

Supplementary Figures and Notes

Discovery of an expansive bacteriophage family that includes the most abundant viruses from the human gut

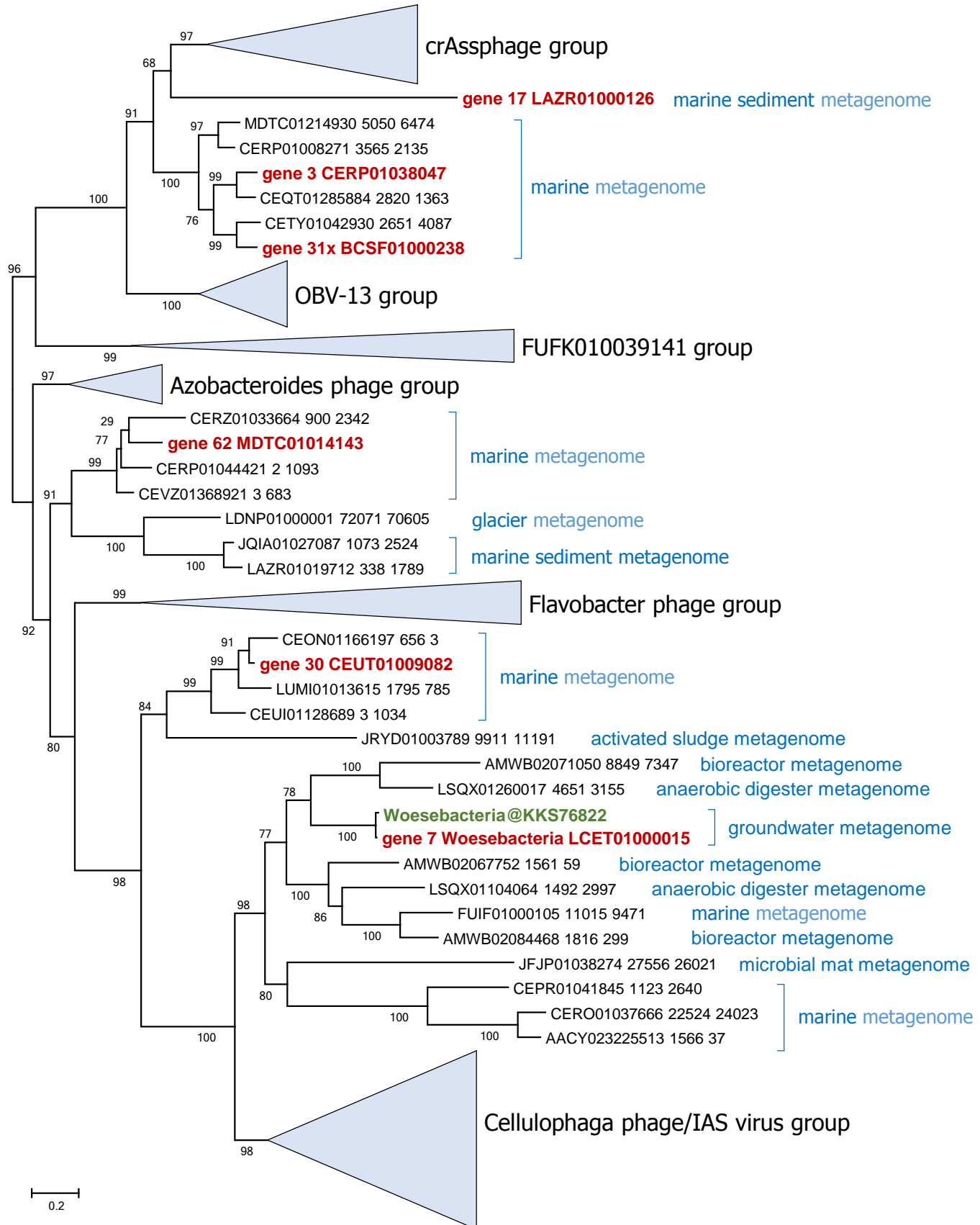
Natalya Yutin¹, Kira S. Makarova¹, Ayal B. Gussow¹, Mart Krupovic², Anca Segall^{1,3},
Robert A. Edwards³, Eugene V. Koonin^{1,*}

