **A Transposon Mutagenesis Library Provides Fine-scale Resolution of Gene** 110 **Essentiality**

111 We generated a transposon insertion library by transforming a *Shigella icsA* mutant [21] with
112 a plasmid containing a mini-Tn10 transposase with decreased hotspot activity ([22]; **Fig. 1A**,
113 inset) inducible by isopropyl- β -D-thiogalactopyranoside (IPTG) [23, 24]. After overnight
114 growth on Tryptic Soy Broth (TSB) agar plates containing IPTG, we harvested
115 approximately 10^6 colonies carrying transposon insertions. We pooled and then split this

116 library of clones into six replicates. Three replicates were subject to additional growth step
117 inside the cytoplasm of HeLa cells for four hours. The resulting cells were then harvested, the
118 replicates were bar coded and libraries were prepared. We sequenced all six of these pools on
119 a single Illumina HiSeq lane. For all the analyses presented in this study, we have pooled the
120 data from all replicates and from both of these treatments, as we are focusing on *Shigella*
121 genes that are essential across any permissive growth conditions.

122 From this pool, we mapped insertions at 131,670 unique positions on the *Shigella*
123 chromosome (with many insertions occurring on both the forward and reverse strands but at
124 the same position), and 12,552 unique positions on the large *Shigella* virulence plasmid (see
125 **Methods**). The median distance between inserts on the chromosome was 17 base pairs (bp);
126 on the plasmid this distance was 9 bp. 95% of all inter-insert distances on the chromosome
127 were less than 107 bp; the corresponding figure for the plasmid was 59 bp (**Figs. 1A and B**).

128 Although the distribution of transposon inserts was relatively even across both the
129 chromosome (**S1 Fig.**), at smaller scales we found many regions in which few or no insertions
130 occurred. Quantitative analyses showed that regions containing no transposon insertions for
131 100bp or more were considerably enriched (see **Methods; Figs. 1C and D**). It is likely that
132 many of these regions are critical for cellular growth in *Shigella*. Indeed, we found that for
133 many of the protein-coding genes in these regions, the orthologous *E. coli* genes are known to
134 be essential (**Figs. 2 and 3**).

135 In contrast to the *Shigella* chromosome, we found that few open reading frames on the
136 virulence plasmid were devoid of insertions. Only six out of 263 plasmid ORFs had no inserts.
137 Two of these were replication proteins (CP0258 and CP0259), and two (CP0217 and CP0218)
138 were located within the plasmid stabilisation region. The remaining two, *mxiH* and *acp*, are
139 both less than 250 bp in length (the cut-off used here to classify ORFs as essential; see below),
140 and thus have a lower likelihood of being hit due simply to their smaller target size. The third
141 replication protein of the plasmid, CP0260, contained a single insert in its 858 bp length. The
142 absence of inserts in the plasmid replication or stabilisation regions is explained by the fact
143 that if such insertions did occur, the plasmid would be lost; such insertions would thus never
144 be sequenced. Thus, this data is consistent with the fact that the *Shigella* plasmid contains no
145 essential genes [25], and suggested that our transposon library provided a fine-scale
146 assessment of which *Shigella* chromosomal ORFs provide critical cellular functions.

147 **Average Distance Between Inserts Clearly Delineates Essential and Non-essential** 148 **ORFs**

149 We next quantified which transposon insertion patterns in the chromosome were good
150 predictors of the essentiality of open reading frames. To do so, we first identified 3,027

151 orthologous open reading frames present in both *E. coli* and *Shigella* for which we also had
152 data on essentiality from both the Keio [1, 26] and the Profiling of the Escherichia coli
153 Chromosome (PEC) [27] studies (**S1 Table**). We considered this combined gene set as a gold-
154 standard of essentiality, for two reasons: it is not subject to artefacts that might exist in Tn-seq
155 dataset, such as insertion biases or biases arising during sequencing library preparation (e.g.
156 [28]); and combining both the Keio and PEC datasets should result in few false positive or
157 false negative essentiality characterizations. This set consisted of 277 orthologues considered
158 essential by both studies, 2,717 genes considered non-essential by both studies, and 33 genes
159 for which the two studies disagree.

160 We next quantified several characteristics for each protein-coding gene in our Tn-seq data set,
161 including the total number of inserts per ORF, the mean distance between inserts, the length
162 of the 5' fraction of the ORF upstream of the first insertion, the largest uninterrupted region in
163 the ORF, and others (**S2 Fig.**). We took as a null hypothesis that generally, genes have
164 maintained their essentiality characteristics since the divergence of *E. coli* and *Shigella*. We
165 then tested which of these characteristics best predicted the essentiality status of their
166 orthologous counterparts in the gold-standard dataset of open reading frame essentiality in *E.*
167 *coli* (the Keio and PEC datasets).

168 We found that the best predictor of essentiality status in *E. coli* was the mean distance
169 between transposon insertions in their *Shigella* orthologues (Materials and Methods; **S2 Fig.**).
170 For the *Shigella* orthologues of the 277 *E. coli* essential genes, only four had a mean distance
171 between inserts of less than 150bp. 17 (6%) had a mean inter-insert distance less than 250bp.
172 In contrast, only 6% of the orthologues of non-essential *E. coli* genes had a mean distance
173 between inserts of greater than 250bp (**Fig. 3A**). We selected this mean inter-insert distance
174 of 250bp as a cut-off for classifying *Shigella* ORFs as essential, as it provided a balance
175 between protein coding genes classified as essential in *E. coli* but non-essential in *Shigella* (a
176 6% false negative rate) versus non-essential in *E. coli* and essential in *Shigella* (a 6% false
177 positive rate). By extension, genes that are less than 250bp in length and in which we do not
178 observe insertions were inferred as essential (26 ORFs in total, of which 12 were ribosomal
179 proteins and five were leader peptides). We note, importantly, that almost all of the predictors
180 we tested performed extremely well (**S2 Fig.**).

181 We next investigated in greater detail the disagreements in essentiality classification between
182 *E. coli* and *Shigella* (**Fig. 3B**). Of the 17 *E. coli*-essential genes that this metric identified as
183 non-essential in *Shigella*, all are likely to be false negatives (i.e. in fact essential in *Shigella*,
184 but not classified as such by our criterion). All 17 have a mean distance between inserts of
185 greater than 100 bp (**Fig. 3A**), and nine are uninterrupted for more than 90% of their reading
186 frame. This suggests, surprisingly, that there are no genes that are essential in *E. coli* but

187 whose *Shigella* orthologues are non-essential. This similarity in essentiality is not due to the
188 fact that we use a characteristic that most closely predicts essentiality in the gold standard
189 dataset – this result is robustly corroborated by any meaningful metric that we used (e.g. using
190 other mean distances between inserts as cut-offs for essentiality, using the total number of
191 inserts, the longest uninterrupted gene fraction, or others (**S2 Fig.**)). Overall, this data gives us
192 a very strong prior that genes have maintained their essentiality status (or near-lethal effects
193 on growth) since the divergence of *E. coli* and *Shigella*.

194 **Many Non-essential *E. coli* orthologues of Essential *Shigella* Genes Exhibit** 195 **Impaired Growth**

196 As a result of this strong prior, we thus expect that many of the discrepancies in essentiality
197 between *E. coli* and *Shigella* are false positives due simply to the *Shigella* mutants being non-
198 essential, but having significantly impaired growth. Indeed, of the 160 discrepant genes
199 classified as non-essential in *E. coli* but essential in *Shigella*, 34% of the orthologous *E. coli*
200 deletion genotypes exhibit low growth yields (less than 0.5 OD600 after 22 hours of growth
201 in LB [1]). This contrasts strongly with the 2557 ORFs we classified as non-essential in
202 *Shigella*: only 3.7% of the orthologous *E. coli* deletion genotypes had low growth yields (**Fig.**
203 **4**). Similar but less striking patterns were observed for growth in glucose minimal MOPS
204 media (**S3 Fig.**).

205 It is important to note that Tn-seq assays have only limited power to differentiate essential
206 genes from those whose deletion results in severe growth deficiencies. During the course of
207 preparing the library for sequencing, we estimate that at least 20 generations of growth
208 occurred. If a mutant has a growth rate even 60% that of the wild type, we would expect it to
209 undergo only 12 doublings in contrast to the 20 of the wild type. This would result in a greater
210 than 200-fold underrepresentation of such a mutant ($2^{12}/2^{20}$). In addition, this calculation does
211 not take into account any effects that the mutations have on the length of the lag time, which
212 might also have significant effects on the relative frequency of some mutants.

213 In light of this limited resolution power; given that our prior expectation is that essentiality
214 status changes only rarely; and because we are specifically interested in genes that may have
215 changed in essentiality status, from this point on we focus our analysis on essential *Shigella*
216 genes whose orthologous *E. coli* deletion genotype exhibits robust growth yields (OD600
217 greater than 0.75 after 22 hours growth in LB (**Fig. 4**)). For these genes, we have relatively
218 high confidence that while their deletion in *E. coli* has few effects on growth, their disruption
219 in *Shigella* is lethal or results in a severe growth deficiency.

220 **Artefacts of the Transposon Screen Explain Some False Positive Discrepancies**

221 35 essential *Shigella* genes have orthologous *E. coli* deletion genotypes with growth yields
 222 higher than 0.75 OD600 after 22 hours growth in LB [1] (**Fig. 4**). Careful inspection
 223 suggested that some genes were present in this set due to differences in the growth conditions
 224 between *E. coli* Keio and PEC collections and our own. For example, *fhuACD*, and *tonB* all
 225 appeared in this set of genes (**Table 1**). All four are involved in iron acquisition, and it is
 226 likely that iron was limiting in the solid agar media [29] used during the preparation of the
 227 Tn-seq library, as compared to the liquid LB used to measure growth in the Keio study.

228

229 **Table 1. Genes artefactually inferred as essential in *Shigella*.**

Gene(s)	Evidence for Artefactual Inference of <i>Shigella</i> Essentiality
<i>acrAB</i> , <i>tolC</i> , <i>ybaB</i> , <i>ksgA</i> , <i>yebC</i> , <i>smpB</i> (0.73) ¹	Affect kanamycin resistance [30-32]
<i>dnaQ</i> , <i>hold</i> , <i>recD</i> , <i>xseA</i> , <i>ruvA</i> (0.58), <i>ruvB</i> (0.60), <i>ruvC</i> (0.61), <i>recB</i> (0.58), <i>recC</i> (0.65)	Likely to affect the transposition process; <i>dnaQ</i> , <i>hold</i> , <i>ruvA</i> , and <i>ruvB</i> inferred as essential using Tn-seq in <i>Salmonella</i> [13]
<i>priB</i>	Deficient in plasmid maintenance [33, 34]
<i>fhuACD</i> , <i>tonB</i>	Involved in iron acquisition which is critical for growth in the iron limited media used in this study [29]
<i>rpsT</i>	S20 ribosomal subunit; new data indicates mutants have poor growth [35] (in conflict with Keio data)
<i>miaA</i>	tRNA dimethyl transferase; previous data indicates <i>E. coli</i> mutants have poor growth [36] (in conflict with Keio data)
<i>ompA</i>	Outer membrane porin; clean knockouts appear viable [37]; mutant forms are frequently lethal [38]
<i>pitA</i>	Metal phosphate transporter with ten transmembrane segments; transposon disruption of substrate transporters is three-fold more likely to be inferred as essential compared to clean deletion (see main text; Fig)
<i>potB</i>	Type I ABC transporter (Putrescine / spermidine transporter)
<i>cysU</i>	Type I ABC transporter (Sulfate / thiosulfate transporter)
<i>sapB</i>	Type I ABC transporter (unknown substrate)
<i>ptsH</i>	Short 306 bp reading frame
<i>ydhR</i>	Short 258 bp reading frame

230 ¹ Gene deletions of orthologous *E. coli* genes with growth levels less than OD600 0.75 have these
 231 levels in parentheses

232

233 For a second set of genes, the discrepancies are likely due to the differences in methodology
234 between the *E. coli* (precise gene deletions) and *Shigella* (transposon inactivation) studies
235 (**Table 1**). We inferred *acrAB* and *tolC* as essential. These genes act together as an efflux
236 pump, and mutations in these genes result in hypersensitivity to antibiotics [39]. Thus, clones
237 with transposon insertions in these genes are unlikely to survive during library growth under
238 kanamycin selection. A similar explanation likely underlies the fact that we inferred *ybaB*,
239 *ksgA*, *yebC*, and *smpB* as essential: these four play role in aminoglycoside resistance [30-32].
240 We also inferred *priB*, *dnaQ*, *hold*, *xseA*, and *recD* as being essential in *Shigella*, although
241 the *E. coli* deletion genotypes exhibit robust growth. All of these are involved in DNA
242 replication, recombination, and double strand break repair, all of which are essential processes
243 in the completion of the transposition process [40]. The related genes *recBC* and *ruvABC*
244 contained a single insert between the five of them, while the *E. coli* deletion genotypes all
245 exhibit only slightly impaired growth of 0.6 OD600 or more (**Table 1**). Certain *recBC*
246 mutants can have considerable effects on the rate of Tn10 excision [41, 42] and we speculate
247 that this may be one reason why we rarely observed insertions in these loci. It has also been
248 speculated that *ruv* mutants inhibit transposition [43]. We propose that after transposition
249 occurs, in order for the event to be successfully resolved, transcription of these genes is often
250 required, and the transposition itself precludes the formation of a proper transcript.
251 Thus, the dispensability of these ten genes in *E. coli*, and the similarity in their function,
252 suggests that they all affect successful transposon insertion rather than having critical effects
253 on growth. Notably, *priB*, *dnaQ*, *hold*, *ruvA*, and *ruvB* were also inferred as essential in the
254 closely related bacterium *Salmonella typhimurium* via a high-throughput transposon assay. In
255 the same study *ruvC*, *recB*, *recC*, and *ybaB*, were inferred as extremely important for growth
256 while *ksgA* and *yebC* were inferred as significantly impairing growth [13]. Again, the
257 majority of these knockouts in *E. coli* exhibit very robust growth (greater than 0.75 OD600
258 after 22h growth in LB). Given the roles that these genes are known to play in transposition
259 and antibiotic resistance, this suggests that the inference of essentiality may be due to
260 artefacts of the transposon screen.

261 For a third set of genes, the literature presents conflicting information on the growth
262 phenotypes, with studies that have individually assessed growth rates suggesting poor growth.
263 These include *rpsT* [35], *miaA* [36], and *ompA* [38] (**Table 1**).

264 There were also two open reading frames that we inferred as differentially essential as they
265 were completely uninterrupted in our data. However, these two open reading frames, *ydhR*
266 and *ptsH*, are very small and less likely to be disrupted, being 306 bp and 258 bp long,

267 respectively. It is probable, then, that this discrepancy is not driven by different physiological
268 roles that they play in *E. coli* as compared to *Shigella*.

269 Finally, we tested for other possible artefactual patterns in the data based on gene function.
270 We asked whether there were specific functional categories in which genes were more likely
271 to be inferred as essential using the transposon mutagenesis screen in *Shigella* as opposed to
272 clean deletions in *E. coli*. We found two functional categories of genes that showed clear
273 enrichment: genes involved in substrate transport and / or active transport, which were 3- and
274 2.1-fold enriched, respectively (**Fig. 5**). We hypothesize that one reason for this enrichment is
275 that truncated versions of these proteins disturb the operation of the *sec* machinery, thereby
276 decreasing or stopping growth. Thus, we propose that the four active transporters we infer as
277 essential in *Shigella* but not *E. coli* (**Table 1**) are artefacts due to the transposition process
278 resulting in truncated proteins.

279 **Genes uniquely essential in *Shigella flexneri***

280 While many differences in essentiality classification between *Shigella* and *E. coli* are likely
281 due to (1) severe growth defects present in both *E. coli* and *Shigella* rather than strict
282 essentiality; and (2) differences in environmental conditions (e.g. iron) between the *E. coli*
283 and *Shigella* assays; and (3) artefacts of the *Shigella* transposon screen that do not occur in
284 the *E. coli* knockout screen, we do find a number of genes which we infer to be uniquely
285 essential to *Shigella*. We expect that the physiological differences between *E. coli* and
286 *Shigella* are driving these differences in gene essentiality (**Table 2**).

287 Among the set of genes essential in *Shigella* but dispensable in *E. coli* is *lysS*: this ORF has a
288 functional homologue in *E. coli* (*lysU* [45]), while in *Shigella flexneri* there is no homologue.
289 Also in this set of genes are *proA*, *proB*, and *proC*. These genes act in proline biosynthesis.
290 Given the rich media the cells were grown in, it is surprising that they would be essential. In
291 addition, as *proB* is involved in the first committed step of proline synthesis, its disruption
292 should not cause accumulation of toxic intermediates. However, the data provide strong
293 evidence that the disruption of any these three genes is either lethal or causes severe growth
294 defects (**Fig. 6**). Interestingly, the active proline transporter *putP* is absent from *Shigella* [46].
295 It is also known that in *Salmonella*, the cryptic proline transporter *proY* is silent [47], and we
296 hypothesize that this may also be true of this transporter in *Shigella*. Thus, inefficient proline
297 transport from the media might necessitate biosynthesis.

