

21 **Abstract**

22 Phage recombinase function units (PRFUs) such as lambda-Red or Rac RecET have been proven
23 to be powerful genetic tools in the recombineering of *Escherichia coli*. Studies have focused on
24 developing such systems in other bacteria as it is believed that these PRFUs have limited efficiency
25 in distant species. However, how the species evolution distance relates to the efficiency of
26 recombineering remains unclear. Here, we present a thorough study of PRFUs to find features that
27 might be related to the efficiency of PRFUs for recombineering. We first identified 59 unique sets
28 of PRFUs in the genus *Corynebacterium* and classified them based on their sequence as well as
29 secondary structure similarities. Then both PRFUs from this genus and other bacteria were chosen
30 for experiment based on sequential and secondary structure similarity as well as species distance.
31 These PRFUs were compared for their ability in mediating recombineering with oligo or double-
32 stranded DNA substrates in *Corynebacterium glutamicum*. We demonstrate that the source of the
33 PRFU is more critical than species distance for the efficiency of recombineering. Our work will
34 provide new ideas for efficient recombineering using PRFUs.

35

36 **Importance**

37 Recombineering using phage recombinase function units (PRFUs) such as lambda-Red or Rac
38 RecET has gained success in *Escherichia coli*, while efforts applying these systems in other
39 bacteria were limited by the efficiency. It is believed that the species distance may be a major
40 reason for the low efficiency. In this study, however, we showed that it is the source of PRFU
41 rather than the species distance that matters for the recombineering in *Corynebacterium*
42 *glutamicum*. Besides, we also showed that the lower transformation efficiency in other bacteria

43 compared to that of *E. coli* could be a major reason for the low performance of heterogeneously
44 expressed RecET. These findings will be helpful for the recombineering using PRFUs.

45 **Keywords:** phage recombinase function unit, RecET, recombineering, efficiency

46

47 **Introduction**

48 Recombineering using a phage recombinase function unit (PRFU), which usually contains one
49 single-stranded annealing protein (SSAP) and one 5'-3' exonuclease (EXO), is powerful in certain
50 bacteria (1-5). Expression of the SSAP alone can mediate recombineering with single-stranded
51 DNA (ssDNA) substrates, while expression of the whole system can use double-stranded DNA
52 (dsDNA) as substrates for genome fragment editing (1, 2).

53 Well-known PRFUs such as lambda-Red or Rac RecET have been used for successful
54 recombineering in *Escherichia coli*. The Red system is phage lambda derived, while the RecET
55 system is from the defective prophage Rac of *E. coli* (1). Lambda-Red contains three adjacent
56 genes *gam*, *bet*, and *exo*, which encode a host nuclease inhibiting the protein Gam, an SSAP Beta,
57 and an EXO, respectively (6). In *E. coli*, it is possible to mediate recombineering using a dsDNA
58 fragment with a homology arm as short as 35 base pairs through the expression of Red-EXO, Red-
59 Beta, and Red-Gam (7). RecET, in contrast, contains only one SSAP RecT and one exonuclease
60 RecE and has been shown to be highly efficient in linear-linear recombination than lambda-Red
61 (2, 8). Both lambda-Red and Rac RecET function by specifically interacting with their respective
62 partners, except for Red-Gam, which is not necessary for dsDNA-mediated recombineering but
63 can improve efficiency (1, 9).

64 Previous studies have shown that PRFU functioned well in related species but performed poorly
65 in distant species, presumably because the exonuclease is more species-specific and could not
66 function well heterogeneously (1, 6). To solve this problem, the PRFUs from closely related
67 species were often searched and developed for specific genetic tools (1, 4, 6, 7). Typically,
68 analogues of Beta or RecT were often searched, and EXOs besides these SSAPs were also
69 compared (10, 11). But how the species evolution distance is related to the efficiency of PRFU
70 remains largely unknown, which makes it difficult to mine efficient PRFUs for recombineering.

71 Here, we present a thorough study of PRFUs in the genomes of the genus *Corynebacterium* to find
72 features that might be related to the efficiency of PRFUs for recombineering. A database to
73 database searching method was first developed to facilitate accurate prediction of novel PRFUs in
74 the genus. Then the identified PRFUs were classified based on both sequential and secondary
75 structure similarity. Typical PRFUs with different sequential and secondary structure similarity as
76 well as host distance to *C. glutamicum* were compared for their ability in mediating
77 recombineering using ssDNA or linear dsDNA substrates. Finally, the recombineering efficiency
78 of the functional PRFUs with their component protein sequence similarity and secondary structure,
79 species relationship, and local genome context were analyzed.

80

81 **Results**

82 **Identification of novel PRFUs in the genus *Corynebacterium***

83 To study the relation of species distance and recombineering efficiency, we first identified the
84 PRFUs in the genus *Corynebacterium*, a taxonomic group with highly diversified species that are
85 of medical, veterinary, or biotechnological relevance (12). Previous studies relating to PRFUs used

86 either reference genomes or nonredundant protein databases for bioinformatics analysis (13-15).
87 As we focused on the relation of species distance and recombineering efficiency, all the bacteria
88 genomes as well as the phage genomes of the genus *Corynebacterium* in the NCBI genome
89 database were used for study.

90 To facilitate the searching process, a database to database method was first developed (Fig. 1). As
91 PRFU is mainly composed of one SSAP and one EXO, which usually appear next to each other,
92 we studied the existence of both among genomes to improve the prediction accuracy. The prophage
93 origin of candidate PRFUs was also verified by showing their location to the prophage region
94 among genomes (16).

95 After iterating for three rounds, we identified candidate PRFUs in 84 of the 429 bacteria genomes
96 and none in the phage genomes studied (Fig. S1a and b). We found that among the 84 bacteria
97 genomes containing PRFUs, 52 of them have PRFUs in the prophage region, while 27 of them
98 have similar proteins coded next to PRFUs as those 52 genomes, and five have no obvious features
99 (Fig. S1c). The PRFU component SSAPs were found in 84 genomes, and EXOs were found in 81
100 genomes. Among them, 54 of the 59 unique SSAPs and 56 of the 57 unique EXOs that were
101 previously annotated as hypothetical proteins were identified.

102 Analysis showed that these SSAPs and EXOs form 59 sets of unique PRFUs in over 23 species
103 (Table 1). Most of these species have been reported to have no PRFU systems or only have SSAP
104 (13, 17, 18). Different from other species, *C. argentoratense* and *C. glutamicum* both separately
105 have one set of incomplete PRFU with only SSAP but no EXO. Not all the genomes sequenced in
106 these 23 species have PRFUs. This indicates that the PRFU occurrence is genome (strain) related
107 but not species related.

