

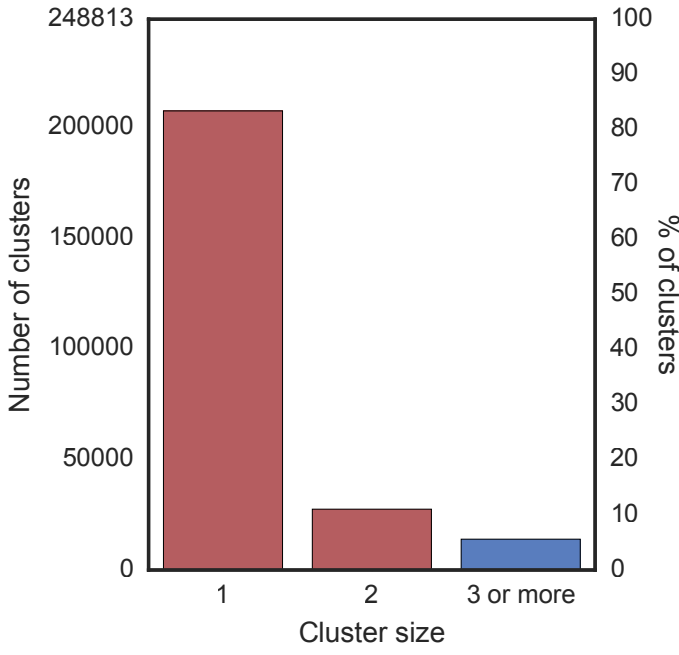
# Discovering viral genomes in human metagenomic data by predicting unknown protein families

Mauricio Barrientos-Somarribas, David N. Messina, Christian Pou, Fredrik Lysholm, Annelie Bjerkner, Tobias Allander, Björn Andersson and Erik L.L. Sonnhammer

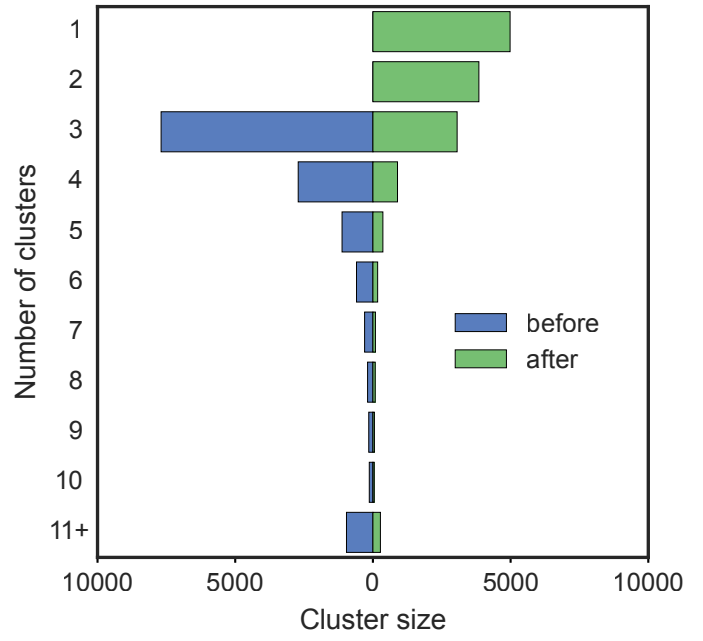
## Supplementary Material

**Supplementary Figure S1.** Detection of high-confidence protein families by clustering and RNAcode calibration. a) Distribution of cluster sizes from MCL. The size distribution was split into three categories: singletons, clusters of size two, and clusters of size 3 and above (candidate for multiple sequence alignment) b) Comparison of cluster size distribution for the filtered MCL clusters ( $3 < \text{size} < 250$ ) before and after subclustering. Cluster size decreases post-subclustering, indicating that the original clusters contained many sequences with low diversity c) Calibration of RNAcode for short sequences. The plot describes the distribution of RNAcode's p-values for simulated alignments of varying lengths from coding and non-coding regions of the *Methanococcus jannaschii* genome. RNAcode separates coding from non-coding sequences up to a length of 50bp with a p-value of 0.15.

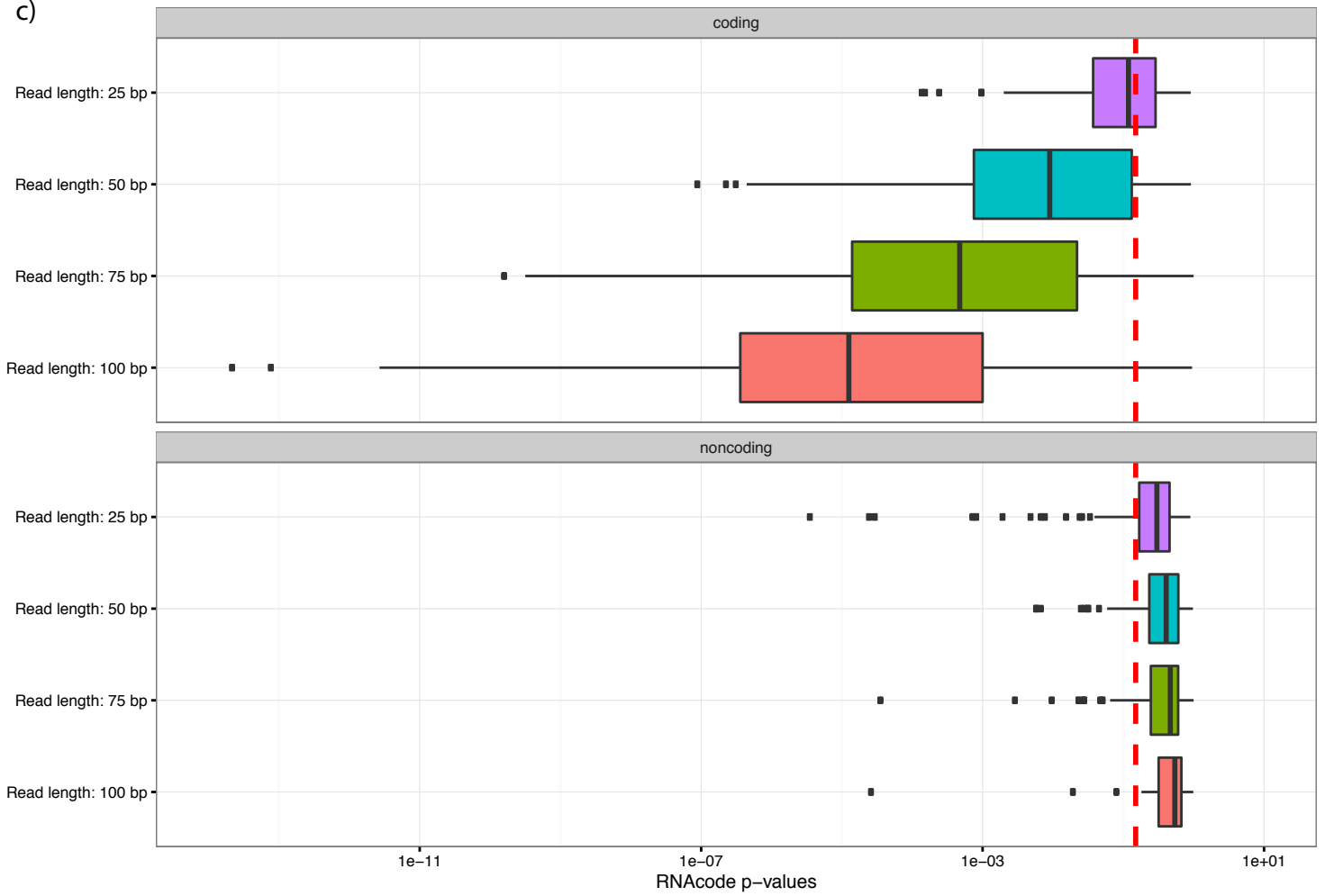
a) MCL cluster size distribution



b) Filtered MCL cluster size before and after subclustering



c)



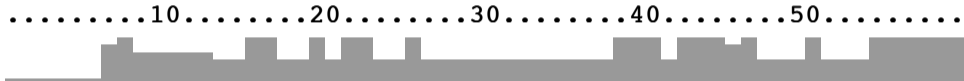
**Supplementary Figure S2.** RNAcode output for the 32 predicted protein families.

cluster113b

Frame +1 p = 0.008

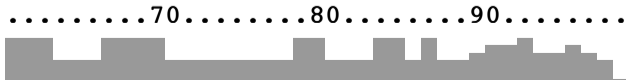
M Y A N R L V N R Q L E K K R I E E

GB3LKKR01C5IFF/1-92 -----**ATGTATGCAAACCGTTTGGTTAATCGTCAACTAGAGAAAAAGAGAATAGAAGAA** 54  
 GB3LKKR01CYB4U/1-92 -----ATGTATGCAAACCGTTTGGTTAATCGTCAACTAGAGAAAAAGAGAATAGAAGAA 54  
 GB3LKKR02GWFIV/1-91 -----ATGTATGCAAACCGTTTGGTTAATCGTCAACTAGAGAAA-AGAGAATAGAAGAA 53  
 GB3LKKR02GSBUP/1-91 -----ATGTATGCAAACCGTTTGGTTAATCGTCAACTAGAGAAAAAGAGAATAGAAGAA 54  
 GB3LKKR01BCTJS/1-96 AAAATAATGTATG--**AATAGATTAATAGCAAAGAGCAAAGATAAAAATCTATCTGAAGAA** 58  
 GB3LKKR02F10X4/1-98 **GAAAAATAATGATGAATAGATTAATAGCAAAGAGCAAAGATAAAAATCTATCTGAAGAA** 60  
 GB3LKKR02HWM0Y/1-98 **GAAAACTAATGATGAATAGATTAATAGCAAAGAGCAAAGATAAAAATCTATCTGAAGAA** 60



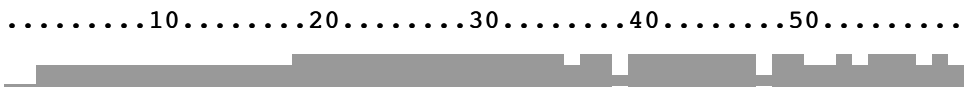
Y K E D H P Q T V E E F

GB3LKKR01C5IFF/1-92 **TATAAAGAAGACCACCCTCAGACAGTAGAAGAATTC** 92  
 GB3LKKR01CYB4U/1-92 TATAAAGAAGACCACCCTCAGACAGTAGAAGAATTC 92  
 GB3LKKR02GWFIV/1-91 TATAAAGAAGACCACCCTCAGACAGTAGAAGAATTC 91  
 GB3LKKR02GSBUP/1-91 TATAAAGAAGACCACCCTCAGACAGTAGAAGAATTC 91  
 GB3LKKR01BCTJS/1-96 TATCGGGAAGCATGTTGG**CAAGTAGCATGCGAAATC** 96  
 GB3LKKR02F10X4/1-98 TATCGGGAAGCATGTTGG**CAAGTAGCA-TACGATAC** 98  
 GB3LKKR02HWM0Y/1-98 TATCGGGAAGCATGTTGG**CAAGTAGCATGCGAAA-A** 98



cluster179a

GB3LKKR01D9539/1-314 257  
 GB3LKKR01B7VFA/1-313 256  
 GB3LKKR01A1AS6/1-174 175  
 GB3LKKR01DZB8X/1-177 136  
 GB3LKKR02JP9RY/1-314 257



GB3LKKR01D9539/1-314 **TAGATAATAAAGGTGGTGGTGCCTTATACTCTTTT---ATGTATCCGCGTT** 197  
 GB3LKKR01B7VFA/1-313 **TAGATAATAAAGGTGGTGGTGCCTTATACTTCTTTT---ATGTATCCGCGTT** 196  
 GB3LKKR01A1AS6/1-174 **TAGATAATAACGGTGTGTCACCGTATACCTTCTTATTATGTATCCGCGTC** 175  
 GB3LKKR01DZB8X/1-177 ----- 76  
 GB3LKKR02JP9RY/1-314 **TAGATAAATCTGCTCAACTTCGTACACTTCTTTT---ATGTATCCGCGTC** 197



Frame -2 p = 0.005  
 R Y C C F F S Q I M P F N H R A A L P N  
 GB3LKKR01D9539/1-314 **CCGGTACTGCTGCTTTTTTCAGTCAGATTATGCCGTTTAATCATCGTGCCTGCACTTCCCTAA** 137  
 GB3LKKR01B7VFA/1-313 CCGTACTGCTGCTTTTTTCAGTCAGATTATGCCGTTTAATCATCGTGCCTGCACTT**CCAAA** 136  
 GB3LKKR01A1AS6/1-174 **CCGGTGT**TGCTGCATTTTCAGTCAGGTTATGCCGTTTAATCATCGTGCCTGCACTT**CCAAA** 127  
 GB3LKKR01DZB8X/1-177 -----TTCTGCATTTTCAGTCAGTTATGCCGTTTAATCATCGTGCCTGCACTT**CCAAA** 16  
 GB3LKKR02JP9RY/1-314 **CCGGTGT**TGCTGCATTT**TTT**AGTCAGGTTATGCCGTTTAATCATCGTGCCTGCACTT**CCAAA** 137



S L M S W F R I A N S P I F N Y S D N P  
 GB3LKKR01D9539/1-314 **TAGTCTGATGTCCTTGGTTTCGTATTGCTAATAGTCCGATTTTAAATTATTCTGATAATCC** 80  
 GB3LKKR01B7VFA/1-313 TAGT**TTG**ATG**TCCG**TGGCTTCGTATTGCTAATAGT**CCT**ATTATCAATTAT**CCTAAGAACAC** 79  
 GB3LKKR01A1AS6/1-174 **TAGGTTGATA**----- 67  
 GB3LKKR01DZB8X/1-177 TAGT**TTG**ATG**TCCG**TGGCTTCGTATTGCTAATAGT**CCT**ATTCTCAATTAT**CCTAAGAATAC** 1  
 GB3LKKR02JP9RY/1-314 TAGT**TTG**ATG**TCCG**TGGCTTCGTATTGCTAATAGT**CCT**ATTATCAATTAT**CCTAAGAACAC** 80



L P S G G V L F K S S P K F L S V N A D  
 GB3LKKR01D9539/1-314 **TTTGCTTCTGGTGGAGTTTTGTTAAGTCGTCGCCGAAGTTTCTTCTGTGAATGCAGA** 23  
 GB3LKKR01B7VFA/1-313 **TGTG**CCTACT**ACTGATGCTCTT**CTCAAGACATCG**GCG**AAGTTTGTTACTGTGAATGCAGA 22  
 GB3LKKR01A1AS6/1-174 ----- 7  
 GB3LKKR01DZB8X/1-177 **TCTG**CCTATT**GATAGTTCTCTT**CTCAAGACAGCG**TCC**AAGTTTCTTACTGTGAATGCAGA 1  
 GB3LKKR02JP9RY/1-314 **TGTG**CCTACT**ACTGATGCTCTT**CTCAAGACATCG**GCG**AAGTTTGTTACTGTGAATGCAGA 23



T Y L G Y W  
 GB3LKKR01D9539/1-314 **TACTTATTTAGGTTATTGG** 1  
 GB3LKKR01B7VFA/1-313 TACTTATTTAGGTTATTGG 1  
 GB3LKKR01A1AS6/1-174 ----- 1  
 GB3LKKR01DZB8X/1-177 TACT----- 1  
 GB3LKKR02JP9RY/1-314 TACTTATTTAGGTTATTGG 1

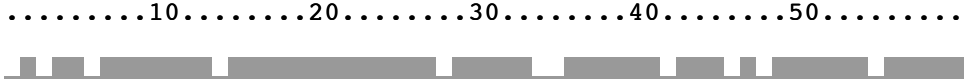


cluster179b

Frame -1 p= 0.010

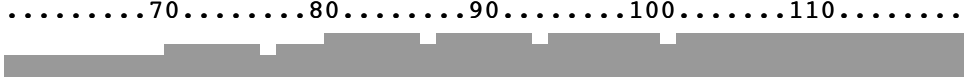
\* I S I C I H S **K K L R R C**

GB3LKKR01C2BL0/1-231 **TAAATAAGTATCTGCATTCACAGTAAGAAACTTCGACGCTG** 172  
 GB3LKKR01EB2BL/1-168 **TAAATAAGTATCTGCATTCACAGTAAGAAACTTCGACGCTG** 169  
 GB3LKKR02ITVHK/1-237 **TAAATAAGTATCTGCATTCACAGTAACAAACTTC**CGCCGATG**** 179  
 GB3LKKR02JWM40/1-158 **TAAATAAGTATCTGC**ACTCACAGAA**AGAAACTTC**CGGCCGAGC**** 159  
 GB3LKKR02IKPFQ/1-162 ----- 103



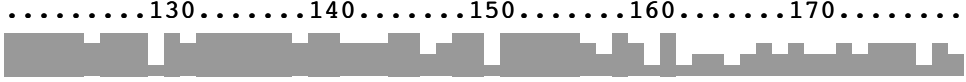
**L E K R T I N R Q S I L R I I E N R T I**

GB3LKKR01C2BL0/1-231 **TCTTGAGAAGAGAACTATCAATAGGCAGAGTATTC**T**TAGGATAA**T**TGAGAATAGGACTAT** 112  
 GB3LKKR01EB2BL/1-168 TCTTGAGAAGAGAACTATCAATAGGCAGAGTATTC**T**TAGGATAA**T**TGAGAATAGGACTAT 119  
 GB3LKKR02ITVHK/1-237 **C**CTTGAGAAGAG**AGCATCAGT**AGTAGGCACAGTGT**T**CTTAGGATAA**T**TGATAATAGGACTAT 119  
 GB3LKKR02JWM40/1-158 **A**CTTAAACAA**AACTCCACC**AGTAGGCACAGTGT**T**CTTAGGATAA**T**TGATAATAGGACTAT 119  
 GB3LKKR02IKPFQ/1-162 -----**ACC**AGAAG**CAAAGG**ATT**ACC**AGAATAAT**T**AAAA**T****CGG**ACTAT 43



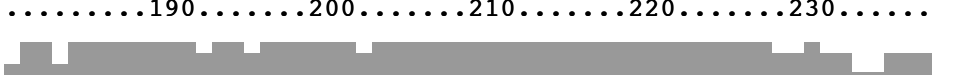
**S N T K P R H Q T I W K C S T M I K R H**

GB3LKKR01C2BL0/1-231 **TAGCAATACGAAGCCACGACATCAA**A**CTATTTGGAAGTGCAGC**C**CGATGAT**T**AAACGGCA** 52  
 GB3LKKR01EB2BL/1-168 TAGCAATACGAAGCCACGACATCAA**A**CTATTTGGAAGTGCACCAG**AT**----- 59  
 GB3LKKR02ITVHK/1-237 TAGCAATACGAAGCCACGACATCAA**A**CTATTTGGAAGTGCAGC**C**CGATGAT**T**AAACGGCA 59  
 GB3LKKR02JWM40/1-158 TAGCAATACGAAGCCACGACATCAA**A**CTATTTGGAAGT**---**----- 59  
 GB3LKKR02IKPFQ/1-162 TAGCAATACGAAGCCA**AGAC**AT**C**AGACTAT**T**AGGAAGTGCAGC**C**CGATGAT**T**AAACGGCA 1



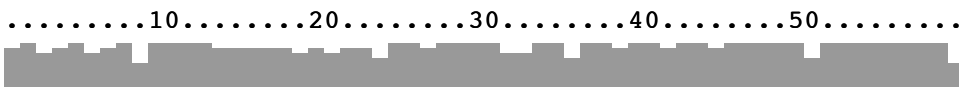
**N L T E K C S N T G T R I H I R S I R**

GB3LKKR01C2BL0/1-231 **TAACCTGACTGAAAAATGCAGCAAC**ACC**GGGACGCGGATACATATAAGAA**GT**ATACGG** 1  
 GB3LKKR01EB2BL/1-168 ----- 1  
 GB3LKKR02ITVHK/1-237 TAACCTGACT**AA**-AAATGCAGCAAC**ACC**GGGACGCGGATACAT**AAA**AGAAGTGT**ACGA** 1  
 GB3LKKR02JWM40/1-158 ----- 1  
 GB3LKKR02IKPFQ/1-162 **T**AATCTGACTGAAAA**AGC**AGCAGTACC**GGA**ACGCGGATACAT**AAA**AGAAGTATA**AGC** 1



cluster182a

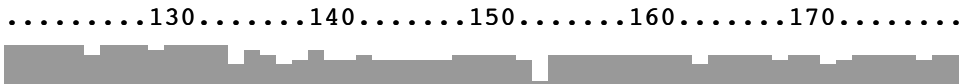
GB3LKKR02H8MH7/1-364 306
GB3LKKR01C3KIQ/1-328 269
GB3LKKR02H78NR/169-302 75
GB3LKKR02I75TX/1-273 214
GB3LKKR01AM7C8/22-109 30
GB3LKKR02IRLPC/358-499 83
contig08842/31-253 164
GB3LKKR02H8NEG/1-364 305
GB3LKKR01EC368/1-365 306
contig04034/238-604 308



GB3LKKR02H8MH7/1-364 247
GB3LKKR01C3KIQ/1-328 209
GB3LKKR02H78NR/169-302 15
GB3LKKR02I75TX/1-273 154
GB3LKKR01AM7C8/22-109 1
GB3LKKR02IRLPC/358-499 23
contig08842/31-253 104
GB3LKKR02H8NEG/1-364 245
GB3LKKR01EC368/1-365 247
contig04034/238-604 248



GB3LKKR02H8MH7/1-364 188
GB3LKKR01C3KIQ/1-328 149
GB3LKKR02H78NR/169-302 1
GB3LKKR02I75TX/1-273 95
GB3LKKR01AM7C8/22-109 1
GB3LKKR02IRLPC/358-499 1
contig08842/31-253 44
GB3LKKR02H8NEG/1-364 185
GB3LKKR01EC368/1-365 188
contig04034/238-604 188



GB3LKKR02H8MH7/1-364
GB3LKKR01C3KIQ/1-328
GB3LKKR02H78NR/169-302
GB3LKKR02I75TX/1-273
GB3LKKR01AM7C8/22-109
GB3LKKR02IRLPC/358-499
contig08842/31-253
GB3LKKR02H8NEG/1-364
GB3LKKR01EC368/1-365
contig04034/238-604

\* R T Q K P
TAACGGACTCAAAAACC 128
TAA CGGACTCAAAAACC 89
-----AAAATCC 1
TAA CGGACTCAAAAACC 35
----- 1
--CGGACCTTAAAAACC 1
AACGGCCTTAAAAAACC 1
TAA CGGACTCAAAAACC 126
TAA CGGACTCAAAAACC 128
TAA CGGTCTTAAAAACC 128



GB3LKKR02H8MH7/1-364
GB3LKKR01C3KIQ/1-328
GB3LKKR02H78NR/169-302
GB3LKKR02I75TX/1-273
GB3LKKR01AM7C8/22-109
GB3LKKR02IRLPC/358-499
contig08842/31-253
GB3LKKR02H8NEG/1-364
GB3LKKR01EC368/1-365
contig04034/238-604

Frame -1 p = 5.2e-04
A K T N A T S M G V R K E G M G L Q Y
AGCAAAAATAATGCAACAAGCATGGGAGTACGAAAAG-AGGGGATGGGTCTACAATACA 68
AGCAAAAATAATGCAACAAGCATGGGAGTACGAAAAGAGGGGATGGGTCTACAATACA 29
AGCAAAAATAATGCAACAAGCATGGGAGTACGAAAAGAGGGGATGGGTCTACAATACA 1
AGCAAAAATAATGCAACAAGCATGGGAGTACGAAAAGAGGGGATGGGTCTACAATACA 1
-----GAAGAAG-TGGGGCTACAATATA 1
AACAAAAATAATGCAACAGGCATGGGAGTATGAAAAGAAAGGATGGGTCTGCAATATA 1
AACAAAAATAATGCAACAAGCATGGGAGTATGAAAAGAAAGGATGGGTCTACAATATA 1
AGCAAAAATAATGCAACAAGCATGGGAGTACGAAAAGAGGGGATGGGTCTACAATACA 67
AGCAAAAATAATGCAACAAGCATGGGAGTACGAAAAG-AGGGGATGGGTCTACAATACA 68
AGCAAAAATAATGAAACAAGCATGGGAATACGAAAAGGAAAGGATGGGCTGCAATACA 68



GB3LKKR02H8MH7/1-364
GB3LKKR01C3KIQ/1-328
GB3LKKR02H78NR/169-302
GB3LKKR02I75TX/1-273
GB3LKKR01AM7C8/22-109
GB3LKKR02IRLPC/358-499
contig08842/31-253
GB3LKKR02H8NEG/1-364
GB3LKKR01EC368/1-365
contig04034/238-604

N Y G Q Q A A D A E Y K R N L O M W K E
ACTATGGGCAACAAGCAGCAGACGCTGAATATAAACGAAATCTACAAATGTGGAAAGAAA 8
ACTATGGGCAACAAGCAGCAGACGCTGAATATAAACGAAATCTACAAATGTGGAAAGAAA 1
ACTATGGGCAACAAGCAGCAGACGCTGAATATAAACGAAATCTACAAATGTGGAAAGAAA 1
ACTATGGGCAACAAGCAGCAGACGCTGAATATAAACGAAATCTACAAATGTGGAAAGAAA 1
ATTACGGACAACAAGCGGCAGATGCAGATGATAAACGAAATCTGAGCAAAAC-AAAGACA 1
ATTACGGACAACAAGCGGCAGATGCAGATGATAAACGAAATCTGAGCAAAATGTGGAAAGACA 1
ATTACGGACAACAAGCGGCAGACGCTGAATATAAACGAAATCTGCAAAATGTGGAAAGACA 1
ACTATGGGCAACAAGCAGCAGACGCTGAATATAAACGAAATCTACAAATGTGGAAAGAAA 7
ACTATGGGCAACAAGCAGCAGACGCTGAATATAAACGAAATCTACAAAGCTGGAAAGAAA 8
ACTATGGGCAACAAGCAGCAGACGCTGAATATAAACGAAATCTACAAATGTGGAAAGATA 8



GB3LKKR02H8MH7/1-364
GB3LKKR01C3KIQ/1-328
GB3LKKR02H78NR/169-302
GB3LKKR02I75TX/1-273
GB3LKKR01AM7C8/22-109
GB3LKKR02IRLPC/358-499
contig08842/31-253
GB3LKKR02H8NEG/1-364
GB3LKKR01EC368/1-365
contig04034/238-604

T N
CCAAC 1
CCAAC 1
CCAAC 1
CCAAC 1
GTAGC 1
CCAAT 1
CCAAC 1
CCAAC 1
CCAAC 1
CCAAC 1
CTAAT 1



cluster182b

GB3LKKR02JSIBW/4-194	30
GB3LKKR02G4CF7/1-191	55
GB3LKKR02GLIMM/1-188	18
GB3LKKR02JJI0V/1-222	54
GB3LKKR02JAAI8/38-264	60
GB3LKKR01EDLON/5-176	59
GB3LKKR02GK10E/263-445	60
GB3LKKR02IIMCY/1-227	59
contig10047/1-228	60
GB3LKKR01DC008/1-71	0
GB3LKKR02J3V10/1-105	0

.....10.....20.....30.....40.....50.....