298 A suite of genes involved in acetate utilization (*aceE*, *aceF*, *ackA*, *pta*, and *pykF*) were all
299 inferred as essential in *Shigella* but dispensable in *E. coli*. The significantly detrimental effect
300 on growth that such mutants have has been noted previously using a completely different
301 approach [21]. The difference in essentiality between these two organisms is most likely due

302 to the absence of acetyl CoA synthetase from *Shigella*, and confirms the sensitivity and
 303 relevance of our transposon mutagenesis assay for assaying differences between *E. coli* and
 304 *Shigella* biology.

305 **Table 2. Genes inferred as uniquely essential in *Shigella*.** All gene deletions in
 306 homologous *E. coli* genes show robust growth in rich media after 22 hours (greater than 0.75
 307 OD600), suggesting that these genes are uniquely essential in *Shigella* as compared to *E. coli*.

Gene(s)	Function [44]	Evidence for Different Physiological Roles in <i>E. coli</i> and <i>Shigella</i>
<i>lysS</i>	Aminoacyl tRNA synthetase, tRNA modification	The <i>lysU</i> functional homologue is absent in <i>Shigella</i> [45]
<i>proABC</i>	Proline biosynthesis	The active proline transporter <i>putP</i> is absent from <i>Shigella</i> [46]. The cryptic transporter <i>proY</i> may be silent, as observed in <i>Salmonella</i> [47], possibly necessitating proline biosynthesis
<i>ackA</i> <i>pta</i> <i>aceEF</i> <i>pykF</i>	acetate kinase phosphotransacetylase pyruvate dehydrogenase pyruvate kinase	All affect acetate accumulation [48] and utilization [21], which is required for robust growth (<i>Shigella</i> lacks the acetyl CoA synthetase present in <i>E. coli</i> K12 [49])
<i>rfbF</i> , <i>rfbG</i> , <i>rfc</i> , and <i>rfbI</i>	sugar nucleotide biosynthesis for LPS	No <i>E. coli</i> K12 orthologues, as this locus has been replaced by the laterally transferred <i>wbb</i> locus [50]
<i>spr</i>	Murein DD-endopeptidase	None known
<i>tufB</i>	Elongation factor EF-Tu	None known

308

309 For only two other orthologous gene pairs is there strong evidence of discrepant essentiality
 310 status: *tufB* (two insert locations; **Fig. 7**) and *spr* (one insert at base pair 543 across 567 bp).

311 For neither of these genes do we have a hypothetical causal explanation. Interestingly, we
 312 also found very few transposon insertions in the *tufB* paralogue *tufA* (three insert locations;
 313 **Fig. 7**), suggesting that this gene, too, is important for *Shigella* growth despite its relative
 314 dispensability in *E. coli* (0.72 OD600 after 22h in LB). We note that these two genes are
 315 nearly identical in their sequence, which creates ambiguities in mapping some reads.

316 However, this does not explain the absence of reads mapping to either of them.

317 Understanding the molecular mechanisms driving these apparent disparities in growth
 318 phenotypes between *Shigella* and *E. coli* is an important topic for future research.

319 Finally, the transposon insertion data indicated that a within single large operon, containing
 320 the ORFs *rfbACEFGI* / *rfc*, four genes completely lacked insertions (*rfbF*, *rfbG*, *rfc*, and *rfbI*)
 321 (**Fig. 8**). Only *rfbA* and *rfbC* in this operon have *E. coli* orthologues. The remaining genes lie
 322 within a commonly laterally transferred region of the *E. coli* chromosome containing
 323 *wbbHIJKL*, *wzxB* (*rfbX*), and *glf*. These were all laterally transferred into the K12 lineage [50],

324 replacing the *Shigella*-like *rfb* operon. The genes in this operon all play a role in sugar
325 nucleotide biosynthesis necessary for O-antigen synthesis and production of the
326 lipopolysaccharide component of the outer membrane [44]. This provides some evidence that
327 specific aspects of this process have become essential in *Shigella*, despite these genes having
328 been replaced by a laterally transferred set in *E. coli* K12.

329 **Conclusions**

330 By exploiting the extremely close evolutionary relationship of *Shigella flexneri* with *E. coli*
331 K12, the bacterial strain that has been the most extensively and carefully characterized for its
332 essential gene complement [1, 27], we were able to develop an objective metric to precisely
333 quantify how the results of the Tn-seq data relate to essentiality.

334 A superficial analysis of our Tn-seq data suggested that a total of 451 ORFs in *Shigella* were
335 essential for cellular growth in rich media. This is very much in line with what other Tn-Seq
336 studies have found, with numbers ranging from 480 in *Caulobacter crescentus* [8] to 447 in *B.*
337 *subvibrioides* to 372 in *Agrobacterium tumefaciens* [3]. However, it is considerably more than
338 the number that had been found in *E. coli* using in-frame gene knockouts, which is on the
339 order of 300 essential genes. In addition, we found that close to 100% of the reading frames
340 that were classified as essential in *E. coli* K12 were also essential in *S. flexneri*, giving us a
341 strong prior expectation that the essentiality classifications should match between these two
342 taxa.

343 A more nuanced analysis suggested four explanations for artefactual discrepancies in
344 essentiality between *E. coli* and *Shigella*: (1) many *Shigella* genes were not strictly essential
345 but instead gene disruption caused severe growth impairment; (2) differences in experimental
346 conditions (i.e. iron availability); (3) many of the genes we inferred as essential were
347 important for antibiotic resistance or successful transposition, and are in fact dispensable for
348 growth; and (4) transposon disruption of specific functional classes of genes may result in
349 systematically different effects as compared to gene deletions, for example due to the
350 production of truncated protein products. By carefully dissecting the functions of discrepant
351 genes that do not appear to be artefactual, we were able to pinpoint several genes for which
352 there is some evidence of differential physiological roles in *E. coli* and *Shigella*. Among
353 others, these included *lysS*, three genes involved in proline biosynthesis, and a suite of genes
354 involved in acetate utilization (**Table 2**). In addition to these, we found one large operon
355 which appears to have an essential role in *Shigella* growth but which is missing completely in
356 *E. coli*. Surprisingly, we found only two additional genes that are differentially essential (*tufB*
357 and *spr*) (**Table 2**).

358 Even after attempting to decrease false positive inferences of gene essentiality in *Shigella*, it
359 appears to be considerably more common for genes to be dispensable for growth in *E. coli*,
360 but critical for growth in *Shigella*. We suggests that one reason *Shigella* more may have a
361 larger complement of essential genes than *E. coli* is that it frequently lives as an intracellular
362 pathogen, and may have lost some of the functional redundancy that is present in *E. coli*. This
363 may occur because host environments provide an abundance of nutrients, or because
364 pathogens requiring a small infectious dose, such as *Shigella* [51], have inherently smaller
365 population sizes and are more subject to genetic drift. A third possibility is that changes in
366 gene function or redundancy may have occurred through selection for increased virulence,
367 which has resulted in the inactivation of certain genes being selectively advantageous. Finally,
368 we note that the discrepancies in essentiality between these two bacteria may be exploited to
369 develop antibiotics that have strain-specific effects [21].

370 **Methods**

371 **Strains**

372 For all experiments, *Shigella flexneri* 2457T Δ *icsA* was provided by M. B. Goldberg was used
373 as the parental strain. This strain is unable to exploit the host actin cytoskeleton for motility
374 and spreading [52]. Bacterial cells were grown in Tryptic Soy Broth (TSB) media. For
375 experiments using eukaryotic cells, HeLa cells were cultured in DMEM supplemented with
376 10 mM HEPES, 25 mM glucose, and 4 mM glutamine. *Shigella* were grown to exponential
377 phase in tryptic soy broth, coated with poly-L-lysine, and added at a multiplicity of infection
378 of 25, resulting in an infection rate of around 60%. *Shigella* was centrifuged onto HeLa cells
379 ($600 \times g$ for 5 min). At 30 min postinfection, we added gentamicin (100 μ g/mL) to kill
380 extracellular bacteria. Bacterial cells were allowed to grow within HeLa cells for a total of 4
381 hours.

382 **Transposon library**

383 Using a Tn10 transposon with a T7 promoter [23, 24] we created a library consisting of
384 approximately 10^6 clones. This library was created by mating *E. coli* strain BW20767
385 containing the pJA1 transposon plasmid with a spontaneous nalidixic acid resistant clone of
386 *Shigella flexneri* 2457T Δ *icsA* for 5 hours. Transposase expression was induced by plating
387 onto TSB plates containing 0.2 mM isopropyl- β -D-thiogalactoside (IPTG). Colonies were
388 allowed to grow at 37°C for 18 hours on TSB agar plates. All colonies from these plates were
389 then pooled and 100 μ l aliquots of the transposon library were stored at -80°C.

390 Three replicate experiments were carried out on different days in which an aliquot of the
391 transposon library was grown for 18 hours in TSB to stationary phase, diluted 1:100 and

392 grown to exponential phase (0.7 OD₆₀₀). This exponential phase culture was split into two:
393 part of the bacterial culture was pelleted and saved and other was used for infecting HeLa
394 cells (as described in [21]). After 3.5 hours, HeLa cells infected with the *Shigella* transposon
395 library were trypsinized and pelleted. Uninfected HeLa cells were also collected and used to
396 spike the original bacterial culture not used for HeLa infection in order to account for HeLa
397 DNA. All resulting DNA was extracted using the Bacterial Genomic Miniprep Kit (Sigma).

398 **Sequence library construction and sequencing**

399 To amplify the transposon region from these pools, we used one top strand primer annealing
400 to the transposon and a pool of three bottom strand primers each of which consisted of 10
401 random nucleotides followed by a pentamer of common nucleotides in *E. coli* [53]:
402 N₁₀GGTGC, N₁₀GATAT, and N₁₀AGTAC, using Phusion pfu (**S4 Fig**). A nested PCR was
403 then performed to add the P7 and P5 Illumina adapters, as well as a barcode. The products
404 from this second PCR were then size selected for inserts between 200bp and 300bp,
405 quantified using a Qubit, and sequenced on an Illumina HiSeq2000 at the D-BSSE
406 Quantitative Genomics Facility resulting in 49bp single end reads. We used a custom
407 sequencing primer on the P5 end of the molecule such that on both ends of the molecule,
408 reads started directly on the chromosome.

409 **Read mapping**

410 In total, we obtained 198,682,954 reads. We found that the number of reads at each location
411 in the genome varied by up to four orders of magnitude. For this reason, we considered only
412 whether an insertion had occurred at a specific location, and not on the number of reads we
413 obtained at a specific location, which is likely to be highly biased due to PCR artefacts. We
414 thus first deduplicated the reads using *tally* [54], and then used bowtie2 [55] to align the reads
415 to the *Shigella flexneri* 2a 2457T genome and the *Shigella flexneri* 2a str. 301 plasmid
416 pCP301. The sequence of the *S. flexneri* 2457T plasmid is not available. However, the *S.*
417 *flexneri* 2457T and 2a str. 301 plasmids are nearly identical in sequence (differing by 30
418 SNPs; see below). Sequence reads were not trimmed for quality as read quality is taken into
419 account in bowtie2. We used the --sensitive-local option to allow soft clipping on the 3' end
420 of the reads (so that reads that contained adapter sequences at the 3' end could map
421 successfully), and required at least 22bp of matching sequence at the 5' end of the read.

422 We checked for single nucleotide polymorphisms (SNPs) on both the chromosome and the
423 plasmid using the samtools mpileup and bcftools utilities [56, 57]. We retained as possible
424 SNPs only those sites that fulfilled the following three criteria: (1) the SNP was inferred as
425 homozygous (necessarily true, as *Shigella* is haploid); (2) the quality score was above 20; and
426 (3) at least three reads on both the reverse and forward strands confirmed the SNP. We found

427 99 SNPs on the chromosome (as compared to the reference *Shigella flexneri* 2457T in NCBI)
428 and 30 SNPs on the plasmid (as compared to the *Shigella flexneri* 2a str. 301 plasmid in
429 NCBI (in addition to 12 and 2 small indels, respectively). These are detailed in **S2 Table** and
430 **S3 Table**, respectively.

431 Within chromosomal protein coding regions, 44% of all SNPs were synonymous, while 32%
432 fell outside of genic regions (i.e. protein coding or RNA genes). These fractions are greater
433 than one would expect if such SNPs were randomly located on the genome. Only 24% of all
434 mutations in chromosomal coding regions are expected to be synonymous (not accounting for
435 mutational biases), and only 28% of the chromosome is annotated as nongenic (including
436 repeat regions, although for many of these regions, the absence of an annotation may be
437 erroneous). Additionally, only 2 of the 12 (17%) small chromosomal indels fell in coding
438 regions. This suggests that there was some selection against nonsynonymous substitutions
439 that occurred during the culturing and derivation of the *Shigella flexneri* 2a 2457T *virG*
440 mutant. More importantly, the small number of SNPs that we found suggests that few, if any,
441 reads remained unmapped due to sequence differences between the strain used in our
442 experiments and the sequenced GenBank strain.

443 In total, the reads mapped to 89,028 unique locations on the forward strand and 83,074 on the
444 reverse strand of the chromosome, for a total of 172,102 insertions. Some of these insertions
445 occurred at identical positions but on opposite strands, so in total, insertions occurred at
446 131,670 unique sites in the chromosome. Correspondingly, the reads mapped to 8,208 unique
447 locations on the forward strand and 8,585 unique locations on the reverse strand of the
448 plasmid, for a total of 12,552 unique sites. During the insertion of the Tn10 transposon, a 9 bp
449 target DNA sequence is duplicated [58]. We accounted for this duplication in calculating the
450 distances between insertions (by moving the inferred site of insertion for one direction (we
451 arbitrarily selected the antisense direction) backward by 9 bp). Similarly, this duplication was
452 accounted for in calculating various statistics of insertions within genes: sense insertions that
453 were inferred as occurring in the last 9 bp of a gene were ignored in calculating the mean
454 number of insertions per gene (as these bp are duplicated upstream of the insertion).
455 Antisense insertions occurring in the first 9 bp of a gene were ignored, as these bp are
456 duplicated downstream of the insertion.

457 Using the read frequencies at all unique insert locations, we found that the transposon
458 insertions occurred in a biased manner, integrating more often at sites similar to the known
459 9bp consensus NGCTNAGC [58], although this bias was relatively weak (**Figs. 1A and B**,
460 insets). This low level of bias is likely due to our using a transposon with reduced hotspot
461 activity [22]. In addition, we found that insertion frequency was slightly influenced by

462 nucleotides further downstream of this 9bp consensus (**Figs. 1A and B**, insets). Sequence
463 logos for this analysis were visualized using the R package seqLogo [59].

464 The median distance between inserts was 17 bp in the chromosome and 9 bp in the plasmid
465 (**Fig. 1B**), suggesting that the transposon libraries yielded a relatively fine-grained map of the
466 essential genomic complement for both the chromosome and the plasmid.

467 As expected given the variation in insertion densities across the chromosome, we found high
468 variance in the distribution of inter-insert distances. The total length of the *S. flexneri* genome
469 is 4,599,354 bp in total. Given that we observed 131,670 inserts, under a model of random
470 insertion, we would expect a median distance between inserts of 35bp, with 95% of all inter-
471 insert distances being less than 107 bp (under the assumption that these distances are
472 distributed in a geometric manner (i.e. a negative binomial with the number of successes set
473 to one). For the plasmid, we observed 12,552 inserts over 221,618 bp, such that we expect a
474 median distance of 18bp between inserts, and that 95% of all inter-insert distances are less
475 than 59 bp. However, as noted above, we found that on average transposons insertions were
476 separated by a median of 17bp on the chromosome and 9bp on the plasmid. Fitting a
477 geometric distribution to the observed data over 99% of the range of the inter-insert distances
478 (i.e. from 1 to 237 bp for the chromosome and from 1 to 78 bp for the plasmid) more exactly
479 quantified this over-dispersion, and showed that uninterrupted regions in the chromosome
480 greater than 100 bp were considerably enriched (**Fig. 1C**).

