

1 **DciA helicase operators exhibit diversity across bacterial phyla**

2 Helen C. Blaine^{1*}, Joseph T. Burke^{2*}, Janani Ravi^{2#}, and Christina L. Stallings^{1#}

3 ¹Department of Molecular Microbiology, Washington University School of Medicine, Saint Louis,
4 Missouri 63110, USA.

5 ²Departments of Pathobiology and Diagnostic Investigation, Microbiology and Molecular
6 Genetics, Michigan State University, East Lansing, Michigan 48824, USA.

7 *Contributed equally.

8 #Correspondence to stallings@wustl.edu and janani@msu.edu.

9

10 **ABSTRACT**

11 A fundamental requirement for life is replication of an organism's DNA. Studies in *Escherichia coli*
12 and *Bacillus subtilis* have set the paradigm for how DNA replication occurs in bacteria. During
13 replication initiation in *E. coli* and *B. subtilis*, the replicative helicase is loaded onto the DNA at the
14 origin of replication by an ATPase helicase loader. However, most bacteria do not encode
15 homologs to the helicase loaders in *E. coli* and *B. subtilis*, raising the question of how helicase
16 activity is facilitated in other bacteria during DNA replication initiation. Recent work has identified
17 the DciA protein as a predicted helicase operator that may perform a function analogous to the
18 helicase loaders in *E. coli* and *B. subtilis*. DciA proteins are defined by the presence of a DUF721
19 domain and are conserved in most bacteria. However, we find that the sequence conservation
20 between DciA proteins across different phyla is very low. Therefore, to comprehensively define
21 the DciA protein family, we took a computational evolutionary approach. These analyses identified
22 diversity in sequence features and domain architectures amongst DciA homologs that were
23 associated with specific phylogenetic lineages. The diversity of DciA proteins elucidated here
24 represents the evolution of helicase operation in bacterial DNA replication, highlights the need for

25 phyla-specific analyses of this fundamental biological process, and is an important example of
26 how research in bacterial DNA replication is necessary in organisms beyond *E. coli* and *B. subtilis*.

27

28 **IMPORTANCE**

29 Despite the fundamental importance of DNA replication for life, this process remains understudied
30 in bacteria outside of *Escherichia coli* and *Bacillus subtilis*. In particular, most bacteria do not
31 encode the helicase loading proteins that are essential in *E. coli* and *B. subtilis* for DNA replication.
32 Instead, most bacteria encode a DciA homolog that likely constitutes the predominant mechanism
33 of helicase operation in bacteria. However, it is still unknown how DciA structure and function
34 compares across diverse phyla that encode DciA proteins. In this study, we perform a
35 computational evolutionary analysis that uncovers tremendous diversity amongst DciA homologs.
36 These studies provide a significant advance in our understanding regarding an essential
37 component of the bacterial DNA replication machinery.

38

39

40 INTRODUCTION

41 DNA replication is a process critical to life for all organisms. The current paradigm for the process
42 of DNA replication in bacteria has primarily been based on studies in *Escherichia coli* and *Bacillus*
43 *subtilis*. Bacterial DNA replication begins with the binding of the replication initiation protein DnaA
44 to specific sequences referred to as DnaA boxes at the origin of replication (*oriC*) (1–7). DnaA
45 binding to double-stranded DNA (dsDNA) triggers DNA unwinding at an AT-rich region of DNA
46 called the DNA unwinding element (DUE), leaving a bubble of single-stranded DNA (ssDNA) (2,
47 3, 6, 8, 9). The ssDNA bubble is coated by single-stranded binding protein (SSB) (10), followed
48 by the concerted loading of two hexameric replicative helicases onto the SSB-coated replication
49 fork. The two helicases translocate along the two sides of the replication fork, unwinding the
50 dsDNA as they move (1, 3, 6, 11–14).

51 Bacterial replicative helicases (DnaB in *E. coli* and DnaC in *B. subtilis*) are Superfamily IV
52 type helicases, which are defined as hexameric RecA ATPases (3, 15, 16) that translocate in the
53 5'-3' direction (11, 12, 17). The bacterial replicative helicase translocates on ssDNA using a
54 “hand-over-hand” mechanism, which is driven by nucleotide hydrolysis (18, 19) (reviewed in
55 (17)). The C-terminus of the bacterial replicative helicase contains the RecA-like fold that is
56 responsible for the ATPase activity, and is connected to an N-terminal scaffolding domain via a
57 linker region (6, 20, 21). The replicative helicase must oligomerize into a double-layered
58 hexameric ring to be active during replication, with one layer made up of the N-termini and the
59 other layer comprised of the C-termini (6, 22, 23). In *E. coli* and *B. subtilis*, the loading of the
60 replicative helicase is performed with the help of a helicase loader, termed DnaC in *E. coli* and
61 DnaI in *B. subtilis* (3, 24–27). *dnaC* and *dnaI* were acquired by *E. coli* and *B. subtilis*, respectively,
62 via domestication of related but distinct phage ATPase-containing genes (28). DnaC and DnaI
63 are both in the ATPases Associated with diverse cellular Activities (AAA+) ATPase family, and

64 the ATPase activity of DnaC is required for its helicase loading function at the origin of replication
65 (29, 30).

66 *E. coli* and *B. subtilis* have long represented the paradigm of helicase loading during
67 bacterial replication. However, the majority of bacteria do not encode ATPase helicase loader
68 homologs to DnaC or DnaI. Instead, most bacteria encode the ancestral protein, DciA (DnaC/I
69 Antecedent) (28, 31), which is defined by the presence of a Domain of Unknown Function (DUF)
70 721. Despite the prevalence of DUF721-containing DciA homologs in bacteria (28), DciA has only
71 been studied in actinobacterial (*Mycobacterium tuberculosis* and *Mycobacterium smegmatis*) and
72 gammaproteobacterial (*Pseudomonas aeruginosa* and *Vibrio cholerae*) species (28, 31–33). DciA
73 homologs interact with the replicative helicase DnaB and are essential for *M. tuberculosis*, *M.*
74 *smegmatis*, and *P. aeruginosa* DNA replication and viability (28, 31). Based on DciA's interaction
75 with the replicative helicase and requirement for DNA replication, DciA has been proposed to
76 perform a function analogous to that of the DnaC/I helicase loaders. However, DciA does not
77 have a predicted ATPase domain and, therefore, cannot be considered a helicase loader like
78 DnaC/I. Instead, DciA is referred to as a predicted helicase operator, although the mechanism of
79 DciA helicase operation is still unknown (28, 31).

80 Beyond the presence of the DUF721, DciA domain architecture and the relationship
81 between DciA homologs across diverse bacterial phyla are yet to be investigated. To address
82 these open questions, we took a computational evolutionary approach and analyzed the
83 phylogenetic distribution, domain architecture, and sequence conservation for DciA homologs. We
84 have discovered low sequence similarity between DciA homologs from different phyla, lineage-
85 specific domain architectures, and divergent evolution of specific DciA homologs, all of which
86 likely have functional consequences. This study provides an evolutionary picture of DciA, reveals
87 key differences between homologs, and generates the framework for mechanistic investigation
88 into different classes of DciA proteins.

89 RESULTS

90 DciA proteins vary considerably in sequence across bacterial phyla

91 Most bacteria encode a DciA homolog, defined by the presence of the DUF721 domain (28).
92 Based on studies in *M. tuberculosis*, the DUF721 is predicted to contain a region of structural
93 homology to the N-terminus of DnaA, which was thus named the DnaA N-terminal-like (DANL)
94 domain (31). The presence of the DANL domain has subsequently been confirmed in *V. cholerae*
95 DciA (32). Our analysis of other DciA homologs indicates that this predicted structural domain is
96 conserved, suggesting that it is important for DciA function. However, beyond the annotation of
97 the DUF721 (henceforth referred to as the DciA domain), it is unclear how different DciA homologs
98 relate to each other. A protein BLAST search for *M. tuberculosis* DciA homologs based on primary
99 amino acid sequence only identifies closely related homologs, all of which are in actinobacteria
100 (**Figure S1**). We found a similar pattern when retrieving homologs with *Pseudomonas aeruginosa*,
101 all homologs are proteobacterial (**Figure S1**). This suggests that the DciA homologs in different
102 phyla are diverse with low conservation in their primary amino acid sequence. We, therefore,
103 needed a more extensive method to investigate the relationship between DciA homologs across
104 different phyla. To achieve this, we used the MolEvolvR web application to comprehensively
105 identify and characterize DciA homologs across all bacterial lineages using molecular evolution
106 and phylogeny (34). Since individual DciA proteins from specific lineages were not successful in
107 retrieving homologs from other distant phyla, we selected a much wider range of starting points.
108 We started with 21 DciA proteins from 11 diverse phyla (28), including representatives from
109 actinobacteria, proteobacteria, and cyanobacteria (**Table S1; Figure 1A**), as query sequences to
110 identify diverse DciA homologs across additional bacterial phyla (**Figure 1A**). Our homology
111 search resulted in identifying ~9K DciA homologs from 15 bacterial phyla (**Figure 1B**). In line with

112 both genome sequencing and publication bias, proteobacteria and actinobacteria were over-
113 represented in both our queries and recovered sequences (35–37) (**Figure 1A,B**).

114 No single DciA protein identified homologs in all other phyla (**Figure 1C**), supporting that
115 there is low sequence conservation between DciA homologs. To quantify the sequence
116 conservation across phyla, we analyzed the pairwise similarity for the 21 DciA protein homologs
117 used as our query set (**Figure 2**). We found a wide range of similarities (~20–60%), with the
118 majority of homologs showing 30-40% similarity. Our query sequence dataset contained multiple
119 species within each class of proteobacteria, so we were able to compare the conservation
120 between DciA homologs within proteobacteria as well as between proteobacteria and other phyla.
121 Similarity was high between DciA proteins within alphaproteobacteria (62.3% between *Brucella*
122 *abortus* and *Mesorhizobium australicum*) and gammaproteobacteria (52.3% between *Proteus*
123 *mirabilis* and *Vibrio cholerae*). Overall, DciA protein queries within proteobacteria have an
124 average of ~30% similarity, and proteobacterial homologs have an average of 31.12% similarity
125 to homologs outside of their phyla (**Figure 2**). Actinobacterial and proteobacterial DciA homologs
126 share between 21–34% sequence similarity. For example, the *M. tuberculosis* DciA protein shares
127 28.3% identity with the *Pseudomonas aeruginosa* DciA (**Figure 2**). DciA shares this trait of low
128 sequence conservation across phyla with its interaction partner, the replicative helicase DnaB,
129 where DnaB proteins in *M. tuberculosis* and *P. aeruginosa* only share 20.4% similarity. Bacterial
130 replication proteins in general have low to moderate sequence conservation across phyla (11–
131 49% similarity across replisome proteins between *E. coli* and *B. subtilis*) (38). Therefore, our data
132 showing low DciA sequence conservation is consistent with other replication initiation proteins.
133 Given the low-level conservation across divergent DciA proteins, it was unsurprising that
134 individual DciA query proteins never returned homologs from all other phyla and emphasizes the
135 need for multiple starting points for analysis (**Figure 1C**).