GB3LKKR02JSIBW/4-194	<b>TAATGCAACAAGCGTGGGAGTATGAAAA-G-AGGAATGGGG--</b>	82
GB3LKKR02G4CF7/1-191	<b>TAA</b> TGCAACAAGCGTGGGAGTATGAAAAAGAAGGAATGGGG--	113
GB3LKKR02GLIMM/1-188	CCAAAACAAT <b>TGCCGAAGCAAATAAAATCGCAGGA-GTGGACA</b>	77
GB3LKKR02JJI0V/1-222	<b>TAA</b> TGCAACAAGCATGGGAGTATGAAAAAGAGGGAAT <b>GGT--</b>	112
GB3LKKR02JAAI8/38-264	<b>TAA</b> TGCAACAAGGCATGGGAGTATGAAAAAGAA <b>GGCATGGGT--</b>	118
GB3LKKR01EDLON/5-176	<b>TAA</b> TGCAACAAGCGTGGGAGTATGAAAAAGAAGGAATGGGG--	117
GB3LKKR02GK10E/263-445	<b>TAA</b> TGCAACAAGCATGGGAGT <b>ACG</b> AAAAAGAG <b>GGGATGGGT--</b>	118
GB3LKKR02IIMCY/1-227	<b>TAA</b> TGCAACAAGCATGGGAGT <b>ACG</b> AAAAAGN <b>GGGATGGGG--</b>	117
contig10047/1-228	<b>TAA</b> TGCAACAAGCATGGGAGT <b>ACG</b> AAAAAGAG <b>GGGATGGGT--</b>	118
GB3LKKR01DC008/1-71	-----	0
GB3LKKR02J3V10/1-105	-----	0

Frame +2 p = 1.5e-04

\* C N K R G S M K R G M G

.....70.....80.....90.....100.....110.....

.....70.....80.....90.....100.....110.....



GB3LKKR02JSIBW/4-194	<b>-CTACAATATAATTACGGACAACAAGCGGCAGACGCTGAATACAAGAG-AATCTGCAAA</b>	140
GB3LKKR02G4CF7/1-191	<b>-CTACAATATAATTACGGACAACAAGCGGCAGACGCTGAATACAAGAGAAATCTGCAAA</b>	172
GB3LKKR02GLIMM/1-188	<b>CAGAAGGGGCAAAACTAGACAACGAATGGAA</b> GACGCTGAATACAAGAGAAATCTGCAAA	137
GB3LKKR02JJI0V/1-222	<b>-CTGCAATACAACCTATGGACAACAAGCGGCAGATGCTGAATACAAGCGGAACCTACA</b> AAAT	171
GB3LKKR02JAAI8/38-264	<b>-CTGCAATATAATTACGGACAACAAGCGGCAGATGCAGAGTATAAACG-AATCTGCAGAT</b>	176
GB3LKKR01EDLON/5-176	<b>-CTACAATATAATTACGGACAACAAGCGGCAGACGCTGAATACAAGAGAAATCTGC----</b>	172
GB3LKKR02GK10E/263-445	<b>-CTACAATACAACCTATGGGCAACAAGCAGCAGACGC-GAAATAAAACGAAATCTAAAATG</b>	176
GB3LKKR02IIMCY/1-227	<b>TCCTACAATACAACCTATGGGCAACAAGCAGCAGACCTGAAATAAAACG-AATCTACA</b> AAAT	176
contig10047/1-228	<b>-CTACAATACAACCTATGGGCAACAAGCAGCAGACGCTGAAATAAAACGAAATCTACA</b> AAAT	177
GB3LKKR01DC008/1-71	----- <b>TATAAACGAAATCTGCAGAT</b>	20
GB3LKKR02J3V10/1-105	<b>-ATATCGCGAGATCATGGGCAACAAGCAGCAGACGCTGAAATAAAACGAAATCTACA</b> AAAT	59

.....130.....140.....150.....160.....170.....



GB3LKKR02JSIBW/4-194	<b>GTGGAAAGACACCAACTTTGGAGCACAAAGAGCCGAGATGGAAAACGCAG</b>	191
GB3LKKR02G4CF7/1-191	GTGGAAAGACACCAACT <b>TTT</b> -----	191
GB3LKKR02GLIMM/1-188	GTGGAAAGACACCAACTTTGGAGCACAAAGAGCCGAGATGGAAAACGCAG	188
GB3LKKR02JJI0V/1-222	GTGGAAAGACACCAACTTTGGAGCACAAAGAGCCGAGATGG <b>AGATGCAG</b>	222
GB3LKKR02JAAI8/38-264	GTGGAAAGACACCAACTTTGGAGCACAAAGAG <b>ACG</b> AAATGGAGAAAGCGG	227
GB3LKKR01EDLON/5-176	-----	172
GB3LKKR02GK10E/263-445	<b>TAGGAAA-</b> -----	183
GB3LKKR02IIMCY/1-227	GTGGAAAGAAACCAACTTTGGAGCGCAAAGGAATGAAATGGAAA <b>AGCAG</b>	227
contig10047/1-228	GTGGAAAGAAACCAACTTTGGAGCGCAAAGGAATGAAATGGAAA <b>AGCAG</b>	228
GB3LKKR01DC008/1-71	GTGGAAAGACACCAACTTTGGAGCACAAAGAG <b>ACG</b> AAATGGAGAAAGCGG	71
GB3LKKR02J3V10/1-105	GTGGAAAG <b>AAA</b> CAACTTTGGAGCGCAAAGGAATGAAATGGAAAA <b>----</b>	105

.....190.....200.....210.....220.....230.....



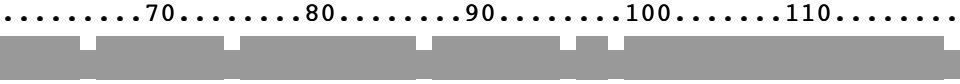
cluster211a

Frame -3 p = 3.9e-06

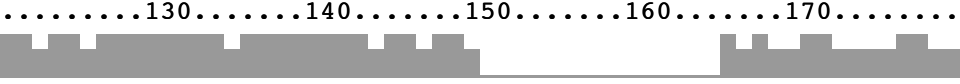
	A	V	A	G	G	T	Y	R	D	W	L	E	T	V	
GB3LKKR01EWCHU/1-169	<b>GCAGTAGCAGGAGGTAC</b>	-----	-----	-----	-----	<b>ATACAGGGATTGGTTGGAAACAGTGT</b>									110
GB3LKKR02IEAKO/31-199	GCAGTAGCAGGAGGTAC	-----	-----	-----	-----	ATACAGGGATTGGTTGGAAACAGTGT									110
GB3LKKR02FP5VJ/14-182	GCAGTAGCAGGAGGTAC	-----	-----	-----	-----	ATACAGGGATTGGTTGGAAACAGTGT									110
GB3LKKR02F78VO/1-176	<b>ATGCTAAACAGAATTGCAGTAAGCGGTGGA</b>	<b>ACTTATAGAGATTGGTTAGAAACAGTAT</b>													125



	Y	T	A	G	K	Y	L	D	R	P	E	T	P	V	F	I	G	G	M	T	
GB3LKKR01EWCHU/1-169	<b>ACACAGCAGGAAAATACCTTGACAGGCCCGAAAACACCTGTATTTATCGGAGGTATGACAC</b>																				50
GB3LKKR02IEAKO/31-199	ACACAGCAGGAAAATACCTTGACAGGCCCGAAAACACCTGTATTTATCGGAGGTATGACAC																				50
GB3LKKR02FP5VJ/14-182	ACACAGCAGGAAAATACCTTGACAGGCCCGAAAACACCTGTATTTATCGGAGGTATGACAC																				50
GB3LKKR02F78VO/1-176	<b>ATACGGCAGGAAAATACCTTGACAGACCGGAAAACA</b>	<b>CCGTATTTATCGGCGGTATGACCC</b>																			65



	Q	Y	I	E	F	D	E	V	I	S	K	S	A	T	E	S	L	O	H	A	
GB3LKKR01EWCHU/1-169	<b>AATATATCGAATTCGACGAAGTGATTTCAAAAAGTGCGACAGAAAGTTTACAACATGCTA</b>																				5
GB3LKKR02IEAKO/31-199	AATATATCGAATTCGACGAAGTGATTTCAAAAAGTGCGACAGAAAGTTTACAACATGCTA																				5
GB3LKKR02FP5VJ/14-182	AATATATCGAATTCGACGAAGTGATTTCAAAAAGTGCGACAGAAAGTTTACAACATGCTA																				5
GB3LKKR02F78VO/1-176	AATATATC <b>GAGTTCGATGAAGTA</b> ATTTCAAAA <b>TCAGCA</b> ACAGAA <b>CAATATACGGC</b> ----																				5



GB3LKKR01EWCHU/1-169	K	
GB3LKKR02IEAKO/31-199	<b>AA</b>	1
GB3LKKR02FP5VJ/14-182	AA	1
GB3LKKR02F78VO/1-176	--	1



cluster211b

GB3LKKR02FZEFS/108-325  
 GB3LKKR02I4GU4/159-256  
 GB3LKKR01A3RF7/9-215  
 GB3LKKR01EE5WR/1-203  
 GB3LKKR02JMT8K/92-304  
 contig03712/344-517  
 GB3LKKR01EP7G6/106-331

\* V N N G R T Q L A  
**TAAGTTAACAATGGACGCACTCAACTTGCAA** 60  
**TAA**GTTAACAATGGACGCACTCAACTTGCAA 60  
**TAA**ATTAACAATGGA**TGC**ACTCAACTTGCAG 60  
 GAAATTGACAATGGACGCGCTGAACCTGCAA 44  
 CAGTCTATCAATGGACGCACTCAACTTAGCA 54  
 AAAATTAACAATGGATGCGCTAAACCTGCAA 60  
 AAGAATT-AAGTGAACAAAGAAAGTTTAGA 59



GB3LKKR02FZEFS/108-325  
 GB3LKKR02I4GU4/159-256  
 GB3LKKR01A3RF7/9-215  
 GB3LKKR01EE5WR/1-203  
 GB3LKKR02JMT8K/92-304  
 contig03712/344-517  
 GB3LKKR01EP7G6/106-331

Frame +3 p = 3.4e-14  
 T K V Y N M L N R I A V S G  
**CAA-AAGTTTACAACATGTTAAACAGGATCGCA**-----**GTAAGTGGT** 101  
 CAAAAAGTT-TCAACATGTTAAACAGGATCGCA-----GTAAGT--- 98  
 CA**AAAA**GTTTACAACATG**CTAAACAGAATT**GC-----GTA**AGC**GGT 102  
 CAGAAAAGTTTACAACATG**CTAAACAGGATCGCA**-----GTA**TCAGGA** 86  
 CAAAAAGTTTACAACATGTTAAAC**AGAATT**GC-----GTAAGT**GGC** 96  
 CAGAAAAGTTTACAACATG**CTAAACAGAATT**GC-----GTA**GCAGGA** 102  
 AAGGAA**TACTACTACATGTTCTATGAATTATAT****ACTAAGAGAATGTCAGCA-GAAGC**GGT 118



GB3LKKR02FZEFS/108-325  
 GB3LKKR02I4GU4/159-256  
 GB3LKKR01A3RF7/9-215  
 GB3LKKR01EE5WR/1-203  
 GB3LKKR02JMT8K/92-304  
 contig03712/344-517  
 GB3LKKR01EP7G6/106-331

G T Y R D W L E T V Y T A G K Y L D R P  
**GGAACATATAGGGATTGGCTGGAAACAGTATATACGGCGGGAAAGTACCTTGACAGGCCG** 161  
 ----- 98  
 GGA**ACT**TAT**AGA**GATTGG**TTA**GAAACAGTATATAC**GCAGGA**AAA**TACCTTGACAGAL-CC** 161  
 G**GTACCTACCGAGACT**TGGCTGGAAACAG**TG**TAT**ACAGCAGGA**AAA**TATCTTGACAGACCC** 146  
 GGAACAT**ACAGAGACT**GG**TTA**GAAACAGT**ATTCACAGGC**GGAGAATACATGGAA**AGATGT** 156  
 G**GTACATACAGGGATTGGTTG**GAAACAG**TGTACACAGCAGGA**AAA**TACCTTGACAGGCC** 162  
 GGA**ACT**TAT**AGAGATTGGTTA**GAAACAGTATATAC**GCAGGA**AAA**TACC**TTGAC**AGACCG** 178



GB3LKKR02FZEFS/108-325  
 GB3LKKR02I4GU4/159-256  
 GB3LKKR01A3RF7/9-215  
 GB3LKKR01EE5WR/1-203  
 GB3LKKR02JMT8K/92-304  
 contig03712/344-517  
 GB3LKKR01EP7G6/106-331

E T P V F I G G M T O Y I E F D E V I  
**GAAACACCGGTATTTATCGGCGGCATGACTCAATACATTGAATTCGATGAAGTGATC** 218  
 ----- 98  
 GAAACACCGGTA-TTATCGGGG-TATGAC-CAAT**TATATCGAG**TTTCGATG----- 207  
 GAAACA**CCTGTG**TTTATC**GGAGGT**ATG**ACACAA****TATATC**GAATTC**GAC**GAAGTGATC 203  
 GAAACA**CCAATATTCGAA**GG**TGGAACA**AGC**CAAGAA**TT**GTATTT**CAAGA**AGTA**ATC 213  
 GAAACA**CCT**GTA----- 174  
 GAAACACCGGTATTTATCGGC**GGT**ATG**ACC**CAAT**TATATCGAG**TT**CGCT**----- 226

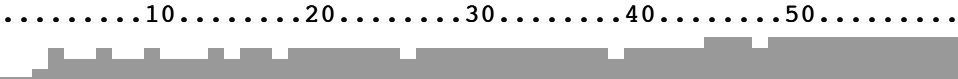


# cluster241a

Frame -3 p = 0.002

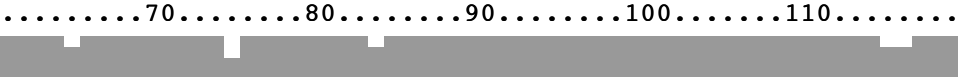
H Y Q C E E P G Y I M G L M A I T P M

GB3LKKR02IEC5V/306-457 **CACTACCAATGCGAAGAGCCGGGATACATTATGGGGCTAATGGCAATCACACCGATGA** 95  
 GB3LKKR01EIO7Z/1-110 CACTACCAAT**TGTGAAGAA**CCGGGCTACATTATGG**GATTA**ATGGCAATCACACCGATGA 95  
 GB3LKKR02HOXTL/22-174 CACTACCAATGCGAAGAGCCGGGATACATTATGGGGCTAATGGCAATCACACCGATGA 95  
 GB3LKKR01DJQRC/358-498 -----**--CGAAGAGCCGGGATACATTATGGGGCTAATGGCAATCACACCGATGA** 82  
 GB3LKKR01CUDIG/48-200 AATTAC**CTA**TGCGAAGAGCCGGGATACATTATGGGGCTAATGGCAATCACACCGATGA 95



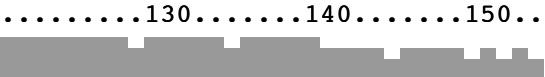
I D Y S Q G N D F D L N L Q T I D D L H

GB3LKKR02IEC5V/306-457 **TTGATTATTCACAAGGAAACGACTTTGATCTAAACTTGCAAACAATAGATGATCTGCATA** 35  
 GB3LKKR01EIO7Z/1-110 TTGATTAT**TCG**CAAGGAAACGACTTTGAT**TTAA**ACTTGCAAACAATGGAT----- 35  
 GB3LKKR02HOXTL/22-174 TTGATTATTCACAAGG**AAAT**GACTTTGAT**TTAA**ACTTGCAAACAATGGATGAT**CTC**CATA 35  
 GB3LKKR01DJQRC/358-498 TTGATTATTCACAAGG**AAAT**GACTTTGAT**TTAA**ACTTGCAAACAATGGATGAT**CTC**CATA 22  
 GB3LKKR01CUDIG/48-200 TTGATTATTCACAAGG**AAAT**GACTTTGAT**TTAA**ACTTGCAAACAATGGATGAT**CTC**CATA 35



K P A L D G I G Y Q D

GB3LKKR02IEC5V/306-457 **AACCGGCATTAGATGGTATCGGCTACCAGGAT** 1  
 GB3LKKR01EIO7Z/1-110 ----- 1  
 GB3LKKR02HOXTL/22-174 AACCGGCA**CTAGATGGAATTGAAGATTA**ACAG 1  
 GB3LKKR01DJQRC/358-498 AACCGGCA**CTAGATGGAATTGGATACCAAGAC** 1  
 GB3LKKR01CUDIG/48-200 AACCGGCA**CTAGATGGAATTGGATACCAAGAC** 1

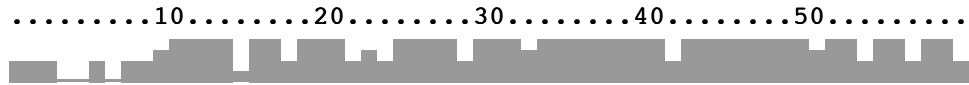


# cluster241b

Frame +1 p = 5.5e-06

G G K P M N S G H I H Y Q C E E P

GB3LKKR01C994Z/1-147 -----GGAGGGAAACCAATGAACAGCGGACATATACACTACCAATGCGAAGAGCCG 51  
 GB3LKKR01AVTUK/1-175 ATCGGAGGAGGAGG-AAACCAACAACAGCGGACATATACACTACCAATGTTGAAGAACCG 59  
 GB3LKKR01E2N9H/1-177 ATCCTACGAGGAGGCAAGCCACTAAACAACGGACATATACATTACCAATGCGAGGAACCA 60  
 GB3LKKR01B1HN4/10-186 ATCGGACGAGGAGGCAAGCCACTAAACAACGGACATATACATTACCAATGCGAGGAACCA 60  
 GB3LKKR02J46T7/1-166 -----GAGGGAAACCAATGAACAGCGGGCATATACACTACCAATGCGAAGAGCCG 50



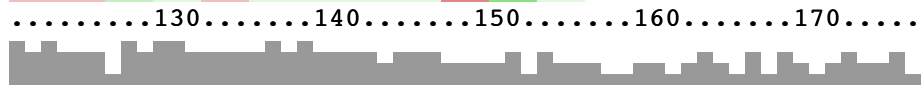
G Y I M G L M A I T P M I D Y S Q G N D

GB3LKKR01C994Z/1-147 GGATACATTATGGGGCTAATGGCAATCACACCGATGATTGATTATTCACAAGGAAATGAC 111  
 GB3LKKR01AVTUK/1-175 GGCTACATTATGGGATTAATGGCAATCACACCGATGATTGATTATTCGCAAGGA-ACGAC 118  
 GB3LKKR01E2N9H/1-177 GGATACATCATGGGACTGATGGCTATCACACCGATGATCGACTATTCGCAAGGAAACGAC 120  
 GB3LKKR01B1HN4/10-186 GGATACATCATGGGACTGATGGCTATCACACCGATGATCGACTATTCGCAAGGCAATGAC 120  
 GB3LKKR02J46T7/1-166 GGATACATTATGGGGCTAATGGCAACACACAGGAAAACAAATTCATACATGATCAGT-AA 109



F D L N L Q T M D D L Q

GB3LKKR01C994Z/1-147 TTTGATTTAAACTTGCAAACAATGGATGATCTCCAA 147  
 GB3LKKR01AVTUK/1-175 TTTGATTTAAACTTGCAAACAATGGATATCTGCAGA 175  
 GB3LKKR01E2N9H/1-177 TTTGATCTAAACTTGCAAACAATAGATGATCTGCAT 177  
 GB3LKKR01B1HN4/10-186 TTTGATCTAAACTTGCAAACAATAGATGATCTACAC 177  
 GB3LKKR02J46T7/1-166 TGTAGGATCAAACCAGACAATCAAAGCCGGTTTCGA 166



cluster258

Frame +3 p= 0.015

K V P F **Y R K K S F W T M L I S I L T**  
**A-AGGTACCGTTTTATCGTAAGAAGAGCTTTTGGACAATGTTGATTTCATTCTTACG** 59  
**A-AAGTACCGTTTTATCGTAAGAAGAGTTTTTGGACATTGTTGATTTCATTCTTACG** 59  
**A--GGTACCGTTTTATCGTAA-AAGAGTTTTTGGACAATGTTGATTTCATTCTTACG** 57  
**-----TATCGTAAGAAGAGTTTTTGGACATTGTTGATTTCATTCTTACG** 45  
**ACTCGTACCGTTTTATCGTAAGAAGAG-TTTTGGACATTGTTGATTTCATTCTTACG** 59  
**AGAAGCTTCGGATTACGTTTCGTTTGTGTTTGAACAGTATCCC-----** 45  
**AGAAGCTTCGGATTACGTTTCGTTTGTGTTTGAACAGTATCCC-----** 45



