



HAL
open science

Exploring Host-Binding Machineries of Mycobacteriophages with AlphaFold2

Christian Cambillau, Adeline Goulet

► **To cite this version:**

Christian Cambillau, Adeline Goulet. Exploring Host-Binding Machineries of Mycobacteriophages with AlphaFold2. *Journal of Virology*, 2023, 97 (3), pp.01793-22. 10.1128/jvi.01793-22. hal-04274513

HAL Id: hal-04274513

<https://hal.science/hal-04274513>

Submitted on 13 Nov 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Exploring host-binding machineries of mycobacteriophages with AlphaFold2

Christian Cambillau^{a,b} and Adeline Goulet^{c,#}

^aSchool of Microbiology, University College Cork, Cork, Ireland

^bAlphaGraphix, Formiguères, 66210, France

^cLaboratoire d'Ingénierie des Systèmes Macromoléculaires (LISM), Institut de Microbiologie, Bioénergies et Biotechnologie (IM2B), CNRS and Aix-Marseille Université UMR7255, Marseille, France

Running title: Structural insights into mycobacteriophages' tail tips

#Address correspondence to Adeline Goulet, adeline.goulet@univ-amu.fr.

Keywords: bacteriophage, *Mycobacteria*, AlphaFold2, host-binding machineries, carbohydrate-binding module, receptor-binding protein, polyglycine helices

17 **Abstract**

18 Although more than 12,000 bacteriophages infecting mycobacteria (or mycobacteriophages) have been
19 isolated so far, there is a knowledge gap on their structure-function relationships. Here, we have explored
20 the architecture of host-binding machineries from seven representative mycobacteriophages of the
21 *Siphoviridae* family infecting *Mycobacterium smegmatis*, *Mycobacterium abscessus* and *Mycobacterium*
22 *tuberculosis*, using AlphaFold2 (AF2). AF2 enables confident structural analyses of large and flexible
23 biological assemblies resistant to experimental methods, thereby opening new avenues to shed light on
24 phage structure and function. Our results highlight the modularity and structural diversity of siphophage
25 host-binding machineries that recognize host-specific receptors at the onset of viral infection. Interestingly,
26 the studied mycobacteriophages' host-binding machineries present unique features as compared with those
27 of phages infecting other Gram-positive actinobacteria. Although they all assemble the classical Dit (Distal
28 tail), Tal (Tail-associated lysin) and receptor-binding proteins, five of them contain two potential additional
29 adhesion proteins. Moreover, we have identified brush-like domains formed of multiple polyglycine
30 helices, which expose hydrophobic residues, as potential receptor-binding domains. These polyglycine-rich
31 domains, which have been observed in only five native proteins, may be a hallmark of mycobacteriophages'
32 host-binding machineries, and more common in nature than expected. Altogether, the unique composition
33 of mycobacteriophages' host-binding machineries indicate they might have evolved to bind to the peculiar
34 mycobacterial cell envelope rich in polysaccharides and mycolic acids. This work provides a rational
35 framework to efficiently produce recombinant proteins or protein domains and test their host-binding
36 function, and hence to shed light on molecular mechanisms used by mycobacteriophages to infect their
37 host.

38

39 **Importance**

40 Mycobacteria include both saprophytes, such as the model system *Mycobacterium smegmatis*, and
41 pathogens, such as *Mycobacterium tuberculosis* and *Mycobacterium abscessus* that are poorly responsive
42 to antibiotic treatments and pose a global public health problem. Mycobacteriophages have been collected
43 at a very large scale over the last decade, and they have proven to be valuable tools for mycobacteria genetic
44 manipulation, rapid diagnostics and infection treatment. Yet, molecular mechanisms used by
45 mycobacteriophages to infect their host remain poorly understood. Therefore, exploring the structural
46 diversity of mycobacteriophages' host-binding machineries is important not only to better understand viral
47 diversity and bacteriophage-host interactions, but also to rationally develop biotechnological tools. With
48 the powerful protein structure prediction software AlphaFold2, which was publicly released a year ago, it
49 is now possible to get structural and functional insights on such challenging assemblies.

50

51 Introduction

52 Mycobacteria are members of the Gram-positive *Actinobacteria* phylum ubiquitously found in terrestrial
53 and aquatic ecosystems (1, 2). They include both saprophytes, such as the model system *Mycobacterium*
54 *smegmatis*, and pathogens, such as *Mycobacterium tuberculosis* and *Mycobacterium abscessus* that are
55 poorly responsive to antibiotic treatments and pose a global public health problem (3). One of the
56 distinguishing features of *Mycobacteria* is their thick, waxy cell wall rich in mycolic acids that gives them
57 many of their biological properties and acts as a protective barrier against harsh environmental conditions
58 (4, 5). Viruses that infect mycobacteria, or mycobacteriophages, have been collected at a very large scale
59 over the last decade thanks to the work of phage hunters enrolled in integrated research-education programs
60 such as SEA-PHAGES (Science Education Alliance Phage Hunters Advancing Genomics and Evolutionary
61 Sciences) (6). To date, more than 12,000 mycobacteriophages have been identified and over 2,000 genomes
62 have been sequenced, representing more than half of phages and genomic sequences available on
63 Actinobacteriophages (<https://phagesdb.org/>, as for September 2022, (7)). Based on their remarkable
64 genomic diversity, mycobacteriophages are classified into 31 clusters (clusters A-Z and AA-AE), further
65 divided into subclusters, and 7 singletons (phages with no other close relatives) (8). This substantial
66 collection of mycobacteriophages is highly valuable to better understand phage diversity and evolution,
67 phage-host interactions, and to develop phage-based strategies to detect and treat mycobacterial infections
68 (9–13). Yet, molecular mechanisms of mycobacteriophage infections are poorly understood. In particular,
69 the architecture of mycobacteriophages' host-binding machineries and their mode of action to bind to the
70 host cell surface at the onset of viral infection are not known.

71 While only one paper reports low-resolution structural data of the host-binding machinery of the *M.*
72 *abscessus* siphophage Araucaria (14), many structural and functional analyses of host-binding machineries
73 from siphophages infecting other Gram-positive bacteria have revealed molecular bases of their
74 architecture and function (15–24). These machineries assemble an hexameric ring of the Dit protein (distal
75 tail protein) bound to the last ring of major tail protein (MTP) and to a trimer of Tal (tail-associated lysin).
76 Dit proteins are made up of a N-terminal domain, the fold of which is conserved in tail proteins such as
77 MTP, head-to-tail joining proteins, and Dit of *Siphoviridae*, *Myoviridae* and phage tail-like system (23–
78 28), and a C-terminal galectin domain (for an example see Fig. 2 DS6A). The Dit N-terminal domain is
79 referred to as the belt domain because of its canonical belt loop assembling the ring. Tal proteins are
80 composed of a N-terminal structural domain of ~350-400 residues, similar to the myophage T4 gp27 protein
81 (29) (for an example see Fig. 3 Abinghost), that can be followed by a C-terminal extension of 1,000 to
82 2,000 residues (18, 19, 30). This extension is believed to play a role in cell wall
83 polysaccharide/peptidoglycan degradation, e.g. the *Lactococcus lactis* P335 phage TP901-1 (31), or host
84 binding, e.g. the *Bacillus subtilis* and *Escherichia coli* phages SPP1 and T5 (32, 33). This Dit-Tal common

85 core may be decorated by receptor-binding proteins (RBP) and other proteins resulting in a variety of host-
86 binding machinery compositions and architectures. Moreover, Dit and Tal can be functionalized with
87 carbohydrate-binding modules, involved in host cell surface polysaccharide binding, and lysin domains
88 (16, 18, 19, 34, 35). In particular, the Tal C-terminal extensions of *Streptococcus thermophilus* and
89 *Oenococcus oeni* phages contain one to five carbohydrate-binding modules (CBM). However, even though
90 siphophages' host-binding machineries share common structural and functional domains, their
91 compositional and architectural diversity leads to different host-binding mechanisms.

92 Since host-binding machineries bind to host-specific cell surface components, we reasoned that those of
93 mycobacteriophages might differ from those of previously-mentioned phages. Therefore, in order to get
94 structural insights into mycobacteriophages' host-binding machineries, we have performed AlphaFold2
95 (AF2) structure predictions of seven representative siphophages infecting *M. smegmatis*, *M. abscessus* and
96 *M. tuberculosis*. AF2 overcomes bottlenecks of other structure prediction methods such as HHpred (36)
97 and I-TASSER (37) that rely on the presence of structural homologs in the Protein Data Bank (PDB), and
98 produces 3D structures of full-length proteins, multimeric assemblies and protein-protein complexes that
99 equally well compare with experimental structures when AF2 confidence scores are carefully considered
100 (38–41). It therefore positions as a powerful technique providing reliable structural information on flexible
101 and multi-domains proteins that cannot be otherwise analysed using experimental techniques such as X-ray
102 crystallography and cryo-electron microscopy (18, 19, 30). We have predicted the structures of Dits and
103 Tals, and of potential RBPs and additional adhesion proteins assembling these machineries. Interestingly,
104 these mycobacteriophage proteins share common features with those of other Gram-positive-infecting
105 siphophages, such as the presence of CBMs, and also present unique features. In particular, brush-like
106 polyglycine-rich domains with surface-exposed hydrophobic residues are present in potential RBPs and/or
107 additional adhesion proteins, and could be involved in the recognition and binding of the mycobacterial
108 cell wall. All in all, these mycobacteriophage protein structures confirm the mosaic and versatile nature of
109 phage host-binding machineries evolved to bind to a specific host and to initiate the viral infection cycle.

110

111

112 **Results**

113 We have selected seven mycobacteriophages that belong to different genomic clusters, infect different
114 *Mycobacterium* species, and assemble host-binding machineries of different composition (Fig. 1A).
115 Abinghost, Badfish, Butters and Wildcat infect *M. smegmatis* and belong to different clusters (clusters N
116 and V for Butters and Wildcat, respectively) and subclusters (subclusters B1 and B3 for Badfish and
117 Abinghost, respectively). Isca infects *M. abscessus* and belongs to the Subcluster A3. Araucaria also infects
118 *M. abscessus* but it has not been assigned to a cluster yet. Lastly, DS6A infects *M. tuberculosis* and is a

119 singleton. Genes coding for proteins assembling siphophage host-binding machineries are, in most cases,
120 comprised in sequence between the long *tmp* (gene coding for the tape measure protein that defines the
121 length of the tail) and the *holin/lysin* cassette (42). Genome analyses in between the *tmp* and the first *lysin*
122 of these phages have revealed different compositions of their host-binding machineries (Fig. 1B). Of note,
123 *lysin A* and *lysin B* in Isca correspond to ORFs 9 and 10 upstream the *tmp*. We have found that three to five
124 proteins of more than 150 residues assemble mycobacteriophages' host-binding machineries including Dit,
125 Tal, potential RBP and potential additional adhesion proteins (Fig. 1A and B).

