

1 **Title page**

2 **Prophages integrating into prophages: a mechanism to accumulate type III secretion effector**
3 **genes and duplicate Shiga toxin-encoding prophages in *Escherichia coli***

4

5 **Short title: Prophages-in-prophage in STEC and EPEC**

6

7 Keiji Nakamura^a, Yoshitoshi Ogura^b, Yasuhiro Gotoh^a, Tetsuya Hayashi^{a*}

8 ^a Department of Bacteriology, Graduate School of Medical Sciences, Kyushu University, Fukuoka,
9 Japan, 812-8582

10 ^b Division of Microbiology, Department of Infectious Medicine, Kurume University School of
11 Medicine, Fukuoka, Japan, 830-0011

12

13 * Corresponding author

14 E-mail: thayash@bact.med.kyushu-u.ac.jp

16 **Abstract**

17 Bacteriophages (or phages) play major roles in the evolution of bacterial pathogens via horizontal
18 gene transfer. Multiple phages are often integrated in a host chromosome as prophages, not only
19 carrying various novel virulence-related genetic determinants into host bacteria but also providing
20 various possibilities for prophage-prophage interactions in bacterial cells. In particular, *Escherichia*
21 *coli* strains such as Shiga toxin (Stx)-producing *E. coli* (STEC) and enteropathogenic *E. coli* (EPEC)
22 strains have acquired more than 10 PPs (up to 21 PPs), many of which encode type III secretion
23 system (T3SS) effector gene clusters. In these strains, some prophages are present at a single locus in
24 tandem, which is usually interpreted as the integration of phages that use the same attachment (*att*)
25 sequence. Here, we present prophages integrating into T3SS effector gene cluster-associated loci in
26 prophages, which are widely distributed in STEC and EPEC. Some of the prophages integrated into
27 prophages are Stx-encoding prophages and have induced the duplication of Stx-encoding phages in
28 a single cell. The identified *att* sequences in prophage genomes are apparently derived from host
29 chromosomes. In addition, two or three different *att* sequences are present in some prophages, which
30 results in the generation of prophage clusters in various complex configurations. These “prophages-
31 in-prophages” represent a medically and biologically important type of inter-phage interaction that
32 promotes the accumulation of T3SS effector genes in STEC and EPEC, the duplication of Stx-
33 encoding prophages in STEC, and the conversion of EPEC to STEC and that may be distributed in
34 other types of *E. coli* strains as well as other prophage-rich bacterial species.

35

36 **Author summary**

37 Multiple prophages are often integrated in a bacterial host chromosome and some are present at a
38 single locus in tandem. The most striking examples are Shiga toxin (Stx)-producing and
39 enteropathogenic *Escherichia coli* (STEC and EPEC) strains, which usually contain more than 10
40 prophages (up to 21). Many of them encode a cluster of type III secretion system (T3SS) effector
41 genes, contributing the acquisition of a large number of effectors (>30) by STEC and EPEC. Here,
42 we describe prophages integrating into T3SS effector gene cluster-associated loci in prophages,

43 which are widely distributed in STEC and EPEC. Two or three different attachment sequences
44 derived from host chromosomes are present in some prophages, generating prophage clusters in
45 various complex configurations. Of note, some of such prophages-in-prophages are Stx-encoding
46 prophages and have induced the duplication of Stx-encoding prophages. Thus, these “prophages-in-
47 prophages” represent an important inter-phage interaction as they can promote not only the
48 accumulation of T3SS effectors in STEC and EPEC but also the duplication of Stx-encoding
49 prophages and the conversion of EPEC to STEC.

50

51 **Introduction**

52 Horizontal gene transfer (HGT) is an important mechanism for generating genetic and
53 phenotypic variations in bacteria [1-3]. Phages are major players in HGT, and many temperate phages
54 that confer virulence potential to host bacteria through the transfer of virulence-related genes have
55 been identified [4]. Most temperate phages integrate their genomes into host chromosomes by site-
56 specific recombination to become a part of the chromosomes as prophages (PPs) and enter a lysogenic
57 cycle. Recombination takes place between the homologous sequences of phage and host DNA (*attP*
58 and *attB*, respectively) and is mediated by a phage-encoded integrase [5]. Many bacterial
59 species/strains contain multiple PPs [6-8], providing various possibilities for PP-PP interactions [9,
60 10]. In particular, *Escherichia coli* strains such as Shiga toxin (Stx)-producing *E. coli* (STEC) strains
61 have acquired more than 10 PPs (up to 21 PPs) [11-14], and some of the PPs are located at the same
62 loci in tandem.

63 STEC strains cause diarrhea and severe illnesses, such as hemorrhagic colitis (HC) and life-
64 threatening hemolytic -uremic syndrome (HUS). Their key virulence factor is Stx. While there are
65 two subtypes (Stx1 and Stx2) with several variants and STEC produces one or more Stx
66 subtypes/variants [15-18], the known *stx* genes are all encoded by PP genomes. In addition, typical
67 STEC strains share the locus of enterocyte effacement (LEE) locus-encoding T3SS with
68 enteropathogenic *E. coli* (EPEC), and more than 30 effectors have been carried into STEC and EPEC
69 by multiple PPs [19-22]. Thus, EPEC strains are generally regarded as progenitors of typical STEC

70 strains. For example, O157:H7 STEC evolved from an ancestral EPEC O55:H7 through the phage-
71 mediated acquisition of *stx* along with a serotype change [23, 24].

72 In this study, we initially analyzed the duplicated Stx2-encoding PPs (referred to as Stx-
73 PPs) in STEC O145:H28, one of the major types of non-O157 STEC [25, 26], and found that one of
74 them is integrated into another PP. We then identified its *att* sequence. By subsequent analyses of
75 PPs carrying similar *att* sequences, we show that PP integration in PP (referred to as PP-in-PP) is a
76 genetic event widely occurring in STEC and EPEC and represents a mechanism underlying the
77 evolution and diversification of these bacteria.

78

79 **Results**

80 **Integration of inducible and packageable Stx2a phages into a PP integrated into the *ompW* locus** 81 **in STEC O145:H28**

82 We previously identified 18 PPs in the finished genome of O145:H28 strain 112648 [27],
83 including two Stx2a-PPs found at the *ompW* (P09) and *yecE* loci (P12). The two Stx2a-PP genomes
84 were identical in sequence; thus, they were considered duplicated PPs. As a lambda-like PP (P08)
85 was also found at *ompW*, we initially thought that P08 and P09 had been integrated in tandem.
86 However, by analyzing the potential *att* sequences of the three PPs, we found that while P08 and P12
87 were integrated into the *ompW* and *yecE* genes with *attL/R* sequences of 121 bp and 21 bp,
88 respectively, P09 was integrated into the P08 genome with a 21-bp *attL/R* sequence similar to that of
89 P12 (Fig 1a). By analyzing the PPs at the *ompW* and *yecE* loci in O145:H28 strains, we identified
90 another strain (12E129) that carries the same set of PPs: a lambda-like PP at *ompW*, an Stx2a-PP in
91 the PP at *ompW*, and another Stx2a-PP at *yecE* (Fig 1b). The potential *att* sequences of the three PPs
92 were identical to those of the corresponding PPs in strain 112648 (S1 Fig). The genomes of the two
93 Stx2a PPs in strain 12E129 were also nearly identical, excepting the left end. Hereafter, PPs integrated
94 into the same locus are collectively referred to as PPxxx (where xxx denotes the integration locus),
95 such as PP $ompW$.

96 To precisely determine the *att* sequences of each PP, we amplified and sequenced the *attP*-
97 flanking regions of excised and circularized genomes of these PPs. Although the two Stx2a-PPs in
98 strain 112648 were indistinguishable, those of strain 12E129 were distinguishable, allowing sequence
99 determination of the *attP*-flanking regions of three PPs from mitomycin C (MMC)-treated cell lysates.
100 This analysis confirmed that the predicted *att* sequences exactly represented those of the three PPs
101 and revealed that these PPs were induced to generate excised and circularized phage genomes by
102 MMC treatment (S1 Fig). The *att* sequences of P08 and P09/P12 in strain 112648 were also confirmed
103 using the same strategy. These results indicate that, in both strains, one of the duplicated Stx2a-PPs
104 has been integrated into PPomp*W*.

105 We further examined the packageability of these PP genomes into phage particles by PCR
106 analysis of DNase-treated culture supernatants of strain 12E129 with or without MMC treatment (Fig
107 1c). This analysis detected DNase-resistant genomic DNA of the two Stx2a-PPs, but did not that of
108 PPomp*W*, indicating that the duplicated Stx2a-PPs were both packaged into the phage particles. In a
109 similar analysis of strain 112648, the packaged genome of Stx2a-PP (P09 and/or P12) was detected.
110 That of P08 (PPomp*W*) was also not detected (data not shown), but the reason is currently unknown.

111

112 **Dynamics of PPomp*W*s, PPs integrated into PPomp*W*, and PPyec*Es* in STEC O145:H28**

113 To investigate the distribution of PPomp*W*s and the *att* sequences found in two PPomp*W*s
114 (referred to as *att*-in-PPomp*W*) among O145:H28 strains, we selected 64 genomes from 239 strains
115 analyzed in our previous study [27]. This set comprised 8 finished and 56 draft genomes and
116 encompassed seven of the eight clades previously identified in the major lineage (sequence type (ST)
117 32) and a minor lineage (ST137/6130) of O145:H28, thus largely representing the overall phylogeny
118 of O145:H28 as shown by a whole genome-based maximum likelihood (ML) tree (Fig 2).

119 PPomp*W*s were present in all 64 strains analyzed, including the two aforementioned strains.
120 All-to-all sequence comparison of the PPomp*W*s from eight finished genomes and 12 PPomp*W*s
121 sequenced in this study revealed that the PPomp*W* genomes were highly conserved, although
122 sequence diversification and segment replacement, probably by recombination, were detected in

123 some parts of several *PPompW*s (S2a Fig). Further analysis of the 20 *PPompW*s revealed that all
124 contained the 21-bp *att-in-PPompW* sequence (Fig 2), with one exception where the *att*-containing
125 region had been replaced by an insertion sequence (IS). These results indicate that a *PPompW*
126 containing *att-in-PPompW* was acquired by an ancestral strain and has been stably maintained in
127 O145:H28.

128 Examination of PP integration into the *att-in-PPompW* and *yecE* loci in the 64 strains
129 revealed that PPs are integrated into the two loci in 14 and 21 strains, respectively, with marked
130 variation in the PP content between strains (Fig 2). At the *att-in-PPompW* locus, Stx2a-PPs were
131 present in 10 strains and non-Stx-PPs in four strains (all belonging to ST32 clade H). More variable
132 PPs were found at *yecE*: Stx1a-PPs in 11 strains, Stx2a-PPs in eight strains, an Stx2d-PP in one strain,
133 and non-Stx-PPs in two strains. Two aforementioned strains (112648 and 12E129) carrying two
134 duplicated Stx2a-PPs belonged to different ST32 clades, indicating that duplication occurred
135 independently.

136 All-to-all sequence comparison of 27 PP genomes integrated into the *att-in-PPompW* (n=8;
137 all were Stx2a-PPs) or *yecE* (n=19; 9 Stx1-PPs, 8 Stx2a-PPs, one Stx2d-PP, and one non-Stx-PP)
138 locus revealed that the Stx1a-PP genomes were relatively well conserved, while regions with 2-3%
139 sequence divergence were probably introduced by recombination (S2b Fig). In contrast, the Stx2a-
140 PP genomes were highly variable except for those in the ST32 clade A/B/C strains. Interestingly,
141 although the Stx2a-PP of strain RM9872 (clade C) was integrated into *yecE*, this PP was similar to
142 the Stx2a-PP at *att-in-PPompW* in clade A/B strains. Considering the high conservation of Stx1a-PPs
143 at the *yecE* locus in these clades, it is likely that the Stx1a-PP at *yecE* has been replaced by the Stx2a-
144 PP originally integrated into *att-in-PPompW* in strain RM9872.