108 **Sequence and second structure similarity of the PRFUs and their components**

109 To study the relation of these PRFUs, their component SSAPs and EXOs were classified into five
110 types and three types separately based on both protein sequence similarity and secondary structure
111 similarity (Fig. 2). Each type is conserved in sequence as well as secondary structure and lies
112 within one branch of the phylogenetic trees. Some identified SSAPs or EXOs from other bacteria
113 were also incorporated to construct the tree.

114 For SSAPs, the residues in helices are conserved in the same type, while the residues in strands
115 tend to vary from sequence to sequence. Multiple sequence alignment result showed that SSAP-2
116 shares much similarity to Beta, whereas SSAPs in SSAP-5 were found to have sequence similarity
117 with RecT. A conserved β - β - β - α fold structure was observed in the predicted structure of SSAP-1,
118 SSAP-3, and SSAP-4, while a deformed β - β - β - β - α fold structure was observed in SSAP-2 and
119 SSAP-5 (Fig. 2c). The same β - β - β - α fold has been observed in human Rad52, which was found to
120 be composed of highly conserved amino acid residues that are responsible for ssDNA and dsDNA
121 binding (19). The structure of the N-terminus of human Rad52 (Rad52₁₋₂₁₂) is shown in Fig. 2d as
122 an example. A deformed β - β - α - β - β - α fold of Rad52₁₋₂₁₂ was observed in all SSAP types except
123 SSAP-1. Similar to the Rad52₁₋₂₁₂, both the C-terminus and the N-terminus of the SSAPs are rich
124 in helices. However, compared to the conserved N-terminus, the residues in the C-terminal region
125 of different SSAP types are irregular. This might be consistent with the fact that the conserved N-
126 terminus of SSAP is responsible for homologous pairing, while the C-terminal interacts with EXO
127 or host factors (20-23).

128 Unlike SSAPs, the second structure of EXOs varies considerably between types (Fig. 2e). The
129 residues in the N-terminal half of EXOs are conserved across types, while the residues in the C-

130 terminal half vary considerably even in the same type. The residue Leu-91 and residue Phe-94 of
131 λ Exo have been found to be important for forming the λ Exo-Beta complex (23). These residues
132 lie within a helix of the α - α - β - β fold of the λ Exo (Fig. 2f), which was also observed in EXO-3.
133 These SSAPs and EXOs form eight types of PRFUs in total, which is less than their theoretical
134 combination (Fig. 2g). Typically, SSAP-1 only appears with EXO-1 and forms PRFU-I, while
135 SSAP-5 and EXO-2 only appear with each other and form PRFU-VII, which suggests the
136 conservation of these PRFUs.

137 **The efficiency of PRFUs for recombineering shows no obvious relation to species distance**

138 Previous studies have indicated that the expression of SSAP alone can promote ssDNA-mediated
139 recombineering, while expression of the full PRFU can facilitate genetic manipulation of the
140 genome with dsDNA (1, 6, 7, 24). To understand the relation of species distance and sequence
141 similarity to recombineering efficiency, PRFUs were chosen based on the species evolutionary
142 relationship as well as the PRFU type. Then both oligo-mediated recombineering and dsDNA-
143 mediated recombineering experiments were conducted in a previously constructed strain C.g-kan
144 (-) (25).

145 An oligo targeting the defective *kan*(-) was first used for ssDNA-mediated recombineering to
146 recover the function of *kan* (18, 25). As a result, among the eight SSAPs chosen, five of them could
147 promote ssDNA-mediated recombineering in *C. glutamicum* (Table 2, Sap1–Sap8). Then a
148 selection marker interspaced by an 800bp homology arm targeting *recA* was used as the dsDNA
149 substrate to compare the recombineering ability of the full PRFUs whose SSAPs were effective.
150 Though no complete PRFU was found in the *C. glutamicum* R strain, the gene next to its SSAP
151 coding gene was also tested. It was found that complete PRFUs whose SSAP were functional also

152 showed an effect in recombineering, except for the PRFU from *C. variabile*, which was not tested
153 due to the failure in strain culturing (Table 3). The incomplete PRFU from *C. glutamicum* showed
154 no activity in the dsDNA recombineering experiment.

155 *C. halotolerans*, *C. aurimucosum*, *C. diphtheriae*, and *C. striarum* all have similar distance to *C.*
156 *glutamicum*, but PRFUs from these species were either functional or nonfunctional, which is also
157 true regarding the PRFUs from *C. variabile* and *C. lubricantis* (Fig. S2, Tables 2 and 3). The
158 efficacy of the functional PRFUs varied regarding their host distance to *C. glutamicum*. Both
159 RecET from the distant *E. coli* and PRFUs from the closely related *C. aurimucosum* showed
160 relative high efficiency in promoting recombineering in *C. glutamicum*, while PRFUs from the
161 closely related *C. halotolerans* and *C. diphtheriae* were more than two order of magnitude less
162 efficient.

163 **Local genome context is related to the function of PRFUs**

164 It has been proven that the products of genes besides the genes coding SSAP or EXO, such as Red-
165 Gam, may also affect the recombineering efficiency (1, 4, 7). These PRFUs have been shown to
166 be prophage derived, while the evolutionary relation of prophages might be different from their
167 hosts (26). These inspired us to further investigate the relationship of the local genome context and
168 the recombineering ability of PRFUs.

169 To compare the local genome context of multi-genomes, we considered the five proteins coded
170 upstream and downstream of SSAP. A value S_{dis} of two local genome contexts was calculated to
171 evaluate their distance. As a result, it was found that four of the five functional SSAPs were from
172 genomes sharing a similar local genome context (Fig. 3a). Through the analysis of the distance of
173 the local genome context of other species to that of *C. glutamicum*, it was found that the local

174 genome contexts of these functional ones have obviously smaller and wider distance range (S_{dis}
175 from 0 to 35) than others (S_{dis} from 35 to 40, Fig. 3b).

176 To test whether the conservation of the local genome context might be related to the
177 recombineering function of PRFUs in *C. glutamicum*, we verified the function of three additional
178 SSAPs whose S_{dis} are smaller than 35 (Table 2, Sap10–Sap12). All the SSAPs in SSAP-5 have S_{dis}
179 less than 35, while RecT and SSAP from *C. aurimucosum* are all highly similar to SSAPs in SSAP-
180 5 and showed high performance in recombineering. Therefore, we also tested four SSAPs from
181 other bacteria who were highly identical to those in SSAP-5 to test whether SSAPs that are highly
182 similar to SSAPs in SSAP-5 may be efficient (Table 2, Sap13–Sap16). As a result, all the selected
183 SSAPs sharing a conserved local genome context were proven to be functional in mediating
184 recombineering, while SSAPs that were highly conserved as SSAP-5 from other bacteria showed
185 no recombineering efficiency. This indicates that the conservation of the local genome contexts
186 might be related to the function of PRFUs.