481 **Paired end read mapping and inference of IS element dynamics**

482 We used 100 bp paired end Illumina sequencing data from this same library to look for
483 structural rearrangements due to IS elements in the genome. However, this analysis was
484 complicated by the fact that many IS elements share close to 100% identity with others
485 around the genome. During these analyses we thus restricted our searches to regions of the
486 genome for which we had *a priori* expectations that they harboured a rearrangement (i.e. if
487 there were no inserts and the orthologous *E. coli* locus was non-essential or absent).
488 Specifically, we followed the following procedure: we extracted a 50 kilobase pair (Kbp)
489 region from the genome surrounding each hypothesized rearrangement (in all cases, this was
490 a deletion). We then used bowtie2 with the paired end option, allowing up to 10 Kbp inserts
491 to map all reads from our 100 bp PE dataset. From these mapped reads, we retained only read
492 pairs that had (1) mapping quality scores greater than 20; (2) at least one read that matched
493 perfectly (i.e. at all 101 bases of the read) to the genome; and (3) were unique in their length
494 at any specific location (thereby excluding artefacts such as PCR doublets). From these paired
495 reads we then inferred the insert size, which is plotted in **S5 Fig**. The vast majority of insert
496 sizes ranged between 100 and 400bp. However, some were much larger (e.g. up to 9,000 bp

497 in **S5B Fig.**). We inferred that these surrounded regions of the genome that must have been
498 deleted.

499 Such deletions would result in the set of genes contained within as being inferred as essential
500 because of their lack of transposon insertions. However, in the vast majority of cases, we
501 found that when large operons lacked insertions but had non-essential orthologous operons in
502 *E. coli*, or were missing entirely from *E. coli*, these operons were in fact missing from the
503 *Shigella* clone that we used, most likely due to the rapid dynamics of IS elements in this
504 bacterium [60]. For example, no sequence reads we obtained mapped to the *yeaKLMNOP*
505 operon, which spans a total of 9,240 bp. Upon further analysis using a paired end genomic
506 data set, we found that this region was clearly missing from our *Shigella* clone (**S6B Fig.**).
507 This was similarly true for several other operons, as well as for single genes. We did not
508 consider any region in which we identified a deletion in our downstream analyses.

509 **Essential open reading frames**

510 We identified 3,027 unambiguously ORFs that were present in both *E. coli* and *Shigella*
511 *flexneri* 2457T [61], and for which we had essentiality data. We used reciprocal shortest
512 distance [62] to find orthologues, with the requirement that the alignment of the two
513 hypothetical orthologues extend over at least 60% of the longer ORF. To establish a gold-
514 standard set of essential genes we combined the data from two studies of the effects of gene
515 deletion on growth in *E. coli* K12: the Keio collection [1] and the PEC study [27]. We
516 retained only those ORFs which we had data on essentiality from both studies. We then
517 quantified which transposon insertion patterns that most closely corresponded with the
518 essentiality delineations in these studies. Specifically, we selected the feature that maximized
519 the number of true positive essential genes (maximizing the sensitivity) while minimizing the
520 number of FP (maximizing specificity) (this metric is a receiver operator characteristic for
521 which we quantified the area under the curve (AUC; **S2 Fig.**)). We selected from eleven non-
522 independent features: (1) the total number of insertions; (2) the mean number of bp between
523 insertions; (3) the median number of bp between insertions; (4) the number of bp in the 5' end
524 preceding the first insertion; (5) the number of bp in the 5' end preceding the first insertion
525 relative to the total bp in the gene; (6) the number of bp in the 5' end preceding the second
526 insertion; (7) the number of bp in the 5' end preceding the second insertion relative to the
527 total bp in the gene; (8) the number of bp in the longest uninterrupted stretch of the gene; (9)
528 the number of bp in the longest uninterrupted stretch of the gene relative to the total length of
529 the gene, and (10) the number of bp in the longest stretch of the gene interrupted by at most
530 one insertion; (11) the number of bp in the longest stretch of the gene interrupted by at most
531 one insertion, relative to the total length of the gene.

532 We found that for both the PEC dataset and the Keio dataset, the two best predictors of
533 essentiality were the mean distance between inserts (AUC = 0.972 for the PEC dataset, 0.952
534 for the Keio dataset, and 0.973 for the genes on which both datasets agreed on the essentiality
535 classifications); and the fraction of the gene that lay in the longest uninterrupted region (AUC
536 = 0.969 for the PEC dataset, 0.955 for the Keio dataset, and 0.971 for the genes on which both
537 datasets agreed on the essentiality classifications) (**S2 Fig.**). We selected mean distance as on
538 average, it marginally outperformed the other statistic on the gold standard data set.

539 We note that for eight of the 14 genes classified as essential solely in the Keio dataset, the
540 orthologous *Shigella* ORFs have mean distances less than 30 bp, suggesting that these genes
541 may be falsely annotated as essential in the Keio study. In contrast, nine of the ten genes
542 inferred as essential solely in the PEC dataset have mean distances greater than 200 bp; the
543 tenth has a mean distance of 189 bp.

544 **tRNA disruptions**

545 We found insertions in 27 out of 99 tRNAs, with tRNAs for certain amino acids being
546 considerably overrepresented (**S2 Table**).

547 **Additional analyses of differentially classified essential genes**

548 We also tested for the enrichment of certain functional categories in the set of genes that were
549 classified as being essential in *Shigella* but not *E. coli*. This differs from the analysis present
550 in **Fig. 5** in that we are asking whether across a broad set of functions, are specific categories
551 enriched for *Shigella*-essential genes. In **Fig. 5** we ask whether *within* a single functional
552 category, is there a much higher fraction of *Shigella* essential genes than we would expect,
553 given the fraction of genes in that functional category that are essential in *E. coli*.

554 Thus, for this analysis, we separated genes by primary functional category and secondary
555 subcategory using the MultiFun designations (e.g. the primary functional category cell
556 processes divided into the secondary subcategories of cell division, SOS, stress, protection,
557 and motility). We then calculated the fraction of *Shigella*-essential genes within each
558 secondary subcategory and compared this to the total fraction of *Shigella*-essential genes
559 within the primary category (e.g. we calculated the fraction of essential genes in Ribosomal
560 Function (the secondary category) and the fraction of *Shigella*-essential genes in all other
561 categories in Cell Structure (the primary category) (**S6 Fig.**). We tested for enrichment
562 (depletion) using a Fisher exact test.

563 We also examined gene conservation. Highly conserved genes were considered to be those
564 present in more than 50% of all gamma-proteobacteria [61]. We found that genes classified as
565 uniquely essential in *Shigella* were much more conserved across gamma-proteobacteria (79%

566 highly conserved) compared to genes that were found non-essential in both *E. coli* and
567 *Shigella* (36% highly conserved; $p=1.0e-33$, Wilcoxon rank sum test).

568 **Availability of supporting data**

569 All read data are in the SRA with accession numbers XXX.

570 **List of abbreviations used**

571 bp - base pairs; *Shigella* – *Shigella flexneri* 2a 2457T; *E. coli* – *Escherichia coli* BW25113;

572 ORF – open reading frame; PEC - Profiling the E coli Chromosome database

573 **Competing interests**

574 The authors declare no competing interests.

575 **Authors' contributions**

576 NEF, DB, and OKS conceived and designed the transposon mutagenesis. NEF performed the

577 mutagenesis and sequencing. OKS analysed the data with input from DB and NEF. NEF and

578 OKS wrote the paper.

579 **Acknowledgements**

580 The authors thank Luise Wolf for input on designing the sequencing protocol.

581 **Endnotes**

582 **References**

- 583 1. Baba T, Ara T, Hasegawa M, Takai Y, Okumura Y, Baba M, et al. Construction of
584 *Escherichia coli* K-12 in-frame, single-gene knockout mutants: the Keio collection. *Mol Syst Biol.*
585 2006;2:2006 0008. Epub 2006/06/02. doi: 10.1038/msb4100050. PubMed PMID: 16738554; PubMed
586 Central PMCID: PMC1681482.
- 587 2. Hutchison CA, Peterson SN, Gill SR, Cline RT, White O, Fraser CM, et al. Global transposon
588 mutagenesis and a minimal mycoplasma genome. *Science.* 1999;286(5447):2165-9. doi: Doi
589 10.1126/Science.286.5447.2165. PubMed PMID: WOS:000084157300055.
- 590 3. Curtis PD, Brun YV. Identification of essential Alphaproteobacterial genes reveals operational
591 variability in conserved developmental and cell cycle systems *Mol Microbiol.* 2014. doi:
592 10.1111/mmi.12686.
- 593 4. Blank D, Wolf L, Ackermann M, Silander OK. The predictability of molecular evolution
594 during functional innovation. *P Natl Acad Sci USA.* 2014;111(8):3044-9. PubMed PMID:
595 WOS:000332180900041.
- 596 5. Liu GW, Yong MYJ, Yurieva M, Srinivasan KG, Liu J, Lim JSY, et al. Gene Essentiality Is a
597 Quantitative Property Linked to Cellular Evolvability. *Cell.* 2015;163(6):1388-99. PubMed PMID:
598 WOS:000366044800003.

- 599 6. Bergmiller T, Ackermann M, Silander OK. Patterns of Evolutionary Conservation of Essential
600 Genes Correlate with Their Compensability. *Plos Genet.* 2012;8(6). PubMed PMID:
601 WOS:000305961000055.
- 602 7. van Opijnen T, Camilli A. Transposon insertion sequencing: a new tool for systems-level
603 analysis of microorganisms. *Nat Rev Microbiol.* 2013;11(7). doi: Doi 10.1038/Nrmicro3033. PubMed
604 PMID: WOS:000320368400013.
- 605 8. Christen B, Abeliuk E, Collier JM, Kalogeraki VS, Passarelli B, Collier JA, et al. The essential
606 genome of a bacterium. *Mol Syst Biol.* 2011;7. doi: Artn 528
607 Doi 10.1038/Msb.2011.58. PubMed PMID: ISI:000294537800007.
- 608 9. van Opijnen T, Bodi KL, Camilli A. Tn-seq: high-throughput parallel sequencing for fitness
609 and genetic interaction studies in microorganisms. *Nat Methods.* 2009;6(10):767-U21. doi: Doi
610 10.1038/Nmeth.1377. PubMed PMID: ISI:000270355200023.
- 611 10. Kamp HD, Patimalla-Dipali B, Lazinski DW, Wallace-Gadsden F, Camilli A. Gene Fitness
612 Landscapes of *Vibrio cholerae* at Important Stages of Its Life Cycle. *Plos Pathog.* 2013;9(12). doi: Artn
613 E1003800
614 Doi 10.1371/Journal.Ppat.1003800. PubMed PMID: WOS:000330535400034.
- 615 11. Klein BA, Tenorio EL, Lazinski DW, Camilli A, Duncan MJ, Hu LDT. Identification of
616 essential genes of the periodontal pathogen *Porphyromonas gingivalis*. *Bmc Genomics.* 2012;13. doi:
617 Artn 578
618 Doi 10.1186/1471-2164-13-578. PubMed PMID: WOS:000314646600001.
- 619 12. van Opijnen T, Camilli A. A fine scale phenotype-genotype virulence map of a bacterial
620 pathogen. *Genome Res.* 2012;22(12):2541-51. doi: Doi 10.1101/Gr.137430.112. PubMed PMID:
621 WOS:000311895500022.
- 622 13. Langridge GC, Phan MD, Turner DJ, Perkins TT, Parts L, Haase J, et al. Simultaneous assay
623 of every *Salmonella Typhi* gene using one million transposon mutants. *Genome Res.*
624 2009;19(12):2308-16. doi: Doi 10.1101/Gr.097097.109. PubMed PMID: ISI:000272273400015.
- 625 14. Gawronski JD, Wong SMS, Giannoukos G, Ward DV, Akerley BJ. Tracking insertion
626 mutants within libraries by deep sequencing and a genome-wide screen for *Haemophilus* genes
627 required in the lung. *P Natl Acad Sci USA.* 2009;106(38):16422-7. doi: Doi 10.1073/Pnas.0906627106.
628 PubMed PMID: WOS:000270071600076.
- 629 15. Goodman AL, McNulty NP, Zhao Y, Leip D, Mitra RD, Lozupone CA, et al. Identifying
630 Genetic Determinants Needed to Establish a Human Gut Symbiont in Its Habitat. *Cell Host Microbe.*
631 2009;6(3):279-89. doi: Doi 10.1016/J.Chom.2009.08.003. PubMed PMID: WOS:000270290700011.
- 632 16. Lee SA, Gallagher LA, Thongdee M, Staudinger BJ, Lippman S, Singh PK, et al. General and
633 condition-specific essential functions of *Pseudomonas aeruginosa*. *P Natl Acad Sci USA.*
634 2015;112(16):5189-94. PubMed PMID: WOS:000353239100085.
- 635 17. Chao MC, Pritchard JR, Zhang YJJ, Rubin EJ, Livny J, Davis BM, et al. High-resolution
636 definition of the *Vibrio cholerae* essential gene set with hidden Markov model-based analyses of

- 637 transposon-insertion sequencing data. *Nucleic Acids Res.* 2013;41(19):9033-48. doi: Doi
638 10.1093/Nar/Gkt654. PubMed PMID: WOS:000326044700026.
- 639 18. Pritchard JR, Chao MC, Abel S, Davis BM, Baranowski C, Zhang YJJ, et al. ARTIST: High-
640 Resolution Genome-Wide Assessment of Fitness Using Transposon-Insertion Sequencing. *Plos Genet.*
641 2014;10(11). PubMed PMID: WOS:000345455200029.
- 642 19. Jacobs MA, Alwood A, Thaipisuttikul I, Spencer D, Haugen E, Ernst S, et al. Comprehensive
643 transposon mutant library of *Pseudomonas aeruginosa*. *P Natl Acad Sci USA.* 2003;100(24):14339-44.
644 PubMed PMID: ISI:000186803800105.
- 645 20. Chao MC, Abel S, Davis BM, Waldor MK. The design and analysis of transposon insertion
646 sequencing experiments. *Nat Rev Microbiol.* 2016;14:119-28.
- 647 21. Kentner D, Martano G, Callon M, Chiquet P, Brodmann M, Burton O, et al. *Shigella* reroutes
648 host cell central metabolism to obtain high-flux nutrient supply for vigorous intracellular growth. *P*
649 *Natl Acad Sci USA.* 2014;111(27):9929-34. doi: Doi 10.1073/Pnas.1406694111. PubMed PMID:
650 WOS:000338514800054.
- 651 22. Kleckner N, Bender J, Gottesman S. Uses of Transposons with Emphasis on Tn10. *Method*
652 *Enzymol.* 1991;204:139-80. PubMed PMID: WOS:A1991GN46900007.
- 653 23. Badarinarayana V, Estep PW, Shendure J, Edwards J, Tavazoie S, Lam F, et al. Selection
654 analyses of insertional mutants using subgenic-resolution arrays. *Nat Biotechnol.* 2001;19(11):1060-5.
655 PubMed PMID: ISI:000172002600022.
- 656 24. Chan K, Kim CC, Falkow S. Microarray-based detection of *Salmonella enterica* serovar
657 typhimurium transposon mutants that cannot survive in macrophages and mice. *Infect Immun.*
658 2005;73(9):5438-49. doi: Doi 10.1128/Iai.73.9.5438-5449.2005. PubMed PMID:
659 WOS:000231460000017.
- 660 25. Sansonetti PJ, Kopecko DJ, Formal SB. Involvement of a Plasmid in the Invasive Ability of
661 *Shigella-Flexneri*. *Infect Immun.* 1982;35(3):852-60. PubMed PMID: WOS:A1982NE19800013.
- 662 26. Yamamoto N, Nakahigashi K, Nakamichi T, Yoshino M, Takai Y, Touda Y, et al. Update on
663 the Keio collection of *Escherichia coli* single-gene deletion mutants. *Mol Syst Biol.* 2009;5. doi: Artn
664 335
665 Doi 10.1038/Msb.2009.92. PubMed PMID: WOS:000273359200008.
- 666 27. Kato JI, Hashimoto M. Construction of consecutive deletions of the *Escherichia coli*
667 chromosome. *Mol Syst Biol.* 2007;3:132. PubMed PMID: ISI:000249223700002.
- 668 28. Lamichhane G, Zignol M, Blades NJ, Geiman DE, Dougherty A, Grosset J, et al. A
669 postgenomic method for predicting essential genes at subsaturation levels of mutagenesis: Application
670 to *Mycobacterium tuberculosis*. *P Natl Acad Sci USA.* 2003;100(12):7213-8. doi: Doi
671 10.1073/Pnas.1231432100. PubMed PMID: WOS:000183493500058.
- 672 29. Inoue T, Shingaki R, Hirose S, Waki K, Mori H, Fukui K. Genome-wide screening of genes
673 required for swarming motility in *Escherichia coli* K-12. *J Bacteriol.* 2007;189(3):950-7. PubMed
674 PMID: WOS:000244112100030.