127 **Predicted structures of Dits**

128 The Dit size of the seven analysed mycobacteriophages varies from 323 residues for the shortest (DS6A)
129 to 573 residues for the longest (Butters) (Fig. 1A). The Dit of DS6A is a classical Dit containing only the
130 belt and galectin domains (Fig. 2). However, the galectin loops are slightly longer than those of Dit from
131 siphophages infecting the lactic acid bacteria (LAB) *L. lactis*, *S. thermophilus*, and *O. oeni* (18, 19, 23).
132 Even though the Dit of Abinghost, Araucaria and Badfish are ~150 residues longer than that of DS6A, their
133 predicted structure show that they are also classical Dits devoid of any CBM insertion. Instead, they possess
134 two long loops in the belt domain (residues 24-107 in the Abinghost's Dit, residues 25-76 in the Araucaria
135 and Badfish's Dit) and the galectin domain (residues 296-343 in the Abinghost's, Araucaria's and Badfish's
136 Dits) (Fig. 2). It is worth noting that these flexible loops are likely to adopt conformations different from
137 those predicted by AF2 (Fig. S1). Moreover, their location at the periphery of the Dit ring indicate that they
138 may be involved in establishing contacts with other components of the host-binding machinery or with the
139 host cell surface. Lastly, the Dit of Isca, Butters and Wildcat present an additional domain inserted in the
140 galectin domain. When submitted to the Dali server, the predicted structure of the Isca's Dit 110-residue
141 inserted domain returned a hit with a domain of the myophage T4 gp11 protein found at the interface
142 between the baseplate and the short tail fibers (Table 1). However, the inserted domain in the longest
143 Butter's and Wildcat's Dits has been identified as a CBM family 66 (CBM66) according to the
144 Carbohydrate-Active EnZymes database (CAZy) nomenclature (43) (Table 1).

146 **Predicted structures of Tals**

147 The Tal size of the seven analysed mycobacteriophages varies from 345 residues for the shortest
148 (Araucaria) to 794 residues for the longest (Wildcat) (Fig. 1A). However, they all assemble bulky trimers
149 devoid of elongated extensions like those observed in the Tals of some LAB-infecting siphophages (19, 21,
150 30) (Fig. 3). Abinghost, Araucaria and Badfish possess short Tals of similar size (347-370 residues, Fig.
151 1A). Although their predicted structures superimpose well onto each other (rmsd of 0.8-1.0 Å on C α atoms)
152 and are similar to that of the Tal of *Neisseria meningitidis* MC58 (Table 1), the Tal of Badfish harbours

153 longer N- and C-termini (Fig. 3). Butters, DS6A and Isca possess Tal of medium size (569-599 residues),
154 yet their predicted structures show that they do not contain extra domains. Instead, they incorporate extra
155 α -helices at the N- and C-terminal ends, as well as extra loops within the conserved backbone (Fig. 3 and
156 S1). Lastly, Wildcat has a long, 794-residue Tal that contains an additional N-terminal 237-residue insertion
157 including a lectin-like domain similar to that of a Trypanosoma sialidase (44) (Table1) a linker, and a long
158 α -helix connected to the Tal central core (Fig. 3).

160 **Predicted structures of potential RBPs and additional adhesion proteins**

161 In siphophages infecting Gram-positive bacteria, *bona fide* RBPs, responsible for irreversible binding
162 of a phage to host-specific receptors, are usually encoded by the *orf* adjacent to the *tal* (42, 45). Yet, the
163 order of genes coding for components assembling siphophages' host-binding machineries may vary when
164 *orfs* of ancillary proteins are inserted between the *tal* and *rbp* (42). Moreover, RBPs typically assemble as
165 trimers. Based on these criteria, we have identified candidate *orfs* coding for long RBPs, from 619 residues
166 in Isca to 855 residues in Wildcat (Fig. 1A), that assemble well-packed trimers (Fig. 4). It should be noted
167 that, in Butters, we propose the *orf20*, and not the *tal*-adjacent *orf19*, as being the RBP-coding *orf*. We have
168 made this choice based on the overall structural similarity between the Butters's ORF20 predicted trimer
169 and other phages' RBP trimers (Fig. 4). In addition to RBPs, all mycobacteriophages, but Araucaria and
170 Isca, also have potential additional adhesion proteins, which are adjacent to each other or separated by small
171 ORFs (Fig. 1B). In order to get structural and functional insights into these proteins, we have predicted
172 their structures as monomers, unless otherwise stated.

173 **Mycobacteriophages' potential RBPs are elongated, multi-domain proteins.** The predicted
174 structures of trimers of potential RBPs form elongated assemblies from 250 Å for the shortest (Isca) to 450
175 Å for the longest (Araucaria) (Fig. S1). They are all composed of axial α -helical segments, giving them
176 their rod shape, with unstructured segments and/or α -helices at their N-terminal end that may be involved
177 in anchoring the RBPs to the central Dit-Tal core. A variety of globular domains are found at their C-
178 terminal end and in their central region, except for the RBP of Araucaria, DS6A and Isca (Fig. 4).

179 **Araucaria and Isca.** Araucaria and Isca possess only one RBP. The predicted structure of Araucaria's RBP
180 is composed of a 281-residue α -helical rod, interrupted by a small β -prism (residues 31-56) and followed
181 by three bulky domains. The first domain (residues 311-487), pointing out from the trimer axis, looks like
182 a painter's brush in which a 10-stranded β -sandwich forms the brush handle, and a bundle of 10 left-handed
183 polyglycine type II (PG_{II}) helices forms the brush hairs (Figs .4, 6 and S2). Interestingly, although PG_{II}
184 helix-containing domains have rarely been observed in nature (46–50), the C-terminal domain of the
185 *Salmonella* phage S16 tail fiber adhesin also contains a brush domain involved in host recognition (51).
186 Since we have identified brush domains in other mycobacteriophages' RBPs and additional adhesion

187 proteins, they will be analysed altogether in a next section. The Araucaria' brush domain is followed by a
188 CBM, similar to a CBM66 from a fructosidase specifically hydrolysing levan (52) (Table 1), which is itself
189 followed by a lectin domain similar to the N-terminal domain of the *Burkholderia cenocepacia* superlectin
190 C (BC2L-C) with dual carbohydrate specificity (53) (Table 1) (Fig. 4). Interestingly, this tandem of C-
191 terminal CBM and lectin domain is similar to that found in a *Klebsellia* phage depolymerase that binds to
192 and degrades host cell surface polysaccharides (54). Of note, the three lectin domains in Araucaria's RBP
193 assemble a tight trimer as that of the *Klebsiella* phage depolymerase, which is critical for the trimer stability
194 (54). These CBM and lectin domain are likely involved in the recognition of mycobacteria surface
195 polysaccharides. Although the 240 Å-long Isca's RBP predicted structure is shorter than that of Araucaria
196 and devoid of brush domain, it is overall similar to that of Araucaria. It starts with a 39-residue unstructured
197 stretch followed by α -helical segments, and it is terminated by a CBM and a lectin domain. The CBM also
198 returned a levanase CBM66 as a significant hit using Dali (Table 1). As for the lectin domain, which also
199 forms a tightly packed trimer, it returned the BC2L-C lectin from *B. cenocepacia* as well as the receptor-
200 binding domain of the *L. lactis* phage 1358 RBP (Table 1) as structural homologs, which is in agreement
201 with a host-binding function for this protein.

202 **DS6A.** DS6A is a singleton, which led to predicted structures with low confidence scores in some regions
203 (Fig. S1). Therefore, we show only well-folded domains without junctions in between. In particular, a CBM
204 and a lectin domain assemble the RBP C-terminal end (Fig. 4). The CBM also returned the levanase CBM66
205 as significant hit, and the lectin domain is similar to that of *L. lactis* phage 1358 RBD (Table 1).

206 **Abinghost, Badfish, Butters and Wildcat.** The RBPs of Abinghost, Badfish, Butters and Wildcat all possess
207 a central CBM similar to a CBM16 (Table 1). In the Abinghost's and Badfish's RBPs, this CBM is
208 followed by α -helical segments connected to a second CBM identified by Dali as a levanase CBM66 (Table
209 1), like that found in the Araucaria's and Isca's RBP. However, in the Abinghost's and Badfish's RBP, the
210 C-terminal end is formed by a tightly packed, trimeric knob assembling small domains with a new fold
211 composed of four β -strands and one α -helix (Fig. 4). On the other hand, the RBP of Butters is terminated
212 by a brush domain (residues 486-738), and not by a lectin domain, the overall fold of which differ from
213 that of the Araucaria RBP's brush domain. In particular, the handle comprises an 8-stranded, 2-layer β -
214 sandwich, and the hairs, assembling 12 PG_{II} helices, possess short β -sheets packed onto each other at their
215 tip (Fig. 4 and 6). As for the RBP of Wildcat, the central CBM is followed by α -helical segments abutting
216 to a well-packed trimer of a second CBM at the apex of the RBP. This second CBM shares structural
217 similarities with a beta-agarase CBM6 (Table 1) (55), and is followed by a brush domain that superimposes
218 well on that of the Araucaria's RBP, except that the latter has slightly longer loops at the brush hair tips.
219 To note, the brush domains are packed against the central α -helices and placed above the second CBMs in
220 the Wildcat RBP predicted structure. However, since a 12-residue linker connects the second CBM to the

brush domain (Fig. 4), its intrinsic flexibility may also place the brush domains below the second CBM at the RBP tip, or in a variety of positions in between being stacked against the helical segments and being aligned along the trimeric axis, pointing towards the RBP C-terminal end.

Mycobacteriophages' potential additional adhesion proteins. The predicted structures of potential additional adhesion proteins contain one or two brush domains, except those of Butters and Wildcat (Fig. 5). However, the potential RBPs of Butters and Wildcat contain a brush domain, as does the RBP of Araucaria that has a host-binding machinery devoid of additional adhesion proteins (Fig. 1A). This indicates that brush domains likely play a role in mycobacteria recognition and binding.

Abinghost, Badfish The predicted structures of Abinghost's and Badfish's additional adhesion proteins (ORF32 and ORF33) are very close to each other (Fig. 5), which was also the case for their RBPs. This is in agreement with the fact that these phages belong to the same cluster B. Their adhesion protein 1 (ORF32) starts with a short N-terminal segment and a β -hairpin, which are followed by a tandem of brush domains. The first N-terminal brush domain comprises an 8-stranded, 2-layer β -sandwich and 8 PG_{II} helices, while the second C-terminal brush domain comprises a 10-stranded, 3-layer β -sandwich and 9 PG_{II} helices. As for their adhesion protein 2 (ORF33), they start with a long, poorly predicted N-terminal extension, likely involved in contacting other components of the host-binding machinery, followed by a β -helix domain commonly found in phage tail fibers (56), tail spikes (57), and RBP (58) and associated to polysaccharide binding, and end with a single brush domain similar to that of the Araucaria's RBP (Fig. 5 and S1).

Butters. The Butter's adhesion protein 1 (ORF19) is composed of a N-terminal β -stranded knob followed by a β -barrel, and its C-terminal domain is similar to the *B. cenocepacia* BC2L-C lectin domain (Table 1), as it is also the case for the C-terminal lectin domains of Araucaria's and Isca's RBPs (Fig. 5). The presence of this C-terminal lectin domain led us to predict the structure of a trimeric assembly, which shows an arrangement of the three lectin domains similar to that observed in siphophages' adhesion proteins (23, 24, 59, 60). As for the Butters's adhesion protein 2 (ORF22), it is overall similar to that of Abinghost and Badfish (ORF33). It contains a β -strand-rich N-terminal region and a C-terminal brush domain (Fig. 5). However, this brush domain superimposes well on the C-terminal brush domain of Abinghost's and Badfish's adhesion protein 1 (ORF32) (rmsd = 1.1 Å on $\text{C}\alpha$ atoms).