**A L S V Y F A S S C S R K L V Y R S S G**  
**GCTTTATCTGTTTATTTTGTCTTTCATGTAGTCGGAAGTTGGTTTATCGTTCTTCCGGT** 119  
**GCATTATCTGTTTATTTTGCAGCTTTCATGTAATCGTAAGATTGTTTATCGTTCTTCCGGT** 119  
**GCTTTATCTGTTTATTTTGTCTGCTTTCATGTACTCGGAAGTTGGTTTATCGTTCTTCCGGT** 117  
**GCATTATCTGTTTATTTTGCAGCTTTCATGTAATCGTAAGATTGTTTATCGTTCTTCCGGT** 105  
**GCATTATCTGTTTATTTTGCAGCTTTCATGTAATCGTAAGATTGTTTATCGTTCTTCCGGT** 119  
**-----GGTATTTATTTTGTCTGCTTTCATGTACTCGGAAGTTGGTTTATCGTTCTTCCGGT** 99  
**-----GGTATTTATTTTGTCTGCTTTCATGTACTCGGAAGTTGGTTTATCGTTCTTCCGGT** 99



**V H C D T V Q L D V R S N L K L P \* Y A**  
**GTTTCATTGTGATACGGTTCAACTTGATGTTCCGGTCTAATTTAAAACCTTCCTTAATATGCC** 179  
**GTTTCATTGTGACACGGTTCAACTTGATGTTCCGGTCTAATTTAAAACCTTCCTTAATATGCC** 179  
**GTTTCATTGTGATACTGTTCAAGTTGATGTTCCGGTCTAATTTAAAGCTTCCTTAATATGCC** 177  
**GTTTCATTGTGACACGGTTCTAATTTGATGTTCCGGTCTAATTTAAAACCTTCCTTAATATGCC** 165  
**GTTTCATTGTGACACGGTTCTACTTGTGATGTTCCGGTCTAATTTAAAACCTTCCTTAATATGCC** 179  
**GTTTCATTGTGATACTGTTCAAGTTGATGTTCCGGTCTAATTTAAAGCTTCCTTAATATGCC** 159  
**GTTTCATTGTGATACTGTTCAAGTTGATGTTCCGGTCTAATTTAAAGCTTCCTTAATATGCC** 159



**W C \***  
**TGGTGCTAG** 204  
**TGGTGCTAG** 204  
**TGCTGCTGG** 202  
**TGCTGCTGG** 190  
**TGCTGCTGG** 204  
**TGCTGCTGG** 184  
**TGCTGCTGG** 184



contig07739/468-671  
 GB3LKKR02GPGDJ/44-247  
 GB3LKKR01B1BAU/42-243  
 GB3LKKR01CQ42S/1-190  
 GB3LKKR02GPY0A/70-273  
 GB3LKKR02HCE5B/72-255  
 GB3LKKR01CZX7K/74-257

contig07739/468-671  
 GB3LKKR02GPGDJ/44-247  
 GB3LKKR01B1BAU/42-243  
 GB3LKKR01CQ42S/1-190  
 GB3LKKR02GPY0A/70-273  
 GB3LKKR02HCE5B/72-255  
 GB3LKKR01CZX7K/74-257

contig07739/468-671  
 GB3LKKR02GPGDJ/44-247  
 GB3LKKR01B1BAU/42-243  
 GB3LKKR01CQ42S/1-190  
 GB3LKKR02GPY0A/70-273  
 GB3LKKR02HCE5B/72-255  
 GB3LKKR01CZX7K/74-257

cluster297a

GB3LKKR02JI9WH/42-359  
 GB3LKKR02J4V1J/161-340  
 GB3LKKR02IA3U3/37-354  
 GB3LKKR02JGGMM/177-494  
 GB3LKKR01AFJJ4/82-399

V K D L G F V W I A S I T D K R L R Y V  
**GTTAAAGATTTAGGTTTTGTTTGGATTGCATCTATTACGGACAAGCGTCTTCGGTATGTT** 60  
 GTT**AAAG**GATTTAGGTTTTGTTTGGATTGCATCTATTACGGACAAGCGTCTTCGGTATGTT 60  
 GTT**AAAG**GATTTAGGTTTTGTTTGGATTGCATCTATTACGGACAAGCGTCTTCGGTATGTT 60  
 GTT**AAAG**GATTTAGGTTTTGTTTGGATTGCATCTATTACGGACAAGCGTCTTCGGTATGTT 60  
 GTT**AAAG**GATTTAGGTTTTGTTTGGATT**TCATCTATTACGGACAAGCGTCTTCGGTATGTT** 60  
 .....10.....20.....30.....40.....50.....

V K Y V G K S V Y M D E R S A D F A K S  
**GTCAAGTATGTTGGCAAGTCTGTTTATATGGATGAGCGTTCGCGAGATTTTGCGAAGTCT** 120  
 GTCAAGTATGTTGGCAAGTCTGTTTATATGGATGAGCGTTCGCA**GCT**TTTGGCAAGTCT 120  
 GTCAAGTATGTTGGCAAGTCTGTTTATATGGATGAGCGTTCGCA**GCT**TTTGGCAAGTCT 120  
 GTCAAGTATGTTGGCAAGTCTGTTTATATGGATGAGCGTTCGCA**GCT**TTTGGCAAGTCT 120  
**GTGAAGTATGTTGGCAAGTCTGTTTATATGGATGAGCGTTCGTCGGT**TTTGGCAAGTCT 120  
 .....70.....80.....90.....100.....110.....

GB3LKKR02JI9WH/42-359  
 GB3LKKR02J4V1J/161-340  
 GB3LKKR02IA3U3/37-354  
 GB3LKKR02JGGMM/177-494  
 GB3LKKR01AFJJ4/82-399

L P I T V G K L K T N L Y D F L Q N S R  
**CTTCCTATTACTGTAGGTAATAAAAACTAATCTTTATGACTTTCCTCAGAATAGTAGA** 180  
 CTTCCTATTACTGTAGGTAATA**TT**AAAACTAATCTTTATGACTTTCCTCAGAATAGTAGA 179  
 CTTCCTATTACTGTAGGTAATAAAAACTAATCTTTATGACTTTCCTCAGAATAGTAGA 180  
 CTTCCTATTACTGTAGGTAATAAAAACTAATCTTTATGACTTTCCTCAGAATAGTAGA 180  
 CTTCCTATTACTGTAGGTAATAAAAACTAATCTTTATGACTTTCCTCAGAATAGT**CGA** 180  
 .....130.....140.....150.....160.....170.....

GB3LKKR02JI9WH/42-359  
 GB3LKKR02J4V1J/161-340  
 GB3LKKR02IA3U3/37-354  
 GB3LKKR02JGGMM/177-494  
 GB3LKKR01AFJJ4/82-399

Frame +1 p = 0.001  
 Y R R K F I S A G V G D Y F G D F K A P  
**TATCGCCGTAAATTTATTTTCGGCAGGTGTTGGCGATTATTTTGGGGATTTTAAAGCTCCC** 240  
 T----- 180  
 TAT**CGT**CGTAAATTTGTT**TCACCA**GGTGT**GGT**GATTATTTT**GGC**GATTTTGAAGCTCCC 240  
 TAT**CGT**CGTAAATTTGTT**TCACCA**GGTGT**GGT**GATTATTTT**GGC**GATTTTGAAGCTCCC 240  
 TATCGCCGTAAATTTATTTTCGGCAGGTGTTGGCGATTATTT**GGGA**GATTTTAAAGCTCCC 240  
 .....190.....200.....210.....220.....230.....

GB3LKKR02JI9WH/42-359  
 GB3LKKR02J4V1J/161-340  
 GB3LKKR02IA3U3/37-354  
 GB3LKKR02JGGMM/177-494  
 GB3LKKR01AFJJ4/82-399

G V T S G L W S Y T D H Q T G V V Y R Y  
**GGTGTACTTCTGGTCTTTGGTCTTACACAGATCATCAGACTGGCGTGTGTTTATCGGTAC** 300  
 ----- 180  
 GGTGTTTCTTCTGGTCTTTGGTCTTACACAGATTATAAG**ACCGGT**GTTGTTTAT**CGT**TAC 300  
 GGTGTTTCTTCTGGTCTTTGGTCTTACACAGATTATAAG**ACCGGT**GTTGTTTAT**CGT**TAC 300  
 AGTGTGCTTCTGGTCTTTGGTCTTACACAGATCATAAG**ACCGGT**GCTGTTTATCGGTAC 300  
 .....250.....260.....270.....280.....290.....

GB3LKKR02JI9WH/42-359  
 GB3LKKR02J4V1J/161-340  
 GB3LKKR02IA3U3/37-354  
 GB3LKKR02JGGMM/177-494  
 GB3LKKR01AFJJ4/82-399

R I P R Y Y  
**CGCATCCCTCGCTATTAC** 318  
 ----- 180  
**CGTATCCCTCGCTACTAC** 318  
**CGTATCCCTCGCTACTAC** 318  
 CGCATCCCTCGCT**ACTAC** 318  
 .....310.....

cluster297b

GB3LKKR02FSPL3/170-455
GB3LKKR01A8VS6/1-192
GB3LKKR01BS4WO/257-467
GB3LKKR01BM7GD/221-500
GB3LKKR01AU64A/59-342

S C S Y N A L R E A V K D L G F V W I A
TCTTGTTCTTATAATGCTCTTCGTGAAG-CTGTTAAGGATTTAGGTTTTGTTTGGATTGC 227
TCTTGTTCTTATAATGCTCTTCGTGAAG-CTGTTAAGGATTTAGGTTTTGTTTGGATTGC 193
-----
-CTTGTTCTTATA-TGC-TTTCGTGAAG-C-GTTAAG---TTAGGTTTTGTTTGGATTGC 221
TCTTGTTCTTATAATGTTATTCGTGAAAGTCTGTTAAGGATTTAGGTTTTGTTTGGATTTC 225
.....10.....20.....30.....40.....50.....



GB3LKKR02FSPL3/170-455
GB3LKKR01A8VS6/1-192
GB3LKKR01BS4WO/257-467
GB3LKKR01BM7GD/221-500
GB3LKKR01AU64A/59-342

S I T D K R L R Y V V K Y V G K S V Y M
ATCTATTACGGACAAGCGTCTTCGGTATGTTGTCAAGTATGTTGGCAAGTCTGTTTATAT 167
ATCTATTACGGACAAGCGTCTTCGGTATGTTGTCAAGTATGTTGGCAAGTCTGTTTATAT 167
-----
---CGTCTTCGGTATGTTGTCAAGTATGTTGGCAAGTCTGTTTATAT 93
ATCTATTACGGACAAGCGTCTTCGGTATGTTGTGAAGTATGTTGGCAAGTCTGTTTATAT 161
ATCTATTACGGACAAGCGTCTTCGGTATGTTGTGAAGTATGTTGGCAAGTCTGTTTATAT 167
.....70.....80.....90.....100.....110.....



GB3LKKR02FSPL3/170-455
GB3LKKR01A8VS6/1-192
GB3LKKR01BS4WO/257-467
GB3LKKR01BM7GD/221-500
GB3LKKR01AU64A/59-342

D E R S A A F A K S L P I T V G K L K
GGATGAGCGTTCTG-CAGCTTTTGGCAAGTCTCTTCCATTACTGTAGGTAAATTAATAA 108
GGATGAGCGTTCTG-CAGCTTTTGGCAAGTCTCTTCCATTACTGTAGGTAAATTAATAA 108
GGATGAGTGTCTG-CAGCTTTTGGCAAGTCTCTTCCATTACTGTAGGTAAATTAATAA 34
GGATGAGCGTTCTGTTCTGATTTTGGCAAGTCTCTTCCATTACTGTAGGTAAATTAATAA 101
GGATGAGCGTTCTG-C-GTTTTGGCAAGTCTCTTCCATTACTGTAGGTAAATTAATAA 109
.....130.....140.....150.....160.....170.....



Frame-1 p = 0.006

GB3LKKR02FSPL3/170-455
GB3LKKR01A8VS6/1-192
GB3LKKR01BS4WO/257-467
GB3LKKR01BM7GD/221-500
GB3LKKR01AU64A/59-342

T N L Y D F L Q N S R Y R R K F V S P G
CTAATCTTTATGACTTTCTTCAGAATAGTAGATATCGTCGTAATTTGTTTCACCAGGTG 48
CTAATCTTTATGACTTTCTTCAGAATAGTAGATATCTTCGTAATTTGTTTCACCAGGTG 48
CTAATCTTTATGACTTTCTTCAGAATAGTAGATATCTTCGTAATTTGTTTCACCAGGTG 1
CTAATCTTTATGACTTTCTTCAGAATAGTAGATATCGATATCGCCGTAATTTATTTCCGGCAGGTG 41
CTAATCTTTATGACTTTCTTCAGAATAGTAGATGTCGATGTCGCCGTAATTTATTTCCGGCAGGTG 49
.....190.....200.....210.....220.....230.....



GB3LKKR02FSPL3/170-455
GB3LKKR01A8VS6/1-192
GB3LKKR01BS4WO/257-467
GB3LKKR01BM7GD/221-500
GB3LKKR01AU64A/59-342

V G D Y F G D F E A P G V S S G
TTGGTGATTATTTTGGCGATTTTGAAGCTCCCGGTGTTTCTTCTGGT 1
TTGGTGATTATTTTGGCGATTTTGAAGCTCCCGGTGTTTCTTCTGGT 1
TTGGTGATTATTTTGGCGATTTTGAAGCTCCCGGTGTTTCTTCTGGT 1
TTGGCGATTATTTGGGAGATTTTAAAGCTCCCGGTGTTTACTTCTGGT 1
TTGGCGATTATTTGGGAGATTTTAAAGCTCCAGTGTTCCTTCTGGT 1
.....250.....260.....270.....280.....





GB3LKKR01EYJBS/82-414  
 GB3LKKR02G6OMM/1-335  
 contig09464/1-118  
 GB3LKKR01E1LJ5/1-175  
 GB3LKKR01E2CDT/1-134  
 contig16546/1-333  
 GB3LKKR01CSEX1/105-316  
 GB3LKKR01D984D/221-413

V S A T I P T S S G G A P S G H G P V A  
**GTTTCTGCTACTATTCCCACTTCCAGCGGTGGTGCCTCCGCTCGGACATGGTCCAGTTGCA** 277  
 GTTTCTGCTACTACT**ACT**CCCACTTCCAGCGGTGGTGCCTCCGCTCGGACATGGTCCAGTTGCA 277  
 GTTTCTGCTACTATTCCCACTTCCAGCGGTGGTGCCTCCGCTCGGACATGGTCCAGTTGCA 119  
 GTTTCTGCTACTATTCCCACTTCCAGCGGTGGTGCCTCCGCTCGGACATGGTCCAGTTGCA 176  
 -----GGATTGATTGAAGTGGGGCGGGTGGTGCCTCCG**TCAGGT**CATGGT**CCGGTTGCT** 135  
 GTTTCTGCTACTATA**CCCACT**TAGCGGT**GGGGCCCA**TC**CGGT**CATGGT**CCTGTGGCT** 277  
 ----- 157  
 ----- 137

.....10.....20.....30.....40.....50.....

GB3LKKR01EYJBS/82-414  
 GB3LKKR02G6OMM/1-335  
 contig09464/1-118  
 GB3LKKR01E1LJ5/1-175  
 GB3LKKR01E2CDT/1-134  
 contig16546/1-333  
 GB3LKKR01CSEX1/105-316  
 GB3LKKR01D984D/221-413

S G S F G G L A A L A G N P S A Y A D I  
**TCTGGTTCCTTTGGTGGTCTTGCTGCTCTTGCTGGTAATCCTTCTGCGTATGCCGATATT** 217  
 TCTGGTTCCTTTGGTGGTCTTGCTGCTCTTGCTGGTAATCCTTCTGCGTATGCCGATATT 217  
 TCTGGTTCCTTTGGTGGTCTTGCTGCTCTTGCTGGTAATCCTTCTGCGTATGCCGAT**A** 119  
 TCTGGTTCCTTTGGTGGTCTTGCTGCTCTTGCTGGTAATCCTTCTGCGTATGCCGATATT 176  
 TCTGGTTC**GGA**GGTGGTCTTGCTGCTCTTG**CC**GGTAATCCTTCTGCGTATGCCGATATT 135  
 TCTGGT**GTCGGT**GGGGTCTTGCTGCTCTTGCTGGTAATCCTTCTGCGTATGCCGATATT 217  
 -----**TT** 97  
 ----- 77

.....70.....80.....90.....100.....110.....

GB3LKKR01EYJBS/82-414  
 GB3LKKR02G6OMM/1-335  
 contig09464/1-118  
 GB3LKKR01E1LJ5/1-175  
 GB3LKKR01E2CDT/1-134  
 contig16546/1-333  
 GB3LKKR01CSEX1/105-316  
 GB3LKKR01D984D/221-413

Q L K D A Q Q D R E R S A A A L N D A E  
**CAGTTGAAAGATGCTCAGCAGGATCGAGAGCGTTCCGGCAGCAGCTCTTAATGATGCTGAA** 157  
 CAGTTGAAAGATGCTCAGCAGGAG**CGT**GAGCGTTC**GGCTGCT**TCT**CTCAATGATGCTGAG** 157  
 ----- 119  
 CAGTTGAAAG**GAC**GCTCAGCAGGAG**CGT**GAGCGTTC**GGCTGCT**TCT**CTCAAT****TCAT**---- 157  
 CAGTTGAAAGATGCTCAGCAG----- 135  
 CAGTTGAAAGATGCTCAGCAGGATCGAGAGCGTTCCGGCAGCAGCTCTTAATGATGCTGAA 157  
 CAGTTG--AAGATGCTCAGCAGGAG**CGT**GAGCGTTC**GGCTGCT**TCTCTTAATGATGCT**GAG** 37  
 -----**AG**GAG**CGT**GAGCGTTC**GGCTGCT**TCTCTTAATGATGCT**GAG** 18

.....130.....140.....150.....160.....170.....

GB3LKKR01EYJBS/82-414  
 GB3LKKR02G6OMM/1-335  
 contig09464/1-118  
 GB3LKKR01E1LJ5/1-175  
 GB3LKKR01E2CDT/1-134  
 contig16546/1-333  
 GB3LKKR01CSEX1/105-316  
 GB3LKKR01D984D/221-413

A E W Y K S Q T L D K D L R E R L M K A  
**GCTGAATGGTATAAGTACACAGACTTTGGATAAAGATTTGCGTGAACGTTTGATGAAGGCC** 97  
 GCTGATTGGTATAGG**TCC**CAGACTTTGGATAAAGATTTGCGTGAACG**CTGATGAAGGCT** 97  
 ----- 97  
 ----- 97  
 ----- 90  
 GCTGAATGGTATAAGTACACAGACTTTGGATAAAGATTTGCGTGAACGTTTGATGAAGGCCG 97  
 GCTGATTGGTATAGG**TCC**CAGACTTTGGATAAAGATTTGCGTGAACG**CTGATGAAGGCT** 1  
 GCTGATTGGTATAGG**TC**-CAGACTTTGGATAAAGATTTGCGTGAACG**CTGATGAAGGCT** 1

.....190.....200.....210.....220.....230.....

GB3LKKR01EYJBS/82-414  
 GB3LKKR02G6OMM/1-335  
 contig09464/1-118  
 GB3LKKR01E1LJ5/1-175  
 GB3LKKR01E2CDT/1-134  
 contig16546/1-333  
 GB3LKKR01CSEX1/105-316  
 GB3LKKR01D984D/221-413

Q A G L A E A G I T E S T S R A S L N T  
**CAGGCAGACTTGCAGAAGCTGGAATTACTGAATCTACATCGCGTGCAAGTTTGAACACT** 37  
 CAG**GCC**GGACTT**GCT**GAAAGCTGGA**ATC**ACTGAATCTGCTTCGCGTGCAAGTTT**GAACTCT** 37  
 ----- 37  
 ----- 37  
 ----- 30  
 CAGGCAGACTTGCAGAAGCTGGAATTACTGAATCTACATCGCGTGCAAGTTTGAACACT 37  
 CAG**GCC**GGACTT**GCT**GAAAGCTGGA**ATC**ACTGAATCTGCTTCGCGTGCAAGTT**TAATGCT** 1  
 CAG**GCC**GGACTT**GCT**GAAAGCTGGA**ATC**ACTGAATCTGCTTCGCGTGCAAGTT**TAATGCT** 1

.....250.....260.....270.....280.....290.....

GB3LKKR01EYJBS/82-414  
 GB3LKKR02G6OMM/1-335  
 contig09464/1-118  
 GB3LKKR01E1LJ5/1-175  
 GB3LKKR01E2CDT/1-134  
 contig16546/1-333  
 GB3LKKR01CSEX1/105-316  
 GB3LKKR01D984D/221-413

A I T L S Y T I D N E  
**GCGATAACTTTGTCT**---**TATACTATTGATAACGAG** 1  
**GCTGCTGTTCTTGGT**CAGTCTGGAGTTGG-TGCGAT 1  
 ----- 1  
 ----- 1  
 ----- 1  
 GCGATAACTTTGTCT---TATACTATTGATAACGAG 1  
 GCGATAACT**TTATCC**---TATTCTATTGATAACGAG 1  
 GCGATAACT**TTATCC**---TATTCTATTGATAACGAG 1

.....310.....320.....330.....

cluster339b

Frame +1 p = 1.1e-06

GB3LKKR01AQQGO/136-455
GB3LKKR01CKF4F/185-387
GB3LKKR01EA4GN/113-379
GB3LKKR02FRAT2/1-232
GB3LKKR01DMUJT/1-233
GB3LKKR01BOPRG/40-233
GB3LKKR02F8KPO/1-236

D D L L E G E L O L L T A R A I Y L K S
GATGATCTGCTTGAAGGAGAGTTGCAGTTGTTGACTGCTCGTGCATTTTATTTGAAGTCT 60
AAACATCTGCTTGAAGGAGAGTTGCAGTTGTTGACTGCTCGTGCATTTTATTTGAAGTCT 60
GATGATCTGCTCGAAGGAGAGTTGCAGTTGTTGACTGCTCGTGCATTTTATTTGAAGTCT 60
----- 0
----- 0
GATGAGCTAATTGAAGCAGAAATTGCAGTTGTTGACCGCTCGGGCGCTTTATTTGAAGTCT 60
---TCAACATAACCGGTACGAGTCTTAAGCTGGTAATACTCTTCCTTAGCCTTACCAAG 56
.....10.....20.....30.....40.....50.....



GB3LKKR01AQQGO/136-455
GB3LKKR01CKF4F/185-387
GB3LKKR01EA4GN/113-379
GB3LKKR02FRAT2/1-232
GB3LKKR01DMUJT/1-233
GB3LKKR01BOPRG/40-233
GB3LKKR02F8KPO/1-236

S A S N O E O L S R V N E L T A D D L
TCT-GCTTCTAATCAGGAGCAG--CTGTCCCGTGTGAATGAATTGACTGCCTGATGATTTG 117
TCT-GCTTCTAATCAGGAGCAG--CTGTCCCGTGTGAATGAATTGACTGCCTGATGATTTG 117
TCT-GCTTCTAATCAGGAGCAG--CTGTCCCGTGTGAATGAATTGACTGCCTGATGATTTG 117
----- 29
----- 29
GTTGAATAAAATTGACGGCAGATGATTTG 29
GTTGAATAAAATTGACGGCAGATGATTTG 29
TCTACTGCTAAATCAGGAGGAGTTTGGACACGTGTGAATAAAATGACGGCAGATGATTTG 120
ATT-AGCCTTAAATCAGGAGCAG--CTGTCCCGTGTGAATGAATTGACTGCCTGATGATTTG 113
.....70.....80.....90.....100.....110.....