- 675 30. Skunca N, Bosnjak M, Krisko A, Panov P, Dzeroski S, Smuc T, et al. Phyletic Profiling with
676 Cliques of Orthologs Is Enhanced by Signatures of Paralogy Relationships. *Plos Comput Biol*.
677 2013;9(1). PubMed PMID: WOS:000314595600012.
- 678 31. Sparling PF. Kasugamycin Resistance . 30s Ribosomal Mutation with an Unusual Location on
679 Escherichia-Coli Chromosome. *Science*. 1970;167(3914):56-&. PubMed PMID:
680 WOS:A1970E985700019.
- 681 32. Corvaisier S, Bordeau V, Felden B. Inhibition of transfer messenger RNA aminoacylation and
682 trans-translation by aminoglycoside antibiotics. *J Biol Chem*. 2003;278(17):14788-97. PubMed PMID:
683 WOS:000182516100028.
- 684 33. Kim PD, Banack T, Lerman DM, Tracy JC, Camara JE, Crooke E, et al. Identification of a
685 novel membrane-associated gene product that suppresses toxicity of a TrfA peptide from plasmid RK2
686 and its relationship to the DnaA host initiation protein. *J Bacteriol*. 2003;185(6):1817-24. PubMed
687 PMID: ISI:000181448900008.
- 688 34. Berges H, Oreglia J, JosephLiauzun E, Fayet O. Isolation and characterization of a priB
689 mutant of Escherichia coli influencing plasmid copy number of Delta rop ColE1-type plasmids. *J*
690 *Bacteriol*. 1997;179(3):956-8. PubMed PMID: WOS:A1997WE44000050.
- 691 35. Bubunenko M, Baker T, Court DL. Essentiality of ribosomal and transcription antitermination
692 proteins analyzed by systematic gene replacement in Escherichia coli. *J Bacteriol*. 2007;189(7):2844-
693 53. PubMed PMID: WOS:000245842000030.
- 694 36. Diaz I, Pedersen S, Kurland CG. Effects of Miaa on Translation and Growth-Rates. *Mol Gen*
695 *Genet*. 1987;208(3):373-6. PubMed PMID: WOS:A1987J026300002.
- 696 37. Ambrosi C, Pompili M, Scribano D, Zagaglia C, Ripa S, Nicoletti M. Outer Membrane
697 Protein A (OmpA): A New Player in Shigella flexneri Protrusion Formation and Inter-Cellular
698 Spreading. *Plos One*. 2012;7(11). PubMed PMID: WOS:000311151900175.
- 699 38. Freudl R, Braun G, Hindennach I, Henning U. Lethal Mutations in the Structural Gene of an
700 Outer-Membrane Protein (Ompa) of Escherichia-Coli-K12. *Mol Gen Genet*. 1985;201(1):76-81.
701 PubMed PMID: WOS:A1985ARK1200013.
- 702 39. Ma D, Cook DN, Alberti M, Pon NG, Nikaido H, Hearst JE. Genes Acra and Acrb Encode a
703 Stress-Induced Efflux System of Escherichia-Coli. *Mol Microbiol*. 1995;16(1):45-55. PubMed PMID:
704 WOS:A1995QU56400005.
- 705 40. Kleckner N. Transposable elements in prokaryotes. *Annu Rev Genet*. 1981;15:341-404.
706 PubMed PMID: Medline:6279020.
- 707 41. Lundblad V, Taylor AF, Smith GR, Kleckner N. Unusual alleles of recB and recC stimulate
708 excision of inverted repeat transposons Tn10 and Tn5. *P Natl Acad Sci USA*. 1984;81(3):824-8.
709 PubMed PMID: Medline:6322169.
- 710 42. Chan SH, Lau A, Lei V, Woo J. Effects of recB, recC and recF mutations on Tn10
711 Transposition in Escherichia coli. *ournal of Experimental Microbiology and Immunology*. 2006;9:75-
712 80.
- 713 43. Attfield PV, Benson FE, Lloyd RG. Analysis of the ruv locus of Escherichia coli K-12 and
714 identification of the gene product. *J Bacteriol*. 1985;164(1):276-81. PubMed PMID: Medline:2995311.

- 715 44. Keseler IM, Mackie A, Peralta-Gil M, Santos-Zavaleta A, Gama-Castro S, Bonavides-
716 Martinez C, et al. EcoCyc: fusing model organism databases with systems biology. *Nucleic Acids Res.*
717 2013;41(D1):D605-D12. PubMed PMID: WOS:000312893300086.
- 718 45. Kawakami K, Ito K, Nakamura Y. Differential Regulation of 2 Genes Encoding Lysyl-
719 Transfer Rna-Synthetases in Escherichia-Coli - Lysu-Constitutive Mutations Compensate for a Lyss
720 Null Mutation. *Mol Microbiol.* 1992;6(13):1739-45. PubMed PMID: WOS:A1992JC20400003.
- 721 46. Zhang XB, Liu H, Yang F, Yang J, Xue Y, Dong J, et al. Comparative genome analysis of
722 deleted genes in Shigella flexneri 2a strain 301. *Chinese Sci Bull.* 2003;48(9):846-52. PubMed PMID:
723 WOS:000183995400002.
- 724 47. Liao MK, Gort S, Maloy S. A cryptic proline permease in Salmonella typhimurium.
725 *Microbiol-Uk.* 1997;143:2903-11. PubMed PMID: WOS:A1997XW07300007.
- 726 48. Wong MS, Wu S, Causey TB, Bennett GN, San KY. Reduction of acetate accumulation in
727 Escherichia coli cultures for increased recombinant protein production. *Metab Eng.* 2008;10(2):97-108.
728 PubMed PMID: WOS:000261595200004.
- 729 49. Monk JM, Charusanti P, Aziz RK, Lerman JA, Premyodhin N, Orth JD, et al. Genome-scale
730 metabolic reconstructions of multiple Escherichia coli strains highlight strain-specific adaptations to
731 nutritional environments. *P Natl Acad Sci USA.* 2013;110(50):20338-43. PubMed PMID:
732 WOS:000328061700082.
- 733 50. Hooper SD, Berg OG. Gene import or deletion: A study of the different genes in Escherichia
734 coli strains K12 and O157 : H7. *J Mol Evol.* 2002;55(6):734-44. PubMed PMID:
735 WOS:000179503100012.
- 736 51. Kothary MH, Babu US. Infective dose of foodborne pathogens in volunteers: A review. *J*
737 *Food Safety.* 2001;21(1):49-73. PubMed PMID: WOS:000169172400004.
- 738 52. Goldberg MB, Theriot JA. Shigella-Flexneri Surface Protein Icsa Is Sufficient to Direct Actin-
739 Based Motility. *P Natl Acad Sci USA.* 1995;92(14):6572-6. PubMed PMID: WOS:A1995RG73600073.
- 740 53. Caetanoanollés G. Amplifying DNA with Arbitrary Oligonucleotide Primers. *Pcr Meth Appl.*
741 1993;3(2):85-94. PubMed PMID: WOS:A1993MD19100001.
- 742 54. Davis MPA, van Dongen S, Abreu-Goodger C, Bartonicek N, Enright AJ. Kraken: A set of
743 tools for quality control and analysis of high-throughput sequence data. *Methods.* 2013;63(1):41-9.
- 744 55. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods.*
745 2012;9(4):357-U54. PubMed PMID: WOS:000302218500017.
- 746 56. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and
747 population genetical parameter estimation from sequencing data. *Bioinformatics.* 2011;27(21):2987-93.
748 doi: Doi 10.1093/Bioinformatics/Btr509. PubMed PMID: WOS:000296099300009.
- 749 57. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence
750 Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25(16):2078-9. doi: Doi
751 10.1093/Bioinformatics/Btp352. PubMed PMID: ISI:000268808600014.
- 752 58. Halling SM, Kleckner N. A Symmetrical 6-Base-Pair Target Site Sequence Determines Tn10
753 Insertion Specificity. *Cell.* 1982;28(1):155-63. doi: Doi 10.1016/0092-8674(82)90385-3. PubMed
754 PMID: WOS:A1982MY42100020.

- 755 59. Bambom O. seqLogo: An R package for plotting DNA sequence logos. 2007.
- 756 60. Wei J, Goldberg MB, Burland V, Venkatesan MM, Deng W, Fournier G, et al. Complete
757 genome sequence and comparative genomics of *Shigella flexneri* serotype 2a strain 2457T (vol 71, pg
758 2775, 2003). *Infect Immun*. 2003;71(7):4223-. PubMed PMID: ISI:000183797200078.
- 759 61. Silander OK, Ackermann M. The constancy of gene conservation across divergent bacterial
760 orders. *BMC Research Notes*. 2009;2:2.
- 761 62. Wall DP, Fraser HB, Hirsh AE. Detecting putative orthologs. *Bioinformatics*.
762 2003;19(13):1710-1. PubMed PMID: ISI:000185310600016.
- 763 63. Serres MH, Riley M. MultiFun, a multifunctional classification scheme for *Escherichia coli*
764 K-12 gene products. *Microb Comp Genomics*. 2000;5(4):205-22. PubMed PMID: 11471834.
- 765

766 **Table and Figure legends**

767 **Table 1. Genes artefactually inferred as essential in *Shigella*.** Listed here are those genes
768 that were likely inferred as essential largely due to the gene disruptions having direct effects
769 on (1) antibiotic resistance, (2) successful transposition events, (3) differences between the
770 growth conditions of the *E. coli* and *Shigella* essentiality studies, or (4) systematically
771 different effects of gene disruption versus gene deletion.

772 **Table 2. Genes inferred as uniquely essential in *Shigella*.** Listed here are genes inferred as
773 essential in *Shigella*, but which have orthologous *E. coli* deletion genotypes that exhibit
774 robust growth (greater than 0.75 OD600 after 22h growth in LB). The genes in the *rfb* operon
775 have no orthologues in *E. coli* K12 (see main text).

776 **Fig 1. Histograms of distances between inserts on the chromosome for (A) the *Shigella***
777 **chromosome and (B) the *Shigella* virulence plasmid.** The median distance between inserts
778 is indicated by the dotted line. The insets in (A) and (B) indicate the slight but detectable
779 biases in transposon insert location using a weight matrix motif. The reverse cumulative plots
780 show the observed fractions of distances between inserts for the chromosome (C) and the
781 plasmid (D). In blue, the observed frequencies are plotted. In black, the expected frequencies
782 are plotted, given a geometric distribution (negative binomial with the number of successes
783 set to one) of inter-insert distances (see main text). For both the chromosome and the plasmid,
784 there are considerably more large regions uninterrupted by transposons than one would expect
785 given the geometric null model, observed as a shift of the curve to the right.

786 **Fig 2. Orthologous genes known to be essential in *E. coli* are also essential in *Shigella*.** A
787 region of the *Shigella* chromosome is shown, with genes whose orthologues are known to be
788 essential for growth in *E. coli* (coloured in white) [1, 27], or non-essential (coloured in grey).
789 The unique locations of transposon insertions are plotted as vertical black segments. In the
790 genome region shown here, none of the genes essential in *E. coli* have orthologues that are
791 interrupted in *Shigella*.

792 **Fig 3. Differences in essentiality classification between *E. coli* K12 and *Shigella*.** (A)
793 Cumulative distributions showing the mean distances between inserts for ORFs depending on
794 whether their orthologues are known to be non-essential (black curve) or essential (blue
795 curve) in *E. coli*. All ORFs that are completely uninterrupted by transposons have been
796 plotted at the very right end of the x-axis. The dotted vertical line indicates the cut-off that we
797 used to delineate essentiality in *Shigella* (a mean distance between transposons of 250 bp or
798 more). The 17 blue points to the left of the dotted vertical line indicate ORFs that are essential
799 in *E. coli* but not *Shigella* by our metric. These are likely to be false positives (i.e. non-
800 essential in both *Shigella* and *E. coli*), as all have inter-insert distances greater than 100 bp

801 (see main text). Black points to the right of the dotted vertical line indicate ORFs that we
802 classify as essential in *Shigella* but not in *E. coli*. Many of these ORFs have *E. coli*
803 orthologues whose deletion genotypes exhibit robust growth, suggesting that their essentiality
804 status has changed. (B) A Venn diagram showing the overlap between essential orthologous
805 ORFs in *E. coli* and *Shigella*.

806 **Fig 4. Orthologous gene pairs that are non-essential in *E. coli* but inferred as essential in**
807 ***Shigella* (blue) tend to exhibit low growth yields in *E. coli*.** ORFs that we infer to be
808 uniquely essential in *Shigella* consistently have *E. coli* orthologues with low growth
809 phenotypes in LB media after 22 hours (apparent as a strong leftward shift in the cumulative
810 curve). For genes inferred as uniquely essential in *Shigella*, 34% of the orthologous *E. coli*
811 deletion genotypes exhibit low growth yields (less than 0.5 OD600 after 22 hours of growth
812 in LB). For genes we classified as non-essential in *Shigella* and *E. coli* only 3.7% exhibit low
813 growth yields. Thus, some genes we infer as essential in *Shigella* may not be strictly essential,
814 but instead be required for robust growth. Despite this enrichment for low-growth phenotypes,
815 there are many genes which we infer as essential in *Shigella*, but which have *E. coli*
816 orthologues whose deletion genotypes exhibit robust growth (OD600 greater than 0.75 after
817 22 hours growth in LB).

818 **Fig 5. Transposon disruption of *Shigella* genes with transport-related functions are**
819 **more likely to be inferred as essential compared to clean deletions of similarly**
820 **functioning genes in *E. coli*.** We classified genes according to function using the MultiFun
821 functional classification system [63]. For any category containing more than ten essential *E.*
822 *coli* genes, we also calculated the number of *Shigella*-essential genes. As expected, most
823 categories show a relative excess of *Shigella*-essential genes, as we inferred approximately
824 50% more genes as being essential in *Shigella* versus *E. coli* (**Fig. 3B**). However, two
825 functional categories show a clear excess above this level: substrate transport and active
826 transport, showing a 3- and 2.1-fold increased probability of inferring a gene as being
827 essential in *Shigella* as opposed to *E. coli*. This provides evidence that genes in these
828 functional categories may be more likely to be inferred as artefactually essential. For each
829 functional category (y-axis), we show the number of genes in that category (to the right of
830 each bar); the number of genes found to be essential in *E. coli* (within each bar); and the level
831 of enrichment of essential genes in *Shigella* (x-axis).

832 **Fig 6. Three genes involved in proline biosynthesis (*proABC*) appear uniquely essential**
833 **in *Shigella*.** The orthologous *E. coli* deletion strains exhibit robust growth (OD600 greater
834 than 0.75 after 22 hours growth in LB), but are essential by our criteria. *proA* and *proC*
835 completely lack transposon insertions, while *proB* contains only two insertions near the 3' end,

836 which leaves approximately 70% of the gene intact, including the entire kinase and substrate-
837 binding domain.

838 **Fig 7. Both elongation factor paralogues *tufA* and *tufB* appear differentially essential in**
839 ***Shigella* as compared to *E. coli*.** The orthologous *E. coli* deletion strains of *tufA* and *tufB*
840 exhibit robust growth (OD600 of 0.72 and 0.78 after 22 hours in LB), but are essential by our
841 criteria. Both genes contain insertions only at the 5' or 3' ends of the genes. Genes that are
842 essential in both *E. coli* and *Shigella* are coloured in white. Those inferred as being essential
843 in *Shigella* but for which the orthologous deletion genotypes exhibit robust growth in *E. coli*
844 are indicated in blue. Genes inferred as essential in *Shigella* and which do not exhibit robust
845 growth in *E. coli* are coloured in light blue. tRNA genes are indicated in dark grey.

846 **Fig 8. The region of the genome containing the *rfb* operon is largely uninterrupted by**
847 **transposon insertions.** *rfbI*, *rfc*, *rfbG*, and *rfbF* are completely uninterrupted by transposon
848 insertions; *rfbE* is uninterrupted over 90% of its length. None of these genes have orthologous
849 counterparts in *E. coli* K12 due to a lateral transfer event that occurred at this locus (see main
850 text). This operon encodes genes active in O-antigen biosynthesis.

851 **S1 Table.** Full table of gene characteristics and orthologue relationships used in the analyses.

852 **S2 Table.** List of chromosomal SNPs and indels observed in the *Shigella* strain used here that
853 differ from the GenBank sequence NC004741.

854 **S3 Table.** List of plasmid SNPs and indels observed in the *Shigella* strain used here that are
855 different from the GenBank sequence of the *Shigella flexneri* 2a strain 301 virulence plasmid
856 pCP301 (NC_004851).

857 **S4 Table.** Table listing tRNA genes and the number of insertions in each.

858 **S1 Fig. Distribution of transposon insertions across the genome.** We observed little bias
859 on the chromosomal level of insert locations.