145

146 **Wide distribution of *PPompW*s and *att-in-PPompW* in *E. coli***

147 We next examined the distribution of *PPompW*s and the *att-in-PPompW* sequence (or
148 sequences similar to it) in the entire *E. coli* lineage by searching for them in 767 publicly available
149 complete *E. coli* genomes. *PPompW* was found in 44% of the *E. coli* strains examined (338 strains of

150 92 serotypes; all but O145:H28 and O26:H11 comprised a single ST). Phylogenetic analysis of the
151 *E. coli* strains representing each of the 92 serotypes showed that PP*ompW*s are widely distributed in
152 *E. coli* (Fig 3a). In contrast, after filtering the *att* sequence in *yecE*, 21-bp sequences identical to the
153 *att*-in-PP*ompW* sequence or with a 1-base mismatch (hereafter, collectively referred to as 21-bp
154 sequences) were detected in 150 strains of 20 serotypes belonging to five different *E. coli* phylogroups
155 (Fig 3a and S1 Table). In 145 of the 150 strains, the 21-bp sequence was present in PP*ompW*s. In 28
156 of the 145 strains (all were serotype O157:H7), two 21-bp sequences were found in two PPs located
157 in tandem at *ompW* (4 strains) or in a PP*ompW* and a PP cluster present at the *mlrA* or *ydfJ* locus (1
158 and 23 strains, respectively). One atypical O157:H7 strain (PV15-279) carried a PP*ompW*, but the
159 21-bp sequence was found in the PP cluster at *ydfJ*. The remaining four strains (all were serotype
160 O177:H25) contained no PP*ompW*s, and their 21-bp sequences were found in a PP or a PP cluster at
161 *ydfJ*.

162 By examining the 145 strains containing the 21-bp *att*-in-PP*ompW* sequence, we identified
163 additional strains carrying PPs integrated in PP*ompW*s in non-O145:H28 lineages: one O157:H7
164 strain and two O145:H25 strains (S2 Table). Moreover, as in the two aforementioned O145:H28
165 strains, the duplication of Stx2-PP and integration of the copies into the *att*-in-PP*ompW* and *yecE* loci
166 occurred in two of the three strains (Stx2d-PP in O157:H7 strain 28RC1 and Stx2a-PP in O145:H25
167 strain CFSAN004176; S2 Table and S3 Fig), although one of the duplicated Stx2a-PPs in the
168 O145:H25 strain contained a large genomic deletion and its *stx2A* gene was inactivated by multiple
169 insertions and deletions in the coding sequence [14].

170

171 **Close association of the *att*-in-PP*ompW* sequence with the PP regions encoding T3SS effector** 172 **genes**

173 Comparison of the PP*ompW* genomes containing the *att*-in-PP*ompW* sequence (S4 Fig)
174 revealed that while the early regions were relatively well conserved, the late regions were highly
175 variable between PPs due to sequence diversification, deletions and IS insertions. In particular, the
176 PP*ompW* genomes of phylogroup A strains have been highly degraded by deletions. However,

177 multiple T3SS effector genes are present just upstream of the *att*-in-PP*ompW* sequence in all
178 PP*ompWs* except for that in an O182:H25 strain, from which effector genes have apparently been
179 deleted (S4 Fig). Thus, the *att*-in-PP*ompW* sequence is closely linked to the T3SS effector-encoding
180 locus located at the very end of PP*ompW* genomes. Such regions of lambda-like phages encoding
181 various T3SS effector genes are called exchangeable effector loci or EELs [20]. The PP*ompWs*
182 containing the *att*-in-PP*ompW* sequence were also apparently lambda-like phages.

183 By analyzing T3SS effector genes in the EELs in the 19 PP*ompW* genomes, we identified
184 seven effector genes belonging to the *nleA*, *nleH*, *nleF*, and *espM* families and three *nleG* subfamilies
185 (G1-3) (Fig 3). Although there were variations in the effector gene repertoire between PP*ompWs* and
186 gene inactivation due to various types of mutations (mostly deletions) was detected in several
187 PP*ompWs*, a similar set of effector genes was found at the PP*ompW* EELs. As one or more IS elements
188 were present at all EELs, the variation in effector gene repertoire was probably generated by IS
189 insertion-associated events. The conservation patterns of effector genes among the 19 PP*ompW*
190 genomes suggest that the EELs of O157:H7 strain Sakai (phylogroup E; E731 in Fig 3) and EPEC
191 O76:H7 strain FORC_042 (phylogroup B1, E398 in Fig 3) represent the ancestral structure encoding
192 seven effector genes.

193 It should be noted that all strains carrying a PP(s) that contained the 21-bp *att* sequence and
194 the associated EEL(s) possessed the *eae* gene, a marker gene of the LEE (S1 Table), indicating that
195 they are all EPEC or typical (LEE-positive) STEC.

196

197 **PP clusters that contained PPs carrying the *att*-in-PP*ompW* sequence and identification of** 198 **additional *att* sites in PP genomes.**

199 In four of the aforementioned 28 O157:H7 strains that contained two 21-bp sequences
200 identical or nearly identical to *att*-in-PP*ompW*, the sequences were each present in the EEL-associated
201 region of two PP*ompW* genomes integrated in tandem (Fig 4a and S5 Fig). In these strains (as
202 represented by FRIK2069 in Fig 4a), while one of the EELs encoded an effector gene set similar to

203 that of other PP*ompW* EELs, the other encoded an *nleG* variant different from the three *nleG* families
204 at other PP*ompW* EELs.

205 In 24 O157:H7 strains, one 21-bp sequence was present in PP*ompW*, and the other was
206 present in PP-in-PP clusters comprising two to four PPs. In one strain (FRIK944; [Fig 4b](#)), the PP
207 cluster was present at *mlrA* (synonyms: *yehV*) and comprised two PPs, an Stx1-PP and a lambda-like
208 PP. By analyzing the *attL/R* sites of each PP, we found that while Stx1-PP is integrated into *mlrA* [9],
209 the lambda-like PP is in Stx1-PP, using the 96-bp *att* sequence (referred to as *att-in-PP_2*; see [S6 Fig](#)
210 for the sequence) associated with an EEL similar to PP*ompW* EELs. The lambda-like PP also
211 contained an *nleG* variant, but the 21-bp sequence was present between *attL* and the integrase gene
212 and was not associated with the *nleG* variant. As it is now known that the 21-bp sequence is present
213 in a PP genome other than PP*ompW* genomes, we hereafter refer to it as *att-in-PP_1*. Intriguingly,
214 between *att-in-PP_1* and the integrase gene of the lambda-like PP, the 121-bp *att* sequence for
215 PP*ompW*s was present. Although PP integration into the 121-bp sequence in PP genomes has yet to
216 be identified, this sequence can serve as a potential *att* site in PP genomes. We therefore refer to it as
217 *att-in-PP_3*.

218 In the remaining 23 strains, PP clusters comprising two to four PPs were present at *ydfJ*
219 ([Fig 4c](#); see [S5 Fig](#) for other strains). In these strains, one or two lambda-like PPs, which carry EELs
220 similar to the PP*ompW* EELs or encode multiple *nleG* variants, were integrated into *ydfJ* (see [S7 Fig](#)
221 for the *att* sequences). The former type of EEL was associated with *att-in-PP_2*, into which another
222 lambda-like PP was integrated. Similar to the PP integrated into PP*mlrA* ([Fig 4b](#)), the PPs integrated
223 into PP*ydfJ* contained the *att-in-PP_1* and *att-in-PP_3* sequences downstream of the integrase gene
224 and encoded *nleG* variants at the opposite PP end. Moreover, in one of the 23 strains (PV15-279, an
225 atypical O157:H7 strain [28]), an Stx2a-PP was integrated into the *att-in-PP_1* of the PP integrated
226 into PP*ydfJ* ([Fig 4c](#)).

227 Among the four aforementioned O177:H25 strains that contained the 21-bp *att-in-PP_1*
228 sequence, a similar but slightly different pattern of PP integration into PP genomes was observed ([Fig](#)
229 [4d](#)). In these strains, the *att-in-PP_1* sequence was found in a PP-like region that probably represents

230 two highly degraded PP genomes integrated in tandem between the *rspB* and *trg* genes. EELs similar
231 to the PP*ompW* EELs, *att-in-PP_2* and *att-in-PP_1* were found in this order, and a lambda-like PP
232 was integrated into *att-in-PP_2*. Moreover, the lambda-like PPs integrated into *att-in-PP_2* contained
233 the *att-in-PP_1* and *att-in-PP_3* sequences and multiple *nleG* variants, as PPs integrated into PP*mlrA*
234 or PP*ydfJs* (Fig 4d). This finding indicates that the distribution of these three *att* sequences in PP
235 genomes is not limited to O157:H7 strains.

236

237 **Origins of *att-in-PP* sequences**

238 Finally, to explore the origins of these *att-in-PP* sequences, we compared their flanking
239 sequences with *E. coli* chromosome sequences. The *att-in-PP_1*-flanking sequences in PP*ompWs* and
240 other PPs (all are integrated into PPs as shown in S5 Fig) were highly conserved, implying that these
241 sequences have a common origin (S8 Fig). Moreover, the 100-bp sequences including the *att-in-PP_1*
242 sequence showed a notable similarity (87% identity) to the corresponding *yecE* region (Fig 5, see S8
243 Fig for sequence alignment), suggesting that the *att-in-PP_1* and its flanking sequence originated
244 from the *yecE* locus.

245 Sequence similarity was also detected between the 96-bp *att-in-PP_2* sequence in
246 PP*ompWs* and the *ykgJ/ecpE* intergenic region of the *E. coli* chromosome (78% identity) (Fig 5). As
247 the homologous sequence extended to 125 bp in PP*mlrA*, we performed an additional search of *E.*
248 *coli* complete genomes and identified seven *att-in-PP_2*-containing PPs, although this search was
249 limited to six STEC genomes fully annotated for PPs (S9 Fig). The identified PPs included the Stx1a-
250 PP (Sp15) at *mlrA* of O157:H7 strain Sakai [11], the aforementioned duplicated Stx2a-PPs of
251 O145:H28 strain 112648, duplicated Stx2a-PPs of the atypical O157:H7 strain PV15-279 (one in
252 PP*ompW* and the other in *yecE*; carrying a T3SS effector gene), and two PPs in O26:H11 and
253 O111:H8 STEC strains [12] (at *ydfJ* and *ssrA*, respectively; the former carries a T3SS effector gene).
254 In these seven PPs, homologous sequences further extended to 309 bp with 84% identity (Fig 5).
255 Contrary to the observation for *att-in-PP_1* and its flanking sequences, there was notable diversity in
256 the *att-in-PP_2* sequence (20/96 polymorphic sites) between the PP*ompWs*, PP*mlrA* in FRIK944, and

257 the other seven PPs (S9 Fig). These findings indicate that the *att*-in-PP₂ and its flanking sequences
258 originated from the *ykgJ/ecpE* intergenic region on the chromosomes of *E. coli* or its close relatives,
259 but acquisition of the sequences by phages might have occurred multiple times.

260 The 121-bp *att*-in-PP₃ sequence was found in many of the PPs-in-PPs identified in this
261 study (Figs 4 and 5, and S5 Fig) and showed 81% identity to the *E. coli ompW*, suggesting that its
262 possible origin is also the chromosome of *E. coli* or its close relatives. Interestingly, PP*ompW*s and
263 many other PPs-in-PPs contained two or three *att*-in-PP sequences in the same order. The sequences
264 between the *att*-in-PP sequences (indicated by green in Fig 5) were also conserved (up to a 5-single
265 nucleotide polymorphism (SNP) difference); however, the location of the *att*-in-PP set in PP*ompW*s
266 was different from that of other PPs integrated in PPs. This finding suggests that the region
267 encompassing three (or two) *att*-in-PP sequences was once acquired by either type of phage and
268 spread to the other by recombination or some other mechanisms.

269

270 Discussion

271 As summarized in Fig 6, we identified various PP integration patterns in STEC and EPEC
272 strains, including PP integration into PPs. Most temperate phages are integrated into host genomes
273 by integrase-mediated recombination between *attP* and *attB*. Tandem PP integration can occur if the
274 two phages share the same *attB* site. In contrast to this traditional view of the mechanism for
275 generating tandem PPs, this study identified many PPs that contain *att* sequences, which allow
276 another PP to be integrated into their genomes, forming a PP-in-PP configuration. The combination
277 of the two integration mechanisms generates more complex PP clusters in host genomes (combination
278 of tandem PPs and PPs-in-PPs). Frequent colocalization of multiple *att*-in-PP sequences potentially
279 generates much more variation than detected in this study. These *att*-in-PP sequences originated from
280 the host chromosome, providing more opportunities for lysogenization to incoming phages and
281 allowing the duplication of PPs encoding medically or biologically important genes, such as *stx*.
282 There may be some previously unrecognized interaction(s) between integrating PPs and their “host”
283 PPs. Analyses of such interactions as well as the mechanisms of incorporating *att*-in-PP sequences

284 from host chromosomes are worthy of future studies to better understand the processes of PP-in-PP
285 formation.

