

Supplementary Information for

Rates of mutation and recombination in *Siphoviridae* phage genome evolution over three decades

Anne Kupczok, Horst Neve, Kun D. Huang, Marc P. Hoepfner, Knut J. Heller, Charles M. A. P. Franz, Tal Dagan

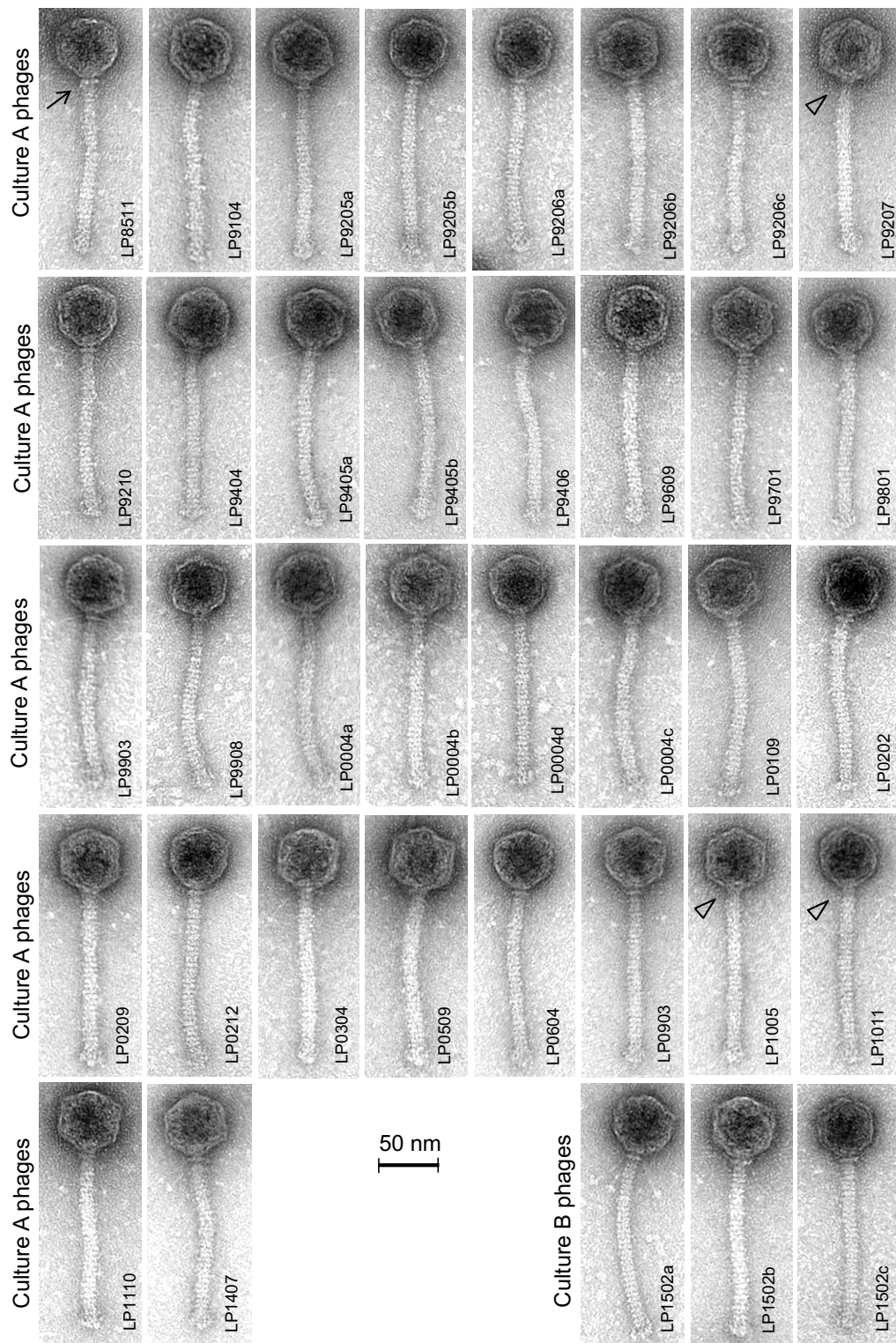


Figure S1: Electron micrographs of phages sequenced in this study. Phage names indicate isolation time (in YYMM format). All 34 culture A and all three phages from culture B are members of the lactococcal 936 group of *Siphoviridae* phages with isometric heads (ca. 55 nm diameter) and flexible, non-contractile tails (ca. 160 nm length). Except for few phages,

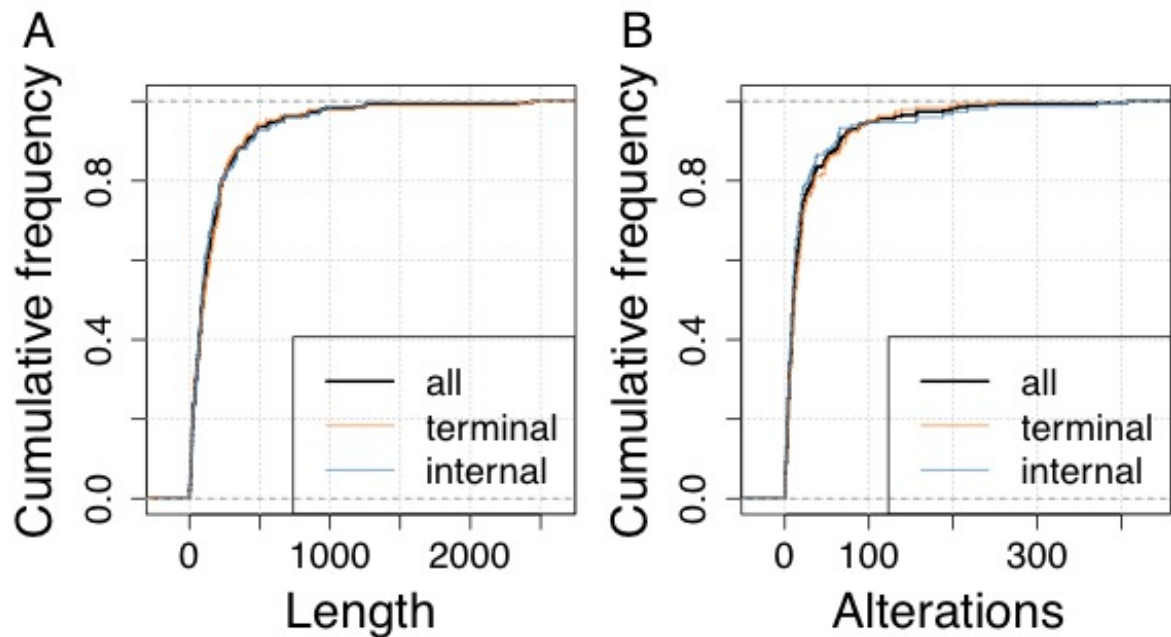


Figure S3: Recombination statistics. (A) Length distribution of 345 detected recombined segments, mean=176.9, median=96, no difference between terminal and internal branches ($p > 0.05$, Anderson-Darling Test). (B) Distribution of number of nucleotide alterations introduced by one recombination event, mean=26.35, median=11, no difference between terminal and internal nodes ($p > 0.05$, Anderson-Darling Test). Sequencing errors in single strains would lead to an overrepresentation of recombination events along terminal branches. Since the distribution of recombination length and altered nucleotides is not significantly different between terminal and internal branches, we can exclude the presence of strain-related artefacts in the data. Posterior mean values of parameters estimated by ClonalFrameML are: ratio of recombination to mutation rate $R/\theta = 1.247$, import length $\delta = 111.2\text{nt}$, nucleotide distance of imports $\nu = 0.1694$ differences per nucleotide. Thus, the relative effect of recombination to mutation is $r/m = (R/\theta) \times \delta \times \nu = 23.50$.

