

Supplementary text 1:

Distribution of Spacer Hits: Contributions of Microbial Diversity and Diversity of CRISPR-Cas systems

The entire set of identified CRISPR arrays, all the matching protospaces and, separately, the subset of virus protospaces were broken down by the CRISPR-Cas system subtype and by the microbial phyla (Table X1-X3). For all 88 phylum/subtype categories with at least 10 matches, the fraction of spacers with matches and the fraction of virus matches among all matches were calculated. Both of these fractions were analyzed as follows.

The calculated fraction $f_{i,j}$ for the spacers in phylum i and belonging to the category j was modeled as the product of a phylum-specific factor p_i and a subtype-specific factor s_j . The deviation of the observed and the predicted fractions was assumed to be log-normally distributed with the expectation of 1. Formally,

$$\log f_{i,j} = \log p_i + \log s_j + e_{i,j}$$

where $e_{i,j} \sim \text{Norm}(0, \sigma)$.

Optimization of vectors \mathbf{p} and \mathbf{s} with respect to Σe^2 provides both the phylum- and subtype-specific contribution factors and the residuals $e_{i,j}$ that allow one to estimate the quality of the fit between the model and the observation. Setting one of the vectors to an arbitrary positive value (e.g. 1) and optimizing the other provides a model with phylum- or subtype-only contributions; the null model replaces both vectors with $\text{mean}(\log f_{i,j})$. The sums of squared residuals can be compared using the Fisher F-test.

Contribution of phylum- and subtype-specific factors to the fraction of spacers with matches

	phylum- and subtype	subtype-	subtype-	None
No. of data points	88	88	88	88
No. of parameters	39	21	18	0
No. of degrees of freedom	47	65	68	85
Σe^2	4.32	7.17	13.02	18.30

Contribution of phylum- and subtype-specific factors to the fraction of spacer matches to virus sequences

	phylum- and subtype	subtype-	subtype-	None
No. of data points	88	88	88	88

No. of parameters	39	21	18	0
No. of degrees of freedom	47	65	68	85
Σe^2	0.24	0.40	0.66	0.79

For both the overall fraction of matches and the fraction of viral matches, both phylum- and subtype-specific factors provide significant and independent contributions to the reduction of residuals (Fisher F-test p-value <0.05 except p=0.07 for the contribution of subtype-specific factors alone to the overall fraction of matches). Taken together, the two factors explain 70-75% of the original variance in the observed fraction.

Thaumarchaeota	0	21	252	0	0	0	0	0	0	0	0	22	0	0	0	41	0	0	0	0	0	0	225
Thermobaculum	0	0	0	0	0	18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	35
Thermococci	59	416	495	0	0	0	0	0	0	0	0	0	153	117	0	0	0	0	0	0	0	0	1970
Thermodesulfobacteria	0	0	187	0	0	0	0	0	0	0	0	0	29	51	0	0	0	0	0	0	0	0	63
Thermoplasmata	13	0	117	112	194	32	0	84	0	0	0	0	63	32	0	0	0	108	0	0	0	0	547
Thermotogae	0	6	1226	0	43	0	0	0	0	0	0	72	212	368	242	50	0	0	0	0	0	0	1512
unclassified	0	0	0	0	0	171	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	87
Acidobacteria																							
unclassified Archaea (miscellaneous)	0	0	0	0	60	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	8
unclassified Bacteria (miscellaneous)	3	0	631	123	408	17	0	0	24	0	237	3	33	196	0	289	0	0	0	0	0	0	850
unclassified Euryarchaeota	0	0	32	0	1	0	0	0	0	0	0	20	21	0	0	0	0	0	0	0	0	0	51
unclassified Proteobacteria	0	0	0	0	0	0	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
uncultured archaeon	0	0	0	55	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Verrucomicrobia	0	0	0	201	0	63	0	119	0	0	67	0	0	123	0	15	0	0	23	0	0	0	80
Zetaproteobacteria	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	5	0	0	0	0	0	0	25

Number of spacers with hits

Phylum	I	I-A	I-B	I-C	I-D	I-E	I-F	I-U	II-A	II-B	II-C	III	III-A	III-B	III-C	III-D	IV-A	V-A	V-B	VI-A	VI-B	VI-C	V-U	?
Acidithiobacillia	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0
Acidobacteria	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
Actinobacteria	56	6	120	429	0	1716	0	287	20	0	145	3	11	1	0	11	63	0	0	0	0	0	0	746
Alphaproteobacteria	0	1	0	45	0	34	0	0	0	0	10	0	0	1	0	2	0	0	0	0	0	0	0	27
Aquificae	0	0	5	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	4
Armatimonadetes	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
Bacteria candidate phyla	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
Bacteroidetes/Chlorobi group	1	0	230	98	0	3	0	4	1	0	319	4	8	16	1	8	0	5	0	0	50	0	0	92
Betaproteobacteria	2	9	0	157	1	72	116	0	0	0	221	2	2	5	2	0	0	0	0	0	0	0	0	319
Candidatus	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2
Marinimicrobia																								
Chloroflexi	0	0	1	7	0	16	0	0	0	0	0	0	2	0	0	2	0	0	0	0	0	0	0	2
Chrysiogenetes	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Cloacimonetes	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
Crenarchaeota	5	69	0	0	3	0	0	0	0	0	0	0	5	4	0	21	0	0	0	0	0	0	0	12
Cyanobacteria/Melainabacteria group	0	0	25	1	35	3	0	1	0	0	0	4	2	16	0	4	0	0	0	0	0	0	1	34
Deinococcus-Thermus delta/epsilon subdivisions	0	5	6	9	0	38	0	4	0	0	0	0	0	1	0	0	1	0	0	0	0	0	0	21
Elusimicrobia	7	0	52	28	0	16	7	2	0	2	178	0	5	1	0	8	2	0	0	0	0	0	0	70
environmental samples	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
Fibrobacteres	0	0	0	0	0	0	0	0	0	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0
Firmicutes	40	6	1380	978	11	338	7	0	2608	0	115	35	94	136	33	34	13	0	7	8	0	0	2	1798
Fusobacteriia	0	0	448	0	0	0	0	0	10	0	1	19	24	0	0	1	1	0	0	0	0	2	0	123
Gammaproteobacteria	286	0	0	529	0	1950	4736	12	0	7	46	0	13	45	0	15	65	0	0	0	0	0	0	2534
Halobacteria	0	0	47	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	26
Methanobacteria	0	0	37	0	0	0	0	0	0	0	0	1	6	0	0	0	0	0	0	0	0	0	0	3
Methanococci	0	1	5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Methanomicrobia	9	0	36	14	9	8	0	1	0	0	0	3	2	31	4	14	0	0	0	0	0	0	0	10
nitrifying bacterium enrichment culture	0	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
Nitrospira	0	0	0	0	0	11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	2
Planctomycetes	1	0	1	9	0	0	0	1	0	0	1	0	0	0	1	0	0	0	0	0	0	0	0	0
Spirochaetia	106	3	268	40	0	263	0	0	15	0	16	0	0	0	0	0	0	0	0	0	0	0	0	462
Synergistia	0	0	0	13	0	6	0	0	0	0	0	0	2	1	0	0	0	0	0	0	0	0	0	0
Tenericutes	0	0	0	0	0	0	0	0	25	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
Thaumarchaeota	0	0	1	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	1
Thermococci	0	2	4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	10
Thermodesulfobacteria	0	0	3	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
Thermoplasmata	0	0	0	0	5	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
Thermotogae	0	0	10	0	0	0	0	0	0	0	0	0	3	1	0	0	0	0	0	0	0	0	0	17
unclassified Bacteria (miscellaneous)	0	0	3	1	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0
unclassified	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Proteobacteria																								
Verrucomicrobia	0	0	0	18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	6

Supplementary text 2

Examples of intergenic hits into unannotated ORFs

Spacer 1:

>JODE01000009.1_62635_14_spacer_63457_32

GACCAGGCGTCCCCGCTGACGCTGATCGCGGT

Hit coordinates:

ContigID: LGDA01000250

Start: 45153

Stop: 45184

Intergenic sequence:

>gi|925387131|gb|LGDA01000250.1|:44655-45964 Streptomyces sp. WM6368 P399contig36.1, whole genome shotgun sequence

TAAGACATGCGGTGTGGCCGGGCACATGGATGCCCGGCATCCGACAGGGTCCTGCGAGCTGCCGCGTG
AGCACAGTGCAGCAGGACAGGCGGAGTCTGGCGGGGAACGGGAGTCGGCTGCGGCATACCCCGGAGGACG
ACTAAGTACAGCTGGATGGTTGCCGGCCACGCCCTCGGTGGACCATCCGACGGGGCTCCTCGTCTGC
GGCGGCGTCGGCCAGGAGCAGGTGCGGGTCTCCAGCCGGAACGGCGCGTACCGCCCCGGCTCCGGCT
CGTACCGGCGATACGGGGCGCTCTCGTAGACGACGGTGGTGTTCGGGCAGTCTCCGCGACCCCGAGCAG
CGCTTGGCGCTCGGCGTCTCGGCGGCAAGGCTCCAGCGCAACTTCGCCGCGGTCCACTGGGCGGCGTAC
TCGCAGTGGTAACCGCTGGCCGGCGGAATCCACTGCGCCGGATCCTTGTGCGCCTTCGACCGGTTGACT
TCGCGGTACCCGCGATCAGCGTCAGCGGGGACGCCTGGTTCGTCGTACGCCTCACGCCGGACCGCGTC
CCACGCGAAGCTGCCCGAGTCGTGGACTTCGGCCAAAGGTACGAAATGGTCTACATCAAGGCCCGCGGGC
TCCGTCAACAACGACGTCGCCGTAAGGGCTCAGCCGCGATCCACCCGACATCTTGACGCCCGCAGCGACGA
CCGGGGCCTCGACGGCCTCGGACAGGATGACTTCCCTGCGGGTATCACAGCCGTCCGGTTCGCGTTCAGCCC
GCGGTTCCAGTGCTTGTACAGGTCCCACTTGTAAACCTCCCGCTGCTCATCAGCGACGGGGATCCGGTTCG
ATGGCCTCGAACAGCGGCAGCGGGCGCACGGACTCCGGGAGCCCCAGCGATACCGGACGAGGCGGGGACGC
CCGCTGGGTGCCGGCCGCGTGCGCCGGGTCGGTGGCGAGCAGGGGCAGGGCAGCGAGAGCGAGCGC
GGGAGGACCACGCTGCAGCAGGTTCTTGATCACGCAGTGGGTTGTAGCGGTGCCCCAGCGGGACAGGGCC
TGGGGCCTCAGCTACCCGACAGGCCGCGGAATCACCGCAGAAACGATCAGCGGCTGGGCCGACCGGC
AGCCGGTCAACGCTCCTCAGCGCGAGATCATCTTTGGCCCGTGGGCAACGGCTCGGTCCCGGGCAACCC
ACACGCATCATCGGGAGTACCAGCAGATTCCAGCGCACGGCAGTTGCAACATCATGGCGTCAGGTGCCT
CGAGTCAGCTGCACATGACGGTCTTTGCGACGTCTCCAGACGAGGGCCG