286 Notably, most *att*-in-PP sequences identified are linked to EELs that encode multiple
287 effector genes for the LEE-encoded T3SS, and PPs integrated into *att*-in-PPs often carry effector
288 genes. Thus, the PP-in-PP system has promoted the accumulation of effector genes in EPEC and
289 STEC strains [21, 29] and can promote further accumulation of these genes, which may increase the
290 pathogenicity of these strains [22, 30]. Furthermore, a significant portion of the PPs integrated in *att*-
291 in-PP_1 (13/18) encoded *stx* genes, indicating that the *att*-in-PP_1 sequence has promoted the
292 acquisition of *stx* genes and thus the conversion of EPEC to typical STEC, even if the *yecE* locus, the
293 origin of *att*-in-PP_1 and one of the integration hot spots of Stx-PP [12, 31-33], has been occupied
294 by another PP.

295 In conclusion, the findings obtained here highlight that PP integration systems are much
296 more complicated than previously recognized and provide additional insights into the evolution of
297 EPEC and STEC and their pathogenicity. It is also possible to find similar PP integration patterns in
298 other types of *E. coli* and other PP-rich species if PP clusters are carefully investigated. Similar
299 integration systems could also be found for genetic elements utilizing integrase-mediated integration
300 mechanisms, such as integrative and conjugative elements (ICEs) [34].

301

302 **Material and Methods**

303 **Bacterial strains**

304 The 64 O145:H28 strains analyzed in this study are listed in [S3 Table](#). Of these, 59 were
305 from our laboratory stock, which were genome-sequenced in our previous study [27], and 5 were
306 completely genome-sequenced stains (the plasmid genome was not finished in strain 2015C-3125),
307 the genome sequences of which were downloaded from the NCBI database. To construct the
308 completely genome-sequenced *E. coli* strain set, a total of 875 complete genomes were downloaded
309 from the database (accessed on the 20th of July 2019). After excluding laboratory, commercial and
310 re-sequenced strains and substrains, the 767 strains listed in [S4 Table](#) were used for analysis.

311 Annotation was carried out using the DDBJ Fast Annotation and Submission Tool (DFAST) [35], if
312 necessary.

313 **Extraction of total cellular and phage DNA**

314 Bacterial cells were grown overnight to the stationary phase at 37°C in lysogeny broth (LB)
315 medium. For prophage induction, cells were grown to the late log phase (0.7-0.9 OD₆₀₀), and MMC
316 was added to the culture to a final concentration of 1 µg/ml. After a 3-hr incubation, aliquots of the
317 culture were isolated, and the cells were collected by centrifugation. Total cellular DNA was extracted
318 from the cells using the alkaline-boiling method and used for PCR analyses. Phage particles were
319 isolated from the culture supernatant after a 3-hr incubation with MMC. The culture was first treated
320 with chloroform, and bacterial cell debris was removed by centrifugation. The supernatant was
321 filtered through a 0.2-µm-pore-size filter (Millipore) and incubated with DNase I (final concentration:
322 400 U/ml, TaKaRa) and RNase A (50 µg/ml, Sigma) at 37°C for 1 hr. After inactivating DNase I by
323 incubation at 75°C for 10 min and adding EDTA (5 mM, Nacalai Tesque), the sample was treated
324 with proteinase K (100 µg/ml; Wako) and used as packaged phage DNA. Total cellular DNA and
325 packaged phage DNA from MMC-untreated cultures were prepared with the same protocol. The
326 primers used in these analyses are listed in [S5 Table](#).

327 **Analyses of PP integration and sequencing of PP genomes**

328 PP integration into the *ompW*, *att-in-PPompW* (later renamed *att-inPP_1*) and *yecE* loci in
329 56 O145:H28 draft genomes was first examined by a BLASTN search as outlined in [S10a Fig](#). The
330 integration of Stx PPs into *att-in-PPompW* and/or *yecE* was determined by long PCR amplification
331 using primers targeting the *stx* genes and sequences adjacent to these integration sites, as
332 schematically shown in [S10b Fig](#). The products of long PCR were used for sequence determination
333 of each PP. The primers used in this analysis are listed in [S6 Table](#).

334 Sequencing libraries were prepared for each product of long PCR (ranging from 15 to 33
335 kb) using the Nextera XT DNA Sample Preparation Kit (Illumina) and sequenced on the Illumina
336 MiSeq platform to generate paired-end (PE) reads (300 bp x 2). PP genomic sequences were obtained
337 by assembling and scaffolding Illumina PE reads using the Platanus_B assembler (v1.1.0)

338 (<http://platanus.bio.titech.ac.jp/platanus-b>) [36]; then, gaps were closed by Sanger sequencing PCR
339 products that spanned the gaps. Annotation of all PP genomes was carried out with DFAST, followed
340 by manual curation using IMC-GE software (In Silico Biology). All sequences have been deposited
341 in the DDBJ/EMBL/GenBank databases under the accession numbers listed in [S3 Table](#).
342 GenomeMatcher (v2.3) [37] was used for genome sequence comparison and to display the results.

343 **Searches for PPompWs and att-in-PPompW sequences in the complete *E. coli* genomes**

344 Serotypes and *eae* subtypes of the 767 complete *E. coli* strains were determined by
345 BLASTN as previously described [29]. Systematic ST determination was performed by a read
346 mapping-based strategy using the SRST2 program [38] with default parameters. Read sequences of
347 the complete genomes were simulated with the ART program (ART_Illumina, version 2.5.8) [39].
348 The genomes whose ST was not precisely defined (possible ST containing a novel allele, an uncertain
349 ST, and no STs in the present database) were reanalyzed using MLST 2.0 with “Escherichia coli #1”
350 schemes [40] (<https://cge.cbs.dtu.dk/services/MLST/>).

351 The presence of PPompW and the att-in-PPompW sequence was examined in the complete
352 genomes by a BLASTN-based search as follows. The presence of PPompW was determined using
353 two query sequences: one was the integrase gene of O145:H28 strain 112648 (EC112648_1574)
354 (thresholds: >90% identity and >90% coverage), and the other was the *ompW*-containing region on
355 the chromosome of *E. coli* K-12 (No. NC_000913; nucleotide positions 1,314,020-1,315,224; no PP
356 integration) to examine the absence of PP insertion into the *ompW* locus (threshold: >85% identity
357 and <60% coverage). When either the integrase gene or the *ompW* locus split by some insertion was
358 detected, we analyzed the gene organization of these regions to determine if PPompW was present.
359 The search for the 21-bp att-in-PPompW sequence (5'-GTCATGCAGTTAAAGTGGCGG-3') ([S1c](#)
360 [Fig](#)) was performed with the blastn-short task option (thresholds: >95% identity and 100% coverage).
361 The 21-bp sequences in the *yecE* gene, which were similar to the att-in-PPompW sequence, were
362 removed.

363 **SNP detection and phylogenetic analysis**

364 The SNP sites (3,277 sites) of the core genomic sequences of the 64 O145:H28 strains were
365 detected by MUMmer [41], followed by filtering recombinogenic SNPs by Gubbins [42], and used
366 for reconstruction of an ML tree in RAxML [43] with the GTR gamma substitution model as
367 previously described [27]. To reconstruct the phylogeny of the *E. coli* strains carrying PPompW, we
368 used 92 *E. coli* strains representing each of the 92 serotypes that contained PPompW-carrying strains.
369 Strains in which *att-in-PPompW* was detected were preferentially selected from the serotypes that
370 contained multiple strains. *Escherichia* cryptic clade I strain TW15838 (No. AEKA01000000) was
371 used as an outgroup. The core genes (n=2,642) of these strains, which were defined as the genes
372 present in 100% of strains, were identified by Roary [44], and their concatenated sequence alignments
373 were generated by the same software. Based on the alignment (109,927 SNP sites in total), an ML
374 tree was constructed using RAxML as described above. Phylogroups of the strains were determined
375 by ClermonTyping [45]. ML trees were displayed and annotated using iTOL [46] or FigTree (v1.4.3)
376 (<http://tree.bio.ed.ac.uk/software/figtree/>).

377

378 **Acknowledgements**

379 This research was supported by AMED under Grant Number 20fk0108065h0803 to T.H., and a
380 KAKENHI from the Japan Society for the Promotion of Science (18K07116) to K.N. We thank M.
381 Horiguchi, M. Kumagai, Y. Nagayoshi, and K. Ozaki for providing technical assistance. We also
382 thank the EHEC working group in Japan for providing O145:H28 strains.

383

384 **Author Contributions**

385 **Conceptualization:** Nakamura K, Hayashi T.

386 **Data curation:** Nakamura K

387 **Formal analysis:** Nakamura K, Ogura Y, Gotoh Y.

388 **Funding acquisition:** Nakamura K, Hayashi T.

389 **Investigation:** Nakamura K

390 **Methodology:** Nakamura K, Ogura Y, Gotoh Y.

391 **Project Administration:** Hayashi T.

392 **Resources:** Nakamura K, Ogura Y.

393 **Visualization:** Nakamura K

394 **Writing - Original Draft Preparation:** Nakamura K

395 **Writing - Review & Editing:** Ogura Y, Gotoh Y, Hayashi T.

396

397 **References**

- 398 1. Brüssow H, Canchaya C, Hardt WD. Phages and the evolution of bacterial pathogens: from
399 genomic rearrangements to lysogenic conversion. *Microbiol Mol Biol Rev.* 2004;68(3):560-
400 602. doi: 10.1128/MMBR.68.3.560-602.2004
- 401 2. Frost LS, Leplae R, Summers AO, Toussaint A. Mobile genetic elements: the agents of open
402 source evolution. *Nat Rev Microbiol.* 2005;3(9):722-732. doi: 10.1038/nrmicro1235
- 403 3. Partridge SR, Kwong SM, Firth N, Jensen SO. Mobile genetic elements associated with
404 antimicrobial resistance. *Clin Microbiol Rev.* 2018;31(4):e00088-17. doi:
405 10.1128/CMR.00088-17
- 406 4. Feiner R, Argov T, Rabinovich L, Sigal N, Borovok I, Herskovits AA. A new perspective on
407 lysogeny: prophages as active regulatory switches of bacteria. *Nat Rev Microbiol.*
408 2015;13(10):641-650. doi: 10.1038/nrmicro3527
- 409 5. Fogg PC, Colloms S, Rosser S, Stark M, Smith MC. New applications for phage integrases. *J*
410 *Mol Biol.* 2014;426(15):2703-2716. doi: 10.1016/j.jmb.2014.05.014
- 411 6. Canchaya C, Proux C, Fournous G, Bruttin A, Brüssow H. Prophage genomics. *Microbiol Mol*
412 *Biol Rev.* 2003;67(2):238-276. doi: 10.1128/MMBR.67.2.238-276.2003
- 413 7. Casjens S. Prophages and bacterial genomics: what have we learned so far? *Mol Microbiol.*
414 2003;49(2):277-300. doi: 10.1046/j.1365-2958.2003.03580.x
- 415 8. Bobay LM, Rocha EP, Touchon M. The adaptation of temperate bacteriophages to their host
416 genomes. *Mol Biol Evol.* 2012;30(4):737-751. doi: 10.1093/molbev/mss279
- 417 9. Asadulghani M, Ogura Y, Ooka T, Itoh T, Sawaguchi A, Iguchi A, et al. The defective