860 **S2 Fig. ROC curves showing the predictive power of various features.** To select a feature
861 that was the best predictor of essentiality in *E. coli* orthologues, used only ORFs that we had
862 data on essentiality from both the Keio and PEC studies. We then selected transposon
863 insertion patterns that most closely match the essentiality delineations in these studies.
864 Specifically, we selected the feature that maximized the number of true positive “essential”
865 genes (maximizing the sensitivity) while minimizing the number of FP (maximizing
866 specificity). We selected from eleven (non-independent) features shown here: (1) the total
867 number of insertions; (2) the mean number of bp between insertions; (3) the median number
868 of bp between insertions; (4) the number of bp in the 5' end preceding the first insertion; (5)
869 the number of bp in the 5' end preceding the first insertion relative to the total bp in the gene;

870 (6) the number of bp in the 5' end preceding the second insertion; (7) the number of bp in the
871 5' end preceding the second insertion relative to the total bp in the gene; (8) the number of bp
872 in the longest uninterrupted stretch of the gene; (9) the number of bp in the longest
873 uninterrupted stretch of the gene relative to the total length of the gene, and (10) the number
874 of bp in the longest stretch of the gene interrupted by at most one insertion; (11) the number
875 of bp in the longest stretch of the gene interrupted by at most one insertion, relative to the
876 total length of the gene. See the **Methods** section for more details of this analysis.

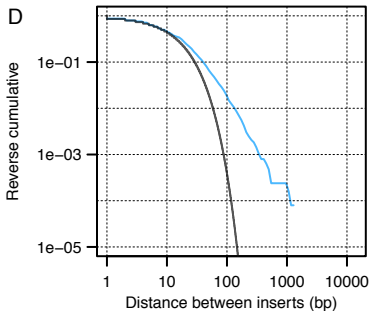
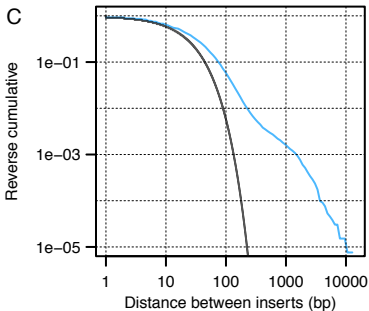
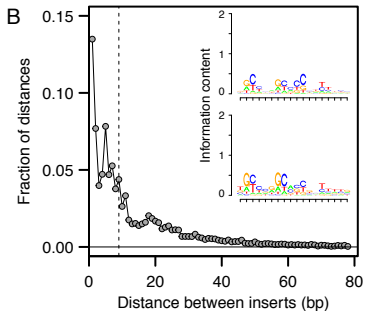
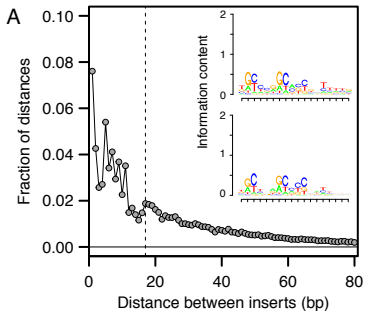
877 **S3 Fig. Analogous plots to that shown in Fig. 3, for growth in minimal glucose MOPS**
878 **media after (A) 24 and (B) 48 hours.** In both cases, we find that the shift is less pronounced
879 than that observed for LB.

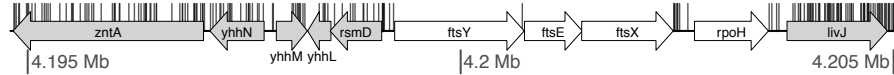
880 **S4 Fig. Schematic of the primer positions used for Illumina sequencing of transposon**
881 **insertions.**

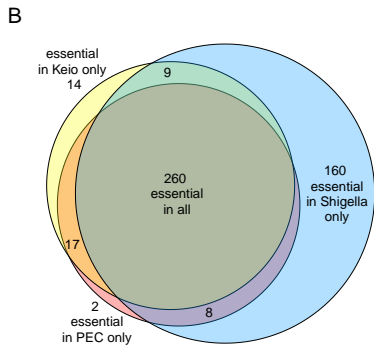
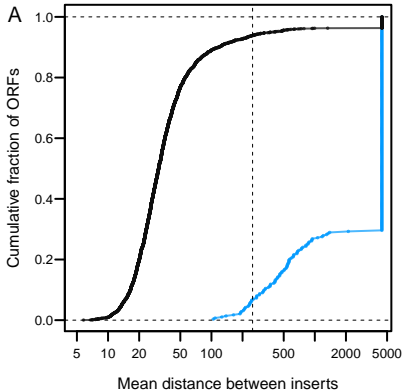
882 **S5 Fig. Inferred fragment lengths of perfectly mapped reads across several genomic**
883 **regions.** For each plot, the inferred fragment lengths are arranged by increasing length
884 (ranked on the y-axis). Thus, very long fragments are present at the top of the y-axis. Most
885 fragments have lengths between 100bp and 400bp; a small number have lengths over 1000bp
886 or more. It is very likely that these are not the true insert sizes, but appear that way because of
887 large scale deletions in our *Shigella* clone compared to the clone present in the NCBI genome
888 database; see **Methods** for more details. (A) A region of the chromosome in which a
889 complicated series of rearrangements has occurred, leading to paired end reads perfectly
890 matching to different locations in this region. 45 mapped read pairs span more than 1.5 Kbp, a
891 size that is not concordant with the majority of insert sizes. (B) A genomic region where an
892 approximately 10Kbp deletion occurred, removing a region containing the *yeaKLMNOP*
893 operon. 92 mapped read pairs span more than 8.5Kbp. This region is flanked by two IS
894 elements. (C) A region where an approximately 4Kbp deletion occurred, removing two genes
895 with no *E. coli* K12 orthologues. 68 mapped read pairs span more than 4 Kbp, and again this
896 region is flanked by two IS elements. (D) A genomic region where an approximately 2Kbp
897 deletion occurred, removing *yhdW*. 244 mapped read pairs span more than 2Kbp, and the
898 region is flanked by two IS elements. (E) A deletion in the region of the chromosome
899 containing *S4145 (viaN)*. 232 mapped read pairs spanned more than 1.8 Kbp, and this region
900 is also flanked by two IS elements. (F) A region of the chromosome containing the *rfb* operon.
901 Most of the genes within this operon are uninterrupted by transposons. However, we find no
902 evidence that this is due to a deletion of this region in our *Shigella* clone, as we find no reads
903 mapping across the region; a small number of reads map within the region; and the closest IS
904 elements are 15 Kb upstream of *rfbJ* and 20 Kb downstream of *rfbA*. The genes in this operon
905 have no orthologues in *E. coli* K12.

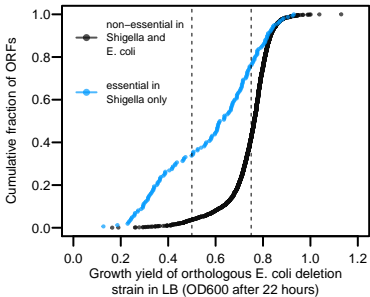
906 **S6 Fig. Functional categories of non-essential *E. coli* genes that are enriched (or**
907 **depleted) for essential *Shigella* genes.** We classified genes according to function using the
908 MultiFun classification system [63]. For this analysis we considered only genes that are non-
909 essential in *E. coli*. We find that genes uniquely essential in *Shigella* are enriched in some
910 functional categories. For example, of the 27 ribosomal proteins identified as non-essential in
911 *E. coli*, we identify approximately 45% as being essential in *Shigella*. In contrast, of the 565
912 membrane proteins identified as non-essential in *E. coli*, we find that less than 10% are
913 essential in *Shigella*. Thus, ORFs uniquely essential in *Shigella* are far more likely to function
914 in the ribosome than one would expect. The number of non-essential *E. coli* genes is indicated
915 above each bar; the probability of finding the level of enrichment (or depletion) that we
916 observe in each secondary category is indicated for cases in which this probability is less than
917 0.025, using a Fisher exact test. (A) All non-essential genes in *E. coli*. (B) An identical
918 analysis excluding all non-essential genes in *E. coli* that exhibit very low growth yields
919 (OD600 less than 0.5 after 22 hours of growth LB). In both cases, the only subcategory
920 notably enriched for essential genes is that containing ribosomal proteins. The only categories
921 appreciably depleted for genes with essential function are genes with function in RNA
922 processing and to some extent, energy production.

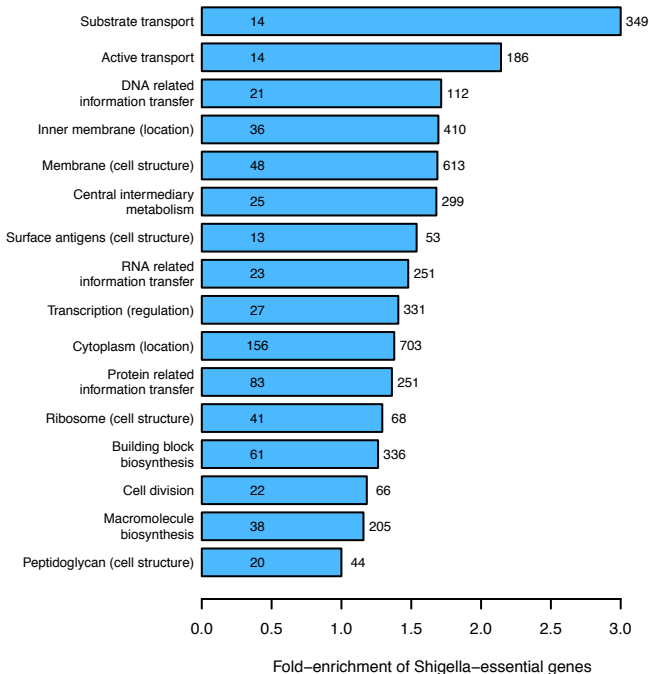
923 **S7 Fig. Transposon insertion locations across the entire *Shigella* chromosome. Each**
924 **insertion site is indicated by a vertical red line.** ORFs are indicated in light green; rRNAs
925 in dark green; and tRNAs in dark blue. The annotation is taken from the GenBank sequence
926 of *Shigella flexneri* 2a 2457T.

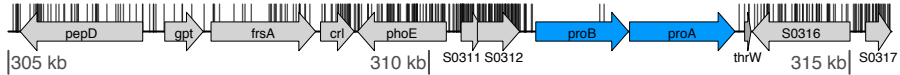


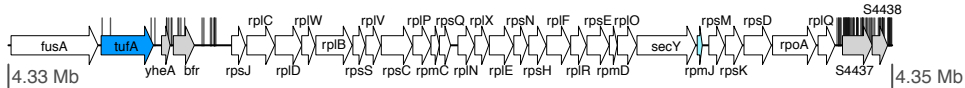
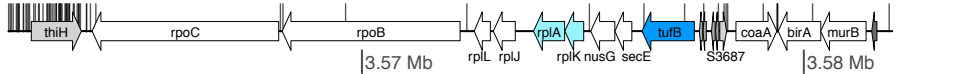




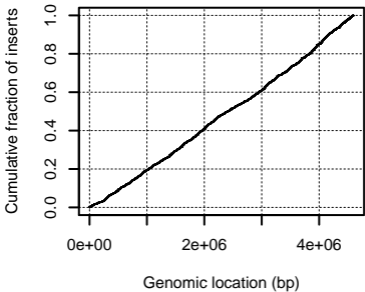


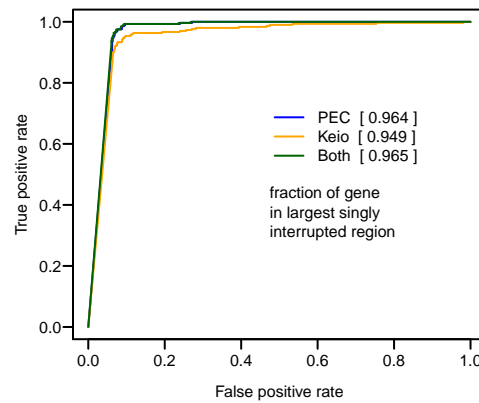
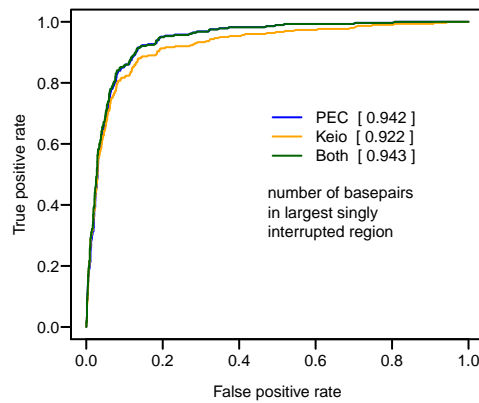
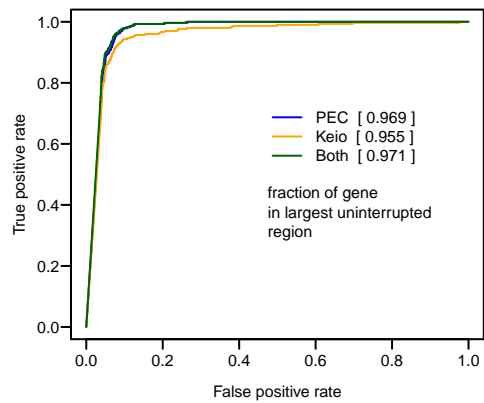
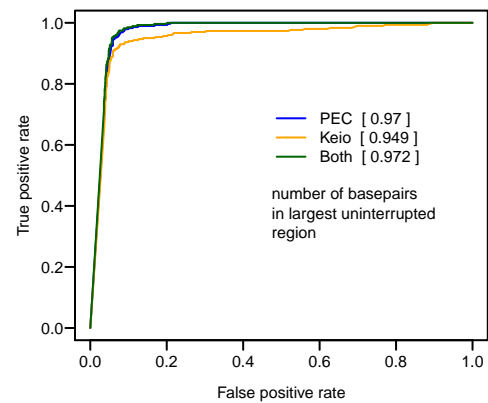
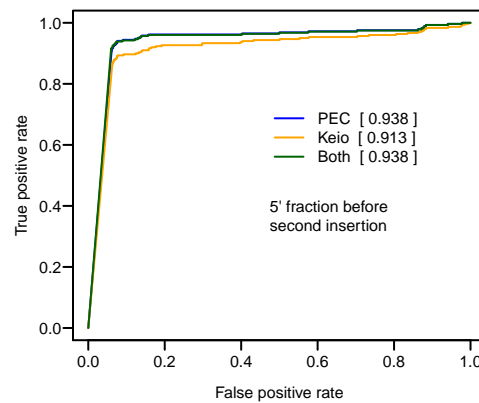
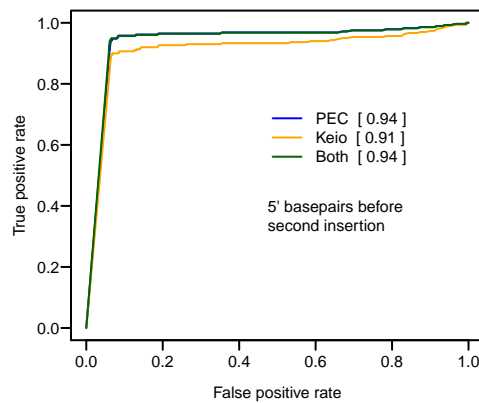
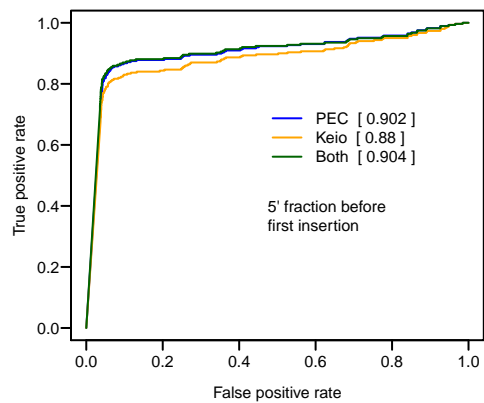
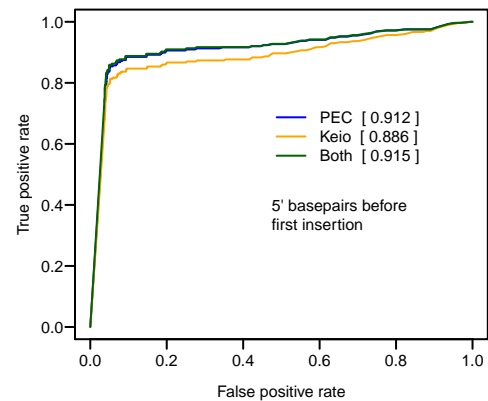
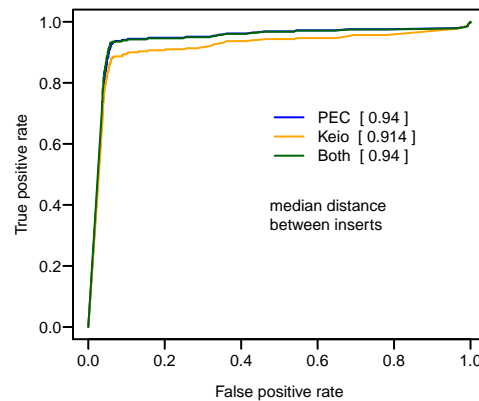
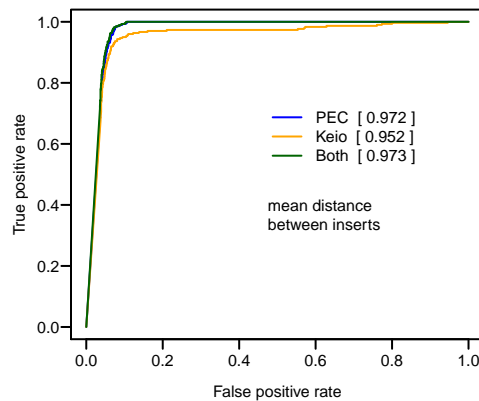
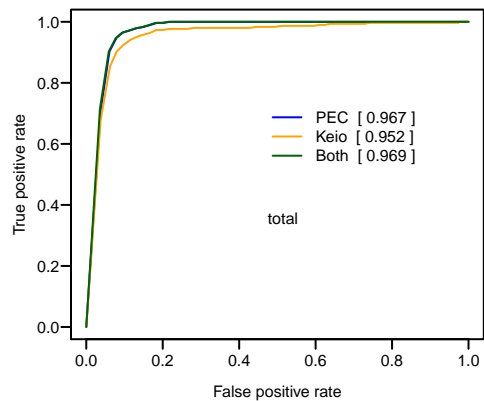