DS6A The DS6A's additional adhesion proteins (ORF25 and ORF27) also contain a C-terminal brush domain, similar to that of Araucaria's RBP (rmsd 1.1 Å on $\text{C}\alpha$ atoms). As for their N-terminal regions, while the N-terminal jelly roll of ORF25 is predicted with high confidence, the confidence on the ORF27 N-terminal β -strand-rich predicted structure is weak (Fig. 5 and S1).

Wildcat. Lastly, the predicted structures of the two Wildcat's additional adhesion proteins show features that distinguish them from the other mycobacteriophages' additional adhesion proteins. The adhesion protein 1 (ORF44) is made up of three domains: the N-terminal domain is a three-stranded β -sheet, the

255 central α/β domain has been identified as a penicillin-binding domain/ β -lactamase/DD-peptidase structural
256 homolog using Dali (61) (Table 1), and the two C-terminal domains are jelly-rolls packed one onto the
257 other (Fig. 5). Canonical catalytic residues of β -lactamase/DD-peptidase are not found in the central
258 domain, indicating that this protein does not have catalytic activity on host cell surface components.
259 Interestingly, the second adhesion protein (ORF45) harbours a long, unstructured N-terminal segment
260 followed by three consecutive domains, each being composed of two packed jelly-rolls, as those observed
261 in the first adhesion protein (Fig. 5). Such domains have no structural homologs in the PDB.

263 **Mycobacteriophages' brush domains present solvent-exposed hydrophobic residues**

264 All mycobacteriophages under study, but Isca, have proteins in their host-binding machineries
265 containing brush domains. These brush domains are formed of a β -sandwich as the brush handle, from
266 which emerge PG_{II} helices making the brush hairs and accounting for 36% to 47% of glycine residue (Fig.
267 6). The brush domain in the Araucaria's RBP is similar to those found in the Abinghost's and Badfish's
268 adhesion protein 2 (ORF33), in both DS6A's adhesion proteins (ORF25 and ORF27), and in the Wildcat's
269 RBP (ORF43). The handle region is composed of a 10-stranded, 3-layer β -sandwich, and the hairs are
270 formed of ten PG_{II} helices (Fig. 6). Yet, the loops at the brush hair tip are of variable lengths from one
271 phage to another. Brush domains with a fold different from that described above are present in Butters's
272 RBP (ORF20) and adhesion protein 2 (ORF22), and in the Abinghost's and Badfish's additional adhesion
273 proteins 1 (ORF32). Variations occur in the handle, which are formed of a 2-layer β -sandwich as in the
274 Butters' RBP and in the Abinghost's and Badfish's additional adhesion protein 1 (for the N-terminal brush
275 domain), and/or in the hairs, which assemble a variable number of PG_{II} helices with different packing. In
276 these brush domains, variations are also observed at the brush hair tip including small parallel and anti-
277 parallel β -sheets (Fig. 6). These structural variations could reflect the adaptation of phages to their host cell
278 surface.

279 A striking feature of these brush domains is the distribution of their bulky hydrophobic residues (Val,
280 Leu, Ile, Met, Phe, Trp). They are found not only buried within the brush handle, as expected, but also
281 exposed at the surface of the brush hair (Fig. 6). They often form clusters of 4 to 11 residues at the brush
282 hair tip (DS6A's ORF27 and ORF25, respectively), which is the region identified as being the host-binding
283 domain of the long tail fibers of phages infecting Gram-negative bacteria and determining host specificity
284 (51, 62, 63). Such solvent-exposed hydrophobic patches are not common in soluble proteins, except in
285 domains interacting with lipids such as those found in phospholipases or lipases and therefore may be
286 involved in host recognition.

288 **Insights into the architecture of whole mycobacteriophages' host-binding machineries**

289 In order to gain insights into the assembly of mycobacteriophages' host-binding machinery, we have
290 selected in the actinobacteriophage database (phagesdb.org) a good quality negative staining electron
291 micrograph (nsEM) of Badfish in which the tail distal end exhibits recognizable features. The tail itself is
292 made of stacked hexamers of the MTP and is terminated by a bulky volume that can ascribed to the stacked
293 Dit and Tal oligomers that are known to interact together (42). The thin and long extension following the
294 Dit-Tal core could correspond to a RBP trimer with its C-terminal CBMs accounting for the enlargement
295 visible at the end of the tail tip (Fig. 7A and B). This implies that the RBP trimer attaches to the Tal trimer
296 via its N-terminal α -helices. To test this hypothesis, we have submitted three Tal monomers and three RBP
297 monomers (residues 1 to 330 because of memory limitation) to AF2 to predict the structure of the Tal-RBP
298 interface. Remarkably, the predicted structure of the Tal-RBP complex shows the RBPs' N-terminal α -
299 helices nestled in the Tal trimer internal cavity (Fig. 7B, C and S1). The positively-charged surface of each
300 RBP N-terminal helix interacts with a 1350-Å² negatively-charged patch of each Tal monomer (Fig. 7C).
301 This predicted structure supports a direct interaction between the Badfish's Tal trimer and a trimer of its
302 potential RBP. In the nsEM image of Badfish, we cannot identify signals accounting for the two potential
303 additional adhesion proteins. Yet, it should be noted that negative staining with uranyl acetate can disrupt
304 the structural integrity of protein complexes. Notably, the RBP of the *S. thermophilus* phage STP1 has not
305 been observed in virions by nsEM for a long time, and hence was thought to be absent in this phage, because
306 of a too high concentration of uranyl acetate (21, 64).

307 We have also selected in the database a nsEM image of Isca (Fig. 7D), which possesses an RBP but is
308 devoid of additional adhesion proteins (Fig. 1). Our predicted structures of Isca's Dit, Tal and RBP match
309 with the overall architecture of its host-binding machinery observed in virions. The two disks after the last
310 MTP ring can be ascribed to the Dit and Tal oligomers, and hence, the thin extension terminated by a cone-
311 shaped structure at the tail tip can be ascribed to a RBP trimer (Fig. 7D and E). Yet, in this case, the Tal-
312 RBP interface could not be predicted by AF2, likely because of the long, unstructured N-terminal stretch
313 of the Isca's RBP (39 residues).

314 315 **Discussion**

316 Deciphering molecular mechanisms used by mycobacteriophages to infect their host is important not
317 only to better understand viral diversity and evolution, phage-host interactions, but also to develop phage-
318 based biotechnological tools. Our predicted structures of mycobacteriophages' host-binding machineries
319 emphasize on the modularity and structural diversity of these macromolecular assemblies that tailed phages
320 have evolved to adapt to the surface of bacteria present in their ecological niche and hence to enhance their
321 ability to initiate contacts with potential hosts. Components of mycobacteriophages' host-binding
322 machineries variably combine structural and functional domains, including CBMs, lectin domains and β -
323 helices that have also been found in tailed phages infecting Gram-negative or Gram-positive bacteria to

324 interact with cell surface polysaccharides. Surprisingly, polyglycine-rich domains, which we named brush
325 domains based on their peculiar shape, appear to be widespread in mycobacteriophages' host-binding
326 machineries and likely play an important role in host binding. Brush domains have rarely been observed in
327 nature, yet they appear as a versatile scaffold with functions related to protein-protein or protein-ligand
328 interactions (46, 48–51). The unique example of a viral brush domain is that of the *Salmonella* phage S16
329 tail fiber adhesin that initiates host adsorption by binding to the outer membrane protein OmpC and
330 lipopolysaccharides of the *Salmonella* cell wall (65). In particular, residues located in highly variable loops
331 at the brush hair tip were shown to be involved in host binding (51). In mycobacteriophages, these brush
332 domains present surface-exposed hydrophobic residues, and we therefore propose that they have likely
333 evolved to interact with lipids of the waxy mycobacterial cell wall in which exceptionally long chain fatty
334 acids, the mycolic acids, form an outer mycomembrane. Interestingly, these brush domains may be more
335 common in nature than expected and could be the hallmark of mycobacteriophages. Further AF2 structure
336 predictions of host-binding machineries of different phages will provide valuable information to better
337 understand the function and evolution of these brush domains as well as to unveil unprecedented structural
338 and functional domains.

339 The mycomembrane is a very unusual feature for bacteria that belong to the *Actinobacteria* phylum of
340 Gram-positive (monoderm) bacteria, which raises questions related to the nature of cell receptors
341 recognized by mycobacteriophages and to their host-binding mechanisms. The mycomembrane inner
342 leaflet, composed of mycolic acids, is linked to arabinogalactan that is in turn covalently attached to
343 peptidoglycan, while its outer leaflet is composed of free phospholipids, glycolipids and lipoglycans (66).
344 The mycomembrane also contains porins and transporters and is covered by a capsule formed of proteins,
345 polysaccharides and small amounts of lipids (66). Therefore, there is a variety of potential receptors in the
346 mycobacterial cell wall for phage binding. Yet, our predicted structures indicate that mycobacteriophages
347 would preferentially use cell surface polysaccharides and lipids rather than proteins to recognize and attach
348 to their hosts. Indeed, CBM, lectin domains, β -helix domains and brush domains, likely involved in
349 carbohydrate and lipid binding, are present in mycobacteriophages' potential RBP and additional adhesion
350 proteins. This composition is in agreement with the fact that, although there are only few papers reporting
351 the identification of mycobacteriophages' receptors, cell wall associated-glycolipids or phospholipids have
352 been identified as cell surface receptors of the mycobacteriophages D4, D29, I3 and Phlei, the lipid or sugar
353 moieties of which being important for host binding depending on the phage (67–69). Moreover, the
354 extended RBPs that certainly plugged to the Tal trimer, like in our Badfish Tal-RBP complex structure
355 prediction, recall the overall architecture and composition of the CBM-rich Tal C-terminal extension of *S.*
356 *thermophilus* and *O. oeni* siphophages that bind to cell surface polysaccharides (18, 19). In the Araucaria
357 low-resolution structure, the long fiber observed at the tail distal end, similar to that of the coliphages

358 Lambda and T5, suggested that this mycobacteriophage could bind to proteinaceous receptors (14). The
359 nsEM image of Isca, which only contains RBP like Araucaria, also shows a single fiber attached to Tal
360 (Fig. 7D). However, the predicted structures of Araucaria's and Isca's RBPs reveal that they encompass
361 CBM, lectin domains and brush domains at the tip of the RBP trimeric assemblies, which does not support
362 the possibility of binding to proteins. Mycobacteriophages would therefore be more similar to siphophages
363 infecting Gram-positive bacteria that generally bind to cell wall polysaccharides like the model *L. lactis*
364 phage p2 (70, 71) than siphophages infecting Gram-negative bacteria that irreversibly bind to membrane
365 proteins like the model *E. coli* phages Lambda and T5 (72). The presence of multiple host-binding domains,
366 either within the same protein (the Araucaria's and Isca's RBPs for instance), or in different host-binding
367 proteins (the DS6A's and Butters' RBP and additional adhesion proteins for instance), offer the advantage
368 to achieve strong interactions with the host in spite of low binding affinity of saccharides and lipids through
369 avidity phenomena, and this could be strengthened with the presence of multiple copies of the additional
370 adhesion proteins and their possible multimerization.