Top 10 BLASTX hits: into HNH endonuclease

Fields: query id, subject id, subject ids, subject length, s. start, s. end, evalue, query seq, subject seq, q. start, q. end, score, subject sci names, subject title

10 hits found

gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |919577089|ref|WP_052876068.1| gi |919577089|ref|WP_052876068.1| 236 21 235 2.97e-126
 PAAHAAGTQAGVPASSGIAGAPGVRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKSNRKADK
 DPAQWIPPASGYHCEYAAQWTAAKLRWSLAADDAERQALLGVAEDCPNTTVVYESA
 PTAHAEPVPTAHSTAGAPSAAGLRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYDDVVVTDAAAGLDVDHDFVPLAEVHDSGGFAWDVRRREAYANDQASPLTLIAVTAKSNRKADK
 DPAQWLPPARGYHCEYAAQWTAATKLRWNLAAADDAERQALLGVAEDCPNTTVVYESA941 297 958 Streptomyces sp. NRRL F-4335HNH endonuclease [Streptomyces sp. NRRL F-4335]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |494479460|ref|WP_007268937.1| gi |494479460|ref|WP_007268937.1| 235 43 235 3.47e-110
 VRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQWTA
 KLRWSLAADDAERQALLGVAEDCPNTTVVYESA
 VRAPLPLFEAIDRITVAEEQREGYTRSLYKHWRNLNATDGCOTRREVILSEALEAPQVTAGCKMSSGLWLSPYDDVTMTDAAGLDVDHDFVPLAEVHDSGGYGDWAARREAYANDQASPLTLIAVTAKSNRKADKPAQWLPPAAGYHCQYAAATWVGT
 KLRWDLAADEAERQALLGLAEDCPNTTVVYEPAP 872 294 852 Streptomyces sp. C HNH endonuclease [Streptomyces sp. C]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |739837920|ref|WP_037688867.1| gi |739837920|ref|WP_037688867.1| 234 1 234 9.32e-110
 LQRGPPALALAALPLLATAPAAHAAGTQAGVPASSGIAGAPGVRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQ
 SPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQWTAAKLRWSLAADDAERQALLGVAEDCPNTTVVYESA MFRGLPALALCALPILA-
 APTAHSATAETAPLAATGVTGPAVGRAPLPLFEAIDRLPVAEHRDGYKRDLYKHWRNLNATDGCOTRREVILAEAVTAPVAAAGCKLTSGSWLSAYDNVVSDAARLDVDHDFVPLAEVYDSEAPWAAARREAYANDQASPLTLIAVSAASNRKADKPAEWLPS
 DDSYHCTYAAASWGTCLRWDLAVDENERQALLGLAEDCPNTTVVYESA 998 294 849 Streptomyces lavendulae HNH endonuclease [Streptomyces lavendulae]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |1154964266|ref|WP_078869587.1| gi |1154964266|ref|WP_078869587.1| 226 1 225 1.18e-105
 LQRGPPALALAALPLLATAPAAHAAGTQAGVPASSGIAGAPGVRAP--
 LPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQWTA
 QALLGVAEDCPNTTVVYESA MLRGLPALVLAAPVPLLATAPAAHA-----
 GTALTPSSAQPVTPLFDVAQDLVADEQREGYQSRSLYKHWRNLNATDGCOTRREVILAEAVLAPTVTAGCRLSGGSWHSAYDDLTVTDAARLDVDHDFVPLAEVHDSGGFAWDVRRREAYANDQASPLTLIAVSAASNRKADKPAEWMPSDGSYHCTYTATWVAT
 KLRWSLAADDERQALLGLAEDCPNTTVVYEPAP 998 297 821 Streptomyces sp. NRRL B-1347HNH endonuclease [Streptomyces sp. NRRL B-1347]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |860602850|ref|WP_048478564.1| gi |860602850|ref|WP_048478564.1| 235 1 235 5.86e-105
 LQRGPPALALAALPLLATAPAAHAAGTQAGVPASSGIAGAPGVRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQ
 SPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQWTAAKLRWSLAADDAERQALLGVAEDCPNTTVVYESA
 MLRGLPALALCTLPLFTAAPAAHSAPTVAAPAAAL TGGPAGLRAPMPLFEAIDRLPVGEHREGYKRDLYKHWRNLNAGDGCOTRREVILAEAVVAPQVAAAGCKLTGGSWRSAYDDVVVTDAAARLDVDHDFVPLAEVYDSEAPWAAARREAYANDQ
 SPDTLIAVSAASNRKADKPAEWLPSDGSYHCTYAAATWVGTCLRWDLAVDENERQALLGLAEDCPNTTVVYEPAP 998 294 818 Streptomyces roseus HNH endonuclease
 [Streptomyces roseus]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |664443183|ref|WP_030965733.1| gi |664443183|ref|WP_030965733.1| 233 1 233 2.10e-102
 LQRGPPALALAALPLLATAPAAHAAGTQAGVPASSGIAGAP-
 GVRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQWTA
 KLRWSLAADDAERQALLGVAEDCPNTTVVYESA MRHGLPALTLALPLLASTPTAHSAPAD---
 PVRGPAAGGPAGVRAPLPLFEAIDRLPVAEHRDGYKRDLYKHWRNLNAGDGCOTRREVILAEAVTAPVAAAGCKLTSGSWRSAYDNVVSDAARLDVDHDFVPLAEVYDSEQTPWTAARREAYANDQASPLTLIAVSAASNRKADKPAEWLPSDGSYHCTY
 AATWVGTCLRWDLAVDENERQALLGLAEDCPNTTVVHETAP 998 294 801 Streptomyces sp. NRRL S-378 HNH endonuclease [Streptomyces sp. NRRL S-378]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |926355608|ref|WP_053686546.1| gi |926355608|ref|WP_053686546.1| 232 46 232 1.85e-101
 LFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQWTA
 AADDAERQALLGVAEDCPNTTVVYESA
 LFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILAEAVIAPVVAAGCKLTGGSWRSAYDDVGLITDAARLDVDHDFVPLAEVHDSGGYAWDAARREAYANDQASPLTLIAVSAASNRKADKPAQWLPPYDSEAPWAAARREAYANDQ
 AADDAERQALLGLAEDCPNTTVVYIAP854 294 794 Streptomyces MULTISPECIES: HNH endonuclease [Streptomyces]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |664052845|ref|WP_030592179.1| gi |664052845|ref|WP_030592179.1| 241 45 241 6.65e-98
 GAPGVRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQ
 WTAAKLRWSLAADDAERQALLGVAEDCPNTTVVYESA
 GAGGLATLPLDAVERIPVAEQRREGYKRELYKHWRNLNAGDGCOTRREVILSEAVEAPVVAAGCKLTGGSWLSAYDDVGTVSDAGLDVDHDFVPLAEVHDSGGYAWDAARREAYANDQASPLTLIAVTAKSNRKADKPAQWLPPAADYRCTYAAE
 WTGKLRWLAADDAERQALLLAGECPNTTVVYGTAP 884 294 772 Streptomyces globisporus HNH endonuclease [Streptomyces globisporus]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |1154960293|ref|WP_078865614.1| gi |1154960293|ref|WP_078865614.1| 238 20 238 2.19e-97
 LPLLATAPAAHAAGTQAGVPASSGIAGAPGVRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKS
 NRKADKPAQWIPPASGYHCEYAAQWTAAKLRWSLAADDAERQALLGVAEDCPNTTVVYESA VPAADPTAAKETTTLQVLP-----
 VQGAANVRAPLALFEAIDRLPVAEHRDGYKRDLYKHWRNLNATDGCOTRREVILAEAVGAPQVAAAGCKLTGTTWRSAYDNLVTDAAARLDVDHDFVPLAEVYDSEQTPWAAARREAYANDQASPLTLIAVSAASNRKADKPAEWLPSDGSYHCTY
 AATWVGTCLRWDLAVDENERQALLGLAEDCPNTTVVYEPAP 962 294 768 Streptomyces roseus HNH endonuclease [Streptomyces roseus]
 gi |925387131|gb|LGDA01000250.1|:44655-45964 gi |494479287|ref|WP_007268764.1| gi |494479287|ref|WP_007268764.1| 233 41 233 1.87e-96
 VRAPLPLFEAIDRIPVADEQREGYKDWLYKHWRNLNATDGCOTRREVILSEAVEAPVVAAGCKMSSGSRSPYGDVVVTDAAAGLDVDHDFVPLAEVHDSGSFAWDVRRREAYANDQASPLTLIAVTAKSNRKADKPAQWIPPASGYHCEYAAQWTA
 KLRWSLAADDAERQALLGVAEDCPNTTVVYESA
 LRTVPVLYEADRILPVAPV
 KLRWNLAVDDAERQALLGLAEDCPNTTVVYETAP 872 294 761 Streptomyces sp. C HNH endonuclease [Streptomyces sp. C]

Spacer 2:

>JQ0001000126.1_6927_10067_28_spacer_9569_35
AGTCTATTTTGCCTGCTAACGGTTTCAGCATCTC

Hit coordinates:

ContigID: AOUX01000192
Start: 3802
Stop: 3836

Intergenic sequence:

