# Comprehensive Genomic and Evolutionary Analysis of Biofilm Matrix Clusters and Proteins in the *Vibrio* Genus

Yiyan Yang[1,*], Jing Yan[2,3], Rich Olson[4], Xiaofang Jiang[1,*]

[1]Intramural Research Program, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA

[2]Department of Molecular, Cellular and Developmental Biology, Yale University, New Haven, CT, USA

[3]Quantitative Biology Institute, Yale University, New Haven, CT, USA

[4]Department of Molecular Biology and Biochemistry, Molecular Biophysics Program, Wesleyan University, Middletown, CT, USA

[*]Corresponding authors: Yiyan Yang, E-mail: yiyan.yang@nih.gov; Xiaofang Jiang, E-mail: xiaofang.jiang@nih.gov

**Abstract**

*Vibrio cholerae* pathogens cause cholera, an acute diarrheal disease resulting in significant morbidity and mortality worldwide. Biofilms in vibrios enhance their survival in natural ecosystems and facilitate transmission during cholera outbreaks. Critical components of the biofilm matrix include the *Vibrio* polysaccharides produced by the *vps*-1 and *vps*-2 gene clusters and the biofilm matrix proteins encoded in the *rbm* gene cluster, together comprising the biofilm matrix cluster. However, the biofilm matrix clusters and their evolutionary patterns in other *Vibrio* species remain underexplored. In this study, we systematically investigated the distribution, diversity, and evolution of biofilm matrix clusters and proteins across the *Vibrio* genus. Our findings reveal that these gene clusters are sporadically distributed throughout the genus, even appearing in species phylogenetically distant from *V. cholerae*. Evolutionary

25  analysis of the major biofilm matrix proteins RbmC and Bap1 shows that they are structurally
26  and sequentially related, having undergone structural domain and modular alterations.
27  Additionally, a novel loop-less Bap1 variant was identified, predominantly represented in two
28  phylogenetically distant *Vibrio cholerae* subspecies clades that share specific gene groups
29  associated with the presence or absence of the protein. Furthermore, our analysis revealed that
30  *rbmB*, a gene involved in biofilm dispersal, shares a recent common ancestor with Vibriophage
31  tail proteins, suggesting that phages may mimic host functions to evade biofilm-associated
32  defenses. Our study offers a foundational understanding of the diversity and evolution of biofilm
33  matrix clusters in vibrios, laying the groundwork for future biofilm engineering through genetic
34  modification.

## Introduction

35

36  *Vibrio cholerae*, the pathogen responsible for cholera, causes an acute diarrheal disease that can
37  lead to hypotonic shock and death. Annually, it infects 3-5 million people, resulting in 100,000–
38  120,000 deaths (1). *V. cholerae* forms biofilms—surface-associated communities encased in a
39  matrix—which enhance survival in ecosystems, and transmission during outbreaks (2, 3), while
40  providing protection from environmental stresses like nutrient scarcity, antimicrobial agents,
41  predation by unicellular eukaryotes, and attack by phages (4–6).

42  The biofilm matrix is primarily comprised of Vibrio polysaccharide (VPS), making up
43  approximately half of its mass and essential for biofilm 3D structural development (7–9). Genes
44  involved in VPS production are organized into two *vps* gene clusters, *vps*-1 and *vps*-2. A gene
45  cluster in this study is defined as a group of closely located genes on a chromosome that are
46  often functionally related and may include multiple operons. The *vps*-1 gene cluster contains 12
47  genes (*vpsU*, VC0916 and *vpsA-K*, VC0917-VC0927) while the *vps*-2 gene cluster is relatively
48  shorter only containing 6 genes (*vpsL-Q*, VC0934-VC0939) (7, 9). Meanwhile, biofilm matrix
49  proteins, such as RbmA, RbmC and Bap1, encoded by *rbmA* (VC0928), *rbmC* (VC0930) and
50  *bap1* (VC1888), respectively, are crucial for preserving the structural integrity of the wild-type
51  biofilm (10, 11), among which RbmA and RbmC are encoded in a *rbm* (rugosity and biofilm
52  structure modulator) gene cluster separating the two *vps* gene clusters. The gene encoding Bap1
53  is distant from the *rbm* gene cluster, yet it also modulates the development of corrugated colonies
54  and is crucial for biofilm formation (11–13). RbmA, as a biofilm scaffolding protein involved in
55  cell-cell and cell-biofilm adhesion, is required for rugose colony formation and biofilm structure
56  integrity in *V. cholerae* (10, 11, 13–15). The other two major biofilm matrix proteins, RbmC and
57  Bap1, are homologues sharing 47% sequence similarity and containing overlapping domains to
58  facilitate their robust adhesion to diverse surfaces (11, 16). Both proteins have a conserved β-
59  propeller domain with eight blades and at least one β-prism domain. RbmC, however, is
60  characterized by two β-prism domains and additional tandem β/γ crystallin domains, known as
61  M1M2 (16, 17). Most notably, Bap1's β-prism contains an additional 57-amino acid (aa)
62  sequence which promotes *V. cholerae* biofilm adhesion to lipids and abiotic surfaces while
63  RbmC mainly mediates binding to host surfaces through recognition of N- and O-glycans and
64  mucins (16). Another interesting gene in the *rbm* gene cluster is *rbmB* (VC0929), which encodes
65  a putative polysaccharide lyase that has been proposed to have a role in VPS degradation and cell

66　detachment (11, 18–20). Other genes included in the *rbm* gene cluster are *rbmDEF* (VC0931-
67　VC0933). Together, the *vps*-1, *rbm* and *vps*-2 gene clusters comprise a functional genetic module
68　— the *V. cholerae* biofilm matrix cluster (*V. cholerae* BMC or VcBMC) (18).

69　The biofilm matrix cluster has primarily been investigated in commonly studied *V. cholerae*
70　strains and a few other *Vibrio* species (21–24). However, it has not yet been systematically
71　studied at the strain level within *V. cholerae* or more extensively across the *Vibrio* genus. Since
72　the biofilm matrix cluster encodes proteins for VPS synthesis and matrix proteins, which are the
73　major components of *Vibrio* biofilms, a systematic genomic analysis of this cluster and the
74　identification of relevant genes across the *Vibrio* genus can provide a prospective and
75　comprehensive view of the genetic basis underlying VPS production and biofilm formation.

76　In this study, we comprehensively annotated the genes involved in the biofilm matrix cluster to
77　explore their distribution, diversity and gene synteny by conducting large-scale comparative
78　genomics and phylogenetic analyses on 6,121 Vibrio genomes spanning 210 species across the
79　entire *Vibrio* genus as well as within the *V. cholerae* species. We observed not only a prevalent
80　presence of this cluster in *V. cholerae* but also in other distantly related species. Our analysis
81　reveals a distinct evolutionary pattern for the *vps*-1 and *vps*-2 gene clusters: genes in the *vps*-2
82　gene cluster often co-located with *rbmDEF* genes, while *vps*-1 genes are commonly adjacent to
83　*rbmABC* genes. This suggests a functional relatedness between them and explains why these two
84　*vps* gene clusters are separated by a *rbm* gene cluster in contemporary *V. cholerae* strains.
85　Additionally, we inferred that the *bap1* genes originated as an ancient duplication of *rbmC* in a
86　clade of species closely related to *V. cholerae*, while *rbmC* genes are present in two major clades
87　and may have undergone structural domain alterations throughout their evolutionary history.
88　Furthermore, a novel loop-less Bap1 variant was identified, predominantly found in two
89　phylogenetically distant *Vibrio cholerae* subspecies clades that share gene groups linked to the
90　presence/absence of the protein. Finally, our findings suggest that RbmB, a putative VPS
91　degradation enzyme, are evolutionarily related to Vibriophage pectin lyase-like tail proteins. The
92　systematic and accurate curation of biofilm matrix clusters and their proteins not only enhances
93　our understanding of *Vibrio* biofilm formation from a genomic view but also offers insights for
94　developing strategies to engineer and control biofilms.

95

96　**Results**


97　**Biofilm matrix clusters are found in phylogenetically distant *Vibrio* species**

98　Leveraging over 6,000 genomes from Genome Taxonomy Database (GTDB r214) (25) across
99　the *Vibrio* genus, we systematically annotated the proteins within the biofilm matrix clusters and
100　depicted an overview of the cluster's gene occurrences spanning 209 *Vibrio* species and seven *V.*
101　*cholerae* subspecies (Fig.1A). We defined a full biofilm matrix cluster if it contains the 12 key
102　*vps* genes (namely *vpsAB*, *vpsDEF*, *vpsIJK*, and *vpsLMNO*) whose deletions have been shown to
103　cause a dramatic reduction in VPS production and biofilm formation (9) and all of the *rbm* genes.