1, National Center for Biotechnology Information, National Library of Medicine, Bethesda, Maryland 20894, USA;

2, Institut Pasteur, Unité Biologie Moléculaire du Gène chez les Extrêmophiles, 25 rue du Docteur Roux, 75015 Paris, France;

3, Viral Information Institute, Department of Biology, San Diego State University, San Diego, California 92182, USA

Supplementary Figure 1

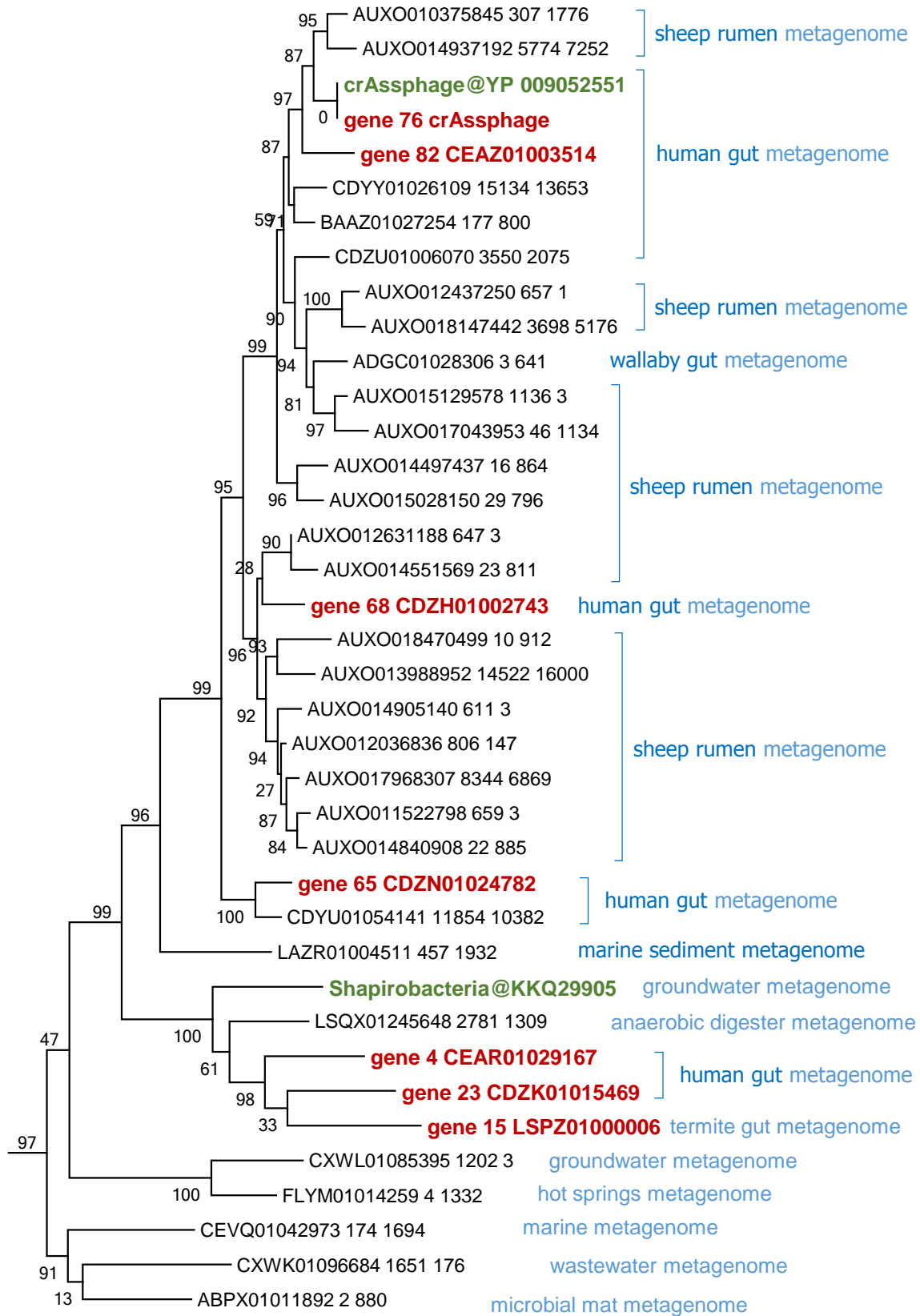


Phylogenetic tree of the MCP for all identified members of the crAss-like family.

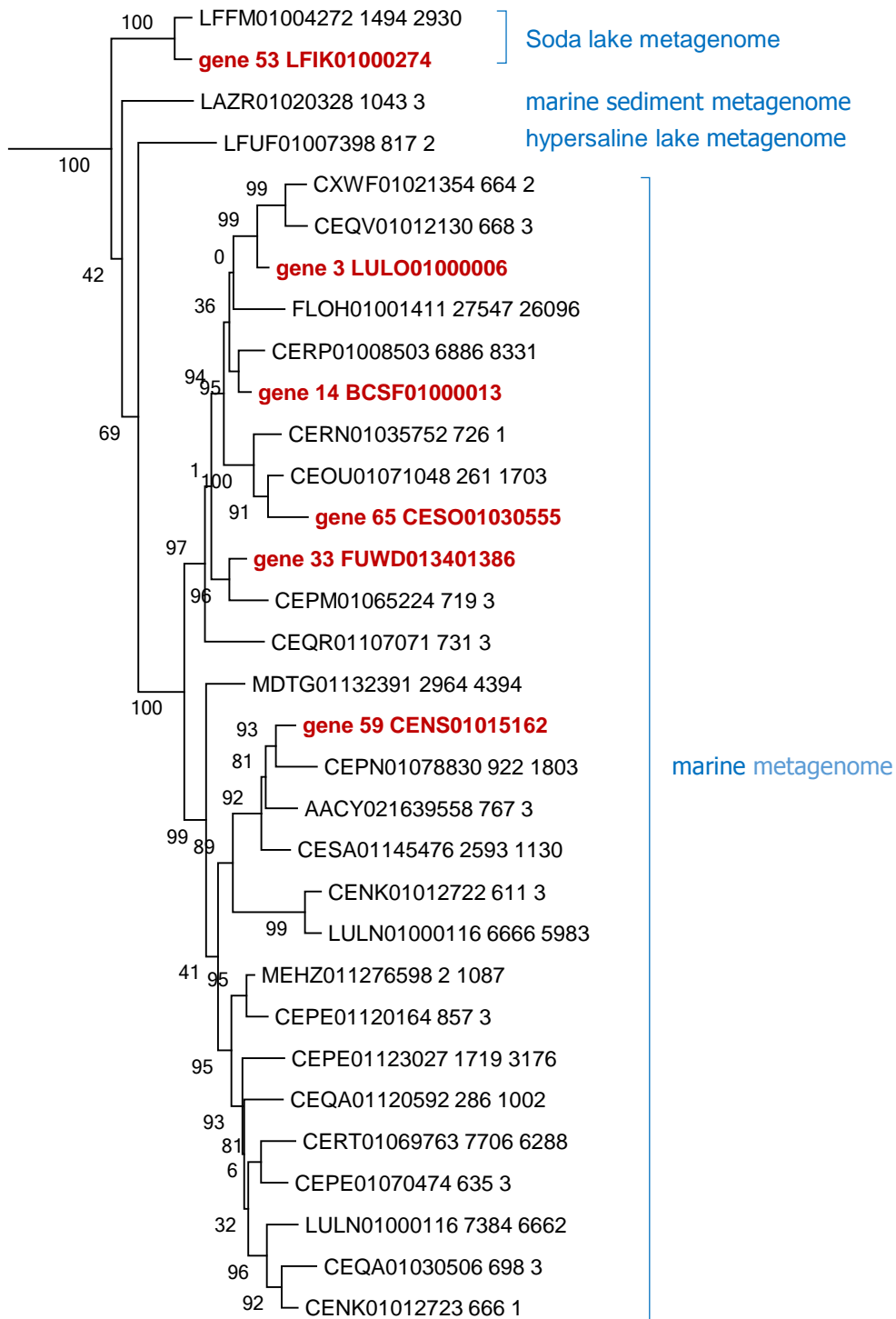
Translated wgs sequences denoted by three numbers: contig ID, orf start and orf end coordinates.