371 Mycobacteriophages of the actinobacteriophage database, but the *M. tuberculosis* phage DS6A, have
372 been isolated on the non-pathogenic *M. smegmatis* mc²155, and the host range of only few of them has
373 been reported (12, 73). There is a correlation between phage clusters and host range, yet
374 mycobacteriophages have the ability to expand or modify their host range through mutations in tail genes
375 (73). However, such mutations in the mycobacteriophages Rosebush (phagesdb.org, cluster B2, ORF32,
376 mutations L297R, W121G, G288D, W293R) and Halo (phagesdb.org, Cluster G, ORF22, mutations
377 A306V, A604E) do not map to CBM, lectin domain or brush domain identified in our predicted structures
378 (Fig. S3). Instead, they are found in positions that likely impact the trimeric assembly of these proteins,
379 thereby changing their ability to interact with mycobacterial cell envelope and modifying host specificity.
380 Interestingly, the Halo mutant that was isolated for its ability to solely infect *M. tuberculosis* did acquire
381 the capacity to bind to *M. tuberculosis* and was still able to strongly bind, but not infect, *M. smegmatis* (73).
382 Generally, host binding is a two-step process that includes the reversible binding of phages that search for
383 specific receptors, followed by irreversible contacts that anchor the whole host-binding machinery to the
384 host in a stable conformation allowing genome injection. In this process, the host-binding proteins and cell
385 receptors mediating the reversible contacts are not necessarily the same as those mediating the irreversible
386 attachment (74), and we hypothesize that the potential additional adhesion proteins, rich in brush domains,
387 would be important for reversible contacts with the host while also strengthening the irreversible contact
388 of the RBP. The observation that the Halo mutant binds both *M. tuberculosis* and *M. smegmatis* but only
389 infects *M. tuberculosis* and our predicted structure of the ORF22 of the Halo host-binding machinery
390 indicate that these mutations would rather impact the dynamics of the whole host-binding machineries
391 leading to the irreversible attachment of phages before DNA ejection than the mere recognition of receptors.

392 The function of CBMs, lectin domains, and brush domains identified in mycobacteriophages needs to
393 be validated experimentally, and our structure predictions provide a rational basis to produce recombinant
394 proteins or protein domains and perform cell-binding assays. Moreover, structural analyses of whole host-
395 binding machineries in virions alone and bound to the host cell wall will be essential to shed light on their
396 overall architecture and on molecular mechanisms used by mycobacteriophages to recognize and bind to
397 their host cell surface. In particular, deciphering the role, at the molecular level, and the dynamics of CBM
398 insertions in the Dit-Tal core and of the CBM, lectin domains and brush domains in the potential RBP and
399 additional adhesion proteins throughout host binding, will be important to better understand phage-host
400 interactions and also offers great perspectives for host-range engineering (17, 63), which may be
401 instrumental to rationally develop diagnostics tools and phage therapy protocols.

402 Lastly, our work illustrates the tremendous potential of AF2 to unveil virus structure and biology. AF2
403 structure predictions can be used to revisit viral genome annotations and to efficiently characterize the
404 overwhelming number of newly discovered viruses, as exemplified by the human gut phageome (75), and
405 the ‘viral dark matter’ - viral proteins of unknown function – that can represent up to 75% of phage encoded
406 proteins.

409 **Material and Methods**

410 **Phage selection**

411 Seven mycobacteriophages that belong to different clusters and infect different hosts were selected for
412 analysis in this study. Abinghost (cluster B3), Badfish (cluster B1), Butters (cluster N) and Wildcat (cluster
413 V) infect *M. smegmatis*. Araucaria (not classified) and Isca (cluster A3) infect *M. abscessus*. DS6A is a
414 singleton that infects *M. tuberculosis*. The Genbank accession numbers are as follow: Abinghost
415 (MN444873), Araucaria (AHAS00000000), Badfish (KJ194580), Butters (KC576783), DS6A (JN698994),
416 Isca (MN586063) and Wildcat (DQ398052).

417 **Protein structure predictions**

418 We used a Github notebook
419 (<https://colab.research.google.com/github/deepmind/alphafold/blob/main/notebooks/AlphaFold.ipynb#scrollTo=XUo6foMQxwS2>) to perform structure predictions of protein monomers and multimers. Due to
420 memory limitations, the number of residues for monomer and multimer structure predictions had to be less
421 than 1400, and we therefore split long sequences into smaller stretches. First, we ran structure predictions
422 of monomers in order to determine sensible boundaries of stretches that were then submitted for structure
423 prediction as trimers (Tal, RBP) or hexamers (Dit) to AlphaFold multimer (40). Moreover, we predicted
424 structures of stretches with overlapping segments to allow assembly of full-length multimers using *Coot*
425

426 (76). The structure prediction of Badfish Tal-RBP complex was performed by submitting three sequences
427 of Tal and three sequences of the potential RBP (residues 1-330). The pLDDT values of predicted structure,
428 stored in the pdb file as B-factors, were plotted using Excel (Fig. S1 and S3). The final predicted protein or
429 domain structures were submitted to the Dali server (77) to identify the closest structural homologs in the
430 PDB. Visual representations of the structures were prepared with ChimeraX (78).

431 432 433 **Author Contributions**

434 Conceptualization: A.G., C.C.; Investigation: C.C.; Resources: C.C.; Formal analysis: A.G. C.C.; Writing
435 original draft: A.G., C.C.; Writing-review and editing: A.G., C.C.; Visualization: A.G., C.C.; Project
436 administration: A.G.

437 **Acknowledgments**

438 This research received no specific grant from any funding agency in the public, commercial, or not-for
439 profit sectors.

440 **Conflict of Interest**

441 C.C. is employee of AlphaGraphix (cambillau.alphagraphix@gmail.com).

442 **Acknowledgments**

443 We acknowledge UCSF ChimeraX for molecular graphics that is developed by the Resource for
444 Biocomputing, Visualization, and Informatics at the University of California, San Francisco, with support
445 from National Institutes of Health R01-GM129325 and the Office of Cyber Infrastructure and
446 Computational Biology, National Institute of Allergy and Infectious Diseases.

447 **Supplemental Material**

448 Supplemental material for this article is provided.

449 **Data Availability**

450 All coordinates of predicted structures will be deposited in the open data repository Zenodo.

Table 1. Dali hits of protein and domain predicted structures.

	Dali hits			
	PDB ID	Z-score*	rmsd (Å)#	lali\$
Abinghost				
Tal (ORF30)	3d37	18.5	4	250/316
RBP (ORF31)				
CBM16	2zew	15	2.6	140/147
CBM66	4azz	16.4	2.6	53/165
Araucaria				
Tal (ORF71)	3d37	18.4	3.8	247/316
RBP (ORF72)				
CBM66	4azz	14.3	2.3	131/165
lectin	2wq4	13	2.3	117/131
Badfish				
Tal (ORF30)	3d37	18	4.2	248/316
RBP (ORF31)				
CBM16	2zey	15.3	2.8	136/145
CBM66	4azz	16.4	2.6	53/165
Butters				
Dit (ORF17)				
CBM66	4azz	15.4	2.8	54/165
RBP (ORF20)				
CBM16	2zew	14.3	2.6	129/144
Additional adhesion protein 1 (ORF 19)				
lectin	2wq4	8.9	2.5	109/134
DS6A				
RBP (ORF23)				
CBM66	4b11	10	3.3	126/162
lectin	4l92	6.3	3.5	121/391
Isca				
Dit (ORF27)				
domain of T4 gp11	1el6	5.4	3.2	83/208
RBP (ORF31)				
CBM66	4azz	15.6	2.6	154/165
lectin	2wq4/4l92	9/8.9	2.7/2.8	114/134 ; 124/391
Wildcat				
Dit (ORF39)				
CBM66	4azz	15.9	2.5	51/165
Tal (ORF40)	1mz5	15.8	2.5	165/622
RBP (ORF43)				
CBM16	2zew	14.3	2.6	129/144
CBM6	2cdo	6.3	2.4	94/138
Additional adhesion protein 1 (ORF 44)				
penicillin-binding domain	1ei5	29.6	3.5	312/518

452 * : the Z-score is a measure of structural similarity

453 #: root mean square deviation on C α atoms

454 \$: lali indicates the number of aligned residues in the predicted structure over the total number of residues in the PDB

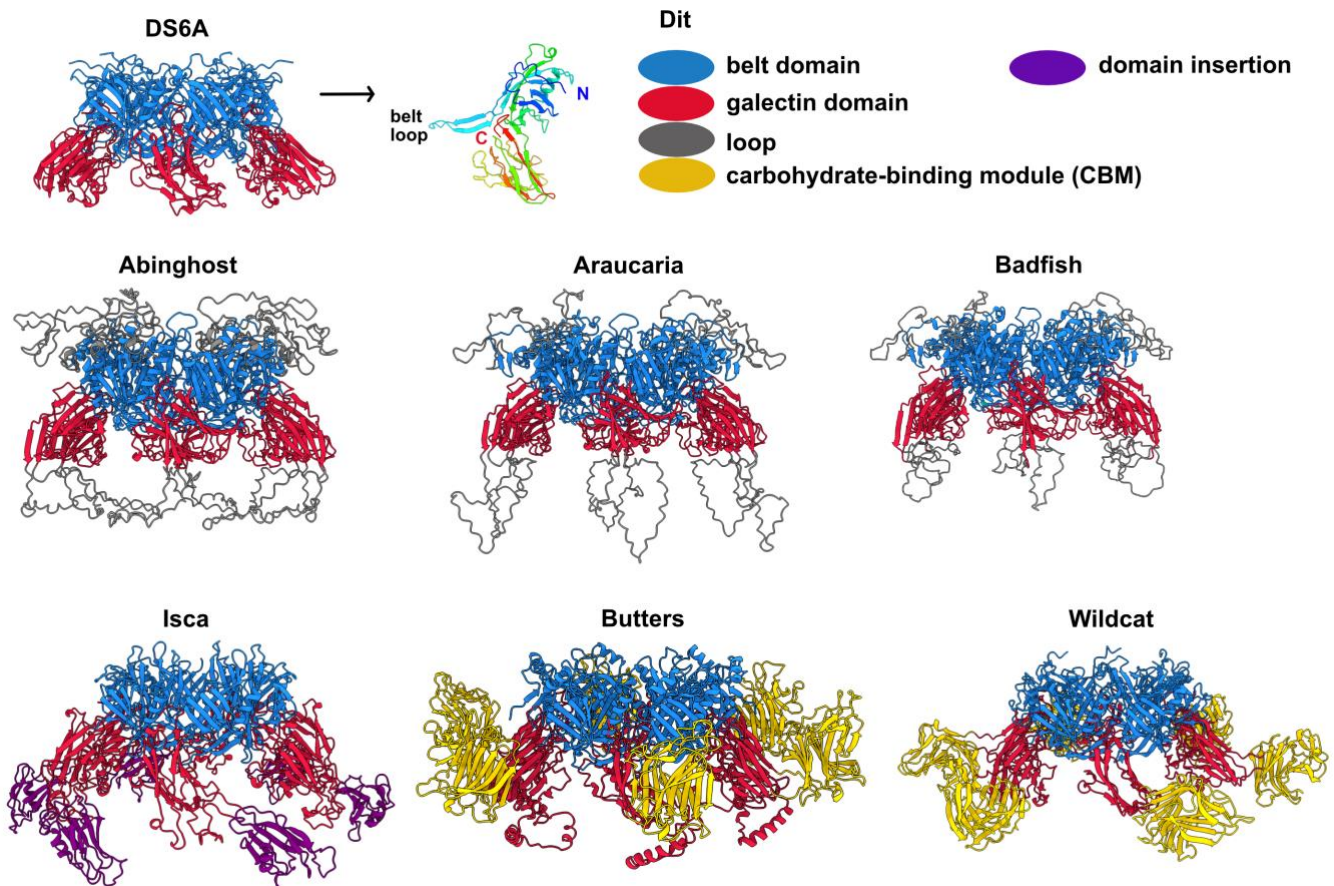
455 hit.