104 We reconstructed a *Vibrio* species tree, which shares a similar topology to that in a previous
105 study (26), and mapped the presence and absence of the key *vps* genes and *rbm* genes to the tree
106 tips. It is interesting to discover that, using this criterion, the full biofilm matrix clusters not only
107 exist in *V. cholerae* and closely related species (such as *V. metoecus* and *V. mimicus*) but are also
108 sporadically distributed across the *Vibrio* genus in distant species like *V. anguillarum*, *V. ordalii*,
109 *V. aestuarianus*, *V. coralliilyticus*, *V. neptunius* and *V. cortegadensis* (Fig.1A). Among all genes,
110 *vps*-2 genes are the most prevalent genes with *vpsL* existing in 50% of the species, *vpsM* in
111 41.2%, *vpsN* in 58.3% and *vpsQ* in 64.4% following by *vps*-1 genes *vpsA* (33.3%) and *vpsB*
112 (33.8%). The higher prevalence of *vps*-2 genes is due to the identification of *vps*-2 similar loci in
113 our data, such as the *cps* (capsular polysaccharide) locus in *Vibrio parahaemolyticus*, the *wcr*
114 (capsular and rugose polysaccharide) locus in *Vibrio vulnificus*, and *vps*-2-like loci in *Aliivibrio*
115 *fischeri*, all of which contain homologs of *vpsLMNO* (Supplementary Figure 1) (27–31). It is
116 important to note that these loci contain genes associated with functions other than VPS
117 production in biofilms, such as capsular polysaccharide synthesis. Therefore, they are less likely
118 to represent true *vps*-2 gene clusters and are instead designated as *vps*-2 similar gene clusters in
119 this study.

120 We next investigated the gene synteny within the biofilm matrix cluster to gain insights on how
121 the *vps-1*, *vps-2* and *rbm* gene clusters have evolved during the speciation of *Vibrio* species
122 (Figure 1B and Supplementary Figure 2). The Vibrio (sub)species clearly form two major clades,
123 Clades A and B, each of which are featured with different patterns in the biofilm matrix clusters
124 (Fig.1B). The examination of the isolation sources and potential hosts of *Vibrio* species in these
125 clades indicates that Clade A species are primarily isolated from marine water and from healthy
126 or diseased invertebrates such as prawns, corals, and bivalve mollusks like clams and oysters
127 (Supplementary Table 1). In contrast, species in Clade B are mostly found in seawater and
128 brackish waters, inhabiting both invertebrate and vertebrate hosts, including fish (such as *V.*
129 *aestuarianus*, *V. ordalii*, and *V. anguillarum*) and humans (such as *V. metoecus*, *V. mimicus*, and
130 *V. cholerae*), and often acting as pathogens (Supplementary Table 1).

131 From Figure 1B, we also observed that *rbmA* genes are absent in seven *Vibrio* species from
132 Clade A (i.e. *V. hepatarius_A*, *V. hepatarius*, *V. sinaloensis*, *V. atypicus*, *V. tubiashii_A*, *V.*
133 *tubiashii*, and *V. bivalvicida*) despite the presence of *rbmD* and *rbmEF* genes in the same operon
134 and the presence of distant *rbmC* genes. Although these species are phylogenetically distant, we
135 observed conservation in the neighborhoods of their *rbmC* genes. These *rbmC* genes are often
136 immediately adjacent to a gene containing a methyl-accepting chemotaxis domain and are close
137 to an operon encoding a system for the uptake and metabolism of disaccharides, suggesting their
138 potential involvement in sugar binding process (Supplementary Figure 3 and Supplementary
139 Table 2). These species typically possess several, but not all, *vps*-2 similar and *vps*-1 similar
140 genes. For genes not annotated as *vps*-like genes, most of them are glycosyltransferases,
141 acyltransferases and polysaccharide biosynthesis proteins, which might be responsible for the
142 synthesis, modification and export of VPS (Supplementary Figure 2 and Supplementary Table 3).

143 Additionally, we observed that *vps*-1 gene clusters tend to co-locate with *rbmABC* genes, while
144 *vps*-2 gene clusters consistently pair with *rbmDEF* genes (see red and blue boxes in Figure 1B).

145     This patten is evident in a sub-lineage of Clade A, which includes *V. coralliilyticus*, *V.*
146     *coralliilyticus_A*, *V. neptunius*, and *V. sp013113835*. In this sub-lineage, *vps*-2 and *rbmDEF*
147     genes are joined but remain distant from the joined *vps*-1 and *rbmABC* genes (Fig.1B). In
148     contrast, Clade B features an intact biofilm matrix cluster, where the *vps*-2 genes and *rbmDEF*
149     genes are consistently adjacent and linked to the joined *vps*-1 and *rbmABC* genes. We also
150     observed that in *V. aestuarianus*, *V. anguillarum*, and *V. ordalii*, the *vps*-2 gene cluster is in the
151     opposite orientation compared to other species within Clade B. Overall, the consistent co-
152     location of the *vps*-1 genes with *rbmABC* and the *vps*-2 genes with *rbmDEF* in several Clade A
153     species and across the whole Clade B suggests their functional associations. This organization
154     may also help explain the intact biofilm matrix clusters commonly observed in *Vibrio cholerae*
155     strains, where the two *vps* gene clusters, separated by a *rbm* gene cluster, could result from the
156     merging of *rbmABC* and *rbmDEF* genes.

157     **Biofilm matrix proteins RbmC and Bap1 experienced structural diversification during**

158     **evolution**

159     RbmC and Bap1, two major biofilm matrix proteins in Vibrio biofilms that share 47% sequence
160     identity, have been shown in previous studies to possess both shared and distinct domains. This
161     suggests that they are functionally and evolutionarily related, with potential domain gain and loss
162     occurring during their evolution (11, 16). Furthermore, we are interested in Bap1-encoded gene
163     due to its higher mutation frequency of 0.0718 compared to all *rbm* genes, implying it may have
164     undergone a stronger positive selection pressure (Supplementary Figure 4). To examine the
165     origin and divergence of these two matrix proteins, we compiled a data set consisting of 2,004
166     *rbmC* and 2,062 *bap1* genes identified across the *Vibrio* genus.

167     There are three extra RbmC variants as well as one Bap1 variant (Supplementary Figures 5 and
168     6). Two of the RbmC variants differ from the standard RbmC protein by having none or only one
169     of the two mucin-binding domains (referred to as M1M2-less and partial M1M2 RbmC,
170     respectively). Most of the M1M2-less RbmC (59%) and partial M1M2 RbmC (85%) proteins
171     were found to have signal peptides, likely still functioning as an intact protein despite losing
172     domains. Five genes from *V. alfacsensis* and *V. sp002608565* genomes represent the third RbmC
173     variant, which show an averaged ~43% and ~52% similarity to the standard *rbmC* and *bap1*,
174     respectively. This variant has β-propeller and β-helix domains, with corresponding genes located
175     in positions typically associated with *rbmC* genes in the biofilm matrix clusters. Therefore, they
176     are labeled as "*rbmC* w/ β-helix" (Fig. 1B). On the other hand, *bap1* genes are exclusively found
177     in *V. cholerae* and its closely related species within Clade B. Upon examining the neighboring
178     genes of *bap1*, we identified a duplicate *bap1* gene, that encodes a Bap1 variant, directly
179     adjacent to the standard *bap1* (Fig.1B). It shares all of the domains with standard Bap1 protein
180     but lacks the 57aa loop in the β-prism domain and is therefore referred to as loop-less Bap1 (16).
181     Taken together, we identified a total of six structural groups representing different protein
182     variants: RbmC with β-helix, M1M2-less RbmC, partial M1M2 RbmC, standard RbmC, standard
183     Bap1 and loop-less Bap1 (Fig.2A).

Next, after sequence redundancy removal, a codon-based phylogenetic tree was constructed. The phylogeny indicates that the RbmC and Bap1 form two distinct clades, and the long branch connecting them suggests their distant divergence. Protein sequences from the same structural group typically cluster together, although there are exceptions. For instance, a group of genes encoding M1M2-less RbmC is exclusively found in *V. cholerae* and nested within the largest standard RbmC clade, while genes for loop-less Bap1 fall into a subclade within the standard Bap1 clade (Fig.2A). Taking this phylogenetic information into consideration, we have further divided all of the protein sequences into eight protein groups: RbmC with β-helix, M1M2-less RbmC, M1M2-less RbmC in *V. cholerae*, partial M1M2 RbmC, RbmC clade 1, RbmC clade 2, Bap1 clade, and loop-less Bap1 (see gene group cartoon illustrations in Fig.2A).

We mapped these protein groups onto the *Vibrio* species tree tips to infer their evolutionary events. The eight protein groups demonstrated distinct patterns between Clades A and B (Fig.2B). Genes encoding RbmC variants are observed across the species in Clade A, but no Bap1 encoded genes are found. We also observed that RbmC has undergone a series of alterations in the M1M2 domains with a β-helix domain replacing the original M1M2 and β-prisms domains in Clade A. Genes encoding standard RbmC are prevalent in Clade B, in contrast to their restricted presence in a subclade of Clade A (see nodes labeled as "standard RbmC in Clade A" and "standard RbmC in Clade B" in Fig.2B). Genes encoding Bap1 are also found exclusively in Clade B, suggesting that Bap1 genes originated at the ancestral node of this clade (see node labeled as "Origination of Bap1" in Fig.2B). The phylogenetic analysis of the β-propeller domains suggests that Bap1 may have diverged from the ancestor of standard RbmC in both Clade A and Clade B (Supplementary Figure 7). It has been reported that the sequence of Bap1's β-prism diverges from the β-prisms in RbmC (17), and our analysis further shows that Bap1's β-prism domains are closer to RbmC's first β-prism domain (β-prism C1) than to the second (β-prism C2), sharing the most recent common ancestor with RbmC's first β-prism domains exclusively in Clade A (Supplementary Figure 8). This observation aligns with previous findings (14, 16). In addition, the genes encoding loop-less Bap1 are likely to originate from a *V. cholerae* lineage within Clade B (see node labels as "Origination of loop-less Bap1 within *V. cholerae*" in Fig.2B).