187

188 Discussion

189 Although genome sequence analysis has grown rapidly these years, few PRFUs have been
190 identified in the genus *Corynebacterium*. This, together with the phage origin of verified PRFUs,
191 led us to investigate all the genomes sequenced in the genus. As each genome represents a strain
192 record, the occurrence of PRFUs in both reference and non-reference genomes and the irregular
193 distribution pattern of PRFUs among genomes within the same species verifies the idea that the
194 existence of PRFUs is strain related but not species related. When searching for PRFUs among all
195 the plasmids of the genus *Corynebacterium*, only one SSAP coding gene was found on the plasmid

196 of the *C. glutamicum* R strain, which may explain why no PRFU was previously found in *C.*
197 *glutamicum* when searching its reference genome (13, 18). These indicate that the prediction of
198 PRFUs using reference genomes might lead to errors in conclusion about their existence in the
199 species, while using all the genomes as this study did or using nonredundant protein databases as
200 in other reports may avoid this problem (1, 7, 27).

201 To compare the ability of PRFUs, we verified several novel identified ones with both ssDNA and
202 linear dsDNA-mediated heterogeneous recombineering in *C. glutamicum*. Some PRFUs tested
203 showed a relatively high ability in conducting recombineering as previously reported tools in other
204 bacteria (1, 6, 7). The recombination efficiency was as high as about 10^{-2} , which means there will
205 be one recombinant over one hundred cells successfully transformed. This suggests the great
206 potential of these PRFUs for recombineering in the genus *Corynebacterium*. However, as we have
207 shown in addition to other studies (6, 15), the relative low transformation efficiency could be a
208 major limit of recombineering efficiency in other bacteria other than the PRFUs used.

209 It has been believed that the efficiency of PRFU for recombineering was related to the species
210 distance, which led to the study of PRFUs in the native species of several bacteria (6, 10, 27). In
211 this study, the host species where most of the tested PRFUs belong to share similar distances to *C.*
212 *glutamicum*, but the tested SSAPs were either functional or non-functional. RecET from *E. coli*
213 showed relatively high ability in promoting recombineering as those from closely related species
214 of *C. glutamicum* at a recombination efficiency as high as about 10^{-2} . The phenomenon that PRFUs
215 from distantly related species performed similarly as those from closely related species (at the
216 same order of magnitude) has also been observed in other reports (1, 7, 18). These suggest an
217 alternative explanation of the relation of the species evolution distance and recombineering
218 efficiency. It was observed in this study that the local genome context conservation in the genus

219 *Corynebacterium* is related to the recombineering function of PRFUs in *C. glutamicum*, but RecET
220 from *E. coli* shows no such conservation. Considering the different performances of PRFUs in this
221 study and previously in *E. coli*(1), it seems that the efficiency of PRFUs for recombineering is
222 dependent on their source.

223 In this study, we developed a database to database searching method to facilitate accurate
224 prediction of novel PRFUs in a species that was reported to have no such system. Analysis of the
225 distribution of PRFUs among genomes indicates that their occurrence is related to strains but not
226 species. This is consistent with their features as mobile elements and may provide new insights for
227 the data mining of mobile elements. Through the analysis of the recombineering efficiency of
228 PRFUs with their component protein sequence similarity, second structure, species relationship,
229 as well as local genome context, we demonstrate that the efficiency of PRFUs for recombineering
230 is dependent on the source of the recombinase. This finding will provide new insights for the study
231 of PRFUs for recombineering.

232

233 **Methods**

234 **Analysis of PRFU across genomes**

235 Genomes of the genus *Corynebacterium*, including chromosomes, plasmids, and phages were
236 downloaded from NCBI (<https://www.ncbi.nlm.nih.gov/>) genome database. The conserved
237 domain data of SSAP and EXO from the Conserved Domains Database of NCBI (Accession
238 number: c104285, c104500, c101936, and pfam09588) were also downloaded. Then the following
239 process was iterated.

240 1. Multiple sequence alignments were done to get the consensus of SSAPs and EXOs and
241 construct database SSAPdb as well as EXOdb.

242 2. Then genomes of the genus *Corynebacterium* were searched for potential SSAP and EXO in
243 the SSAPdb and EXOdb separately through PSI-BLAST for one iteration.

244 3. Potential SSAP and EXO were extracted and added to the SSAPdb and EXOdb, respectively.

245 This process was iterated for three times until no new potential SSAP or EXO was found. Potential
246 SSAPs were chosen from the result of PSI-BLAST using a cut off E value of 10^{-6} . Potential EXOs
247 were chosen by setting a cut off E value as 10^{-5} .

248 Prophage was predicted by PHASTER (<http://phaster.ca/>) (16). The genome location of PRFUs as
249 well as the local genome context were used for the verification of their prophage origin.

250 **Protein sequence and structure analysis**

251 Multiple protein sequence alignment was performed with MUSCLE (28). The alignment results
252 were then edited manually to get the consensus parts and shown as the default settings in ClustalX
253 (29).

254 All of the second structures of the alignments were predicated by Jpred4
255 (<http://www.compbio.dundee.ac.uk/jpred4/index.html>) (30).

256 **Alignment of the local genome context**

257 The adjacent five proteins coded upstream and downstream of the SSAP of each genome were
258 extracted to evaluate the local genome context similarity if available. First, the similarity of these
259 proteins was analyzed through PSI-BLAST. Protein pairs with an E value less of than 10^{-5} were

260 grouped together. The groups with proteins of similar annotation or the proteins that were not
261 grouped in the previous step but had similar annotation were merged into one group. The proteins
262 that were not grouped yet were each assigned a group separately.

263 A similarity matrix of the local genome context was then calculated through the similarity score
264 S_{sim} between each pair of them. The similarity score S_{sim} contained two parts, S_{sam} and S_{loc} . S_{sam}
265 was calculated as the sum number of protein pairs from the two genomes that are of same group.
266 S_{loc} was the sum of S_i , which was calculated in the following way: for each protein coded upstream
267 of SSAP and downstream of SSAP, its position was assigned as “site.” Site i has a range from 1 to
268 11, corresponding to the position from the farthest upstream position of SSAP to the farthest
269 downstream. If the two proteins of both genomes at the site i belong to the same group, then S_i was
270 assigned 1, 2, 3, 4, 5, 6, 5, 4, 3, 2, or 1 for each i from 1 to 11 accordingly. This can be written as

$$271 \quad S_{sim} = \sum_{i=1}^{11} S_i + \sum S_{sam}. \quad (1)$$

272 The distance matrix of these genomes was obtained from the similarity matrix, which used a score
273 S_{dis} . S_{dis} that was calculated as

$$274 \quad S_{dis} = 47 - S_{sim}. \quad (2)$$

275 The number 47 was S_{sim} for a pair of completely same local genome context that had an ideal S_{dis}
276 value of 0. As a different site was assigned a different weight, we called these site-weighted scores.