>gi|456831403|gb|AOUX01000192.1|:3205-4402 *Leptospira interrogans* serovar Naam str. Naam ctg718000007680, whole genome shotgun sequence

```
TTAGCTCCGCATCTGTCCAACCTTTTTAGGGCGTTTTTTTTAATAGTTTTTTGGTCTGAATTACGTTGGT
CGCGCGATCAGGATAATATTTTCAACCAAATCAACCGTATACGGTGGTAAGACTTTTGACTTCTCGTCC
GTAATTTCAATTTCTTGCGAATCGATCCTTTAATGGTAGCAAAAGGCAATTCAACAGTTTTTTTATTTA
CAAATATTAATTTTTATTATTCAGAATATATTTACGTATTTGTTGTCCAAGCTTTCTTTCAAAAAC
TAATCTCGTTCAGTATTTGTAATTTAGAACGAACAGTTTCAACTCCGCGTTTTTTGTCCCCGTCGATG
CGTTTTAGATCGGATTCGATCTCGATTAATCTTTGAAGTCGTTATTAAGATCCTCTAATGTTTTAGGCG
GTCGGGTTGAACCTCATCGACAGTCTCCTTATTGATTTTTAATCTACTTGCCTCGTTGTCGACTAATA
GTAATTTACGTTTCCCTTTTGGGACCTTTGTTCCCTTTTTTTTCGATCGCCTCAATAACAAGGAGCGATT
AGGTCCTGGGATTAATCTGATTGTTTTATTGTAATCGGAGATGCTGAAACCGTTAGCAGTGCAAAATAGA
CTCATACGGCTATATAATTTGCCTATATGCGCCTTCAGTTCTTCTCGGTGGCCTCCGAAAGTAATCTTT
CTGACTTTTGGGTTTGTGTTGCTCATTGTACATCTATAAATTATAATTACTACTATGGTTAGGCATAAT
AAACGTATTATAAAAAGCTCAAGCATCTAGATCCCAGAAAATATAAAACACAGTTACCCAAATAGCAGCA
GCTAAAAAGATCAAGAATGCTTGTAAAAAGAGATACGCTGTCATAACTTTACAACCCCAAGTGGATT
AACGGAACCCCTGTGATGATGAAGGGCTATAGCGATCAATGAGGCGACGCACTTAAGGCTTGTTTTTTT
GCGTCGGACCGATCCATGTGGAGCTCGTCACCATCACCATCTTTCTGAAAAAATAGATCTGAATGCTCTA
TAGTTTGTAAAGCACGTCTGAAAAGATCATATATCGGCCTCCACCAAACCTGTTTATTGTTCCAAATC
TTAACTGCGTCCGCAATTGATTGCGCAATATTGTGGTTACGCGACAAGAGTGACATCTAACGGCATAGG
CCTTGTTTC
```

Top 10 BLASTX hits: into hypothetical protein

```
# Fields: query id, subject id, subject ids, subject length, s. start, s. end, evalue, query seq, subject seq, q. start, q. end, score, subject sci names, subject title
# 10 hits found
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|446021850|ref|WP_000099705.1| gi|446021850|ref|WP_000099705.1| 176 1 145 8.09e-94
MSSTRPPKTLLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
MSSTRPPKTLLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL 436
2 735 Leptospira interrogans hypothetical protein [Leptospira interrogans]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|488052634|ref|WP_002124031.1| gi|488052634|ref|WP_002124031.1| 176 1 145 3.59e-93
MSSTRPPKTLLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
```

MSSTRPPKLTLEDLNKQLRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL 436
2 730 Leptospira interrogans Gam-like protein [Leptospira interrogans]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|516461609|ref|WP_017850447.1| gi|516461609|ref|WP_017850447.1| 176 1 145 8.23e-87
MSSTRPPKLTLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
MSLTQPPKLTLEDLNKRLQRLIKIDSDLKRIDGDKNAEVEVTRSKFTDTERDLVFERESLDKQVREYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKTALKSWTGAEL 436
2 688 Leptospira interrogans hypothetical protein [Leptospira interrogans]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|490915838|ref|WP_004777741.1| gi|490915838|ref|WP_004777741.1| 176 1 145 1.75e-79
MSSTRPPKLTLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
MSSTQPPKLTLEDLNKRLQRLIEIESDLKRIDGDKNTEVESVRSRFDTERELVFEKENLDKQVREYILNNKDIIFVHKKTVELPFATIKRIDSQEIEITDEKSKTLPPTYVDLVEKYYSRPNVVIQTKKTIKKAALKNWTNAEL 436
2 640 Leptospira kirschneri Gam-like protein [Leptospira kirschneri]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|487897388|ref|WP_001970854.1| gi|487897388|ref|WP_001970854.1| 178 1 147 1.00e-69
LSMSSTRPPKLTLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
MSTTQNPQPKLTLEDLNKMQRLVEIESDLKRIEAGEKNSEVESVRSRFDTERDLVFEKEELDKQVRDFVQNKDITLFAHRKTVELPFATIKRIDSQEIEITDEKSKELPPYSVDLIEKFYPERANNAIQIKKTVKKTALKSWTDAEL 442
2 575 Leptospira interrogans Gam-like protein [Leptospira interrogans]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|1001638526|ref|WP_061216057.1| gi|1001638526|ref|WP_061216057.1| 176 1 145 4.86e-69
MSSTRPPKLTLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
MNAATQPKLTLEDLNARLQRLVEIEANLKRIGGDKNTEVESVRSRFDTERDLVFEKEQLDRQVREFVLQNKDITLFAHRKTIELPFATIKRIDSQEIEITDEKSKTLPYSVDLIEKFYPERASDAIQIKKSIKKAALKSWTDAEL 436
2 571 Leptospira santarosai Gam-like protein [Leptospira santarosai]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|490641537|ref|WP_004506532.1| gi|490641537|ref|WP_004506532.1| 178 1 147 1.55e-67
LSMSSTRPPKLTLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
MNRAPNPQPKLTLEDLNKMQRLIEIKSDLKRIEAGEKNTEVESIRSRFVGRDLVFEKEELDKQVREFVQNKDITLFAHRKTIELPFATIKRIDSQEIEITDEKSKTLPYSVDLIEKFYPERANNAIQIKKSVKKTALKSWTDAEL 442
2 561 Leptospira weilii Gam-like protein [Leptospira weilii]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|515129273|ref|WP_016758123.1| gi|515129273|ref|WP_016758123.1| 177 8 146 1.70e-65
PKTLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
PKSLEDLNIRMQRLEIESDLKRIEAGEKNTEVEDTRSHFVEVERDLVFEKERLDKQIRAFVMMENKNTLFAHRKTIELPFATIKRIDSQEIEITDEKSKTLPYSVDLIEKFYPERANNAIQIKKSVKKTALKSWTDAEL 418 2
547 Leptospira MULTISPECIES: hypothetical protein [Leptospira]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|516465861|ref|WP_017854324.1| gi|516465861|ref|WP_017854324.1| 103 1 103 6.90e-64
MSKQTQKSERLLSEATEKKLKAHIGKLYSRMSLFCTANGFSISDYNKTIIRLIPGPKSLLVIEAIEKKGTKVPGKRRKLLLDNERASRLKINKETVDEFNPTA 725 417 529 Leptospira
interrogans hypothetical protein [Leptospira interrogans]
gi|456831403|gb|AOUX01000192.1|:3205-4402 gi|516474782|ref|WP_017863226.1| gi|516474782|ref|WP_017863226.1| 176 1 145 1.48e-61
MSSTRPPKLTLEDLNKRLQRLIEIESDLKRIDGDKNAEVEVTRSKFTNTERELVFERESLDKQIRKYILNNKNLIFVNKKTVELPFATIKRIDSQEIEITDEKSKVLPPYTVDLVEKYYPDRATNVIQTKKTIKKNALKSWTDAEL
MNAATQPKLTLEDLNARLQRLVEIEANLKRIDGDKNTEVESVRSRFDTERDLVFEKEQLDRQVREFVLQNKDITLFAHRKTIELPFATIKRIDSQEIEITDEKSKTLPYSVDLIEKFYPERVSDAIQIKKLIKTTALKSWTDAEL 436
2 521 Leptospira santarosai hypothetical protein [Leptospira santarosai]

Spacer 3:

>HG916826.1_842637_846204_26_spacer_844192_32
TACATCAGCAAGGACGGCCCGGACATGATCCC

Hit coordinates:

ContigID: LMMB01000004

Start: 84179

Stop: 84210

Intergenic sequence:

>gi|944790866|gb|LMMB01000004.1|:84017-84790 Pseudomonas sp. Leaf83 contig_12, whole genome shotgun sequence
TGTACAAGCCAGTATTTCCGGGGTTC TAGCCGGATCGCTTCTCAGATGGGAGCGTTTTGGGAACACTTT
GGGAATGGTGGAACGAAAAAGGCAGCCTGCGAGGGCTGCCTTTTGTCTGTTTCTAGGGCCTGGATTACAG
CTTTAGCGCGTCTCCAGGAGCGGGATCATGTCCGGGCCGTCTTGCTGATGTA CTGACTTGCCGTAATGACGT
TCGAGCATTGCGATTGTCTGTGCTTCCATACGGCTGACGCAGGCTGGGCAGAGGCAC TTGTAACCGCCG
CTGTTGGTGGTCATGGCTGGGGCTTCTGTGCAGGGCAGGATGGCGCCCACTTCGGCGTGGCGCGCGGC
GATGTGCTGGCGGTGGGCGGTGATGATGCTGGCCAGCTCGGCTTCGTCGATTTGCCCGTCTTCCAGGGCC
TTGAGGATCATCGCTCGACCTTGCCCGCTTTACCGCCGTGCCATGGCGCGGGCCAGCAGGTCAATGT
TGTCCTGTCTGT CAGCTCGGGCAGGGGACGTAACCGCCGTTGTAGAGGCTGGTGAGGTAGTCGGGCAGA
TAGCTGGTGC CGGCACCTTCTCCAGCTGGCGCACTTGTTCGTCCGTCAGCGGGCGGTTGCCGGGTTCT
CGTACAGCTTGTGTGCAACTGCTTGAGGTCCAGCCCCGGGCAGGTGGCGGCGCATTCGCGGCCCCCTGG
GAAGGCGCCCCGGTGAGGGGGAATGCGAGGGCGATATCGGCAAGGCCGTGCTGGTGGGCGATCCGTTGTG
AATC