456



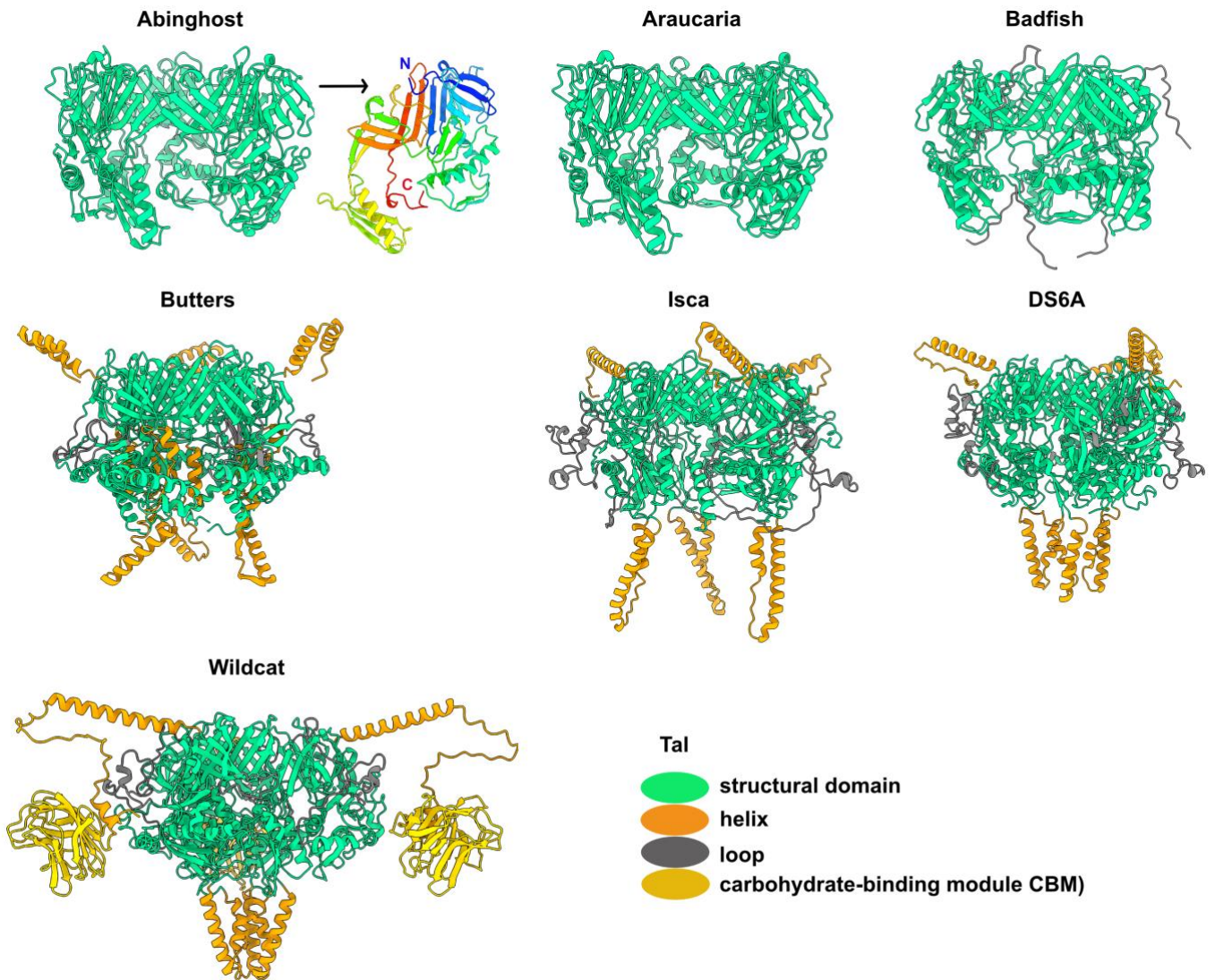
458

459 **Fig. 1. Diversity of mycobacteriophages analysed in this study. A.** Genomic properties (clusters), host
 460 spectrum, and host-binding machinery composition of the seven analysed mycobacteriophages. **B.**
 461 Schematics of the genomic region coding for proteins assembling host-binding machineries.



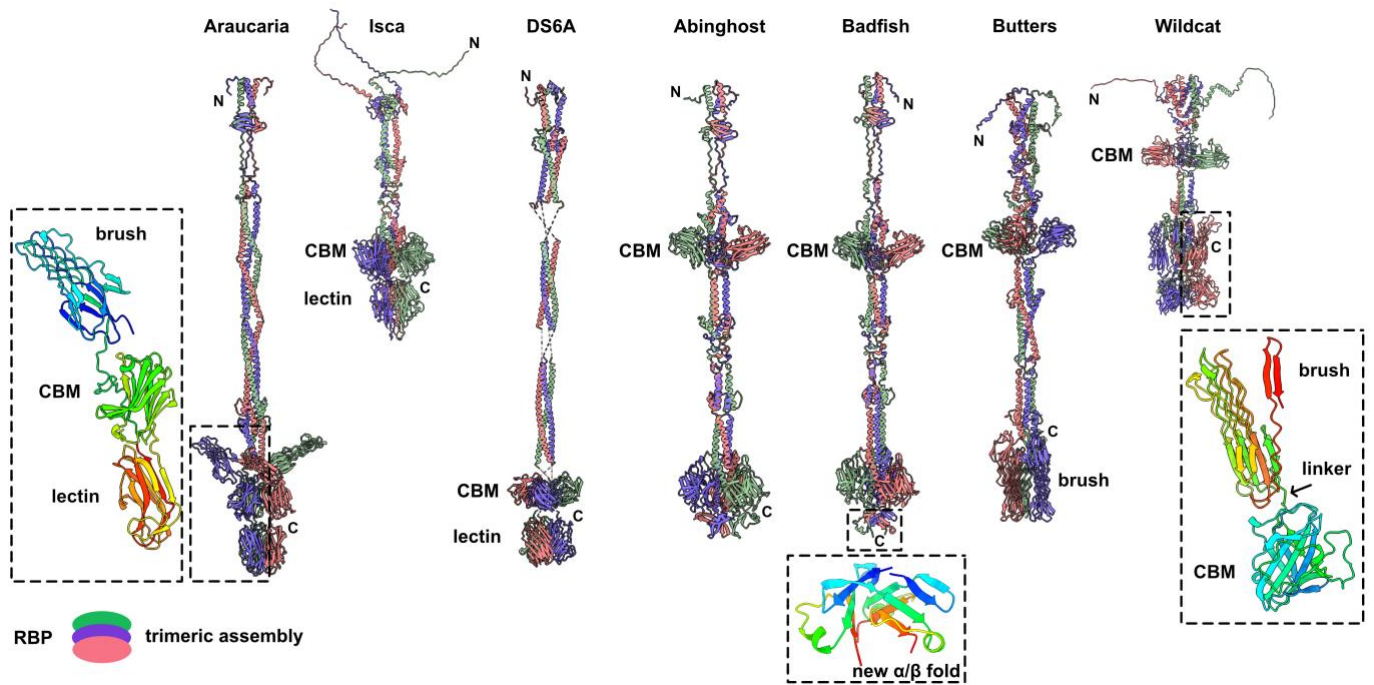
462

463 **Fig. 2. DIT predicted structures.** Ribbon representations of DIT predicted structures. One rainbow-colored
 464 DS6A DIT monomer is shown for clarity. While the DIT of Abinghost, Araucaria, Badfish and DS6A contain
 465 only the belt and galectin domains, the DIT of Butters, Isca and Wildcat have an extra domain inserted in
 466 the galectin domain. This extra-domain has been identified as a CBM for the DIT of Butters and Wildcat.
 467 The color code is indicated.



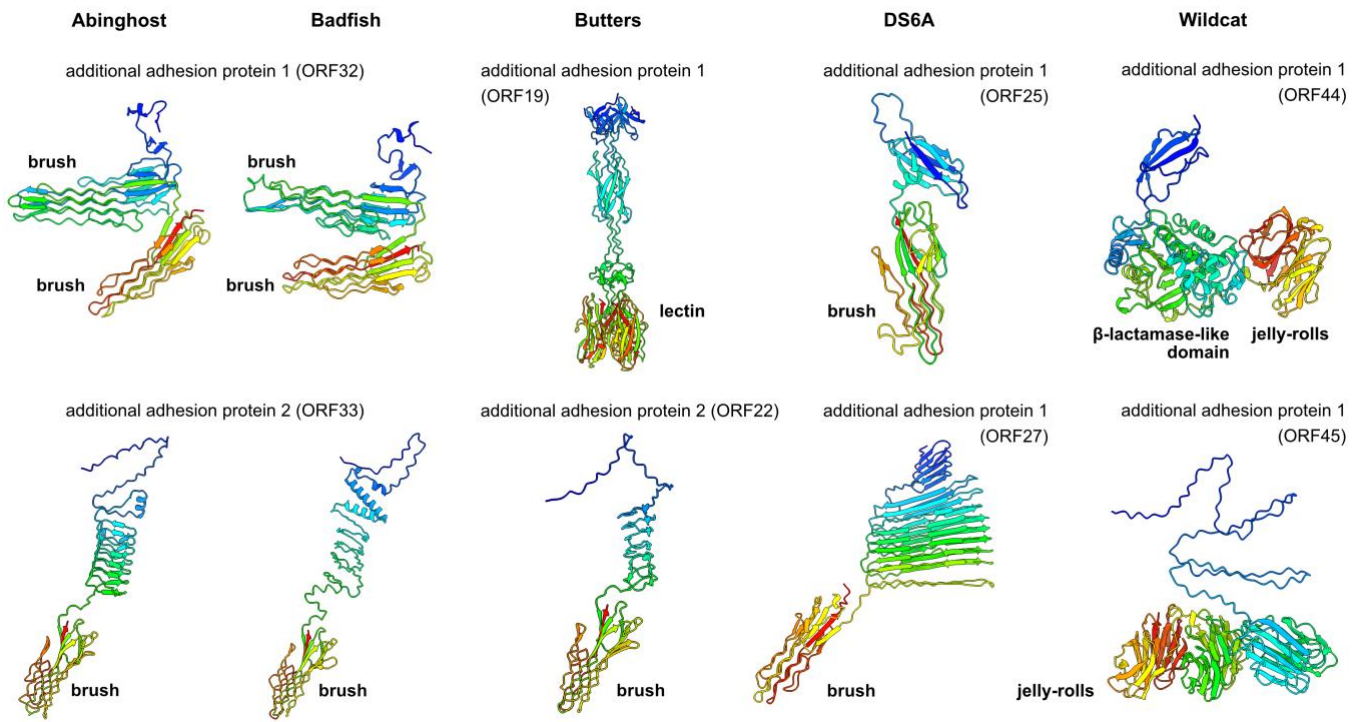
468

469 **Fig. 3. Tal predicted structures.** Ribbon representations of the Tal predicted structures. One rainbow-
 470 colored Abinghost Tal monomer is shown for clarity. The color code is indicated.