277 **Phylogenetic analysis**

278 A genome phylogenetic tree was constructed using the 16S rRNA of the representative genomes.
279 A protein phylogenetic tree was constructed using the multiple sequence alignment results of
280 SSAPs and EXOs. For the construction of a local genome context phylogenetic tree, a distance

281 matrix was used as the input. The phylogenetic trees were constructed by MEGA6 using the
282 Neighbor-Joining method with the default settings except for the Interior-branch Test value, which
283 was set as 1000 (31).

284 **Strain cultivation**

285 The *C. glutamicum* strain C.g-kan (-) was cultured in Luria-Bertani (LB) or brain heart infusion-
286 supplemented (BHIS) medium. For the preparation of the electrocompetent cells, the strains were
287 first cultured in LB medium overnight at 30°C. The cultures were then added to 100 ml BHIS
288 medium at a start optical density of 0.3 and cultured at 30°C until the optical density reached 0.9,
289 with isopropyl β -D-1-thiogalactopyranoside (IPTG) and antibiotics added according to the
290 circumstance. The culture was then chilled on ice for 30 min and centrifuged at 4°C to collect the
291 cells. The sediments were washed four times using 30 ml 10% (v/v) glycerol and resuspended in
292 1 ml 10% (v/v) glycerol. Then a volume of 100 μ l was divided into each tube and stored at -80°C
293 for later experiments. To maintain the plasmids, a final concentration of 100 μ g/ml spectinomycin,
294 15 μ g/ml kanamycin, or 10 μ g/ml chloromycetin was added to the corresponding cultures. For
295 induction purposes, a final concentration of 0.25 mmol/L IPTG was used (25, 32).

296 **Statistics**

297 The ratio of the number of sequence finished genomes with PRFUs and sequence finished genomes
298 as well as the ratio of the number of sequence unfinished genomes with PRFUs and sequence
299 unfinished genomes were used for a t-test (Table S1).

300 **Recombineering with PRFUs**

301 The plasmid, ssDNA, and dsDNA substrates were transformed into *C. glutamicum*
302 electrocompetent cells by electroporation. Different substrates were pre-chilled and added to the
303 electrocompetent cells according to the experiment. For recombineering with ssDNA, 2 µg oligo
304 was added to the competent cells. For plasmid and dsDNA, 500 ng or 1 µg substrate was used. The
305 mixture was then used for electroporation at 1.8 kV/mm, 200 Ω, 25 µF. Then 1 ml 46°C BHIS
306 medium was instantly added. The mixture was heat shocked at 46°C for 6 min and then cultured
307 for outgrowth at 30°C, 200 rpm up to 1 h for strain construction purposes or 30°C, 150 rpm for 5
308 h for the transformation of ssDNA, dsDNA, and control plasmid. Cultures were then plated directly
309 after collecting the cells for most of the time. To count the number of viable cells, 100 µl of 10⁻⁴
310 dilution was plated on BHIS with chloromycetin added (15, 18).

311 Fragments were prepared through overlap extension PCR using TransStart FastPfu Fly DNA
312 Polymerase (TransGen, AP231-03) and collected using a DNA Gel Extraction Kit (TsingKe,
313 GE101-200). Results were verified by agarose gel electrophoresis and Sanger sequencing.

314

315 **Acknowledgments**

316 This work was supported by grants from the National Natural Science Foundation of China
317 [31730003, 31670077] and Natural Science Foundation of Shandong Province [ZR2017ZB0210].

318 We thank Haiying Jin for her help in the recombineering with ssDNA experiment. We also thank
319 Qilong Qin for reviewing the manuscript.

320 We thank LetPub (www.LetPub.com) for its linguistic assistance during the preparation of this
321 manuscript.

322

323 **Author contributions**

324 YZ designed and conducted the experiment, prepared and revised the manuscript. QS designed,
325 supervised the study and revised the manuscript. QW revised the manuscript. TY participated the
326 experiment.

327

328 **References**

- 329 1. Datta S, Costantino N, Zhou X, Court DL. 2008. Identification and analysis of recombineering
330 functions from Gram-negative and Gram-positive bacteria and their phages. *Proc Natl Acad Sci U*
331 *S A* 105:1626-31.
- 332 2. Zhang Y, Buchholz F, Muyrers JP, Stewart AF. 1998. A new logic for DNA engineering using
333 recombination in *Escherichia coli*. *Nat Genet* 20:123-8.
- 334 3. Murphy KC. 1998. Use of bacteriophage lambda recombination functions to promote gene
335 replacement in *Escherichia coli*. *J Bacteriol* 180:2063-71.
- 336 4. Yang P, Wang J, Qi Q. 2015. Prophage recombinases-mediated genome engineering in
337 *Lactobacillus plantarum*. *Microb Cell Fact* 14:154.
- 338 5. Vellani TS, Myers RS. 2003. Bacteriophage SPP1 Chu is an alkaline exonuclease in the SynExo
339 family of viral two-component recombinases. *J Bacteriol* 185:2465-74.
- 340 6. van Kessel JC, Hatfull GF. 2007. Recombineering in *Mycobacterium tuberculosis*. *Nat Methods*
341 4:147-52.
- 342 7. Yin J, Zhu H, Xia L, Ding X, Hoffmann T, Hoffmann M, Bian X, Muller R, Fu J, Stewart AF, Zhang Y.
343 2015. A new recombineering system for *Photorehabdus* and *Xenorhabdus*. *Nucleic Acids Res*
344 43:e36.
- 345 8. Fu J, Bian X, Hu S, Wang H, Huang F, Seibert PM, Plaza A, Xia L, Muller R, Stewart AF, Zhang Y.
346 2012. Full-length RecE enhances linear-linear homologous recombination and facilitates direct
347 cloning for bioprospecting. *Nat Biotechnol* 30:440-6.
- 348 9. Muyrers JP, Zhang Y, Buchholz F, Stewart AF. 2000. RecE/RecT and Redalpha/Redbeta initiate
349 double-stranded break repair by specifically interacting with their respective partners. *Genes*
350 *Dev* 14:1971-82.
- 351 10. Swingle B, Bao Z, Markel E, Chambers A, Cartinhour S. 2010. Recombineering using RecTE from
352 *Pseudomonas syringae*. *Appl Environ Microbiol* 76:4960-8.
- 353 11. Van Pijkeren JP, Neoh KM, Sirias D, Findley AS, Britton RA. 2012. Exploring optimization
354 parameters to increase ssDNA recombineering in *Lactococcus lactis* and *Lactobacillus reuteri*.
355 *Bioengineered* 3:209-17.