A horizontal gene transfer event (HGT) of genes encoding M1M2-less RbmC was observed from *V. cortegadensis* species in Clade B to *V. aestuarianus* species in Clade A. We inferred this to be a result of horizontal gene transfer (see yellow dashed line in Fig.2B) because the genes encoding M1M2-less RbmC, while phylogenetically closest (Fig.2A), are found in two distantly related species in the *Vibrio* species tree (Fig.2B). Interestingly, the biofilm matrix clusters in the genomes of these two species are also highly similar and only slightly differ in the direction and location of the *rbmABC* genes (Fig.1B). In terms of the absence of M1M2 domains in RbmC proteins within *Vibrio cholerae* Clade 1, it is likely the result of a domain loss event in the standard RbmC proteins. This is supported by them forming into a distinct subclade within the standard RbmC Clade 1 on the gene tree (Fig. 2A).

**Loop-less Bap1 proteins are predominantly found in two distant subspecies clades of *V.***

*cholerae*

The standard Bap1 protein and the loop-less variant are predicted to be highly similar in both structures (TM-score=0.8020) and sequences (identity=78.5%) (Supplementary Figure 5E-F). In addition, a 22aa signal peptide was predicted at the N-terminus of loop-less Bap1, which differs in sequence pattern and peptide length from that of the standard Bap1, whose signal peptide is 26aa (Fig.2A). Therefore, despite the absence of a loop and the likely loss of adhesion ability (16), the presence of a signal peptide in the loop-less Bap1 suggests that the protein is still likely to be expressed and secreted.

We next examined the distribution of loop-less Bap1 in the *V. cholerae* subspecies phylogenetic tree. The phylogeny reveals that *V. cholerae* is divided into seven distinct subspecies clades, with loop-less Bap1-encoding genes predominantly enriched in Clades 2 and 3, and a few scattered in Clade 5 (Fig.3). Clades 2 and 3 are phylogenetically distant (Fig.3A), suggesting that the predominant presence of loop-less Bap1 in these clades may reflect selective pressure acting on the protein. While all strains in Clade 2 are classified as *Vibrio cholerae* by GTDB, 38.75% are identified as *Vibrio paracholerae* according to NCBI taxonomy (noting that GTDB does not recognize *Vibrio paracholerae* as a species), indicating its close relationship to *V. paracholerae*. In Clade 3, no clear taxonomic patterns were observed, though the strain *Vibrio albensis* ATCC 14547 (reclassified as *V. cholerae* species in GTDB) is included. Clade 5 is noteworthy, as it includes strains associated with 7PET (7[th] pandemic El Tor) or a putative 7PET lineage (Supplementary Figure 9 and Supplementary Table 4).

To further investigate the genomic traits and potential roles in metabolic pathways of strains harboring loop-less Bap1, we conducted a comparative genomic analysis using Evolink (33) (see details in Methods), which identified gene groups predominantly present and absent in loop-less Bap1-positive strains, referred to as positively and negatively associated gene groups, respectively. We identified five positively and seven negatively associated gene groups (Fig.3B and Supplementary Table 5). Among the positively associated groups, group_3468 is annotated as putative diguanylate cyclase (DGC) with a GGDEF domain and located close to methyl-accepting chemotaxis-related proteins (Supplementary Figure 10A). Among the negatively associated groups, group_1552 encodes an HlyD family secretion protein and is part of *ybhGFSR* system with other two negatively associated gene groups, *ybhF* and *ybhS* (34). Group_2125 contains a methyl-accepting chemotaxis protein signaling domain and is adjacent to a gene encoding a chitinase (*chiA*) (Supplementary Figure 10B). Despite having opposing associations, group_3045 and group_971, are both predicted to function as histidine kinases involved in signal transduction and positioned near group_525, which is annotated as a c-di-GMP phosphodiesterase (Supplementary Figure 10C). These associations highlight potential regulatory pathways and signaling mechanisms that may influence biofilm formation in loop-less Bap1 positive strains in *V. cholerae*.

**RbmB is evolutionarily related to Vibrio prophage pectin lyase-like tail proteins**

We further studied *rbmB* due to its key role in VPS degradation, which regulates biofilm dispersal and cell detachment in *Vibrio cholerae* (11, 19, 20, 35). Its high mutation frequency (0.0604) among all *rbm* genes also suggests a strong positive selection pressure on it (Supplementary Figure 4), highlighting its adaptive significance in promoting biofilm dispersal and its potential as a target for infection control.

By integrating both gene synteny and structural information, we confidently identified 1,760 *rbmB* genes (see details in Methods). We also identified 7,532 genes encoding the pectin lyase-like domain across the Vibrio genus. However, *rbmB* genes account for only 23.4% of these, raising our curiosity about the source and relationships of *rbmB* with other genes. Particularly, given the well-documented role of pectin lyase-like domains in breaking down polysaccharides (36) and their presence in certain phage tail depolymerases that facilitate biofilm degradation (37–39), we explored the possibility that RbmB is evolutionarily related to Vibriophage proteins. To address the abovementioned questions, we constructed a gene tree for Vibrio proteins predicted to have the single-stranded right-handed β-helix (RBH)/pectin lyase-like domains (Fig.4A and Supplementary Figure 11). We observed that over half of the genes (56.1%) are unidentified non-RbmB-encoded genes, and 28.2% are putative pectate lyases. The third largest gene group comprises RbmB-encoded genes (N=319), forming a monophyly in the gene tree (highlighted in red in Fig.4A). The top five species to which these genes belong are *V. cholerae* (N=225), *V. mimicus* (N=20), *V. coralliilyticus* (N=19), *V. metoecus* (N=15) and *V. anguillarum* (N=12) species. Genes in this group have a median length of 408 amino acids and possess signal peptides. This group is closely related to a sister group consisting of 21 non-RbmB-encoded genes (highlighted in yellow in Fig.4A). Together, the two groups are part of a larger clade that includes a large outgroup of 143 non-RbmB-encoded genes (highlighted in blue in Fig.4A). Both groups of 21- and 143-non-RbmB-encoded genes exhibit high structural similarity and moderate sequence similarity to those of the RbmB group, suggesting their close evolutionary relationship (Fig.4B-C).

The 21 non-RbmB encoded genes belong to *V. cholerae* (N=9), *V. anguillarum* (N=6), *V. hepatarius* (N=2), *V. hepatarius_A* (N=2) and *V. mimicus* (N=2) species, with a median gene length of 374 amino acids and possessing signal peptides. Thirteen out of the 21 genomes containing these genes also host confidently curated *rbmB* genes, located hundreds of genes away, and all of these genomes additionally contains *rbmC* genes. Taken together, we believe that these genes encode secretory proteins that are functionally different from the real *rbmB* and are named *rbmB*-like genes in this study. The 21 genes likely play distinct roles across different species due to their involvement in varying gene contexts (see *rbmB*-like genes in Fig.1B and Supplementary Table 6). Interestingly, while *V. hepatarius* and *V. hepatarius_A* species lack true *rbmB* genes, they possess putative polysaccharide lyases with β-jelly roll domains located near *vps*-2 similar genes, which might serve as RbmB alternatives for polysaccharide degradation or biofilm dispersal (Supplementary Figure 2 and Supplementary Table 6).

300    On the other hand, the majority of the 143 non-RbmB encoded genes are from *V. cholerae*
301    (N=124), while the remaining are from *V. mimicus* (N=8), *V. anguillarum* (N=6), *V. metoecus*
302    (N=4) and *V. sp000176715* (N=1) species, with a median gene length of 834 amino acids and
303    lacking signal peptides. One hundred and twenty-six of the 143 genomes containing these genes
304    possess confidently curated *rbmB* genes, which are far from these genes, and all the genomes,
305    except for one, also host *rbmC* or *bap1* genes. Strikingly, we found that 142 of 143 the genes are
306    in the prophage regions. For the only gene not detected in any prophage region in the same
307    contig, it is likely due to the fact that it is the sole gene in the contig that is a relatively short
308    contig that is only 2,667 base pairs long. Gene synteny analysis demonstrated the similarity in
309    the locations of the genes in the 15 representative prophage genomes, where they are situated
310    between two head and packing function-related genes and close to a tail protein (Fig.4D). In
311    addition, BLASTp results showed that all of the 143 genes' best hits (41) share around 30%
312    identity with the tail fiber protein in Vibrio phage vB_VchM_Kuja (GeneBank accession:
313    MN718199) when queried against the Infrastructure for a PHAge Reference Database
314    (INPHARED, accessed on August 15[th], 2024) (42), suggesting these genes may also function as
315    part of the phage tail fibers (Supplementary Table 7). Based on the phylogenetic relationships
316    between RbmB, RbmB-like, and prophage pectin lyase-like proteins, we infer that they are
317    derived from a common ancestor, with the prophage proteins diverging before the split of the
318    RbmB and RbmB-like proteins. Overall, our finding marks the first time that RbmB has been
319    demonstrated to be evolutionarily related to Vibriophage pectin lyase-like tail proteins, thus
320    expanding our understanding of their genetic and functional connections.