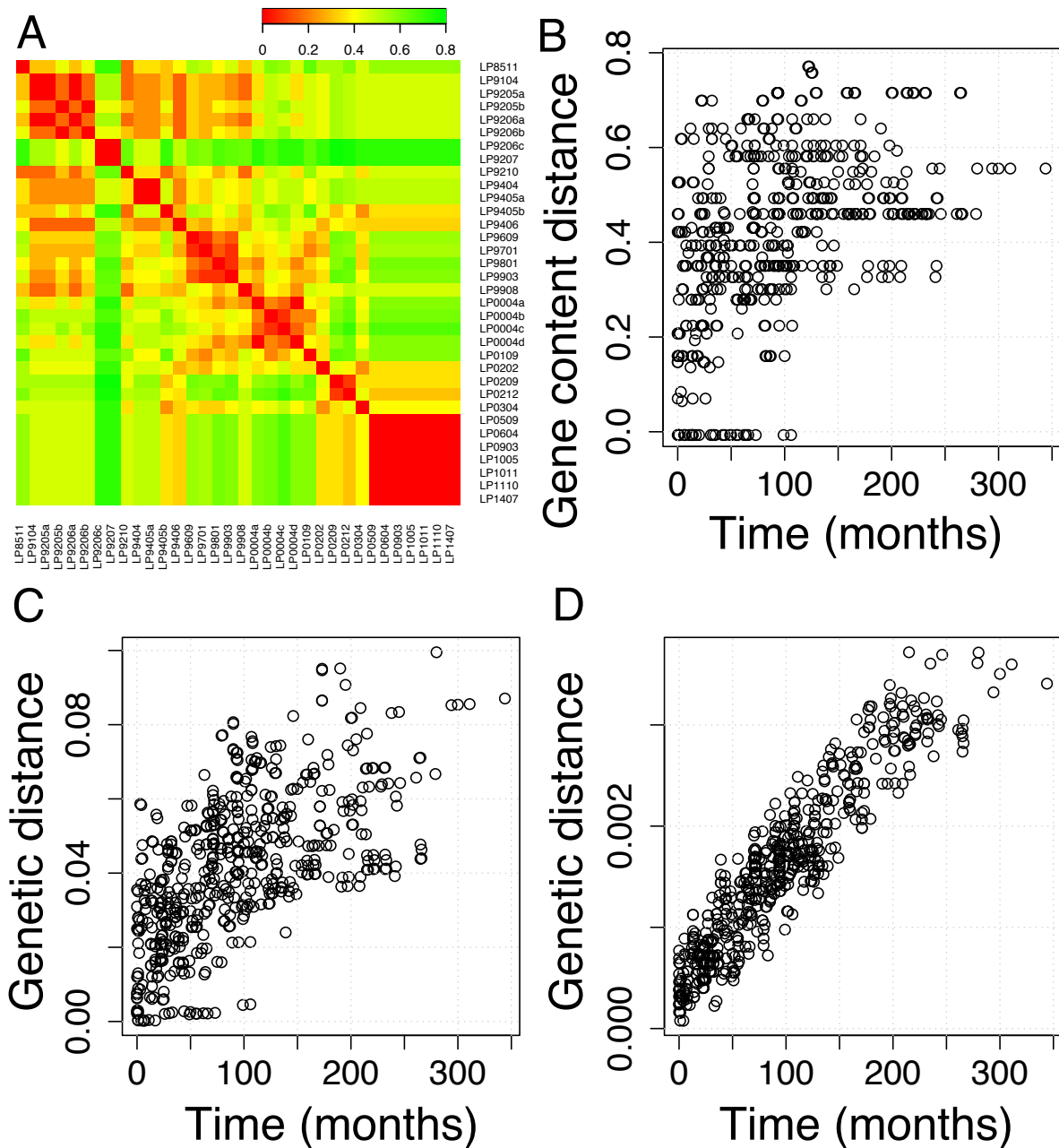


Figure S5: Gene indel family evolution over time. (A) Heatmap of pairwise gene content distance based on 25 gene indel families. We observe that the gene content distance is lower between isolates close in time, although there appears to be no gradual change of the gene content distance with time. (B) Pairwise gene content distance against pairwise isolation time differences ($r^2=0.19$). In accordance with a, we observe a weak correlation of pairwise gene content distance with time distance. (C) Pairwise genetic distances based on whole-genome alignment against pairwise isolation time differences ($r^2=0.39$). (D) Genetic distances based on whole-genome alignment after recombination masking against pairwise isolation time differences ($r^2=0.87$). In comparison to the gene content distances, the genetic distances based on the whole-genome alignment have a stronger correlation with time, and this correlation even increased when recombinations were masked.

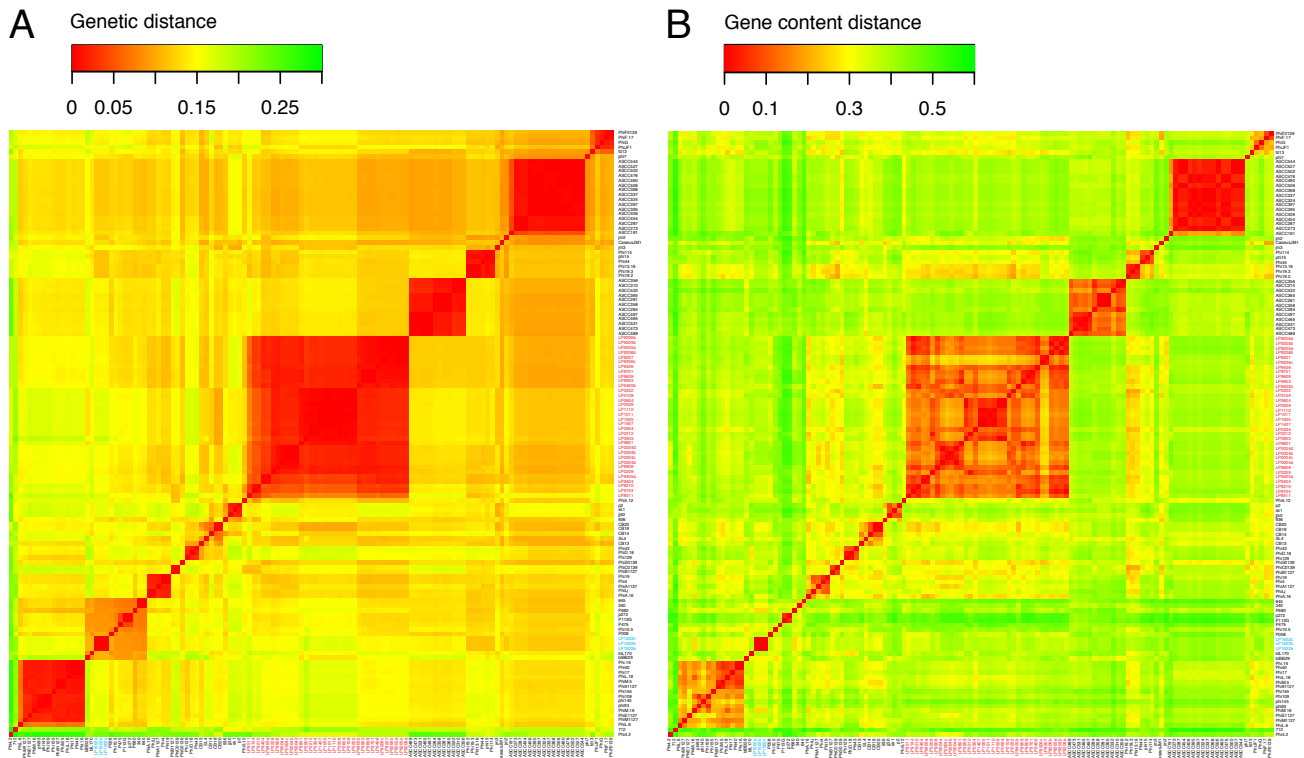


Figure S6: Heatmap of (A) genetic distances, and (B) gene content distances based on Jaccard index, for all phages of the 936 group. Genome order is the same in both matrices, based on clustering of genetic distances. Culture A phages are highlighted in red and culture B phages are highlighted in blue.