136 While we identified a diverse set of DciA homologs, including those with moderate to low
137 sequence conservation from most lineages (**Figure 1**), we were still missing multiple phyla
138 previously reported to contain DciA homologs. Therefore, we expanded our search further by
139 including 66 DciA query proteins from 20 bacterial phyla as query sequences (**Figure 3A; Table**
140 **1**). This hugely diversified our results identifying >13K unique DciA homologs from 22 bacterial
141 phyla (**Figure 3B,C**). Consistent with the previous smaller set of DciA query proteins (**Figure 1**),
142 we found that most DciA proteins identified homologs within their corresponding phylum as well
143 as ones within actinobacteria and proteobacteria (**Figure 3C**). However, a few DciA proteins
144 recovered homologs only within their own phylum. Specifically, most actinobacterial (10/12) and
145 a few proteobacterial (10/30) DciA proteins only recovered homologs within their respective
146 phylum (**Figure 1C, 3C**). In addition, the DciA proteins from nitrospirae, gemmatimonadetes,
147 fibrobacteres, chlorobi, chlamydiae, and bacteroidetes identified homologs from actinobacteria
148 and not proteobacteria, suggesting that these DciA proteins are closer evolutionarily to
149 actinobacterial DciA homologs than proteobacterial DciA. In contrast, acidobacteria,
150 cyanobacteria, and thermodesulfobacteria recovered homologs from proteobacteria and not
151 actinobacteria. The evolution of DciA seems to mirror the phylogenetic distance of these species;
152 gemmatimonadetes, chlamydiae, and bacteroidetes are more closely evolutionary related to
153 actinobacteria than proteobacteria based on the 16S rRNA gene (28, 39). Together these data
154 suggest that the DciA homologs in actinobacteria and proteobacteria likely represent the most
155 evolutionarily divergent repertoire of DciA homologs. This analysis also indicates that DciA
156 proteins from distinct phyla carry lineage-specific signatures, likely co-evolving with other phylum-
157 specific protein families involved in DNA replication.

158

159 **DciA domain architecture varies in a lineage-specific manner**

160 The evolutionary divergence in DciA proteins prompted us to take a closer look at the sequence-
161 structure features of these homologs, including sequence alignment and domain architectures.
162 As a first step, we aligned the 66 DciA starting point protein sequences from 20 phyla (**Figure 3**;
163 **Table 1**) using the MolEvolvR web application (34) and examined their domain architectures
164 (**Figure 4**; *see Methods*). Tracing the phylogenetic tree of DciA sequences confirms that the
165 evolution of DciA proteins roughly corresponds to bacterial phylogeny, where DciA homologs tend
166 to cluster with their own phyla (*e.g.*, proteobacteria and actinobacteria) (**Figure 4 center**). Only
167 DciA homologs that encoded the signature DUF721 domain were included in our alignment, and
168 could be classified into one of four distinct groups based on where the DUF721 domain occurs
169 within the protein (**Figures 4, 5, and S2**). We describe the group memberships and descriptions
170 of our 66 DciA query proteins in the following sections (**Table 1**).

171 In Group 1 DciA proteins, the DUF721 spans at least 70% of the protein sequence, with
172 ≤ 25 amino acids on either side of the DUF721 (**Figures 4, 5A, S2; Table 1**, See Group 1 example:
173 *R. rickettsii* DciA). Group 1 DciA proteins are present in acidobacteria, bacteroidetes, chlamydiae,
174 chlorobi, gemmatimonadetes, elusimicrobia, fibrobacteres, fusobacteria, thermotogae,
175 planctomycetes, spirochaetes, and proteobacteria (**Table 1**). The DUF721 domain spanning the
176 entire protein in this group suggests that the DUF721 is likely sufficient for DciA function in these
177 bacteria.

178 Group 2 DciA homologs have ≤ 25 amino acids C-terminal to the DUF721 and an N-
179 terminal extension of > 25 amino acids and no known domains (**Figures 4, 5A, S2; Table 1**, See
180 Group 2 example: *M. tuberculosis* DciA). Group 2 DciA homologs are present in actinobacteria
181 and verrucomicrobia (**Figure 5A; Table 1**). Group 2 DciA queries in the mycobacteriales order
182 are predicted to have intrinsically disordered regions by MobiDBLite software (also noted in (40))
183 in the N-terminal extension (**Figure 4, left**). This predicted intrinsically disordered region is not

184 present in bifidobacteriales, another order within actinobacteria, or in other Group 2 proteins.
185 Expression of the *M. tuberculosis* DciA DUF721 alone is not sufficient to support viability in
186 mycobacteria (31), suggesting that the N-terminal extension is essential for DciA function in
187 mycobacteria, although its function remains unknown. The requirement for the sequence N-
188 terminal to the DUF721 in mycobacteria and the absence of this N-terminal sequence in Group 1
189 and 3 DciA proteins demonstrates divergent evolution of DciA homologs in bacteria and the
190 potential for functional variation. In addition to the N-terminal extension in all Group 2 DciA
191 homologs, 4 actinobacterial proteins also encode a short predicted region of disorder C-terminal
192 to the DUF721 (19aa in *M. tuberculosis*; **Figure 4, left, Table 1**).

193 Group 3 DciA homologs have ≤ 25 amino acids N-terminal to the DUF721 and a C-terminal
194 extension >25 amino acids long (**Figures 4, 5A, S2; Table 1**, See Group 3 example: *V. cholerae*
195 DciA). Group 3 DciA homologs are present in acidobacteria, spirochaetes, deferribacteres,
196 cyanobacteria, dictyoglomi, nitrospirae, thermodesulfobacteria, and proteobacteria. The DciA
197 homologs in *T. palladium*, *Nitrospirae moscoviensis*, and *V. parahemolyticus* are the only Group
198 3 proteins that are predicted to be intrinsically disordered in the region C-terminal to the DUF721
199 domain using the MobiDBLite predictor (**Figure 4, left; Table 1**). However, the C-terminal
200 sequence of *V. cholerae* DciA was predicted to be intrinsically disordered in prior studies using
201 PONDR and IUPred2A software, which was further confirmed by small-angle X-ray scattering
202 (32). When we analyzed the *V. cholerae* sequence using the JRONN disorder prediction track,
203 we were able to identify an intrinsically disordered region C-terminal to the DUF721 (**Figure 4**).
204 Therefore, it is possible that other Group 3 proteins could have intrinsically disordered domains
205 that are not predicted by the tools we are using. The C-terminal intrinsically disordered sequence
206 of *V. cholerae* DciA is required for its interaction with DnaB as well as for the enhancement of the
207 association between *V. cholerae* DnaB and ssDNA *in vitro* (32). The absence of the C-terminal
208 extension in Groups 1 and 2 DciA proteins further supports evolutionary divergence in DciA

209 protein sequence and structure, which likely impacts specialized lineage-specific protein function
210 or mechanism of interaction with the replicative helicase. The one actinobacterial DciA protein
211 from our query set that falls into Group 3 is encoded by *S. seoulensis*. This protein encodes a
212 YspA domain (Pfam: YAcAr/PF10686) C-terminal to the DUF721 (**Figure 4, 5A, S2; Table 1**).
213 YspA domains within proteins typically have fusions to domains that process nucleotide-derived
214 ligands such as ADP-ribose and may function as sensors of these ligands and nucleic acids (41).

215 In Group 4 DciA proteins, the DUF721 domain falls > 25 amino acids away from both
216 termini (**Figure 4, 5A, S2; Table 1**, See Group 4 example: *C. mirabilis* DciA). In our query set,
217 Group 4 consists of DciA proteins in proteobacteria, actinobacteria, and synergistetes. The Group
218 4 DciA homologs in *F. fastidiosum*, *S. coelicor*, and *S. avermitilis* have regions of disorder both
219 N-terminal and C-terminal to the DUF721 (**Figure 4, left**).

220 Analysis of how the DciA groups are distributed in different bacterial phyla for our starting
221 set of DciA homologs (**Figure 4**) reveals that 75% of homologs we queried from actinobacteria
222 fall into Group 2, and 66% of proteobacterial DciA proteins from our starting set fall into Group 3.
223 (**Figure 5A; Table 1**). These two phyla had the most representatives in our set of query proteins,
224 and they form distinct clusters on the phylogenetic tree (**Figure 4, right**). Further examining the
225 characteristics of each group revealed that 75% of DciA proteins from gram-positive bacteria
226 included in the query fall into Group 2, and 54% of DciA starting points from gram-negative
227 bacteria fall into Group 3 (**Figure 5B**). By contrast, no Group 1 DciA proteins are from gram-
228 positive bacterial queries, and Group 1 makes up 37% of gram-negative bacterial queries (**Figure**
229 **5B**). In addition to classifying each DciA protein from our query set into one of four groups based
230 on the position of the DUF721, we found that the alphaproteobacteria DciA proteins in the
231 Rickettsiales and Hyphomicrobiales orders harbored an insertion within the DUF721 that is not
232 present in any other phyla (**Figure 4, right; Figure S2**). This could indicate further functional
233 divergence of the DciA homologs in these orders.

234 Overall, these analyses reveal the diversity of domain architectures amongst our selected
235 66 DciA homologs. A more comprehensive analysis of all DciA homologs would be required to
236 fully understand the distribution of diverse domain architectures in different phyla. Nonetheless,
237 these data demonstrate that DciA homologs have diverged significantly in sequence structure
238 based on the position and sequence of the DUF721, which may impact their function and
239 interaction with the DNA replication machinery.

240

241 **The DciA domain proximity network**

242 To further investigate the possible functions of each of the ~13K putative DciA homologs that we
243 have identified using our diverse query set, we interrogated their domain architectures (*see*
244 *Methods*). A striking majority of these homologs showed no variation, carrying a single DciA
245 domain (the Pfam DUF721 domain) (**Figure 6**). Less than 1% of the proteins exhibited novel
246 fusions with the DciA domain. For example, the DciA protein from *S. seoulensis* identified other
247 YAcAr/PF10686 Pfam members in Streptomyces, as noted above in **Figure 4**. These
248 actinobacterial homologs carry the YspA/YAcAr-like domain known to be associated with NAD
249 utilization and ADP-ribosylation domains (41) (**Figure 6**). We find only a few streptomyces DciA
250 homologs that share this domain architecture (**Figure 4, S3**), suggesting that DciA in this genus
251 might have evolved this unique function. In addition, some proteobacteria also showed variation
252 in domains associated with the DUF721 domain: i) Reyranela species carry a C-terminal
253 thioredoxin domain, with possible redox function, and ii) the pseudomonas genus has a rare
254 instance of a DciA dyad (two DciA domains) (**Figure 6**). We also note that there are ABC
255 transporter-like proteins within proteobacteria with ~30% similarity to the query DciA protein in
256 acidobacteria, but none of the other 65 query proteins (**Figure 6**). Finally, we found that one query
257 DciA from deferribacteres identifies peptidase-like proteins in proteobacteria, again with no
258 similarity to any of the other divergent DciA proteins. The novel fusions and alternate domain

259 architectures identified within the bacterial DciA homologs have been summarized in the form of
260 a network of domain architectures reconciling all DciA homologs, with domains as nodes and co-
261 occurrence within proteins as edges, and their co-occurrences have been further quantified
262 (**Figure 6**). The association of the DUF721 and DciA proteins with other functional domains could
263 shed light on the activities for DciA in these bacteria.