- 418 prophage pool of *Escherichia coli* O157: prophage-prophage interactions potentiate horizontal
419 transfer of virulence determinants. PLoS Pathog. 2009;5(5):e1000408. doi:
420 10.1371/journal.ppat.1000408
- 421 10. De Paepe M, Hutinet G, Son O, Amarir-Bouhram J, Schbath S, Petit MA. Temperate phages
422 acquire DNA from defective prophages by relaxed homologous recombination: the role of
423 Rad52-like recombinases. PLoS Genet. 2014;10(3):e1004181. doi:
424 10.1371/journal.pgen.1004181
- 425 11. Hayashi T, Makino K, Ohnishi M, Kurokawa K, Ishii K, Yokoyama K, et al. Complete genome
426 sequence of enterohemorrhagic *Escherichia coli* O157:H7 and genomic comparison with a
427 laboratory strain K-12. DNA Res. 2001;8(1):11-22. doi: 10.1093/dnares/8.1.11
- 428 12. Ogura Y, Ooka T, Iguchi A, Toh H, Asadulghani M, Oshima K, et al. Comparative genomics
429 reveal the mechanism of the parallel evolution of O157 and non-O157 enterohemorrhagic
430 *Escherichia coli*. Proc Natl Acad Sci USA. 2009;106(42):17939-17944. doi:
431 10.1073/pnas.0903585106
- 432 13. Kyle JL, Cummings CA, Parker CT, Quinones B, Vatta P, Newton E, et al. *Escherichia coli*
433 serotype O55:H7 diversity supports parallel acquisition of bacteriophage at Shiga toxin phage
434 insertion sites during evolution of the O157:H7 lineage. J Bacteriol. 2012;194(8):1885-1896.
435 doi: 10.1128/JB.00120-12
- 436 14. Lorenz SC, Gonzalez-Escalona N, Kotewicz ML, Fischer M, Kase, A, et al. Genome
437 sequencing and comparative genomics of enterohemorrhagic *Escherichia coli* O145:H25 and
438 O145:H28 reveal distinct evolutionary paths and marked variations in traits associated with
439 virulence & colonization. BMC Microbiol. 2017;17(1):183. doi: 10.1186/s12866-017-1094-3
- 440 15. Teel LD, Melton-Celsa AR, Schmitt CK, O'Brien AD. One of two copies of the gene for the
441 activatable Shiga toxin type 2d in *Escherichia coli* O91:H21 strain B2F1 is associated with an
442 inducible bacteriophage. Infect Immun. 2002;70(8):4282-4291. doi: 10.1128/IAI.70.8.4282-
443 4291.2002
- 444 16. Muniesa M, Blanco JE, de Simon M, Serra-Moreno R, Blanch AR, Jofre J. Diversity of *stx2*

- 445 converting bacteriophages induced from Shiga-toxin-producing *Escherichia coli* strains
446 isolated from cattle. *Microbiology*. 2004;150(9):2959-2971. doi: 10.1099/mic.0.27188-0
- 447 17. Dallman TJ, Ashton PM, Byrne L, Perry NT, Petrovska L, Ellis R, et al. Applying
448 phylogenomics to understand the emergence of Shiga-toxin-producing *Escherichia coli*
449 O157:H7 strains causing severe human disease in the UK. *Microb Genom*. 2015;1(3):e000029.
450 doi: 10.1099/mgen.0.000029
- 451 18. Ogura Y, Gotoh Y, Itoh T, Sato MP, Seto K, Yoshino S, et al. Population structure of
452 *Escherichia coli* O26:H11 with recent and repeated *stx2* acquisition in multiple lineages.
453 *Microb Genom*. 2017;3(11):e000141. doi: 10.1099/mgen.0.000141
- 454 19. Deng W, Puente JL, Gruenheid S, Li Y, Vallance BA, Vázquez A, et al. Dissecting virulence:
455 Systematic and functional analyses of a pathogenicity island. *Proc Natl Acad Sci USA*.
456 2004;101(10):3597-3602. doi: 10.1073/pnas.0400326101
- 457 20. Tobe T, Beatson SA, Taniguchi H, Abe H, Bailey CM, et al. An extensive repertoire of type
458 III secretion effectors in *Escherichia coli* O157 and the role of lambdoid phages in their
459 dissemination. *Proc Natl Acad Sci USA*. 2006;103(40):14941-14946. doi:
460 10.1073/pnas.0604891103
- 461 21. Ingle DJ, Tauschek M, Edwards DJ, Hocking DM, Pickard DJ, Azzopardi KI, et al. Evolution
462 of atypical enteropathogenic *E. coli* by repeated acquisition of LEE pathogenicity island
463 variants. *Nat Microbiol*. 2016;1:15010. doi: 10.1038/nmicrobiol.2015.10
- 464 22. Hazen TH, Sonnenberg MS, Panchalingam S, Antonio M, Hossain A, Mandomando I, et al.
465 Genomic diversity of EPEC associated with clinical presentations of differing severity. *Nat*
466 *Microbiol*. 2016;1:15014. doi: 10.1038/nmicrobiol.2015.14
- 467 23. Wick LM, Qi W, Lacher DW, Whittam TS. Evolution of genomic content in the stepwise
468 emergence of *Escherichia coli* O157:H7. *J. Bacteriol*. 2005;187(5):1783-1791. doi:
469 10.1128/JB.187.5.1783-1791.2005
- 470 24. Feng PCH, Monday SR, Lacher DW, Allison L, Siitonen A, Keys C, et al. Genetic diversity
471 among clonal lineages within *Escherichia coli* O157:H7 stepwise evolutionary model.

- 472 Emerging Infect Dis. 2007;13(11):1701-1706. doi: 10.3201/eid1311.070381
- 473 25. Karmali MA, Mascarenhas M, Shen S, Ziebell K, Johnson S, Reid-Smith R, et al. Association
474 of genomic O island 122 of *Escherichia coli* EDL 933 with verocytotoxin-producing
475 *Escherichia coli* seropathotypes that are linked to epidemic and/or serious disease. J Clin
476 Microbiol. 2003;41(11):4930-4940. doi: 10.1128/JCM.41.11.4930-4940.2003
- 477 26. European Food Safety Authority and European Centre for Disease Prevention and Control.
478 The European Union summary report on trends and sources of zoonoses, zoonotic agents and
479 food-borne outbreaks in 2017. EFSA Journal. 2018;16(12):e05500. doi:
480 10.2903/j.efsa.2018.5500
- 481 27. Nakamura K, Murase K, Sato MP, Toyoda A, Itoh T, Mainil JG, et al. Differential dynamics
482 and impacts of prophages and plasmids on the pangenome and virulence factor repertoires of
483 Shiga toxin-producing *Escherichia coli* O145:H28. Microb Genom. 2020;6(1):e000323. doi:
484 10.1099/mgen.0.000323
- 485 28. Ogura Y, Seto K, Morimoto Y, Nakamura K, Sato MP, Gotoh Y, et al. Genomic
486 characterization of β -glucuronidase-positive *Escherichia coli* O157:H7 producing Stx2a.
487 Emerging Infect Dis. 2018;24(12):2219-2227. doi: 10.3201/eid2412.180404
- 488 29. Arimizu Y, Kirino Y, Sato MP, Uno K, Sato T, Gotoh Y, et al. Large-scale genome analysis
489 of bovine commensal *Escherichia coli* reveals that bovine-adapted *E. coli* lineages are serving
490 as evolutionary sources of the emergence of human intestinal pathogenic strains. Genome Res.
491 2019;29(9):1495-1505. doi: 10.1101/gr.249268.119.
- 492 30. Dean P, Kenny B. The effector repertoire of enteropathogenic *E. coli*: ganging up on the host
493 cell. Curr Opin Microbiol. 2009;12(1):101-109. doi: 10.1016/j.mib.2008.11.006
- 494 31. Bonanno L, Loukiadis E, Mariani-Kurkdjian P, Oswald E, Garnier L, Michel V, et al. Diversity
495 of Shiga toxin-producing *Escherichia coli* (STEC) O26:H11 strains examined via *stx* subtypes
496 and insertion sites of Stx and EspK bacteriophages. Appl Environ Microbiol.
497 2015;81(11):3712-3721. doi: 10.1128/AEM.00077-15
- 498 32. Cointe A, Birgy A, Mariani-Kurkdjian P, Liguori S, Courroux C, Blanco J, et al. Emerging

- 499 multidrug-resistant hybrid pathotype Shiga toxin-producing *Escherichia coli* O80 and related
500 strains of clonal complex 165, Europe. *Emerging Infect Dis.* 2018;24(12):2262-2269. doi:
501 10.3201/eid2412.180272
- 502 33. Yara DA, Greig DR, Gally DL, Dallman TJ, Jenkins C. Comparison of Shiga toxin-encoding
503 bacteriophages in highly pathogenic strains of Shiga toxin-producing *Escherichia coli*
504 O157:H7 in the UK. *Microb Genom.* 2020;6(3):e000334. doi: 10.1099/mgen.0.000334
- 505 34. Bellanger X, Payot S, Leblond-Bourget N, Guédon G. Conjugative and mobilizable genomic
506 islands in bacteria: evolution and diversity. *FEMS Microbiol Rev.* 2014;38(4):720-760. doi:
507 10.1111/1574-6976.12058
- 508 35. Tanizawa Y, Fujisawa T, Nakamura Y. DFAST: a flexible prokaryotic genome annotation
509 pipeline for faster genome publication. *Bioinformatics.* 2018;34(6):1037-1039. doi:
510 10.1093/bioinformatics/btx713
- 511 36. Kajitani R, Yoshimura D, Ogura Y, Gotoh Y, Hayashi T, Itoh T. Platanus_B: an accurate *de*
512 *novo* assembler for bacterial genomes using an iterative error-removal process. *DNA Res.*
513 2020;27(3):dsaa014. doi: 10.1093/dnares/dsaa014
- 514 37. Ohtsubo Y, Ikeda-Ohtsubo W, Nagata Y, Tsuda M. GenomeMatcher: A graphical user
515 interface for DNA sequence comparison. *BMC Bioinformatics* 2008;9(1):1-9. doi:
516 10.1186/1471-2105-9-376
- 517 38. Inouye M, Dashnow H, Raven LA, Schultz MB, Pope BJ, Tomita T, et al. SRST2: Rapid
518 genomic surveillance for public health and hospital microbiology labs. *Genome Med.*
519 2014;6(11):90. doi: 10.1186/s13073-014-0090-6
- 520 39. Huang W, Li L, Myers JR, Marth GT. ART: a next-generation sequencing read simulator.
521 *Bioinformatics.* 2012;28(4):593-594. doi: 10.1093/bioinformatics/btr708
- 522 40. Wirth T, Falush D, Lan R, Colles F, Mensa P, Wieler LH, et al. Sex and virulence in
523 *Escherichia coli*: an evolutionary perspective. *Mol Microbiol.* 2006;60(5):1136-1151. doi:
524 10.1111/j.1365-2958.2006.05172.x
- 525 41. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and

- 526 open software for comparing large genomes. *Genome Biol.* 2004;**5(2)**:R12. doi: 10.1186/gb-
527 2004-5-2-r12
- 528 42. Croucher NJ, Page AJ, Connor TR, Delaney AJ, Keane JA, Bentley SD, et al. Rapid
529 phylogenetic analysis of large samples of recombinant bacterial whole genome sequences
530 using Gubbins. *Nucleic Acids Res.* 2015;43(3):e15. doi: 10.1093/nar/gku1196
- 531 43. Stamatakis A. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with
532 thousands of taxa and mixed models. *Bioinformatics* 2006;22(21):2688-2690. doi:
533 10.1093/bioinformatics/btl446
- 534 44. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MT, et al. Roary: rapid large-
535 scale prokaryote pan genome analysis. *Bioinformatics* 2015;31(22):3691-3693. doi:
536 10.1093/bioinformatics/btv421
- 537 45. Beghain J, Bridier-Nahmias A, Le Nagard H, Denamur E, Clermont O. ClermonTyping: an
538 easy-to-use and accurate *in silico* method for *Escherichia* genus strain phylotyping. *Microb*
539 *Genom.* 2018;4(7):e000192. doi: 10.1099/mgen.0.000192
- 540 46. Letunic I, Bork P. Interactive tree of life (iTOL) v3: an online tool for the display and
541 annotation of phylogenetic and other trees. *Nucleic Acids Res.* 2016;44(W1):W242–W245.
542 doi: 10.1093/nar/gkw290

543

544 **Supporting information**

545 **S1 Fig. Determination of the *attP* sequences of PPs integrated in *ompW*, *PPompW*, and *yecE* in**
546 **O145:H28 strains 12E129 and 112648.** (a) Schematic representation of the PCR strategy used to
547 amplify the *attP*-flanking region (left panel) and the locations of PCR primers used for each PP (right
548 panel). (b) PCR detection of excised and circularized PP genomes. Total cellular DNA isolated from
549 MMC-treated (+) or MMC-untreated (-) cells was analyzed. A chromosome backbone (CB) region
550 was amplified as a positive control. (c) The *att* sequences of the three PPs in strain 12E129. The *attP*-
551 containing sequences obtained by sequencing the PCR products shown in S1b Fig were aligned with
552 the *attR*-, *attL*-, *attB*-containing sequences to define the *att* sequences of each PP. The *ompW*

553 sequences of strain K-12 MG1655 (accession No. NC_000913) and the *yecE* sequence of O145:H28
554 strain 122715 (accession No. AP019708), in which no PPs were integrated, were used as the *attB*
555 sequences, respectively (indicated by a dagger and a section mark, respectively). The defined *att*
556 sequences are indicated by uppercase letters.