GB3LKKR01AQQGO/136-455
GB3LKKR01CKF4F/185-387
GB3LKKR01EA4GN/113-379
GB3LKKR02FRAT2/1-232
GB3LKKR01DMUJT/1-233
GB3LKKR01BOPRG/40-233
GB3LKKR02F8KPO/1-236

E N W F D V N W N T O V E V P I I D E K
GAGAATTGGTTTGTGATGTGAATTGGAATACGCAAGTTGAAGTTCCCTATTATCGACGAGAAA 177
GAGAATTGGTTTGTGATGTGAATTGGAATACGCAAGTTGAAGTTCCCTATTATCGACGAGAAA 177
GAGAAC TGGTTTGTGATGTGAATTGGAATACGCAAGTTGAAGTTCCCTATTATCGACGAGAAA 177
GAGAAC TGGTTTGTGATGTGAATTGGAATACGGAGTTGAGGTTCCCTATTATCAACGAGAAA 88
GAGAAC TGGTTTGTGATGTGAATTGGAATACGGAGTTGAGGTTCCCTATTATCAACGAGAAA 89
GAGAAC TGGTTTGTGATGTGAATTGGAATACGGAGTTGAGGTTCCCTATTATCAACGAGAAA 180
GAGAAT TGGTTTGTGATGTGAATTGGAATACGCAAGTTGAAGTTCCCTATTATCGACGAGAAA 173
.....130.....140.....150.....160.....170.....



GB3LKKR01AQQGO/136-455
GB3LKKR01CKF4F/185-387
GB3LKKR01EA4GN/113-379
GB3LKKR02FRAT2/1-232
GB3LKKR01DMUJT/1-233
GB3LKKR01BOPRG/40-233
GB3LKKR02F8KPO/1-236

G K V E R T V K M T G K E I R R E Y M K
GGAAAGGTTGAGCGTACGGTCAAGATGACCGGTAAGGAAATTCGCAGAGAATATATGAAA 237
GGA-GGTTGAGCGTACGGTCAAGATG----- 203
GGAAAGGTTGAGCGTACGGTCAAGATGACCGGTAAGGAAATTCGCAGAGAATATATGAAA 237
GGAAGGATTGAGCGTACGATTAAGATGACCGGC AAGGAAATTCGCAGAGAATATATGAAA 148
GGAAGGATTGAGCGTACGATTAAGATGACCGGC AAGGAAATTCGCAGAGAATATATGAAA 149
GGAAGGATTGAGCG----- 194
GGAAGGTTGAGCGTACGGTCAAGATGACCGGTAAGGAAATTCGCAGAGAATATATGAAA 233
.....190.....200.....210.....220.....230.....



GB3LKKR01AQQGO/136-455
GB3LKKR01CKF4F/185-387
GB3LKKR01EA4GN/113-379
GB3LKKR02FRAT2/1-232
GB3LKKR01DMUJT/1-233
GB3LKKR01BOPRG/40-233
GB3LKKR02F8KPO/1-236

L D L Q D F Q Y D M Y T N R W A L R T E
CTTGATTTGCAGGATTTTCAATATGATATGTATACTAATCGTTGGGCGCTTCGCACCTGAG 297
----- 203
CTTGATTTGCAGGATTTTTCATCCACACGAT----- 267
CTTGATTTGCAGGATTTTCAATATGATATGTATACTAATCGTTGGGAGCTTCGCTCTGAA 208
CTTGATTTGCAGGATTTTCAATATGATATGTATACTAATCGTTGGGAGCTTCGCTCTGAA 209
----- 194
CTT----- 236
.....250.....260.....270.....280.....290.....



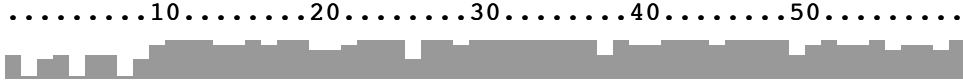
GB3LKKR01AQQGO/136-455
GB3LKKR01CKF4F/185-387
GB3LKKR01EA4GN/113-379
GB3LKKR02FRAT2/1-232
GB3LKKR01DMUJT/1-233
GB3LKKR01BOPRG/40-233
GB3LKKR02F8KPO/1-236

K N R S A I V
AAGAATCGTTTCG-GCTATAGTA 320
----- 203
----- 267
AAGAATCGATTTGGCTATAGTG 232
AAGAATCGATTTGGCTATAGTG 233
----- 194
----- 236
.....310.....320..

cluster348

Frame +3 p= 1.4e-08

GB3LKKR02JPI9I/6-242 **GAGATCCTGGGTCTTAACTCTACTGAGCAAGAAGAAAAAGATAATATTTCTTCGAGCTG** 60  
 contig03211/1-228 -----TGGGTCTTAACTCTACTGAGCAAGAAGAAAAAGATAATATTTCTTCGAGCTG 51  
 GB3LKKR02GT5I1/17-126 **GTTATTGTGGGTCTTAAATCTACC**GAGCAAGAAGAAAAAGATAAT**ATA**CTTCGAG**T**TG 60  
 GB3LKKR01A2JGN/183-393 **GTT**ATCATGGGTCTTAACTCTACTGAGCAAGAAG**GAG**AAAGATAATATTTCTTCGAGCTG 60  
 GB3LKKR02I0PFL/44-279 **TTA**TCATGGGGTCTTAACTCTACTGAGCAAGAAGAAAAAGATAATATTTCTTCGAGCTG 60  
 GB3LKKR02H3Y6W/1-172 -----ATGGGTCTTAACTCTACTGAGCAAGAAG**GAG**AAAGATAATATTTCTTCGAGCTG 52  
 GB3LKKR02FQWKY/1-228 **TTCA**TTATGGGTCTTAA**CAGTACA**GAGCAAGAAG**GAG**AAAGATAAT**ATC**CTT**CGTTCAG** 60  
 GB3LKKR01AYE55/195-429 **GTT**ATCATGGGTCTTAACTCT**ACC**GAGCAAGAAGAAAAAGATAATATTTCTTCGAGCTG 60  
 GB3LKKR01A5HAT/6-242 **GTA**ATTATGGGGT**TAAATGCTACAGAA**CAAGAAGAAAGGGATGAT**ATATTAAGGGCCG** 60  
 contig03850/257-493



GB3LKKR02JPI9I/6-242 **AGATAGTAGTACAAAGGGATGGTCTTCCCTCCGGTAGAGCATTATTCAAATGCGATGTCCG** 120  
 contig03211/1-228 AGATAGTAGTACAAAGGGATGGTCTTCCCTCCGGTAGAGCATTATTCAAATGCGATGTCCG 111  
 GB3LKKR02GT5I1/17-126 AGATGGTAGTACAAAGG**GAC**GGTCTT**TCTTCA**GGTAGA**GCC**TTA**TTTTAAA**----- 110  
 GB3LKKR01A2JGN/183-393 AGATAGTAGTACAAAGGGATGGTCTTCCCTCCGGTAGAGCATTATTCAAATGCGATGTCCG 120  
 GB3LKKR02I0PFL/44-279 AGATAGTAGTACAAAGGGATGGTCTTCCCTCCGGTAGAGCATTATTCAAATGCGATGTCCG 120  
 GB3LKKR02H3Y6W/1-172 -----GTAGTACAAAGG**GAC**GGTCTT**TCTTCA**GGTAGA**GCC**TTA**TTTTAAA****TGT**GAT**GTT**G 55  
 GB3LKKR02FQWKY/1-228 AGATAGTAGTACAAAGGGATGGTCTTCCCTCCGGTAGAGCATTATTCAAATGCGATGTCCG 112  
 GB3LKKR01AYE55/195-429 **AGATT**GTAGTACAAAGGGATGGTCTTCC**TCT**GGTAGA**GCC**TTA**TTTTAGGTTGT**GAT**GTAG** 120  
 GB3LKKR01A5HAT/6-242 AGATAGTAGTACAAAGGGATGGTCTTCCCTCCGGTAGAGCATTATTCAAATGCGATGTCCG 120  
 contig03850/257-493 AGATA**GTGGTT**CAAAGGGATGGT**CTACCATCA**GGTAGA**GCTCTTTTT**TAAATGCTCT**ACAG** 120



GB3LKKR02JPI9I/6-242 **AAAGGCAAAGATGTACGGAAATTTACAAGGGAAACAACGTAAACAATATGATGGAGTGTATG** 180  
 contig03211/1-228 AAAGGCAAAGATGTACGGAAATTTACAAGGGAAACAACGTAAACAATATGATGGAGTGTATG 171  
 GB3LKKR02GT5I1/17-126 ----- 110  
 GB3LKKR01A2JGN/183-393 AAAGGCAAAGATGTACGGAAATTTACAAGGGAAACAACGTAAACAATATGAT**GAA**GTATATG 180  
 GB3LKKR02I0PFL/44-279 AAAGGCAAAGATGTACGGAAATTTACAAGGGAAACAACGTAAACAATATGATGGAGTGTATG 180  
 GB3LKKR02H3Y6W/1-172 **AAAGA**CAAAGATGT**ACA**GAATTTACA**AGAGA**ACAACGTAAACAATATGAT**GAG**GTATATG 115  
 GB3LKKR02FQWKY/1-228 AAAGGCAAAGATGTACGGAAATTTACAAGGGAAACAACGTAAACAATATGAT**GAA**GTATATG 172  
 GB3LKKR01AYE55/195-429 AAAGGCAAAGATGT**ACAGAG**TT**ACT**AAA**GAACAACA**CATATCAA**UAC**CCA**UAC**TAT**CA**-G 179  
 GB3LKKR01A5HAT/6-242 AAAGGCAAAGATGTACGGAAATTTACAAGGGAAACAACGTAAACAATATGAT**GAA**GTGTATT 180  
 contig03850/257-493 **AAACT**CAAC**CGT**GT**ACC**GAATTT**ACT**AGA**GAACAACA**CGA**AAAGA**ATATGAT**GAA**GTGTATT 180



GB3LKKR02JPI9I/6-242 **GTAGTAAGTTGGATGAACAATTTAAAAAGAATACTAACCCGGATGCGGATTCTAAG** 237  
 contig03211/1-228 GTAGTAAGTTGGATGAACAATTTAAAAAGAATACTAACCCGGATGCGGATTCTAAG 228  
 GB3LKKR02GT5I1/17-126 ----- 110  
 GB3LKKR01A2JGN/183-393 GTAGTAAGTTGGATGAACAATTA**ATAACGTA**- 211  
 GB3LKKR02I0PFL/44-279 GTAGTAAGTTGGATGAACAATTTAAAAAGAAAG**ATAATCCA**-**AT****GCT**GAT**TCCA**G 236  
 GB3LKKR02H3Y6W/1-172 **GTA**AAAAATTTGGATGAACAATTT**AAGA**AGGC**ACTAATCCAG**ATGCGGATTCTAAG 172  
 GB3LKKR02FQWKY/1-228 GTAGTAAGTTGGATGAACAATTTAAAAAGAATACTAA-AGGTATCTATACATTTGC 228  
 GB3LKKR01AYE55/195-429 **GTE**-**TAAA**TATAATATGATAAGG**ACTATAA**AC**CA**T**TC**TA**ACC**CTT**ATCA**AG**ATC** 235  
 GB3LKKR01A5HAT/6-242 CTGGAGT**ATTA**GATTTCTATGAT**GAA**GAGT**TCT**AAA**GAT**AA**TCCC**CT**GCA**AATA**AAA** 237  
 contig03850/257-493 CTGGAGT**ATTA**GATTTCTATGAT**GAA**GAGT**TCT**AAA**GAT**AA**TCCC**CT**GCA**AATA**AAA** 237



cluster375a

Frame-1 p= 0.013

GB3LKKR02GFL4E/48-362  
 GB3LKKR01A3N1M/1-255  
 GB3LKKR01B2UZU/181-440  
 GB3LKKR02IIYC9/37-316  
 GB3LKKR02GS4KL/162-448  
 GB3LKKR02GHLLT/28-339

F S I T P G I I Y P V R I O F V N A R D  
**TTTTCTATTACGCCGGGTATTATTTATCCGGTTCGTATTCAGTTTGTCAATGCTCGTGAT** 256  
 TTTTCTATTACGCCGGGTATTATTTATCCG**GTG**CGTATTCAGTTTGTCAATGCTCGTGAT 256  
 TTTTCTATTACGCCGGGTATTATTTATCCG**GTG**CGTATTCAGTTTGTCAATGCTCGTGAT 201  
 TTTTCTATTACGCCGGGTATTATTTATCCG**GTG**CGTATTCAGTTTGTCAATGCTCGTGAT 255  
 TTTTCTATTACGCCGGGTATTATTTATCCG**GTG**CGTATTCAGTTTGTCAATGCTCGTGAT 230  
 TTTTCTGTTACGCCGGGTATTATTT**TACCCG****GTG**CGTATTCAGTTTGTCAATGCTCGTGAT 254  
 .....10.....20.....30.....40.....50.....



GB3LKKR02GFL4E/48-362  
 GB3LKKR01A3N1M/1-255  
 GB3LKKR01B2UZU/181-440  
 GB3LKKR02IIYC9/37-316  
 GB3LKKR02GS4KL/162-448  
 GB3LKKR02GHLLT/28-339

R V T L H Q G I D V R L N P L G V P S F  
**CGGGTTACTTTGCATCAAGGTATTGATGTCGGTTTAAATCCTTTGGGTGTTCCGTCGTTT** 196  
**CGAG**TACTTTGCATCAAGGTATTGATGTCGGT**TCA**AATCCTTTGGGTGTTCCGTCGTTT 196  
**CGAG**TACTTTGCATCAAGGTATTGATGTCGGT**TCA**AATCCTTTGGGTGTTCCGTCGTTT 141  
**CGAG**TACTTTGCATCAAGGTATTGATGTCGGT**TCA**AATCCTTTGGGTGTTCCGTCGTTT 195  
 CGGGT**TACGCTG**CATCAAGGTATTGATGTCGGT**TCA**AAT**CCC**TGGGTGTTCCGTCGTTT 170  
 CGGGT**TACGCTG**CATCAAGGT**T**TGATGTCGGT**TCA**AATCCTTTGGGTGTT**CCA**TCGTTT 194  
 .....70.....80.....90.....100.....110.....



GB3LKKR02GFL4E/48-362  
 GB3LKKR01A3N1M/1-255  
 GB3LKKR01B2UZU/181-440  
 GB3LKKR02IIYC9/37-316  
 GB3LKKR02GS4KL/162-448  
 GB3LKKR02GHLLT/28-339

N P Y V L R L H R F W V P L O L Y H P E  
**AATCCTTATGTACTTTCGGTTGCATCGGTTTGGGTTCCCTTTCAGTTGTATCATCCGGAA** 136  
 AATCCTTATGTACTTTCGGTTGCATCGGTTTGGGTTCCCTTTCAGTTGTATCATCATCCGGAA 136  
 AATCCTTATGTACTTTCGGTTGCAT**CG**TTTTGGGTTCCCTTTCAGTTGTATCATCCGGAA 82  
 AATCCTTATGTACTTTCGGTTGCATCGGTTTGG**GTCC**TTCAGCAGT**ACC**GGAA**CG**-**CG** 136  
**AAC**CCTTATGTACTT**CGACTG**CATCGGTTTGGGTTCCCTATGCAGTTGTATCAT**CCT**GAA 110  
**AAC**CCTTATGTACTTTCGG**CTG**CATCGGTTTGGGTTCC**T**-**G**CAGTTGTATCAT**CCT**GAA 135  
 .....130.....140.....150.....160.....170.....



GB3LKKR02GFL4E/48-362  
 GB3LKKR01A3N1M/1-255  
 GB3LKKR01B2UZU/181-440  
 GB3LKKR02IIYC9/37-316  
 GB3LKKR02GS4KL/162-448  
 GB3LKKR02GHLLT/28-339

M R V N S S K F D M N D L T Y N F I P G  
**ATGCGTGTAAATTCGTCTAAGTTCGATATGAATGATTTGACGTATAACTTTATTTCCCGGT** 76  
 ATGCGTGTAAATTCGTCTAAGTTCGATATGAATGATTTGACGTATAACTTTATTTCCCGGT 76  
 ATGCGTGTAAATTCGTCTAAGTTCGATATGAATGATTTGACGTATAACTTTATTTCCCGGT 22  
**GATACA****TAA**AAGAAGTATAAGCACCACC**CTTT**ATTTATCTACAA**CACCGGGAA****TAA**AGT 76  
 AT**CGA**GTTAAATTCGTCTAAG**TTT**GATATGA**AAC**GATTTGACGTTTAACTTTATTTCCCGGT 50  
 AT**CGA**GTTAAATTCGTCTAAG**TTT**GATATGA**AAC**GATTTGACGTTTAACTTTATTTCCCGGT 76  
 .....190.....200.....210.....220.....230.....



GB3LKKR02GFL4E/48-362  
 GB3LKKR01A3N1M/1-255  
 GB3LKKR01B2UZU/181-440  
 GB3LKKR02IIYC9/37-316  
 GB3LKKR02GS4KL/162-448  
 GB3LKKR02GHLLT/28-339

V V D N K G G G A Y T S F M Y P R S G T  
**GTTGTAGATAATAAAGGTGGTGGTGCCTTATACTTCTTTTATGTATCCGCGTTCGGTACT** 16  
 GTTGTAGATAATAAAA----- 16  
 GTTGTAGATAATAAAGGTGGTGGTGCCTTATACTTCTTTTATGTATCCGCGTTCGGTACT 1  
**TAT**ACGTCAAAT**CAT**T**CA**T**TCA**ACT**TAG****ACGA**T**TAA****CA**----- 16  
**TGT**GTAGATAAAT**TCTGGTCCA**ACT**T**CG**TAC**ACTTCTTTTATGTATCCGCGT**CAAC**GATCG 1  
**TGT**GTAGATAATAACGGT**GCTGC**ACC**G**TATACT**T****TC**-TTTATGTATCCGCGT**CCC**GGTGT 16  
 .....250.....260.....270.....280.....290.....

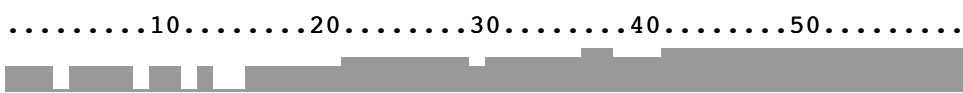


GB3LKKR02GFL4E/48-362  
 GB3LKKR01A3N1M/1-255  
 GB3LKKR01B2UZU/181-440  
 GB3LKKR02IIYC9/37-316  
 GB3LKKR02GS4KL/162-448  
 GB3LKKR02GHLLT/28-339

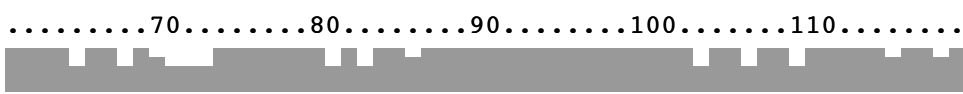
A A F F S  
**GCTGCTTTTTTTCAGT** 1  
 ----- 1  
 GCTGCTTTTTTTCAGT 1  
 ----- 1  
**GGTCTTTTATGTA**-- 1  
 GCT**GCA**TTTTTTCAGT 1  
 .....310...



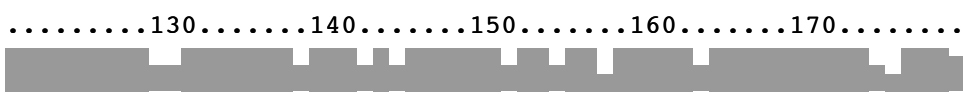
GB3LKKR02JAOTH/146-515 312
GB3LKKR02JRCYP/15-455 383
GB3LKKR02JUV5A/141-499 300
contig01623/1-405 384
GB3LKKR02F42F9/12-452 382
FTSPZO101CX1PY/23-293 233



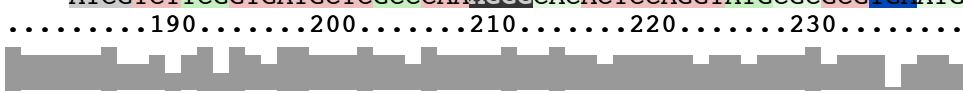
GB3LKKR02JAOTH/146-515 252
GB3LKKR02JRCYP/15-455 323
GB3LKKR02JUV5A/141-499 240
contig01623/1-405 324
GB3LKKR02F42F9/12-452 322
FTSPZO101CX1PY/23-293 174



GB3LKKR02JAOTH/146-515 192
GB3LKKR02JRCYP/15-455 263
GB3LKKR02JUV5A/141-499 180
contig01623/1-405 264
GB3LKKR02F42F9/12-452 263
FTSPZO101CX1PY/23-293 114



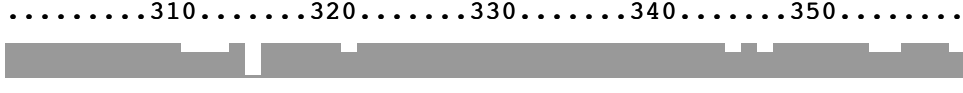
GB3LKKR02JAOTH/146-515
GB3LKKR02JRCYP/15-455
GB3LKKR02JUV5A/141-499
contig01623/1-405
GB3LKKR02F42F9/12-452
FTSPZO101CX1PY/23-293
\* L G Y S R S G F Y E K I O N R S F
TA-ATTAGGATATAGCAGATCGGGG-TTCTATGAGAAGATACAAAACAGGAGTTTT 133
TAAATTAGGATATAGCAGATCGGGG-TTCCATAAGAAGATACAAAACAGGAGTTTT 205
TAAATTAGGATATAGCAGATTCGGGGTTC-TATAAGAAGATACAAAACAGGAGTTTT 120
CAAGTTAGGGTATAGTAGTCCGG-ATTCTATAAAAATAAAGAACAGGAGTTTT 205
TAAATTAGGATATAGCAGAT-CGGGGTTC-TATAAGAAGATACAAAACAGGAGTTTT 204
ATCGTCTTCGGTGATGCTCGCCCAAAGGGCACACTCCAGGTATGCGCGCGTGAATG 55



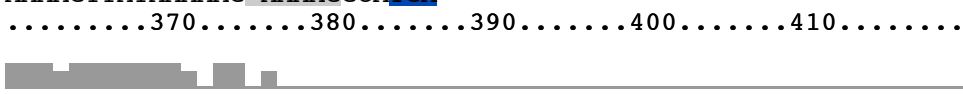
GB3LKKR02JAOTH/146-515
GB3LKKR02JRCYP/15-455
GB3LKKR02JUV5A/141-499
contig01623/1-405
GB3LKKR02F42F9/12-452
FTSPZO101CX1PY/23-293
Frame -2 p = 0.001
N I R E L A Q I F D A I I N F K D Q D W
AATATCCGGGAAGCTAGCTCAGATATTCGATGCGATCATCAACTTCAAAGATCAAGATTGG 75
AATATCCGGGAAGCTAGCTCAGATATTCGATGCGATCATCAACTTCAAAGATCAAGATTGG 145
AATATCCGGGAAGCTAGCTCAGATATTCGACACGATCATCAATTTCAAAGATCAAGATTGG 60
AATATCCGAGAGCTGCTCAGATATTCGATACGATCATCAACTTCAAAGATCAAGATTGG 145
AATATCCGGGAGCT-CTGCAGATAT-CGAGACGATCATCAATTTCAAAGATCAAGATTGG 145
TGTCCAAGTTACCG-GTCCCTTCGTTTCGACACGATCATCAACTTCAAAGATCAAGATTGG 1



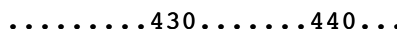
GB3LKKR02JAOTH/146-515
GB3LKKR02JRCYP/15-455
GB3LKKR02JUV5A/141-499
contig01623/1-405
GB3LKKR02F42F9/12-452
FTSPZO101CX1PY/23-293
T E G G K I N R L K R Y R A M S L M E F N
ACTGAGGGTAAGATTAATAGGCTTAAAAGATATAGGGCTATGAGCCTTATGGAGTTCAAT 18
ACTGAGGGTAAGATTAATAGGCTTAAAAGATATAGGGCTATGAGCCTTATGGAGTTCAAT 86
GCGGAGAGTAAGATATAGATAGGCTTAAAAGATATAGAGCCATAAGCCTTATGGAGTTCAAT 2
ACTGAGGGTAAGATTAATAGGCTTAAAAGATATAGGGCTATGAGCCTTATGGAGTTCAAT 86
GCGGAGAGTAAGATATAGATAGGCTTAAAAGATATAGAGCCATAAGCCTTATGGAGTTCAAT 85
ACTGAGGGTAAGATTAATAGGCTTAAAAGGATATAGGGCTATGAGCCTTATGGAGTTCAAT 1