## Discussion

322    Bacterial biofilms play a vital role as a lifestyle niche for bacteria in natural environments. They
323    also represent a significant health hazard due to their contribution to persistent infections and the
324    contamination of medical equipment (43–46). Despite their importance in bacterial survival and
325    the challenges they pose in clinical settings, the organization and evolution of the genes encoding
326    the components in biofilm-related clusters have not been extensively studied. A deeper genomic
327    and phylogenetic understanding of these clusters and genes is crucial for the development of
328    innovative genetic engineering strategies that target biofilm-surface interactions and offer
329    alternatives to antibiotic treatments. In this study, using *Vibrio cholerae*—the causative agent of
330    pandemic cholera and a model organism for biofilm studies (18, 47) as well as other related
331    species in the *Vibrio* genus as examples, we propose a framework that integrates comparative
332    genomics, phylogeny, gene synteny analysis and structure prediction to thoroughly characterize
333    biofilm matrix clusters and related proteins, a methodology that can be extended to the study of
334    the biofilm associated clusters and proteins in other bacterial species including important
335    pathogens. This approach has also allowed us to identify domain and modular changes in
336    proteins across their evolutionary timelines, revealing the commonality of domain alterations in
337    *Vibrio* biofilm matrix proteins and their potential implications for biofilm development.

338    As an alternative to combating antibiotic resistance and biofilm formation in Vibrio pathogens,
339    phage therapies are increasingly attracting attention. Notably, phage host-receptor binding

proteins, typically tail fibers or tailspikes, are recognized for their ability to cleave polysaccharides (48–52). Concurrently, *rbmB* genes, encoding RbmB proteins involved in biofilm disassembly, demonstrate significant potential for controlling biofilms and potentially serve as a promising approach to combat Vibrio infections. Interestingly, RbmB proteins and phage tail proteins both feature a common domain—the single-stranded RBH/pectin lyase-like domain—suggesting a potential functional link. However, the evolutionary relationship between these proteins has remained elusive. Here, we reveal that RbmB proteins, along with a group of RbmB-like proteins, share a more recent common ancestor with prophage pectin lyase-like tail proteins than with other pectin lyase-like domain-containing proteins. This suggests that phage tail proteins may be distant homologs of biofilm dispersal proteins, not only highlighting the involvement of phages in biofilm-associated protein evolution but also offering further evidence of phages mimicking host functions to circumvent bacterial defenses. More importantly, the comprehensive annotation of RbmB in *Vibrio* species, combined with insights into Vibrio prophage pectin lyase-like tail proteins, establishes a foundation for a potential biofilm degrader pool. These resources could pave the way for the development of novel protein-based therapies to effectively and precisely target biofilms in emerging Vibrio pathogens.

Our findings clarify many aspects of the *Vibrio* biofilm matrix cluster while also raising new questions. Although we have conducted a comprehensive search for the cluster in the existing genomes across the *Vibrio* genus, it is important to note that this biofilm-associated cluster is VPS-dependent. For *Vibrio* genomes lacking both the *vps*-1 and *vps*-2 genes, it is highly likely that biofilm formation via the clusters curated in this study is not feasible, and these organisms may instead rely on other VPS-independent mechanisms, such as the *syp* loci in *V. parahaemolyticus* and *V. vulnificus* (31). In genomes containing only *vps*-2 similar genes but lacking *vps*-1 genes, it is plausible that the *vps*-2 similar genes are instead integral part of alternative gene operons, such as the *cps* locus in *V. parahaemolyticus* and the *wcr* locus in *V. vulnificus* (31). Therefore, additional VPS-independent biofilm-associated clusters remain to be explored and annotated. For the associated genes identified in loop-less Bap1-positive strains, further experimental validation is required to determine whether they function cooperatively and how they influence bacterial biofilms and behaviors. It is also interesting to explore whether there are polysaccharide lyases or glycosidic hydrolases, aside from RbmB, that could help bacterial cells escape from the biofilm during dispersal. For instance, while RbmB-like proteins are present in *V. hepatarius_A* and *V. hepatarius*, their effectiveness in biofilm disassembly is questionable due to their remote location from other *vps* and *rbm* genes. Instead, polysaccharide lyases containing the β-jelly roll domain may assume this role. It would also be intriguing to uncover how variations in RbmC and Bap1 influence biofilm assembly and to determine the extent to which changes in a single domain or module affect Vibrio phenotypes. We anticipate that these unresolved questions will be addressed through more detailed genomic annotations and experimental studies in the future.

## Methods

### Curation of the biofilm matrix cluster

We downloaded 6,121 genomes classified by GTDB r214 (Genome Taxonomy Database) (25) as Vibrio and Vibrio_A species from NCBI assembly database (53) (accessed on February 18[th], 2024) (Supplementary Data 1). Genomes were annotated by Prokka v1.14.6 (54) with default parameters. KofamScan (https://github.com/takaram/kofam_scan) (55) and InterProScan v5.63-95.0 (56) (with options "-t p -iprlookup --goterms --pathways" and chunksize of 400) were applied to assign KEGG ortholog and predict domains for the genes with default parameters. These genomes along with their gene protein files (.faa), annotation files (.gff) and kofam annotation files (.kofam.tsv) were used as inputs for ProkFunFind (https://github.com/nlm-irp-jianglab/ProkFunFind) (57) to detect potential biofilm matrix clusters. To prepare the queries for the biofilm matrix protein encoded genes, we have collected a set of KEGG orthologs (i.e. KOfam) covering all *vps* genes as well as the *rbmA* gene from the Kyoto Encyclopedia of Genes and Genomes (KEGG) database (https://www.genome.jp/kegg/) (58). We have also composed a hmm profile for all the *rbm* genes. Any clusters of genes containing more than four of the *vps* or *rbm* genes (with option cluster_min_samples=4) and with a gene neighborhood radius of 18 (with option cluster_eps=18) were assigned a cluster ID as a potential biofilm matrix-associated cluster by ProkFunFind. The 18-gene threshold was determined based on the sum of 12 key *vps* genes and 6 *rbm* genes. The *rbmA*, *rbmB*, *rbmC* and *bap1* as well *vpsE* and *vpsF* genes in an output gene annotation file (.gff) was further recognized and curated in the following section, to generate a refined gene annotation file. The configuration file for ProkFunFind, KOFam list and hmm profile files are provided at https://github.com/nlm-irp-jianglab/ProkFunFind and https://zenodo.org/doi/10.5281/zenodo.11509588. The refined gene annotation output obtained from ProkFunFind is available in Supplementary Data 2.

### Curation and classification of the biofilm matrix proteins RbmC and Bap1

Since Bap1 shares over 40% sequence identity with RbmC, traditional sequence-based computational approaches often perform poorly to distinguish them. Furthermore, these two proteins are usually annotated as hemolysin-like proteins by the NCBI genome annotation pipeline, yet they only share less than 40% identity in the single β-prism domain with hemolysins. Another example lies in the initial scanning of ProkFunFind where both *rbmC* and *bap1* genes have been identified as *rbmC* using hmm profile-based search. Nevertheless, RbmC and Bap1 consist of well-studied domains, which inspires us to leverage structural information to distinguish them. First, 4,066 potential RbmC and Bap1 encoded sequences were obtained by querying WP_000200580.1 (RbmC) and WP_001881639.1 (Bap1) against all protein sequences in *Vibrio* genomes using BLASTp v2.15.0+ (41), with a criteria of > 40% identity, > 250 bit score, and > 200 amino acids in aligned length. Next, to better perform a multiple sequence alignment (MSA), after removing sequence redundancy we excluded the five RbmCs with β-helix encoded genes and only selected high-quality RbmC and Bap1 encoded genes. High-

417 quality genes are those with ≥ 80% identity with a Bap1 query and ranging from 650-700aa in
418 length or with ≥ 80% identity with a RbmC query and ranging from 950-1000aa in length, both
419 with bit scores > 900, while the remaining are classified as low-quality genes. We applied
420 MAFFT v7.475 (59) to align high-quality protein sequences with options "--maxiterate 1000 --
421 localpair" and aligned low-quality protein sequences by adding them to the previously aligned
422 high-quality genes using MAFFT with option "-add". The aligned protein sequences were
423 mapped back to the nucleotide sequences to align by codons using PAL2NAL v14 (60). Finally,
424 a codon-based phylogenetic tree was built with the aligned nucleotide sequences using RAxML
425 v8.2.12 (61) by providing a partition file ("-m GTRGAMMA -q dna12_3.partition.txt"), based
426 on which the encoded genes were initially classified as RbmC or Bap1. The detailed structural
427 classification was performed according to the presence and absence of domains in both
428 sequences and structures (Supplementary Data 3-4). The domain boundaries were manually
429 determined by investigating the MSA in Geneious Prime v2023.1.2 (https://www.geneious.com)
430 and double checked with the predicted structures obtained from ESMfold v2.0.0 (62)
431 (Supplementary Data 5). All gene syntenies were annotated using Clinker v0.0.28 (63).