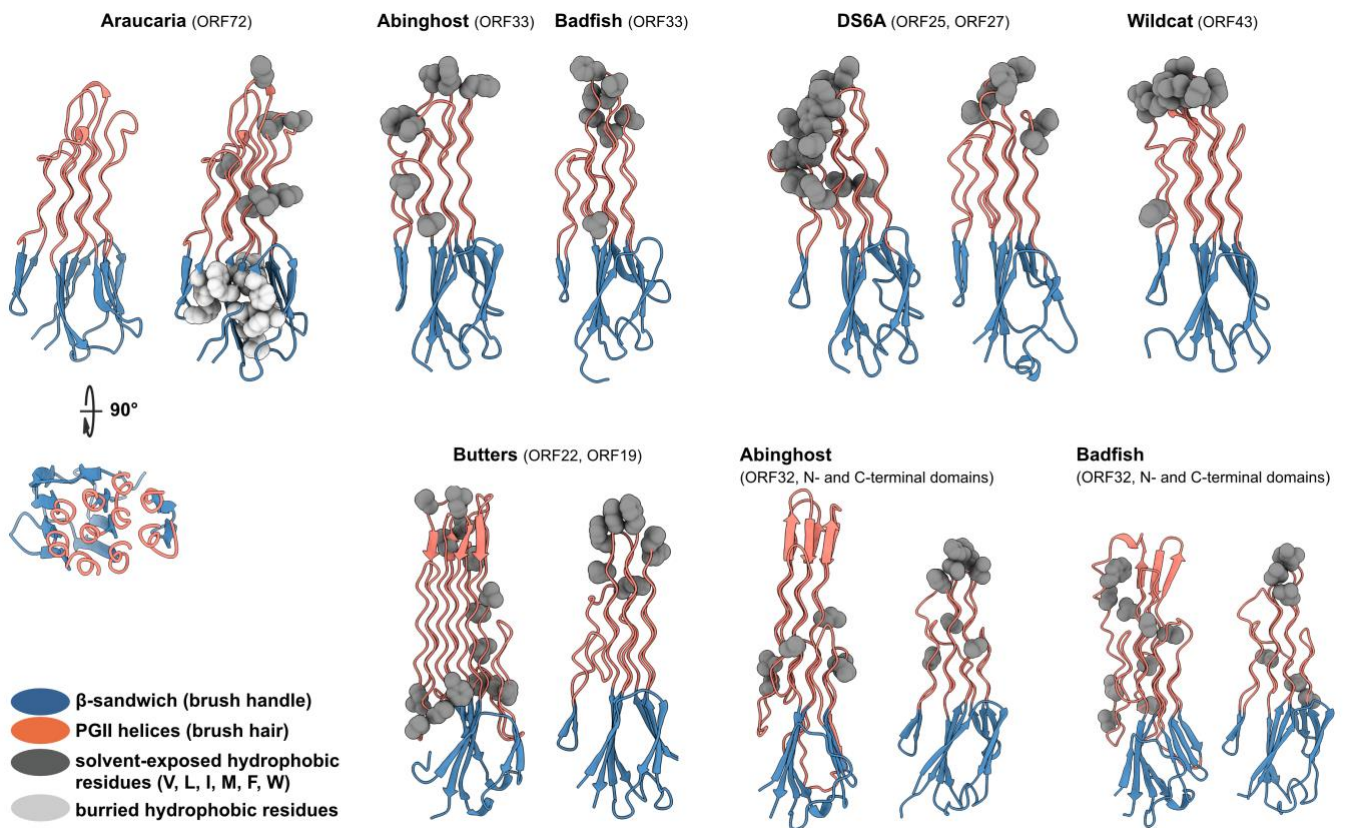
471

472 **Fig. 4. Predicted structures of potential RBPs reveal multi-domain assemblies.** Ribbon representations
 473 of the RBP predicted structures (the N- and C-termini are indicated.). Trimeric assemblies are colored by
 474 chain. The inset on the left-hand side highlights the brush domain, CBM and lectin domain at the
 475 Araucaria's RBP C-terminal end (rainbow colored from the N-terminus to the C-terminus). The inset below
 476 the Badfish's RBP highlights the trimeric knob (rainbow colored) at its C-terminal end. The inset on the
 477 right-hand side highlights the brush domain and CBM, connected by a flexible linker, at the Wildcat's RBP
 478 C-terminal end (rainbow colored).



479

480 **Fig. 5. Additional adhesion protein predicted structures.** Ribbon representations of the additional
 481 adhesion protein predicted structures (rainbow colored from the N-terminus to the C-terminus).



482

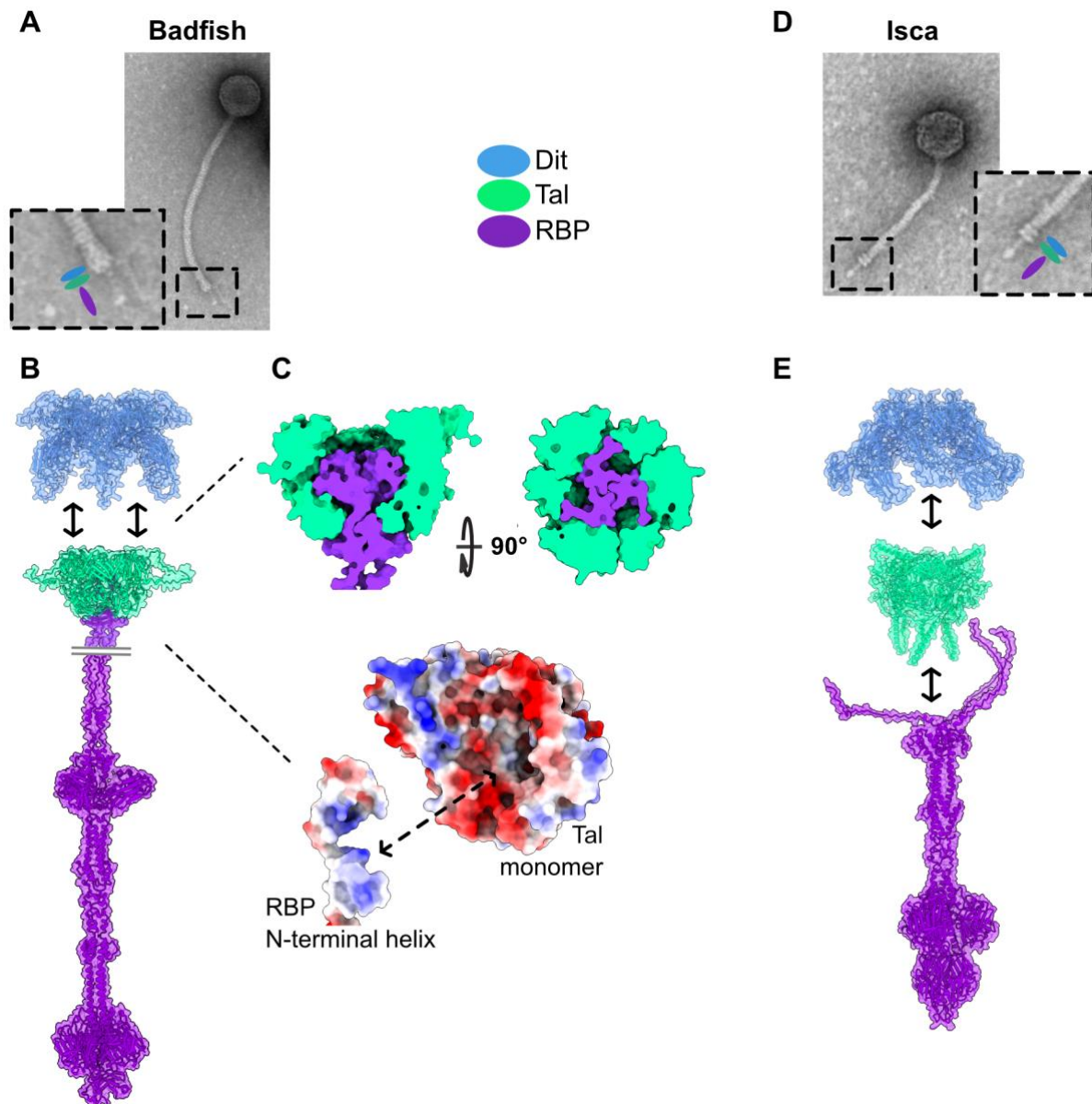
483

484

485

486

Fig. 6. Mycobacteriophages' brush domains present solvent-exposed hydrophobic residues. Ribbon representations of the brush domains identified in the mycobacteriophages' RBPs and additional adhesion proteins. Solvent-exposed and buried hydrophobic residues are shown as spheres. The color code is indicated.



487

488 **Fig. 7. Architectures of mycobacteriophages' host-binding machineries.** **A.** nsEM picture of Badfish
 489 (from phagesdb.org) with its host-binding machinery at the tail distal end highlighted in a dotted box. In
 490 the inset, the colored circles indicate the relative position of Dit, Tal and RBP (the color code is indicated).
 491 **B.** Surface and ribbon representation of the predicted structures of the Badfish's Dit ring, Tal-RBP<sub>(residues 1-
 492 330)</sub> complex, and full-length RBP trimer. These predicted structures are compatible with the overall
 493 architecture of its host-binding machinery. **C.** Top. Orthogonal, slabbed views of the Tal-RBP complex
 494 showing three RBP N-terminal helices plugged into the Tal trimer internal cavity. Bottom. A positively-
 495 charged surface of one RBP N-terminal helix interacts with a negatively-charged surface of a Tal monomer.
 496 **D.** nsEM picture of Isca (from phagesdb.org) with its host-binding machinery at the tail distal end
 497 highlighted in a dotted box. In the inset, the colored circles indicate the relative position of Dit, Tal and
 498 RBP (the color code is indicated). **E.** Surface and ribbon representation of the predicted structures of the
 499 Isca's Dit, Tal, and RBP. These predicted structures are compatible with the overall architecture of its host-
 500 binding machinery.