GB3LKKR02JAOTH/146-515
GB3LKKR02JRCYP/15-455
GB3LKKR02JUV5A/141-499
contig01623/1-405
GB3LKKR02F42F9/12-452
FTSPZO101CX1PY/23-293
K S Y K K K K A \*
AAAAGTTATAAAAAGAAAAGGCATGA 1
AAAAGTTATAAAAAGAAAAGGCATGA 26
AAGAATTATAAACGAAAAGAAAGCATGA 1
AAAAGTTATAAAAAGAAAAGGCATGA 26
AAGAATTATAAACGAAAAGAAAGCATGA 25
AAAAGTTATAAAAAG-AAAAGGCATGA 1



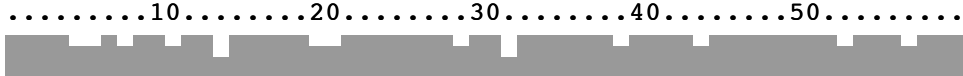
GB3LKKR02JAOTH/146-515 1
GB3LKKR02JRCYP/15-455 1
GB3LKKR02JUV5A/141-499 1
contig01623/1-405 1
GB3LKKR02F42F9/12-452 1
FTSPZO101CX1PY/23-293 1



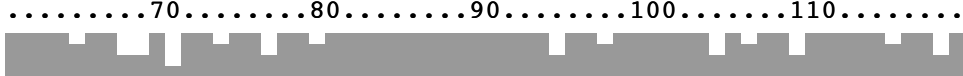
cluster457

Frame +2 p = 1.3e-08

	E	K	G	V	D	I	A	G	Q	Y	S	P	V	V	G	A	G	L	R	
FTSPZO101BPBHB/81-291	<b>GAA</b>	<b>-AAAGGCGTTGACATTGCAGGGCAGTATAGTCCTGTTGTGCGGTGCTGGTTTGC</b>	<b>GA</b>																	59
FS22EC101CA3FG/4-214	GAAGAAAAGGC	<b>G</b> TAGACATTGCAGGG <b>CAATAC</b> AGTCCCTGTTGTGCGGTGCT <b>GGC</b> TTGCGTA																		60
FTSPZO101B9AWK/7-218	GAAGAAAAGGC	<b>G</b> TAGACATTACAGGG <b>CAATAC</b> AGTCCCTGTTGTGCGGTGCT <b>GGC</b> TTGCGTA																		60
FTSPZO101B6C5R/133-266	GAAGAAAAGGC	<b>G</b> TAGACATTGCAGGG <b>CAA</b> TATAGTCCTGTT <b>GTT</b> GGTGCT <b>GGC</b> TTGCGTA																		60
FTSPZO101DMWHB/4-137	GAAGCAG <b>GGT</b> GTTGAC <b>ATA</b> GCAGGG <b>CAA</b> TATAGTCCTATTGTGCGGTGCT <b>GGC</b> TTG <b>TGT</b> A																			60



	T	V	P	A	A	I	G	I	A	M	G	A	K	P	A	L	O	A	G	R	
FTSPZO101BPBHB/81-291	<b>CTGTACCTGCGGCGATTGGTATTGCTATGGGTGCAAAACCTGCGCTACAGGCAGGTCGAC</b>																				119
FS22EC101CA3FG/4-214	CTGTAC <b>CCA</b> CCG <b>GCTAT</b> CGGTATTGCTATGGGT <b>GCT</b> AAACCTGCG <b>TTG</b> CAGGCAGGGCGAC																				120
FTSPZO101B9AWK/7-218	CTGTAC <b>CAATGCTGCTATTGGC</b> ATTGCTATGGGTGCA <b>AAG</b> CC <b>TGCGTTACA</b> AGCAGGGCGGC																				120
FTSPZO101B6C5R/133-266	CTGTAC <b>CCA</b> CCG <b>GCTAT</b> CGGTATTGCTATGGGT <b>GCT</b> AAACCTGCGCTAC <b>CA</b> AGCAGGGCGAC																				120
FTSPZO101DMWHB/4-137	CT <b>GTG</b> CCTACT <b>GCTAT</b> CGGTATTGCTATGGGT <b>GCT</b> AAACCTGCGCTACAGGCAGGGCGGC																				120



	Q	V	S	E	G	L	G	A	I	Q	G	R	M	V	A	N	A	N	A	P	
FTSPZO101BPBHB/81-291	<b>AAGTAAGCGAGGGATTAGGGGCGATACAAGGGCGTATGGTTGCTAATGCTAACGCACCAC</b>																				179
FS22EC101CA3FG/4-214	AA <b>GTT</b> AGCG <b>GA</b> AGGATTAGGG <b>-CT</b> ATGCAAGGGCGTATGATT <b>GCC</b> AATGCTACTGCACCAA																				179
FTSPZO101B9AWK/7-218	AA <b>GTG</b> AGTAAACAG <b>CT</b> AGGGCGCAATGCAATCGCGTATGATT <b>GCA</b> AATGCT <b>AAT</b> GCACCAA																				180
FTSPZO101B6C5R/133-266	AA <b>GTT</b> AGCAAAG <b>CT</b> -----																				134
FTSPZO101DMWHB/4-137	AAGTAGCAAAG <b>GCT</b> -----																				134



	R	T	L	N	T	G	Y	M	G	Q	R	
FTSPZO101BPBHB/81-291	<b>GAACATTAAATACGGGTTATATGGGGCAACGT</b>											211
FS22EC101CA3FG/4-214	<b>GCACGCTAAACACAGGATTTAG</b> GGG <b>CAG</b> CGT											211
FTSPZO101B9AWK/7-218	<b>GGGCGTTAAACACT</b> GGTTATATGGGGCAACGT											212
FTSPZO101B6C5R/133-266	-----											134
FTSPZO101DMWHB/4-137	-----											134



cluster502

Frame +1 p = 3.6e-07

D W L E T V Y T A G K Y L D R P E T P V

GB3LKKR01BXXY4/1-146 **GATTGGTTAGAAACAGTATATACGGCAGGAAAATACCTTGACAGACCGGAAACACCGGTA** 60

contig17153/1-146 GATTGG**TTG**GAAACAGTATATACGGCAGGAA**AAG**TACCTTGACAGACCGGAAACACCGGTA 60

GB3LKKR02FQIJ5/1-127 -----**AT**ACGGCAGGAAAATACCTTGACAGACCGGAAACACCGGTA 41

GB3LKKR02FUJR7/1-146 GATTGG**TTG**GAAACAGTATAT**ACAGCAGGT**AAATACCTTGACAGAC**CCAGAGACA****CCAGTG** 60

.....10.....20.....30.....40.....50.....



F I G G M T Q Y I E F D E V I S K S A T

GB3LKKR01BXXY4/1-146 **TTTATCGGCGGTATGACCCAATATATCGAGTTCGATGAAGTAATTTCAAAATCAGCAACA** 120

contig17153/1-146 **TTCATC****GGT**GGTATGACCCAATAT**ATT**GAGTTCGATGAAGTA**ATC**TCAAAATC**AGCGACA** 120

GB3LKKR02FQIJ5/1-127 TTTATC**GGAG**GTATG**ACA**CAATATATC**GAA**TTC**GAC**GAA**GTGATC**TCAAAA**AGTGCGACA** 101

GB3LKKR02FUJR7/1-146 TTT**ATTGGAG**GTATG**ACACAGTACATTGAA**TTC**GAC**GAA**GTTATC**TCAAAA**AGTGCGACA** 120

.....70.....80.....90.....100.....110.....



E T I Y G S Q P

GB3LKKR01BXXY4/1-146 **GAAACAATATACGGCTCACAACCA** 146

contig17153/1-146 GATACAATAT**TAT**GGCTCACAACCA 146

GB3LKKR02FQIJ5/1-127 GAAACAG**CATAC****GGAAGT**CAACCA 127

GB3LKKR02FUJR7/1-146 GAAACAG**CATAC****GGAAGT**CAACCA 146

.....130.....140.....



cluster532

Frame -1 p = 0.001

V A S S V D A V V E M D S E P E M V E I

contig08608/212-370 **GTCGCTTCTTCGGTGGACGCAGTGGTGGAGATGGATAGCGAGCCGGAGATGGTGGAGATA** 101  
 GB3LKKR01BKK2R/1-158 **GTCGCTTCTTCGGTGGACGCAGTGGTGGAGATGGATAGCGAGCCGGAGATGGTGGAGATA** 101  
 GB3LKKR02HC1UT/118-244 -----GATAGCGAGCCGGAGATGGTGGAGATA 69  
 GB3LKKR01CXJZO/1-159 **GTA**GCTTCTTCGGTGGACGCAGTGGTGGAGATGGATAGCGAG**CCAG**AGATGGTGGAGATA 101  
 GB3LKKR01BBV60/1-159 **GTG**GCT**TCCTCTGCAGATATC**GTGGTGGAGATGGATAGCGAGCCGGAGATGGTGGAGATA 101

.....10.....20.....30.....40.....50.....



E V E T G E F Y K T G A K K G Q P K T G

contig08608/212-370 **GAGGTCGAGACAGGAGAGTTCTATAAGACGGGAGCTAAGAAAGGTCAGCCTAAGACGG-G** 41  
 GB3LKKR01BKK2R/1-158 **GAGGTCGAGACAGGAGAGTTCTATAAGACGGGAGCTAAGAAAGGTCAGCCTAAGACG--G** 41  
 GB3LKKR02HC1UT/118-244 **GAGGTCGAGACAGGAGAGCTCTACAAGACTGGAGGCCAAGAAA**GGCCAA**CCTAAGACAG-G** 9  
 GB3LKKR01CXJZO/1-159 **GAGGTCGAGACAGGAGAGCTCTACAAGACTGGAGGCCAAGAAA**GGCCAA**CCTAAGACAGAG** 41  
 GB3LKKR01BBV60/1-159 **GAGGTCGAGACAGGAGAGCTCTACAAGACTGGAGGCCAAGAAA**GGCCAA**CCTAAGACAGAA** 41

.....70.....80.....90.....100.....110.....



K R R \*

contig08608/212-370 **AAAAAGAAGATAA** 1  
 GB3LKKR01BKK2R/1-158 **AAAAAGAAGATAA** 1  
 GB3LKKR02HC1UT/118-244 **AAAAAGAAGATAA** 1  
 GB3LKKR01CXJZO/1-159 **AAAAAGAAGATAA** 1  
 GB3LKKR01BBV60/1-159 **AAAAAGAAGATAA** 1

.....130.....140.....150.....1





cluster540

FTSPZO101CLSF2/1-245  
 ER8QEOW01CIRYO/1-228  
 DBA-SLE c2415/1-183  
 FC8LRL301AU4EL/48-245  
 is\_serum c2894/36-273  
 cerebrospinal\_rep c831/82-323  
 FS22EC101AWJ2M/34-274

F K P C V L G W O P F L K K S L L L I V  
**TTCAAGCCGTGTGTGCTAGGGTGGCAACCCTTCTTAAAGAAAAGTCTGTTATTGATAGT** 60  
 -----CAACGAGCGAGG-AGT**TGATGAATGAGCAATCAGAAATTGGAAGT** 44  
 ----- 0  
**TACCCGCCGATTACCAACAAGCGAGGAGTTGATGAATGAGCAATCT---GTTGATAGT** 57  
**TACCTGCCGATTACCGACAAGCGAGGAGTTGATGAATGAGCAATCT---ATTGATAGT** 57  
**TACCTGCCGATTACCGACAAGCGAGGAGTTATAGATGAGCAATCT---ATTGATAGT** 57  
**TACCTGCCGATTACCAACAAGCGAGGAGTTGATGAATGAGCAATCA---ATTGATAGT** 57  
 .....10.....20.....30.....40.....50.....



Frame +2 p= 3.8e-04

FTSPZO101CLSF2/1-245  
 ER8QEOW01CIRYO/1-228  
 DBA-SLE c2415/1-183  
 FC8LRL301AU4EL/48-245  
 is\_serum c2894/36-273  
 cerebrospinal\_rep c831/82-323  
 FS22EC101AWJ2M/34-274

E K V S K E G D Y H Q W L S N P V T K A M  
**GGAGAAGGTGAGCAAGGCGGATTACCTGCAATGGTTGAGCAATCCAGTAACCAAGGCAAT** 120  
**GGAGAAGGTGAGCAAGGCGGATTACCTGCAATGGTTGAGCAATCCAGTAACCAAGGCAAT** 104  
**--AGAAGGTGAGTATAGCGGATTACCACCAATGGTTGAGCAATCCAGTAACCAAGGCAAT** 58  
**GGAGAAGGTGAGTATAGCGGATTACCACCAATGGTTGAGCAATCCAGTAACCAAGGCAAT** 117  
**GGAGAAGGTGAGCAAGGCGGATTACCA-CA-CAATGGT-GGACAATCCAGTGACCAAAGCGAT** 113  
**GGAGAAGGTGAGTATAGCGGATTACCACCAATGGTTGAGCAATCCAGTGACCAAAGGCRAT** 117  
**GGAGAAGGTGAGTATAGCGGATTATCACCAATGGTTGAGCAATCCAGTAACCAAGGCGAT** 117  
 .....70.....80.....90.....100.....110.....



FTSPZO101CLSF2/1-245  
 ER8QEOW01CIRYO/1-228  
 DBA-SLE c2415/1-183  
 FC8LRL301AU4EL/48-245  
 is\_serum c2894/36-273  
 cerebrospinal\_rep c831/82-323  
 FS22EC101AWJ2M/34-274

T A Q L T N E L M S L N S L T S L P N P  
**GACGGCGCAATTGACCAACGAGTTGATGAGTCTCAACAGCCTAACGAGCCTACCCAATCC** 180  
**GACGGCGCAATTGACCAACGAGTTGATGAGTCTCAACAGCCTAACGAGCCT-CCCACCC** 163  
**GACAGCTCACTTGACCAACGAATGTTTGGAGCCTGAACCAACTGAGCAGCCTACCGAACCC** 118  
**GACGGCGCAATTGACCAACGAGTTGATGAGTCTCAACAGCCTAACGAGCCTCCCAATCC** 177  
**GACGGCGCAATTGACCAACGAGTTGATGAGTCTCAACAGTTTAAACGAGCCTACCCAATCC** 173  
**GACGGCGCAATTGACCAACGAGTTGATGAGCCTCAACAGCCTAACGAGCCTCCCAATCC** 177  
**GACGGCAATTGACCAACGAGTTGATGAGTCTCAACAGCTTAAACGAGCCTACCCAATCC** 177  
 .....130.....140.....150.....160.....170.....



FTSPZO101CLSF2/1-245  
 ER8QEOW01CIRYO/1-228  
 DBA-SLE c2415/1-183  
 FC8LRL301AU4EL/48-245  
 is\_serum c2894/36-273  
 cerebrospinal\_rep c831/82-323  
 FS22EC101AWJ2M/34-274

L T L E K V H Y Q L G K V Q G L T A V L  
**GTTGACATTGGAGAAGGTGCATTACCAATTGGGCAAGGTGCAGGGATTGACGGCGGTGTT** 240  
**ACAGACATTGGAGAAGGTACATTACCAATTGGGCAAGGTGCAGGGATTGACGGCGGTGCT** 223  
**ATTGACATTGGAGCGGTGAACCTACCAACGGGGCAAGGTGTCGGGCTTGATGGCGGTCTT** 178  
**ATTGACATTGGAGAAGGTACA-** 198  
**ATTGACATTGGAGAAGGTGCATTACCAATTGGGCAAGGTGCAGGGATTGACGGCGGTGCT** 233  
**ATTGGCATTGGAGAAGGTACATTACCAATTGGGCAAGGTGCAGGGATTGACGGCGGTGCT** 237  
**ATTAACATTGGAGAAGGTGCATTACCAATTGGGCAAGGTGTTTCGATCGACTTGTGGCC** 237  
 .....190.....200.....210.....220.....230.....



FTSPZO101CLSF2/1-245  
 ER8QEOW01CIRYO/1-228  
 DBA-SLE c2415/1-183  
 FC8LRL301AU4EL/48-245  
 is\_serum c2894/36-273  
 cerebrospinal\_rep c831/82-323  
 FS22EC101AWJ2M/34-274

D  
**GGAT** 245  
**GGAT** 228  
**GGAT** 183  
**----** 198  
**GGAT** 238  
**GGAT** 242  
**GGG-** 241  
 .....



cluster562

Frame +1 p = 4.6e-07

E A F E Q F D A L R V A N A L G L E Y G  
**GAAGCGTTTTGAACAGTTCGACGCATTCGCGAGTTGCTAATGCACTCGGTCTTGAATATGGC** 60  
 GAAG**GC**ATTTGAACAGTTCGACGC**CACTACGCGTGGCTAAC**GCACCTCGGTCTTGAATATGGC 60  
 -----  
 GAAG**GC**ATTTGAACAGTTAGAC**GCGTTGCGCGTGGCA**AATGCACCT**GGCCT**AGAATAT**GAA** 60  
 GAAG**GC**ATTTGAACAGTTCGACGC**CACTACGCGTGGCMAAC**GCACCTCGGTCTTGAATATGGC 60  
 -----  
 -----**CACTT**GGTCTTGAATATGGC 20  
 -----  
 -----  
 -----  
 .....10.....20.....30.....40.....50.....



V V C K W R D R D Q I P A Y W R V K L  
**GTTGTTTG-CAAA-TGGCGTGACCGTGATCAAATTCCTGCTTACTGGCGGTTAAACTTG** 118  
 GTTGTTTG-CAAA-TGGCGTGACCGTGATCAAATTCCTGCTTACTGGCGGTTAAATTTG 118  
**-TT**GTTTG-CAAA-TGGCGT**GAT**CGTGAGCAAATTCCTGCTTACTGGCGGTTAAATTTG 57  
 ACTGTTTG-CAAA-TGGCGTGACCGTGATCAAATTCCTGCTTACTGGCGGTTAAATTTG 118  
 GTTGTTTG-CAAA-TGGCGTGACCGTGAGCAAATTCCT**GC**ATACTGGCGCACAGCGTTTG 118  
 GTTGTT**TGTC**AAA**TAGG**CGTGACCGTGAGTCGAT**CCAGCTTAT**TGGCGGTTAAATTT**█** 79  
 -----  
 -----GAGCTCT**TGCAGATAT-CC**CTGTG**CAAG**TTTG 30  
 -----  
 -----  
 .....70.....80.....90.....100.....110.....



V N L M N H H N V A I S L H D L A G W I  
**TTAATTTGATGAATCATCATAACGTAGCGATTTCTTTGCACGACTTAGCAGGATGGATT** 177  
 TTAATTTGATGAATCAT**CACGGCGTT**CTATTTCT**TTA**CAC**GATTTG**GCAGGGTGGATT 177  
**TCAACTTAAT**GAATCATCATAACGTAG**CA**ATTAC**CTA**CAC**GATTTG**GCAGGGTGGATT 116  
 TTAATTTGATGAATCATCATAACGTAG**CAATCAAGCTAC**ATGAC**CTG**GCAGGGTGGATT 177  
 TTAAT**TTAAT**GAATAACCATAACGTAG**CAATCAAGCTAC**ATGAC**CTG**GCAGGGTGGATT 177  
**GTAACTT**GATGAATCATCATAACGTAG**CAATCAAGCTAC**ATGAC**CTG**GCAGGGTGGATT 138  
 TTAATTTGATGAATCATCATAACGTAG**CAATCAAGCTAC**ATGAC**CTG**GCAGGGTGGATT 89  
 -----  
 -----  
 .....130.....140.....150.....160.....170.....



cluster565b

FTSPZO101A2LOY/1-134  
is serum c150/232-376  
FSZ2EC10ICV6FM/1-150  
FCPU0RF01DWGD1/74-139

K M Y K K N L F N R F W D L Y D P K E  
**AAGATGTACAAGAAGAATTTGTTCAATCGGTTTTGGGA-TCTGTA-CGAT-CCGAA-GGA** 86  
AAGATGTACAAGAAGAATTTGTTCAATCGGTTTTGGGA-TCTGTA-CGAT-CCGAA-GGA 86  
AAGATGT**TATA**AGAAGAATTTGT**TACAGCAAGTTTTTT****GAATTTG****TACC**GATCCCGAAAGGA 91  
----- 7  
.....10.....20.....30.....40.....50.....



FTSPZO101A2LOY/1-134  
is serum c150/232-376  
FSZ2EC10ICV6FM/1-150  
FCPU0RF01DWGD1/74-139

Frame -1 p = 0.003  
H T H V M M E D **L D H E A V E K L S I**  
**G-CATACGCACGTGATGATGGAAGATTTGGACCACGAGGCAGTGGAGAAGTTGAGCATCA** 31  
G-CATACGCACGTGATGATGGAAGATTTGGACCACGAGGCAGTGGAGAAGTTGAGCATCA 31  
**GCCATACGCAT**GTGATGATGGAAGATTTGGAC**CAT**GAGGCAGTGGAGAAGTT**GAGTATCA** 31  
-----**ATCTGGGC****CATGAA**GCAGTGGAGAAG**CTCAGTATCA** 1  
.....70.....80.....90.....100.....110.....