**Curation of RbmB, RbmA, VpsE and VspF proteins**

433 We composed a confident set of *rbmB* genes by first including any genes within a seven-gene
434 distance of either a curated *rbmC* or a putative *rbmA* gene that possess a single-stranded RBH
435 domain (SUPERFAMILY: SSF51126) or are annotated as *rbmB* by a hidden Markov model
436 (HMM) search. The gene distance threshold of seven was determined based on our observation
437 of the maximum number of genes located between *rbmB* and *rbmA* in the current data. Since
438 *rbmA* genes haven't been thoroughly curated, the neighboring *vps* and *rbm* genes of identified
439 *rbmB* genes adjacent only to a putative *rbmA* gene were manually reviewed to determine if they
440 are real *rbmB* genes. Additionally, ten *rbmB* genes were added to the set because they share over
441 60% sequence identity and cover more than 90% of the alignment with *rbmB* genes in the
442 confident set. The gene context and the presence of *rbmC* in the same genomes were examined to
443 support the likelihood that these genes are real *rbmB* genes but are not connected to other *rbm*
444 genes due to poor genome assembly and sequencing quality.

445 Likewise, we curated genes as *rbmA* genes if they are within an eight-gene distance of either a
446 curated *rbmB* or a curated *rbmC* gene, as confirmed in previous sections, that possess two
447 fibronectin type III domains (Gene3D: 2.60.40.3880) or are annotated as *rbmA* by hidden
448 Markov model (HMM) search. The gene distance threshold of eight was determined based on
449 our observation of the maximum number of genes located between *rbmA* and *rbmC* in the
450 current data. For genes located distantly from any *rbmB* or *rbmC* genes but having two
451 fibronectin type III domains, we only included them to the *rbmA* gene set if they, as well as the
452 *rbmB* or *rbmC* genes in the same genomes, are on the edge of contigs, indicating a break in the
453 contig. Regarding genes possessing fewer than two fibronectin type III domains but close to a
454 *rbmB* or *rbmC*, we annotated them as *rbmA* only if they are split into multiple smaller genes or
455 fragmented due to poor genome assembly.

456  We have cautiously annotated *vpsE* and *vpsF*, as they encode the Wzy-polymerase (VpsE) and
457  Wzx-flippase (VpsF) in the *vps*-1 cluster (64), indicating their important roles in the Wzy/Wzx-
458  dependent VPS synthesis pathway. Any genes within a *vps* gene context that are predicted to be
459  polysaccharide biosynthesis proteins (Pfam: PF13440) and have a polysaccharide biosynthesis
460  C-terminal domain (Pfam: PF14667) or are identified as VpsF family polysaccharide
461  biosynthesis proteins (NCBIfam: NF038256), are regarded as *vpsE* or *vpsF*, respectively. Split
462  and fragmented genes, which only have part or none of the domains, were manually annotated
463  and added if they are close to a well-annotated *vpsF*/*vpsE*.

464  The gene sequences and typing information in this section are provided as Supplementary Data
465  6-9.

### Calculating mutation frequency

467  Given a MSA of nucleotide sequences with N sequences and L positions, the mutations at
468  position *i* are:

$$M_i = \sum_L Number\ of\ sequences\ where\ residue \neq residue\ in\ the\ concensus\ sequence$$

469  This includes any substitution or gap ("-") that is not the same as the reference residue.

470  The mutation frequency is defined as the proportion of mutations relative to all possible
471  positions:

$$Mutation\ Frequency = \frac{\sum_i M_i}{N \times L}$$

### Selection of *Vibrio* species representative genomes

473  We didn't simply use the GTDB representative genomes for the 210 *Vibrio* species in this study.
474  Although the representative genomes generally have high completeness and low contamination,
475  they might have fragmented biofilm matrix clusters and don't necessarily have the matrix
476  proteins due to genome assembly issues. To take this into consideration, we developed a strategy
477  to pick representative genomes which have maximally reflected the biofilm matrix cluster status
478  at the *Vibrio* species levels. For the 23 species whose genomes possess *rbmC* and/or *bap1* genes,
479  we manually selected the representative genomes that have the most intact biofilm matrix
480  proteins as well as the untruncated RbmC/Bap1 proteins and are representative of the gene
481  synteny of the biofilm matrix cluster in the species. For 73 species in which no biofilm matrix
482  cluster associated proteins was detected, their GTDB representative genomes were used. For the
483  remaining 114 species, 76 of them have multiple genomes. We ranked the genomes in each
484  species higher if they have 1) fewer contigs, implying they have less fragmented contigs, 2) more
485  key *vps*-1 and *vps*-2 genes in the same gene cluster, and 3) more curated *rbm* or *bap1* genes. The
486  genomes meeting these criteria best were selected as the representatives, while the genomes in
487  the 38 single-genome species were picked as species representatives. The final 216

488 representative genomes for *Vibrio* species and *V. cholerae* subspecies are provided as
489 Supplementary Data 10.

**Pan-genome analysis of *Vibrio cholerae***

491 A total of 194 core genes were detected and aligned in 1663 *V. cholerae* genomes by pan-
492 genome analysis using the Roary v3.13.0 with options "-i 90 -cd 90 -g 500000 -s -e --mafft" (32).
493 The core gene alignment of a subset of 273 representative genomes with completeness > 90%
494 and contamination < 5% was leveraged to build a phylogenomic tree using FastTree v2.1.11 with
495 default options (65) (Supplementary Data 11). The seven clade representative genomes within V.
496 cholerae species, which have intact biofilm matrix clusters and rbmC/bap1 genes as well as
497 fewer contigs and larger genome lengths, were manually picked for the corresponding clades.
498 The 7PET (7th pandemic El Tor) and putative 7PET lineages were identified by calculating the
499 genomic distance and detecting marker genes with the refence (N16961) using "is-it-7pet" tool
500 (https://github.com/amberjoybarton/is-it-7pet).

**Construction of phylogenomic *Vibrio* species tree**

502 We applied PIRATE v1.0.5 to the 209 *Vibrio* species representative genomes (excluding *V.*
503 *cholerae*) and seven *V. cholerae* subspecies representative genomes to obtain genus-wise marker
504 genes (with options "-k '--diamond'") (66). PIRATE can rapidly create pangenomes from coding
505 sequences over a wide range of amino acid identity thresholds, thus recognizing the most robust
506 set of core genes. The core gene nucleotide alignment provided by PIRATE was used to build
507 the *Vibrio* species tree using FastTree v2.1.11 with options "-gtr -nt" (Supplementary Data 12).
508 According to GTDB, the *Vibrio_A* genus is more distantly related to the *Vibrio* genus and can
509 serve as a reference group for determining the evolutionary relationships within the *Vibrio* genus.
510 Consequently, the genome of *Vibrio_A stylophorae* was selected as the outgroup to root the tree.

**Identification of loop-less Bap1 positive strains associated gene groups**

512 Given the *V. cholerae* phylogenetic tree, the presence and absence of the gene groups defined by
513 Roary (Supplementary Data 13) and the existence of loop-less Bap1 in genomes, we ran Evolink
514 with default parameters (https://github.com/nlm-irp-jianglab/Evolink) to find five positively and
515 seven negatively associated gene groups related to loop-less Bap1 presence. Evolink is a method
516 for rapid identification of associated genotypes provided a trait of interest and uses phylogenetic
517 approaches to adjust for the population structure in microbial data (33). To confirm that the
518 associated genes identified could be reproduced using alternative methods, we repeated the
519 association analysis with Pyseer using a linear mixed model (67). Pyseer identified 592
520 significantly associated genes (adjusted p-value < 0.05), 11 of which overlapped with the 12
521 genes identified by Evolink. The only exception was group_2326, which had an adjusted p-value
522 of 0.108.

**Signal peptide detection**

Signal peptides were predicted for RbmC and Bap1-related proteins using SignalP6.0 server (https://services.healthtech.dtu.dk/services/SignalP-6.0/) (68). The signal peptides were aligned with MAFFT v7.475 (59) and visualized as sequence logo using WebLogo server (https://weblogo.berkeley.edu/logo.cgi) (69) (Supplementary Data 14).

**Construction of gene and domain trees**

After removing sequence redundancy, single-stranded RBH domain containing protein sequences were aligned using MAFFT-DASH (70) to take structural alignment into consideration. The multiple sequence alignment was next trimmed using TrimAl v1.2rev59 (71) with the -gt 0.2 option to obtain cleaner alignment and used to reconstruct their phylogeny using FastTree v2.1.11 with default options (65).

The β-propeller and β-prism domains sequences were extracted based on domain segmentation of RbmC and Bap1 proteins. The alignment using MAFFT v7.475 (59) were used to build trees using FastTree v2.1.11 with default options (65). All trees were visualized and annotated with iTOL v6 server (https://itol.embl.de/) (72).

The tree files were provided as Supplementary Data 15-17.

**Prophage regions identification**

Prophage regions in genomes were detected using VirSorter v2.2.4 (73) with options "--min-length 1000" (Supplementary Data 18). Phage genes within the determined prophage regions were annotated and categorized using Pharokka v1.3.2 (74).