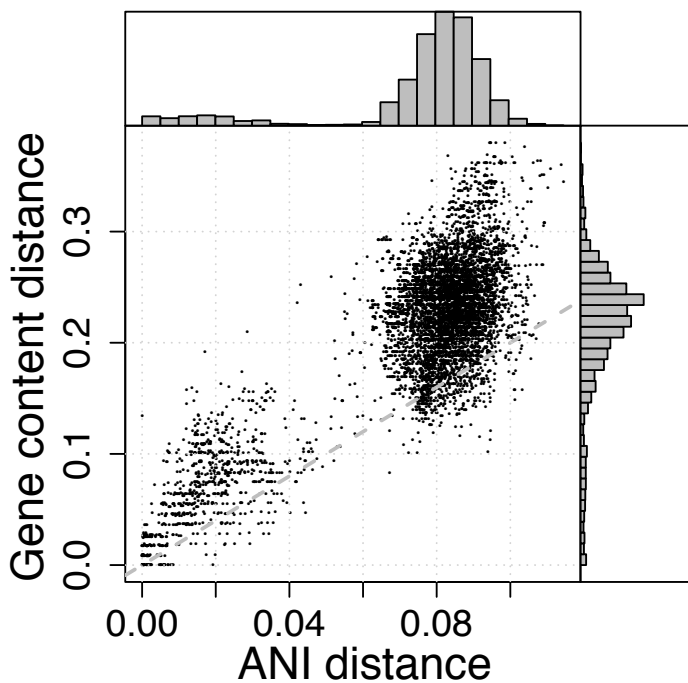


Figure S7: Association between pairwise gene content distance and ANI distance for all *L. lactis* phages of the 936 group. Gene content distances are based on the average proportion of shared genes. These statistics are similar to the ones used by (Mavrigh and Hatfull 2017). The clustering pattern observed in Fig. 6 is also observed with this statistic. For an ANI distance cutoff of 0.05, low distance pairs