FTSPZO101A2LOY/1-134  
is serum c150/232-376  
FSZ2EC10ICV6FM/1-150  
FCPU0RF01DWGD1/74-139

**N F I K T V**  
**ACTTTATCAAGACGGTA** 1  
ACTTTATCAAGACGGTA 1  
ACTTTATCAAG**ACA**ATT 1  
**AT**TTTATCAAGACGATT 1  
.....130.....140.....1



cluster589

Frame -3 p = 8.4e-04

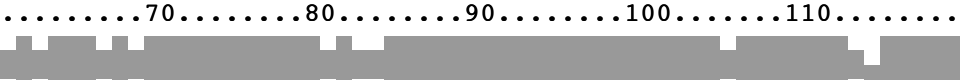
GB3LKKR01DFSJ1/79-314  
GB3LKKR02FHHQ3/61-296  
GB3LKKR02FG5RN/1-236  
contig16113/1-235

**E S Y I N K L Y R L F L H V Y K F F R**  
**GAGTCTTATATTAATAAGTTGTATAGGCTATTCTTACATGTCTATAAATTCTTCCGGA** 177  
GAGTCTTATATTAATAAGTTGTATAGGCTATTCTTACATGTCTATAAATTCTTCCGGA 177  
GAGTCTTATATTAATAAGTTGTATAGGCTATTCTTACATGTCTATAAATTCTTCCGGA 177  
**GAA**TCTTATATTGGTAAG**TTA**TATAGG**TTG**TTCTTATATGTC**TCT**AAATTCTTCC**CGTA** 176



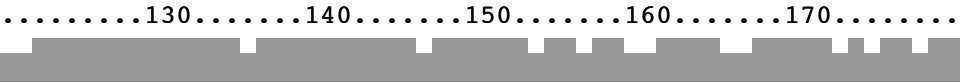
GB3LKKR01DFSJ1/79-314  
GB3LKKR02FHHQ3/61-296  
GB3LKKR02FG5RN/1-236  
contig16113/1-235

**N W H L P E F G S D L S A Y S G R I M F**  
**ATTGGCATTTCCTGAGTTCGGCTCTGATCTTAGTGCTTATTCGGGTCGTATTATGTTTA** 117  
ATTGGCATTTCCTGCGTTCGGCTCTGATCTTAGTGCTTATTCGGGTCGTATTATGTTTA 117  
ATTGGCATTTCCTGAGTTCGGCTCTGATCTTAGTGCTTATTCGGGTCGTATTATGTTTA 117  
**ACT-G**CATT**TA**CCTGACT**TTTGGT**TCTGATATTAGTGCT**TAC**TCCGGTTCGTATT**AA**TTTTA 116



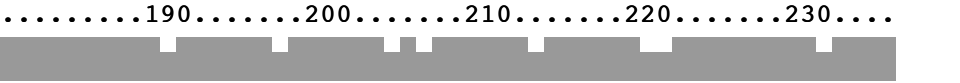
GB3LKKR01DFSJ1/79-314  
GB3LKKR02FHHQ3/61-296  
GB3LKKR02FG5RN/1-236  
contig16113/1-235

**I L L K T G I E Y E K K N Y E S L R D V**  
**TTCTTAAAACGGGTATAGAGTATGAAAAGAAAAGAATTATGAAAGTCTACGAGACGTAT** 57  
TTCTTAAAACGGGTATAGAGTATGAAAAGAAAAGAATTATGAAAGTCTACGAGACGTAT 57  
TTCTTAAAACGGGTATAGAGTATGAAAAGAAAAGAATTATGAAAGTCTACGAGACGTAT 57  
**TC**ATTAAA**ACA**GGTATAGAGTATGAAAAGAAA**CG**GATTATGAAAGTATGCGAAATGTAT 57



GB3LKKR01DFSJ1/79-314  
GB3LKKR02FHHQ3/61-296  
GB3LKKR02FG5RN/1-236  
contig16113/1-235

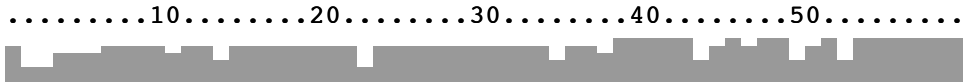
**F N I R S A N P D I S D C M F A L P A**  
**TCAACATACGTTCCGCTAACCCGGACATATCGGATTGTATGTTTGCCTGCTGCG** 1  
TCAACATACGTTCCGCTAACCCGGACATATCGGATTGTATGTTTGCCTGCTGCG 1  
TCAACATACGTTCCGCTAACCCGGACATATCGGATTGTATGTTTGCCTGCTGCG 1  
**ATG**CCCTTCGTTCC**CAATACTCC**GAAC**TCT**CGGATTGTATGTTTGCCTGCG 1



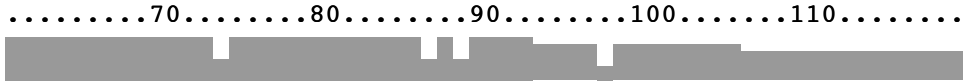
cluster612a

Frame +3 p = 0.004

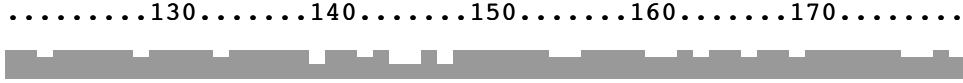
GB3LKKR01DCM8C/36-229 **TTATGTGCATTCAAGATGTCACAGAACTTGAACATTACAAGAAAAGA**CTTGAACAGA 60  
 GB3LKKR01CL8OT/1-157 -----TACAAGAAAAGAACCTTGAACAGA 23  
 GB3LKKR01CR3JI/368-459 TTATGTGCATTTAAGATGTCACAGAACTTGAACATTACAAGAAA**GA**CTTGAACAGA 59  
 GB3LKKR01ESSK1/178-371 TTATGTGCATTTAAGATGTCACAGAACTTGAACATTACAAGAAAAGA**CT**TGAACAGA 60  
 GB3LKKR01E2YMD/41-145 **CTA**TGTGCATTCAAGAT**GTCT**CAGA**ACTTGAATATCACAAGGAA**-**GAGCTA**GAACAGA 59  
 GB3LKKR01CLFJW/1-193 **CGA**TGTGCATTCAAGAT**GTCT**CAGA**ACTTGAATATCACAAGGAA****GAGCTA**GAACAGA 60  
 GB3LKKR01BDX16/1-193 **CTTAGTGC**CTTCAAGAT**GTCC**CAGA**ACTTGAACATTACAAGGAA**GA**TTA**GAACAGA 60



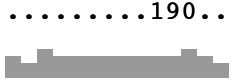
GB3LKKR01DCM8C/36-229 **AGAATCAATACAAGGATTTGGATGTA**CTTGAAGAT**ACCGTTGGTATCAAGAAATACATCC** 120  
 GB3LKKR01CL8OT/1-157 AGAATCAATACAAGGATTTGGATGTA**CTTGAAGATACCGTTGGTATCAAGAAATACATCC** 83  
 GB3LKKR01CR3JI/368-459 AGAATCAATACAAGGATTTGGATGTA**CTTGAAG**----- 92  
 GB3LKKR01ESSK1/178-371 AGAATCAATACAAGGATTTGGATGTA**CTTGAAGATACCGTTGGTATCAAGAAATACATCC** 120  
 GB3LKKR01E2YMD/41-145 AGAATCAATACA**AAA**GATTTGGATGTA**TTA**GAAGAT**ACCGTTGGTAT**----- 105  
 GB3LKKR01CLFJW/1-193 AGAATCAATACA**AAA**GATTTGGATGTA**TTA**GAAGAT**ACCGTTGGTATCAAGAAATACATCC** 120  
 GB3LKKR01BDX16/1-193 AGAATCAATACA**AAA**GATTTGGATGTA**TTA**GAAGAT**ACCGTTGGTATCAAGAAATACATCC** 120



GB3LKKR01DCM8C/36-229 **GAGGCTTGATTTCGAAAGGCAAAGGGAAAAAAGAAAT**TATCTCAAAGACAGTTGAGA**AA**T 180  
 GB3LKKR01CL8OT/1-157 GAGGCTTGATTTCGAAAGGCAAAGGGAAAAAAGAAAT**TATCTCAAAGACAGTTGAGA**AA**T** 143  
 GB3LKKR01CR3JI/368-459 ----- 92  
 GB3LKKR01ESSK1/178-371 **GATGCT**TTGATTTCGAAAGGCAA**AG**AAAAAAGAAAT**TATCTCAAAGACAGTTGAGA**AA**T** 180  
 GB3LKKR01E2YMD/41-145 ----- 105  
 GB3LKKR01CLFJW/1-193 GAGGCTTGGTT**CG**-AAA**GGG**AAA**GAC**AGAAAAGAAAT**TATCTCAAAGACAGTTGAGA**AA**T** 179  
 GB3LKKR01BDX16/1-193 GAGGCTTGATTTCGAA**AGGAA**-**GAT**AGAAA**GAGGTTATTACC**AAA**CCGTTGAGCTAA** 179



GB3LKKR01DCM8C/36-229 **TCGGAGTACCAAGA** 194  
 GB3LKKR01CL8OT/1-157 TC**GAA**GTACCAAGA 157  
 GB3LKKR01CR3JI/368-459 ----- 92  
 GB3LKKR01ESSK1/178-371 **TTGAA**GTACCAAGA 194  
 GB3LKKR01E2YMD/41-145 ----- 105  
 GB3LKKR01CLFJW/1-193 TC**GAA**GTAC**CTAAG** 193  
 GB3LKKR01BDX16/1-193 **GGGAGCGTAAAAG** 193



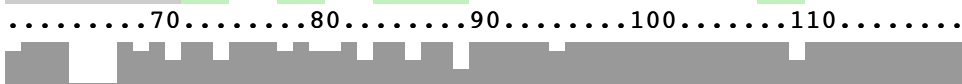
cluster787b

Frame +2 p= 0.016

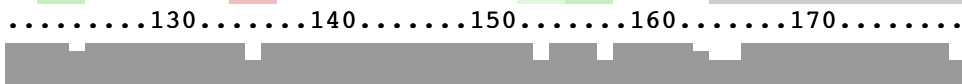
	H N Q L T C Q P I T N I T Y E S	
GB3LKKR01CCRLP/11-280	<b>CACAACCAACTCACATGCCAACCAATCACCAACATCACCTA-CGAAAGCA</b>	50
GB3LKKR01A9PZ6/1-317	-----AACTCACATGCCAACCAATCACCAACATCACCTA-CGAAAGCA	42
GB3LKKR01B6CHJ/4-336	TACT <b>TG</b> CCAACTC <b>ACCT</b> GCC <b>CAG</b> CCAATCACCAACATCACCTA-CGAAAGCA	59
GB3LKKR01D9CHJ/12-353	TACT <b>TGT</b> CAACTC <b>ACCTGT</b> C <b>AG</b> CCAATCACCGACATCACCTA-CGAA <b>TCAA</b>	59
contig03861/84-406	TACT <b>TGC</b> CAACTC <b>ACCT</b> GCCAACCAATCACCGACATCACCTA-CCAG <b>TCAA</b>	59
GB3LKKR02JV51K/1-304	-----ACCC <b>GGCCAGCCAACCGAATACGCAAGC</b> -CGTACC <b>GGCCAC</b>	41



	I R T H P A T A I D L D A I D H E G H	
GB3LKKR01CCRLP/11-280	<b>TCCG---CACCCACCCAGCCACAGCAATAGACCTAGACGCCATCGACCACGAAGGCCACT</b>	107
GB3LKKR01A9PZ6/1-317	TCCG---CACCCACCCAGCCACAGCAATAGACCTAGACGCCATCGACCACGAAGGCCACT	99
GB3LKKR01B6CHJ/4-336	TCCG---CACCCACCCAGCCT <b>TCCGCAATCGACCTC</b> GACGCCATCGACCACGAAGGCCACT	116
GB3LKKR01D9CHJ/12-353	TCCG---CACCCACCCAGCA <b>ACCGCA</b> ATAGACCTAGACGCCATCGACCACGAAGGCCACT	116
contig03861/84-406	TCCG--- <b>CACACAT</b> CCAGCCAC <b>AGCCATC</b> GACCTAGACGCCATCGAC <b>CATGA</b> AGGCCACT	116
GB3LKKR02JV51K/1-304	GCCGCTTC-CA <b>CATCCA</b> ACCACAG <b>CCATC</b> GACCTAGACGCCATCGAC <b>CATGA</b> AGGCCACT	100



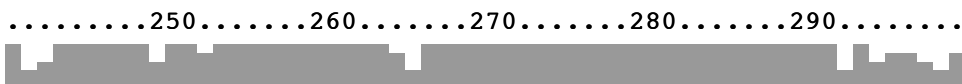
	W V K A S A S K S T T N S T P Q H C K I	
GB3LKKR01CCRLP/11-280	<b>GGGTAAAAGCATCCGCTTCCAAATCGACTACAAACTCAACCCACACAACACTGTAAAATAA</b>	167
GB3LKKR01A9PZ6/1-317	GGGTAAAAGCATCCGCTTCCAAATCGACTACAAACTCAACCCCA- <b>AA</b> CACTGTAAAATAA	158
GB3LKKR01B6CHJ/4-336	GGGTAAAAGCATCCGCTTCCAAATCGACTACAAACTCAACCCACACAACACTGTAAAATAA	176
GB3LKKR01D9CHJ/12-353	GGGTAAAAGCATCCGCTTCCAAATCGACTACAAACTCAACCC <b>CCG</b> CAACACTGTAAAATAA	176
contig03861/84-406	GGGTAAAAGCATCC <b>GATT</b> TCCAAATCGACTACAAG <b>CTCG</b> ACCCACGACACTGTAAAATA-	175
GB3LKKR02JV51K/1-304	GG <b>GTG</b> AAAAGCATCC <b>GAT</b> TCCAAATCGACTACAAG <b>CTCG</b> ACCCCA-GACACTGTAAAATA-	158



	K T A *	
GB3LKKR01CCRLP/11-280	<b>AAACGGCATAA</b>	215
GB3LKKR01A9PZ6/1-317	AA-CGGCATAA	203
GB3LKKR01B6CHJ/4-336	AAACGGCATA <b>A</b>	222
GB3LKKR01D9CHJ/12-353	AAACGGCATA <b>A</b>	229
contig03861/84-406	-----	210
GB3LKKR02JV51K/1-304	-----	192



GB3LKKR01CCRLP/11-280		270
GB3LKKR01A9PZ6/1-317		263
GB3LKKR01B6CHJ/4-336		282
GB3LKKR01D9CHJ/12-353		288
contig03861/84-406		269
GB3LKKR02JV51K/1-304		250



GB3LKKR01CCRLP/11-280		270
GB3LKKR01A9PZ6/1-317		317
GB3LKKR01B6CHJ/4-336		333
GB3LKKR01D9CHJ/12-353		342
contig03861/84-406		323
GB3LKKR02JV51K/1-304		304

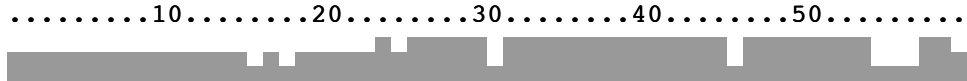


cluster956b

Frame -1 p = 0.003

F Q T G N S T Q F G Y N S Y K P E E N K

```
is_serum_c7126/2-160 TTTCAAACTGGTAACAGCACCCCAGTTTGGATACAATTCCTACAAAACCAGAAAGAAACAAA 100
FCPU0RF0ICTE83/27-185 TTCAAAACTTGGTAACAGCACCCCAGTTTGGGTACAACTTCCTACAAGCCAGACATAAACAAC 100
is_serum_c1111/1-136 TTTCAAACTGGTAACAGCACCCCAGTTTGGCTACAACTCATACAATCCACAAAACAACCAA 100
is_serum_c3480/1-159 TTCAAAACTGGTAACAGCACCCCAGATTTGGGTACAACTTCCTACAAGCCAGACACATCCAAC 100
```



I K Q A N A A Y W K A L T Q Q N D Q A T

```
is_serum_c7126/2-160 ATAAACAAGCAAATGCAGCATACTGGAAGGCACTAACCAAAACGACCAAGCAACA 40
FCPU0RF0ICTE83/27-185 ATAAAGCAAGCAAATGAAGCATACTGGACAGCACTCACAAAAAGCAACGACCAGGCAACA 40
is_serum_c1111/1-136 ATACAAAATGCAAATAACGCATACTGGACAGCACTCACAAAAACGACCAAGCAACA 40
is_serum_c3480/1-159 ATAAAGCAAGCAAATGAAGCATACTGGACAGCACTCACAAAAAGCAACGACCAGGCAACA 40
```

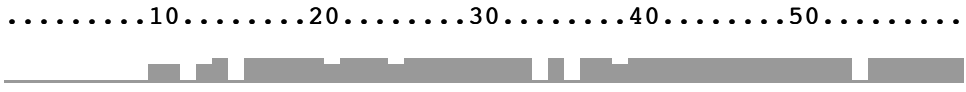


Q I G Q A R A Q Q F E Y H

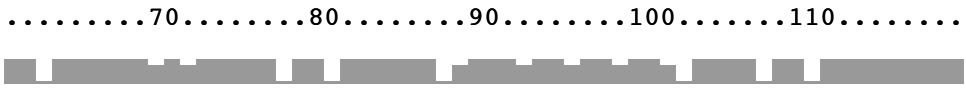
```
is_serum_c7126/2-160 CAAATAGGACAAGCTCGTGCCCAACAATTTGAAATACCAC 1
FCPU0RF0ICTE83/27-185 CAAATAGGCAAGCTCGTGCCCAAGCAATTTGAAATACCAC 1
is_serum_c1111/1-136 CAAATAGGCAAGCCC----- 1
is_serum_c3480/1-159 CAAATAGGCAAGCTCGTGCCCAAGCAATTTGAAATACCAC 1
```



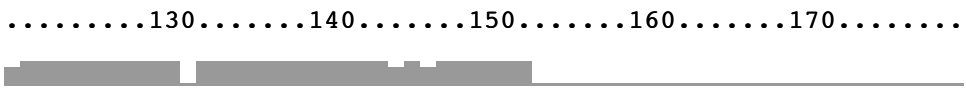
R Y Q S F N G W L D R L O O M A O T G F  
 FTSPZO101CXZAB/1-265 **AGATATCAATCGTTTAATGGTTGGCTAGACCGTCTGCAACAAATGGCACAGACAGGCTTT** 60  
 FTSPZO101C1IV6/1-308 -----ATCGTTTAATGGTTGGTTAGACCGTT**TTACAGCAAATGGCACAGACAGGCTTT** 52  
 FTSPZO101BQG61/1-318 **GCTAATCAGTCTTTCAATGGATGGTTAGACCGTTTACAACAAATGGCACAGACAGGCTTT** 60  
 FTSPZO101D4HEW/1-224 **AGCTCTGCAGATATCAATGGTTGGTTAGACCGTCTGCAACAAATGGCACAG**ACGGGCTTT**** 60  
 FS22EC101BJEV2/1-317 **GCTAATCAGTCTTCAATGGTTGGTTAGACCGTTTACAACAAATGGCACAG**ACTGGCTTT**** 60  
 FSTRRC101DG34E/1-94 ----- 0  
 FTSPZO101A9YRG/1-105 ----- 0  
 FS22EC101BSKRG/1-127 ----- 0



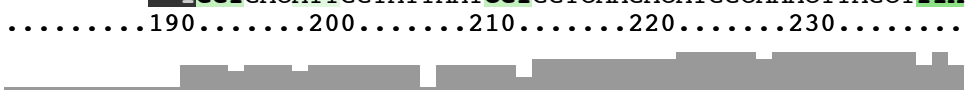
N A Q S N L S O A O N N L T T G R T N A  
 FTSPZO101CXZAB/1-265 **AACGCTCAATCTAACTTATCTCAAGCTCAAAATAATCTTACGACAGGCCGCACAAATGCG** 120  
 FTSPZO101C1IV6/1-308 AACGCTCAAACAAACT**TTGTGCG**CAAGCT**CAAAACAACCTAACCACAGGCCGT**TACAAATGCG 111  
 FTSPZO101BQG61/1-318 **AATGCTCAATCTAACTTTGCTCAAGCTCAAAATAACCTTACGGCA**GGCGGT****TACAAATGCG 120  
 FTSPZO101D4HEW/1-224 **AATGCTCAATCTAACTTATCACAAGCTCAAAATAACCTTACGGCA**GGT****CGCACAAATGCG 120  
 FS22EC101BJEV2/1-317 **AATGCTCAATCTAACTTTGTGCG**CAAGCTCAAAATAACCTTACGGCA**GGT**CGCACAAATGCG 120  
 FSTRRC101DG34E/1-94 ----- 0  
 FTSPZO101A9YRG/1-105 ----- 0  
 FS22EC101BSKRG/1-127 ----- 0



L Q Y G D T S A L D V G M G K D V L G I  
 FTSPZO101CXZAB/1-265 **CTACAGTATGGCGATACAAGTGCCTTGATGTTGGCATGGGCAAAGACGTTCTAGGTATC** 180  
 FTSPZO101C1IV6/1-308 CTACAGTATGGCGATACAAGTGCCT**TTAGATGTTGGCTTGGTAAAGATGTTT**TTAGGCATC 171  
 FTSPZO101BQG61/1-318 **TTACAGTATGGT**GATACAAGTGCCTTGATGTTGGCATGGGC**AAGGACGTTT**TTAGGCATT 180  
 FTSPZO101D4HEW/1-224 CTACAGTATGGCGATACAAGTGCCTTGATGTT----- 153  
 FS22EC101BJEV2/1-317 CTACAGTAT**GGT**GATACAAGTGCCTTGATGTTGGCATG**GGT**TAAAGACGTT**TTAGGC**ATC 180  
 FSTRRC101DG34E/1-94 ----- 0  
 FTSPZO101A9YRG/1-105 ----- 0  
 FS22EC101BSKRG/1-127 ----- 0



N Q O R G D V G I N R G O D M A N L A L  
 FTSPZO101CXZAB/1-265 **AATCAACAACGTTGGCGACGTTGGTATTAATCGCGGTCAAGACATGGCAAACCTTAGCTCTT** 240  
 FTSPZO101C1IV6/1-308 AATCAACAACGTT**GGT**GACATTTGGT**TATAAATCGCGGTCAAGACATGGCAAACCTTAGCTCTT** 231  
 FTSPZO101BQG61/1-318 AATCAA**CAG**CGT**GGT**GACATTTGGTATTAAT**CGT**GGTCAAGACATGGCAAACCTTAGCT**TTA** 240  
 FTSPZO101D4HEW/1-224 ----- 153  
 FS22EC101BJEV2/1-317 AATCAACAACGTT**GGT**GACATTTGGT**TATAAATCGT**GGTCAAGACATGGCAAACCTTAGCT**CTA** 240  
 FSTRRC101DG34E/1-94 -----ATGGCAAACCTTAGCT**CTA** 18  
 FTSPZO101A9YRG/1-105 -----CGGTCAAGACATGGC-AACTTAGCT**CTA** 27  
 FS22EC101BSKRG/1-127 -----**TGGT**GACATTTGGTATTAAT**CGT**GGTCAAGACATGGCAAACCTTAGCT**TTA** 49



G R G T I G A N  
 FTSPZO101CXZAB/1-265 **GGTCGTGGTACGATTGGTGCTAAT** 265  
 FTSPZO101C1IV6/1-308 GGTTCGTGGTACGATTGGTGCTAAT 290  
 FTSPZO101BQG61/1-318 GGTTCGTGGTACGATTGGTGCTAAT 300  
 FTSPZO101D4HEW/1-224 -----**GGCACGATACTGCAGAG** 206  
 FS22EC101BJEV2/1-317 GGTTCGTGGTACGATTGGTGCTAAT 299  
 FSTRRC101DG34E/1-94 GGTTCGTGGTACGATTGGTGCTAAT 78  
 FTSPZO101A9YRG/1-105 GGTTCGTGGTACGATTGGTGCTAAT 87  
 FS22EC101BSKRG/1-127 GGTTCGTGGTACGATTGGTGCTAAT 109



FTSPZO101CXZAB/1-265 265  
 FTSPZO101C1IV6/1-308 308  
 FTSPZO101BQG61/1-318 318  
 FTSPZO101D4HEW/1-224 224  
 FS22EC101BJEV2/1-317 317  
 FSTRRC101DG34E/1-94 94  
 FTSPZO101A9YRG/1-105 105  
 FS22EC101BSKRG/1-127 127



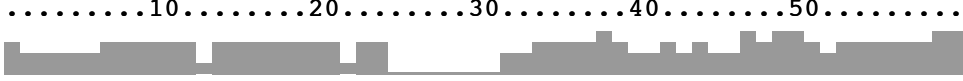


cluster1073

Frame -3 p = 1.7e-10

K T P E L S K A K K A L D A Y L K

contig05414/1-269 **AAGACTCCGGAGCTTTCAAAAGCAAAAAGGCATTGGATGCTTATCTTAA** 219  
 GB3LKKR01DOS1W/1-272 ---AT**CCAGACCTCACAAAGGCAAGAAAGCTTTGGACGCCTTATCTCAA** 213  
 contig07331/1-285 AAGACT**CCAGAACTCACAAAGGCAAGAAAGCTTTGGACGCCTTATCTCAA** 226  
 GB3LKKR01CLC55/1-255 AAGACT**CCAGAACTCACAAAGGCAAGAAAGCTTTGGACGCCTTATCTCAA** 226  
 GB3LKKR01A150I/1-276 AAGACTCCG**GAAGCTCACCAAGGCAAGAAAGCTTTGATGCTTACCTTAA** 226



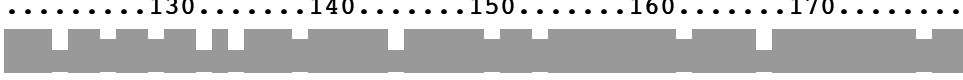
E N K L D P O K D W T K D K K H G K K V

contig05414/1-269 **GAAAACAAGTTGGACCCCTCAAAGGATTGGACCAAAGACAAGAAACACGGTAAGAAGGTT** 159  
 GB3LKKR01DOS1W/1-272 **GAGAACAAGTTGGACCCA**ACT**TAAGAT**TGGACC**AAGGACAAGAAA**CATGGTAAGAAGATT 153  
 contig07331/1-285 **GAGAACAAGTTGGACCCA**ACT**TAAGAT**TGGACC**AAGGACAAGAAA**CATGGTAAGAAGATT 166  
 GB3LKKR01CLC55/1-255 **GAGAACAAGTTGGACCCA**ACT**TAAGAT**TGGACC**AAGGACAAGAAA**CATGGTAAGAAGATT 166  
 GB3LKKR01A150I/1-276 **GAGAACAAGTTGGACCCA**ACT**TAAGAT**TGGACC**AAGGACAAGAAA**CATGGTAAGAAGTT 166