# Data and code availability

The data underlying this article can be accessed through Zenodo (https://zenodo.org/doi/10.5281/zenodo.11509588). All scripts utilized throughout the publication can be accessed through the main branch on the GitHub repository (https://github.com/YiyanYang0728/Vibrio_biofilm_matrix_cluster).

# Acknowledgments

## Conflicts of interest

556   The authors declare that there are no conflicts of interest.

557

## Reference

559   1.   Charles RC, Ryan ET. 2011. Cholera in the 21st century: Current Opinion in Infectious Diseases
560        24:472–477.

561   2.   Donlan RM, Costerton JW. 2002. Biofilms: Survival Mechanisms of Clinically Relevant
562        Microorganisms. Clin Microbiol Rev 15:167–193.

563   3.   Colwell RR, Huq A, Islam MS, Aziz KMA, Yunus M, Khan NH, Mahmud A, Sack RB, Nair GB,
564        Chakraborty J, Sack DA, Russek-Cohen E. 2003. Reduction of cholera in Bangladeshi villages by
565        simple filtration. Proc Natl Acad Sci USA 100:1051–1055.

566   4.   Gupta P, Mankere B, Chekkoora Keloth S, Tuteja U, Pandey P, Chelvam KT. 2018. Increased
567        antibiotic resistance exhibited by the biofilm of Vibrio cholerae O139. Journal of Antimicrobial
568        Chemotherapy 73:1841–1847.

569   5.   Matz C, McDougald D, Moreno AM, Yung PY, Yildiz FH, Kjelleberg S. 2005. Biofilm formation
570        and phenotypic variation enhance predation-driven persistence of *Vibrio cholerae*. Proc Natl Acad
571        Sci USA 102:16819–16824.

572   6.   Beyhan S, Yildiz FH. 2007. Smooth to rugose phase variation in *Vibrio cholerae* can be mediated by
573        a single nucleotide change that targets c-di-GMP signalling pathway. Molecular Microbiology
574        63:995–1007.

575   7.   Yildiz FH, Schoolnik GK. 1999. *Vibrio cholerae* O1 El Tor: Identification of a gene cluster required
576        for the rugose colony type, exopolysaccharide production, chlorine resistance, and biofilm
577        formation. Proc Natl Acad Sci USA 96:4028–4033.

578   8.   Yildiz F, Fong J, Sadovskaya I, Grard T, Vinogradov E. 2014. Structural Characterization of the
579        Extracellular Polysaccharide from Vibrio cholerae O1 El-Tor. PLoS ONE 9:e86751.

580   9.   Fong JCN, Syed KA, Klose KE, Yildiz FH. 2010. Role of Vibrio polysaccharide (vps) genes in VPS
581        production, biofilm formation and Vibrio cholerae pathogenesis. Microbiology 156:2757–2769.

582   10.  Fong JCN, Karplus K, Schoolnik GK, Yildiz FH. 2006. Identification and Characterization of
583        RbmA, a Novel Protein Required for the Development of Rugose Colony Morphology and Biofilm
584        Structure in *Vibrio cholerae*. J Bacteriol 188:1049–1059.

585   11.  Fong JCN, Yildiz FH. 2007. The *rbmBCDEF* Gene Cluster Modulates Development of Rugose
586        Colony Morphology and Biofilm Formation in *Vibrio cholerae*. J Bacteriol 189:2319–2330.

587   12.  Moorthy S, Watnick PI. 2005. Identification of novel stage-specific genetic requirements through
588        whole genome transcription profiling of *Vibrio cholerae* biofilm development. Molecular
589        Microbiology 57:1623–1635.

590   13.  Berk V, Fong JCN, Dempsey GT, Develioglu ON, Zhuang X, Liphardt J, Yildiz FH, Chu S. 2012.
591        Molecular Architecture and Assembly Principles of *Vibrio cholerae* Biofilms. Science 337:236–239.

14. Absalon C, Van Dellen K, Watnick PI. 2011. A Communal Bacterial Adhesin Anchors Biofilm and Bystander Cells to Surfaces. PLoS Pathog 7:e1002210.

15. Maestre-Reyna M, Wu W-J, Wang AH-J. 2013. Structural Insights into RbmA, a Biofilm Scaffolding Protein of V. Cholerae. PLoS ONE 8:e82458.

16. Huang X, Nero T, Weerasekera R, Matej KH, Hinbest A, Jiang Z, Lee RF, Wu L, Chak C, Nijjer J, Gibaldi I, Yang H, Gamble N, Ng W-L, Malaker SA, Sumigray K, Olson R, Yan J. 2023. Vibrio cholerae biofilms use modular adhesins with glycan-targeting and nonspecific surface binding domains for colonization. Nat Commun 14:2104.

17. De S, Kaus K, Sinclair S, Case BC, Olson R. 2018. Structural basis of mammalian glycan targeting by Vibrio cholerae cytolysin and biofilm proteins. PLoS Pathog 14:e1006841.

18. Teschler JK, Zamorano-Sánchez D, Utada AS, Warner CJA, Wong GCL, Linington RG, Yildiz FH. 2015. Living in the matrix: assembly and control of Vibrio cholerae biofilms. Nat Rev Microbiol 13:255–268.

19. Bridges AA, Fei C, Bassler BL. 2020. Identification of signaling pathways, matrix-digestion enzymes, and motility components controlling *Vibrio cholerae* biofilm dispersal. Proc Natl Acad Sci USA 117:32639–32647.

20. Díaz-Pascual F, Hartmann R, Lempp M, Vidakovic L, Song B, Jeckel H, Thormann KM, Yildiz FH, Dunkel J, Link H, Nadell CD, Drescher K. 2019. Breakdown of Vibrio cholerae biofilm architecture induced by antibiotics disrupts community barrier function. Nat Microbiol 4:2136–2145.

21. Lilburn TG, Gu J, Cai H, Wang Y. 2010. Comparative genomics of the family Vibrionaceae reveals the wide distribution of genes encoding virulence-associated proteins. BMC Genomics 11:369.

22. Guo Y, Rowe-Magnus DA. 2011. Overlapping and unique contributions of two conserved polysaccharide loci in governing distinct survival phenotypes in *Vibrio vulnificus*. Environmental Microbiology 13:2888–2990.

23. Chodur DM, Rowe-Magnus DA. 2018. Complex Control of a Genomic Island Governing Biofilm and Rugose Colony Development in Vibrio vulnificus. J Bacteriol 200.

24. Gao C, Garren M, Penn K, Fernandez VI, Seymour JR, Thompson JR, Raina J-B, Stocker R. 2021. Coral mucus rapidly induces chemokinesis and genome-wide transcriptional shifts toward early pathogenesis in a bacterial coral pathogen. The ISME Journal 15:3668–3682.

25. Parks DH, Chuvochina M, Rinke C, Mussig AJ, Chaumeil P-A, Hugenholtz P. 2022. GTDB: an ongoing census of bacterial and archaeal diversity through a phylogenetically consistent, rank normalized and complete genome-based taxonomy. Nucleic Acids Research 50:D785–D794.

26. Lin H, Yu M, Wang X, Zhang X-H. 2018. Comparative genomic analysis reveals the evolution and environmental adaptation strategies of vibrios. BMC Genomics 19:135.

27. Smith AB, Siebeling RJ. 2003. Identification of Genetic Loci Required for Capsular Expression in *Vibrio vulnificus*. Infect Immun 71:1091–1097.

28. Güvener ZT, McCarter LL. 2003. Multiple Regulators Control Capsular Polysaccharide Production in *Vibrio parahaemolyticus*. J Bacteriol 185:5431–5441.

29. Grau BL, Henk MC, Garrison KL, Olivier BJ, Schulz RM, O'Reilly KL, Pettis GS. 2008. Further Characterization of *Vibrio vulnificus* Rugose Variants and Identification of a Capsular and Rugose Exopolysaccharide Gene Cluster. Infect Immun 76:1485–1497.

633 30. Darnell CL, Hussa EA, Visick KL. 2008. The Putative Hybrid Sensor Kinase SypF Coordinates
634    Biofilm Formation in *Vibrio fischeri* by Acting Upstream of Two Response Regulators, SypG and
635    VpsR. J Bacteriol 190:4941–4950.

636 31. Yildiz FH, Visick KL. 2009. Vibrio biofilms: so much the same yet so different. Trends in
637    Microbiology 17:109–118.

638 32. Page AJ, Cummins CA, Hunt M, Wong VK, Reuter S, Holden MTG, Fookes M, Falush D, Keane
639    JA, Parkhill J. 2015. Roary: rapid large-scale prokaryote pan genome analysis. Bioinformatics
640    31:3691–3693.

641 33. Yang Y, Jiang X. 2023. Evolink: a phylogenetic approach for rapid identification of genotype–
642    phenotype associations in large-scale microbial multispecies data. Bioinformatics 39:btad215.

643 34. Feng Z, Liu D, Wang L, Wang Y, Zang Z, Liu Z, Song B, Gu L, Fan Z, Yang S, Chen J, Cui Y. 2020.
644    A Putative Efflux Transporter of the ABC Family, YbhFSR, in Escherichia coli Functions in
645    Tetracycline Efflux and Na+(Li+)/H+ Transport. Front Microbiol 11:556.