Top 10 BLASTX hits: into hypothetical protein

Fields: query id, subject id, subject ids, subject length, s. start, s. end, evalue, query seq, subject seq, q. start, q. end, score, subject sci names, subject title
20 hits found

Query ID	Subject ID	Subject IDs	Subject Length	s. start	s. end	evalue	Query Seq	Subject Seq	q. start	q. end	score	Subject Sci Names	Subject Title
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 980973412 ref WP_059392489.1	gi 980973412 ref WP_059392489.1	156	83	156	1.47e-63	LPELTDNDIDLLARAMGTAVKRGKVDAMILKALEDDGQIDEAELASII TAHRQHIAARHAEVGAIALLHRKPQP	MPEPADTDNDIDLLARAMGTAVKRGKVDAMILKALEDDGQIDEAELASII TAHRQHIAARHAEVGAIALLHRKPQP	515	294	355	Pseudomonas toytomiensis	hypothetical protein [Pseudomonas toytomiensis]
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 980973412 ref WP_059392489.1	gi 980973412 ref WP_059392489.1	156	18	85	1.47e-63	GAFPGGRECAATCPGLDLKQFDNKL YENPGNRPLTDEQVRQLEKVAGTSYLPDYLTSLYNGVYVPCPS	AAFPGGRECAATCLGLDLKQFDNKL YENPGHRPLTDEQVRQLEKVAGTSYLPDYITGLYNGVVFAMPE	709	506	327	Pseudomonas toytomiensis	hypothetical protein [Pseudomonas toytomiensis]
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 545127478 ref WP_021488921.1	gi 545127478 ref WP_021488921.1	155	82	155	2.10e-63	LPELTDNDIDLLARAMGTAVKRGKVDAMILKALEDDGQIDEAELASII TAHRQHIAARHAEVGAIALLHRKPQP	MPELTDNDIDLLARAMGTAVKRGKVDAMILRALEDGQIDEAELASII TAHRQHIAARHAEVGAIALLHRKPQP	515	294	367	Pseudomonas mendocina	hypothetical protein [Pseudomonas mendocina]
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 545127478 ref WP_021488921.1	gi 545127478 ref WP_021488921.1	155	18	84	2.10e-63	GAFPGGRECAATCPGLDLKQFDNKL YENPGNRPLTDEQVRQLEKVAGTSYLPDYLTSLYNGVYVPCPS	AAFPGGRECAA-CLGLDLKQFDNKL YENPGHRPLTDEQVRQLEKVAGTYYLPDYLTSLYNGVVFAMPE	709	506	313	Pseudomonas mendocina	hypothetical protein [Pseudomonas mendocina]
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 917193888 ref WP_051800600.1	gi 917193888 ref WP_051800600.1	156	83	156	2.14e-62	LPELTDNDIDLLARAMGTAVKRGKVDAMILKALEDDGQIDEAELASII TAHRQHIAARHAEVGAIALLHRKPQP	MPEPADTDNDIDLLARAMGTAVKRGKVDAMILKALEDDGQIDEAELASII TAHRQHIAARHAEVGAIALLHRKPQP	515	294	355	Pseudomonas oleovorans	hypothetical protein [Pseudomonas oleovorans]
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 917193888 ref WP_051800600.1	gi 917193888 ref WP_051800600.1	156	18	85	2.14e-62	GAFPGGRECAATCPGLDLKQFDNKL YENPGNRPLTDEQVRQLEKVAGTSYLPDYLTSLYNGVYVPCPS	AAFPGGRECAAACLGLDLKQFDNKL YENPGHRPLTDEQVRQLEKVAGTSYLPDYITGLYNGVVFAMPE	709	506	317	Pseudomonas oleovorans	hypothetical protein [Pseudomonas oleovorans]

gi 944790866 gb LMMB01000004.1 :84017-84790	gi 835622920 ref WP_047589232.1	gi 835622920 ref WP_047589232.1	156	83	156	8.82e-60
LPELTDNIDLLARAMGTAVKRGKVDAMILKALEDGQIDEAELASIIIAHRQHIAARHAEVGAIALHRKPQP	MPELADTDNIDLLARAMGTTIKRGTVDAEILKALEDGEISEAELASIIAAHRQHIAARHAEVGAIALHRKPQP					515
294 331 Pseudomonas mendocina	hypothetical protein [Pseudomonas mendocina]					
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 835622920 ref WP_047589232.1	gi 835622920 ref WP_047589232.1	156	18	84	8.82e-60
GAFPGGRECAATCPGLDLKQFDNKLLENPGRPLTDEQVRQLEKQVAGTSYLPDYLTSLYNGVYVPCP	AAFPGGRECAATCLGLDLKQFDNKLLENPGRPLTDEQVLQLEKQVAGTSYLPDYISGLYNGVYVAMP					709
509 318 Pseudomonas mendocina	hypothetical protein [Pseudomonas mendocina]					
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 503480628 ref WP_013715289.1	gi 503480628 ref WP_013715289.1	159	18	84	3.53e-53
GAFPGGRECAATCPGLDLKQFDNKLLENPGRPLTDEQVRQLEKQVAGTSYLPDYLTSLYNGVYVPCP	AAFPGGRECAATWGLDLKQFDNKLLENPGRPLTDEQVLQLEKQVAGTSYLPDYISGLYNGVYVAMP					709
509 313 Pseudomonas mendocina	hypothetical protein [Pseudomonas mendocina]					
gi 944790866 gb LMMB01000004.1 :84017-84790	gi 503480628 ref WP_013715289.1	gi 503480628 ref WP_013715289.1	159	83	159	3.53e-53
LPELTDNIDLLARAMGTAVKRGKVDAMILKALEDGQIDEAELASIIIAHRQHIAARHAEVGAIALHRKPQP*PP	MPELAELDNIDLLERAMTTTIKRGTVDAMILTALKDGEINEAELASIIAHRQHMAARHAEVSSILALHSKRQEPKP		515	285	279	Pseudomonas mendocina
protein [Pseudomonas mendocina]						hypothetical

Spacer 4:

>LFQC01000004.1_87106_87459_1_spacer_87136_35
GTAACAGATGATATTAATTCAAAGAATTAAGCTT

Hit coordinates:

ContigID: AOSX01000029

Start: 620311

Stop: 620345

Intergenic sequence:

>gi|727535586|gb|AOSX01000029.2|:619946-620968 Clostridium botulinum Af84 Contig_29, whole genome shotgun sequence
AAATTCGGATAGAAACATAGATATTTCTCATATAGATTTTTTTGATTAATTATATGCACTGATGCTATGA
ATAAAAACAAAAAATTCTAAAAGTTTTCCAACCGCATTAAATTTTAAAACCTATCCTAAGCATTGCAA
TATAAGGCTTAAGATAGGTTTTGTTTATGTTTTCTGAAAAATCACAAATGGCTGGAAAAATTTTTTAA
TCATTGTAAAATCAAATGCATATGCTATTCCTTAAAATTAAGGAATGGCTAATTTACTATGGATGAACAT
TAACATAGGATGATTTAAATTTGGAAAAGCTTACTAGATTTCTTACCTATATCGGGTTGAACAATAACAT
GAGATGCATTTAAATAAGTTAATTCCTTTGAATTAATATCATCTGTTACACTGAACAATAACATAAGAT
GTATTATTAATAAATGGTAGTTGTAGCTTATGTTATAGCTATTATTTTTCAATATTTAGAGGAATTTTTT
ACATATGTAGAATATTTGTATATAGTCCCTTAATTCATTATCAAATCTTCTTGAATAAAAAGGAGCCT
GTGTATGCAGTGGGGCTCCTTTTTTATATTATTTGAAGGAATTTTTAACATTTATAGAATATTAATA
TAACGGCTTCCCAATAGGTTGATTATAGACCTCTTGCTTAAGCAAAAAGAACCCCAATAAAAAGGGGTT
CTTTTTGTATGGATTTATATTGGTGCTATGTATTGGTCTATTTTTACTATATCCAGATGTAGAAAAGT
TAATCAGTAAGAAATATATAAGATGTAATTTTTAAATGTATTTCTATAATGAGAGTACCTTATTTTTTA
TTTTTTTGATAAAAATTTAAGAAAATTAGCAATTTATTAATAAAAATTTAAAATAATTATTTAATTTCA
TTTTAAAGTTGTATCAGAATGTATAAAAAGTAAAAGATGTAGGATAAACAAGAAGGAATTGAATTA
TTGTAAGAATATAAATATTATGGTAAAAATATAAATAAATA

Top 10 BLASTX hits: (1 hit found) into hypothetical protein

Fields: query id, subject id, subject ids, subject length, s. start, s. end, evalue, query seq, subject seq, q. start, q. end, score, subject sci names, subject title