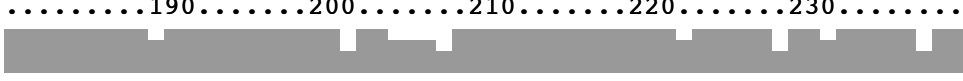
T E L V N K L N K E R D K V A A E Y P E

contig05414/1-269 **ACCGAACTTGTTAACAAGCTCAACAAGGAACGGGATAAAGTCGCTGCTGAATATCCTGAA** 99  
 GB3LKKR01DOS1W/1-272 ACCGAACTTGTTAACAAG**CTGAACAAGGAAAGAGACAAA**GTTCGCTGCT**GCC**TAT**CCG**GAA 93  
 contig07331/1-285 ACCGAACTTGTTAACAAG**CTGAACAAGGAAAGAGACAAA**GTTCGCTGCT**GCC**TAT**CCG**GAA 106  
 GB3LKKR01CLC55/1-255 ACCGAACTTGTTAACAAG**CTGAACAAGGAAAGAGACAAA**GTTCGCTGCT**GCC**TAT**CCG**GAA 106  
 GB3LKKR01A150I/1-276 **ACAGA**ACT**TTGTA**AATAAGCTCAAC**AAAGAACGGGACAAGGTA**GCTGCT**GCTTACCA**GAA 106



K D L K N E A K L V K M K E K K N A E K

contig05414/1-269 **AAGGATTTGAAGAATGAGGCTAAACTCGTTAAGATGAAAGAGAAGAAGAATGCCGAAAAG** 39  
 GB3LKKR01DOS1W/1-272 AAGGAT**GCCGACAACAACAA**GAAG**TTGGTAAAA**CTCAAAGAGAAGAAGGAT**AAA**GAAAA 33  
 contig07331/1-285 AAGGAT**GCCGACAACAACAA**GAAG**TTGGTAAAA**CTCAAAGAGAAGAAGGAT**AAA**GAAAA 46  
 GB3LKKR01CLC55/1-255 AAGGAT**GCCGACAACAACAA**GAAG**TTGGTAAAA**CTCAAAGAGAAGAAGGAT**AAA**GAAAA 46  
 GB3LKKR01A150I/1-276 **GGT**GAC**AAA**GAGAAT**ACCA**AAAA**TTGGTAAAA**CTCAGTAAAGAA**AAAGGC**AAGAAAGAG 46



T E K K K E K K E K L L

contig05414/1-269 **ACTGAAAAGAAA**-----**AAAGAGAAGAAAGAAAAGCTTCTT** 1  
 GB3LKKR01DOS1W/1-272 **TCCGGCAAGAAA****GAAGACAAGAAG**GAGAAG**AAA**GAAAAGAAAT**CT** 1  
 contig07331/1-285 **TCCGGCAAGAAA****GAAGACAAGAAG**GAGAAG**AAA**GAAAAGAAAT**CT** 1  
 GB3LKKR01CLC55/1-255 **TCCGGCAAGGAAGAT**-----**AAAGAGAAGAAAGAAAAGAAATCT** 1  
 GB3LKKR01A150I/1-276 **GAATCTGAAACC**-----**AAAGAGAAGAAAGAAAAGAAATCT** 1



cluster1085

Frame +2 p = 0.006

\* T H K A A

FTSPZO101A5XSA/49-223  
 ER8QEOW01BCS2I/1-68  
 FTSPZO101A2ZKG/1-94  
 FCPU0RF01CACI4/1-115  
 FTSPZO101C8F3X/19-176  
 FS22EC101A98SE/1-175  
 FSTRRC101CPTN7/37-211

**TGAACACACAAAGCGGC** 60  
 ----- 0  
 ----- 0  
 -----TGC 3  
 CGAACGCACAAAGCGGC 60  
**TGAACACATAAAGCGGC** 60  
**TGAACGCACAAAGCGGC** 60

.....10.....20.....30.....40.....50.....



FTSPZO101A5XSA/49-223  
 ER8QEOW01BCS2I/1-68  
 FTSPZO101A2ZKG/1-94  
 FCPU0RF01CACI4/1-115  
 FTSPZO101C8F3X/19-176  
 FS22EC101A98SE/1-175  
 FSTRRC101CPTN7/37-211

N V S Q F T R A E M A A T D S G L D V Q

**CAACGTCTCGCAGTTTACACGCGCAGAAATGGCAGCAACAGACAGTGGCCTTGTATGTGCA** 120  
 -----GCCTTGTATGTGCA 13  
 -----CGCGTAG-ATGGCAGCAACAGACAGTGGCCTT-ATGTGCA 39  
**AAACGTCTCGCAG--ACACGCGCAGAAATGGCCACAACAGATAGCGGCCTTGTATGTTTCG** 60  
**CAATGTCTCGCAGTTCACACGCGCAGAAATGGCAGCAACAGATAGTGGTCTTGTATGTGCA** 120  
**AAACGTCTCGCAATTCACACGCGCAGAGATGGCAGCAACAGACAGCGGCCTTGTATGTGCG** 120  
**CAATGTCTCGCAGTTCACACGCGCAGAGATGGCAGCGACAGATAGTGGTCTTGTATGTGCA** 120

.....70.....80.....90.....100.....110.....



FTSPZO101A5XSA/49-223  
 ER8QEOW01BCS2I/1-68  
 FTSPZO101A2ZKG/1-94  
 FCPU0RF01CACI4/1-115  
 FTSPZO101C8F3X/19-176  
 FS22EC101A98SE/1-175  
 FSTRRC101CPTN7/37-211

I E W I S T N D G R V R D S H R S T

**AATCGAATGGATTAGCACGAATGACGGCAGGGTTAGAGACTCACACAGAAGCACA** 175  
 AATCGAATGGATTAGCACGAATGACGGCAGGGTTAGAGACTCACACAGAAGCACA 68  
 AATCGAATGGATTAGCACGAATGATGGCAGGGTTAGAGACTCGCACCGAAGCGTC 94  
**CATTGAGTGGATTAGCACAAATGATAGTCGGGTTAGAGACTCACACAGAAGCACA** 115  
 AATCGAATGGATTAGCACGAATGACGGCAGGGTTTCAT----- 158  
 AATC**GAGT**GGATTAGCACGAATGATGGCAGGG**GTC**AGAGACTCA**CATCGA**GCAGTC 175  
 AATC**GAGT**GGATTAGCACGAATGATGGCAGGG**GTC**AGAGAC**CCATCGC**AATGTA 175

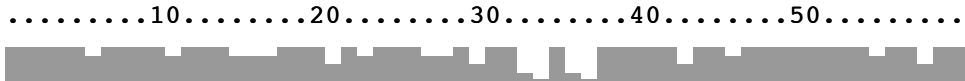
.....130.....140.....150.....160.....170...



cluster1217

Frame -3 p = 7.7e-10

\* S T F G G Q F K T Y K M T N  
**TAGTCT-AC**TTTTGGTGGT**CAGTTCAAGACTTACAAGATGACTAAC** 207  
 -----  
**TAGTCT-AC**TTTTGGTGGT**CAGTTCAAGACTTACAAGATGACTAAC** 207  
**TAGTCTTACT**TTTTGGTGGT**CAGTTCAAGACTTACAAGATGACTAAC** 208  
 -GG**CTTTGGT**TTTTGGTGGT**CAGTTTGTAAACATACAAGATGACAAAT** 207  
**TAAC**TTGACT**TTTCGGTGGT**CAGTT**CAAGACTTATAAGATGACTAAT** 182



GB3LKKR01DTQ2E/16-280  
 GB3LKKR01DOG04/1-117  
 GB3LKKR01AL7BN/1-265  
 GB3LKKR02FUFQO/1-266  
 contig04236/16-280  
 GB3LKKR02GUTL4/19-199

G I E L T L K Y F P L Y D D T T Y N R M  
**GGTATCGAACTTACTCTGAAATATTTCCATTGTATGATGATACCACCTTATAACCGTATG** 147  
 -----  
**GGTATCGAACTTACTCTGAA**TATTTCCATTGTATGATGAT**ACTACTTATAAATCGTATG** 147  
**GGTATCGAACTTACTCTGAA**TATTTCCATTGTATGATGAT**ACTACTTATAAATCGTATG** 148  
**GGC**AT**CGAGTTGACATTG**AAA**CATTTCCCGT**TTGATGATGAT**ACTACTTATAAATCGTTTG** 147  
**GGTATCGAA**TT**GACG**CTGAAAT**TAC**TTCCATT**TTCC**CCAG**GATGTATACAGTATAACCGT** 148  
 .....70.....80.....90.....100.....110.....



GB3LKKR01DTQ2E/16-280  
 GB3LKKR01DOG04/1-117  
 GB3LKKR01AL7BN/1-265  
 GB3LKKR02FUFQO/1-266  
 contig04236/16-280  
 GB3LKKR02GUTL4/19-199

L H P I T L K P L E S Y R L T F L D L G  
**TTGCATCCTATCACACTGAAACCTCTGGAATCATATCGTTTGACATTCCTTGATCTGGGT** 87  
 -----  
 -----TCATATCGTATGACATTCCTTGATCTGGGT 1  
**CTGCATCCG**ATCACACTGAAACCTCTGGAATCATATCGTATGACATTCCTTGATCTGGGT 87  
**CTGCATCCG**ATCACACTGAAACCTCTGGAATCATATCGTATGACATTCCTTGATCTGGGT 88  
**TTACACCCG**GTATCTGGTAA**CCACTGGAATCTT**TAT**AGA**ATGACATTC**TTG**GAT**CTT**GGT 87  
**CCTTACCATTAG**ATCT**CA**GAGTAT**CTTTAG**ATTT**AG**CATATCCAGACGGAGCAACAGCAC 88  
 .....130.....140.....150.....160.....170.....



GB3LKKR01DTQ2E/16-280  
 GB3LKKR01DOG04/1-117  
 GB3LKKR01AL7BN/1-265  
 GB3LKKR02FUFQO/1-266  
 contig04236/16-280  
 GB3LKKR02GUTL4/19-199

R R D G E A N I V K V V R K D R E F V T  
**AGACGTGATGGTGAAGCTAACATTGTTAAGGTAGTTTCGTAAAGATCGTGAATT-CGTTAC** 27  
 AGACGTGATGGTGAAGCTAACATTGTTAAGGTAGTTTCGTAAAGATCGTGAATTCGGTTAC 1  
 AGACGTGATGGTGAAGCTAAC**ATCGTAAAA**GTAGTTTCGTAAAG**GAC**CGTGA-GTTT**GTAAC** 27  
 AGACGTGATGGTGAAGCTAAC**ATCGTAAAA**GTAGTTTCGTAAAG**GAC**CGTGA-GTTT**GTAAC** 28  
 AGACGTGATGGTCAAGCT**AATATCGTTAAGGTT**GTTCGT**AAG**GATCGT**GAGAT-GGTTAT** 27  
**C--**----- 28  
 .....190.....200.....210.....220.....230.....



GB3LKKR01DTQ2E/16-280  
 GB3LKKR01DOG04/1-117  
 GB3LKKR01AL7BN/1-265  
 GB3LKKR02FUFQO/1-266  
 contig04236/16-280  
 GB3LKKR02GUTL4/19-199

W Y T G G A V A  
**TTGGTACTGGTGGTGCTGTTGCT** 1  
 TTGGTATACTGGTGGTGCTGTTGCT 1  
 TTGGT**TAC**ACTGGTGGTGCTGTTGCT 1  
 TTGGT**TAC**ACTGGTGGTGCTGTTGCT 1  
**CTGGAAT**ACTT**CAGGTTCTGTAGCT** 1  
 ----- 1  
 .....250.....260.....

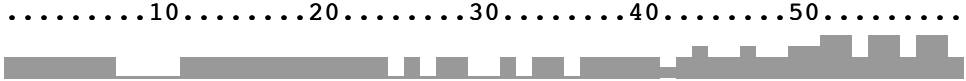


GB3LKKR01DTQ2E/16-280  
 GB3LKKR01DOG04/1-117  
 GB3LKKR01AL7BN/1-265  
 GB3LKKR02FUFQO/1-266  
 contig04236/16-280  
 GB3LKKR02GUTL4/19-199

cluster1423

Frame +1 p= 1.0e-04

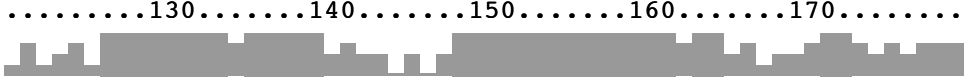
	<b>Q R V Q Y I S P L E D G S Y K V T L V S</b>	
DBA-SLE c4340/1-188	<b>CAGAGGGTGCAGTATATTAGTCCTTTAGAAAGATGGAAGCTACAAAGTTACATTAGTCAGT</b>	60
FTSPZO101BASEK/1-188	CAGAGGGGAAGGTATATTAGTCCTTTAGAAAGATGGAAGCTACAAAGTTACATTAGTCAGT	60
DBA-SLE c13402/1-183	CAGAGGGTGCAGTATATTAGTCCTATTGAC <b>AAGGGCAGCTATAAAGTTACATTGGTTAGC</b>	60
FTSPZO101CIWNM/1-146	----- <b>CTGCAGATATCATTGGTTAGC</b>	21
FTSPZO101DPI37/1-130	----- <b>TTGGTTAGC</b>	9



	<b>G W E M K L T D K L Y S R E K L I G L W</b>	
DBA-SLE c4340/1-188	<b>GGGTGGGAAATGAAGCTTACAGATAAGTTGTACAGTAGAGAGAAGTTGATTGGGTTGTGG</b>	120
FTSPZO101BASEK/1-188	GGGTGGGAAATGAAGCTTACAGATAAGTTGTACAGTAGAGAGAAGTTGATTGGGTTGTGG	120
DBA-SLE c13402/1-183	GGGTGGGAAATG <b>CCGTTGACT</b> GATAAGTTGTACAGTAGAGAGAAGTTGATT <b>GGACAGTGG</b>	120
FTSPZO101CIWNM/1-146	GGGTGGGAAATGAAG <b>TTGACGGGTGAGGTGAT</b> AGTAGAGAGAAGTT <b>GACTGGGTTGTGG</b>	81
FTSPZO101DPI37/1-130	GGGTGGGAAATG <b>CCGTTGACT</b> GATAAGTTGTACAGTAGAGAGAAGTTGATTGGGTTGTGG	69



	<b>G V L Q N K *</b>	
DBA-SLE c4340/1-188	<b>GGAGTTTTGCAGAATAAATAG</b>	180
FTSPZO101BASEK/1-188	GGAGTTTTGCAGAATAAA <b>TAG</b>	180
DBA-SLE c13402/1-183	<b>CGCATGTTGCAGAATAAATA-</b>	173
FTSPZO101CIWNM/1-146	<b>AAGGGGTTGCAGAAACAAATAA</b>	140
FTSPZO101DPI37/1-130	<b>CGTATATTGCAGAATAAA<b>TAG</b></b>	122

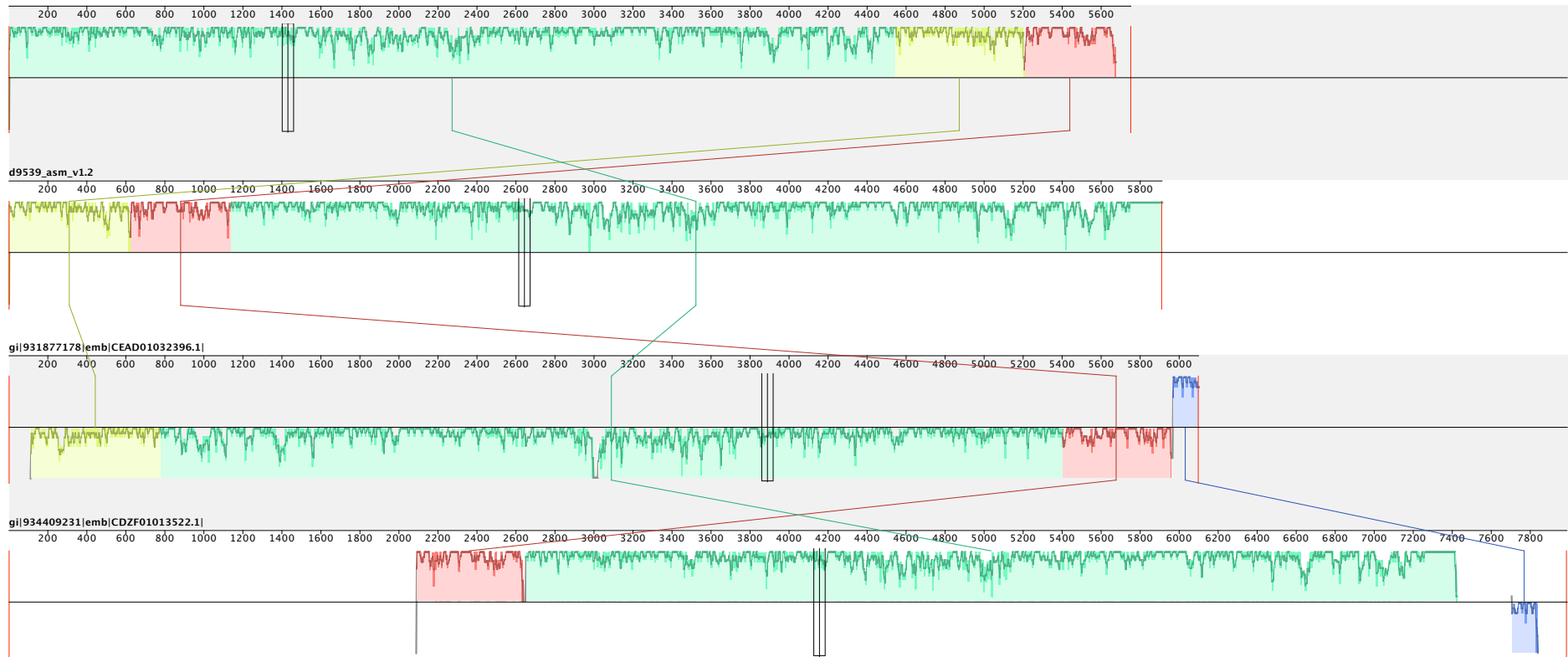


DBA-SLE c4340/1-188	188
FTSPZO101BASEK/1-188	188
DBA-SLE c13402/1-183	183
FTSPZO101CIWNM/1-146	146
FTSPZO101DPI37/1-130	130

.....1



**Supplementary Figure S3.** Whole genome comparison of the bacteriophage HFM assembly with three metagenomic contigs from the NCBI env\_nt database using Mauve<sup>1</sup>. The sequence of the selected env\_nt contigs corresponds mostly to our assembled circular genome. Assuming the contigs are derived from closely related bacteriophages, differences in the order of the conserved blocks (in different colors) may be attributed to different cut points during the linearization of the circular genome, rearrangements, or possible misassemblies. The middle sections of the contigs are well conserved between the different contigs, but the flanking regions differ considerably. Corresponding segments in different contigs are connected with lines in the same colour.

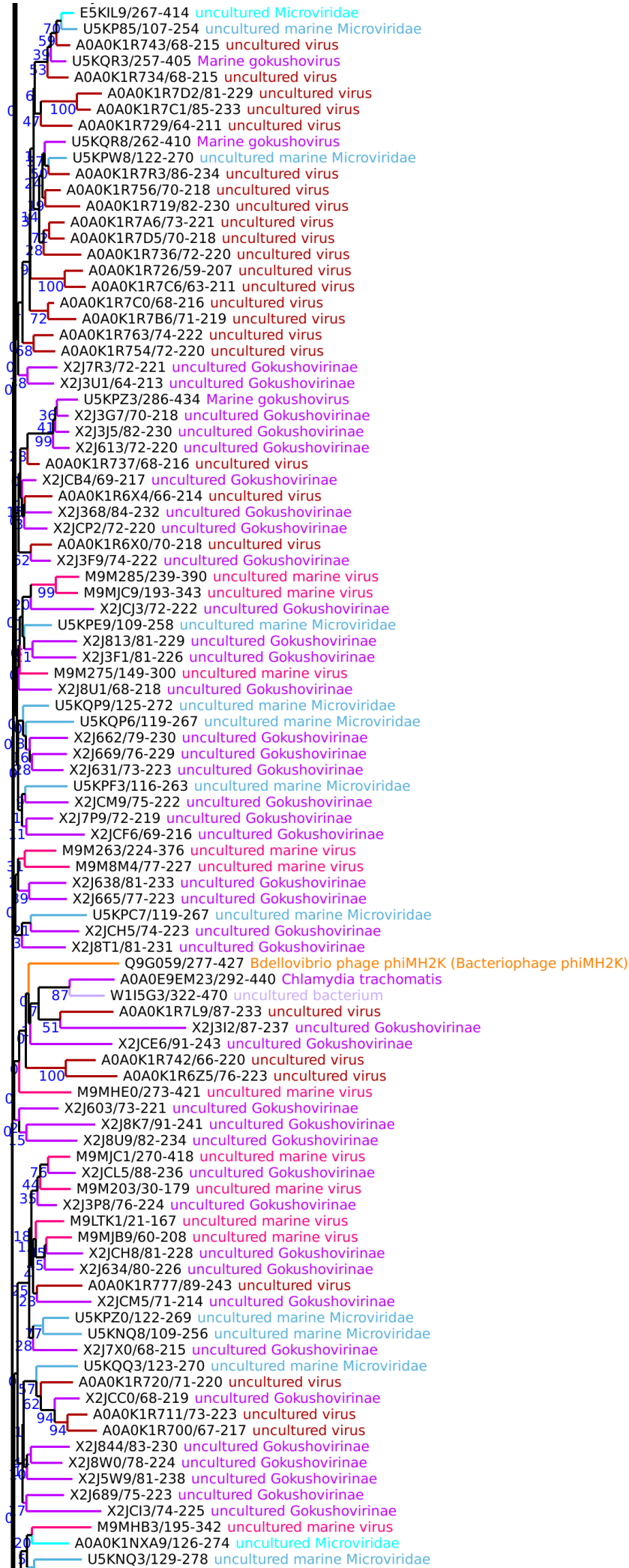


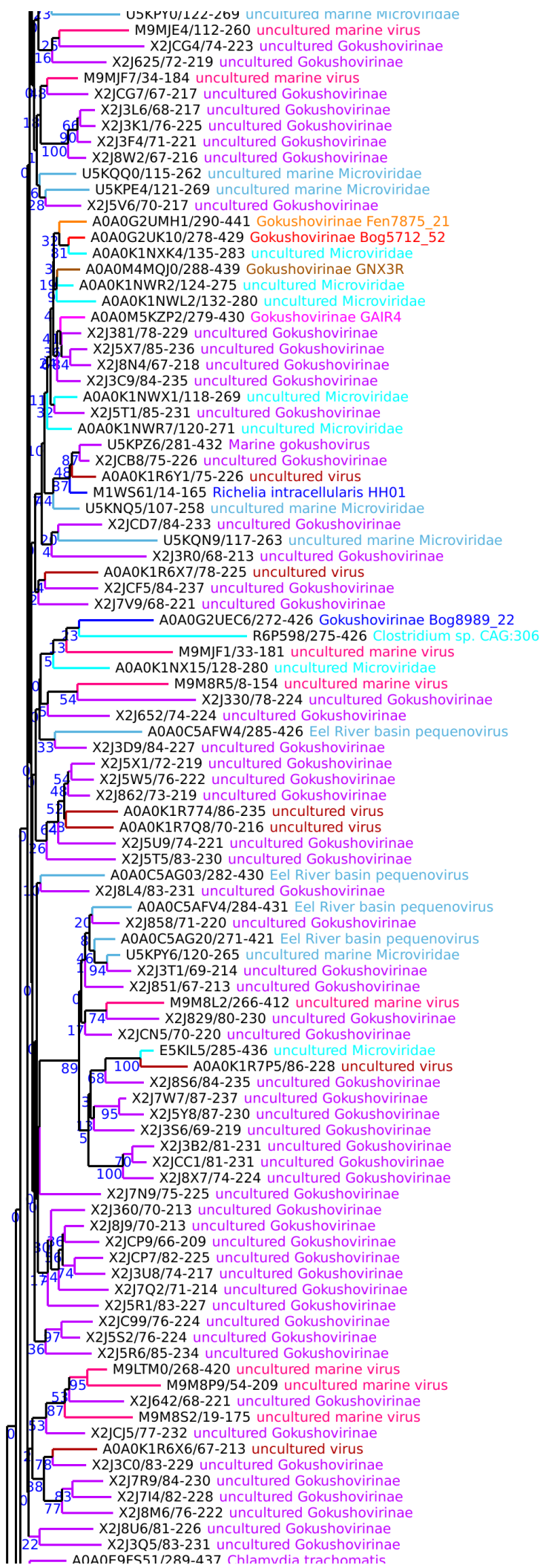
gi|935523605|emb|CDTY01044545.1

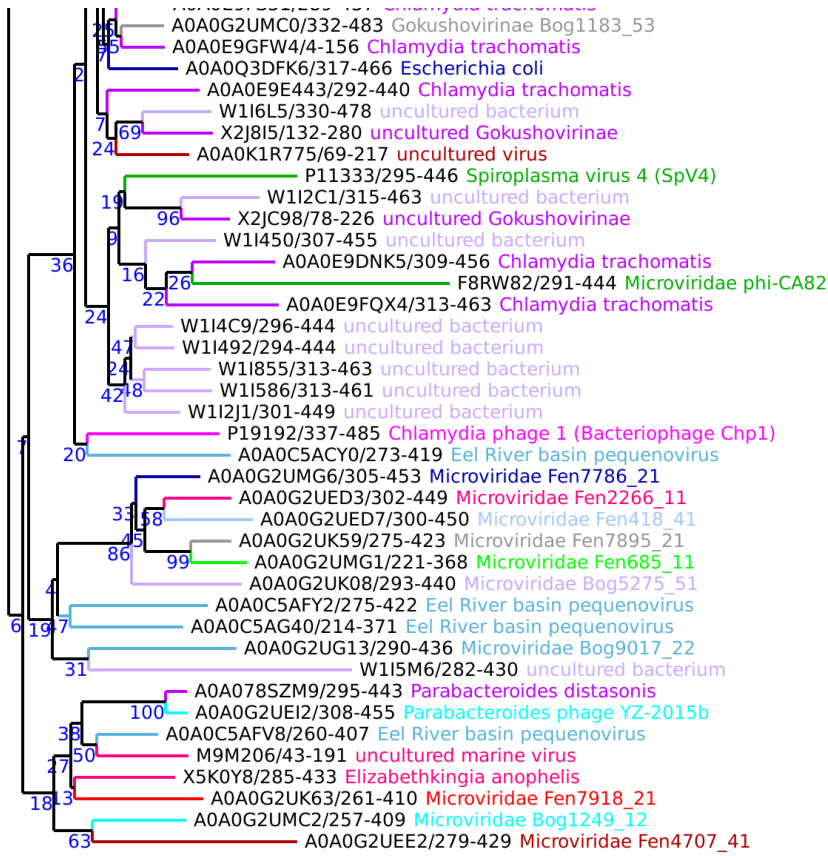
**Supplementary Figure S4.** Neighbour-joining tree of 331 phage capsid F proteins from Uniprot, as well as bacteriophage HFM ORF6 (highlighted in gray). Tree nodes are marked with bootstrap support as percentages. The species as listed in Uniprot is shown to the right and colour coded. Note however that in many cases the listed species is the host in which the phage was found. This will happen for proteins from e.g. “Unassembled WGS sequence” entries.