646 35. Weerasekera R, Moreau A, Huang X, Nam K-M, Hinbest AJ, Huynh Y, Liu X, Ashwood C, Pepi
647    LE, Paulson E, Cegelski L, Yan J, Olson R. 2024. Vibrio cholerae RbmB is an α-1,4-polysaccharide
648    lyase with biofilm-disrupting activity against Vibrio polysaccharide (VPS). PLoS Pathog
649    20:e1012750.

650 36. Burnim AA, Dufault-Thompson K, Jiang X. 2024. The three-sided right-handed β-helix is a
651    versatile fold for glycan interactions. Glycobiology 34:cwae037.

652 37. Latka A, Drulis-Kawa Z. 2020. Advantages and limitations of microtiter biofilm assays in the model
653    of antibiofilm activity of Klebsiella phage KP34 and its depolymerase. Sci Rep 10:20338.

654 38. Hughes KA, Sutherland IW, Clark J, Jones MV. 1998. Bacteriophage and associated polysaccharide
655    depolymerases – novel tools for study of bacterial biofilms. Journal of Applied Microbiology
656    85:583–590.

657 39. Verma V, Harjai K, Chhibber S. 2010. Structural changes induced by a lytic bacteriophage make
658    ciprofloxacin effective against older biofilm of *Klebsiella pneumoniae*. Biofouling 26:729–737.

659 40. Abramson J, Adler J, Dunger J, Evans R, Green T, Pritzel A, Ronneberger O, Willmore L, Ballard
660    AJ, Bambrick J, Bodenstein SW, Evans DA, Hung C-C, O'Neill M, Reiman D, Tunyasuvunakool K,
661    Wu Z, Žemgulytė A, Arvaniti E, Beattie C, Bertolli O, Bridgland A, Cherepanov A, Congreve M,
662    Cowen-Rivers AI, Cowie A, Figurnov M, Fuchs FB, Gladman H, Jain R, Khan YA, Low CMR,
663    Perlin K, Potapenko A, Savy P, Singh S, Stecula A, Thillaisundaram A, Tong C, Yakneen S, Zhong
664    ED, Zielinski M, Žídek A, Bapst V, Kohli P, Jaderberg M, Hassabis D, Jumper JM. 2024. Accurate
665    structure prediction of biomolecular interactions with AlphaFold 3. Nature
666    https://doi.org/10.1038/s41586-024-07487-w.

667 41. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
668    BLAST+: architecture and applications. BMC Bioinformatics 10:421.

669 42. Cook R, Brown N, Redgwell T, Rihtman B, Barnes M, Clokie M, Stekel DJ, Hobman J, Jones MA,
670    Millard A. 2021. INfrastructure for a PHAge REference Database: Identification of Large-Scale
671    Biases in the Current Collection of Cultured Phage Genomes. PHAGE 2:214–223.

672 43. Donlan RM. 2016. Microbial Biofilms, Second Edition. Emerg Infect Dis 22:1142–1142.

673 44. Hall-Stoodley L, Costerton JW, Stoodley P. 2004. Bacterial biofilms: from the Natural environment
674    to infectious diseases. Nat Rev Microbiol 2:95–108.

45. Costerton JW, Stewart PS, Greenberg EP. 1999. Bacterial Biofilms: A Common Cause of Persistent Infections. Science 284:1318–1322.

46. Flemming H-C, Wingender J, Szewzyk U, Steinberg P, Rice SA, Kjelleberg S. 2016. Biofilms: an emergent form of bacterial life. Nat Rev Microbiol 14:563–575.

47. Nelson EJ, Harris JB, Glenn Morris J, Calderwood SB, Camilli A. 2009. Cholera transmission: the host, pathogen and bacteriophage dynamic. Nat Rev Microbiol 7:693–702.

48. Yen M, Cairns LS, Camilli A. 2017. A cocktail of three virulent bacteriophages prevents Vibrio cholerae infection in animal models. Nat Commun 8:14187.

49. Jensen MA, Faruque SM, Mekalanos JJ, Levin BR. 2006. Modeling the role of bacteriophage in the control of cholera outbreaks. Proc Natl Acad Sci USA 103:4652–4657.

50. Bhandare S, Colom J, Baig A, Ritchie JM, Bukhari H, Shah MA, Sarkar BL, Su J, Wren B, Barrow P, Atterbury RJ. 2019. Reviving Phage Therapy for the Treatment of Cholera. The Journal of Infectious Diseases 219:786–794.

51. Barman RK, Chakrabarti AK, Dutta S. 2022. Screening of Potential Vibrio cholerae Bacteriophages for Cholera Therapy: A Comparative Genomic Approach. Front Microbiol 13:803933.

52. Yang Y, Dufault-Thompson K, Yan W, Cai T, Xie L, Jiang X. 2024. Large-scale genomic survey with deep learning-based method reveals strain-level phage specificity determinants. GigaScience 13:giae017.

53. Kitts PA, Church DM, Thibaud-Nissen F, Choi J, Hem V, Sapojnikov V, Smith RG, Tatusova T, Xiang C, Zherikov A, DiCuccio M, Murphy TD, Pruitt KD, Kimchi A. 2016. Assembly: a resource for assembled genomes at NCBI. Nucleic Acids Res 44:D73–D80.

54. Seemann T. 2014. Prokka: rapid prokaryotic genome annotation. Bioinformatics 30:2068–2069.

55. Aramaki T, Blanc-Mathieu R, Endo H, Ohkubo K, Kanehisa M, Goto S, Ogata H. 2020. KofamKOALA: KEGG Ortholog assignment based on profile HMM and adaptive score threshold. Bioinformatics 36:2251–2252.

56. Jones P, Binns D, Chang H-Y, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell A, Nuka G, Pesseat S, Quinn AF, Sangrador-Vegas A, Scheremetjew M, Yong S-Y, Lopez R, Hunter S. 2014. InterProScan 5: genome-scale protein function classification. Bioinformatics 30:1236–1240.

57. Dufault-Thompson K, Jiang X. 2024. Annotating microbial functions with ProkFunFind. mSystems 9:e00036-24.

58. Kanehisa M, Furumichi M, Tanabe M, Sato Y, Morishima K. 2017. KEGG: new perspectives on genomes, pathways, diseases and drugs. Nucleic Acids Res 45:D353–D361.

59. Katoh K. 2002. MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform. Nucleic Acids Research 30:3059–3066.

60. Suyama M, Torrents D, Bork P. 2006. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. Nucleic Acids Research 34:W609–W612.

61. Stamatakis A. 2006. RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models. Bioinformatics 22:2688–2690.

713 62. Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, Smetanin N, Verkuil R, Kabeli O, Shmueli Y, Dos
714     Santos Costa A, Fazel-Zarandi M, Sercu T, Candido S, Rives A. 2023. Evolutionary-scale prediction
715     of atomic-level protein structure with a language model. Science 379:1123–1130.

716 63. Gilchrist CLM, Chooi Y-H. 2021. clinker & clustermap.js: automatic generation of gene cluster
717     comparison figures. Bioinformatics 37:2473–2475.

718 64. Schwechheimer C, Hebert K, Tripathi S, Singh PK, Floyd KA, Brown ER, Porcella ME, Osorio J,
719     Kiblen JTM, Pagliai FA, Drescher K, Rubin SM, Yildiz FH. 2020. A tyrosine phosphoregulatory
720     system controls exopolysaccharide biosynthesis and biofilm formation in Vibrio cholerae. PLoS
721     Pathog 16:e1008745.

722 65. Price MN, Dehal PS, Arkin AP. 2010. FastTree 2 – Approximately Maximum-Likelihood Trees for
723     Large Alignments. PLoS ONE 5:e9490.

724 66. Bayliss SC, Thorpe HA, Coyle NM, Sheppard SK, Feil EJ. 2019. PIRATE: A fast and scalable
725     pangenomics toolbox for clustering diverged orthologues in bacteria. GigaScience 8:giz119.

726 67. Lees JA, Galardini M, Bentley SD, Weiser JN, Corander J. 2018. pyseer: a comprehensive tool for
727     microbial pangenome-wide association studies. Bioinformatics 34:4310–4312.

728 68. Teufel F, Almagro Armenteros JJ, Johansen AR, Gíslason MH, Pihl SI, Tsirigos KD, Winther O,
729     Brunak S, Von Heijne G, Nielsen H. 2022. SignalP 6.0 predicts all five types of signal peptides
730     using protein language models. Nat Biotechnol 40:1023–1025.

731 69. Crooks GE, Hon G, Chandonia J-M, Brenner SE. 2004. WebLogo: A Sequence Logo Generator.
732     Genome Res 14:1188–1190.

733 70. Rozewicki J, Li S, Amada KM, Standley DM, Katoh K. 2019. MAFFT-DASH: integrated protein
734     sequence and structural alignment. Nucleic Acids Research gkz342.

735 71. Capella-Gutiérrez S, Silla-Martínez JM, Gabaldón T. 2009. trimAl: a tool for automated alignment
736     trimming in large-scale phylogenetic analyses. Bioinformatics 25:1972–1973.

737 72. Letunic I, Bork P. 2024. Interactive Tree of Life (iTOL) v6: recent updates to the phylogenetic tree
738     display and annotation tool. Nucleic Acids Research gkae268.