1 hits found

gi 727535586 gb AOSX01000029.2 :619946-620968	gi 1119667520 ref WP_072587154.1	gi 1119667520 ref WP_072587154.1	80	49	
79	2.91e-07VTLNNNIRCIKNGSCSLCYSYYFSIFRGIF LKLNNNIRCIKNGSCSLCYSYYFSLFRRI	396	488	136	Clostridium botulinum
	hypothetical protein [Clostridium botulinum]				

gi 393068133 gb AKOR01000006.1 :34315-34513	gi 446797338 ref WP_000874594.1			
gi 446797338 ref WP_000874594.1 2399	258	323	3.26e-31	
QTSFNQGTYNFSNSATLSFNNSNFNQGTYHFNSAQSTFENSNFNQGTYNFNDNTSFNNDTFNQGAY				
QTSFNQGTYNFSNSATLSFNNSNFNQGTYHFNSAQSTFENSNFNQGTYNFNDNTSFNNDTFNQGTY	1	198	307	
Helicobacter pylori	toxin-like outer membrane protein [Helicobacter pylori]			
gi 393068133 gb AKOR01000006.1 :34315-34513	gi 446797381 ref WP_000874637.1			
gi 446797381 ref WP_000874637.1 2529	258	323	1.63e-30	
QTSFNQGTYNFSNSATLSFNNSNFNQGTYHFNSAQSTFENSNFNQGTYNFNDNTSFNNDTFNQGAY				
QTSFNQGTYNFSNSATLSFNNSNFNQGTYHFNSTQSTFENSNFNQGTYNFNDNTSFNNDTFNQGTY	1	198	302	
Helicobacter pylori	toxin outer membrane protein [Helicobacter pylori]			
gi 393068133 gb AKOR01000006.1 :34315-34513	gi 446797342 ref WP_000874598.1			
gi 446797342 ref WP_000874598.1 2399	258	323	2.63e-30	
QTSFNQGTYNFSNSATLSFNNSNFNQGTYHFNSAQSTFENSNFNQGTYNFNDNTSFNNDTFNQGAY				
QTSFNQGTYDFNSNSATLSFNNSNFNQGTYHFNSAQSTFENSNFNQGAYNFNDNTSFNNDTFNQGTY	1	198	301	
Helicobacter pylori	toxin outer membrane protein [Helicobacter pylori]			
gi 393068133 gb AKOR01000006.1 :34315-34513	gi 1159760846 ref WP_079307344.1			
gi 1159760846 ref WP_079307344.1	2397	259	324	3.32e-30
QTSFNQGTYNFSNSATLSFNNSNFNQGTYHFNSAQSTFENSNFNQGTYNFNDNTSFNNDTFNQGAY				
QTSFNQGTYNFSNSATLSFGNSNFNQGTYHFNSAQSTFENSNFNQGTYNFNDNVSFNNDTFNQGTY	1	198	300	
Helicobacter pylori	toxin [Helicobacter pylori]			

Predicted CRISPR-Array: 2

ContigID: KK099646.1 (Staphylococcus_aureus_DAR3163_GCA_000610585.1)
Start: 60277
Stop: 60684
<https://www.ncbi.nlm.nih.gov/projects/sviewer/?id=KK099646.1&v=60277..60684>
Predicted in: serine-aspartate repeat-containing protein

CRISPR-Finder prediction:

CRISPR id : tmp_1_PossibleCrispr_1

- CRISPR start position : 1 ----- CRISPR end position : 132 ----- CRISPR length : 131
- DR consensus : AGCGATTCAGACTCAGATAGCGACTCAGA
- DR length : 29 Number of spacers : 2

```
1 AGTGACTCAGATTCGGACAGCGATTCAGA CTCAGATAGCGACTCAGATTCAGAT 54
55 AGTGACTCAGACTCAGATAGCGACTCAGA CTCAGATAGCGACTCAGAC 102
103 TCAGATAGCGACTCAGACAGCGACTCAGA 132
```

CRISPR id : tmp_1_PossibleCrispr_2

- CRISPR start position : 278 ----- CRISPR end position : 408 ----- CRISPR length : 130
- DR consensus : AGCGATTCAGACTCAGATAGCGACTCAGA
- DR length : 29 Number of spacers : 2

```
278 AGCGATTCAGATTCAGACAGCGACTCAGA TTCAGATAGCGACTCAGACTCAGAC 331
332 AGCGATTCAGACTCAGATAGCGACTCAGA CAGCGATTCAGATTCGGAT 379
380 AGCGATTCAGATTCAGATGCAGGTAACA 408
```

Top 10 BLASTX results in RefSeq:

Fields: query id, subject id, subject ids, subject length, s. start, s. end, evalue, query seq, subject seq, q. start, q. end, score, subject sci names, subject title
587 hits found

```
gi|601598167|gb|KK099646.1|:60277-60684 gi|1093956063|ref|WP_070956199.1|
  gi|1093956063|ref|WP_070956199.1| 1426 1225 1281 7.54e-09
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  aureus hydrolase [Staphylococcus aureus] 171 150 Staphylococcus
gi|601598167|gb|KK099646.1|:60277-60684 gi|1093956063|ref|WP_070956199.1|
  gi|1093956063|ref|WP_070956199.1| 1426 1227 1283 7.54e-09
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  aureus hydrolase [Staphylococcus aureus] 171 150 Staphylococcus
gi|601598167|gb|KK099646.1|:60277-60684 gi|1093956063|ref|WP_070956199.1|
  gi|1093956063|ref|WP_070956199.1| 1426 1229 1285 7.54e-09
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  aureus hydrolase [Staphylococcus aureus] 171 150 Staphylococcus
gi|601598167|gb|KK099646.1|:60277-60684 gi|1093956063|ref|WP_070956199.1|
  gi|1093956063|ref|WP_070956199.1| 1426 1231 1287 7.54e-09
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  SDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDSDS
  aureus hydrolase [Staphylococcus aureus] 171 150 Staphylococcus
```



```

gi|451810799|gb|ANAE01000140.1|:57689-58634      gi|485835026|ref|WP_001446733.1|
gi|485835026|ref|WP_001446733.1|;gi|218706248|ref|YP_002413767.1| 172      2      167      1.36e-11
RCSAGARAGP--SPRARGSRHPHGHLHLHRSIPACAGLTHILALVGTSPVHPRVRGAHAHLGVGGDFATGSPRARGSP-----
TSWRWWGLRHRSSIPACAGLT--HILALVGTSPVHPRVRGAHPHLGVGGDFATGSPRARGSPTSWRWWGLRHRSSIPACAGLTHILA
KWSSTATVGPGLSPLARGTRWWRGQTEQSRRFIPAGAGNTLTYNTENHTLTVYPRWRGEHGHSPKQTQEVNGLSPLARGTPGFHLSKCMW-----
RFIPAGAGNTTPDIFATVVT--PVYPRWRGEHRNQTQRANGVAGLSPLARGTLALFLSGESRHRFIPAGAGNTYFLS      223      714      174
Escherichia coli;Escherichia coli UMN026 hypothetical protein [Escherichia coli]
gi|451810799|gb|ANAE01000140.1|:57689-58634      gi|485835026|ref|WP_001446733.1|
gi|485835026|ref|WP_001446733.1|;gi|218706248|ref|YP_002413767.1| 172      12      167      2.32e-10
GSPRARGSRAGHRWRVAGRSIPACAGLTRALALVGTSPVHPRVRGAHPHLGVGGDFATGSPRARGSRHRTVLGSCASRSIPACAGLTSSTWPPP
TPTPVHPRVRGAHPHLGVGGDFATGSPRARGSRHRTVLGSCASRSIPACAGLTHILA
GLSPLARGTRWWRGQTEQSRRFIPAGAGNTLTYNTENHTLTVYPRWRGEHGHSPKQTQEVNGLSPLARGTPGFHLSKCMWRFIPAGAGNTTPDIFAT
VVTPVYPRWRGEHRNQTQRANGVAGLSPLARGTLALFLSGESRHRFIPAGAGNTYFLS      3      470      165      Escherichia
coli;Escherichia coli UMN026 hypothetical protein [Escherichia coli]
gi|451810799|gb|ANAE01000140.1|:57689-58634      gi|485835026|ref|WP_001446733.1|
gi|485835026|ref|WP_001446733.1|;gi|218706248|ref|YP_002413767.1| 172      12      167      1.15e-09
GSPRARGSPTSWRWW---GLRHRSSIPACAGLTHILALVGTSPVHPRVRGAHPHLGVGGDFATGSPRARGSP-----
TSWRWWGLRHRSSIPACAGLT--HILALVGTSPVHPRVRGAHPHLGVGGDFATGSPRARGSPTSWRWWGLRHRSSIPACAGLTHILA      GLSPLARGT----
RWWRGQTEQSRRFIPAGAGNTLTYNTENHTLTVYPRWRGEHGHSPKQTQEVNGLSPLARGTPGFHLSKCMW-----RFIPAGAGNTTPDIFATVVT--
PVYPRWRGEHRNQTQRANGVAGLSPLARGTLALFLSGESRHRFIPAGAGNTYFLS      430      897      160      Escherichia
coli;Escherichia coli UMN026 hypothetical protein [Escherichia coli]
gi|451810799|gb|ANAE01000140.1|:57689-58634      gi|485835026|ref|WP_001446733.1|
gi|485835026|ref|WP_001446733.1|;gi|218706248|ref|YP_002413767.1| 172      12      167      1.82e-08
GSPRARGSLVPRWW---GLRHRSSIPACAGLTHILALVGTSPVHPRVRGAHVIGRCSAGARAGSPRARGSRHPHGHLHLHRSIPACAGLT--
HILALVGTSPVHPRVRGAHAHLGVGGDFATGSPRARGSPTSWRWWGLRHRSSIPACAGLTHILA      GLSPLARGT----
RWWRGQTEQSRRFIPAGAGNTLTYNTENHTLTVYPRWRGEHGHSPKQTQEVNGLSPLARGTPGFHLSKCMWRFIPAGAGNTTPDIFATVVT--
PVYPRWRGEHRNQTQRANGVAGLSPLARGTLALFLSGESRHRFIPAGAGNTYFLS      64      531      151      Escherichia
coli;Escherichia coli UMN026 hypothetical protein [Escherichia coli]
gi|451810799|gb|ANAE01000140.1|:57689-58634      gi|485835026|ref|WP_001446733.1|
gi|485835026|ref|WP_001446733.1|;gi|218706248|ref|YP_002413767.1| 172      12      163      3.34e-08
GSPRARGSPTSWRWW---GLRHRSSIPACAGLTHILALVGTSPVHPRVRGAHPHLGVGGDFATGSPRARGSP-----
TSWRWWGLRHRSSIPACAGLT--HILALVGTSPVHPRVRGAHPHLGVGGDFATGSPRARGSPTSWRWWGLRHRSSIPACAGLT      GLSPLARGT----
RWWRGQTEQSRRFIPAGAGNTLTYNTENHTLTVYPRWRGEHGHSPKQTQEVNGLSPLARGTPGFHLSKCMW-----RFIPAGAGNTTPDIFATVVT--
PVYPRWRGEHRNQTQRANGVAGLSPLARGTLALFLSGESRHRFIPAGAGNT      491      946      149      Escherichia
coli;Escherichia coli UMN026 hypothetical protein [Escherichia coli]
gi|451810799|gb|ANAE01000140.1|:57689-58634      gi|485668974|ref|WP_001310069.1|
gi|485668974|ref|WP_001310069.1|;gi|218706249|ref|YP_002413768.1| 126      12      122      6.39e-06
GSPRARGSRAGHRWRVAGRSIPACAGLTRALALVGTSPVHPRVRGAHPHLGVGGDFATGSPRARGSRHRTVLGSCASRSIPACAGLTSSTWPPP
TPTPVHPRVRGAH
GLSPLARGTRIFYRRMSRGNRFIPAGAGNTSAAPANVSAVTVYPRWRGEHPRPSIISCASFGLSPLARGTHQEADESARHARFIPAGAGNTQFHMSRR
AGSSVYPRWRGEH      3      335      130      Escherichia coli;Escherichia coli UMN026 hypothetical protein
[Escherichia coli]
# BLAST processed 1 queries