264

265

266 **DISCUSSION**

267 The recent discovery of DciA as a predicted helicase operator in bacteria (28, 31) has begun to
268 shed light on a long-standing open question of how the majority of bacteria facilitate helicase
269 activity during DNA replication in the absence of the ATPase helicase loaders expressed by *E.*
270 *coli* and *B. subtilis*. The wide distribution of DciA in diverse bacterial phyla indicates that these
271 proteins likely represent the predominant paradigm for helicase operation in bacteria, despite not
272 being conserved in *E. coli* and *B. subtilis*, the organisms typically used as a model for bacterial
273 replication. DciA proteins are defined by the presence of the DUF721 domain and prior
274 phylogenetic analysis indicates that *dnaC* and *dnaI* homologs were acquired through evolution at
275 the expense of *dciA* (named for *dna[CI]* antecedent) (28), suggesting that DciA and DnaC/DnaI
276 perform a common function. In addition, it has been shown that DciA interacts with the replicative
277 helicase and is required for DNA replication and viability in the limited organisms it has been
278 studied in (28, 31, 33). However, the mechanism by which DciA mediates replication initiation is
279 still unknown. Our study has revealed immense diversity in DciA proteins, where there is low
280 sequence similarity between homologs in different phyla (**Figures 1–3**), there are at least 4
281 distinct classes based on the positioning of the DUF721 domain in the protein (**Figures 4–5**),
282 there exist lineage-specific insertion sequences in the DUF721 domain of some proteobacterial
283 species (**Figures 4, S2**), and some DciA proteins have evolved as fusions to other functional

284 domains (**Figure 6**). These data suggest that DciA proteins have been divergently evolving and
285 the mechanism of helicase operation conferred by DciA may have distinct features dependent on
286 the bacteria.

287 Biochemical and genetic studies have only been performed with *M. tuberculosis*, *P.*
288 *aeruginosa*, and *V. cholerae* DciA proteins (28, 31–33). Our analyses demonstrate that the DciA
289 homologs from these species have a number of distinct features, in addition to low sequence
290 similarity, raising the question of how conserved their mechanisms of action will be. In particular,
291 mycobacterial DciA is a Group 2 DciA protein with sequences predicted to be intrinsically
292 disordered both N-terminal and C-terminal to the DUF721 (**Figures 4 and 5**). In contrast, *V.*
293 *cholerae* and *P. aeruginosa* DciA proteins are classed as Group 3, with a long sequence extension
294 C-terminal of the DUF721 (**Figure S2**).

295 The one feature conserved in all DciA homologs is the presence of the DUF721, which
296 contains the DANL domain and is predicted to structurally resemble the N-terminus of DnaA (31).
297 The N-terminus of DnaA is critical for the interaction of DnaA with the helicase and other
298 regulators (42, 43), however, the role of the DciA DANL domain in the interaction with DnaB has
299 yet to be established. A tryptophan residue conserved in the DANL domains of many DciA
300 homologs has structural similarity to a phenylalanine residue in the DnaA N-terminus that has
301 been predicted to have a key role in making contacts between DnaA and its interacting partners,
302 including DnaB (31, 44). Mutation of the conserved tryptophan in the DANL domain of *M.*
303 *tuberculosis* DciA results in slow growth and decreased DNA replication (31). This supports that
304 the conserved tryptophan within the DANL domain plays a key role for DciA function *in vivo*,
305 however that precise role has yet to be elucidated. It is also important to note that not all DciA
306 homologs encode this tryptophan residue within their DANL domain (32) (**Figure S2**). In fact,
307 there is considerable diversity in the DUF721 sequences between DciA homologs from different
308 phyla, including an insertion in the DUF721 of some alphaproteobacteria DciA proteins, which is

309 not observed in the DciA proteins analyzed here from other classes (**Figure 4, Figure S2**).
310 Therefore, even the defining DUF721 feature of DciA proteins has evolved, likely reflecting either
311 lineage-specific adaptation in mechanism of action or mode of interaction with the replicative
312 helicase.

313 There are multiple DciA homologs predicted to encode intrinsically disordered regions N-
314 terminal and/or C-terminal to the DUF721 (**Figure 4**). The intrinsically disordered sequence C-
315 terminal to the DUF721 in the *V. cholerae* Group 3 DciA protein enhances the association
316 between DnaB and ssDNA and truncation of this intrinsically disordered sequence results in loss
317 of the interaction between *V. cholerae* DciA and the DnaB helicase (32). Although it is unknown
318 how much of this mechanism will be conserved in other bacterial species encoding DciA proteins
319 with divergent domain architectures and many DciA proteins do not encode predicted intrinsically
320 disordered regions, there is precedent for roles of intrinsically disordered domains in other
321 bacterial DNA replication proteins. For example, the intrinsically disordered linker (IDL) within the
322 C-terminus of SSB is important for its cooperative ssDNA binding, as well as the displacement of
323 SSB from ssDNA (45, 46) (reviewed in (47)). The IDL has also been proposed to be important for
324 SSB protein-protein interactions, such as the interaction between SSB and the DNA repair protein
325 RecG (48). The intrinsically disordered C-terminus of the replication restart helicase Rep is also
326 important for the interaction between Rep and its regulator PriC, as well as between Rep and the
327 replicative helicase DnaB (49, 50). In addition, the helicase loaders DnaC and DnaI as well as the
328 replication initiation protein DnaA have been predicted to encode intrinsically disordered domains,
329 currently of unknown function (32).

330 A lot of unknowns still remain regarding DciA proteins and bacterial DNA replication. The
331 computational evolutionary analysis described herein highlights the complexities and diversity
332 that have evolved in the fundamental process of DNA replication, where no single species of
333 bacteria will be able to represent a central dogma that holds true throughout the Kingdom. These

334 studies provide a framework for researchers to consider the evolutionary variation while dissecting
335 the mechanistic basis for helicase operation in bacteria.

336

337 **METHODS**

338 **Query selection**

339 We selected our small and full set of DciA query proteins based on the DUF721 location defined
340 by on Pfam annotation, using a variety of DciA containing phyla with annotated DciA sequences.
341 The DciA domain was annotated using Pfam annotation and subsequently confirmed using a
342 multiple sequence alignment (**Figure S2**). The full list of starting points is listed in **Table 1**. The
343 only DciA-containing phyla excluded from our set of 66 query proteins were Deinococcus-
344 Thermus, Chrysiogenetes, and Firmicutes. Firmicutes and Deinococcus-Thermus were
345 subsequently recovered in our MolEvolvR searches.

346

347 **Analysis using MolEvolvR**

348 We used MolEvolvR (34) to determine and characterize all DciA query proteins and their
349 homologs across the bacterial kingdom. We first identified all the homologs for each of the query
350 proteins in RefSeq (51) genomes, and reconciled the comprehensive set of DciA homologous
351 proteins. Next, we characterized each of the query proteins and their homologs in terms of domain
352 architectures (including Pfam (52), Gene3D (53)), localization (using Phobius (54), TMHMM (55)),
353 and disorder (using MobiDB (56)). The domain architectures of these homologs were analyzed
354 by lineage, quantified with Upset plots, and reconciled using domain proximity networks. We then
355 performed phylogenetic analysis including phyletic spreads (sunburst, heatmap), multiple
356 sequence alignment, and tree construction using MolEvolvR and custom R scripts. The MSA for
357 subset of the sequences with representatives from the 4 DciA groups shown in Figure SY was
358 generated using Kalign (57) and Jalview (58). All our data, analyses, and visualizations

359 summarizing the DciA homologs across the bacterial kingdom along with their domain
360 architectures and phyletic spreads are available at https://github.com/JRaviLab/dcia_evolution.

361

362 **Pairwise Similarity Analysis**

363 The similarity matrix was designed using the MatGat application (59). We compared each DciA
364 query protein of the 21 starting points (**Figure 1**) to each other in order to calculate pairwise
365 percent similarities. DnaB from *M. tuberculosis* and *P. aeruginosa* sequence similarity were
366 compared using the EMBOSS Needle pairwise similarity tool (57).

367

368 **FIGURE LEGENDS**

369 **Figure 1. Query of DciA homologs using 20 DciA protein starting points reveals diversity**
370 **across bacterial phyla. A. Lineages of query DciA proteins.** Sunburst plot showing the
371 lineages of the 21 query DciA proteins. In each plot, the inner ring corresponds to the kingdom
372 (bacteria, in this case), and the outer ring represents the distribution of phyla. **B. Lineages of**
373 **DciA homologs.** Sunburst plot showing the phyletic spread of all the DciA homologs generated
374 using the 21 starting points. **C. Phyletic spread of the DciA homologs by query.** The heatmap
375 shows the presence/absence of homologs of DciA across bacterial lineages (columns) for each
376 query DciA (rows). The color gradient indicates the highest number of homologs in a particular
377 lineage. *Note: The sunburst plots only display lineages of >0.1% fraction of total proteins. The
378 heatmap gives the full picture.

379

380 **Figure 2. Pairwise similarity analysis of DciA proteins.** Pairwise percentage similarities for 21
381 query DciA proteins across 11 phyla were computed using MatGat (59) and the standard
382 BLOSUM62 matrix for similarity metric calculation.

383

384 **Figure 3. Retrieving DciA homologs using the extended query set of 66 DciA proteins. A,**
385 **B. Lineages of query and homologous DciA proteins.** See Figure 1 for details. The three phyla
386 excluded in the query searches were Deinococcus-Thermus, Chrysiogenetes, and Firmicutes. **C.**
387 **Phyletic spread of the DciA homologs by query.** See Figure 1 for details. Deinococcus-
388 Thermus and Firmicutes were both recovered in the resulting set of homologs. Our queries did
389 not recover the DciA homologs in chrysiogenetes, suggesting that the homologs in this phylum
390 are the most divergent in sequence from the query sequences.

391
392 **Figure 4. Characterizing the full list (66) of DciA query proteins.**
393 **The domain architectures and disorder predictions (left panel) of the 66 DciA query**
394 **proteins are overlaid with the multiple sequence alignment (right), and phylogenetic tree**
395 **(middle).** Each DciA protein is marked with the kingdom (B, bacteria), phylum (first 6 letters),
396 Genus, and species (represented as '*Gspecies*'). The Pfam and MobiDB annotations for each
397 domain prediction are shown in the legend (top). The colors in the multiple sequence alignment
398 depiction correspond to different amino acids (bottom legend). The data was generated and
399 visualized using the MolEvolvR web application. In addition, JronnWS (60) disorder predictions
400 were performed within Jalview (58) (not shown here), where other DciA proteins such as *V.*
401 *cholerae* show disorder regions.

402
403 **Figure 5. DciA groups and their distribution within our query sequences.**

404 **A. Example domain architectures of each of the 4 groups of DciA homologs.** Group 1 DciA
405 homologs have ≤ 25 aa on either side of the DUF721 (top, blue), Group 2 homologs have ≤ 25 aa
406 C-terminal to the DUF721 (second, pink), Group 3 DciA homologs have ≤ 25 aa amino acids N-
407 terminal to the DUF721 (third, orange), and Group 4 DciA homologs have > 25 aa both N- and C-
408 terminal to the DUF721 (bottom, teal). Graphics created using BioRender.com. **B. Distribution**

409 **of groups within Gram-positive and Gram-negative DciA query proteins.** Pie charts
410 comparing the number of gram-positive and gram-negative bacteria that have DciA homologs in
411 one of the 4 groups. Gram-positive (left), gram-negative (right). Percentages rounded to the
412 nearest whole number. Details of group and Gram stain assignments of each DciA homolog are
413 found in **Table 1**.

414
415 **Figure 6. DciA partners. A. Domain proximity network.** The network visualizes co-occurring
416 domains within all bacterial DciA homologs generated with our 66 starting points (**Figure 1**; **Table**
417 **1**). The nodes and edges correspond to domains and co-occurrence of domains within a protein;
418 the size corresponds to the frequency of occurrence with a minimum scaling factor. The full data
419 can be accessed at https://github.com/JRaviLab/dcia_evolution. **B. Frequencies of co-**
420 **occurring domains in DciA homologs.** Upset plot of the DciA homologs are shown. Blue
421 histogram: Distribution of the predominant domains. Dots and connections: Combinations in
422 which these domains come together in DciA domain architectures. Red histogram: Frequency of
423 occurrences of domain architectures. Of these only the DciA containing domain architectures
424 were used for alignments and phylogenetic trees.