0.1  
 Tree balance = 0.0

**Supplementary Figure S5.** PHACTS<sup>2</sup> prediction on the bacteriophage HFM predicted ORFs. PHACTS uses a machine learning approach based on protein sequences to classify if a phage is more likely to be temperate or lytic and whether it has a gram-positive or gram-negative host preference. The results do not yield information about the lifestyle of the phage, but they suggest that the phage is likely to invade gram-negative hosts

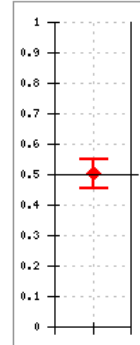
✓ Lifestyle:

### Analysis Statistics

Ten iterations of PHACTS were performed using the default settings. The phage was **non-confidently** predicted as having a **Temperate** lifestyle.

Predicted Class	Averaged Probability	Standard Deviation
Temperate	0.503	0.048

Probability that phage is **Temperate**



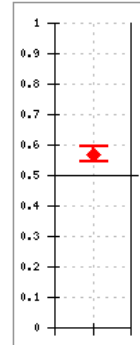
✓ Gram-stain of host:

### Analysis Statistics

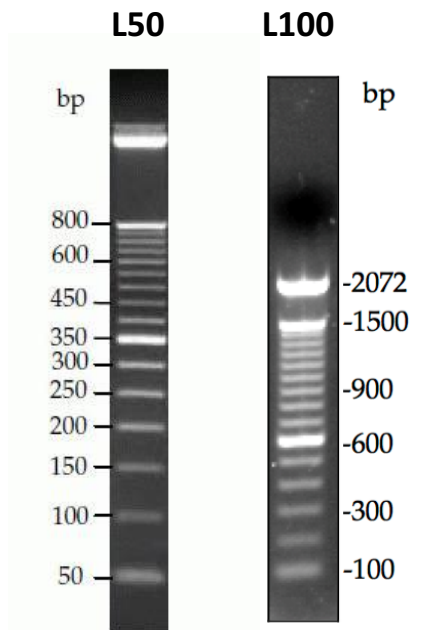
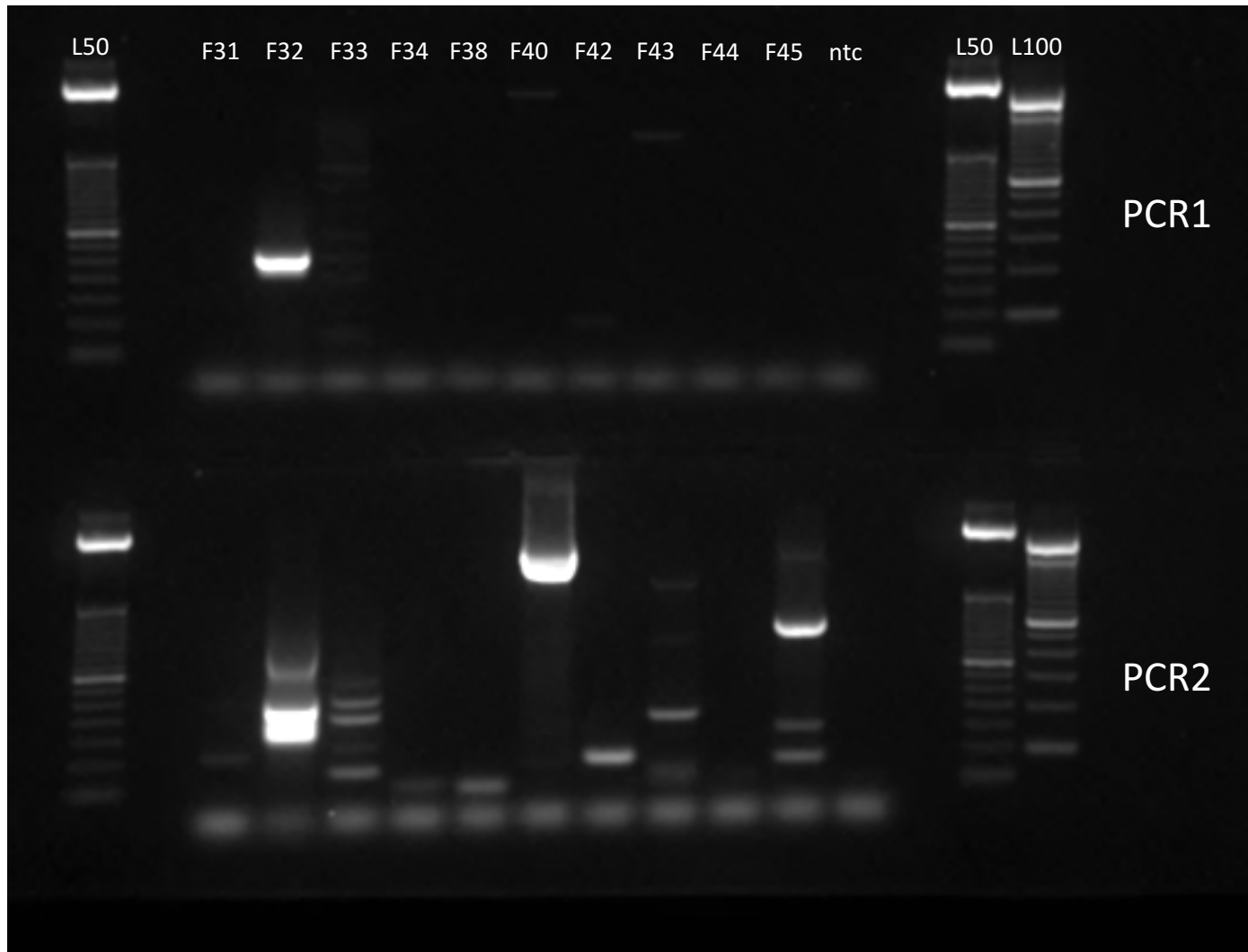
Ten iterations of PHACTS were performed using the default settings. The phage was **confidently** predicted as infecting a **Gram Negative** host.

Predicted Class	Averaged Probability	Standard Deviation
Negative	0.57	0.025

Probability that phage is **Negative**



**Supplementary Figure S6.** Results from two separate Bacteriophage HFM PCRs from the same set of samples. The main image shows one complete gel. The two smaller gel images show the sizes of the bands in the two size markers used in the main gel, see below for description. We have chosen an image where all bands are clearly distinguishable, despite variations in intensities. The PCR amplifications are of bacteriophage HFM in individual FESC D samples. Bands of the expected size were only observed in sample F32, where DNA fragments of 267 bp and 155 bp were obtained in PCR 1 and 2, respectively. Ntc, L50 and L100 refer to non-template control, 100 bp DNA Ladder (ThermoFisher Scientific, [cat.no. 15628050](#)) and 50 bp DNA Ladder (ThermoFisher Scientific, [cat.no. 10416014](#)).



## Supplementary Table S1 - Description of viral metagenomics libraries

library_id	sample_type	nucleic_acid	platform	reads	avg_read_len
ER8QEOW01	nasopharyngeal swabs	DNA	GS	160,120	99.6
ER8QEOW02	nasopharyngeal swabs	RNA	GS	138,941	102.8
FC8LRL301	nasopharyngeal swabs	DNA	GS FLX	165,715	185.9
FC8LRL302	nasopharyngeal swabs	RNA	GS FLX	239,014	208.1
FCPU0RF01	serum	DNA	GS FLX	177,184	143.7
FCPU0RF02	serum	RNA	GS FLX	113,542	140.5
FPFNBIF01	serum	DNA	GS FLX	25,767	170.1
FPFNBIF02	serum	RNA	GS FLX	305,191	199.1
FSTRRC101	serum	DNA	GS FLX	28,217	177.9
FS22EC101	serum	DNA	Titanium	256,310	193.2
FS22EC102	serum	RNA	Titanium	189,704	187.6
FTSPZO101	cerebrospinal fluid	DNA	Titanium	158,065	220.1
FTSPZO102	cerebrospinal fluid	RNA	Titanium	51,683	134.5
GB7HT0Z01	nasopharyngeal & throat swabs	DNA	Titanium	199,388	258.8
GB7HT0Z02	nasopharyngeal & throat swabs	RNA	Titanium	233,531	264.4
GB3LKKR01	feces	DNA	Titanium	723,577	327.0
GB3LKKR02	feces	RNA	Titanium	736,239	323.7
			<b>OVERALL</b>	3,902,188	



**Supplementary Table S2** - Summary of predicted families hits to different databases

Cluster	sample_origin	RNAcode_pvalue	composite score	hits_to
1217	feces	7.70E-10	7.57	env_nt;metahit_2014_cds;metahit_2014_pep;nr;env_nr
339b	feces	1.07E-06	7.52	metahit_2014_cds;env_nt;metahit_2014_pep
211b	feces	3.36E-14	7.49	env_nt;metahit_2014_cds;metahit_2014_pep
562	CSF;serum	4.59E-07	7.39	
297a	feces	0.001	7.38	metahit_2014_cds;env_nt;metahit_2014_pep;
457	CSF;serum	1.35E-08	7.32	env_nt
348	feces	1.42E-08	7.32	env_nt;metahit_2014_cds;metahit_2014_pep;
502	feces	3.60E-07	7.3	metahit_2014_cds;env_nt;metahit_2014_pep;
540	CSF;mucus;serum	3.76E-04	7.24	env_nt
241b	feces	5.54E-06	7.24	env_nt;metahit_2014_cds;metahit_2014_pep;
179a	feces	0.005	7.23	metahit_2014_cds;env_nt;metahit_2014_pep;
406b	CSF;feces	0.001	7.21	env_nt;metahit_2014_cds;metahit_2014_pep
1073	feces	1.66E-10	7.19	metahit_2014_cds;env_nt;metahit_2014_pep
211a	feces	3.90E-06	7.17	metahit_2014_cds;env_nt;metahit_2014_pep;
589	feces	8.36E-04	7.12	metahit_2014_cds;env_nt;metahit_2014_pep;
375a	feces	0.013	7.09	metahit_2014_cds;env_nt;metahit_2014_pep;
182a	feces	5.23E-04	7.05	env_nt;metahit_2014_cds;metahit_2014_pep;
1423	CSF;serum	1.00E-04	7.01	
258	feces	0.015	6.99	env_nt;metahit_2014_cds;metahit_2014_pep;
297b	feces	0.006	6.99	metahit_2014_cds;env_nt;metahit_2014_pep;
1057	CSF;serum	3.44E-06	6.99	
241a	feces	0.002	6.94	env_nt;metahit_2014_cds;metahit_2014_pep;
956b	serum	0.003	6.94	nr
321b	feces	1.88E-07	6.93	metahit_2014_cds;env_nt;metahit_2014_pep;
182b	feces	1.54E-04	6.9	metahit_2014_cds;env_nt;metahit_2014_pep;
532	feces	0.001	6.79	metahit_2014_cds;env_nt;metahit_2014_pep;nr
787b	feces	0.016	6.78	env_nt;metahit_2014_pep;

565b	CSF;serum	0.003	6.74 nr
1085	CSF;mucus;serum	0.006	6.73 metahit_2014_pep
612a	feces	0.004	6.56 metahit_2014_cds;env_nt;metahit_2014_pep;
179b	feces	0.01	6.54 metahit_2014_cds;env_nt;
113b	feces	0.008	6.32 metahit_2014_cds;env_nt;metahit_2014_pep;

**Supplementary Table S3** - Summary of hits of protein families to MetaHIT database. "n" stands for the CDS version, "p" stands for the protein version of the database

Cluster	composite score	sample_origin	MetaHIT
1217	7.57	feces	2014_n;2014_p
339b	7.52	feces	2014_n;2014_p
211b	7.49	feces	2014_n;2014_p
562	7.39	CSF;serum	
297a	7.38	feces	2014_n;2014_p
457	7.32	CSF;serum	
348	7.32	feces	2014_n;2014_p
502	7.3	feces	2014_n;2014_p
241b	7.24	feces	2014_n;2014_p
540	7.24	CSF;mucus;serum	
179a	7.23	feces	2014_n;2014_p
406b	7.21	CSF;feces	2014_n;2014_p
1073	7.19	feces	2014_n;2014_p
211a	7.17	feces	2014_n;2014_p
589	7.12	feces	2014_n;2014_p
375a	7.09	feces	2014_n;2014_p
182a	7.05	feces	2014_n;2014_p
1423	7.01	CSF;serum	
1057	6.99	CSF;serum	
297b	6.99	feces	2014_n;2014_p
258	6.99	feces	2014_n;2014_p
241a	6.94	feces	2014_n;2014_p
956b	6.94	serum	
321b	6.93	feces	2014_n;2014_p
182b	6.9	feces	2014_n;2014_p

532	6.79 feces	2014_n;2014_p
787b	6.78 feces	2014_p
565b	6.74 CSF;serum	
1085	6.73 CSF;mucus;serum	2014_p
612a	6.56 feces	2014_n;2014_p
179b	6.54 feces	2014_n
113b	6.32 feces	2014_n;2014_p

**Supplementary Table S4** - Summary of hits of protein families to the NCBI environmental databases. (n) stands for env\_nt,(p) for env\_nr

Cluster	composite score	sample_origin	env_hits
1217	7.57	feces	marine sediment (p);gut (n)
339b	7.52	feces	gut (n)
211b	7.49	feces	gut (n)
562	7.39	CSF;serum	
297a	7.38	feces	gut (n)
457	7.32	CSF;serum	wastewater (n)
348	7.32	feces	gut (n)
502	7.3	feces	Human gut (n)
241b	7.24	feces	gut (n)
540	7.24	CSF;mucus;serum	wastewater (n)
179a	7.23	feces	gut (n)
406b	7.21	CSF;feces	gut (n)
1073	7.19	feces	gut (n)
211a	7.17	feces	gut (n)
589	7.12	feces	gut (n)
375a	7.09	feces	gut (n)
182a	7.05	feces	Human gut (n)
1423	7.01	CSF;serum	
1057	6.99	CSF;serum	
297b	6.99	feces	gut (n)
258	6.99	feces	gut (n)
241a	6.94	feces	gut (n)
956b	6.94	serum	
321b	6.93	feces	gut (n)
182b	6.9	feces	Human gut (n)

532	6.79 feces	Human gut (n)
787b	6.78 feces	Human gut (n)
565b	6.74 CSF;serum	
1085	6.73 CSF;mucus;serum	
612a	6.56 feces	gut (n)
179b	6.54 feces	gut (n)
113b	6.32 feces	Gut (n)

## Supplementary methods

### Protein family construction

We employed RNACode's evolutionary model to distinguish between coding and non-coding sequences. However, because RNACode depends on long, high-quality multiple alignments for its predictions, it has previously been impossible to use this method for gene prediction from short read metagenomics datasets. We have addressed this issue by adding several filtering steps. In short, we generated candidate protein families by building high-quality multiple sequence alignments from clusters of similar sequences. Subsequently, we calibrated the RNACode thresholds to detect coding sequences in short reads, and combined these with sequence-based measures like ORF length and sequence complexity and to select high-confidence candidate protein families. These steps are outlined in Figure 1.

### Clustering

To build the multiple sequence alignments, we began by clustering sequences by similarity with the MCL Algorithm<sup>3</sup>. This procedure yielded a total of 248,813 raw clusters, of which 13,861 had between 3 and 250 sequences (Supplementary Figure 1a) and were selected for building multiple sequence alignments. However, closer examination of these clusters revealed that many of them contained sequences with very high identity. Considering that RNACode has been reported to perform poorly on alignments with low sequence diversity<sup>4</sup>, we opted to sub-cluster sequences with more than 95% identity within each cluster, using a single-linkage approach, and replaced them with a single representative copy. We hypothesized that these highly similar sequences had likely originated from the same source, and thus could confound RNACode, which would consider them two independent proteins. Supplementary Figure 1b shows the cluster size distributions before and after sub-clustering. The increase in the number of small clusters, and specifically the large fraction (64%) of clusters were reduced to one or two sequences after sub-clustering, strengthens our hypothesis that most initial clusters included nearly identical sequences, thus validating our strategy. It is possible that some of the singletons are true ORFans, but this cannot be determined without additional homologous sequences.

### Multiple sequence alignment of clusters

It is well known that multiple sequence alignment quality has a large effect on bioinformatics tools that use them as input<sup>5</sup>. RNACode is no exception, and unfortunately, viral-enriched metagenomics data is close to a worst-case scenario for a multiple sequence alignment. Factors such as variable coverage, population variability, and sequence quality of next-generation sequencing (NGS) data could lead to clusters composed of sequences of different lengths, which in turn will result in alignments with a large number of gaps. Under these circumstances, RNACode could erroneously predict mostly-gapped regions to have high coding potential, due to the low effective number of triplets that can be scored (S. Washietl, personal communication). In order to overcome this problem, we used a codon-aware aligner, MACSE, which takes possible frameshifts and stop codons into account<sup>6</sup>, coupled with manual inspection and removal of low-coverage regions and poorly-aligning sequences. In the absence of an accurate automated solution, we found this to be the most effective way to ensure that RNACode would not erroneously predict coding potential due to poor alignment quality. To ensure high quality, we here focus on clusters of more than 8

sequences. 456 clusters were aligned, of which 209 low quality alignments (e.g. poor overlap) were discarded.

### **Selection of high-confidence protein families**

As RNACode was originally designed to work with longer alignments, we carried out a calibration procedure using simulated reads of varying lengths from coding and non-coding regions to optimize it for shorter alignments. The results of this calibration procedure suggest that the tool can accurately separate true ORFs from non-coding regions for alignment lengths as short as 50 nt, using a p-value threshold of 0.15 (Supplementary Figure 1c). Based on these findings, we applied the  $P \leq 0.15$  threshold to the remaining clusters, resulting in the exclusion of 154 families predicted to not have significant coding potential.

RNACode does not take sequence complexity into account when assessing coding potential, and thus may yield low p-values for short tandem amino-acid repeats. Since such repeats often appear in non-coding regions, we applied the *seg* algorithm<sup>7</sup> to identify and exclude low-complexity candidates, complementing RNACode's predictions. In total, an additional 61 clusters were excluded based on low complexity.

Finally, to rank our resulting families in terms of coding potential, we calculated a composite score that combines the RNACode score, sequence complexity, and ORF length. Although *ad hoc*, in our experience the composite score is a convenient heuristic that was more useful than simply evaluating the p-value alone, since it integrates external, biologically relevant information that RNACode does not consider.

From the original set of 456 clusters selected for further analysis, our approach identified 32 high-confidence candidate novel ORFan protein families, which we ranked by our ad-hoc composite score. Their final trimmed and filtered alignments were composed of a median of 6 sequences and have a predicted ORF length between 30 and 130 amino acids (Figures 3a and 3b). The coding potential was assessed using RNACode as described in the previous sections, and the RNACode p-value distribution for the 32 clusters ranged from very confident predictions with p-values around  $10^{-10}$  to values close to 0.02 (Figure 2a). The RNACode output for cluster 457 is shown in Figure 3d, and RNACode predictions for all 32 families are available in Supplementary Fig. S2.

### **Phylogenetic tree construction**

The ORF6 protein sequence and all 1657 sequences in the Pfam full alignment of Phage\_F (PF02305) were aligned to the Pfam HMM with *hmmalign*. Unalignable regions in the N- and C-terminus were removed, keeping a conserved core region of about 150 residues. Sequences with >50% gap residues were removed, and sequences >90% identical to other sequences were removed, leaving 331 representative and non-redundant sequences. The tree was built with the Neighbor Joining method using Scoredist distance correction using *Belvu*<sup>8</sup>, with 1000 bootstraps



## References

1. Darling, A. C. E., Mau, B., Blattner, F. R. & Perna, N. T. Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome Res.* **14**, 1394–403 (2004).
2. McNair, K., Bailey, B. A. & Edwards, R. A. PHACTS, a computational approach to classifying the lifestyle of phages. *Bioinformatics* **28**, 614–618 (2012).
3. Enright, A. J., Van Dongen, S. & Ouzounis, C. A. An efficient algorithm for large-scale detection of protein families. *Nucleic Acids Res.* **30**, 1575–1584 (2002).
4. Washietl, S. *et al.* RNAcode: robust discrimination of coding and noncoding regions in comparative sequence data. *RNA* **17**, 578–94 (2011).
5. Blackburne, B. P. & Whelan, S. Measuring the distance between multiple sequence alignments. *Bioinformatics* **28**, 495–502 (2012).
6. Ranwez, V., Harispe, S., Delsuc, F. & Douzery, E. J. P. MACSE: Multiple alignment of coding SEquences accounting for frameshifts and stop codons. *PLoS One* **6**, (2011).
7. Wootton, J. C. & Federhen, S. Analysis of compositionally biased regions in sequence databases. *Methods Enzymol.* **266**, 554–571 (1996).
8. Sonnhammer, E. L. & Hollich, V. Scoredist: a simple and robust protein sequence distance estimator. *BMC Bioinformatics* **6**, 108 (2005).