```

Predicted CRISPR-Array: 4

ContigID: KQ087587.1 (Enterobacter_aerogenes_UCI97_GCA_001030185.1)

Start: 84822

Stop: 86069

<https://www.ncbi.nlm.nih.gov/projects/sviewer/?id=KQ087587.1&v=84822..86069>

PilerCR predicted Array

>gi|844486473|gb|KQ087587.1|:84822-86069 Enterobacter aerogenes strain UCI97 genomic scaffold aeuqd-supercont1.1, whole genome shotgun sequence

GCGCCTTATCCGGGCTACGGTTCGGTATTTGGTTGGTAGCCCCGGTAAGCGTAAGCGCCACCGGGGAGG
ATTCCCGGATAGCGGCGGTAGCGCCTTATCCGGGCTACGGTTCGGTATGCGGTTGGTAGCCCCGGTAAG
CGTAAGCGCCACCGGGGAGGGTTCGGGATGGCGGCGCATAGCGCCTTATCCGGGCTACGGTTCGGTAG
TCGGTTGGTAGCCCCGGTAAGCGTAAGCGCCACCGGGGAGGGTTCGGGATGGCGGCGCATGGCGCCTTA
TCCGGGCTACGGTTCGGTATTCGGTCGGTAGCCCCGGTAAGCGTAAGCGCCACCGGGGAGGATTCGGGA
TAGCGGCGCATCGCGCCTTATCCGGGCTACGATTTCGGTATTCGGTTGGTAGCTCCGGTAAGCGTAAGCGC
CACCGGGGAGGGTTCGGGATGGCGGCGCNN
NNNTAAACTGCGC
GCGTGGGTCGCCAGCTTCGGCCTGCTGTTGGTATTGGCTGCCCGGATCCGCAGCGTCTGCGGCGAAGCTT
GTCAGCGCCAGCAGTAAATAGAGGTATTTCAATTTTTTATTATCAATAAGGGAGACCCCGCAGTATAAA
GGCGGGCGGAAAAGAAGAAAATGGCGGGAATTAGCGGGTGAATACCCGGATAGCGGCGCATAGCGCCTT
ATCCGGGCTACGGTTCGGTATTTGGTTGGTAGCCCCGGTAAGCGTAAGCGCCACCGGGGAGGATTCGGG
GATAGCGGCGCGTACGCGCCTTATCCGGGCTACGGTTCGGTATTCGGTTGGTAGCCCCGGTAAGCGTAAGC
GCCACCGGGGAGGGTTCGGGATGGCGGCGCATAGCGCCTTATCCGGGCTACGGTTCGGTAGTCGGTTG
GTAGCCCCGGTAAGCGTAAGCGCCACCGGGGAGGGTTCGGGATGGCGGCGCATGCGCCTTATCCGGGCT
TACGGTTCGGTATTCGGTCGGTAGCCCCGGTAAGCGTAAGCGCCACCGGGGAGGATTCGGGATAGCGGC
GCATCGCGCCTTATCCGGGCTACGATTTCGGTATTCGGTTGGTAGCTCCGGTAAGCGTAAGCGCCACCGG
GAGGGTTCGGGATGGCGGCGCATAGCGCCTTATCCGGGCTACGGTTCGGTAGTCGGT

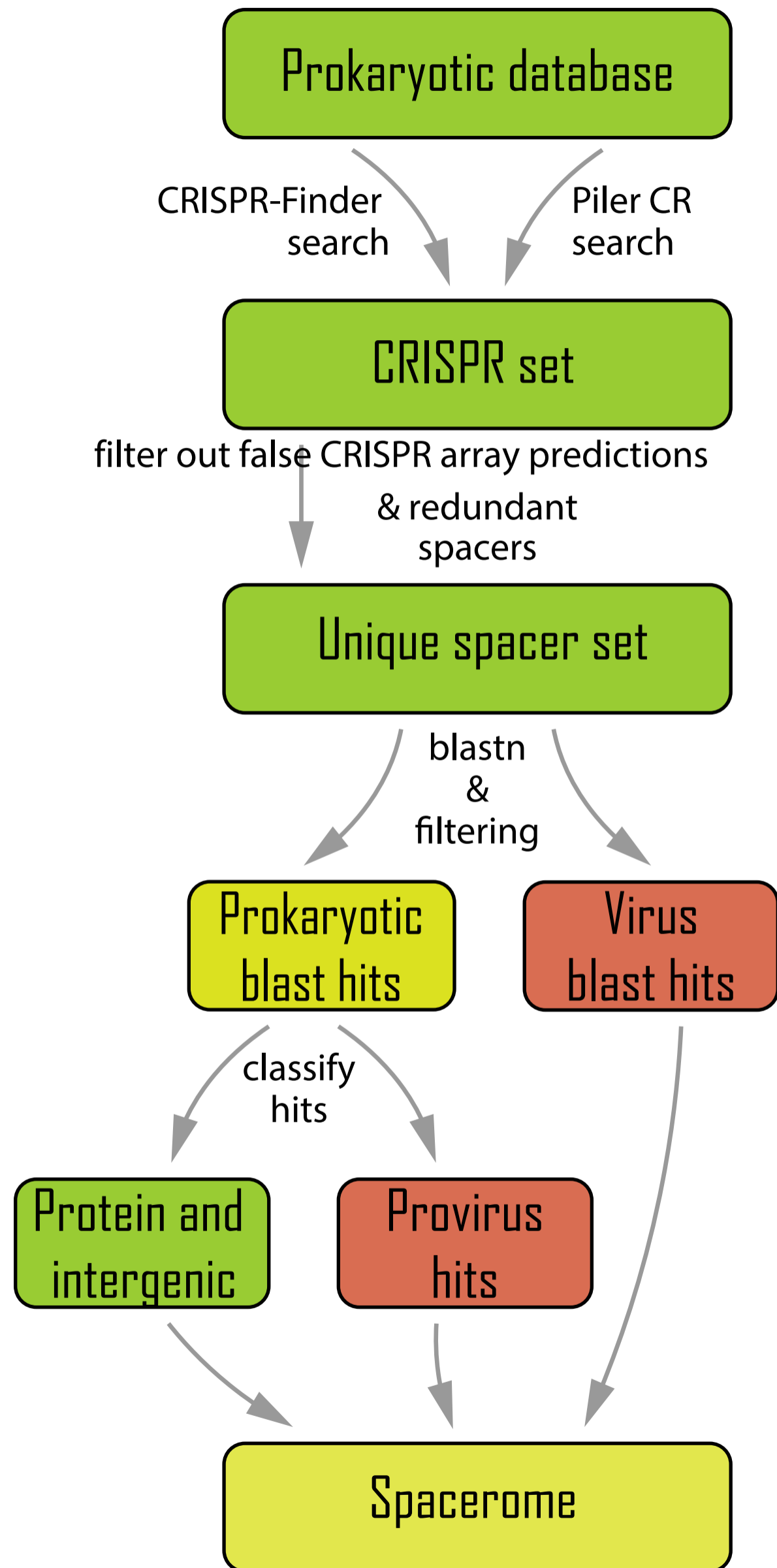
Top 10 BLASTX results in RefSeq:

Fields: query id, subject id, subject ids, subject length, s. start, s. end, evalue, query seq, subject seq, q. start, q. end, score, subject sci names, subject title