- 356 12. Oliveira A, Oliveira LC, Aburjaile F, Benevides L, Tiwari S, Jamal SB, Silva A, Figueiredo HCP,
357 Ghosh P, Portela RW, Azevedo VAD, Wattam AR. 2017. Insight of Genus *Corynebacterium*:
358 Ascertaining the Role of Pathogenic and Non-pathogenic Species. *Frontiers in Microbiology* 8.
359 13. Resende BC, Rebelato AB, D'Afonseca V, Santos AR, Stutzman T, Azevedo VA, Santos LL, Miyoshi
360 A, Lopes DO. 2011. DNA repair in *Corynebacterium* model. *Gene* 482:1-7.
361 14. Garriss G, Poulin-Laprade D, Burrus V. 2013. DNA-damaging agents induce the RecA-
362 independent homologous recombination functions of integrating conjugative elements of the
363 SXT/R391 family. *J Bacteriol* 195:1991-2003.
364 15. Huang Y, Li L, Xie S, Zhao N, Han S, Lin Y, Zheng S. 2017. Recombineering using RecET in
365 *Corynebacterium glutamicum* ATCC14067 via a self-excisable cassette. *Sci Rep* 7:7916.
366 16. Arndt D, Grant JR, Marcu A, Sajed T, Pon A, Liang Y, Wishart DS. 2016. PHASTER: a better, faster
367 version of the PHAST phage search tool. *Nucleic Acids Res* 44:W16-21.
368 17. Nakamura Y, Nishio Y, Ikeo K, Gojobori T. 2003. The genome stability in *Corynebacterium* species
369 due to lack of the recombinational repair system. *Gene* 317:149-55.
370 18. Binder S, Siedler S, Marienhagen J, Bott M, Eggeling L. 2013. Recombineering in
371 *Corynebacterium glutamicum* combined with optical nanosensors: a general strategy for fast
372 producer strain generation. *Nucleic Acids Res* 41:6360-9.
373 19. Kagawa W, Kurumizaka H, Ishitani R, Fukai S, Nureki O, Shibata T, Yokoyama S. 2002. Crystal
374 structure of the homologous-pairing domain from the human Rad52 recombinase in the
375 undecameric form. *Mol Cell* 10:359-71.
376 20. Smith CE, Bell CE. 2016. Domain Structure of the Redbeta Single-Strand Annealing Protein: the C-
377 terminal Domain is Required for Fine-Tuning DNA-binding Properties, Interaction with the
378 Exonuclease Partner, and Recombination in vivo. *J Mol Biol* 428:561-78.
379 21. Lloyd JA, McGrew DA, Knight KL. 2005. Identification of residues important for DNA binding in
380 the full-length human Rad52 protein. *J Mol Biol* 345:239-49.
381 22. Ploquin M, Bransi A, Paquet ER, Stasiak AZ, Stasiak A, Yu X, Cieslinska AM, Egelman EH, Moineau
382 S, Masson JY. 2008. Functional and structural basis for a bacteriophage homolog of human
383 RAD52. *Curr Biol* 18:1142-6.
384 23. Caldwell BJ, Zakharova E, Filsinger GT, Wannier TM, Hempfling JP, Chun-Der L, Pei D, Church GM,
385 Bell CE. 2019. Crystal structure of the Redbeta C-terminal domain in complex with lambda
386 Exonuclease reveals an unexpected homology with lambda Orf and an interaction with
387 *Escherichia coli* single stranded DNA binding protein. *Nucleic Acids Res*
388 doi:10.1093/nar/gky1309.
389 24. Lopes A, Amarir-Bouhram J, Faure G, Petit MA, Guerois R. 2010. Detection of novel
390 recombinases in bacteriophage genomes unveils Rad52, Rad51 and Gp2.5 remote homologs.
391 *Nucleic Acids Res* 38:3952-62.
392 25. Su TY, Jin HY, Zheng Y, Zhao Q, Chang YZ, Wang Q, Qi QS. 2018. Improved ssDNA recombineering
393 for rapid and efficient pathway engineering in *Corynebacterium glutamicum*. *Journal of*
394 *Chemical Technology and Biotechnology* 93:3535-3542.
395 26. Pratama AA, Chaib De Mares M, van Elsas JD. 2018. Evolutionary History of Bacteriophages in
396 the Genus *Paraburkholderia*. *Front Microbiol* 9:835.
397 27. van Pijkeren JP, Britton RA. 2012. High efficiency recombineering in lactic acid bacteria. *Nucleic*
398 *Acids Res* 40:e76.
399 28. Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
400 *Nucleic Acids Res* 32:1792-7.
401 29. Thompson JD, Gibson TJ, Plewniak F, Jeanmougin F, Higgins DG. 1997. The CLUSTAL_X windows
402 interface: flexible strategies for multiple sequence alignment aided by quality analysis tools.
403 *Nucleic Acids Res* 25:4876-82.

- 404 30. Drozdetskiy A, Cole C, Procter J, Barton GJ. 2015. JPred4: a protein secondary structure
405 prediction server. *Nucleic Acids Res* 43:W389-94.
- 406 31. Tamura K, Stecher G, Peterson D, Filipski A, Kumar S. 2013. MEGA6: Molecular Evolutionary
407 Genetics Analysis version 6.0. *Mol Biol Evol* 30:2725-9.
- 408 32. Yu X, Jin H, Liu W, Wang Q, Qi Q. 2015. Engineering *Corynebacterium glutamicum* to produce 5-
409 aminolevulinic acid from glucose. *Microb Cell Fact* 14:183.

410

411

412 **Supplemental material**

413 **Figure S1. PRFUs in the genus *Corynebacterium*.** a Genomes with or without PRFUs. b

414 Genomes in species with or without PRFUs. c The phage origin of predicted PRFUs. F: genomes

415 with PRFUs, NF: genomes without PRFUs. CL: PRFUs found in genomes that have been

416 classified into species, UCL: PRFUs found in genomes that have not been classified into any

417 species, F: genomes with PRFUs, NF: genomes without PRFUs. Numbers indicate the genome

418 number.

419 **Figure S2. Species 16sRNA distance to *C. glutamicum*.**

420 **Table S1. The occurrence of PRFU and genome assemble level.**

421 **Table S2. Strains and plasmids used in the study.**

422

423

424

425

Table 1. Distribution of PRFUs in the genus *Corynebacterium*

Species	Genome ^a	SSAP ^a	EXO ^a	PRFU ^b
<i>C. argentoratense</i>	3	1	0	1
<i>C. aurimucosum</i>	6	2	2	1
<i>C. diphtheriae</i>	63	22	22	5
<i>C. falsenii</i>	2	1	1	1
<i>C. freneyi</i>	1	1	1	1
<i>C. glutamicum</i>	20	1	0	1
<i>C. halotolerans</i>	3	1	1	1
<i>C. humireducens</i>	2	2	2	1
<i>C. jeikeium</i>	21	5	5	5
<i>C. kroppenstedtii</i>	4	4	4	3
<i>C. lactis</i>	1	1	1	1
<i>C. lubricantis</i>	1	1	1	1
<i>C. minutissimum</i>	3	1	1	1
<i>C. pyruviciproducens</i>	1	1	1	1
<i>C. simulans</i>	3	1	1	1
<i>C. striatum</i>	10	6	6	4
<i>C. timonense</i>	2	2	2	2
<i>C. ulcerans</i>	17	4	4	4
<i>C. ulceribovis</i>	1	1	1	1
<i>C. variabile</i>	3	1	1	1

<i>C. bouchesdurhonense</i>	1	1	1	1
<i>C. phoceense</i>	1	1	1	1
<i>C. provencense</i>	2	1	1	1
<i>C. sp</i> ^b	115	22	21	22

426 ^a Number of SSAP or EXO identified or the number of genomes. Species without PRFUs are

427 not listed.