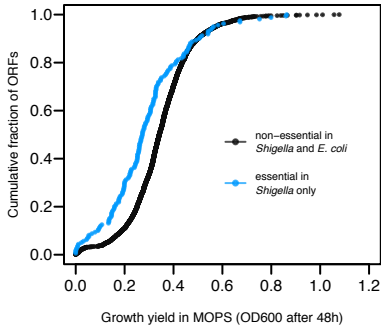
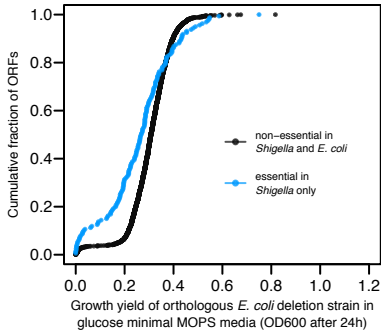


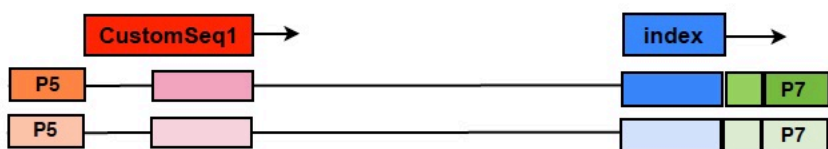
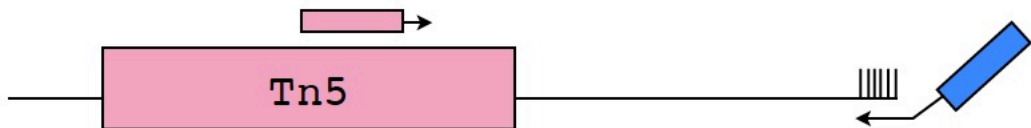












double stranded PCR product,
gel purify band/smear at (300-400bp)
(insert at 200-300bp)

1: Arbitrary PCR

using phusion PCR enzyme, so that no A's added at 3' end

1_transp_deep_seq_F1: TCTATCGCCTTCTTGACGAG

1_transp_deep_seq_R2: GTGACTGGAGTTCAGACGTGT GCTCTTCCGATCT NNNNNNNNNggtgc
 1_transp_deep_seq_R3: GTGACTGGAGTTCAGACGTGT GCTCTTCCGATCT NNNNNNNNNgatata
 1_transp_deep_seq_R4: GTGACTGGAGTTCAGACGTGT GCTCTTCCGATCT NNNNNNNNNnagtac

2: Nested PCR (addition of P7 and P5)

using phusion PCR enzyme, so that no A's added at 3' end

2_transp_deep_seq_FP5: AATGATACGGCGACCACCGAGATCT actggtcggcg CATTAGGGGATTCATCAG

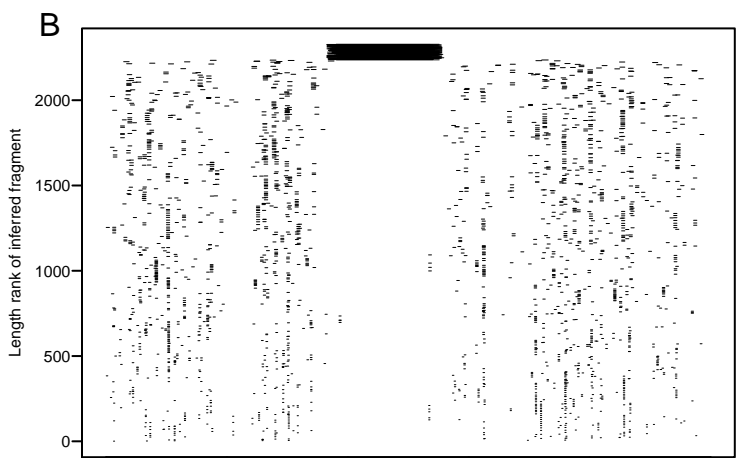
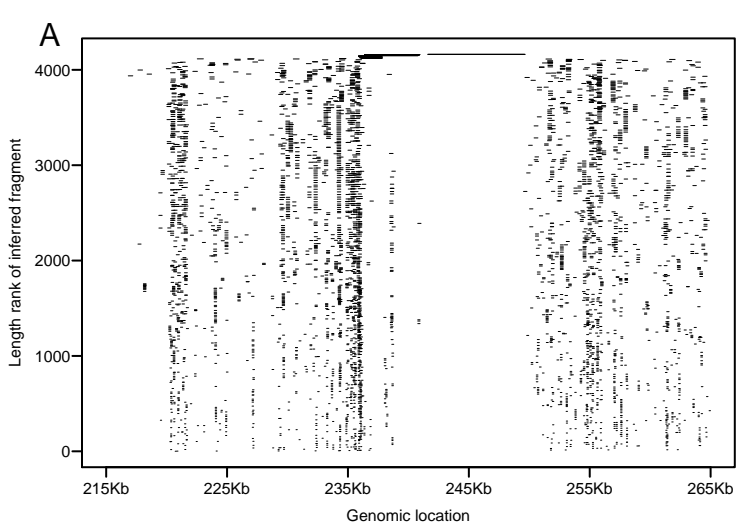
Rev_P7 6_02
 CAAGCAGAAGACGGCATAACGAGAT agtctt GTGACTGGAGTTCAGACGTGT
 CAAGCAGAAGACGGCATAACGAGAT gtatct GTGACTGGAGTTCAGACGTGT
 CAAGCAGAAGACGGCATAACGAGAT tcatgg GTGACTGGAGTTCAGACGTGT
 CAAGCAGAAGACGGCATAACGAGAT cgcgac GTGACTGGAGTTCAGACGTGT
 CAAGCAGAAGACGGCATAACGAGAT acgata GTGACTGGAGTTCAGACGTGT
 CAAGCAGAAGACGGCATAACGAGAT tataga GTGACTGGAGTTCAGACGTGT

3: Sequencing on Illumina chip

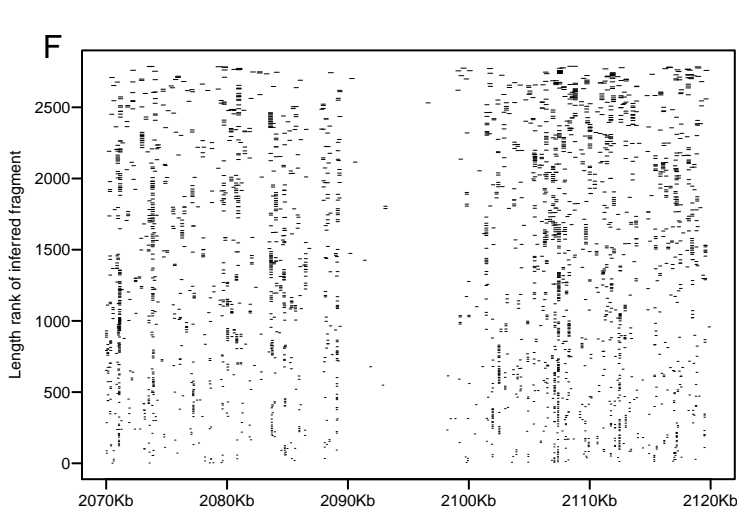
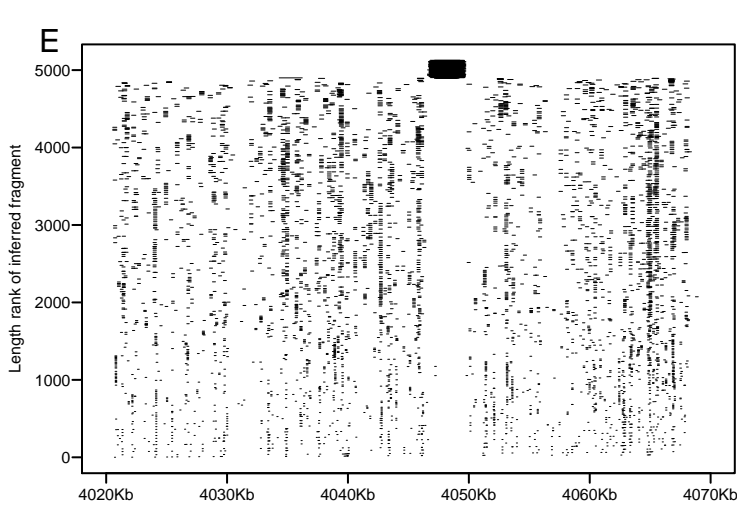
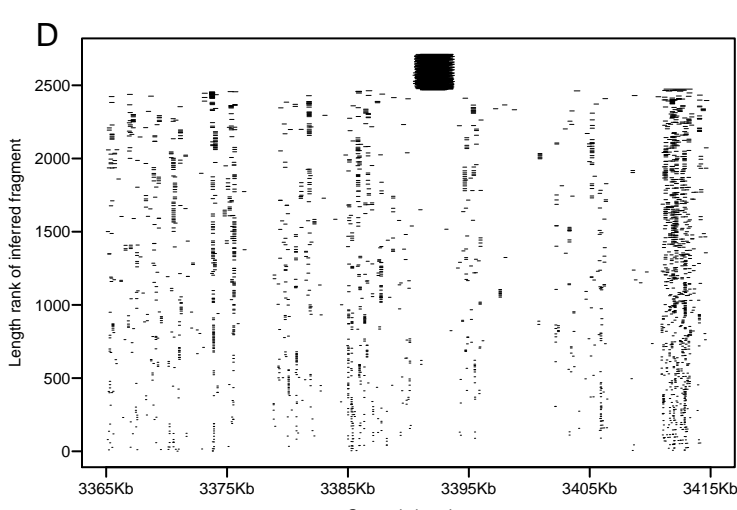
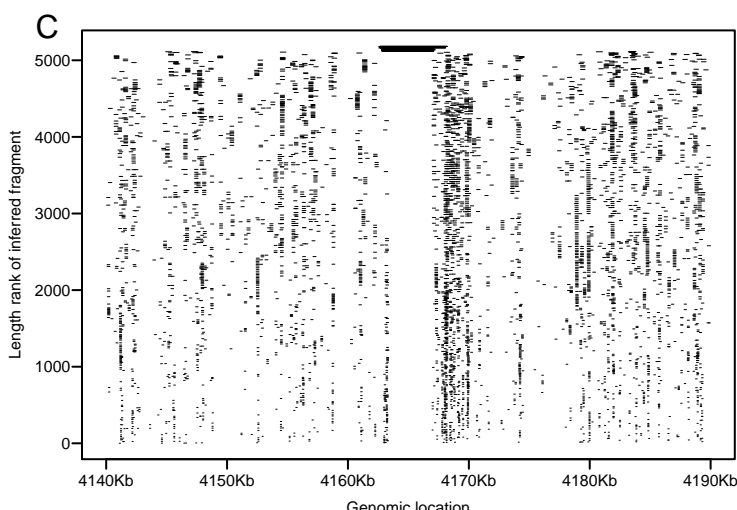
transp_deep_seq_Custom seq 1: actggtcggcgCATTAGGGGATTCATCAG

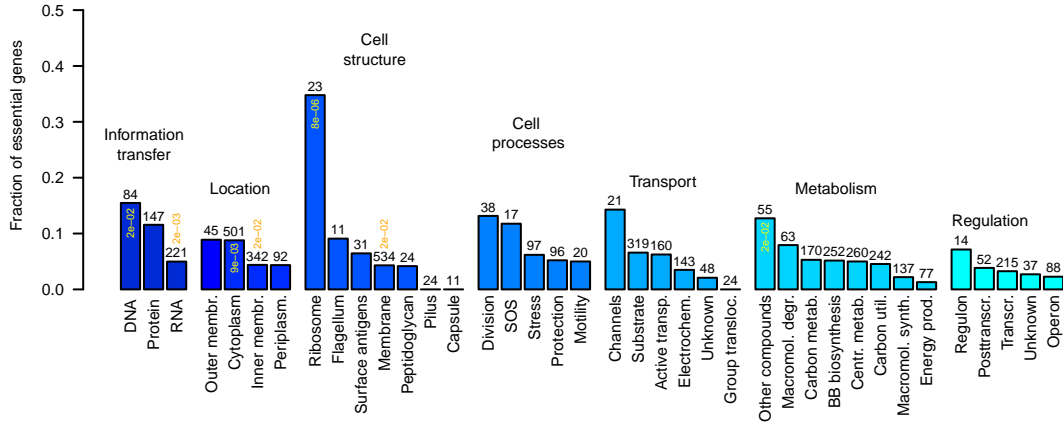
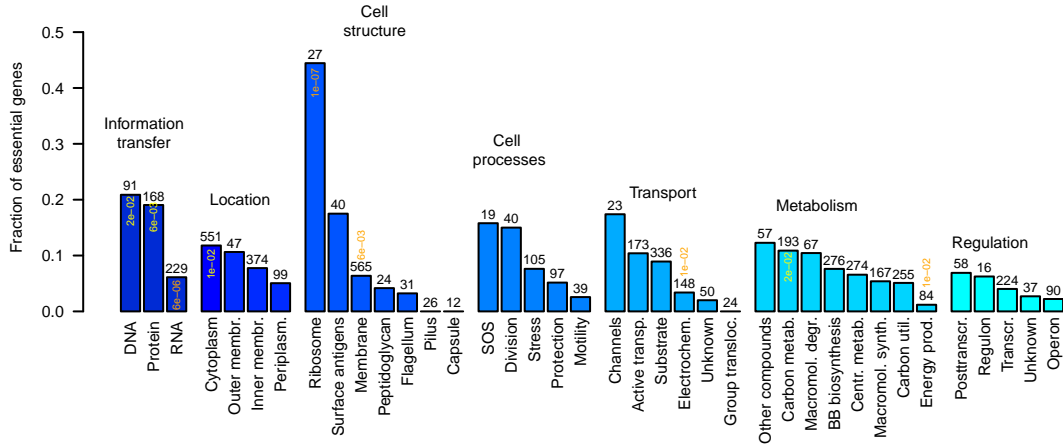
4: Barcode is read on Illumina chip

transp_deep_seq_Index seq: GATCGGAAGAGCACACGTCTGAACTCCAGTCAC

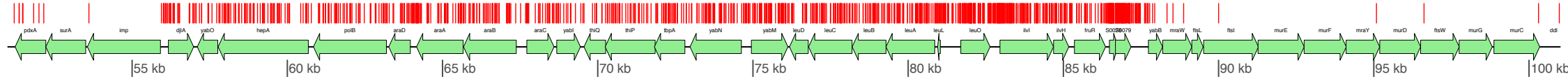


bioRxiv preprint doi: <https://doi.org/10.1101/038869>; this version posted January 11, 2016. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