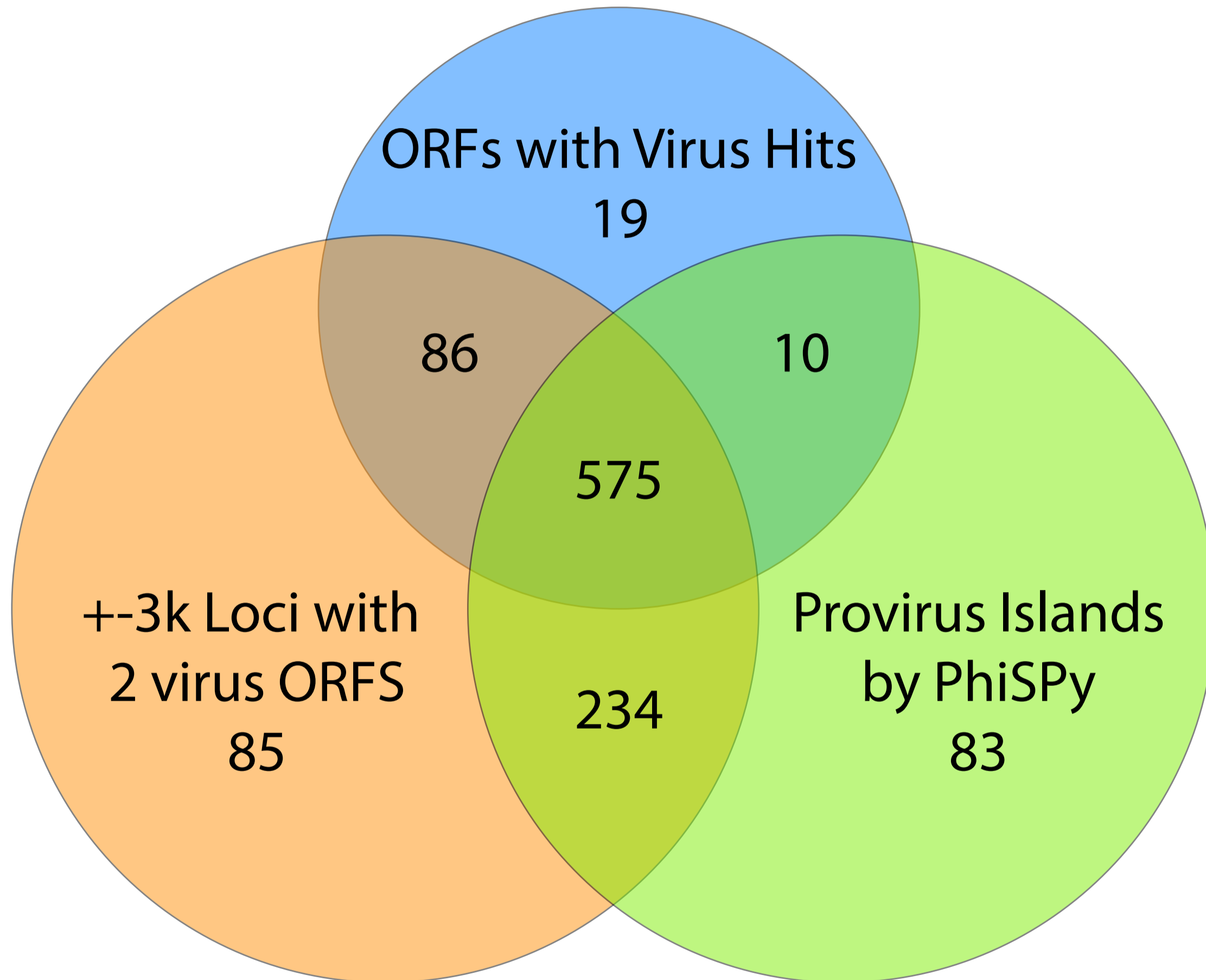
99 hits found

gi|844486473|gb|KQ087587.1|:84822-86069 gi|1172898100|ref|WP_080967648.1|
gi|1172898100|ref|WP_080967648.1| 135 11 135 2.47e-77
RLPNRSPDKALCAAIREPSPVALTLTGATNRIIPNRSPDKARCAAIREPSPVALTLTGATDRIPNRSPDKAPCAAIREPSPVALTLTGATNRLPNRSPD
KALCAAIREPSPGGAYAYRGYQPHTEP
RLPNRSPDKALCAAIREPSPVALTLTGATNRIIPNRSPDKARCAAIREPSPVALTLTGATDRIPNRSPDKAPCAAIREPSPVALTLTGATNRLPNRSPD
KALCAAIREPSPGGAYAYRGYQPHTEP 1246 872 623 Enterobacter aerogenes hypothetical protein
[Klebsiella aerogenes]
gi|844486473|gb|KQ087587.1|:84822-86069 gi|1172898100|ref|WP_080967648.1|
gi|1172898100|ref|WP_080967648.1| 135 23 135 2.63e-68
AAIREPSPVALTLTGATNRIIPNRSPDKARCAAIREPSPVALTLTGATDRIPNRSPDKAPCAAIREPSPVALTLTGATNRLPNRSPDKALCAAIREPSP
GGAYAYRGYQPHTEP
AAIREPSPVALTLTGATNRIIPNRSPDKARCAAIREPSPVALTLTGATDRIPNRSPDKAPCAAIREPSPVALTLTGATNRLPNRSPDKALCAAIREPSP
GGAYAYRGYQPHTEP 447 109 563 Enterobacter aerogenes hypothetical protein [Klebsiella
aerogenes]
gi|844486473|gb|KQ087587.1|:84822-86069 gi|1172898100|ref|WP_080967648.1|
gi|1172898100|ref|WP_080967648.1| 135 11 135 1.27e-18
RLPELPTFYRIVARIRRDAPLSGNPPRWRLRPLPTEYRTVARIRRHAPPSGNPPRWRLRPLPTEYRTVARIRRYAPPSGNPPVALTLTGATNR
IPNRSPDKALRAAIREPSPGGAYAYRGYQPHTEP
RLPNRSPDKALCAAIREPSPVALTLTGATNRIIPNRSPDKARCAAIREPSPVALTLTGATDRIPNRSPDKAPCAAIRE-----
PSPVALTLTGATNRLPNRSPDKALCAAIREPSPGGAYAYRGYQPHTEP 413 18 225 Enterobacter aerogenes
hypothetical protein [Klebsiella aerogenes]
gi|844486473|gb|KQ087587.1|:84822-86069 gi|1172898100|ref|WP_080967648.1|
gi|1172898100|ref|WP_080967648.1| 135 11 135 1.60e-18
RLPELPTFYRIVARIRRDAPLSGNPPRWRLRPLPTEYRTVARIRRHAPPSGNPPRWRLRPLPTEYRTVARIRRYAPPSGNPPVALTLTGATNR
IPNRSPDKALRAAIREPSPGGAYAYRGYQPHTEP
RLPNRSPDKALCAAIREPSPVALTLTGATNRIIPNRSPDKARCAAIREPSPVALTLTGATDRIPNRSPDKAPCAAIRE-----
PSPVALTLTGATNRLPNRSPDKALCAAIREPSPGGAYAYRGYQPHTEP 1176 781 224 Enterobacter aerogenes
hypothetical protein [Klebsiella aerogenes]
gi|844486473|gb|KQ087587.1|:84822-86069 gi|1172898100|ref|WP_080967648.1|
gi|1172898100|ref|WP_080967648.1| 135 27 61 1.06e-13

NPPPVALTLTGATNQIPNRSPDKALCAAIRVFTPL	EPSPVALTLTGATNRIPNRSPDKARCAAIRESSPV	836	732
130	Enterobacter aerogenes hypothetical protein [Klebsiella aerogenes]		
gi 844486473 gb KQ087587.1 :84822-86069	gi 1172898100 ref WP_080967648.1		
gi 1172898100 ref WP_080967648.1	135 1 30 1.06e-13		
VALTLTGATNRIPNRSPDKALRAAIRESSP	MALTLTGATNRLPNRSPDKALCAAIREPSP	915	826 121
130	Enterobacter aerogenes hypothetical protein [Klebsiella aerogenes]		
gi 844486473 gb KQ087587.1 :84822-86069	gi 1172898100 ref WP_080967648.1		
gi 1172898100 ref WP_080967648.1	135 1 30 2.18e-10		
VALTLTGATNRIPNRSPDKALRAAIRESSP	MALTLTGATNRLPNRSPDKALCAAIREPSP	152	63 121
130	Enterobacter aerogenes hypothetical protein [Klebsiella aerogenes]		
gi 844486473 gb KQ087587.1 :84822-86069	gi 1172898100 ref WP_080967648.1		
gi 1172898100 ref WP_080967648.1	135 27 50 2.18e-10		
EPSPVALTLTGATNRIPNRSPDKA 73	2 101	Enterobacter aerogenes hypothetical protein	
[Klebsiella aerogenes]			
gi 844486473 gb KQ087587.1 :84822-86069	gi 1172898100 ref WP_080967648.1		
gi 1172898100 ref WP_080967648.1	135 87 120 1.13e-05		
NPPPVALTLTGATNQIPNRSPDKALCAAIRVFTPL	EPSPVALTLTGATNRLPNRSPDKALCAAIREPSP	836	735
130	Enterobacter aerogenes hypothetical protein [Klebsiella aerogenes]		
gi 844486473 gb KQ087587.1 :84822-86069	gi 1172898104 ref WP_080967652.1		
gi 1172898104 ref WP_080967652.1	105 1 105 9.68e-63		
VALTLTGATNRIPNRSPDKARCAAIRESSPV	ALTLTGATDRIPNRSPDKAPCAAIREPSP	VALTLTGATNRLPNRSPDKALCAAIREPSP	GGAYAYRG
YQPHTEP			
MALTLTGATNRIPNRSPDKARCAAIRESSPV	ALTLTGATDRIPNRSPDKAPCAAIREPSP	VALTLTGATNRLPNRSPDKALCAAIREPSP	GGAYAYRG
YQPHTEP 423	109 523	Enterobacter aerogenes hypothetical protein [Klebsiella aerogenes]	