have an average gene content distance of 0.06947 and r^2 of 0.3925 ($n=881$, $p<10^{-6}$), and high distance pairs have an average gene content distance of 0.2331 and r^2 of 0.07739 ($n=7120$, $p<10^{-6}$). We cannot distinguish between phages of high or low gene flux modes (Mavrigh and Hatfull 2017) due to sampling density.

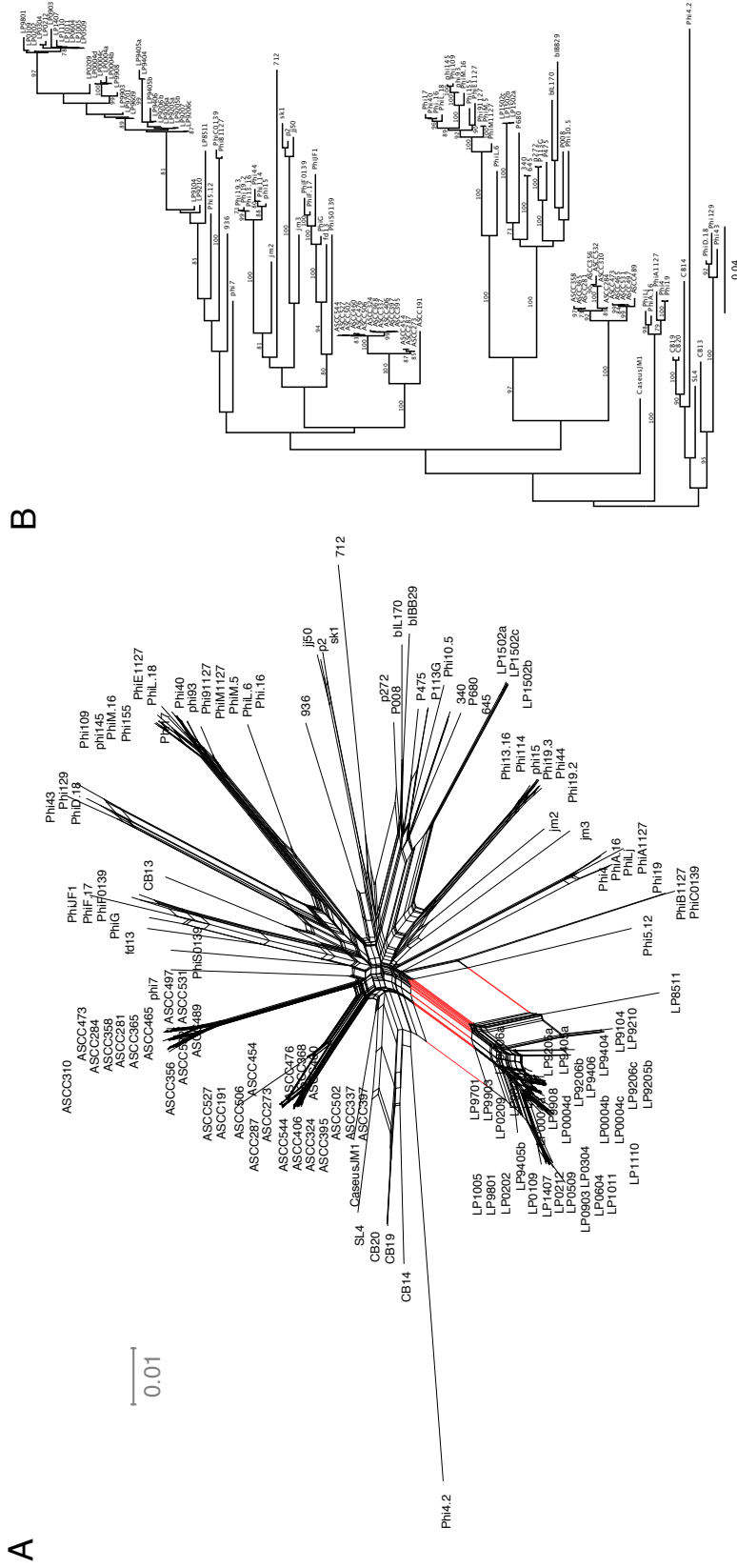
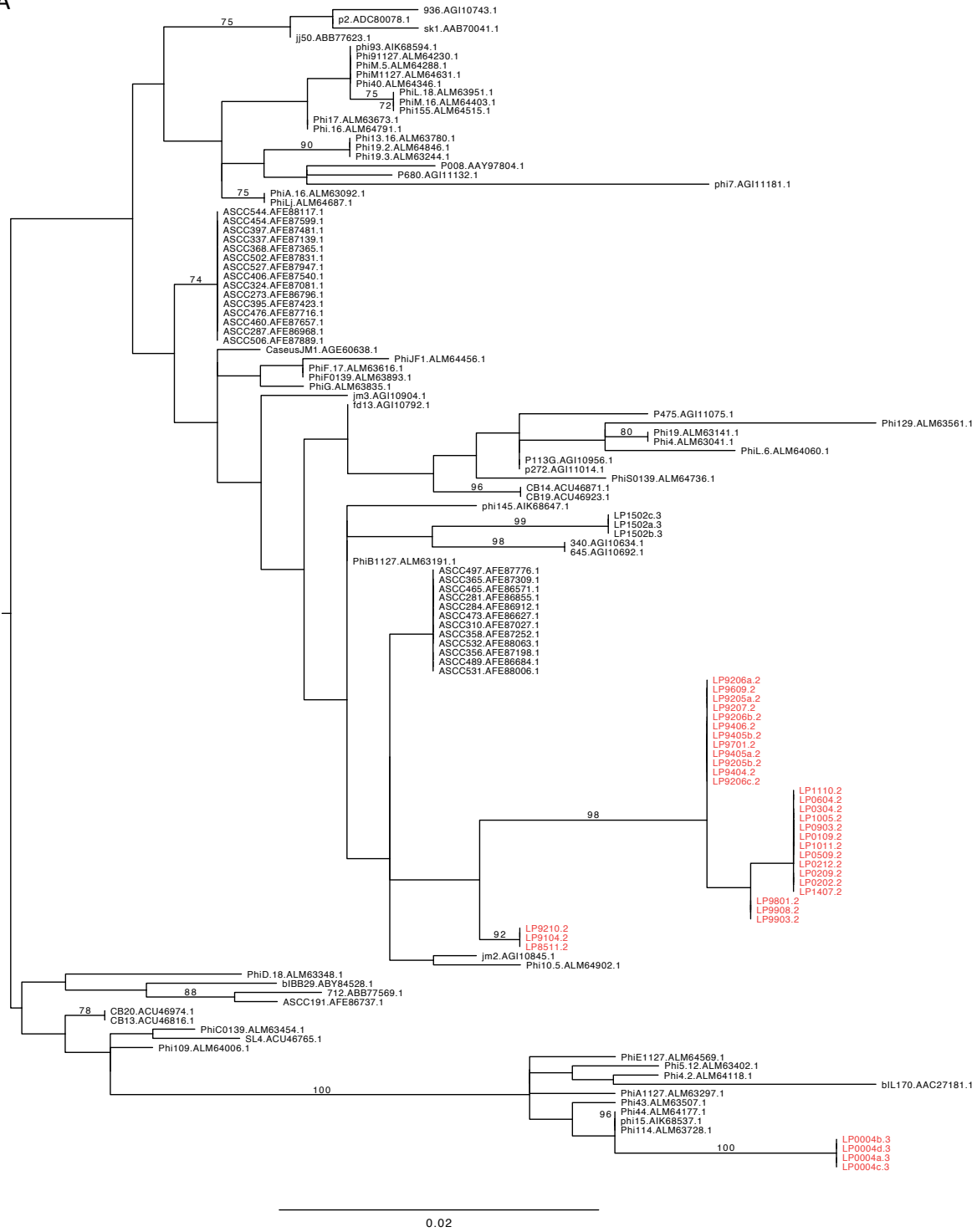
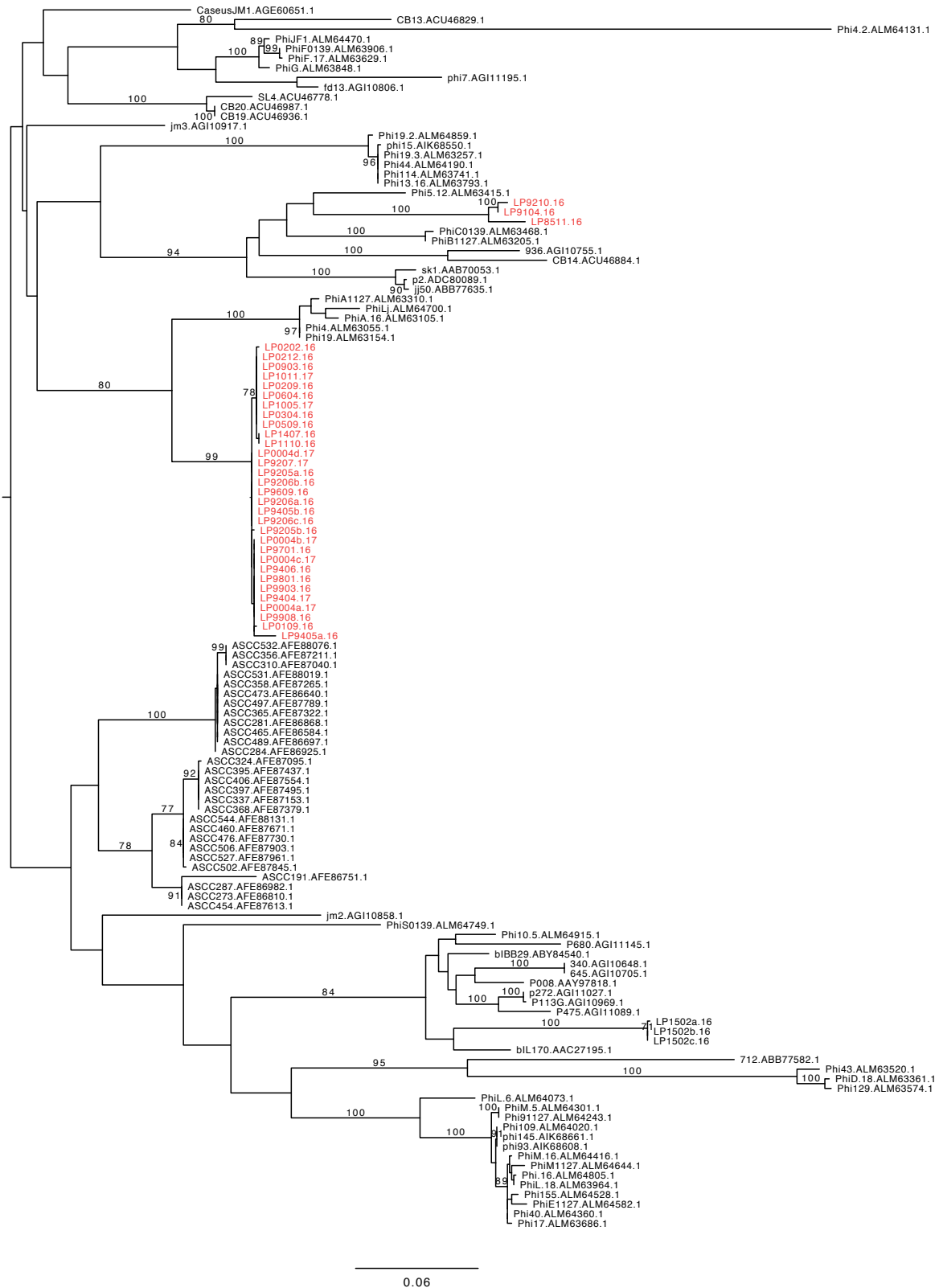