739 73. Guo J, Bolduc B, Zayed AA, Varsani A, Dominguez-Huerta G, Delmont TO, Pratama AA, Gazitúa
740     MC, Vik D, Sullivan MB, Roux S. 2021. VirSorter2: a multi-classifier, expert-guided approach to
741     detect diverse DNA and RNA viruses. Microbiome 9:37.

742 74. Bouras G, Nepal R, Houtak G, Psaltis AJ, Wormald P-J, Vreugde S. 2023. Pharokka: a fast scalable
743     bacteriophage annotation tool. Bioinformatics 39:btac776.

744

## Figure legends

**Figure 1. The distribution of biofilm matrix clusters across the *Vibrio* genus.** (A) The phylogenomic tree with the presence and absence of important genes in biofilm matrix clusters mapped to tips representing 216 Vibrio (sub)species. The tree was rooted with the representative genome of *Vibrio_A stylophorae* species (NCBI Assembly accession=GCA_921293875.1). (B) Gene syntenies for biofilm matrix clusters in 29 (sub)species that possess biofilm matrix protein encoding genes (*rbmC* and/or *bap1*) are illustrated using the same color palette as in panel A and the phylogenomic tree displayed is a subtree derived from the tree in panel A. The clusters are aligned with each other using the *rbmC* gene as the anchor. Genes that are not concatenated are located on different contigs, whereas genes separated by the "//" symbol are found in the same contig but are hundreds of genes away from each other. The red boxes highlight the proximity of the *vps*-1 and *rbmABC* genes within the genome, while the blue boxes indicate the close genomic location of the *vps*-2 and *rbmDEF* genes. The *rbmE* and *rbmF* genes are combined under the single gene name *rbmEF* due to overlaps in their gene sequences and frequent annotations as a single gene. Similarly, the *vpsC* and *vpsG* genes are merged into one gene name, *vpsCG*, as they both share a highly similar domain. PS: Polysaccharide.

**Figure 2. The gene tree and evolutionary analysis for RbmC and Bap1 proteins.** (A) The gene tree was built with non-redundant codon sequences of 514 RbmC and 483 Bap1 proteins, which is rooted at the midpoint. The outer circle indicates the species of origin, while the inner circle indicates the protein structural features with grey representing truncated proteins. The cartoons at the bottom demonstrate the domain composition for the corresponding structures. Color ranges indicate different protein groups based on both structural features and phylogenetic relationships, whose legend was put under the corresponding structural features. Note that the RbmC with a β-helix domain was omitted from the gene tree due to it causing a poor multiple sequence alignment. The sequence logos for the signal peptides are shown for the Bap1 clade and loop-less Bap1 clade. (B) The distribution of 9 protein groups along the phylogenetic tree suggests evolutionary events for *rbmC* and *bap1* genes. The tree replicates the one in Fig.1B while retaining the outgroup species. The species and protein group colors are consistent with those in panel A.

**Figure 3. Loop-less Bap1 encoded genes are predominantly found in two distant *V. cholerae* clades, which share specific gene groups associated with the presence/absence of the protein.** (A) Unrooted phylogenomic tree of *V. cholerae* species (N=273), with bootstrap values displayed at clade ancestral nodes and nodes representing clade divergence. (B) The phylogenomic tree for *V. cholerae* species was built with protein sequences from the core genes found by Roary (32). The tree was rooted at Clade 1. The presence and absence of RbmC/Bap1 variants (inner circles, using the same palette in Fig.2) and gene groups either positively (dark red) or negatively (dark blue) associated with loop-less Bap1-positive strains (outer circles) are mapped to the tips.
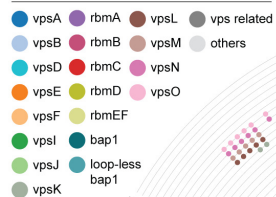
**Figure 4. Single-stranded right-handed β-helix (RBH) domain containing gene tree suggests an association between RbmB and prophage proteins.** (A) The gene tree was built with non-

785 redundant protein sequences containing single-stranded RBH domains (SUPERFAMILY:
786 SSF51126) and was rooted at the midpoint. Encoded proteins are annotated as colored dots at
787 tips. The inner circle represents the associations of the genes with the prophages found in the
788 same contigs, while the outer circle represents the gene lengths. Bootstrap values are shown at
789 three key internal nodes. The color ranges highlight the clades for RbmB encoded genes (red),
790 RbmB-like encoded genes (yellow) and prophage-related genes (blue). (B-C) Pairwise
791 superimposition of predicted protein structures. The structures displayed are for RbmB (colored
792 red, gene accession: GCA_013111535.1_02619), RbmB-like (colored yellow, gene accession:
793 GCA_002284395.1_03257), and prophage proteins (colored blue, gene accession:
794 GCA_002097735.1_02038). The signal peptides were removed from RbmB and RbmB-like
795 proteins and the structures were predicted by AlphaFold3 (40). (D) Gene syntenies for the 15
796 representative prophages that possess single-stranded RBH domain containing genes. Each gene
797 synteny is accompanied by the genome accessions from which the prophage fragment was found.
798 Genes encoding the single-stranded RBH domain are colored red, while other genes are colored
799 according to phage functional categories. AlgG: Mannuronan C5-epimerase; NosD: Putative
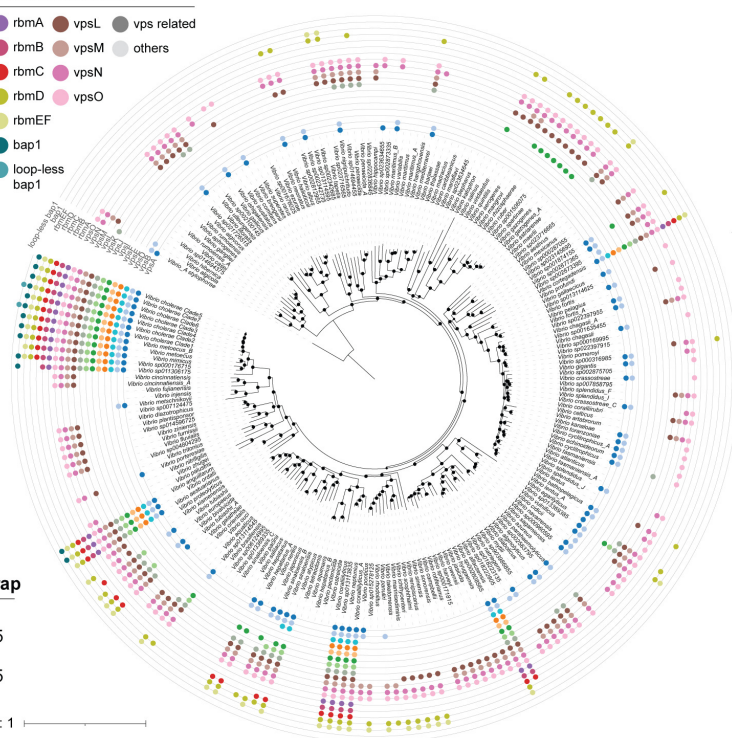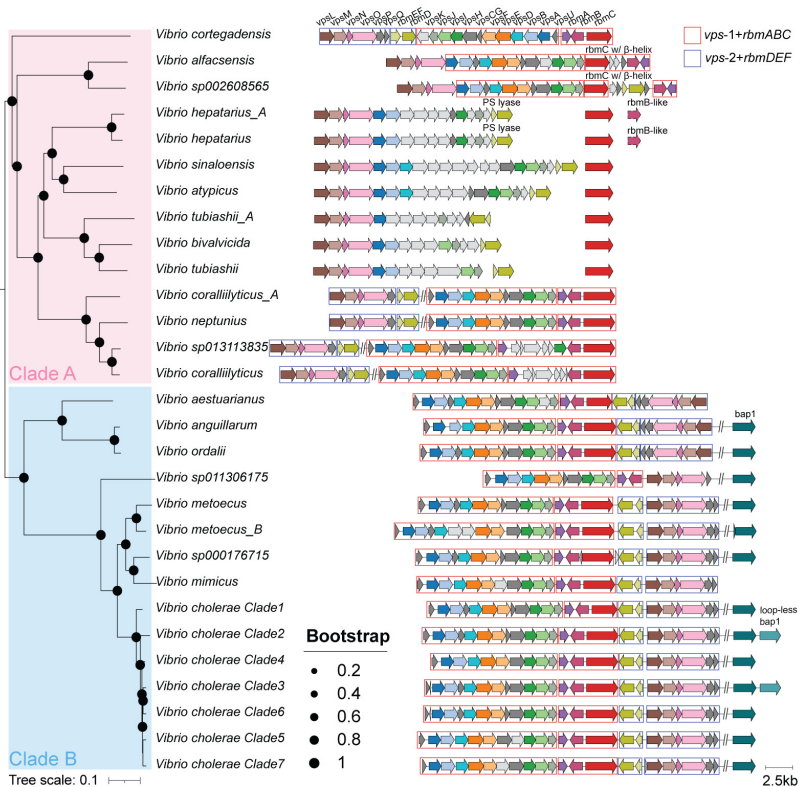800 ABC transporter binding protein.

801