425
426
427 **SUPPLEMENTARY FIGURES**

428 **Figure S1. Mycobacterium tuberculosis and Pseudomonas aeruginosa DciA proteins only**
429 **identify DciA homologs in their respective phylum.** Heatmap description as in Figure 1.

430
431 **Figure S2. Multiple sequence alignment of select DciA proteins.** Alignment of the 66 DciA
432 proteins used in the full query set. Numbering of residues across the top of the alignment is based
433 on the annotation *M. tuberculosis* DciA. Based on the numbering of *M. tuberculosis* residues, The

434 DUF721 domain falls between 75–162 amino acids (red box), the conserved tryptophan residue
435 in the DANL domain is at position 133 (red star), and the insertion present in some
436 alphaproteobacterial DciA proteins occurs after position 118 on the alignment. The multiple
437 sequence alignment was generated with Kalign (57, 61) and visualized using Jalview (58); (color
438 scheme: Clustalx).

439

440 **Figure S3. Full list representative homolog characterization (with DciA)**

441 The domain architectures, multiple sequence alignment, and phylogenetic tree were generated
442 using representative DciA homologs (one per domain architecture per lineage). See Figure 4 for
443 details.

444 **TABLES**

445 **Table 1. DciA query proteins used to identify homologs across the bacterial kingdoms.**

446 Protein, domain architecture, group, and lineage-related metadata for each of the 66 diverse
447 starting points of DciA proteins across the bacterial kingdom are shown in this table. The full
448 homolog data, analyses, and figures can be found here:
449 https://github.com/jravailab/dcia_evolution.

450

451

452 **ACKNOWLEDGEMENTS**

453 We are very grateful to the Midwest Microbial Pathogenesis Conference (MMPC) 2021 organizers
454 for providing HCB a travel award and opportunity to present the DciA story, and for providing an
455 interactive venue for JR and JTB to start this collaboration with CLS and HCB.

456

457 **FUNDING**

458 CLS is supported by a Burroughs Wellcome Fund Investigator in the Pathogenesis of Infectious
459 Disease Award. HCB is supported by the Sondra Schlesinger Student Fellowship in Molecular
460 Microbiology. JR is supported by Michigan State University (MSU) College of Veterinary Medicine
461 Endowed Research Funds and MSU start-up funds.

462

463 **DATA AVAILABILITY AND REUSE**

464 All the data, analyses, and visualizations are available in our GitHub repository,
465 https://github.com/JRaviLab/dcia_evolution. Text, figures, and data are licensed under Creative
466 Commons Attribution CC BY 4.0.

467

468 **REFERENCES**

- 469 1. Kaguni JM. 2011. Replication initiation at the Escherichia coli chromosomal origin. Current
470 opinion in chemical biology, 2011/08/18 ed. 15:606–613.
- 471 2. Mott ML, Berger JM. 2007. DNA replication initiation: mechanisms and regulation in
472 bacteria. Nature Reviews Microbiology 5:343–354.
- 473 3. Jameson KH, Wilkinson AJ. 2017. Control of Initiation of DNA Replication in Bacillus subtilis
474 and Escherichia coli. Genes 8:22–22.
- 475 4. Fuller RS, Kornberg A. 1983. Purified dnaA protein in initiation of replication at the
476 Escherichia coli chromosomal origin of replication. Proceedings of the National Academy
477 of Sciences 80:5817 LP – 5821.
- 478 5. Fuller RS, Funnell BE, Kornberg A. 1984. The dnaA protein complex with the E. coli
479 chromosomal replication origin (oriC) and other DNA sites. Cell 38:889–900.

- 480 6. Chodavarapu S, Kaguni JM. 2016. Replication Initiation in Bacteria. *Enzymes* 39:1–30.
- 481 7. Schaper S, Messer W. 1995. Interaction of the initiator protein DnaA of *Escherichia coli* with
482 its DNA target. *The Journal of biological chemistry* 270:17622–17626.
- 483 8. O'Donnell M, Langston L, Stillman B. 2013. Principles and concepts of DNA replication in
484 bacteria, archaea, and eukarya. *Cold Spring Harbor perspectives in biology* 5:a010108–
485 a010108.
- 486 9. Bramhill D, Kornberg A. 1988. Duplex opening by dnaA protein at novel sequences in
487 initiation of replication at the origin of the *E. coli* chromosome. *Cell* 52:743–755.
- 488 10. Meyer RR, Laine PS. 1990. The single-stranded DNA-binding protein of *Escherichia coli*.
489 *Microbiological reviews* 54:342–380.
- 490 11. LeBowitz JH, McMacken R. 1986. The *Escherichia coli* dnaB replication protein is a DNA
491 helicase. *The Journal of biological chemistry* 261:4738–4748.
- 492 12. Baker TA, Funnell BE, Kornberg A. 1987. Helicase action of dnaB protein during replication
493 from the *Escherichia coli* chromosomal origin in vitro. *Journal of Biological Chemistry*
494 262:6877–6885.
- 495 13. Lewis JS, Jergic S, Dixon NE. 2016. The *E. coli* DNA Replication Fork. *The Enzymes* 39:31–
496 88.
- 497 14. Oakley AJ. 2019. A structural view of bacterial DNA replication. *Protein science: a*
498 *publication of the Protein Society* 28:990–1004.
- 499 15. Gorbalenya AE, Koonin EV. 1993. Helicases: amino acid sequence comparisons and
500 structure-function relationships. *Current Opinion in Structural Biology* 3:419–429.

- 501 16. Leipe DD, Aravind L, Grishin NV, Koonin EV. 2000. The bacterial replicative helicase DnaB
502 evolved from a RecA duplication. *Genome research* 10:5–16.
- 503 17. Fernandez AJ, Berger JM. 2021. Mechanisms of hexameric helicases. *Critical Reviews in*
504 *Biochemistry and Molecular Biology* 56:621–639.
- 505 18. Spinks RR, Spenkelink LM, Stratmann SA, Xu Z-Q, Stamford NPJ, Brown SE, Dixon NE,
506 Jergic S, van Oijen AM. 2021. DnaB helicase dynamics in bacterial DNA replication resolved
507 by single-molecule studies. *Nucleic acids research* 49:6804–6816.
- 508 19. Itsathitphaisarn O, Wing RA, Eliason WK, Wang J, Steitz TA. 2012. The hexameric helicase
509 DnaB adopts a nonplanar conformation during translocation. *Cell* 151:267–277.
- 510 20. Nakayama N, Arai N, Kaziro Y, Arai K. 1984. Structural and functional studies of the dnaB
511 protein using limited proteolysis. Characterization of domains for DNA-dependent ATP
512 hydrolysis and for protein association in the primosome. *The Journal of biological chemistry*
513 259:88–96.
- 514 21. Sakamoto Y, Nakai S, Moriya S, Yoshikawa H, Ogasawara N. 1995. The *Bacillus subtilis*
515 dnaC gene encodes a protein homologous to the DnaB helicase of *Escherichia coli*.
516 *Microbiology (Reading, England)* 141 (Pt 3:641–644.
- 517 22. Liu B, Eliason WK, Steitz TA. 2013. Structure of a helicase–helicase loader complex reveals
518 insights into the mechanism of bacterial primosome assembly. *Nat Commun* 4:2495.
- 519 23. Arias-Palomo E, Puri N, O’Shea Murray VL, Yan Q, Berger JM. 2019. Physical Basis for the
520 Loading of a Bacterial Replicative Helicase onto DNA. *Molecular cell* 74:173-184.e4.

- 521 24. Kobori JA, Kornberg A. 1982. The *Escherichia coli* dnaC gene product. II. Purification,
522 physical properties, and role in replication. *Journal of Biological Chemistry* 257:13763–
523 13769.
- 524 25. Wahle E, Lasken RS, Kornberg A. 1989. The dnaB-dnaC replication protein complex of
525 *Escherichia coli*. II. Role of the complex in mobilizing dnaB functions. *J Biol Chem*
526 264:2469–2475.
- 527 26. Koonin EV. 1992. DnaC protein contains a modified ATP-binding motif and belongs to a
528 novel family of ATPases including also DnaA. *Nucleic acids research* 20:1997–1997.
- 529 27. Bruand C, Farache M, McGovern S, Ehrlich SD, Polard P. 2001. DnaB, DnaD and Dnal
530 proteins are components of the *Bacillus subtilis* replication restart primosome. *Molecular*
531 *microbiology* 42:245–255.
- 532 28. Brézellec P, Vallet-Gely I, Possoz C, Quevillon-Cheruel S, Ferat J-L. 2016. DciA is an
533 ancestral replicative helicase operator essential for bacterial replication initiation. *Nat*
534 *Commun* 7.
- 535 29. Davey MJ, Fang L, McInerney P, Georgescu RE, O'Donnell M. 2002. The DnaC helicase
536 loader is a dual ATP/ADP switch protein. *EMBO J* 21:3148–3159.
- 537 30. Ioannou C, Schaeffer PM, Dixon NE, Soultanas P. 2006. Helicase binding to Dnal exposes
538 a cryptic DNA-binding site during helicase loading in *Bacillus subtilis*. *Nucleic Acids Res*
539 34:5247–5258.

- 540 31. Mann KM, Huang DL, Hooppaw AJ, Logsdon MM, Richardson K, Lee HJ, Kimmey JM,
541 Aldridge BB, Stallings CL. 2017. Rv0004 is a new essential member of the mycobacterial
542 DNA replication machinery. *PLoS Genet* 13:e1007115.
- 543 32. Chan-Yao-Chong M, Marsin S, Quevillon-Cheruel S, Durand D, Ha-Duong T. 2020.
544 Structural ensemble and biological activity of DciA intrinsically disordered region. *Journal*
545 *of Structural Biology* 212:107573.
- 546 33. Marsin S, Adam Y, Cargemel C, Andreani J, Baconnais S, Legrand P, Li de la Sierra-Gallay
547 I, Humbert A, Aumont-Nicaise M, Velours C, Ochsenbein F, Durand D, Le Cam E, Walbott
548 H, Possoz C, Quevillon-Cheruel S, Ferat J-L. 2021. Study of the DnaB:DciA interplay reveals
549 insights into the primary mode of loading of the bacterial replicative helicase. *Nucleic Acids*
550 *Research* gkab463.
- 551 34. Joseph T Burke*, Samuel Z Chen*, Lauren M Sosinski*, John B Johnson, Janani Ravi.
552 MolEvolvR: A web-app for characterizing proteins using molecular evolution and
553 phylogeny. In preparation.
- 554 35. Blackwell GA, Hunt M, Malone KM, Lima L, Horesh G, Alako BTF, Thomson NR, Iqbal Z.
555 2021. Exploring bacterial diversity via a curated and searchable snapshot of archived DNA
556 sequences. *PLoS biology* 19:e3001421–e3001421.
- 557 36. Land M, Hauser L, Jun S-R, Nookaew I, Leuze MR, Ahn T-H, Karpinets T, Lund O, Kora G,
558 Wassenaar T, Poudel S, Ussery DW. 2015. Insights from 20 years of bacterial genome
559 sequencing. *Functional & integrative genomics* 15:141–161.