428 ^b Unique PRFU numbers in each species, including incomplete ones.

429 ^c Genomes in the genus *Corynebacterium* that have not been classified.

Table 2. Recombineering with ssDNA using different SSAPs

Strain	Substrate ^a			Viable cells (x 10 ⁷) ^a	Transformation efficiency	Recombination efficiency	SSAP type	SSAP source
	ssDNA-	ssDNA+	plasmid (x 10 ⁴)					
Sap0	1 ± 1	0	0.67 ± 0.25	41.90 ± 13.11	1.36 x 10 ⁻⁵	0	-	Empty vector
Sap1	2 ± 3	444 ± 40	2.13 ± 1.16	9.93 ± 3.93	2.15 x 10 ⁻⁴	2.08 x 10 ⁻²	SSAP-5	<i>Corynebacterium aurimucosum</i> DSM 44827
Sap2	0	61 ± 30	3.77 ± 1.93	4.03 ± 1.20	9.34 x 10 ⁻⁴	1.62 x 10 ⁻³	SSAP-4	<i>Corynebacterium diphtheriae</i> HC07
Sap3	0	136 ± 55	4.21 ± 0.67	16.13 ± 6.38	2.61 x 10 ⁻⁴	3.23 x 10 ⁻³	SSAP-3	<i>Corynebacterium glutamicum</i> R
Sap4	0	14 ± 3	2.43 ± 0.41	30.63 ± 7.03	7.92 x 10 ⁻⁵	5.76 x 10 ⁻⁴	SSAP-3	<i>Corynebacterium halotolerans</i> 307 CHAL
Sap5	0	0	1.20 ± 0.71	18.10 ± 9.92	6.65 x 10 ⁻⁵	0	SSAP-1	<i>Corynebacterium kroppenstedtii</i> DNF00591
Sap6	2 ± 2	0	3.26 ± 1.23	41.73 ± 8.27	7.82 x 10 ⁻⁵	0	SSAP-2	<i>Corynebacterium lubricantis</i> DSM 45231
Sap7	0	0	3.26 ± 0.60	34.00 ± 9.64	9.60 x 10 ⁻⁵	0	SSAP-2	<i>Corynebacterium striatum</i> 963 CAUR
Sap8	1 ± 2	23 ± 15	2.70 ± 0.43	16.50 ± 2.83	1.64 x 10 ⁻⁴	8.52 x 10 ⁻⁴	SSAP-2	<i>Corynebacterium variabile</i> DSM 44702
SmrecT	1 ± 1	20 ± 16	1.19 ± 0.30	22.07 ± 0.98	5.40 x 10 ⁻⁵	1.68 x 10 ⁻³	SSAP-5	<i>Escherichia coli</i> MG1655
Sap10	1 ± 1	133 ± 47	3.60 ± 0.85	0.16 ± 0.04	2.30 x 10 ⁻²	3.69 x 10 ⁻³	SSAP-3	<i>Corynebacterium sp</i> HMSC070E08
Sap11	0 ± 0	5 ± 3	8.53 ± 0.92	10.67 ± 1.67	8.00 x 10 ⁻⁴	5.86 x 10 ⁻⁵	SSAP-5	<i>Corynebacterium humireducens</i> DSM 45392

Sap12	4 ± 7	16 ± 5	9.33 ± 1.15	10.20 ± 1.40	9.15 x 10 ⁻⁴	1.71 x 10 ⁻⁴	SSAP-5	<i>Corynebacterium timonense</i> 5401744
Sap13	-	-	-	-	-	-	SSAP-5	<i>Nocardia farcinica</i> IFM 10152
Sap14	0	0	8.40 ± 0.40	9.50 ± 3.25	8.84 x 10 ⁻⁴	0	SSAP-5	<i>Proteus mirabilis</i> HI4320
Sap15	0	0	9.53 ± 0.50	9.10 ± 1.74	1.05 x 10 ⁻³	0	SSAP-5	<i>Aneurinibacillus aneurinilyticus</i> ATCC 12856
Sap16	0	0	6.67 ± 1.15	8.50 ± 1.56	7.84 x 10 ⁻⁴	0	SSAP-5	<i>Thermosipho melanesiensis</i> BI429

430 ^aCell numbers are shown as the average ± S.E.M. Values indicated are the average of three or more experiments. Transformation efficiency is
431 the number of cells transformed with plasmid of the viable cells. Recombination efficiency is the recombinants of the successfully transformed
432 cells (cells transformed with plasmid). Sap0 expressing empty vector was used as the negative control. SmrecT expressing RecT was used as
433 positive control. Sap1–Sap8 were tested to verify the function of SSAPs. Sap10–Sap16 were additionally tested to study the relation of SSAP
434 function and local genome context. ssDNA-: no ssDNA control. ssDNA+: addition of 2 µg oligo substrate. plasmid: addition of 500 ng plasmid
435 substrate. -: not available.