501 **References**

- 502 1. Prasanna AN, Mehra S. 2013. Comparative phylogenomics of pathogenic and non-pathogenic
503 mycobacterium. *PLoS One* 8:e71248.
- 504 2. Barka EA, Vatsa P, Sanchez L, Gaveau-Vaillant N, Jacquard C, Meier-Kolthoff JP, Klenk H-P,
505 Clément C, Ouhdouch Y, van Wezel GP. 2016. Taxonomy, Physiology, and Natural Products of
506 Actinobacteria. *Microbiol Mol Biol Rev* 80:1–43.
- 507 3. Broncano-Lavado A, Senhaji-Kacha A, Santamaría-Corral G, Esteban J, García-Quintanilla M.
508 2022. Alternatives to Antibiotics against Mycobacterium abscessus. *Antibiotics (Basel)* 11:1322.
- 509 4. Daffé M, Marrakchi H. 2019. Unraveling the Structure of the Mycobacterial Envelope. *Microbiol*
510 *Spectr* 7.
- 511 5. Batt SM, Minnikin DE, Besra GS. 2020. The thick waxy coat of mycobacteria, a protective layer
512 against antibiotics and the host's immune system. *Biochemical Journal* 477:1983–2006.
- 513 6. Hanauer DI, Graham MJ, SEA-PHAGES, Betancur L, Bobrownicki A, Cresawn SG, Garlena RA,
514 Jacobs-Sera D, Kaufmann N, Pope WH, Russell DA, Jacobs WR, Sivanathan V, Asai DJ, Hatfull GF.
515 2017. An inclusive Research Education Community (iREC): Impact of the SEA-PHAGES program on
516 research outcomes and student learning. *Proc Natl Acad Sci U S A* 114:13531–13536.
- 517 7. Russell DA, Hatfull GF. 2017. PhagesDB: the actinobacteriophage database. *Bioinformatics*
518 33:784–786.
- 519 8. Hatfull GF. 2020. Actinobacteriophages: Genomics, Dynamics, and Applications. *Annu Rev Virol*
520 7:37–61.
- 521 9. Hatfull GF, Sarkis GJ. 1993. DNA sequence, structure and gene expression of mycobacteriophage
522 L5: a phage system for mycobacterial genetics. *Mol Microbiol* 7:395–405.
- 523 10. Hatfull GF. 2018. Mycobacteriophages. *Microbiol Spectr* 6.
- 524 11. Sarkis GJ, Jacobs WR, Hatfull GF. 1995. L5 luciferase reporter mycobacteriophages: a sensitive
525 tool for the detection and assay of live mycobacteria. *Mol Microbiol* 15:1055–1067.
- 526 12. Dedrick RM, Guerrero-Bustamante CA, Garlena RA, Russell DA, Ford K, Harris K, Gilmour KC,
527 Soothill J, Jacobs-Sera D, Schooley RT, Hatfull GF, Spencer H. 2019. Engineered bacteriophages for
528 treatment of a patient with a disseminated drug-resistant Mycobacterium abscessus. *Nat Med* 25:730–733.
- 529 13. Little JS, Dedrick RM, Freeman KG, Cristinziano M, Smith BE, Benson CA, Jhaveri TA, Baden
530 LR, Solomon DA, Hatfull GF. 2022. Bacteriophage treatment of disseminated cutaneous Mycobacterium
531 chelonae infection. *Nat Commun* 13:2313.
- 532 14. Sassi M, Bebeacua C, Drancourt M, Cambillau C. 2013. The First Structure of a
533 Mycobacteriophage, the Mycobacterium abscessus subsp. bolletii Phage Araucaria. *Journal of Virology*
534 87:8099–8109.
- 535 15. Biemann R, Habann M, Eugster MR, Lurz R, Calendar R, Klumpp J, Loessner MJ. 2015.
536 Receptor binding proteins of *Listeria monocytogenes* bacteriophages A118 and P35 recognize serovar-
537 specific teichoic acids. *Virology* 477:110–118.
- 538 16. Dieterle M-E, Spinelli S, Sadvovskaya I, Piuri M, Cambillau C. 2017. Evolved distal tail
539 carbohydrate binding modules of *Lactobacillus* phage J-1: a novel type of anti-receptor widespread
540 among lactic acid bacteria phages. *Mol Microbiol* 104:608–620.
- 541 17. Dunne M, Rupf B, Tala M, Qabrati X, Ernst P, Shen Y, Sumrall E, Heeb L, Plückthun A,
542 Loessner MJ, Kilcher S. 2019. Reprogramming Bacteriophage Host Range through Structure-Guided
543 Design of Chimeric Receptor Binding Proteins. *Cell Reports* 29:1336-1350.e4.
- 544 18. Goulet A, Joos R, Lavelle K, Van Sinderen D, Mahony J, Cambillau C. 2022. A structural
545 discovery journey of streptococcal phages adhesion devices by AlphaFold2. *Front Mol Biosci* 9:960325.
- 546 19. Goulet A, Cambillau C. 2021. Structure and Topology Prediction of Phage Adhesion Devices
547 Using AlphaFold2: The Case of Two *Oenococcus oeni* Phages. *10. Microorganisms*, 2021/10/24 ed. 9.
- 548 20. Kizziah JL, Manning KA, Dearborn AD, Dokland T. 2020. Structure of the host cell recognition
549 and penetration machinery of a *Staphylococcus aureus* bacteriophage. *PLoS Pathog* 16:e1008314.
- 550 21. Lavelle K, Goulet A, McDonnell B, Spinelli S, van Sinderen D, Mahony J, Cambillau C. 2020.
551 Revisiting the host adhesion determinants of *Streptococcus thermophilus* siphophages. *Microb Biotechnol*
552 13:1765–1779.
- 553 22. Legrand P, Collins B, Blangy S, Murphy J, Spinelli S, Gutierrez C, Richet N, Kellenberger C,

554 Desmyter A, Mahony J, van Sinderen D, Cambillau C. 2016. The Atomic Structure of the Phage Tuc2009
555 Baseplate Tripod Suggests that Host Recognition Involves Two Different Carbohydrate Binding
556 Modules. *mBio* 7:e01781-15, /mbio/7/1/e01781-15.atom.

557 23. Sciarra G, Bebeacua C, Bron P, Tremblay D, Ortiz-Lombardia M, Lichière J, van Heel M,
558 Campanacci V, Moineau S, Cambillau C. 2010. Structure of lactococcal phage p2 baseplate and its
559 mechanism of activation. *Proc Natl Acad Sci USA* 107:6852–6857.

560 24. Veesler D, Spinelli S, Mahony J, Lichiere J, Blangy S, Bricogne G, Legrand P, Ortiz-Lombardia
561 M, Campanacci V, van Sinderen D, Cambillau C. 2012. Structure of the phage TP901-1 1.8 MDa
562 baseplate suggests an alternative host adhesion mechanism. *Proceedings of the National Academy of
563 Sciences of the United States of America*, 2012/05/23 ed. 109:8954–8.

564 25. Pell LG, Kanelis V, Donaldson LW, Lynne Howell P, Davidson AR. 2009. The phage λ major tail
565 protein structure reveals a common evolution for long-tailed phages and the type VI bacterial secretion
566 system. *Proc Natl Acad Sci USA* 106:4160–4165.

567 26. Cardarelli L, Pell LG, Neudecker P, Pirani N, Liu A, Baker LA, Rubinstein JL, Maxwell KL,
568 Davidson AR. 2010. Phages have adapted the same protein fold to fulfill multiple functions in virion
569 assembly. *Proc Natl Acad Sci USA* 107:14384–14389.

570 27. Veesler D, Robin G, Lichière J, Auzat I, Tavares P, Bron P, Campanacci V, Cambillau C. 2010.
571 Crystal structure of bacteriophage SPP1 distal tail protein (gp19.1): a baseplate hub paradigm in gram-
572 positive infecting phages. *J Biol Chem* 285:36666–36673.

573 28. Bárdy P, Füzik T, Hrebík D, Pantůček R, Thomas Beatty J, Plevka P. 2020. Structure and
574 mechanism of DNA delivery of a gene transfer agent. *Nat Commun* 11:3034.

575 29. Kanamaru S, Leiman PG, Kostyuchenko VA, Chipman PR, Mesyanzhinov VV, Arisaka F,
576 Rossmann MG. 2002. Structure of the cell-puncturing device of bacteriophage T4. *Nature* 415:553–557.

577 30. Goulet A and C C. 2022. Present impact of AlphaFold2 revolution on structural biology, and an
578 illustration with the structure prediction of the bacteriophage J-1 host adhesion device. *Frontiers in
579 Molecular Biosciences* in press.

580 31. Stockdale SR, Mahony J, Courtin P, Chapot-Chartier M-P, van Pijkeren J-P, Britton RA, Neve H,
581 Heller KJ, Aideh B, Vogensen FK, van Sinderen D. 2013. The lactococcal phages Tuc2009 and TP901-1
582 incorporate two alternate forms of their tail fiber into their virions for infection specialization. *J Biol
583 Chem* 288:5581–5590.

584 32. Linares R, Arnaud C-A, Degroux S, Schoehn G, Breyton C. 2020. Structure, function and
585 assembly of the long, flexible tail of siphophages. *Curr Opin Virol* 45:34–42.

586 33. Sao-Jose C, Lhuillier S, Lurz R, Melki R, Lepault J, Santos MA, Tavares P. 2006. The
587 ectodomain of the viral receptor YueB forms a fiber that triggers ejection of bacteriophage SPP1 DNA. *J
588 Biol Chem*, 2006/02/17 ed. 281:11464–70.

589 34. Hayes S, Vincentelli R, Mahony J, Nauta A, Ramond L, Lugli GA, Ventura M, van Sinderen D,
590 Cambillau C. 2018. Functional carbohydrate binding modules identified in evolved dits from siphophages
591 infecting various Gram-positive bacteria: Functional carbohydrate binding modules identified in evolved
592 dits. *Mol Microbiol* 110:777–795.

593 35. Hayes S, Mahony J, Vincentelli R, Ramond L, Nauta A, van Sinderen D, Cambillau C. 2019.
594 Ubiquitous Carbohydrate Binding Modules Decorate 936 Lactococcal Siphophage Virions. *Viruses* 11.

595 36. Zimmermann L, Stephens A, Nam S-Z, Rau D, Kübler J, Lozajic M, Gabler F, Söding J, Lupas
596 AN, Alva V. 2018. A Completely Reimplemented MPI Bioinformatics Toolkit with a New HHpred
597 Server at its Core. *J Mol Biol* 430:2237–2243.

598 37. Yang J, Yan R, Roy A, Xu D, Poisson J, Zhang Y. 2015. The I-TASSER Suite: protein structure
599 and function prediction. *Nat Methods* 12:7–8.

600 38. Akdel M, Pires DEV, Pardo EP, Jänes J, Zalevsky AO, Mészáros B, Bryant P, Good LL,
601 Laskowski RA, Pozzati G, Shenoy A, Zhu W, Kundrotas P, Serra VR, Rodrigues CHM, Dunham AS,
602 Burke D, Borkakoti N, Velankar S, Frost A, Basquin J, Lindorff-Larsen K, Bateman A, Kajava AV,
603 Valencia A, Ovchinnikov S, Durairaj J, Ascher DB, Thornton JM, Davey NE, Stein A, Elofsson A, Croll
604 TI, Beltrao P. 2022. A structural biology community assessment of AlphaFold2 applications. *Nat Struct
605 Mol Biol* 29:1056–1067.

606 39. Gao M, Nakajima An D, Parks JM, Skolnick J. 2022. AF2Complex predicts direct physical
607 interactions in multimeric proteins with deep learning. *Nat Commun* 13:1744.