557 **S2 Fig. All-to-all genome sequence comparison of PPompWs and that of PPatt-in-PPompWs and**
558 **PPyecEs found in 64 O145:H28 strains.** Dot plot matrixes of the concatenated sequences of the 20
559 PPompW genomes (a) and 27 PPatt-in-PPompW and PPyecE genomes (b) found in 64 O145:H28
560 strains are shown. Strain names and information on the ST and ST32 clade of each strain are indicated.
561 Sequence identities are indicated by a heatmap. In panel A, the nucleotide sequences between *att*-in-
562 PPompW and *attR* (approximately 428 bp in length) were excluded from this analysis because the
563 sequences of this region in three strains (499, EH1910, and KIH15-140) were not determined. In
564 panel B, the subtype of Stx encoded by each PP and the integration site of each PP are indicated. PP
565 groups sharing similar genomic sequences are framed by boxes.

566 **S3 Fig. Variation in the sequences of the Stx2-encoding PPatt-in-PPompW and PPyecE genomes**
567 **found in *E. coli* strains analyzed in this study.** Dot plot matrixes of the concatenated sequences of
568 the PPatt-in-PPompW and PPyecE genomes encoding Stx2 are shown. The genome sequences of a
569 pair of PPatt-in-PPompW and PPyecE found in two O145:H28 strains, that in an O157:H7 strain and
570 that in an O145:H25 strain were compared. The names of host strains, Stx2 subtypes, and integration
571 site of each PP are indicated. PPs in the same strain are framed by boxes. Sequence identities are
572 indicated by a heatmap.

573 **S4 Fig. Comparison of the PPompW genomes containing the *att*-in-PPompW sequence.** In the
574 left panel, along with the same ML tree as shown in Fig 3, *E. coli* strain ID, phylogroup (PG), and
575 the presence (colored) or absence (open) of the *att*-in-PPompW sequence in each *E. coli* are indicated.
576 In the right panel, the genome structures of PPompWs containing *att*-in-PPompW are drawn to scale.
577 A large chromosome inversion resulting from recombination between PPompW and another PP in
578 two strains is indicated by an asterisk (E473 and E471), and relevant PP regions are shown.
579 Homologous regions and sequence identities are depicted by shading with a color gradient. The

580 Stx2a-PPs integrated into the *att*-in-PP*ompW* locus in strains E474 and E118 are schematically
581 indicated.

582 **S5 Fig. Variation in the PP integration patterns in the PP clusters that contained PPs carrying**
583 **potential *att* sites.** In the left table, a list of 33 strains that possessed PP clusters that contained PPs
584 carrying the 21-bp sequence identical or nearly identical to the *att*-in-PP*ompW* sequence is provided.
585 In the right panel, the patterns of PP integration are schematically illustrated. Strains showing each
586 pattern are also indicated in the left table. CDSs shown by colored triangles include pseudogenes.
587 The 21-bp sequence (renamed *att*-in-PP_1) and other *att* sequences are indicated. Among these
588 sequences, the two indicated by an asterisk are truncated by IS insertion. Several *att* sequences are
589 missing because of deletions. The T3SS effector set (light green triangles) consists of any of the seven
590 effector family/subfamily genes that are encoded by the PP*ompW* EELs shown in Fig 3. More detailed
591 genomic structures of four PP clusters (indicated in bold in the left table) are presented in Fig 4. Types
592 a, c and d include a minor variation; homologous recombination between the second PP*ompW* and
593 the first PPydfJ (type a2), integrase-deficient PPydfJs with or without additional PP integration in
594 tandem (types c2 and c3, respectively), and a region comprising two degraded PPs integrated in
595 tandem between the *rspB* and *trg* genes without PP integration into the *att*-in-PP_2 locus (type d2)
596 are shown.

597 **S6 Fig. Variation in the PP integration patterns in the PP clusters that contained PPs carrying**
598 **potential *att* sites.** (a) Locations of the *att*-in-PP_2 sequences in representative PP genomes. (b)
599 Comparison of the nucleotide sequence of *att*-in-PP_2 among the PPs shown in panel A.

600 **S7 Fig. Variation in the PP integration patterns in the PP clusters that contained PPs carrying**
601 **potential *att* sites.** (a) Schematic representation of the *ydfJ*-flanking region and the PP clusters
602 present at the *ydfJ* locus in three *E. coli* strains. Because both integrase genes of the PPydfJs in strain
603 PV15-279 (PPydfJ-L and PPydfJ-R) have been inactivated by IS insertion, the PPydfJ-L of O26:H11
604 strain 11368 was used for sequence determination of the *attP*-flanking region of PPydfJ by
605 sequencing a PCR amplicon obtained with two primers (indicated by red and blue arrows). (b) The
606 *att* sequences of the four PPydfJs. The *attP*-containing sequence of the PPydfJ-L of strain 11368 was

607 aligned with the *attR*-, *attL*-, and *attB*-containing sequences to define the *att* sequences of each PP.
608 The *ydfJ* sequence of O104:H4 strain C227-11, in which no PPs were not integrated, was used as the
609 *attB* sequence. The 18- or 19-bp *att* sequence that we defined is indicated by uppercase letters.

610 **S8 Fig. The *att*-in-PP_1 and its flanking sequences in PPs and comparison with the *E. coli yecE***
611 **sequence.** (a) The locations of the *att*-in-PP_1 (initially called *att*-in-PP*ompW*) sequences in the
612 genomes of six PP*ompWs* and three other PPs integrated in PPs and in the *yecE* locus of *E. coli*
613 O145:H28 strain 122715. The 21-bp *att*-in-PP_1 sequence and the additional 79-bp sequence
614 homologous to the *yecE* gene are indicated by red and purple, respectively. The *att*-in-PP_2 and *att*-
615 in-PP_3 are also indicated by blue and orange, respectively. The sequences of the two regions
616 indicated by green are conserved between PPs with up to 5 SNPs. The lengths of the two regions are
617 185 bp (left) and 228 bp (right). (b) Alignment of the 100-bp sequences homologous to the *yecE* locus
618 in the nine PPs shown in panel A with the corresponding sequence of the *yecE* locus of strain *E. coli*
619 O145:H28 strain 122715. The 21-bp *att*-in-PP_1 sequence is indicated by uppercase letters. The 100-
620 bp sequences of these PPs were 87% identical to the *yecE* sequence.

621 **S9 Fig. The *att*-in-PP_2 sequence and its flanking sequences.** (a) The locations of the *att*-in-PP_2
622 sequences (blue) in seven PP genomes and on the chromosome of *E. coli* K-12 strain MG1655. The
623 96-bp *att*-in-PP_2 sequences and their flanking sequences (184 bp and 29 bp in length) homologous
624 to the *ykgJ-ecpE* region on the *E. coli* MG1655 chromosome are indicated by blue, pink, and dark
625 brown, respectively. The presence of *stx* and T3SS effector genes in each PP is also indicated. (b)
626 Alignment of the *att*-in-PP_2 and its flanking sequences in the PPs shown in panel a with the
627 corresponding sequence of the *ykgJ-ecpE* region on the *E. coli* MG1655 chromosome. Only the PP
628 genomic regions homologous to the *ykgJ-ecpE* region are shown. The 184-bp regions (pink) of PPs
629 show 83% sequence identity with the *ykgJ-ecpE* region. Note that the 96-bp *att*-in-PP_2 (blue;
630 indicated by uppercase letters) contained 20 SNPs.

631 **S10 Fig. Procedures used to determine the PP integration into the *ompW*, *att*-in-PP*ompW* (later**
632 **in the manuscript, renamed *att*-in-PP_1) and *yecE* loci.** (a) Analysis of PP integration by a
633 BLASTN search. Draft genomes of O145:H28 (n=56) were searched by BLASTN, using the

634 sequences of the *attL*- and *attR*-containing regions of the PPs at *ompW*, *att*-in-*PPompW* and *yecE* in
635 strain 112648 (P08L/R, P09L/R and P12L/R, respectively) as queries. Each query sequence was
636 composed of the sequences from the host chromosome and PP (60 bp each) with the *att* sequence
637 determined in this study (121 bp for P08 and 21 bp for P09/P12) located between them. PP integration
638 at each locus was considered positive when *attL*- and *attR*-containing sequences were both detected
639 (identity threshold: >95%). PP integration in all but two genomes was determined by this analysis.
640 In strains EH1910 and H27V05, although PPs integrated into *yecE* (PP*yecE*) were detected, PP*ompW*
641 was not detected. Unexpectedly, however, the P09L/R sequences (corresponding to the *attL*- and
642 *attR*-containing sequences of the PP-in-PP*ompW*) were detected in EH1910, and a partial P09 *attL*
643 sequence (74.5% coverage) was detected in H27V05. Therefore, the *ompW* and *att*-in-PP*ompW* loci
644 of the two genomes were defined as ‘Others’, and subjected to long PCR analysis along with the
645 identified PPs. (b) Long PCR analysis and sequence determination of PPs. Strategies for five types
646 of analysis are shown. Type I analysis: The genomes of PP*ompW*s that did not contain PPs were
647 divided into three segments and amplified by three long PCRs to obtain the PCR products for genomic
648 sequence determination. Note that the left and right segments included the left and right PP*ompW*-
649 chromosome junctions, respectively (the same strategy was employed in Types II-V analyses). Type
650 II analysis: The genomes of PP*ompW*s that contained an Stx-PP were amplified together with the Stx-
651 PP genomes using 5 or 6 primer pairs to confirm the presence of these PPs and to obtain the PCR
652 products for genome sequence determination. Two primers targeted the *stx* gene (*stx1* or *stx2*). As we
653 detected recombination between the Stx-PP and a PP located at the *ydfJ* locus in two strains (EH1910
654 and 499), a different primer (the leftmost one) was used, thus labeled Type IIb. Type III analysis: In
655 four strains, in which the PP*ompW* contained an Stx-PP, the genome of PP*ompW* and the early region
656 of the Stx-PP were amplified using 4 primer pairs, and only these genomic regions were sequenced.
657 Type IV and V analyses: The genomes of PP*yecE*s were amplified using 2 or 3 primer pairs to obtain
658 the PCR products for genomic sequence determination. When the PP*yecE* contained the *stx* gene
659 (Type IV), two *stx*-targeting primers were used as in Type II analysis. For the PP*yecE* in strain
660 H27V05 (Type Va), only the early region was amplified and sequenced.

661 **S1 Table. *E. coli* strains containing the *att*-in-PPomp*W* sequence at non-*yecE* loci.**

662 **S2 Table. *E. coli* strains containing the *att*-in-PPomp*W* sequence at non-*yecE* loci.**

663 **S3 Table. *E. coli* O145:H28 strains analyzed in this study.**

664 **S4 Table. Complete *E. coli* genomes downloaded from the NCBI database.**

665 **S5 Table. Primers used for PCR amplification for prophage regions.**

666 **S6 Table. Primers used for long PCR analysis.**

667

668 **Figure legends**

669 **Fig 1. Integration sites of the inducible and packageable duplicated Stx2a phages in two STEC**

670 **O145:H28 strains.** (a) The duplicated Stx2a phages and their *att* sequences in strain 112648. The

671 genome structures of three PPs (P08, P09 and P12) are drawn to scale. The *att* sites of each PP are

672 indicated by open (*attL*) or filled (*attR*) symbols (P08, rhombus; P09, circle; P12, square). The *att*

673 sequences of the Stx2a phages (P09 and P12) are shown in the inset. (b) The genome structures of

674 two Stx2a-PPs and a PP integrated into *ompW* (PPomp*W*) in strain 12E129. Sequence homology

675 between the two Stx2a-PPs is also shown, with their integration sites indicated in parentheses.

676 Homologous regions are indicated by shading with different colors according to sequence identity.

677 The genes that were targeted by the PCR primers used in Fig 1c are indicated by asterisks. (c)

678 Detection of packaged DNA of the three PPs in the DNase-treated lysates of strain 12E129 with (+)

679 or without (-) MMC treatment. The chromosome backbone (CB) region was amplified as a negative

680 control.