436

437

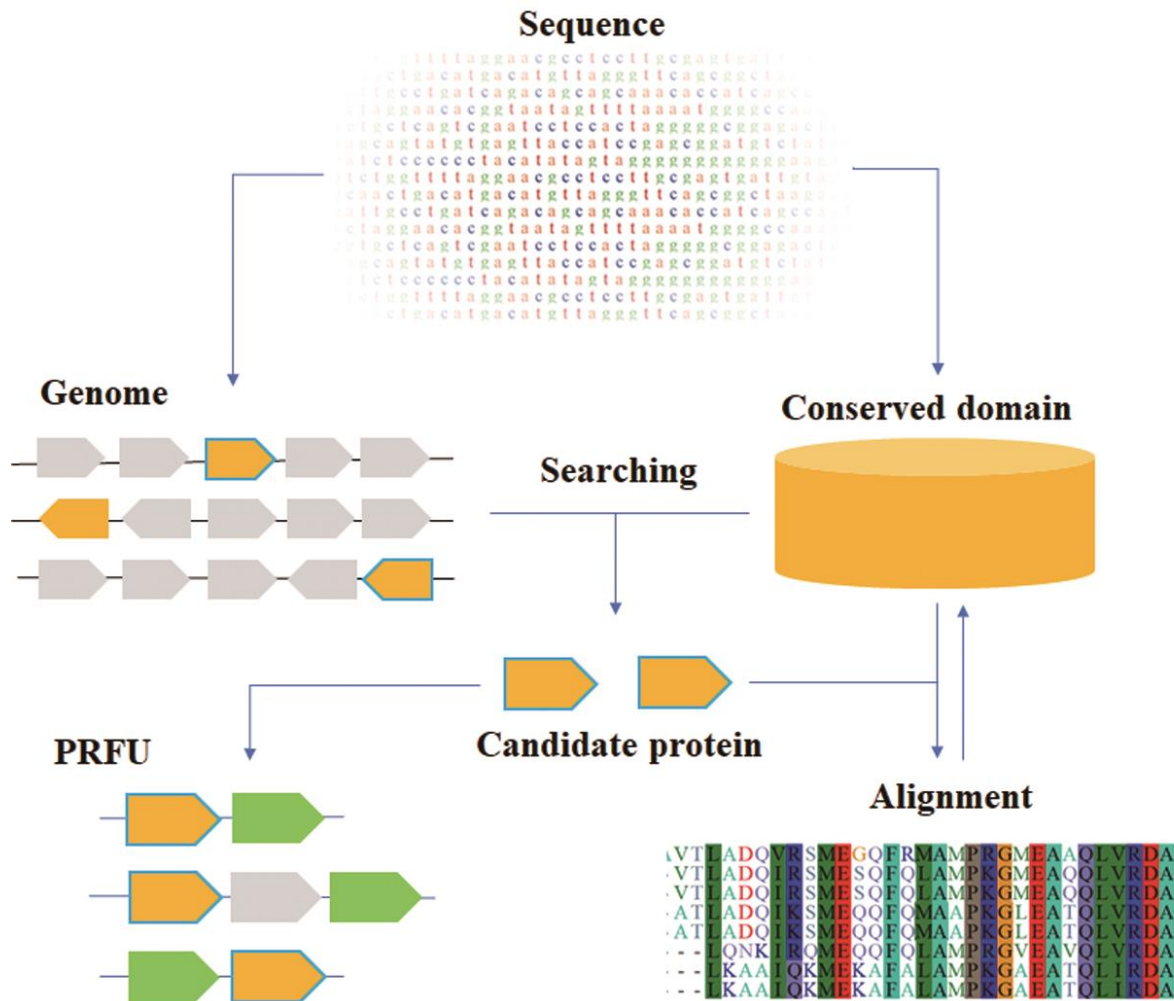
438

439

Table 3. Recombineering with dsDNA using different PRFUs

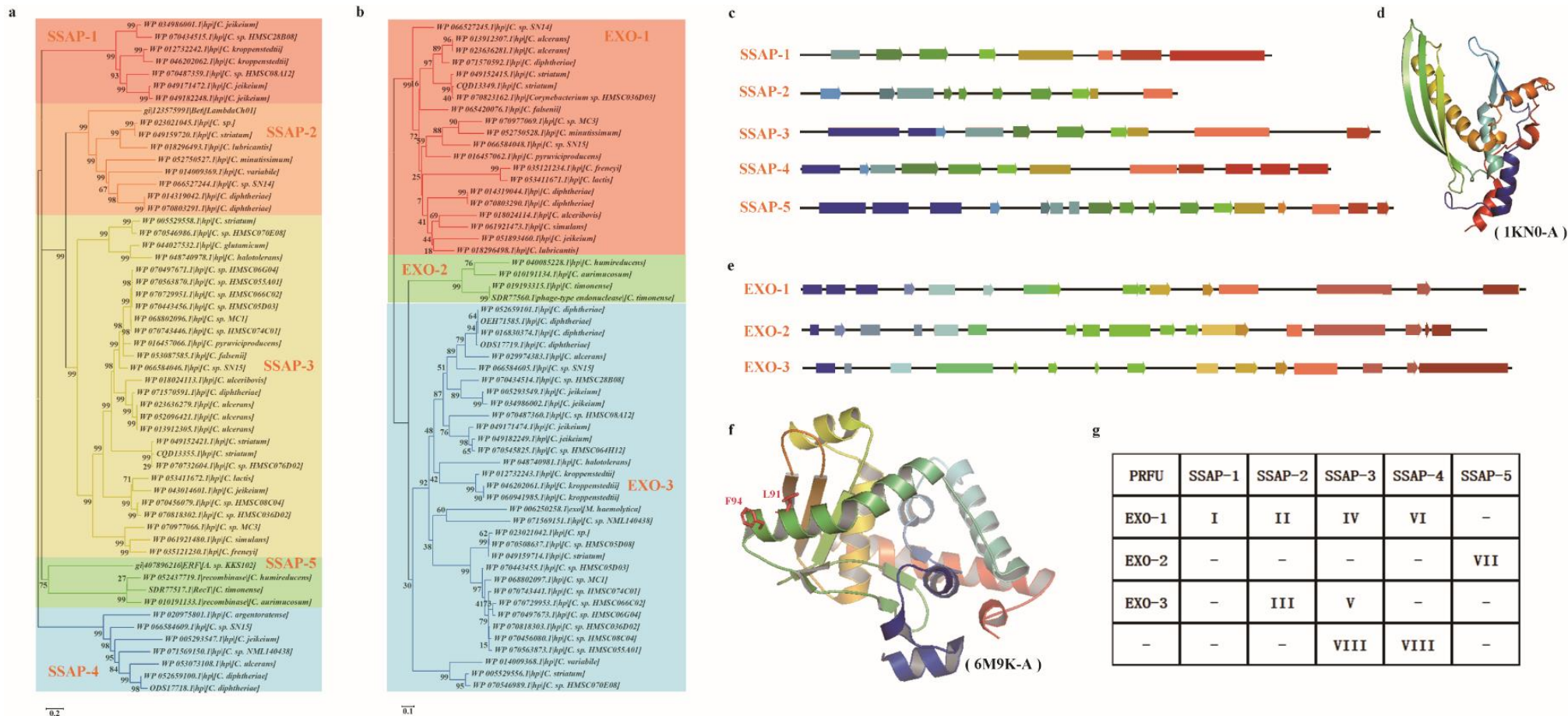
Strain	Substrate ^a			Viable cells (x 10 ⁵) ^a	Transformation efficiency	Recombination frequency	PRFU type	PRFU source
	dsDNA-	dsDNA+	plasmid (x 10 ⁵)					
ET1	0	2365 ± 822	1.52 ± 0.09	3.73 ± 1.44	4.08 x 10 ⁻¹	1.56 x 10 ⁻²	VII	<i>Corynebacterium aurimucosum</i> DSM 44827
ET2	0	4 ± 2	3.11 ± 2.07	6.10 ± 1.87	5.1 x 10 ⁻¹	1.29 x 10 ⁻⁵	VI	<i>Corynebacterium diphtheriae</i> HC07
ET3	0	0	1.53 ± 0.50	8.5 ± 1.97	1.8 x 10 ⁻¹	0	VIII	<i>Corynebacterium glutamicum</i> R
ET4	0	85 ± 73	0.66 ± 0.23	11.43 ± 8.44	5.73 x 10 ⁻¹	1.3 x 10 ⁻³	IV	<i>Corynebacterium halotolerans</i> 307 CHAL
ET8	-	-	-	-	-	-	III	<i>Corynebacterium variabile</i> DSM 44702
ET9	1	2456 ± 391	1.02 ± 0.17	6.03 ± 1.42	1.69 x 10 ⁻¹	2.41 x 10 ⁻²	-	<i>Escherichia coli</i> MG1655