- 608 40. Evans R et al. 2021. Protein complex prediction with AlphaFold-Multimer. BioRxiv
609 <https://doi.org/10.1101/2021.10.04.463034>;
- 610 41. Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K, Bates
611 R, Židek A, Potapenko A, Bridgland A, Meyer C, Kohl SAA, Ballard AJ, Cowie A, Romera-Paredes B,
612 Nikolov S, Jain R, Adler J, Back T, Petersen S, Reiman D, Clancy E, Zielinski M, Steinegger M,
613 Pacholska M, Berghammer T, Bodenstein S, Silver D, Vinyals O, Senior AW, Kavukcuoglu K, Kohli P,
614 Hassabis D. 2021. Highly accurate protein structure prediction with AlphaFold. *Nature* 596:583–589.
- 615 42. Goulet A, Spinelli S, Mahony J, Cambillau C. 2020. Conserved and Diverse Traits of Adhesion
616 Devices from Siphoviridae Recognizing Proteinaceous or Saccharidic Receptors. *Viruses*, 2020/05/10 ed.
617 12.
- 618 43. Lombard V, Golaconda Ramulu H, Drula E, Coutinho PM, Henrissat B. 2014. The carbohydrate-
619 active enzymes database (CAZy) in 2013. *Nucleic Acids Res* 42:D490-495.
- 620 44. Buschiazzi A, Tavares GA, Campetella O, Spinelli S, Cremona ML, París G, Amaya MF, Frasca
621 AC, Alzari PM. 2000. Structural basis of sialyltransferase activity in trypanosomal sialidases. *EMBO J*
622 19:16–24.
- 623 45. Dunne M, Hupfeld M, Klumpp J, Loessner M. 2018. Molecular Basis of Bacterial Host
624 Interactions by Gram-Positive Targeting Bacteriophages. *Viruses* 10:397.
- 625 46. Buglino J, Shen V, Hakimian P, Lima CD. 2002. Structural and biochemical analysis of the Obg
626 GTP binding protein. *Structure* 10:1581–1592.
- 627 47. Mus F, Eilers BJ, Alleman AB, Kabasakal BV, Wells JN, Murray JW, Nocek BP, DuBois JL,
628 Peters JW. 2017. Structural Basis for the Mechanism of ATP-Dependent Acetone Carboxylation. *Sci Rep*
629 7:7234.
- 630 48. Pentelute BL, Gates ZP, Tereshko V, Dashnau JL, Vanderkooi JM, Kossiakoff AA, Kent SBH.
631 2008. X-ray structure of snow flea antifreeze protein determined by racemic crystallization of synthetic
632 protein enantiomers. *J Am Chem Soc* 130:9695–9701.
- 633 49. Warkentin E, Weidenweber S, Schühle K, Demmer U, Heider J, Ermler U. 2017. A rare
634 polyglycine type II-like helix motif in naturally occurring proteins. *Proteins* 85:2017–2023.
- 635 50. Weidenweber S, Schühle K, Demmer U, Warkentin E, Ermler U, Heider J. 2017. Structure of the
636 acetophenone carboxylase core complex: prototype of a new class of ATP-dependent
637 carboxylases/hydrolases. *Sci Rep* 7:39674.
- 638 51. Dunne M, Denyes JM, Arndt H, Loessner MJ, Leiman PG, Klumpp J. 2018. Salmonella Phage
639 S16 Tail Fiber Adhesin Features a Rare Polyglycine Rich Domain for Host Recognition. *Structure*
640 26:1573-1582.e4.
- 641 52. Cuskin F, Flint JE, Gloster TM, Morland C, Baslé A, Henrissat B, Coutinho PM, Strazzulli A,
642 Solovyova AS, Davies GJ, Gilbert HJ. 2012. How nature can exploit nonspecific catalytic and
643 carbohydrate binding modules to create enzymatic specificity. *Proc Natl Acad Sci U S A* 109:20889–
644 20894.
- 645 53. Sulák O, Cioci G, Delia M, Lahmann M, Varrot A, Imberty A, Wimmerová M. 2010. A TNF-like
646 trimeric lectin domain from Burkholderia cenocepacia with specificity for fucosylated human histo-blood
647 group antigens. *Structure* 18:59–72.
- 648 54. Squeglia F, Maciejewska B, Łątka A, Ruggiero A, Briers Y, Drulis-Kawa Z, Berisio R. 2020.
649 Structural and Functional Studies of a Klebsiella Phage Capsule Depolymerase Tailspike: Mechanistic
650 Insights into Capsular Degradation. *Structure* 28:613-624.e4.
- 651 55. Henshaw J, Horne-Bitschy A, van Bueren AL, Money VA, Bolam DN, Czjzek M, Ekborg NA,
652 Weiner RM, Hutcheson SW, Davies GJ, Boraston AB, Gilbert HJ. 2006. Family 6 carbohydrate binding
653 modules in beta-agarases display exquisite selectivity for the non-reducing termini of agarose chains. *J*
654 *Biol Chem* 281:17099–17107.
- 655 56. Garcia-Doval C, Castón JR, Luque D, Granell M, Otero JM, Llamas-Saiz AL, Renouard M,
656 Boulanger P, van Raaij MJ. 2015. Structure of the Receptor-Binding Carboxy-Terminal Domain of the
657 Bacteriophage T5 L-Shaped Tail Fibre with and without Its Intra-Molecular Chaperone. *Viruses* 7:6424–
658 6440.
- 659 57. Steinbacher S, Seckler R, Miller S, Steipe B, Huber R, Reinemer P. 1994. Crystal Structure of P22
660 Tailspike Protein: Interdigitated Subunits in a Thermostable Trimer. *Science* 265:383–386.
- 661 58. Hawkins NC, Kizziah JL, Hatoum-Aslan A, Dokland T. 2022. Structure and host specificity of

662 Staphylococcus epidermidis bacteriophage Andhra. *Sci Adv* 8:eade0459.

663 59. Dunne M, Rupf B, Tala M, Qabrati X, Ernst P, Shen Y, Sumrall E, Heeb L, Pluckthun A,
664 Loessner MJ, Kilcher S. 2019. Reprogramming Bacteriophage Host Range through Structure-Guided
665 Design of Chimeric Receptor Binding Proteins. 5. *Cell Rep*, 2019/10/31 ed. 29:1336-1350 e4.

666 60. Farenc C, Spinelli S, Vinogradov E, Tremblay D, Blangy S, Sadovskaya I, Moineau S, Cambillau
667 C. 2014. Molecular insights on the recognition of a *Lactococcus lactis* cell wall pellicle by the phage
668 1358 receptor binding protein. *J Virol* 88:7005–7015.

669 61. Bompard-Gilles C, Remaut H, Villeret V, Prangé T, Fanuel L, Delmarcelle M, Joris B, Frère J,
670 Van Beeumen J. 2000. Crystal structure of a D-aminopeptidase from *Ochrobactrum anthropi*, a new
671 member of the “penicillin-recognizing enzyme” family. *Structure* 8:971–980.

672 62. Drexler K, Riede I, Montag D, Eschbach M-L, Henning U. 1989. Receptor specificity of the
673 *Escherichia coli* T-even type phage Ox2. *Journal of Molecular Biology* 207:797–803.

674 63. Trojet SN, Caumont-Sarcos A, Perrody E, Comeau AM, Krisch HM. 2011. The gp38 adhesins of
675 the T4 superfamily: a complex modular determinant of the phage’s host specificity. *Genome Biol Evol*
676 3:674–686.

677 64. Lavelle K, Murphy J, Fitzgerald B, Lugli GA, Zomer A, Neve H, Ventura M, Franz CM,
678 Cambillau C, van Sinderen D, Mahony J. 2018. A Decade of *Streptococcus thermophilus* Phage
679 Evolution in an Irish Dairy Plant. *Appl Environ Microbiol* 84:e02855-17.

680 65. Marti R, Zurfluh K, Hagens S, Pianezzi J, Klumpp J, Loessner MJ. 2013. Long tail fibres of the
681 novel broad-host-range T-even bacteriophage S16 specifically recognize *Salmonella* OmpC: T4-like
682 *Salmonella* phage S16. *Molecular Microbiology* 87:818–834.

683 66. Chiaradia L, Lefebvre C, Parra J, Marcoux J, Burlet-Schiltz O, Etienne G, Tropis M, Daffé M.
684 2017. Dissecting the mycobacterial cell envelope and defining the composition of the native
685 mycomembrane. *Sci Rep* 7:12807.

686 67. Chen J, Kriakov J, Singh A, Jacobs WR, Besra GS, Bhatt A. 2009. Defects in glycopeptidolipid
687 biosynthesis confer phage I3 resistance in *Mycobacterium smegmatis*. *Microbiology* 155:4050–4057.

688 68. Furuchi A, Tokunaga T. 1972. Nature of the receptor substance of *Mycobacterium smegmatis* for
689 D4 bacteriophage adsorption. *J Bacteriol* 111:404–411.

690 69. Imaeda T, Blas FS. 1969. Adsorption of Mycobacteriophage on Cell-wall Components. *Journal of*
691 *General Virology* 5:493–498.

692 70. Bebeacua C, Tremblay D, Farenc C, Chapot-Chartier M-P, Sadovskaya I, van Heel M, Veessler D,
693 Moineau S, Cambillau C. 2013. Structure, adsorption to host, and infection mechanism of virulent
694 lactococcal phage p2. *J Virol* 87:12302–12312.

695 71. Ainsworth S, Sadovskaya I, Vinogradov E, Courtin P, Guerardel Y, Mahony J, Grard T,
696 Cambillau C, Chapot-Chartier M-P, van Sinderen D. 2014. Differences in lactococcal cell wall
697 polysaccharide structure are major determining factors in bacteriophage sensitivity. *mBio* 5:e00880-
698 00814.

699 72. Goulet A, Spinelli S, Mahony J, Cambillau C. 2020. Conserved and Diverse Traits of Adhesion
700 Devices from Siphoviridae Recognizing Proteinaceous or Saccharidic Receptors. *Viruses* 12.

701 73. Jacobs-Sera D, Marinelli LJ, Bowman C, Broussard GW, Guerrero Bustamante C, Boyle MM,
702 Petrova ZO, Dedrick RM, Pope WH, Science Education Alliance Phage Hunters Advancing Genomics
703 And Evolutionary Science Sea-Phages Program, Modlin RL, Hendrix RW, Hatfull GF. 2012. On the
704 nature of mycobacteriophage diversity and host preference. *Virology* 434:187–201.

705 74. Bertozzi Silva J, Storms Z, Sauvageau D. 2016. Host receptors for bacteriophage adsorption.
706 *FEMS Microbiology Letters* 363:fnw002.

707 75. Shkoporov AN, Hill C. 2019. Bacteriophages of the Human Gut: The “Known Unknown” of the
708 Microbiome. *Cell Host Microbe* 25:195–209.

709 76. Emsley P, Lohkamp B, Scott WG, Cowtan K. 2010. Features and development of Coot. *Acta*
710 *Crystallogr D Biol Crystallogr* 66:486–501.

711 77. Holm L. 2020. Using Dali for Protein Structure Comparison. *Methods Mol Biol* 2112:29–42.

712 78. Pettersen EF, Goddard TD, Huang CC, Meng EC, Couch GS, Croll TI, Morris JH, Ferrin TE.
713 2021. UCSF ChimeraX: Structure visualization for researchers, educators, and developers. *Protein Sci*
714 30:70–82.

715