681 **Fig 2. Variation in the PP content at the *ompW*, *att*-in-PPomp*W* and *yecE* loci in STEC**

682 **O145:H28.** In the left panel, an ML tree of 64 O145:H28 strains is shown. Completely sequenced

683 strains are indicated in bold (plasmids were not finished for strain 2015C-3125). The tree was

684 constructed based on the recombination-free SNPs (3,277 sites) identified on the conserved

685 chromosome backbone (CB) (3,961,936 bp in total length) by RAxML using the GTR gamma

686 substitution model [43]. The reliabilities of the tree's internal branches were assessed using

687 bootstrapping with 1,000 pseudoreplicates. Along with the tree, the geographic and ST/clade

688 information of strains, the presence or absence of PPs at three loci (*ompW*, *att-in-PPompW* and *yecE*)
689 and the types of PPs at the *att-in-PPompW* and *yecE* loci are shown. PPs sequenced in this study and
690 those in the finished genomes are indicated by asterisks and daggers, respectively. Note that the *att-*
691 *in-PPompW* sequence is missing from the *PPompW*s of strains EH2246 and 12E109; a deletion in the
692 latter stain was detected in its draft genome assembly. The bar indicates the mean number of
693 nucleotide substitutions per site. In the right panel, the patterns of the PP content at the three loci are
694 schematically presented. Strains showing each pattern are also indicated in the left panel by diagrams.
695 Note that we detected recombination between the Stx2a-PP at *att-in-PPompW* and a PP present at the
696 *ydfJ* locus that induced a large chromosome inversion in three strains (10942, 499, and EH1910).

697 **Fig 3. Phylogenetic positions of *E. coli* strains carrying *PPompW* and the genome structures of**
698 **their EELs associated with *att-in-PPompW*.** In the upper panel, an ML tree of 92 complete genomes
699 of *E. coli* strains that carry *PPompW* is shown. The tree was constructed based on 109,927 SNP sites
700 in 2,642 core genes and rooted by cryptic Escherichia clade I strain TW15838 (No. AEKA01000000)
701 used as an outlier. Along with the tree, strain IDs used in this paper (see Dataset S2 for more details),
702 phylogroups, and the presence (colored) or absence (open) of *att-in-PPompW* in each strain are
703 indicated. The bar indicates the mean number of nucleotide substitutions per site. In the lower panel,
704 the repertoires of T3SS effector genes that were encoded by the effector exchangeable loci (EELs) in
705 the *PPompW*s containing *att-in-PPompW* are shown. The genomic structures of EELs are drawn to
706 scale. All effector genes were aligned using BLASTN, and orthologous genes (sequence identity;
707 >90%, coverage; >90%) are indicated by the same color. Genes with over 90% identity but less than
708 90% coverage and those containing indels and nonsense mutations in the sequence alignment to intact
709 genes are indicated by asterisks.

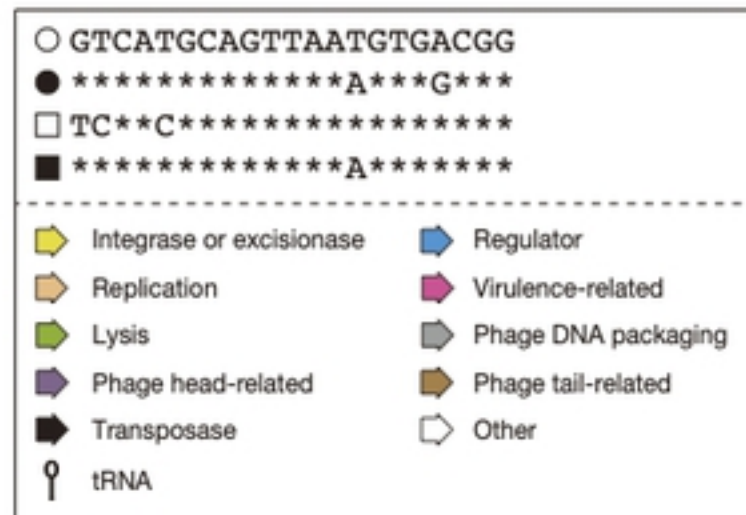
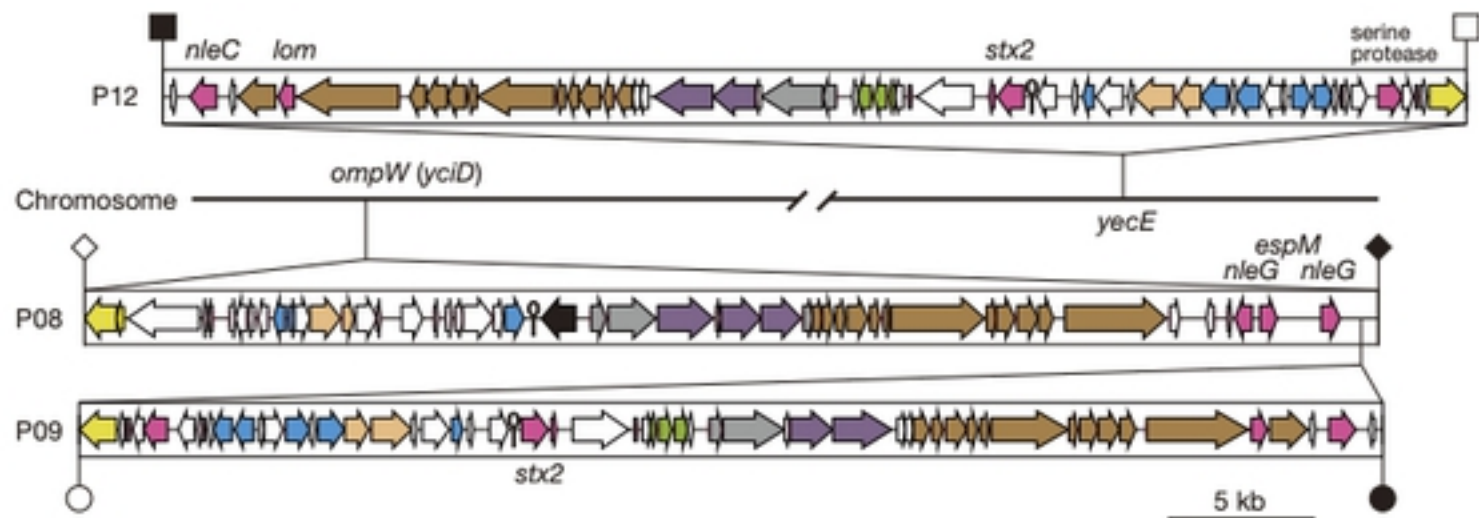
710 **Fig 4. PP clusters that contained PP carrying potential *att* sequences in O157:H7 and O177:H25**
711 **strains.** The genomic structures of three representative PP clusters of the 33 PP clusters found in
712 O157:H7 and that of O177:H25 strains are shown (A, strain FRIK2069; B, strain FRIK944; C,
713 atypical O157:H7 strain PV15-279; D, O177:H25 strain SMN152S1). The identified *att* sequences,
714 coding sequences (CDSs) (including pseudogenes), and ISs in each PP are indicated. T3SS effector

715 genes found in the PP*ompW* EELs (Fig 3) and other effector genes (*nleG* variants) are distinguished
716 by different colors. In panel C, the *att* sequence indicated by an asterisk is truncated by an IS insertion,
717 and integration of an Stx2a-PP into the *att*-in-PP_1 site is schematically presented. The genome
718 structures of all PP clusters identified in this analysis are illustrated in S5 Fig.

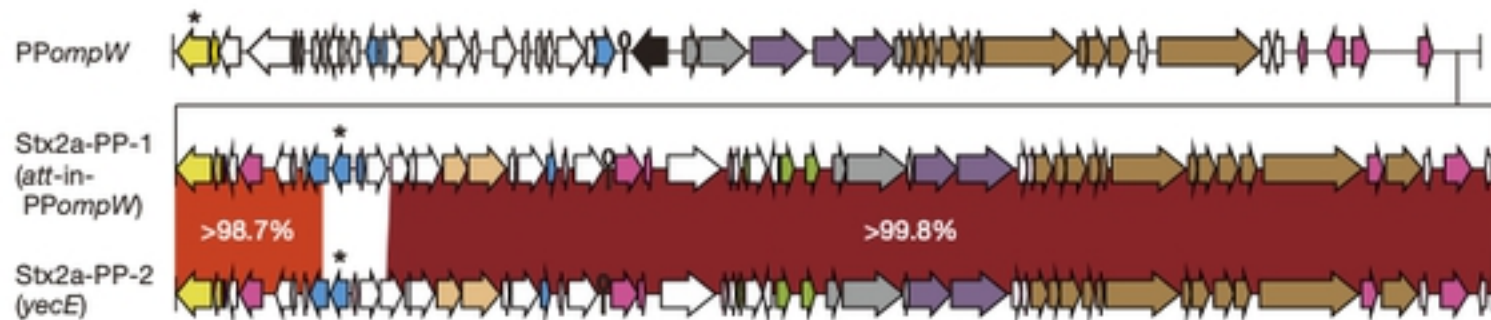
719 **Fig 5. Locations of the *att*-in-PP sequences in PPs and the PP genome regions homologous to *E.***
720 ***coli* chromosome regions.** Three loci in the *E. coli* chromosome showing sequence homology to
721 three identified *att*-in-PP sequences and their flanking sequences are shown at the top. The left- and
722 right-end regions of representative PPs that contained the *att*-in-PP sequences are shown below.
723 Homologous sequences are indicated by the same color. The color used for each *att*-in-PP sequence
724 is the same as that used in Fig 4. See Fig 4 for the details of “PPs-in-PPs” and “PP*mlrA*” and S9 Fig
725 for information on “Sp15 and other PPs”. Alignments of the *att*-in-PP_1 and *att*-in-PP_2 sequences
726 and their flanking sequences with corresponding chromosome sequences are shown in S8 and S9 Figs,
727 respectively.

728 **Fig 6. Summary of the variable PP integration patterns found in this study.**