- 560 37. Tatusova T, Ciufu S, Fedorov B, O'Neill K, Tolstoy I. 2014. RefSeq microbial genomes
561 database: new representation and annotation strategy. *Nucleic Acids Research* 42:D553–
562 D559.
- 563 38. Robinson A, Causer RJ, Dixon NE. 2012. Architecture and conservation of the bacterial
564 DNA replication machinery, an underexploited drug target. *Current drug targets* 13:352–
565 372.
- 566 39. Hug LA, Baker BJ, Anantharaman K, Brown CT, Probst AJ, Castelle CJ, Butterfield CN,
567 Hermsdorf AW, Amano Y, Ise K, Suzuki Y, Dudek N, Relman DA, Finstad KM, Amundson R,
568 Thomas BC, Banfield JF. 2016. A new view of the tree of life. *Nature Microbiology* 1:16048–
569 16048.
- 570 40. Necci M, Piovesan D, Dosztányi Z, Tosatto SCE. 2017. MobiDB-lite: fast and highly specific
571 consensus prediction of intrinsic disorder in proteins. *Bioinformatics (Oxford, England)*
572 33:1402–1404.
- 573 41. Burroughs AM, Zhang D, Schäffer DE, Iyer LM, Aravind L. 2015. Comparative genomic
574 analyses reveal a vast, novel network of nucleotide-centric systems in biological conflicts,
575 immunity and signaling. *Nucleic acids research*, 2015/11/20 ed. 43:10633–10654.
- 576 42. Sutton MD, Carr KM, Vicente M, Kaguni JM. 1998. *Escherichia coli* DnaA protein. The N-
577 terminal domain and loading of DnaB helicase at the *E. coli* chromosomal origin. *J Biol*
578 *Chem* 273:34255–34262.
- 579 43. Abe Y, Jo T, Matsuda Y, Matsunaga C, Katayama T, Ueda T. 2007. Structure and function
580 of DnaA N-terminal domains: specific sites and mechanisms in inter-DnaA interaction and
581 in DnaB helicase loading on oriC. *The Journal of biological chemistry* 282:17816–17827.

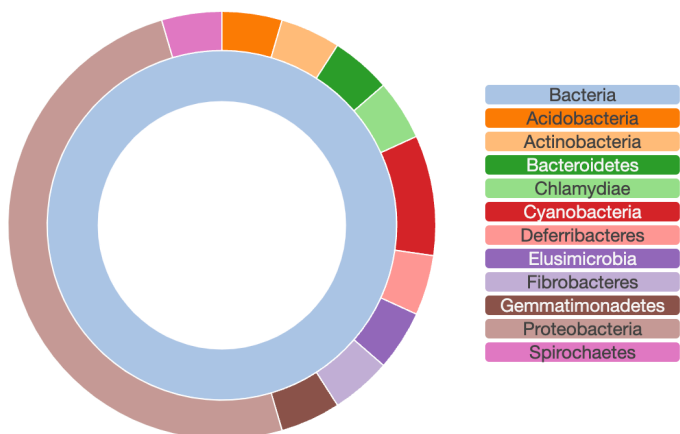
- 582 44. Keyamura K, Abe Y, Higashi M, Ueda T, Katayama T. 2009. DiaA dynamics are coupled
583 with changes in initial origin complexes leading to helicase loading. *The Journal of*
584 *biological chemistry* 284:25038–25050.
- 585 45. Tan HY, Wilczek LA, Pottinger S, Manosas M, Yu C, Nguyenduc T, Bianco PR. 2017. The
586 intrinsically disordered linker of *E. coli* SSB is critical for the release from single-stranded
587 DNA. *Protein science: a publication of the Protein Society* 26:700–717.
- 588 46. Kozlov AG, Weiland E, Mittal A, Waldman V, Antony E, Fazio N, Pappu RV, Lohman TM.
589 2015. Intrinsically Disordered C-Terminal Tails of *E. coli* Single-Stranded DNA Binding
590 Protein Regulate Cooperative Binding to Single-Stranded DNA. *Journal of Molecular*
591 *Biology* 427:763–774.
- 592 47. Antony E, Lohman TM. 2019. Dynamics of *E. coli* single stranded DNA binding (SSB)
593 protein-DNA complexes. *Seminars in cell & developmental biology* 86:102–111.
- 594 48. Bianco PR, Pottinger S, Tan HY, Nguyenduc T, Rex K, Varshney U. 2017. The IDL of *E. coli*
595 SSB links ssDNA and protein binding by mediating protein-protein interactions. *Protein*
596 *science: a publication of the Protein Society* 26:227–241.
- 597 49. Guy CP, Atkinson J, Gupta MK, Mahdi AA, Gwynn EJ, Rudolph CJ, Moon PB, van
598 Knippenberg IC, Cadman CJ, Dillingham MS, Lloyd RG, McGlynn P. 2009. Rep Provides a
599 Second Motor at the Replisome to Promote Duplication of Protein-Bound DNA. *Molecular*
600 *Cell* 36:654–666.
- 601 50. Nguyen B, Shinn MK, Weiland E, Lohman TM. 2021. Regulation of *E. coli* Rep helicase
602 activity by PriC. *Journal of Molecular Biology* 433:167072–167072.

- 603 51. O’Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, Rajput B, Robbertse
604 B, Smith-White B, Ako-Adjei D, Astashyn A, Badretdin A, Bao Y, Blinkova O, Brover V,
605 Chetvernin V, Choi J, Cox E, Ermolaeva O, Farrell CM, Goldfarb T, Gupta T, Haft D, Hatcher
606 E, Hlavina W, Joardar VS, Kodali VK, Li W, Maglott D, Masterson P, McGarvey KM, Murphy
607 MR, O’Neill K, Pujar S, Rangwala SH, Rausch D, Riddick LD, Schoch C, Shkeda A, Storz
608 SS, Sun H, Thibaud-Nissen F, Tolstoy I, Tully RE, Vatsan AR, Wallin C, Webb D, Wu W,
609 Landrum MJ, Kimchi A, Tatusova T, DiCuccio M, Kitts P, Murphy TD, Pruitt KD. 2016.
610 Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and
611 functional annotation. *Nucleic acids research* 44:D733-45.
- 612 52. Mistry J, Chuguransky S, Williams L, Qureshi M, Salazar GA, Sonnhammer ELL, Tosatto
613 SCE, Paladin L, Raj S, Richardson LJ, Finn RD, Bateman A. 2021. Pfam: The protein families
614 database in 2021. *Nucleic Acids Research* 49:D412–D419.
- 615 53. Nowotny J, Ahmed S, Xu L, Oluwadare O, Chen H, Hensley N, Trieu T, Cao R, Cheng J.
616 2015. Iterative reconstruction of three-dimensional models of human chromosomes from
617 chromosomal contact data. *BMC bioinformatics* 16:338–338.
- 618 54. Käll L, Krogh A, Sonnhammer ELL. 2004. A combined transmembrane topology and signal
619 peptide prediction method. *Journal of molecular biology* 338:1027–1036.
- 620 55. Möller S, Croning MD, Apweiler R. 2001. Evaluation of methods for the prediction of
621 membrane spanning regions. *Bioinformatics (Oxford, England)* 17:646–653.
- 622 56. Piovesan D, Necci M, Escobedo N, Monzon AM, Hatos A, Mičetić I, Quaglia F, Paladin L,
623 Ramasamy P, Dosztányi Z, Vranken WF, Davey NE, Parisi G, Fuxreiter M, Tosatto SCE.

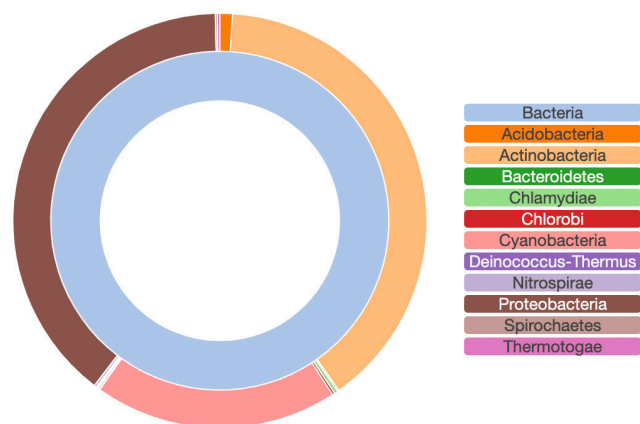
- 624 2021. MobiDB: intrinsically disordered proteins in 2021. *Nucleic acids research* 49:D361–
625 D367.
- 626 57. Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, Basutkar P, Tivey ARN,
627 Potter SC, Finn RD, Lopez R. 2019. The EMBL-EBI search and sequence analysis tools
628 APIs in 2019. *Nucleic acids research* 47:W636–W641.
- 629 58. Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009. Jalview Version 2—
630 a multiple sequence alignment editor and analysis workbench. *Bioinformatics* 25:1189–
631 1191.
- 632 59. Campanella JJ, Bitincka L, Smalley J. 2003. MatGAT: An application that generates
633 similarity/identity matrices using protein or DNA sequences. *BMC Bioinformatics* 4:29–29.
- 634 60. Troshin PV, Procter JB, Barton GJ. 2011. Java bioinformatics analysis web services for
635 multiple sequence alignment—JABAWS:MSA. *Bioinformatics* 27:2001–2002.
- 636 61. Lassmann T. 2020. Kalign 3: multiple sequence alignment of large datasets. *Bioinformatics*
637 36:1928–1929.
- 638

Figure 1

A Query distribution



B Homolog distribution



C Phyletic spread of homologs (lineages)

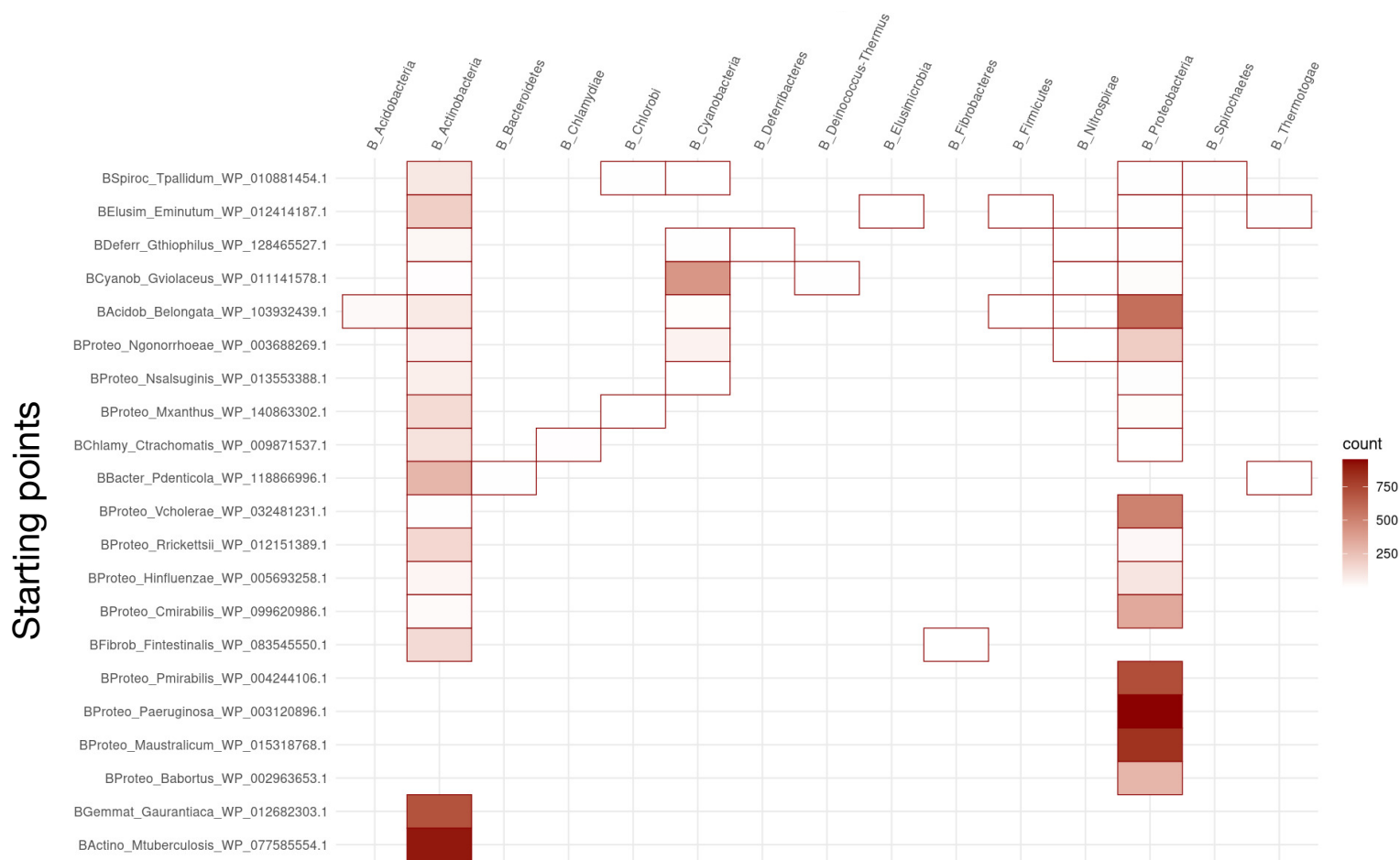
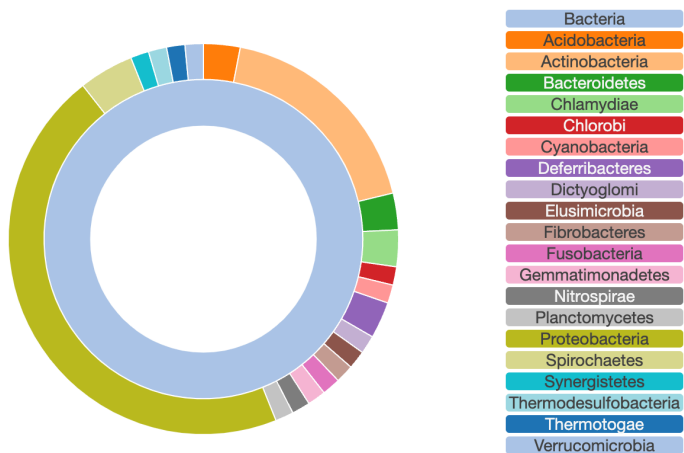