Figure S8: Relationships of *L. lactis* phages of the 936 group based on core gene alignment. (A) Neighbor-net based on ordinary least squares distances calculated with SplitsTree (Huson and Bryant 2006). The split that supports the monophyly of culture A phages is marked in red. (B) ML phylogeny, only bootstrap values above 70 are displayed. Note that this is an unrooted tree. Culture A phages are monophyletic in the tree whereas culture B phages group in a notably different region.

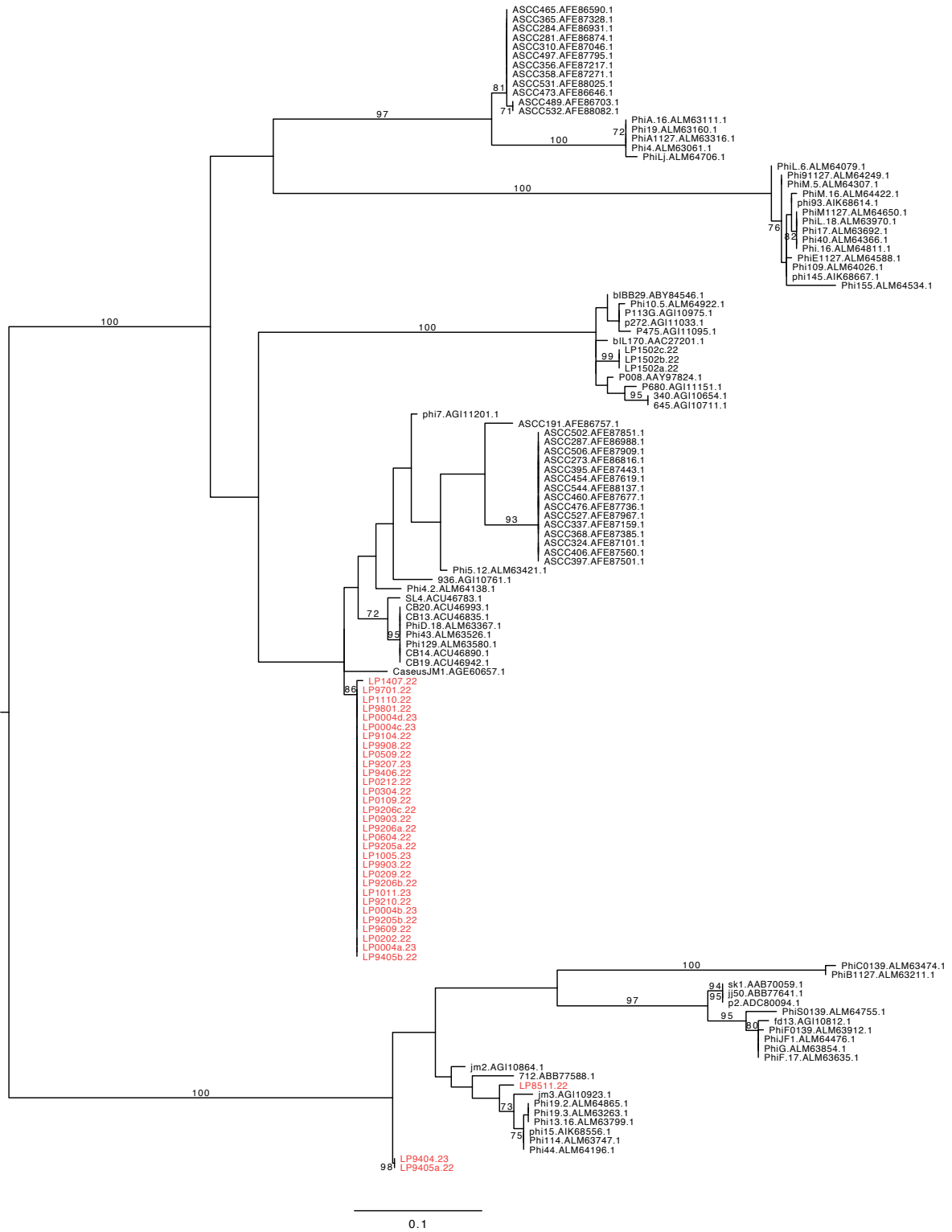
A



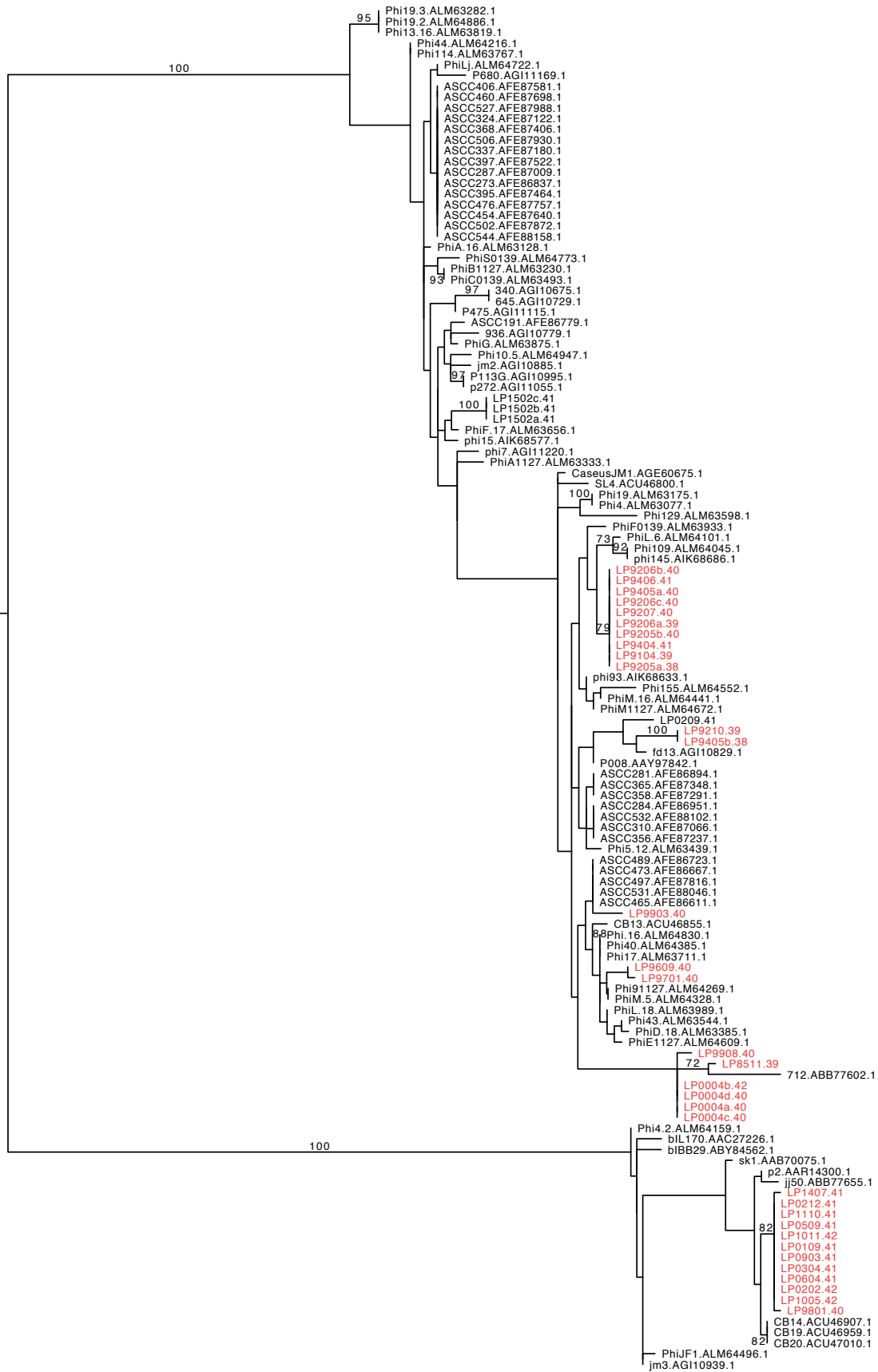
B



C

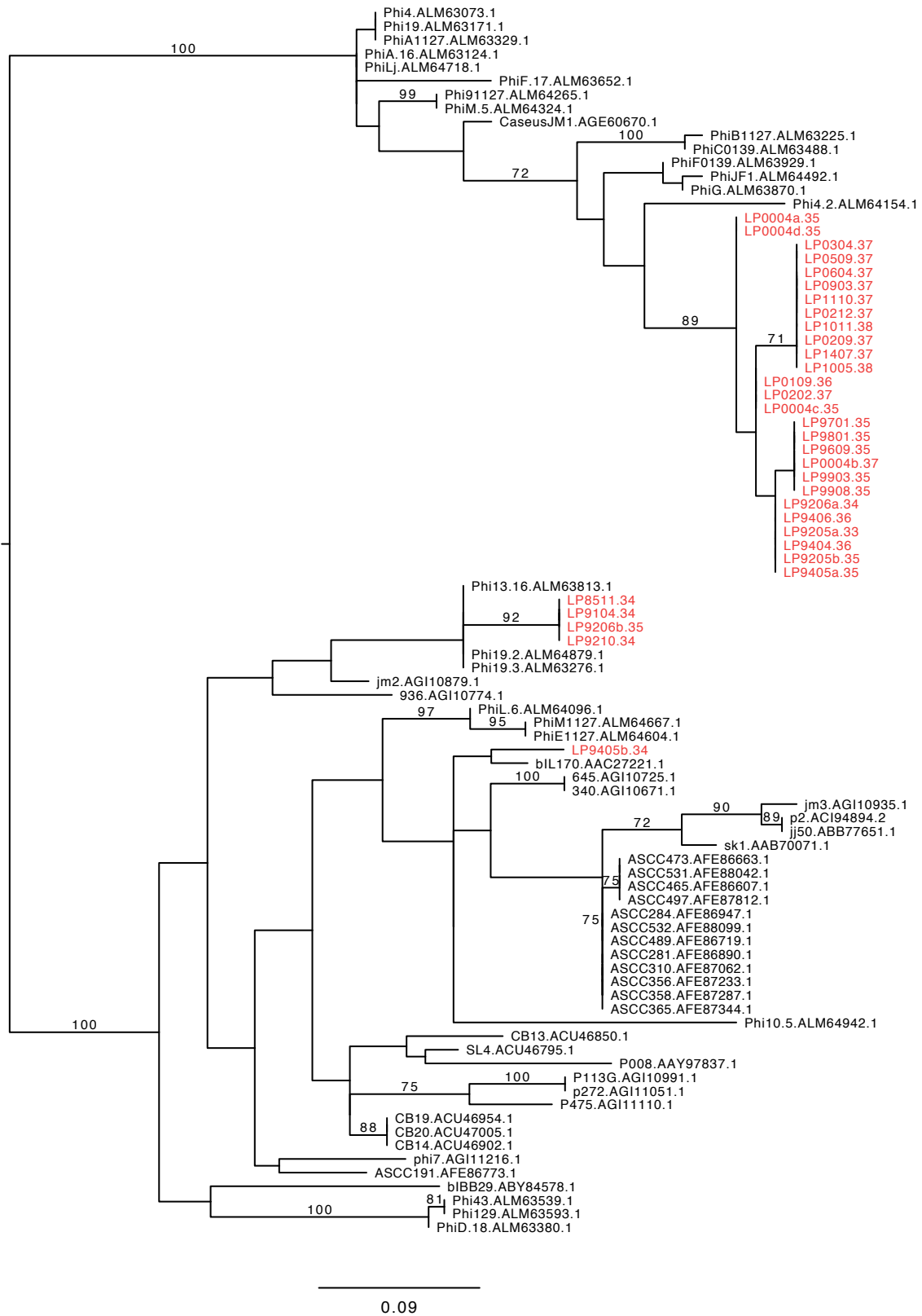


D

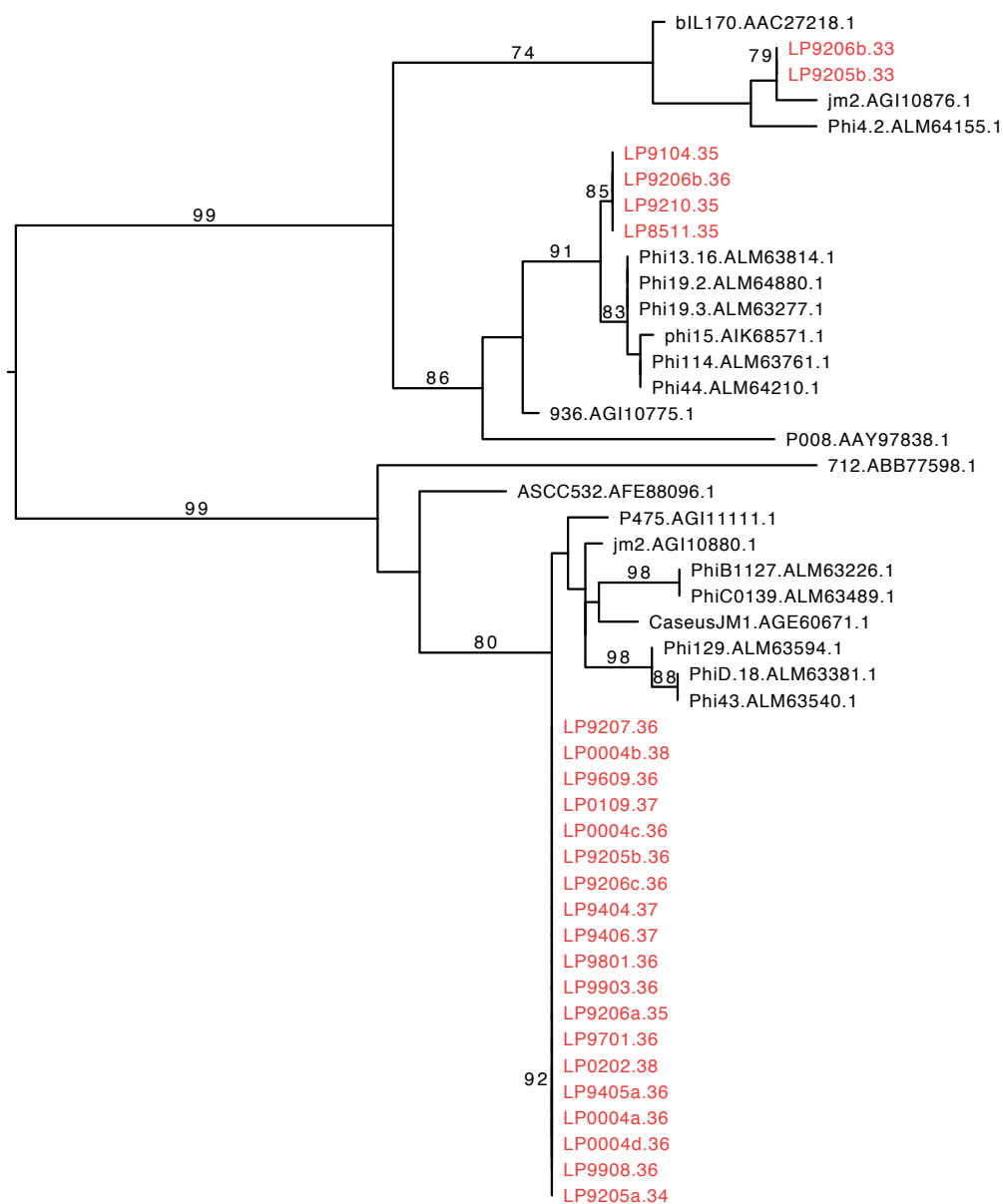


0.2

E



F



0.08

G

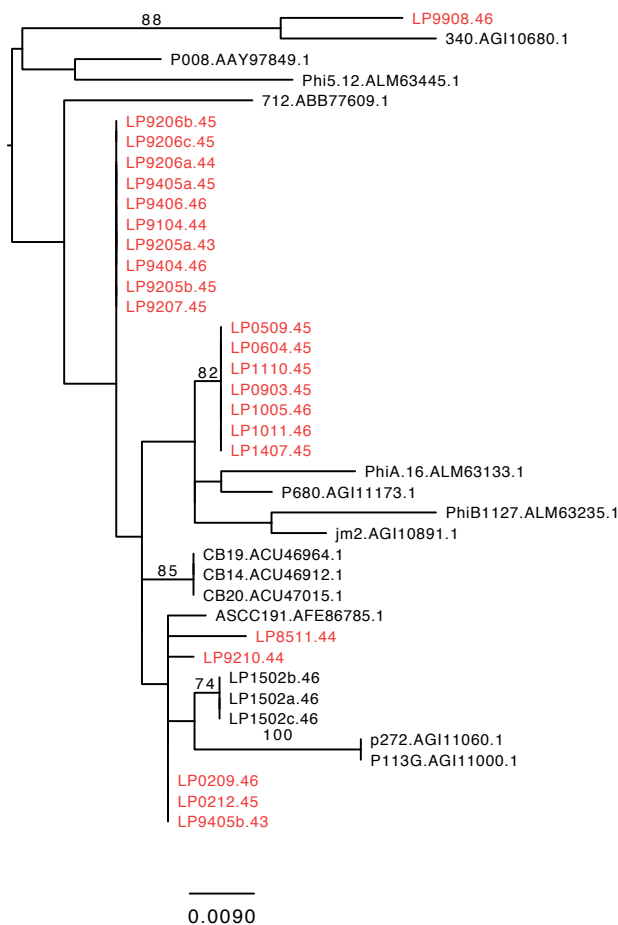


Figure S9: Gene phylogenies that conflict the monophyly of culture A lactococcal phages.

All trees are unrooted. Proteins from culture A and culture B phages are labelled by the name and CDS number. Remaining proteins are labelled by name and protein ID. Culture A phages are highlighted in red. Only bootstrap support values above 70 are displayed. All phylogenies are shown where the monophyly of culture A phages is contradicted by a branch with bootstrap support of more than 70. (A) Family 2 – Phage terminase, large subunit. (B) Family 16 – Phage tail length tape measure protein. (C) Family 22 – Phage lysin. (D) Family 31 – Sak3. Family 31 shows similarity to ORF35 from phage p2 (blastp e-value $\leq 2 \times 10^{-64}$) that was described as Sak3 (Bouchard and Moineau 2004). Sak3 is a DNA-single strand annealing protein involved in homologous recombination (Bouchard and Moineau 2004; Scaltriti et al. 2011). Sequences from family 31 were also detected as homologous to the eukaryotic recombination protein Rad52 using virfam (Lopes et al. 2010). Proteins from the Rad52 family promote recombination between phages and prophages by facilitating homologous recombination between divergent sequences (De Paepe et al. 2014). The alignment exhibits a high amount of conflict as displayed by the split network (Fig. S10A). The conflict can be traced back to different regions of the protein (Fig. S10). (E) Family 60 – conserved hypothetical protein. (F) Family 61 – HNH homing endonuclease. (G) Family 66 – DNA polymerase. Two protein families (64,66) are predicted to function as DNA polymerases, but they show less than 50% global protein identity. Exactly one gene from each family is present in each isolate of culture A phages at a conserved genomic location. In addition, family 66 shows a conflicting phylogeny. This implies the frequent exchange of multiple versions of DNA polymerase between phages of the 936 group.