A



B



C

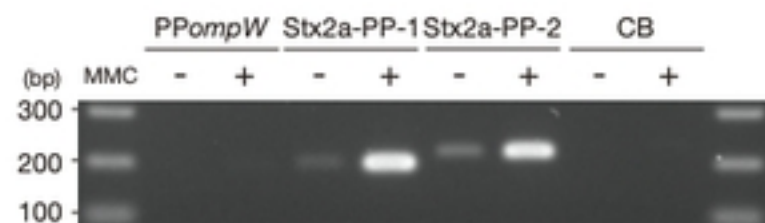


fig1

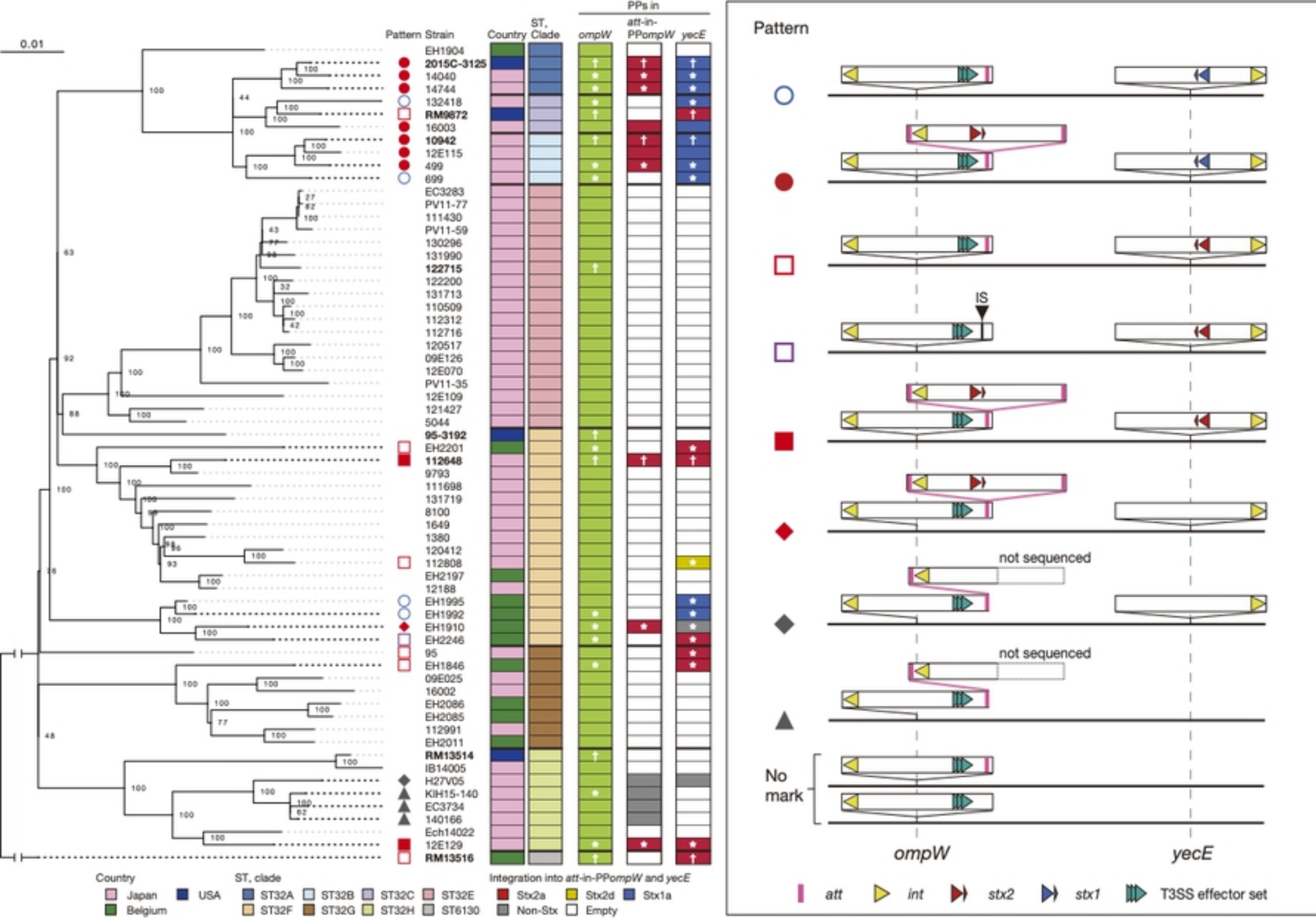
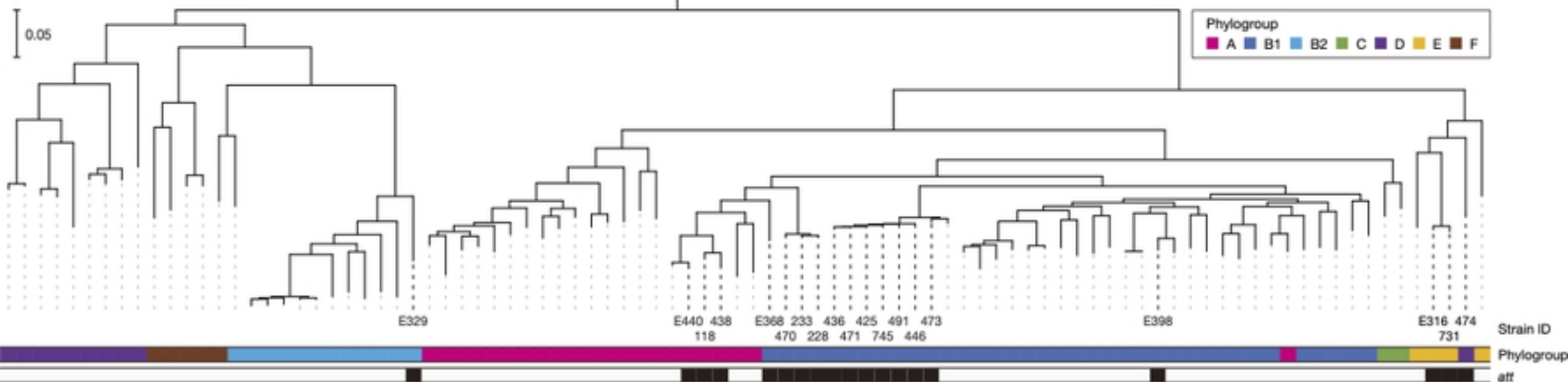
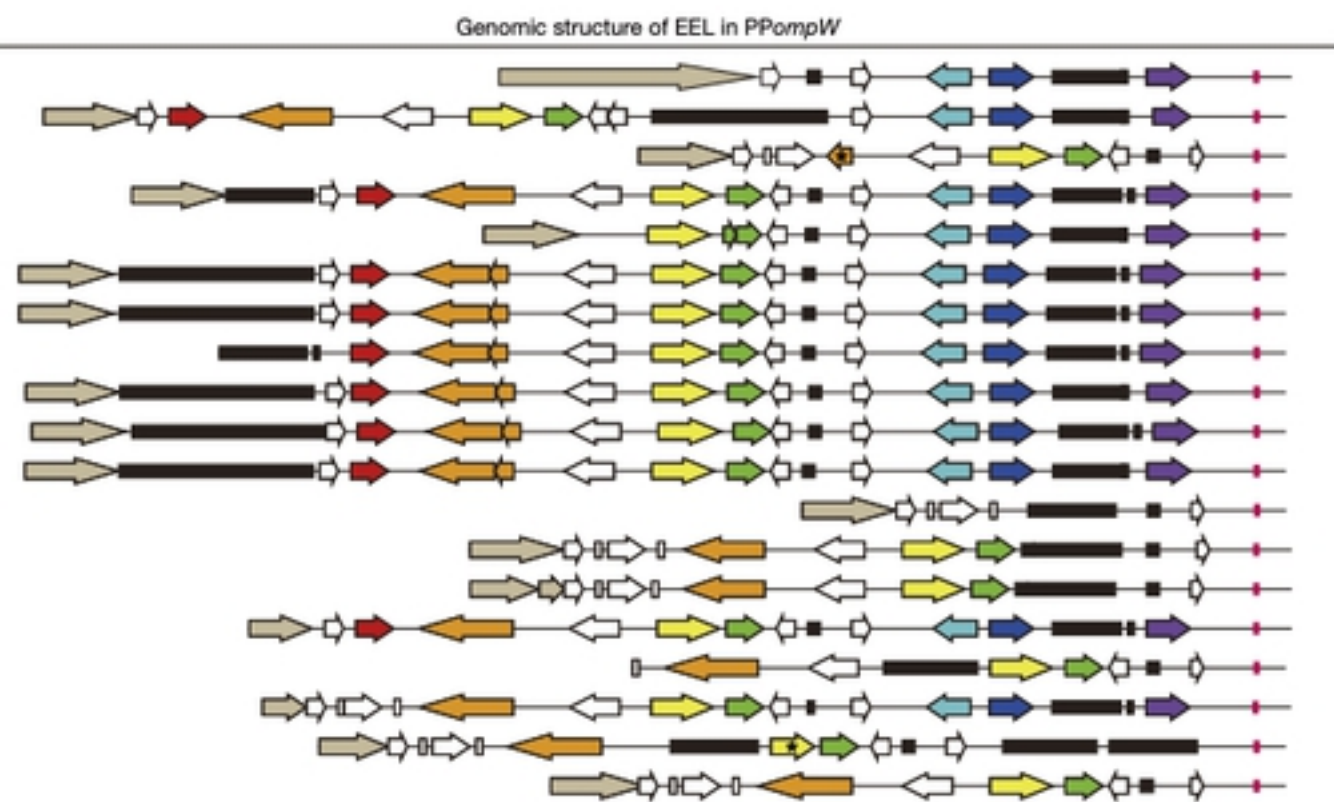


fig2



Group	ID	T3SS effector gene						
		<i>nleG1</i>	<i>nleA</i>	<i>nleH</i>	<i>nleF</i>	<i>nleG2</i>	<i>espM</i>	<i>nleG3</i>
D	E474	-	-	-	-	+	+	+
E	E731	+	+	+	+	+	+	+
E	E316	-	*	+	+	-	-	-
B1	E398	+	+	+	+	+	+	+
B1	E473	-	-	+	*	+	+	+
B1	E446	+	*	+	+	+	+	+
B1	E491	+	*	+	+	+	+	+
B1	E745	+	*	+	+	+	+	+
B1	E425	+	*	+	+	+	+	+
B1	E471	+	*	+	+	+	+	+
B1	E436	+	*	+	+	+	+	+
B1	E228	-	-	-	-	-	-	-
B1	E233	-	+	+	+	-	-	-
B1	E470	-	+	+	+	-	-	-
B1	E368	+	+	+	+	+	+	+
A	E438	-	+	+	+	-	-	-
A	E118	-	+	+	+	+	+	+
A	E440	-	+	*	+	-	-	-
B2	E329	-	+	+	+	-	-	-



1 kb

T3SS effectors

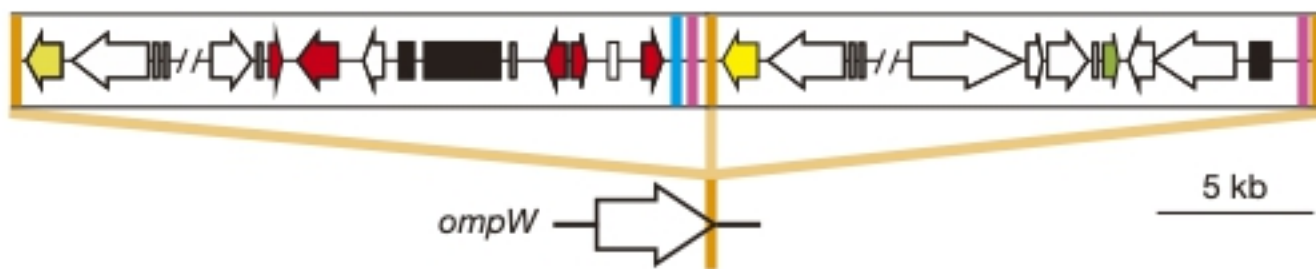
- nleG1* (red arrow)
- nleA* (orange arrow)
- nleH* (yellow arrow)
- nleF* (green arrow)
- nleG2* (light blue arrow)
- espM* (dark blue arrow)
- nleG3* (purple arrow)

Other gene/element

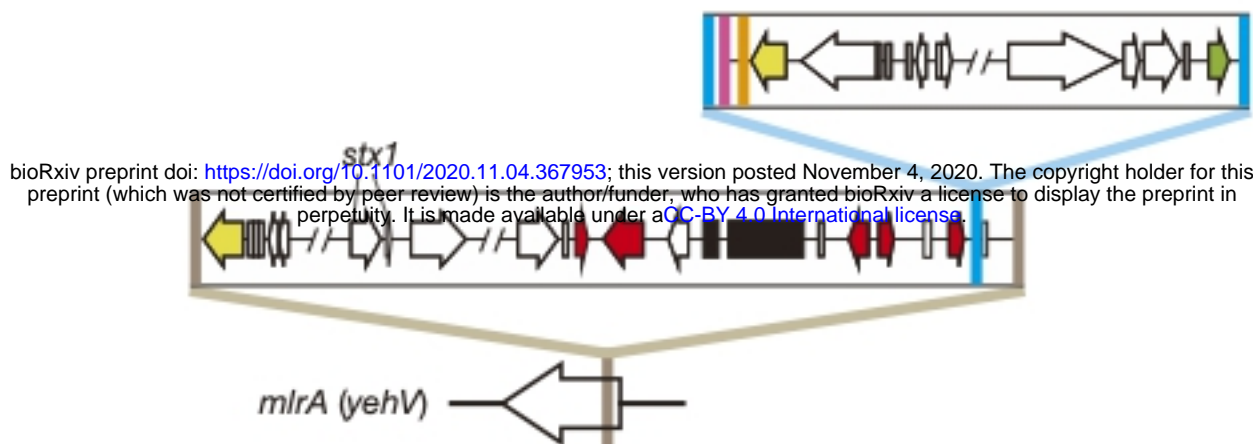
- att site (pink square)
- ISs (black rectangle)
- phage tail genes (grey arrow)
- others (white arrow)