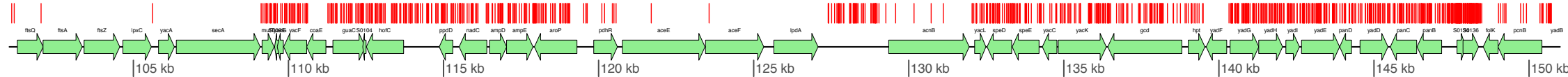


51 – 101 Kb



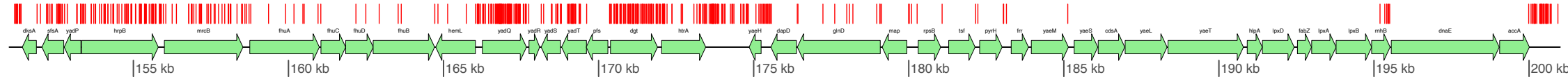
5 kb

101 – 151 Kb



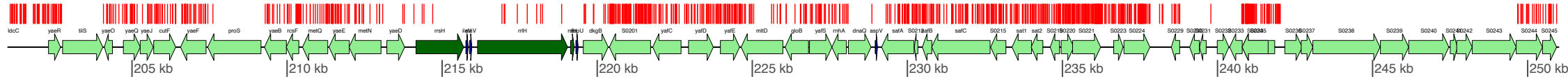
5 kb

151 – 201 Kb



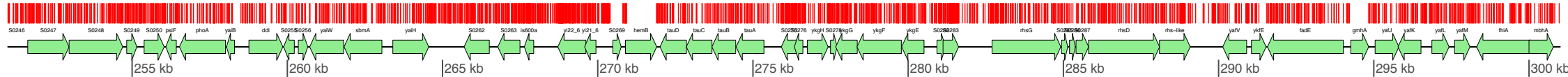
5 kb

201 – 251 Kb



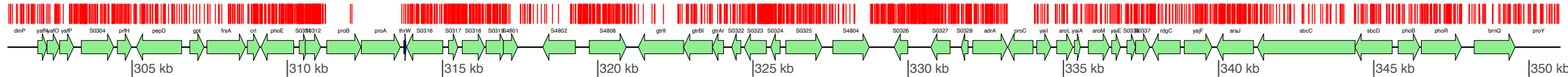
5 kb

251 – 301 Kb



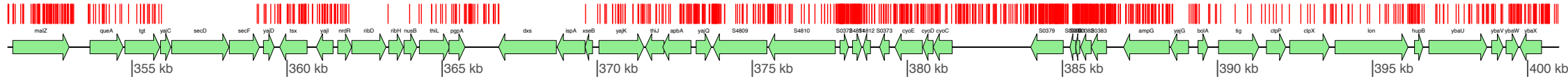
5 kb

301 – 351 Kb



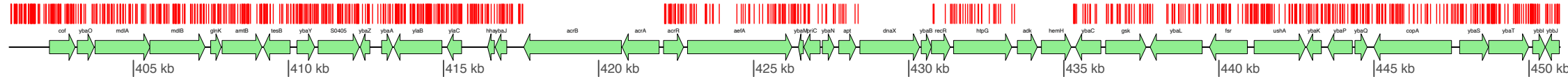
5 kb

351 – 401 Kb



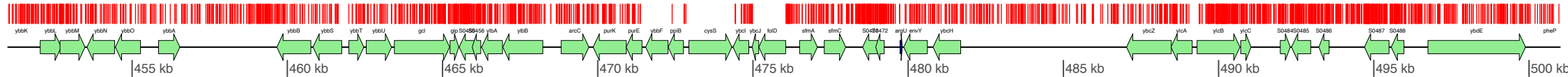
5 kb

401 – 451 Kb



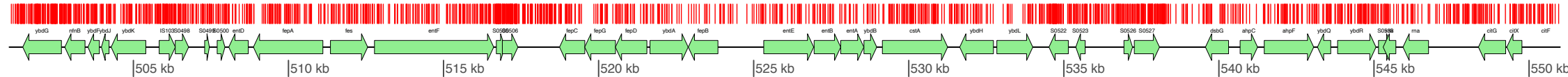
5 kb

451 – 501 Kb



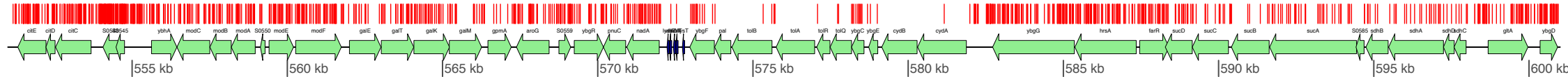
5 kb

501 – 551 Kb



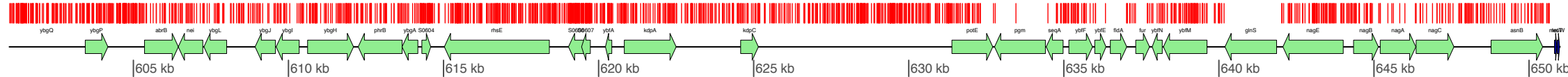
5 kb

551 – 601 Kb



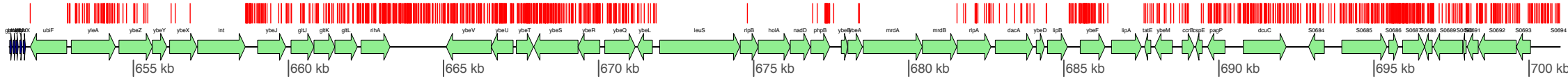
5 kb

601 – 651 Kb



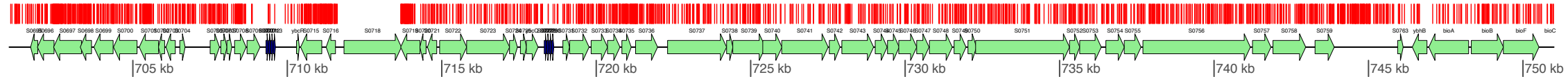
5 kb

651 – 701 Kb



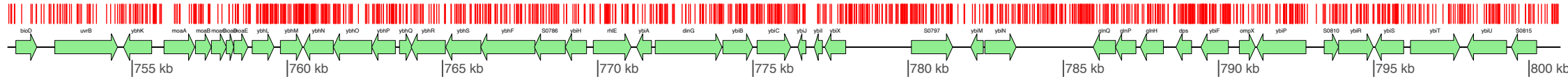
5 kb

701 – 751 Kb



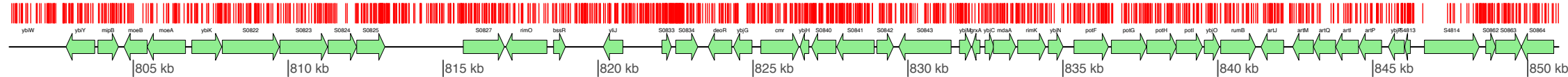
5 kb

751 – 801 Kb



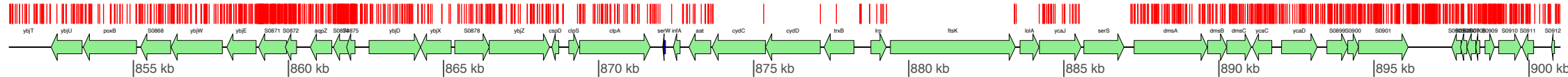
5 kb

801 – 851 Kb



5 kb

851 – 901 Kb



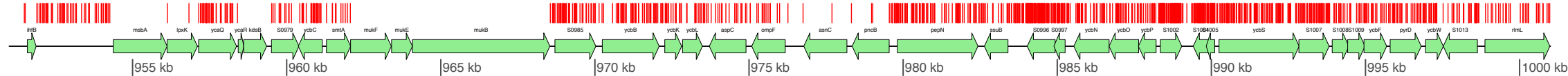
5 kb

901 – 951 Kb



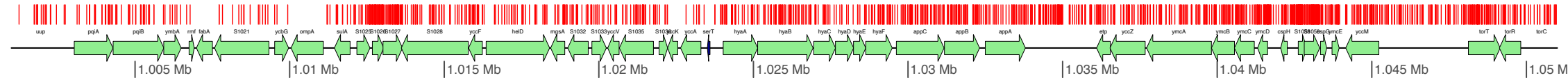
5 kb

951 – 1001 Kb



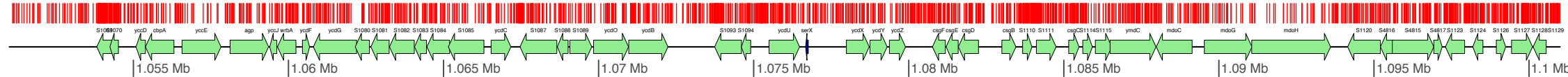
5 kb

1001 – 1051 Kb



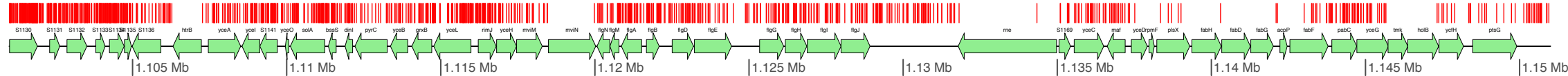
5 kb

1051 – 1101 Kb



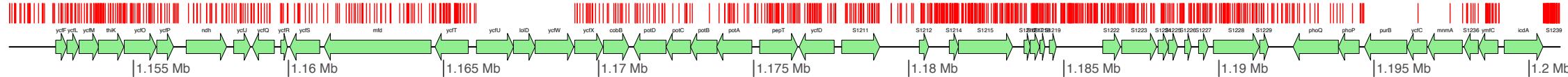
5 kb

1101 – 1151 Kb



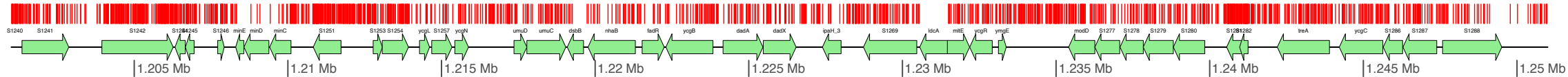
5 kb

1151 – 1201 Kb



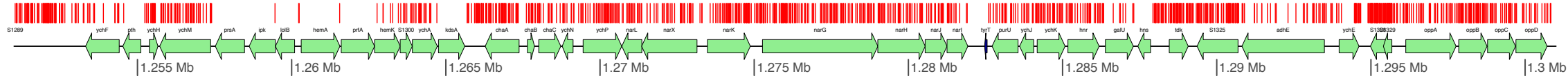
5 kb

1201 – 1251 Kb



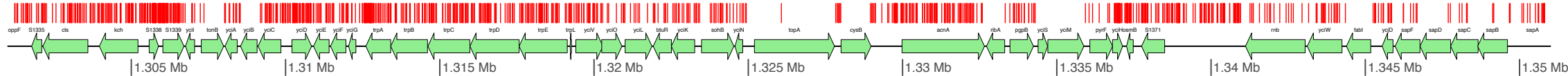
5 kb

1251 – 1301 Kb



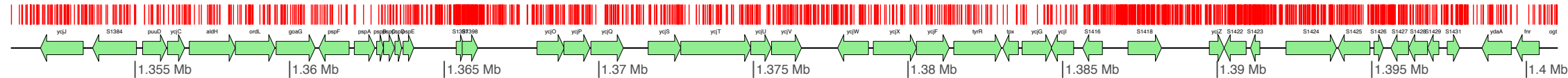
5 kb

1301 – 1351 Kb



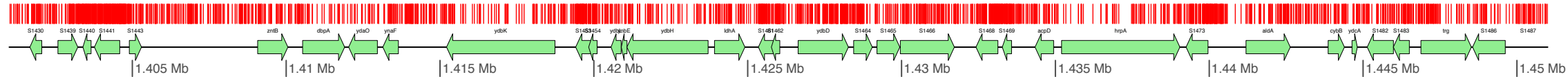
5 kb

1351 – 1401 Kb



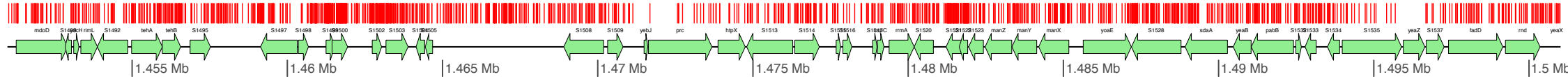
5 kb

1401 – 1451 Kb



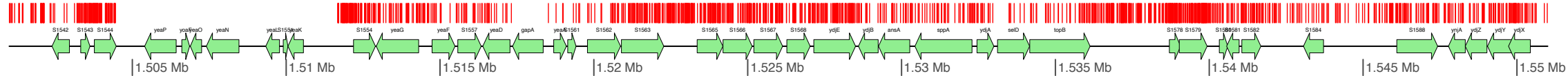
5 kb

1451 – 1501 Kb



5 kb

1501 – 1551 Kb



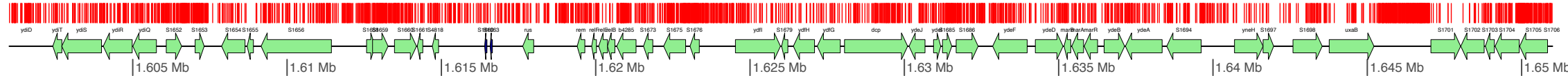
5 kb

1551 – 1601 Kb



5 kb

1601 – 1651 Kb



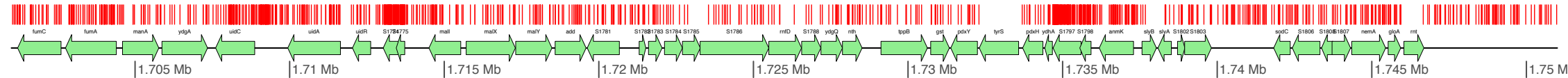
5 kb

1651 – 1701 Kb



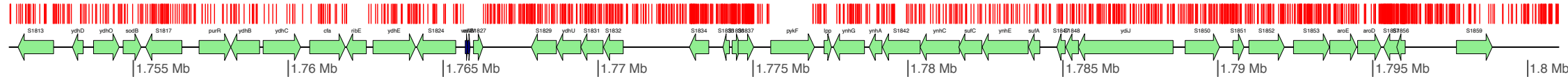
5 kb

1701 – 1751 Kb



5 kb

1751 – 1801 Kb



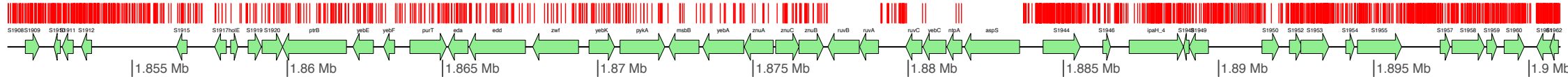
5 kb

1801 – 1851 Kb



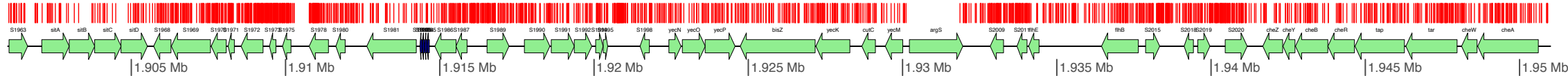
5 kb

1851 – 1901 Kb



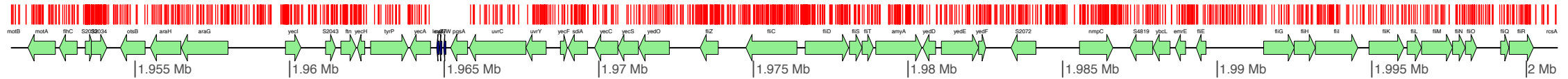
5 kb

1901 – 1951 Kb



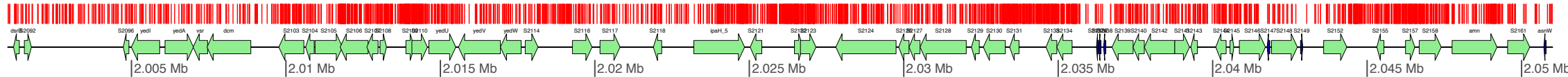
5 kb

1951 – 2001 Kb



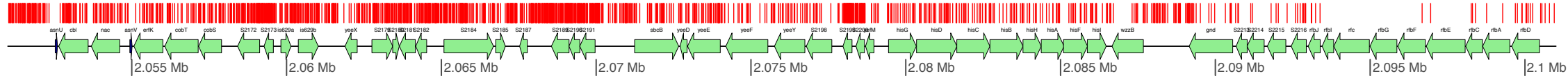
5 kb

2001 – 2051 Kb



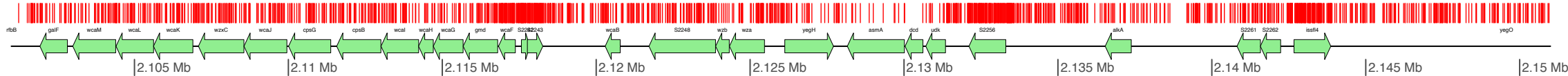
5 kb

2051 – 2101 Kb



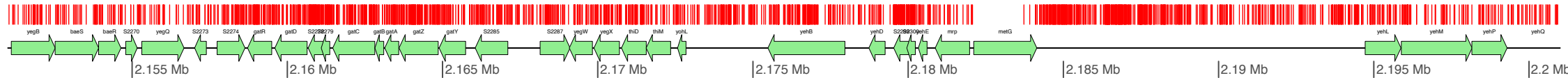
5 kb

2101 – 2151 Kb



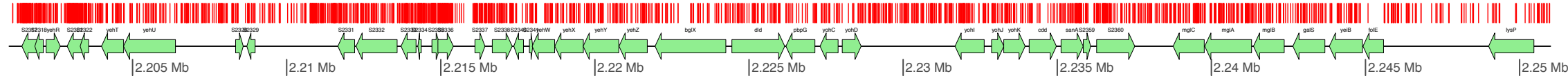