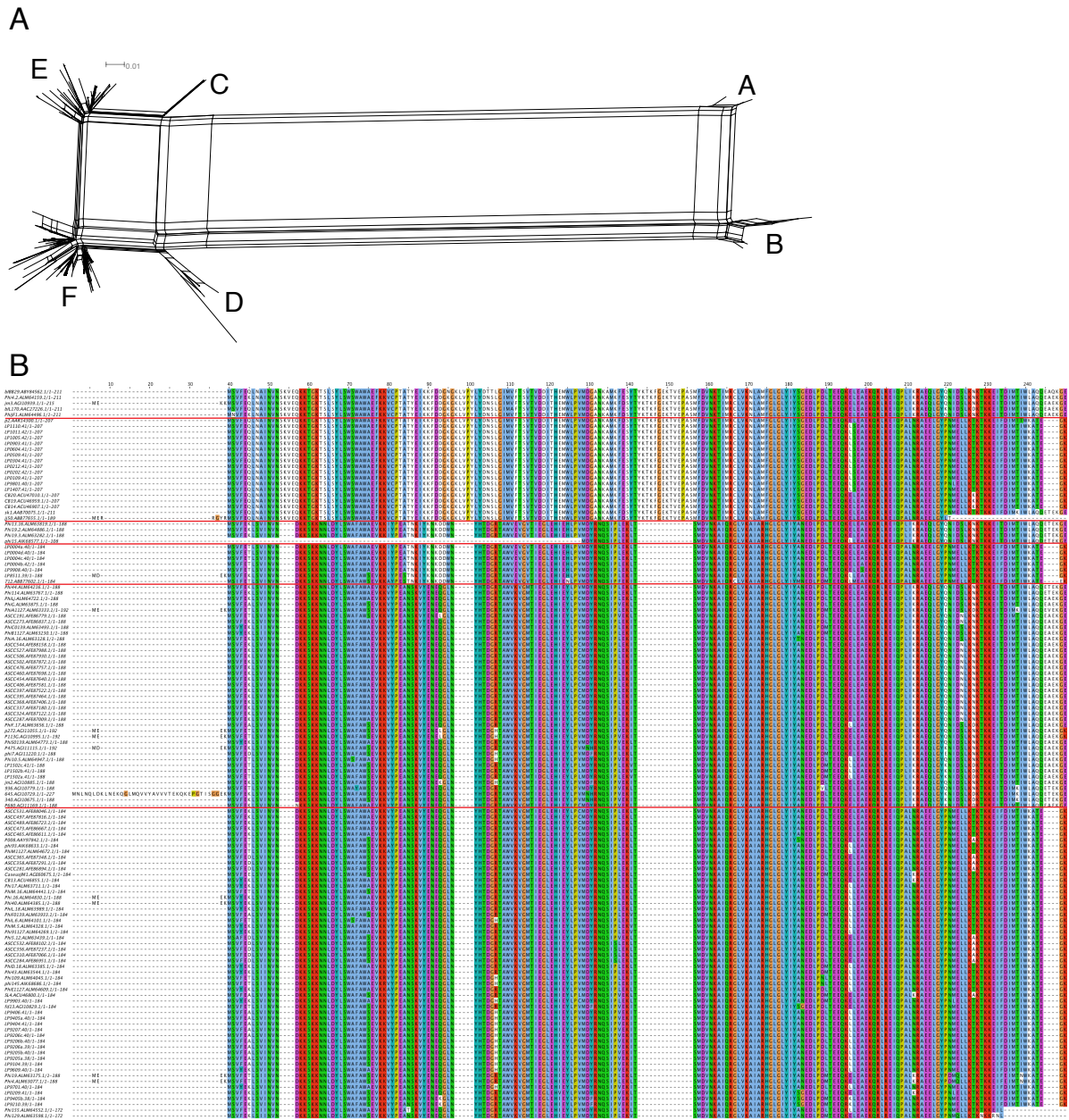


Figure S10: Alignment of Sak3 proteins from phages of the 936 group. (A) Neighbor-net based on ordinary least squares distances calculated with SplitsTree (Huson and Bryant 2006). (B) Alignment delimited into the taxon groups A, B, C, D, E, and F (top to bottom) as marked in the neighbor-net. Differences that split taxa A and B from the rest can be found in the N-terminal region (alignment positions 40-184), whereas differences that split taxa A, C, and E from B, D, and F can be found in the C-terminal region (alignment positions 212-250). ORF35 from phage p2 belongs to group B. Its N-terminal region contains the domains for oligomerization, DNA binding, and the motifs for ATPase activity that are important for the stimulation of RecA, whereas its C-terminal region might constitute the interaction domain with RecA (Scaltriti et al. 2011).

Table S1: Protein families.

Protein Family	Samples	Alignment length	% Recombination	Type	Annotation	Annotation evidence
1	34	174	7.5%	core (all)	Phage terminase, small subunit	RAST
2	34	540	72.0%	core (all)	Phage terminase, large subunit	RAST
3	34	94	0.0%	core (all)	Phage-associated homing endonuclease	RAST
4	34	378	0.0%	core (all)	Phage portal protein	RAST
5	34	178	100.0%	core (all)	Putative capsid maturation protein	Similarity to YP_762516.1
6	34	393	100.0%	core (all)	Major capsid protein	Similarity to ALM64064.1
7	34	87	0.0%	core (all)	Conserved hypothetical protein	
8	34	104	0.0%	core (all)	Probable head completion protein 1	Similarity to HCP1_BPLSK
9	34	113	0.0%	core (all)	Probable head completion protein 2	Similarity to HCP2_BPLSK
10	34	121	0.0%	core (all)	Probable tail terminator protein	Similarity to TTTP_BPLSK
11	34	656	33.8%	core	Neck passage structure	RAST, Similarity to ANS02693.1
12	34	301	72.1%	core (all)	Major tail protein	Similarity to ALM63411.1
13	34	228	71.5%	core	Tail protein extension	Similarity to ALM63412.1
14	34	91	0.0%	core (all)	Probable chaperone for the tape measure protein	Similarity to GP12_BPLSK
15	34	170	100.0%	core (all)	Conserved hypothetical protein	
16	34	919	44.0%	core (all)	Phage tail length tape measure protein	RAST
17	34	460	0.0%	core (all)	Distal tail protein	Similarity to ALM64750.1
18	34	379	0.0%	core	Tail associated lysin	Similarity to ALM63688.1
19	34	97	0.0%	core	Conserved hypothetical protein	
20	34	264	0.0%	core	Putative receptor binding protein	Similarity to AAT81500.1
21	34	117	0.0%	core (all)	Phage holin	RAST
22	34	263	3.8%	core (all)	Phage lysin	RAST
23	34	89	1.1%	core (all)	Putative transglycolase	Phagocyte based on similarity to sk1p21
24	34	115	58.3%	core (all)	Conserved hypothetical protein	
25	34	62	100.0%	core	Conserved hypothetical protein	
26	34	102	100.0%	core	SaV protein	Similarity to SAV_BPLSK
27	34	39	100.0%	core	Conserved hypothetical protein	
28	34	172	100.0%	core	Conserved hypothetical protein	
29	34	71	100.0%	core (all)	Conserved hypothetical protein	
30	34	119	100.0%	core (all)	Single stranded DNA-binding protein, phage associated	RAST
31	34	211	77.7%	core (all)	Sak3	Similarity to ALM64385.1
32	34	94	100.0%	core	Conserved hypothetical protein	
33	34	96	100.0%	core	Conserved hypothetical protein	
34	34	90	100.0%	core	Conserved hypothetical protein	
35	34	73	100.0%	core	Conserved hypothetical protein	

Protein Family	Samples	Alignment length	% Recombination	Type	Annotation	Annotation evidence
36	34	56	100.0%	core	Conserved hypothetical protein	
37	34	75	100.0%	core	Conserved hypothetical protein	
38	34	43	0.0%	core (all)	Conserved hypothetical protein	
39	34	160	16.3%	core (all)	Putative holliday junction endonuclease	Similarity to YP_762566.1
40	34	55	0.0%	core	Conserved hypothetical protein	
41	4	153	0.0%	gene indel	HNH homing endonuclease	RAST
42	1	-	-	nonsense	Neck passage structure	RAST
43	1	-	-	nonsense	Neck passage structure	RAST
44	1	-	-	nonsense	Neck passage structure	RAST
45	1	-	-	nonsense	Neck passage structure	RAST
46	1	-	-	gene indel	Putative lytic transglycosylase	Similarity to YP_002004010.1
47	31	53	0.0%	gene indel	Putative lytic transglycosylase	Similarity to B3TJQ7_9CAUD
48	19	59	0.0%	nonsense	Conserved hypothetical protein	
49	33	56	100.0%	nonsense	Conserved hypothetical protein	
50	7	82	100.0%	gene indel	Conserved hypothetical protein	
51	31	86	100.0%	gene indel	Conserved hypothetical protein	
52	13	81	0.0%	gene indel	Conserved hypothetical protein	
53	32	220	0.0%	gene indel	Methylase	Similarity to WP_062806946.1
54	2	107	0.0%	gene indel	Conserved hypothetical protein	
55	2	116	0.0%	gene indel	Putative excisionase	Similarity to A5GZ54_9CAUD
56	2	45	0.0%	nonsense	Methylase	Similarity to ALM64768.1
57	22	325	70.2%	gene indel	Methylase	Similarity to ALM64939.1
58	30	82	37.8%	gene indel	Conserved hypothetical protein	
59	32	57	0.0%	gene indel	Conserved hypothetical protein	
60	32	109	100.0%	gene indel	Conserved hypothetical protein	
61	23 (25)*	147	100.0%	gene indel	HNH homing endonuclease	RAST
62	11	169	100.0%	gene indel	HNH homing endonuclease	RAST
63	20	60	33.3%	gene indel	Conserved hypothetical protein	
64	11	316	2.8%	gene indel	DNA Polymerase	Similarity to AIK68583.1
65	5	145	0.0%	gene indel	HNH homing endonuclease	RAST
66	23	298	100.0%	gene indel	DNA polymerase	Similarity to YP_008320131.1
67	7	153	0.0%	gene indel	Phage anti-repressor protein	RAST
68	4	54	0.0%	nonsense	Phage anti-repressor protein	RAST
69	7	81	0.0%	gene indel	Conserved hypothetical protein	
70	5	42	0.0%	nonsense	Conserved hypothetical protein	
71	18	48	0.0%	gene indel	Conserved hypothetical protein	
72	1	-	-	gene indel	HNH homing endonuclease	RAST
73	32	62	100.0%	gene indel	Conserved hypothetical protein	
74	2	163	0.0%	gene indel	HNH homing endonuclease	RAST

Alignment length is the length of the protein alignment. Type marks core families of culture A phages and core families of all phages of the 936 group are denoted as "core (all)". One family (marked *) occurs two times in two strains, thus the total number of sequences in the family is given in brackets.