440 ^a Cell numbers are shown as the average ± S.E.M. Values indicated are the average of three or more experiments. Transformation efficiency
441 was calculated as the number of cells transformed with plasmid of the viable cells. Recombination efficiency was calculated as the
442 recombinants of the successfully transformed cells (cells transformed with plasmid). ET9 expressing RecET was used as the positive control.
443 dsDNA-: no dsDNA control. dsDNA+: addition of 1 µg dsDNA substrate. plasmid: addition of 500 ng plasmid substrate. -: not available.

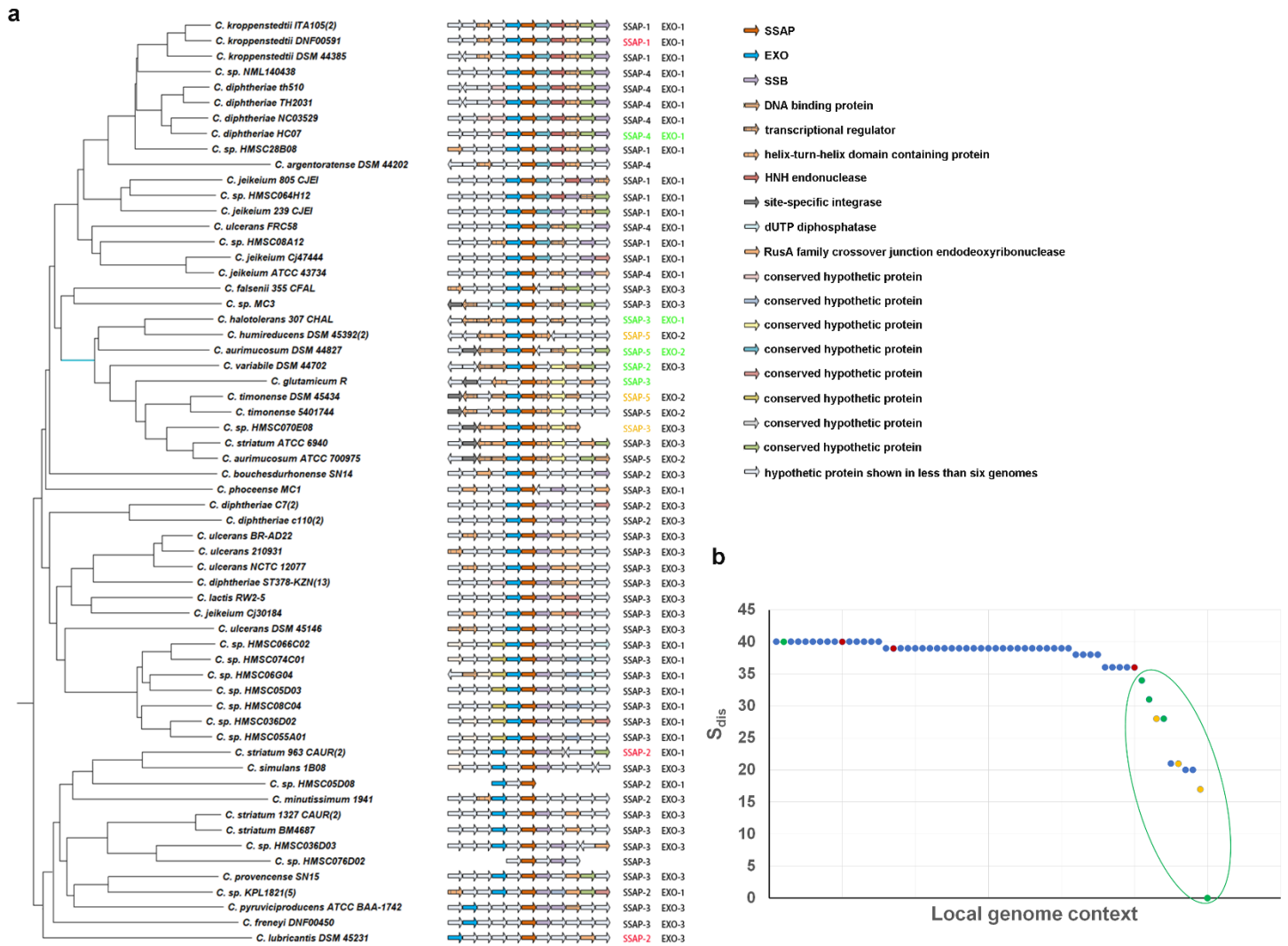


445

446 **Fig. 1.** Workflow of finding phage recombinase function units (PRFUs). The genome protein
447 sequence and the conserved domain database sequence were first downloaded from NCBI
448 (<https://www.ncbi.nlm.nih.gov/>). The genome protein sequence database was then searched in the
449 conserved domain database using the PSI-BLAST method. Candidate single-stranded annealing
450 proteins (SSAPs) or EXOs were then aligned in the conserved domain database to get the
451 consensus parts and form the new conserved domain database. The searching and alignment step
452 were iterated until no new candidates were found. All the candidate SSAPs and EXOs were then
453 collected and paired as PRFUs.



455 **Fig. 2.** Classification of the PRFUs. **a** Phylogenetic tree and classification of the SSAPs identified. **b** Phylogenetic tree and classification
 456 of the EXOs identified. Each type is shown in a different color. The labels are shown as the accession number of the protein, the note
 457 for the protein, and the source of the protein. Phylogenetic trees were constructed by MEGA6 using the Neighbor-Joining method. **c**
 458 Schematic second structure of each SSAP type. **d** N-terminal structure of human SSAP Rad52. PDB: 1KN0. **e** Schematic second structure
 459 of each EXO type. **f** Structure of λ Exo. PDB: 6M9K. **g** Theoretical PRFU type and observed PRFU type. -: not detected.



460

461 **Fig. 3.** Local genome context of the PRFUs. **a** A phylogenetic tree of local genome contexts of PRFUs. **b**
 462 The local genome context distance (S_{dis}) of different species with PRFUs compared to that of *C.*
 463 *glutamicum*. PRFUs functioning in the recombining experiment were colored in green. Those that
 464 showed no function were colored in red. Additional verified functional SSAPs were colored in orange.