5 kb

2151 – 2201 Kb



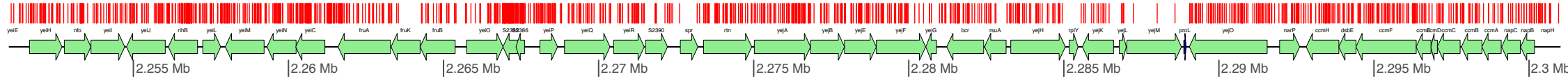
5 kb

2201 – 2251 Kb



5 kb

2251 – 2301 Kb



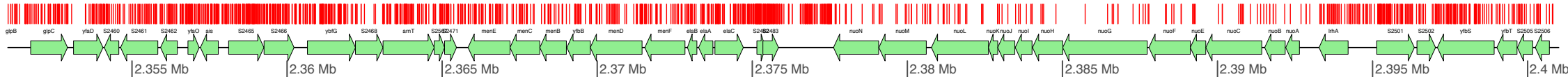
5 kb

2301 – 2351 Kb



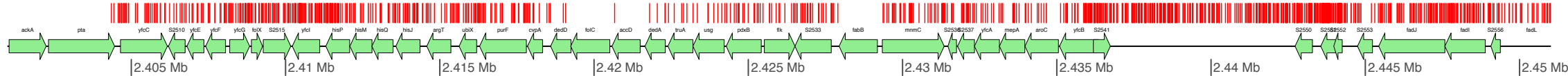
5 kb

2351 – 2401 Kb



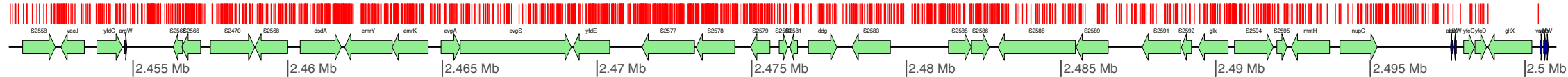
5 kb

2401 – 2451 Kb



5 kb

2451 – 2501 Kb



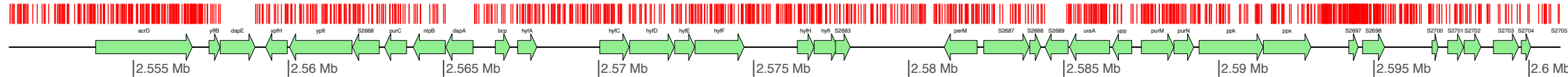
5 kb

2501 – 2551 Kb



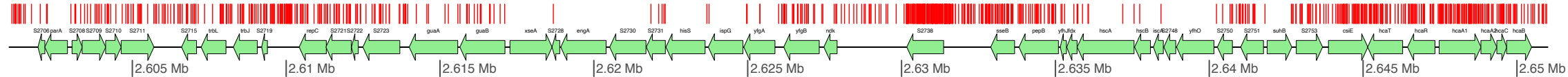
5 kb

2551 – 2601 Kb



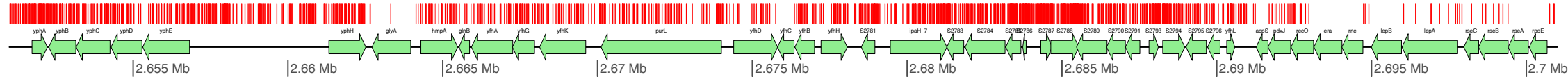
5 kb

2601 – 2651 Kb



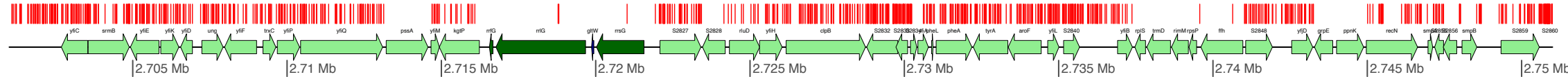
5 kb

2651 – 2701 Kb



5 kb

2701 – 2751 Kb



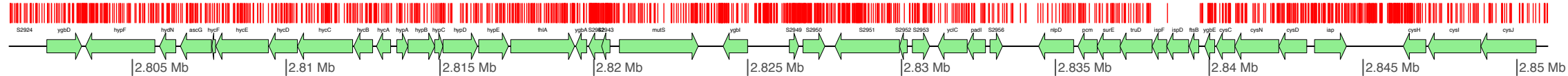
5 kb

2751 – 2801 Kb



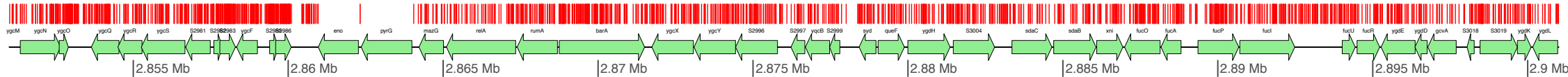
5 kb

2801 – 2851 Kb



5 kb

2851 – 2901 Kb



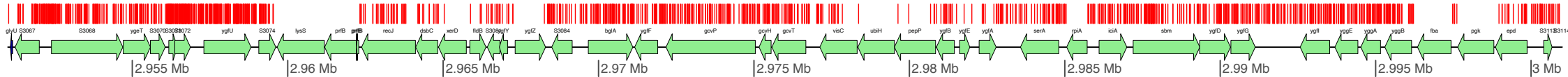
5 kb

2901 – 2951 Kb



5 kb

2951 – 3001 Kb



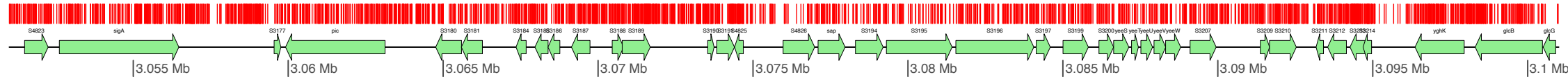
5 kb

3001 – 3051 Kb



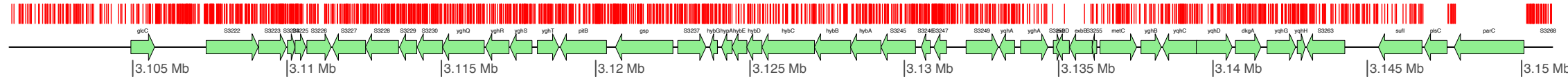
5 kb

3051 – 3101 Kb



5 kb

3101 – 3151 Kb



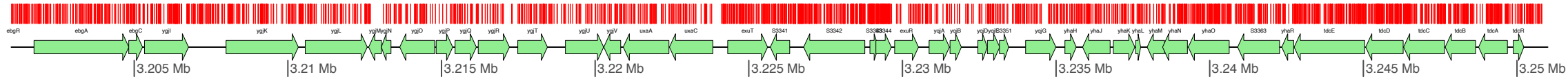
5 kb

3151 – 3201 Kb



5 kb

3201 – 3251 Kb



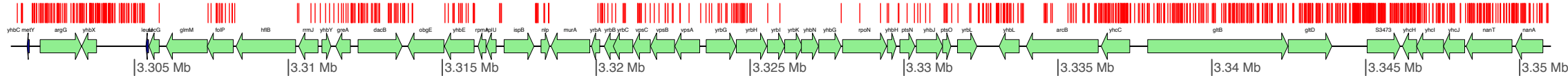
5 kb

3251 – 3301 Kb



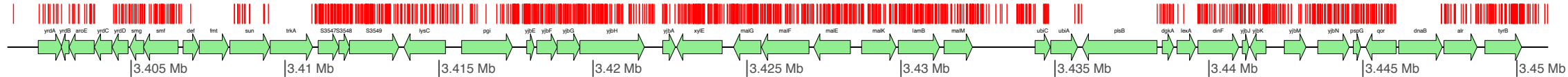
5 kb

3301 – 3351 Kb



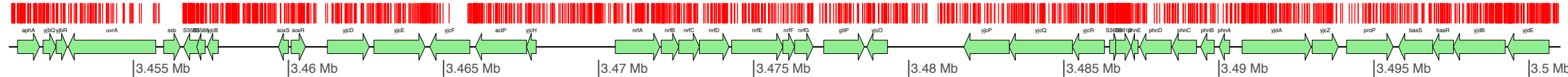
5 kb

3401 – 3451 Kb



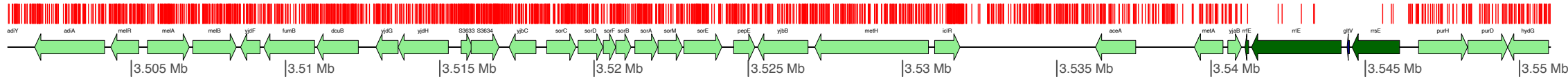
5 kb

3451 – 3501 Kb



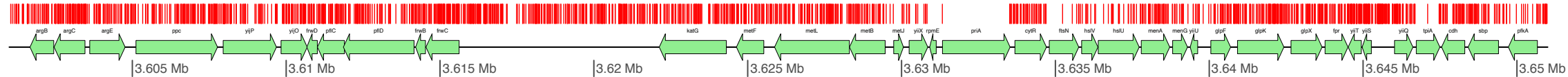
5 kb

3501 – 3551 Kb



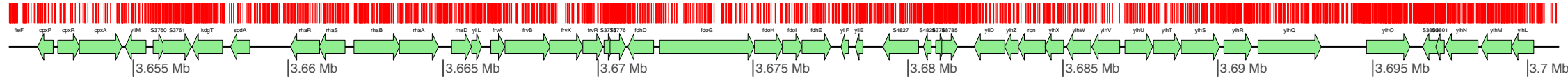
5 kb

3601 – 3651 Kb



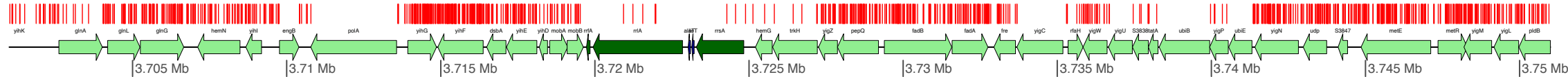
5 kb

3651 – 3701 Kb



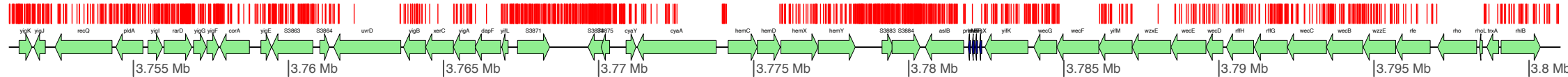
5 kb

3701 – 3751 Kb



5 kb

3751 – 3801 Kb



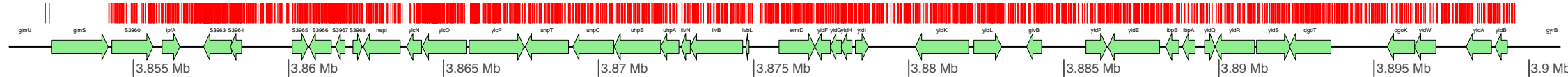
5 kb

3801 – 3851 Kb



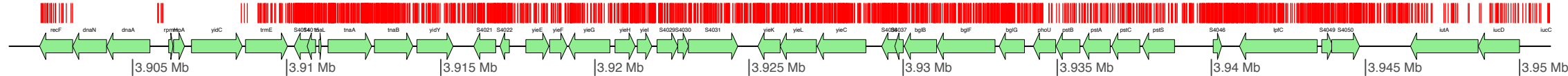
5 kb

3851 – 3901 Kb



5 kb

3901 – 3951 Kb



5 kb

3951 – 4001 Kb



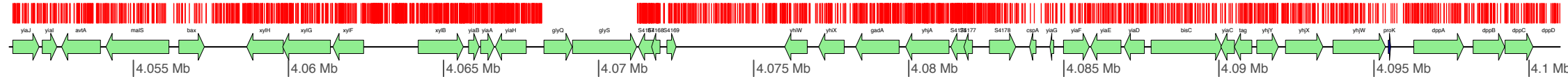
5 kb

4001 – 4051 Kb



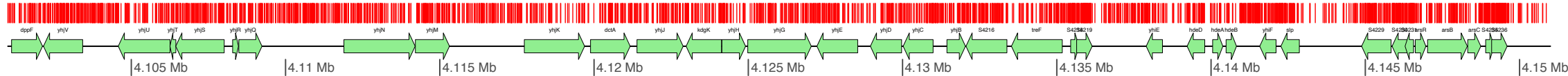
5 kb

4051 – 4101 Kb



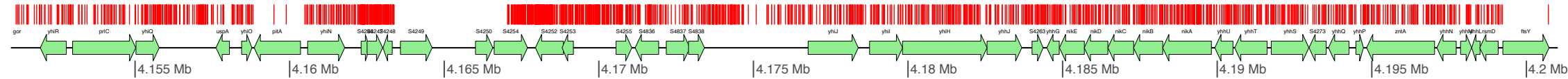
5 kb

4101 – 4151 Kb



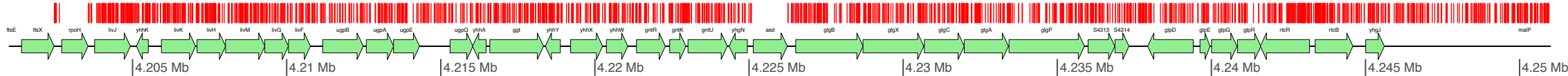
5 kb

4151 – 4201 Kb



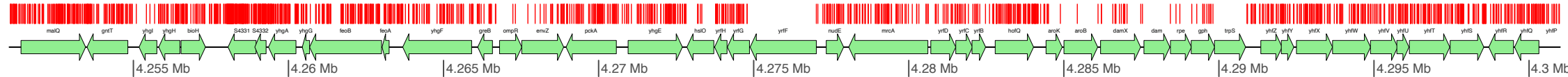
5 kb

4201 – 4251 Kb



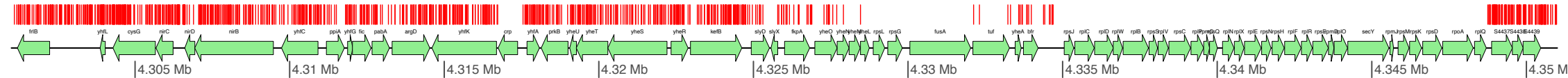
5 kb

4251 – 4301 Kb



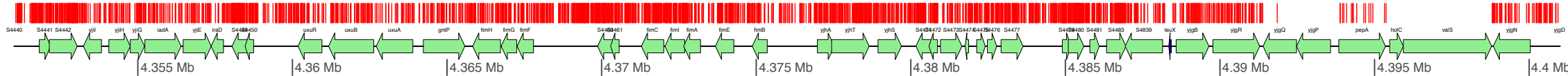
5 kb

4301 – 4351 Kb



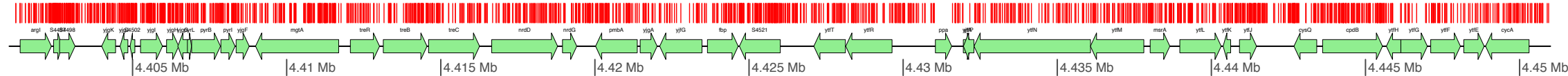
5 kb

4351 – 4401 Kb



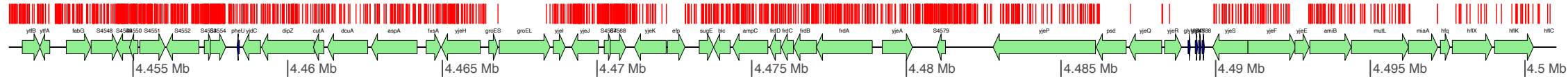
5 kb

4401 – 4451 Kb



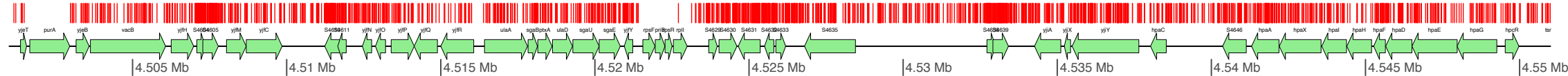
5 kb

4451 – 4501 Kb



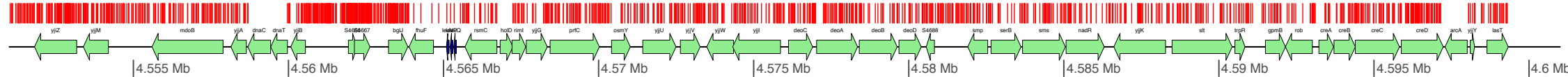
5 kb

4501 – 4551 Kb



5 kb

4551 – 4601 Kb



5 kb