Table S2: Substitution rate estimation with BEAST.

Data set	Model	Clock model	Population model	Substitution rate estimate [95% HPD interval] in substitutions/site/year	Marginal likelihood	Bayes factor
Masked whole-genome alignment	1	strict	constant	1.888 [1.470-2.317] x 10 ⁻⁴	-22798.77	0
Masked whole-genome alignment	2	strict	exponential	1.890 [1.489-2.330] x 10 ⁻⁴	-22800.48	-1.715
Masked whole-genome alignment	3	strict	Bayesian Skyline	1.900 [1.464-2.317] x 10 ⁻⁴	-22798.89	-0.1258
Masked whole-genome alignment	4	strict	Extended Bayesian Skyline	1.889 [1.460-2.331] x 10 ⁻⁴	-22798.29	0.4780
Masked whole-genome alignment	5	Lognormal	constant	1.915 [1.486-2.371] x 10 ⁻⁴	-22803.45	-4.683
Masked whole-genome alignment	6	Random local clock	constant	1.889 [1.479-2.342] x 10 ⁻⁴	-22821.26	-22.49
Core codon alignment	1	strict	constant	1.737 [1.329-2.143] x 10 ⁻⁴	-	-
Core codon alignment, position 1	1	strict	constant	1.494 [0.9759-2.087] x 10 ⁻⁴	-	-
Core codon alignment, position 2	1	strict	constant	0.8751 [0.5001-1.264] x 10 ⁻⁴	-	-
Core codon alignment, position 3	1	strict	constant	2.887 [2.107 - 3.770] x 10 ⁻⁴	-	-
Core codon alignment, position 1	6	Random local clock	constant	1.525 [0.9690-2.105] x 10 ⁻⁴	-	-
Core codon alignment, position 2	6	Random local clock	constant	0.8892 [0.5160-1.301] x 10 ⁻⁴	-	-
Core codon alignment, position 3	6	Random local clock	constant	2.927 [2.148 - 3.823] x 10 ⁻⁴	-	-

Substitution rate estimate is mean rate for lognormal and random local clock models. Bayes factors compare H1 to H0 with H0 being the null model of a strict clock and a constant population size. Complex population size models have an absolute Bayes factor difference of less than two to the basic model. For the exponential growth model (model 2), the growth rate is estimated as -0.032 [-0.11,0.040]. This interval includes 0, thus the data is compatible with constant population size. Non-strict clock models are clearly rejected in comparison to the strict clock model (Bayes Factor > 4). For the random local clock model (model 6), the rate indicator estimate is 0.48 [0,2] and the highest posterior probability is at zero rate changes, thus the data is compatible with a strict molecular clock. The random local clock model based on a partition of the core codon alignment supports a strict molecular clock for each codon position, since the rate indicator estimates are 0.88 [0,2] for position 1, 0.78 [0,2] for position 2 and 0.50 [0,2] for position 3.

Table S3: Publicly available genomes for *L. lactis* phages of the 936 group.

Accession	Phage
AF009630	bIL170
AF011378	sk1
DQ054536	P008
DQ227763	712
DQ227764	jj50
EU221285	bIBB29
FJ848881	SL4
FJ848882	CB13
FJ848883	CB14
FJ848884	CB19
FJ848885	CB20
GQ979703	p2
JQ740787	ASCC191
JQ740788	ASCC273
JQ740789	ASCC281
JQ740790	ASCC284
JQ740791	ASCC287
JQ740792	ASCC310
JQ740793	ASCC324
JQ740794	ASCC337
JQ740795	ASCC356
JQ740796	ASCC358
JQ740797	ASCC365
JQ740798	ASCC368
JQ740799	ASCC395
JQ740800	ASCC397
JQ740801	ASCC406
JQ740802	ASCC454
JQ740803	ASCC460
JQ740804	ASCC465
JQ740805	ASCC473
JQ740806	ASCC476
JQ740807	ASCC489
JQ740808	ASCC497
JQ740809	ASCC502
JQ740810	ASCC506
JQ740811	ASCC527
JQ740812	ASCC531
JQ740813	ASCC532
JQ740814	ASCC544
KC182542	340
KC182543	645
KC182544	936
KC182545	fd13
KC182546	jm2
KC182547	jm3
KC182548	P113G

Accession	Phage
KC182549	p272
KC182550	P475
KC182551	P680
KC182552	phi7
KC522412	CaseusJM1
KM091442	phi15
KM091443	phi93
KM091444	phi145
KP793101	Phi4
KP793102	PhiA.16
KP793103	Phi19
KP793104	PhiB1127
KP793105	Phi19.3
KP793106	PhiA1127
KP793107	PhiD.18
KP793108	Phi5.12
KP793109	PhiC0139
KP793110	Phi43
KP793111	Phi19.2
KP793112	Phi129
KP793113	PhiF.17
KP793114	Phi17
KP793115	Phi114
KP793116	Phi13.16
KP793117	PhiG
KP793118	PhiF0139
KP793119	Phi10.5
KP793120	PhiL.18
KP793121	Phi109
KP793122	PhiL.6
KP793123	Phi4.2
KP793124	Phi44
KP793125	Phi91127
KP793126	PhiM.5
KP793127	Phi40
KP793128	PhiM.16
KP793129	PhiJF1
KP793130	Phi155
KP793131	PhiE1127
KP793132	PhiM1127
KP793133	PhiLj
KP793134	PhiS0139
KP793135	Phi.16

References

- Bouchard JD, Moineau S. 2004. Lactococcal phage genes involved in sensitivity to AbiK and their relation to single-strand annealing proteins. *J. Bacteriol.* 186:3649–3652.
- De Paepe M, Hutinet G, Son O, Amarir-Bouhram J, Schbath S, Petit M-A. 2014. Temperate phages acquire DNA from defective prophages by relaxed homologous recombination: the role of Rad52-like recombinases. *PLOS Genet.* 10:e1004181.
- Hecht A, Glasgow J, Jaschke PR, Bawazer LA, Munson MS, Cochran JR, Endy D, Salit M. 2017. Measurements of translation initiation from all 64 codons in *E. coli*. *Nucleic Acids Res.* 45:3615–3626.
- Huson DH, Bryant D. 2006. Application of phylogenetic networks in evolutionary studies. *Mol. Biol. Evol.* 23:254–67.
- Lopes A, Amarir-Bouhram J, Faure G, Petit M-A, Guerois R. 2010. Detection of novel recombinases in bacteriophage genomes unveils Rad52, Rad51 and Gp2.5 remote homologs. *Nucleic Acids Res.* 38:3952–62.
- Mavrich TN, Hatfull GF. 2017. Bacteriophage evolution differs by host, lifestyle and genome. *Nat. Microbiol.* 2:17112.
- Murphy J, Bottacini F, Mahony J, Kelleher P, Neve H, Zomer A, Nauta A, van Sinderen D. 2016. Comparative genomics and functional analysis of the 936 group of lactococcal *Siphoviridae* phages. *Sci. Rep.* 6:21345.
- Scaltriti E, Launay H, Genois M-M, Bron P, Rivetti C, Grolli S, Ploquin M, Campanacci V, Tegoni M, Cambillau C, et al. 2011. Lactococcal phage p2 ORF35-Sak3 is an ATPase involved in DNA recombination and AbiK mechanism. *Mol. Microbiol.* 80:102–116.