Figure 2

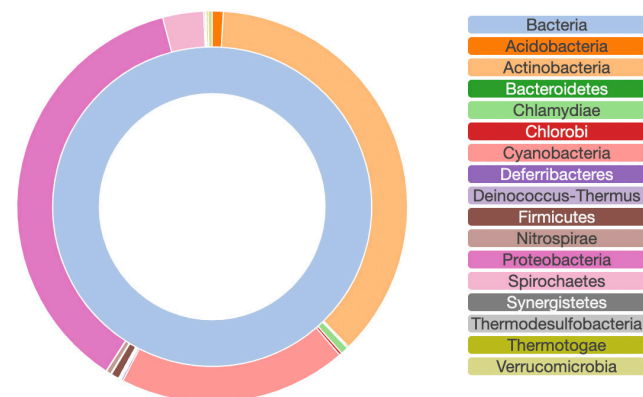
Phylum	Species	B. elo	M. tub	P. den	C. tra	G. vio	G. thio	E. min	F. int	G. aur	B. abo	P. aer	N. gon	P. mir	H. inf	R. ric	N. sal	M. aus	V. cho	C. mir	M. xan	T. pal	
Acidobacteria	<i>Bryocella elongata</i>																						
Actinobacteria	<i>Mycobacterium tuberculosis</i>	23.0																					
Bacteroidetes	<i>Prevotella denticola</i>	42.3	26.2																				
Chlamydiae	<i>Chlamydia trachomatis</i>	31.9	28.9	41.4																			
Cyanobacteria	<i>Gloeobacter violaceus</i>	26.7	33.2	26.1	28.9																		
Deferribacteres	<i>Geovibrio thiophilus</i>	32.9	28.9	27.4	30.8	34.4																	
Elusimicrobia	<i>Elusimicrobium minutum</i>	38.5	21.4	38.8	39.7	26.1	29.5																
Fibrobacteres	<i>Fibrobacter intestinalis</i>	36.3	27.3	41.6	42.2	26.7	30.8	35.4															
Gemmatimonadetes	<i>Gemmatimonas aurantiaca</i>	39.4	28.3	43.7	41.4	28.3	33.6	38.8	34.5														
Proteobacteria	<i>Brucella abortus</i>	25.1	31.0	26.3	29.1	33.3	32.6	20.6	25.7	33.1													
Proteobacteria	<i>Pseudomonas aeruginosa</i>	29.8	28.3	31.3	31.3	30.0	34.2	30.5	31.3	29.8	36.0												
Proteobacteria	<i>Neisseria gonorrhoeae</i>	30.7	29.9	29.3	29.3	32.8	37.0	25.7	28.6	29.3	34.9	39.3											
Proteobacteria	<i>Proteus mirabilis</i>	26.2	29.9	26.7	32.0	38.3	39.0	29.1	32.0	24.4	40.0	44.2	38.4										
Proteobacteria	<i>Haemophilus influenzae</i>	35.6	23.0	40.4	38.8	23.3	29.5	42.3	39.8	36.5	21.1	28.2	26.4	26.7									
Proteobacteria	<i>Rickettsia rickettsii</i>	30.8	25.7	43.9	39.7	26.1	31.5	38.3	39.8	36.4	26.9	30.5	30.0	26.7	32.7								
Proteobacteria	<i>Nitratifractor salsuginis</i>	24.8	26.7	24.8	22.9	30.6	32.7	28.8	29.4	26.1	29.7	27.5	30.7	32.6	27.5	33.3							
Proteobacteria	<i>Mesorhizobium australicum</i>	25.9	32.1	24.7	27.1	37.2	34.3	23.5	27.1	27.7	62.3	37.3	35.5	37.2	29.5	27.1	28.9						
Proteobacteria	<i>Vibrio cholerae</i>	27.4	27.8	26.1	32.5	32.2	29.9	26.1	30.6	36.3	33.1	42.0	38.9	52.3	30.6	32.5	35.7	34.9					
Proteobacteria	<i>Caulobacter mirabilis</i>	28.0	34.8	26.9	27.5	40.7	34.1	22.0	25.3	28.0	47.8	28.6	35.2	34.1	23.1	27.5	33.0	46.2	35.7				
Proteobacteria	<i>Myxococcus xanthus</i>	38.5	28.3	42.6	32.8	26.7	31.5	34.7	40.7	41.7	26.9	31.3	29.3	26.2	27.9	40.2	25.5	30.1	22.9	26.9			
Spirochaetes	<i>Treponema pallidum</i>	31.0	31.0	28.3	29.7	31.1	39.0	28.3	26.2	32.4	29.7	36.6	31.7	38.4	34.5	29.7	35.3	34.3	37.6	31.9	26.9		

Figure 3

A Query distribution



B Homolog distribution



C Phyletic spread of homologs (lineages)