fig3

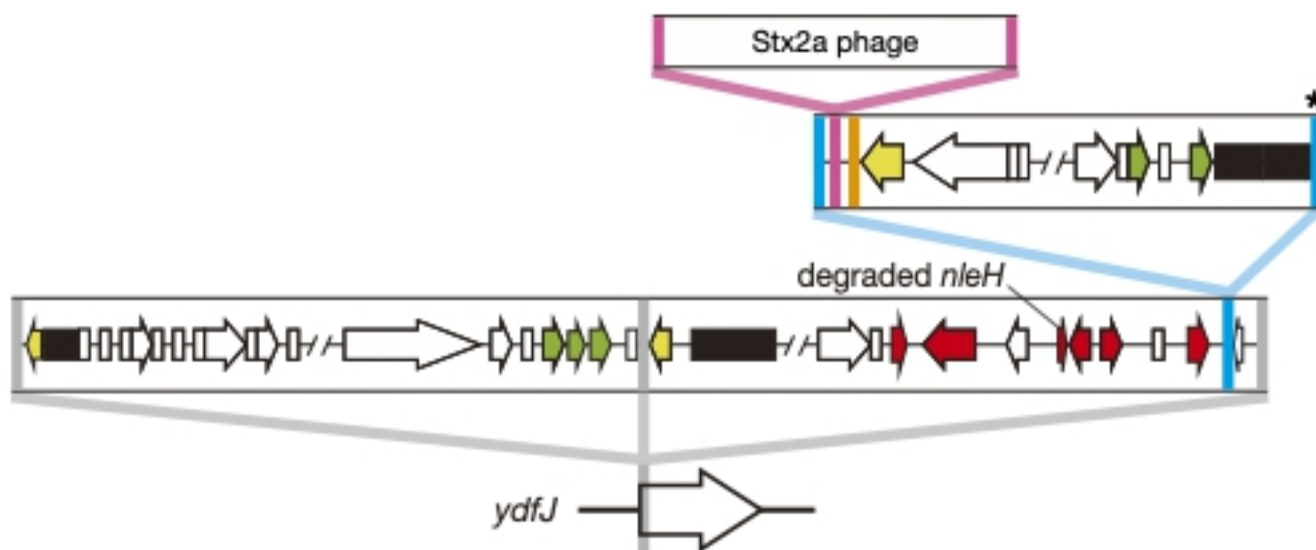
A



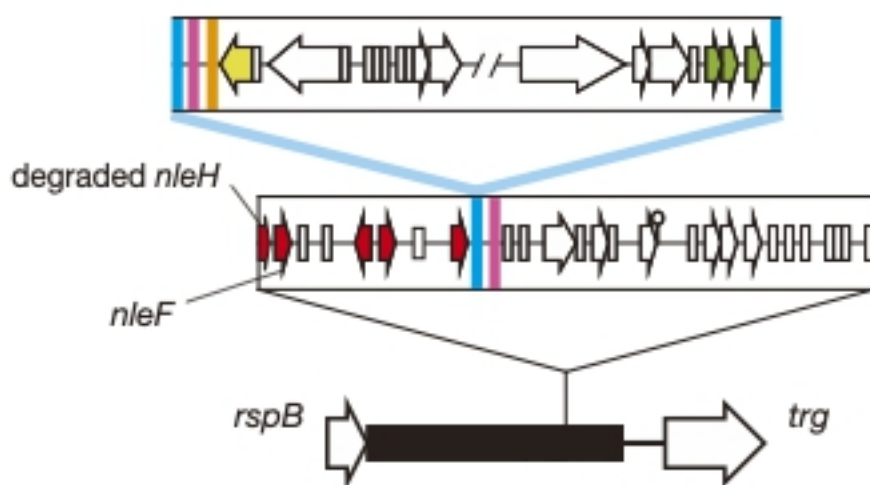
B



C



D



att site

- att-in-PP_1 (*att-in-PPompW*)
- att-in-PP_2
- att-in-PP_3 (*att* for PPompW)
- att for PPmlrA
- att for PPydfJ

gene/element

- ➡ T3SS effector set (*nleG1/A/G2/espM/nleG3*)
- ➡ *nleG* variant
- ➡ integrases
- ⏏ other genes
- IS or IS clusters
- ⚭ tRNA

fig4

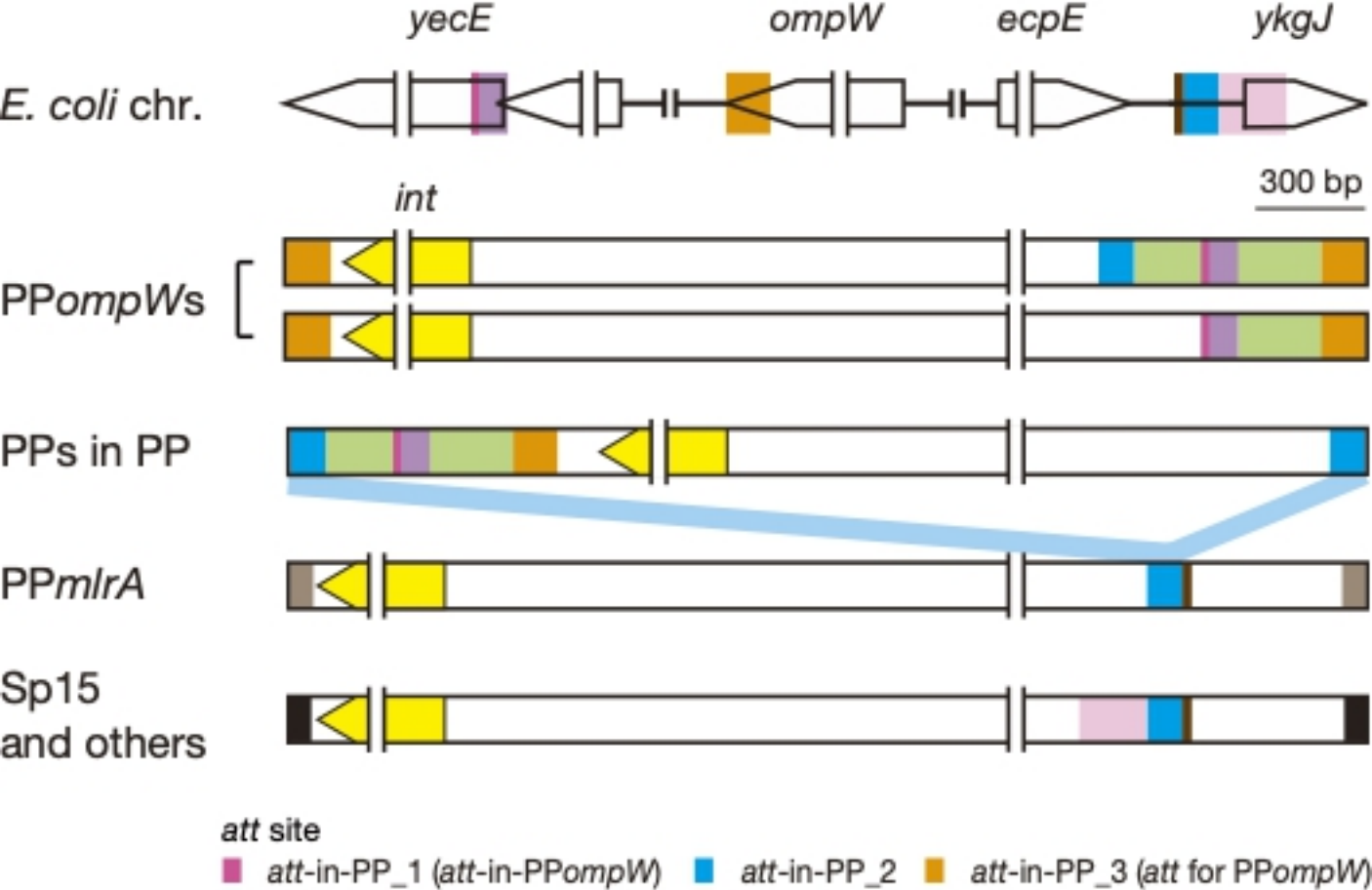


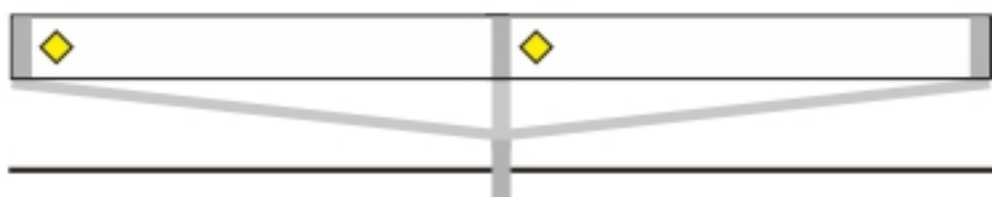
fig5

PP in host



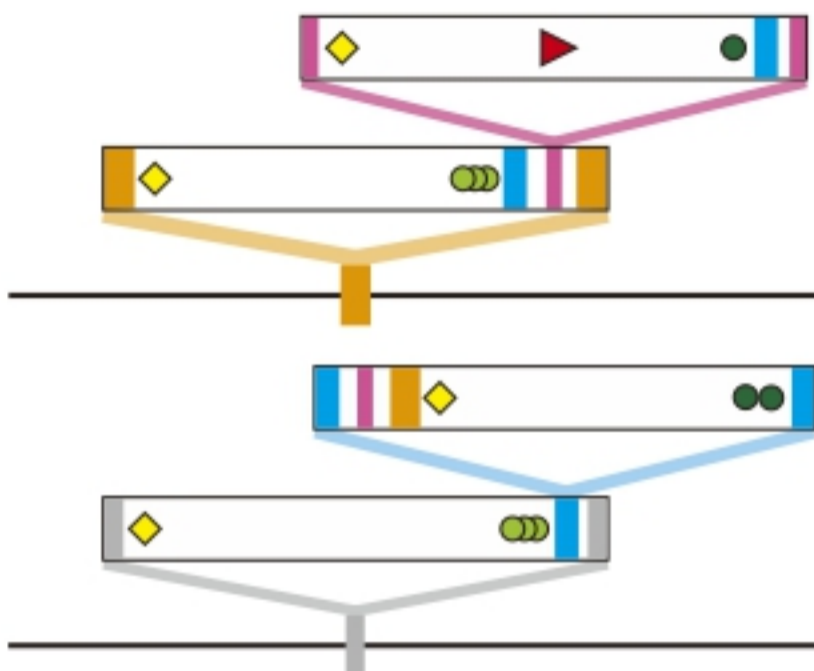
Chromosome

tandem PPs

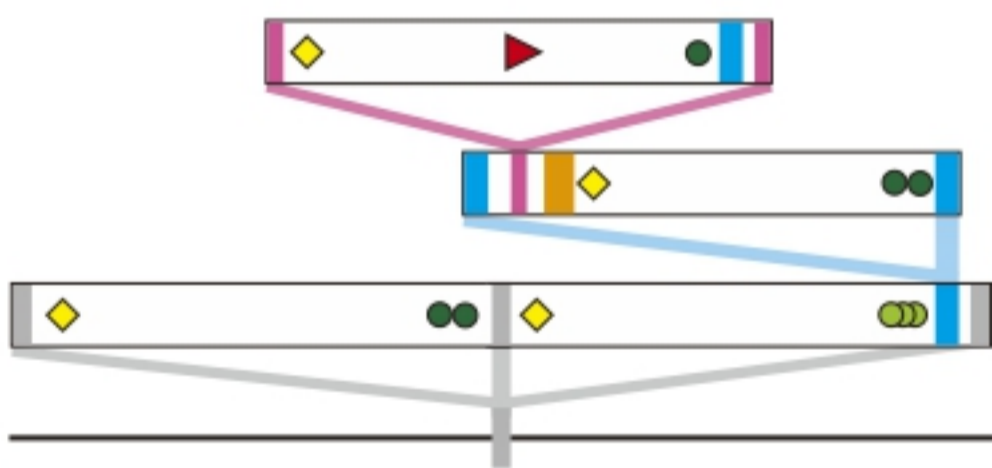


PP-in-PP

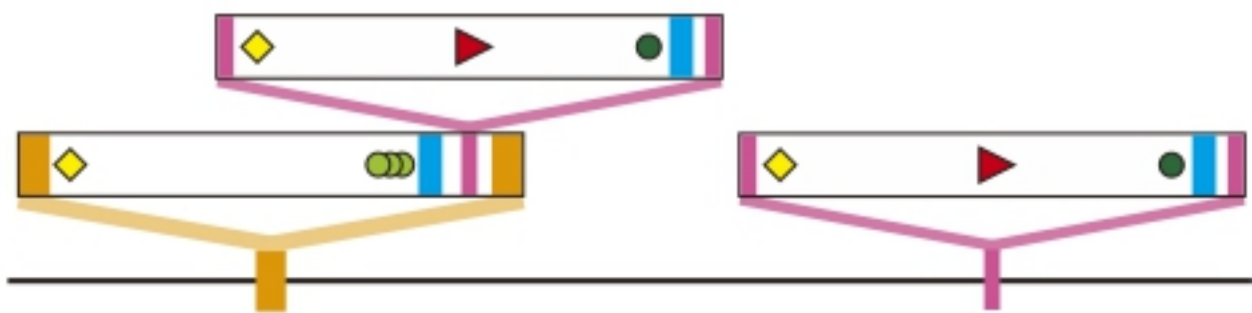
bioRxiv preprint doi: <https://doi.org/10.1101/2020.11.04.367953>; this version posted November 4, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.



combination of 'tandem PPs' and 'PP-in-PP'



intra-chromosome PP duplication



■ *att-in-PP_1* (*att-in-PPompW*) ■ *att-in-PP_2* ■ *att-in-PP_3* (*att* for *PPompW*)
■ *att* for PP ◆ integrase ► *stx2* ● T3SS effector ○ T3SS effector set

fig6