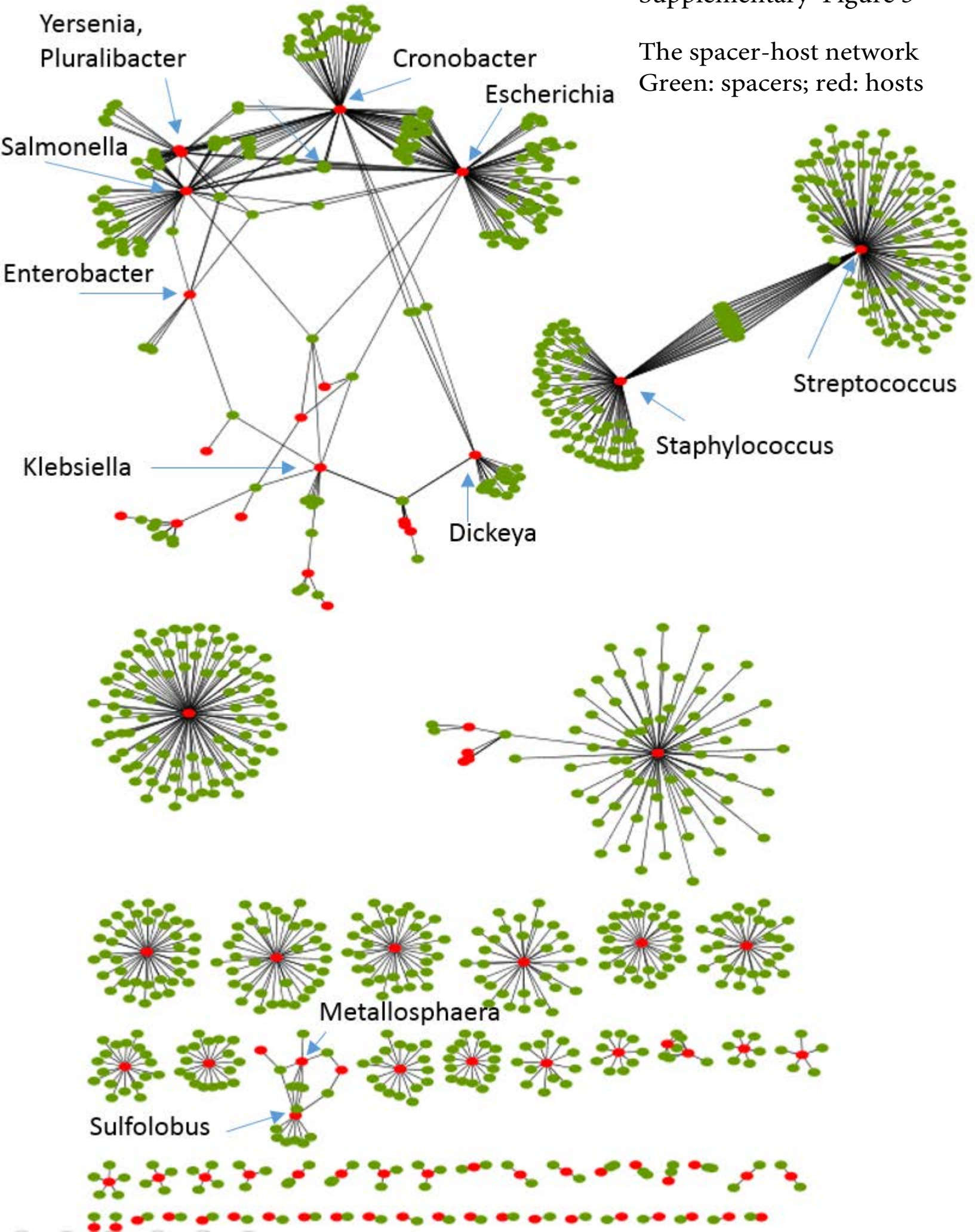


Supplementary Figure 1. The computational pipeline for CRISPR spacer analysis

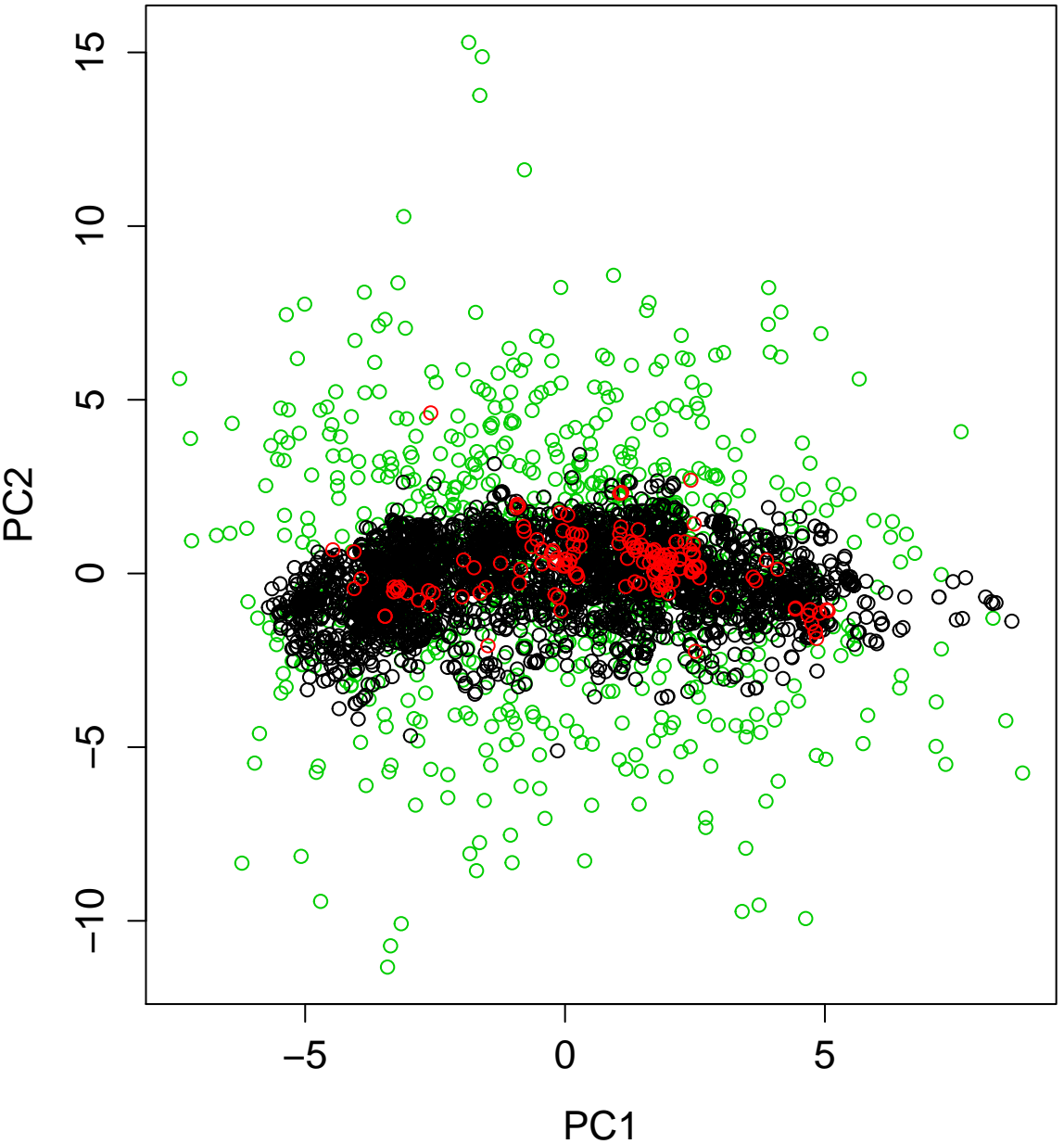


Supplementary Figure 2. Identification of prophage regions with different approaches

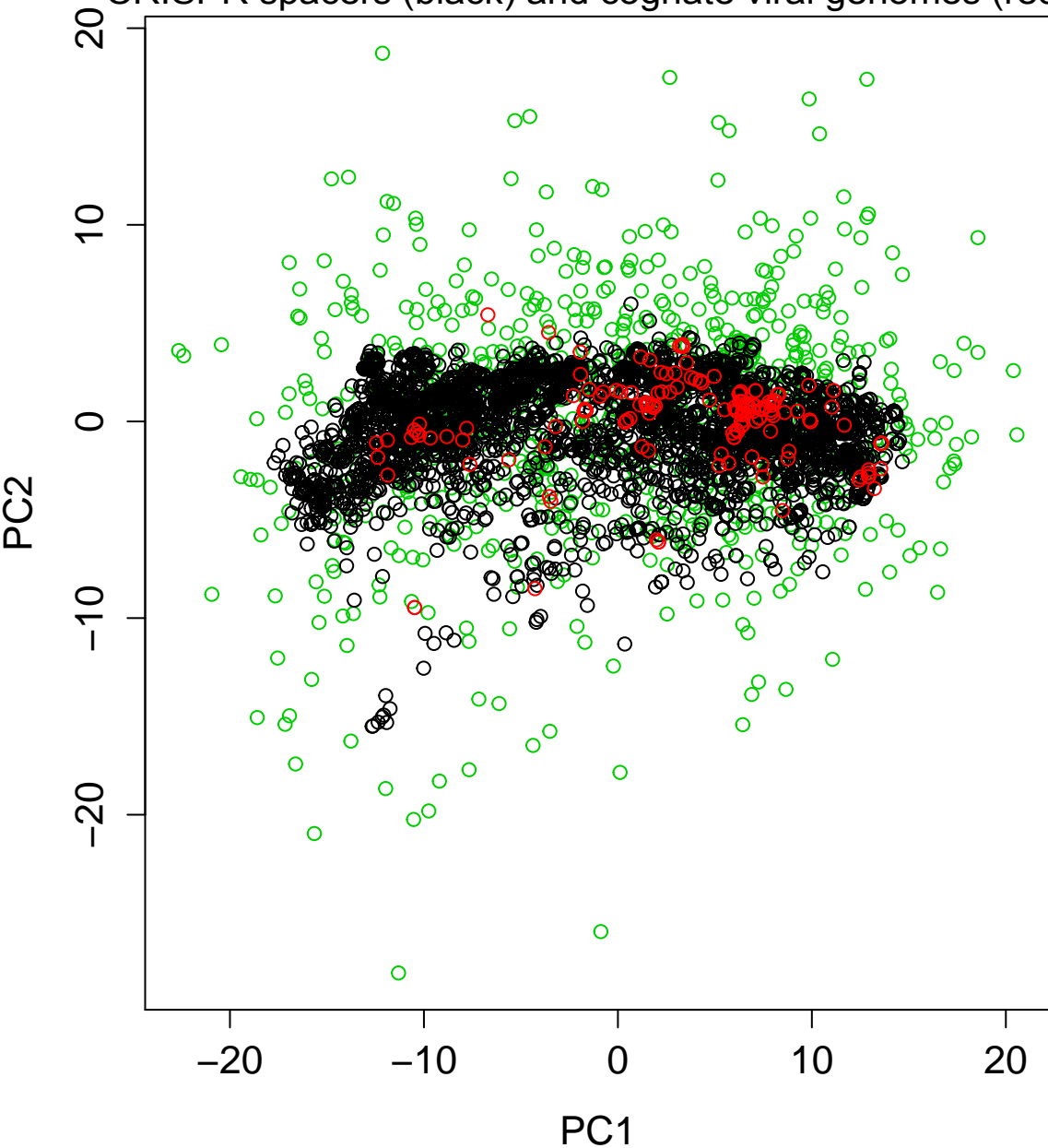
The spacer-host network
Green: spacers; red: hosts



Supplementary figure 4: Principal component analysis of the dinucleotide frequencies in microbial genomes (green), CRISPR spacers (black) and cognate virus genomes (red)



Supplementary figure 5: principal component analysis of tetranucleotide frequencies in microbial genomes (green), CRISPR spacers (black) and cognate viral genomes (red)



Supplementary figure 6 : Defining the threshold for CRISPR spacer identification