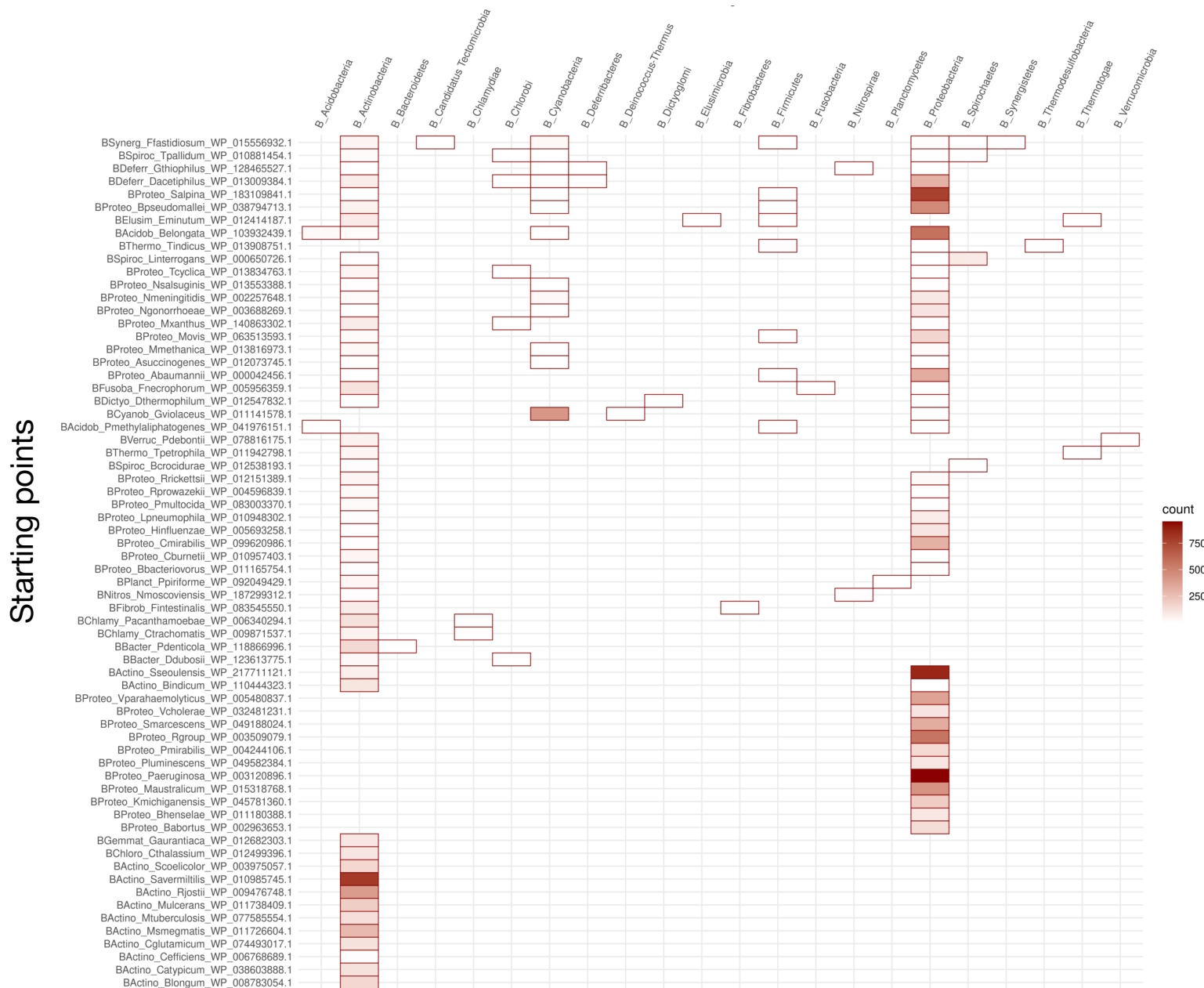


Figure 4

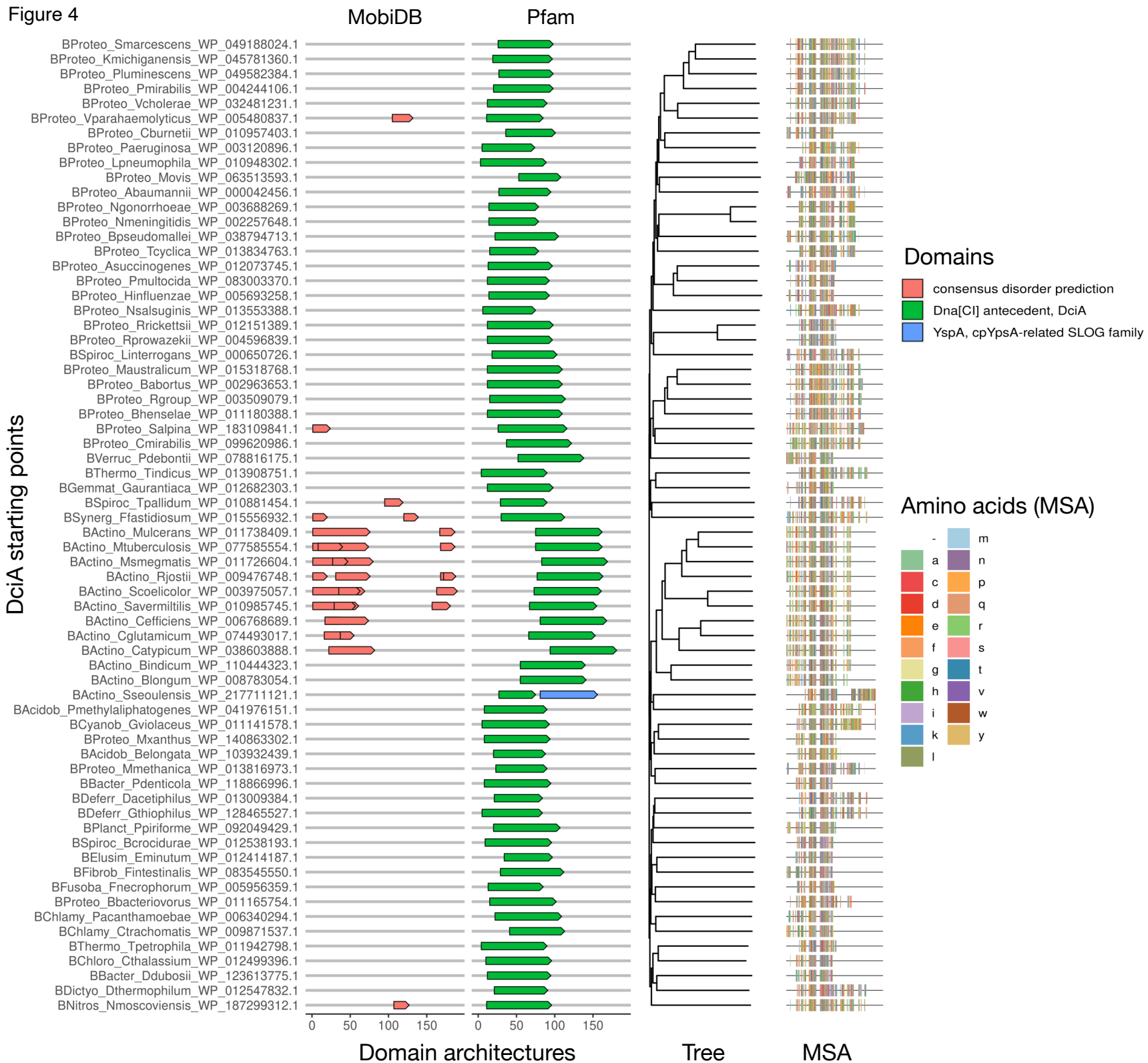
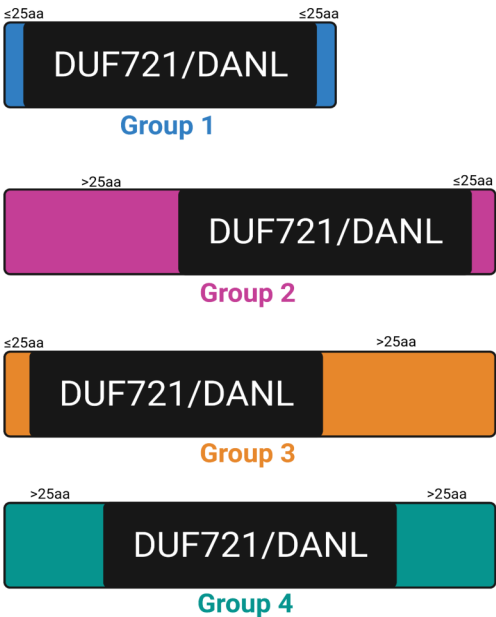
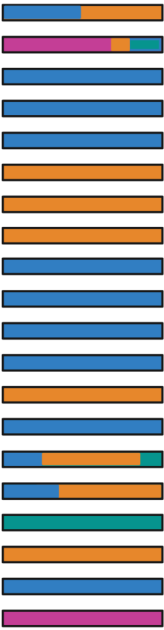


Figure 5

A



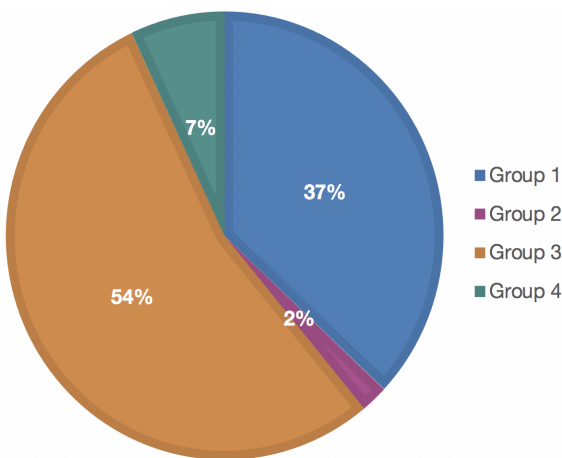
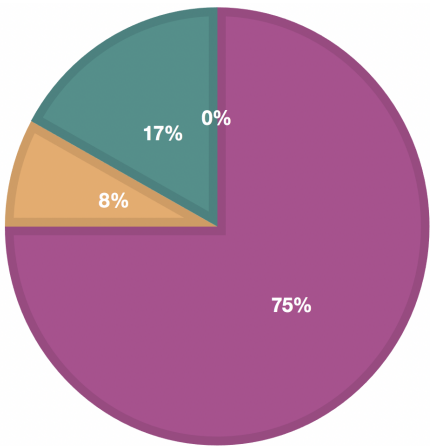
- Acidobacteria (2)
- Actinobacteria (12)
- Bacteroidetes (2)
- Chlamydiae (2)
- Chlorobi (1)
- Cyanobacteria (1)
- Deferribacteres (2)
- Dictyoglomi (1)
- Elusimicrobia (1)
- Fibrobacteres (1)
- Fusobacteria (1)
- Gemmatimonadetes (1)
- Nitrospirae (1)
- Planctomycetes (1)
- Proteobacteria (30)
- Spirochaetes (3)
- Synergistetes (1)
- Thermodesulfobacteria (1)
- Thermotogae (1)
- Verrucomicrobiae (1)



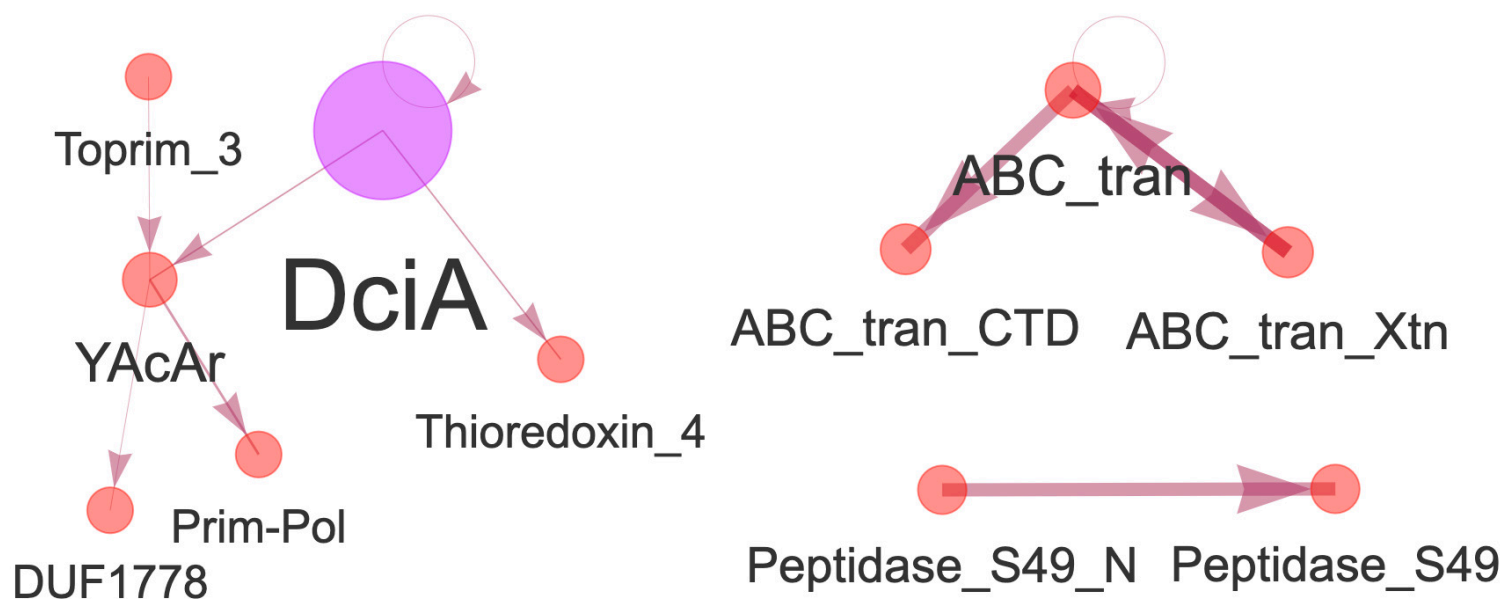
Gram Positive Bacteria (12)

Gram Negative Bacteria (54)

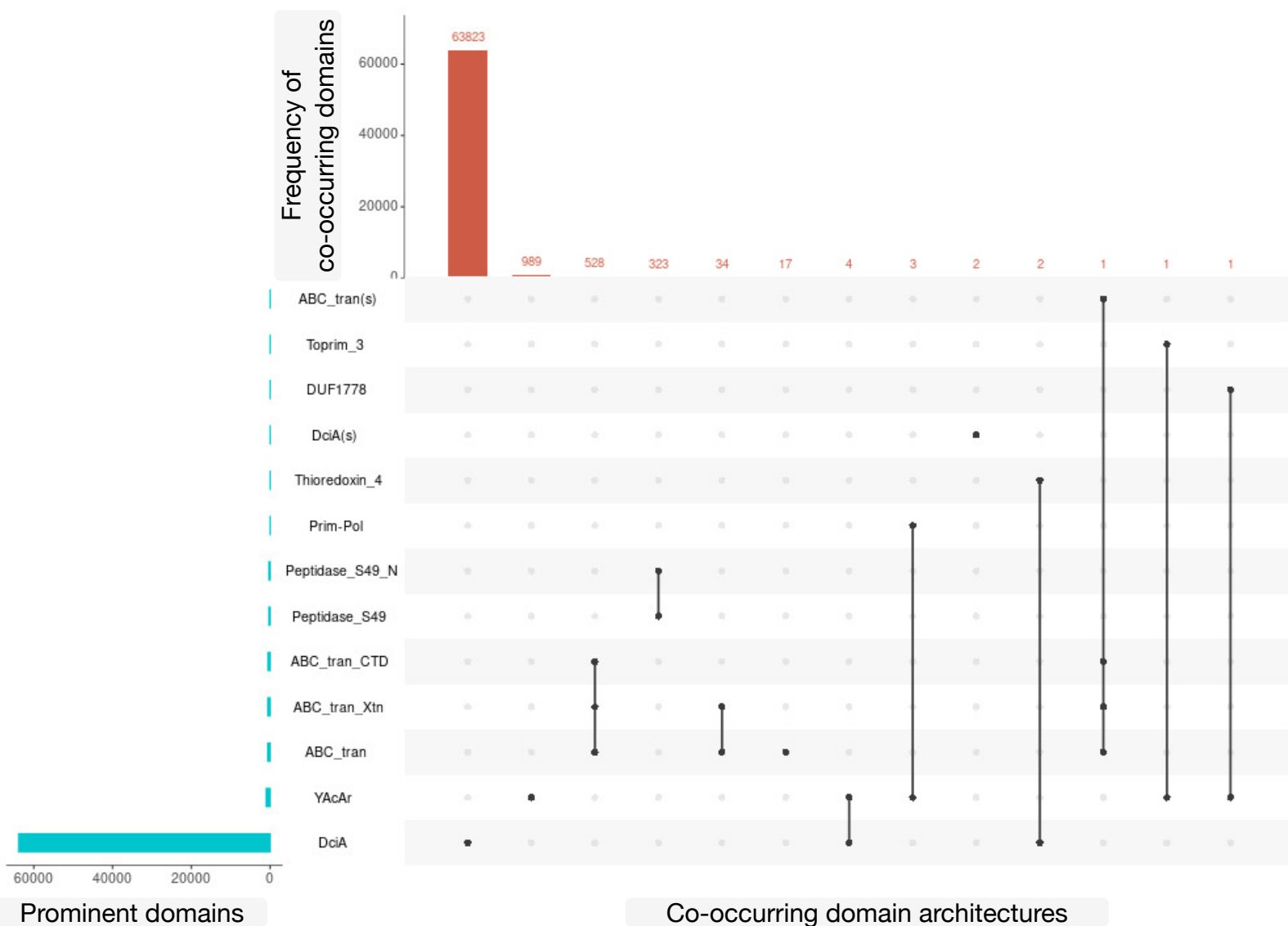
B



A Domain proximity network



B Domain co-occurrence frequency



Accession	Gene	Species	Strain	Accession	Gene	Species	Strain
Phage T400000001	MP101000001	Parvovirus sp.	001	MP101000001	NS	001	001
Phage T400000002	MP101000002	Parvovirus sp.	002	MP101000002	NS	002	002
Phage T400000003	MP101000003	Parvovirus sp.	003	MP101000003	NS	003	003
Phage T400000004	MP101000004	Parvovirus sp.	004	MP101000004	NS	004	004
Phage T400000005	MP101000005	Parvovirus sp.	005	MP101000005	NS	005	005
Phage T400000006	MP101000006	Parvovirus sp.	006	MP101000006	NS	006	006
Phage T400000007	MP101000007	Parvovirus sp.	007	MP101000007	NS	007	007
Phage T400000008	MP101000008	Parvovirus sp.	008	MP101000008	NS	008	008
Phage T400000009	MP101000009	Parvovirus sp.	009	MP101000009	NS	009	009
Phage T400000010	MP101000010	Parvovirus sp.	010	MP101000010	NS	010	010
Phage T400000011	MP101000011	Parvovirus sp.	011	MP101000011	NS	011	011
Phage T400000012	MP101000012	Parvovirus sp.	012	MP101000012	NS	012	012
Phage T400000013	MP101000013	Parvovirus sp.	013	MP101000013	NS	013	013
Phage T400000014	MP101000014	Parvovirus sp.	014	MP101000014	NS	014	014
Phage T400000015	MP101000015	Parvovirus sp.	015	MP101000015	NS	015	015
Phage T400000016	MP101000016	Parvovirus sp.	016	MP101000016	NS	016	016
Phage T400000017	MP101000017	Parvovirus sp.	017	MP101000017	NS	017	017
Phage T400000018	MP101000018	Parvovirus sp.	018	MP101000018	NS	018	018
Phage T400000019	MP101000019	Parvovirus sp.	019	MP101000019	NS	019	019
Phage T400000020	MP101000020	Parvovirus sp.	020	MP101000020	NS	020	020
Phage T400000021	MP101000021	Parvovirus sp.	021	MP101000021	NS	021	021
Phage T400000022	MP101000022	Parvovirus sp.	022	MP101000022	NS	022	022
Phage T400000023	MP101000023	Parvovirus sp.	023	MP101000023	NS	023	023
Phage T400000024	MP101000024	Parvovirus sp.	024	MP101000024	NS	024	024
Phage T400000025	MP101000025	Parvovirus sp.	025	MP101000025	NS	025	025
Phage T400000026	MP101000026	Parvovirus sp.	026	MP101000026	NS	026	026
Phage T400000027	MP101000027	Parvovirus sp.	027	MP101000027	NS	027	027
Phage T400000028	MP101000028	Parvovirus sp.	028	MP101000028	NS	028	028
Phage T400000029	MP101000029	Parvovirus sp.	029	MP101000029	NS	029	029
Phage T400000030	MP101000030	Parvovirus sp.	030	MP101000030	NS	030	030
Phage T400000031	MP101000031	Parvovirus sp.	031	MP101000031	NS	031	031
Phage T400000032	MP101000032	Parvovirus sp.	032	MP101000032	NS	032	032
Phage T400000033	MP101000033	Parvovirus sp.	033	MP101000033	NS	033	033
Phage T400000034	MP101000034	Parvovirus sp.	034	MP101000034	NS	034	034
Phage T400000035	MP101000035	Parvovirus sp.	035	MP101000035	NS	035	035
Phage T400000036	MP101000036	Parvovirus sp.	036	MP101000036	NS	036	036
Phage T400000037	MP101000037	Parvovirus sp.	037	MP101000037	NS	037	037
Phage T400000038	MP101000038	Parvovirus sp.	038	MP101000038	NS	038	038
Phage T400000039	MP101000039	Parvovirus sp.	039	MP101000039	NS	039	039
Phage T400000040	MP101000040	Parvovirus sp.	040	MP101000040	NS	040	040
Phage T400000041	MP101000041	Parvovirus sp.	041	MP101000041	NS	041	041
Phage T400000042	MP101000042	Parvovirus sp.	042	MP101000042	NS	042	042
Phage T400000043	MP101000043	Parvovirus sp.	043	MP101000043	NS	043	043
Phage T400000044	MP101000044	Parvovirus sp.	044	MP101000044	NS	044	044
Phage T400000045	MP101000045	Parvovirus sp.	045	MP101000045	NS	045	045
Phage T400000046	MP101000046	Parvovirus sp.	046	MP101000046	NS	046	046
Phage T400000047	MP101000047	Parvovirus sp.	047	MP101000047	NS	047	047
Phage T400000048	MP101000048	Parvovirus sp.	048	MP101000048	NS	048	048
Phage T400000049	MP101000049	Parvovirus sp.	049	MP101000049	NS	049	049
Phage T400000050	MP101000050	Parvovirus sp.	050	MP101000050	NS	050	050
Phage T400000051	MP101000051	Parvovirus sp.	051	MP101000051	NS	051	051
Phage T400000052	MP101000052	Parvovirus sp.	052	MP101000052	NS	052	052
Phage T400000053	MP101000053	Parvovirus sp.	053	MP101000053	NS	053	053
Phage T400000054	MP101000054	Parvovirus sp.	054	MP101000054	NS	054	054
Phage T400000055	MP101000055	Parvovirus sp.	055	MP101000055	NS	055	055
Phage T400000056	MP101000056	Parvovirus sp.	056	MP101000056	NS	056	056
Phage T400000057	MP101000057	Parvovirus sp.	057	MP101000057	NS	057	057
Phage T400000058	MP101000058	Parvovirus sp.	058	MP101000058	NS	058	058
Phage T400000059	MP101000059	Parvovirus sp.	059	MP101000059	NS	059	059
Phage T400000060	MP101000060	Parvovirus sp.	060	MP101000060	NS	060	060
Phage T400000061	MP101000061	Parvovirus sp.	061	MP101000061	NS	061	061
Phage T400000062	MP101000062	Parvovirus sp.	062	MP101000062	NS	062	062
Phage T400000063	MP101000063	Parvovirus sp.	063	MP101000063	NS	063	063
Phage T400000064	MP101000064	Parvovirus sp.	064	MP101000064	NS	064	064
Phage T400000065	MP101000065	Parvovirus sp.	065	MP101000065	NS	065	065
Phage T400000066	MP101000066	Parvovirus sp.	066	MP101000066	NS	066	066
Phage T400000067	MP101000067	Parvovirus sp.	067	MP101000067	NS	067	067
Phage T400000068	MP101000068	Parvovirus sp.	068	MP101000068	NS	068	068
Phage T400000069	MP101000069	Parvovirus sp.	069	MP101000069	NS	069	069
Phage T400000070	MP101000070	Parvovirus sp.	070	MP101000070	NS	070	070
Phage T400000071	MP101000071	Parvovirus sp.	071	MP101000071	NS	071	071
Phage T400000072	MP101000072	Parvovirus sp.	072	MP101000072	NS	072	072
Phage T400000073	MP101000073	Parvovirus sp.	073	MP101000073	NS	073	073
Phage T400000074	MP101000074	Parvovirus sp.	074	MP101000074	NS	074	074
Phage T400000075	MP101000075	Parvovirus sp.	075	MP101000075	NS	075	075
Phage T400000076	MP101000076	Parvovirus sp.	076	MP101000076	NS	076	076
Phage T400000077	MP101000077	Parvovirus sp.	077	MP101000077	NS	077	077
Phage T400000078	MP101000078	Parvovirus sp.	078	MP101000078	NS	078	078
Phage T400000079	MP101000079	Parvovirus sp.	079	MP101000079	NS	079	079
Phage T400000080	MP101000080	Parvovirus sp.	080	MP101000080	NS	080	080
Phage T400000081	MP101000081	Parvovirus sp.	081	MP101000081	NS	081	081
Phage T400000082	MP101000082	Parvovirus sp.	082	MP101000082	NS	082	082
Phage T400000083	MP101000083	Parvovirus sp.	083	MP101000083	NS	083	083
Phage T400000084	MP101000084	Parvovirus sp.	084	MP101000084	NS	084	084
Phage T400000085	MP101000085	Parvovirus sp.	085	MP101000085	NS	085	085
Phage T400000086	MP101000086	Parvovirus sp.	086	MP101000086	NS	086	086
Phage T400000087	MP101000087	Parvovirus sp.	087	MP101000087	NS	087	087
Phage T400000088	MP101000088	Parvovirus sp.	088	MP101000088	NS	088	088
Phage T400000089	MP101000089	Parvovirus sp.	089	MP101000089	NS	089	089
Phage T400000090	MP101000090	Parvovirus sp.	090	MP101000090	NS	090	090
Phage T400000091	MP101000091	Parvovirus sp.	091	MP101000091	NS	091	091
Phage T400000092	MP101000092	Parvovirus sp.	092	MP101000092	NS	092	092
Phage T400000093	MP101000093	Parvovirus sp.	093	MP101000093	NS	093	093
Phage T400000094	MP101000094	Parvovirus sp.	094	MP101000094	NS	094	094
Phage T400000095	MP101000095	Parvovirus sp.	095	MP101000095	NS	095	095
Phage T400000096	MP101000096	Parvovirus sp.	096	MP101000096	NS	096	096
Phage T400000097	MP101000097	Parvovirus sp.	097	MP101000097	NS	097	097
Phage T400000098	MP101000098	Parvovirus sp.	098	MP101000098	NS	098	098
Phage T400000099	MP101000099	Parvovirus sp.	099	MP101000099	NS	099	099
Phage T400000100	MP101000100	Parvovirus sp.	100	MP101000100	NS	100	100