Representative sequences are shown in red, nr proteins (denoted by their source and protein ID) are shown in green. Branches corresponding to large groups are collapsed into triangles. The next 6 panels show the collapsed branches expanded. The tree was constructed using FastTree as described under Methods. Support values were obtained using 100 bootstrap replications; values greater than 50% are shown.

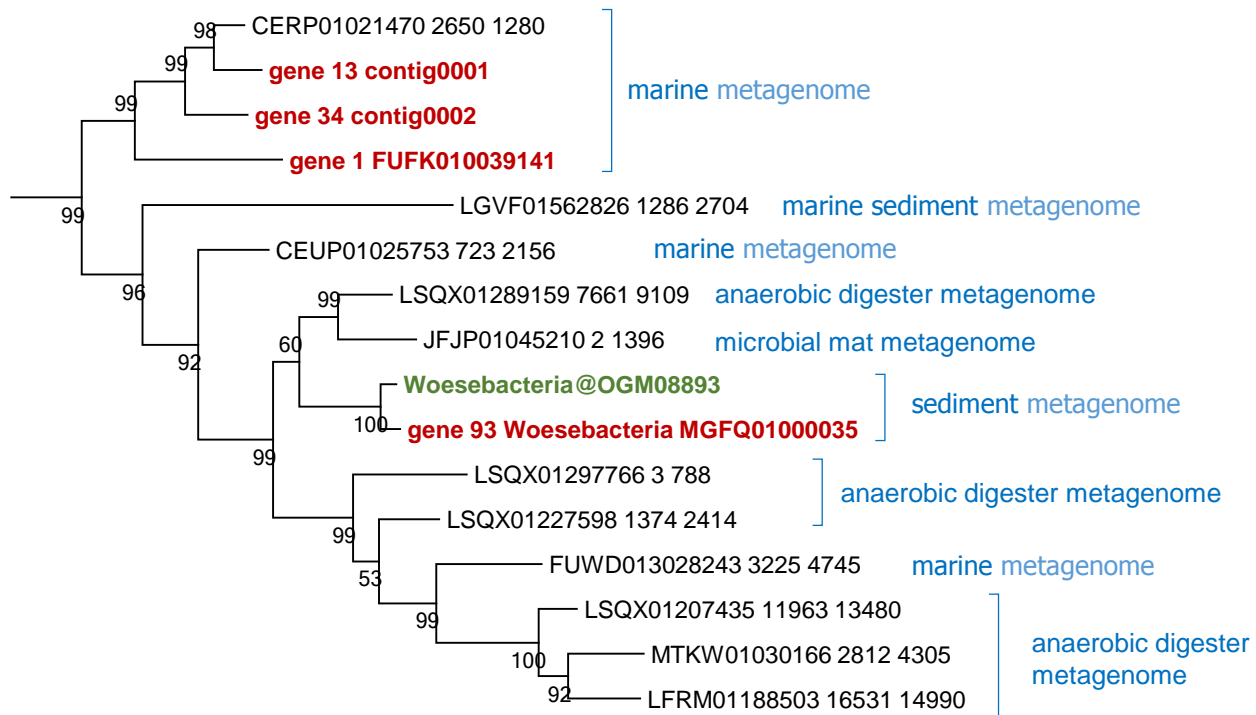
CrAssphage group



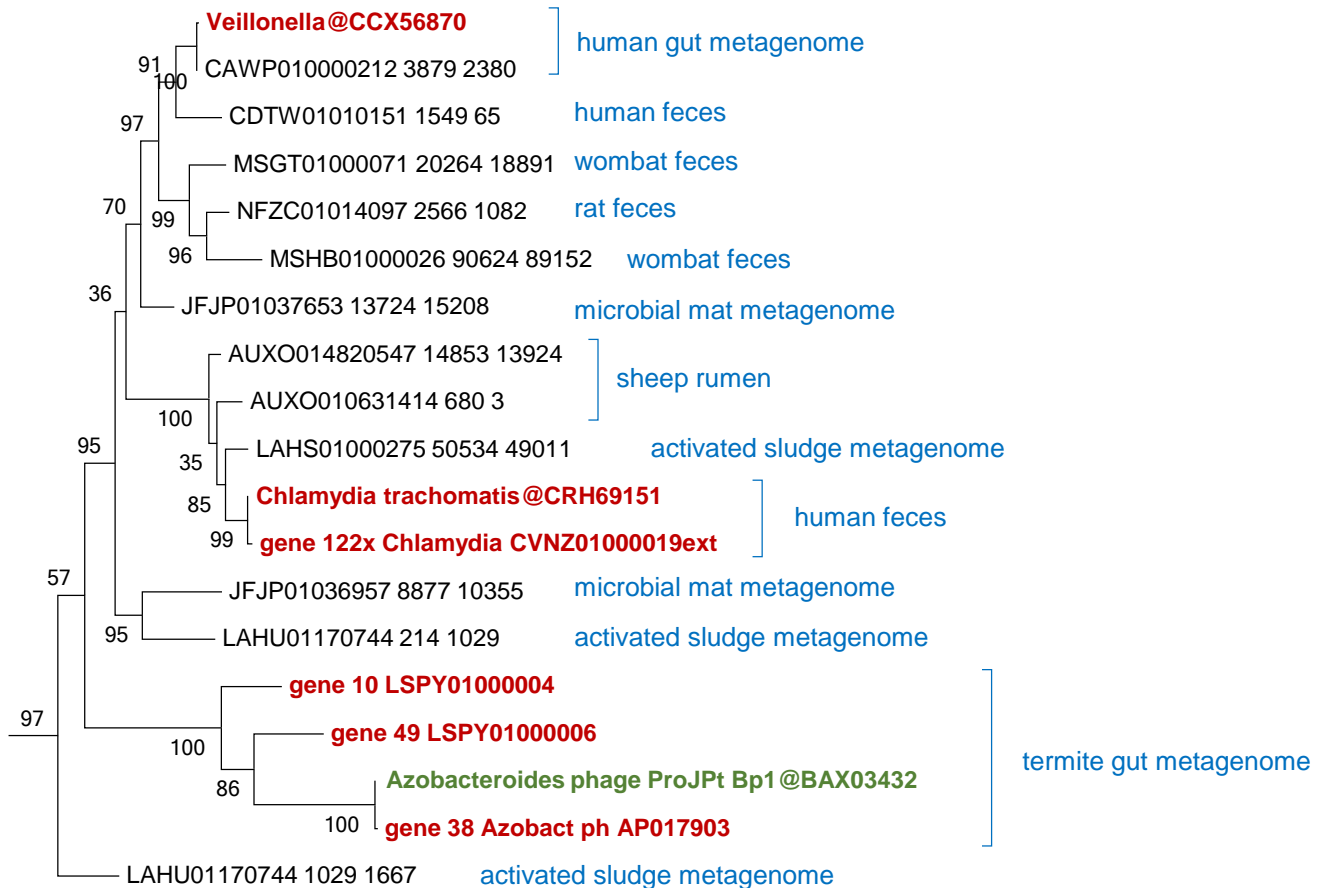
OBV-13 group



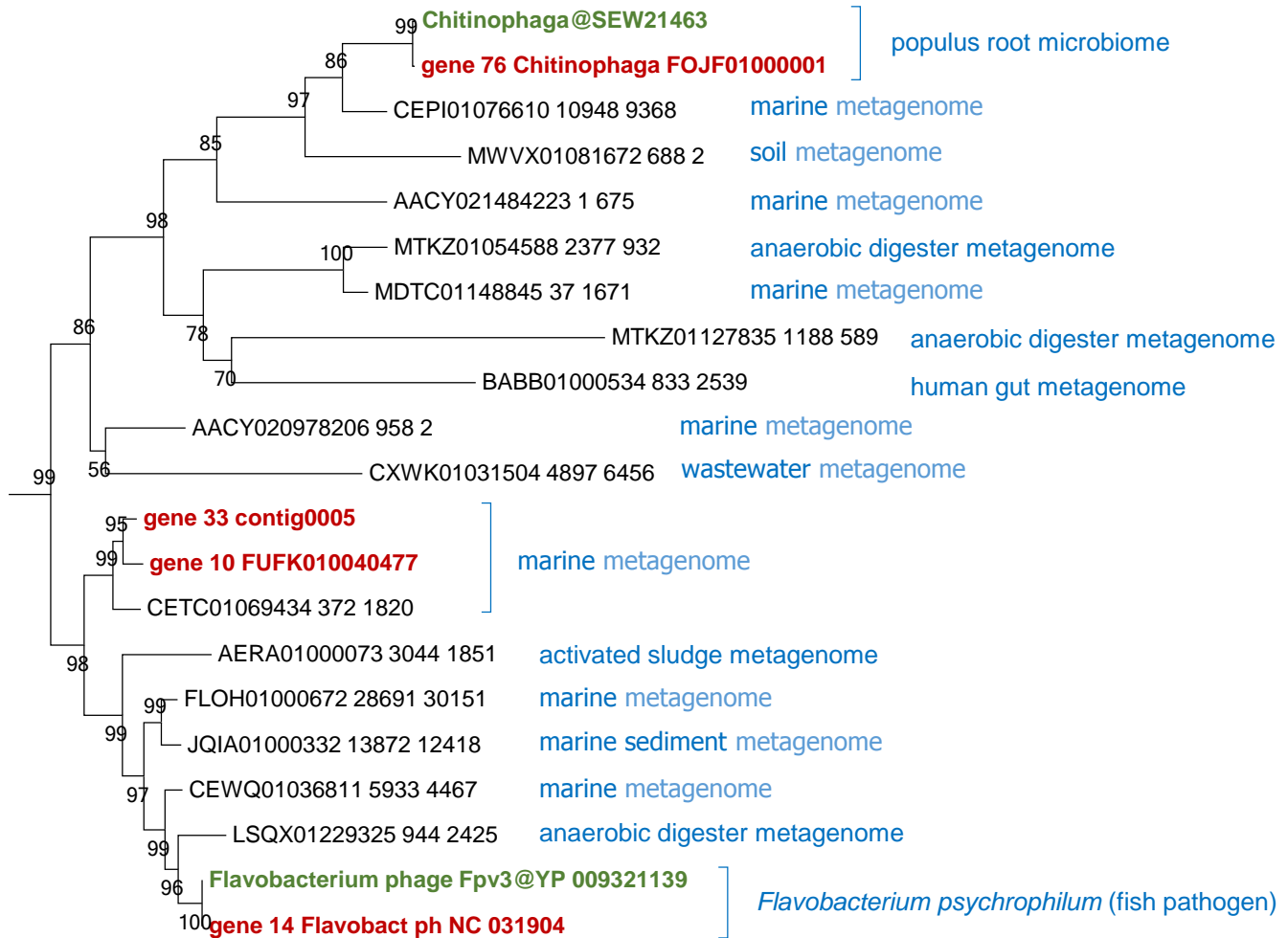
FUFK010039141 group



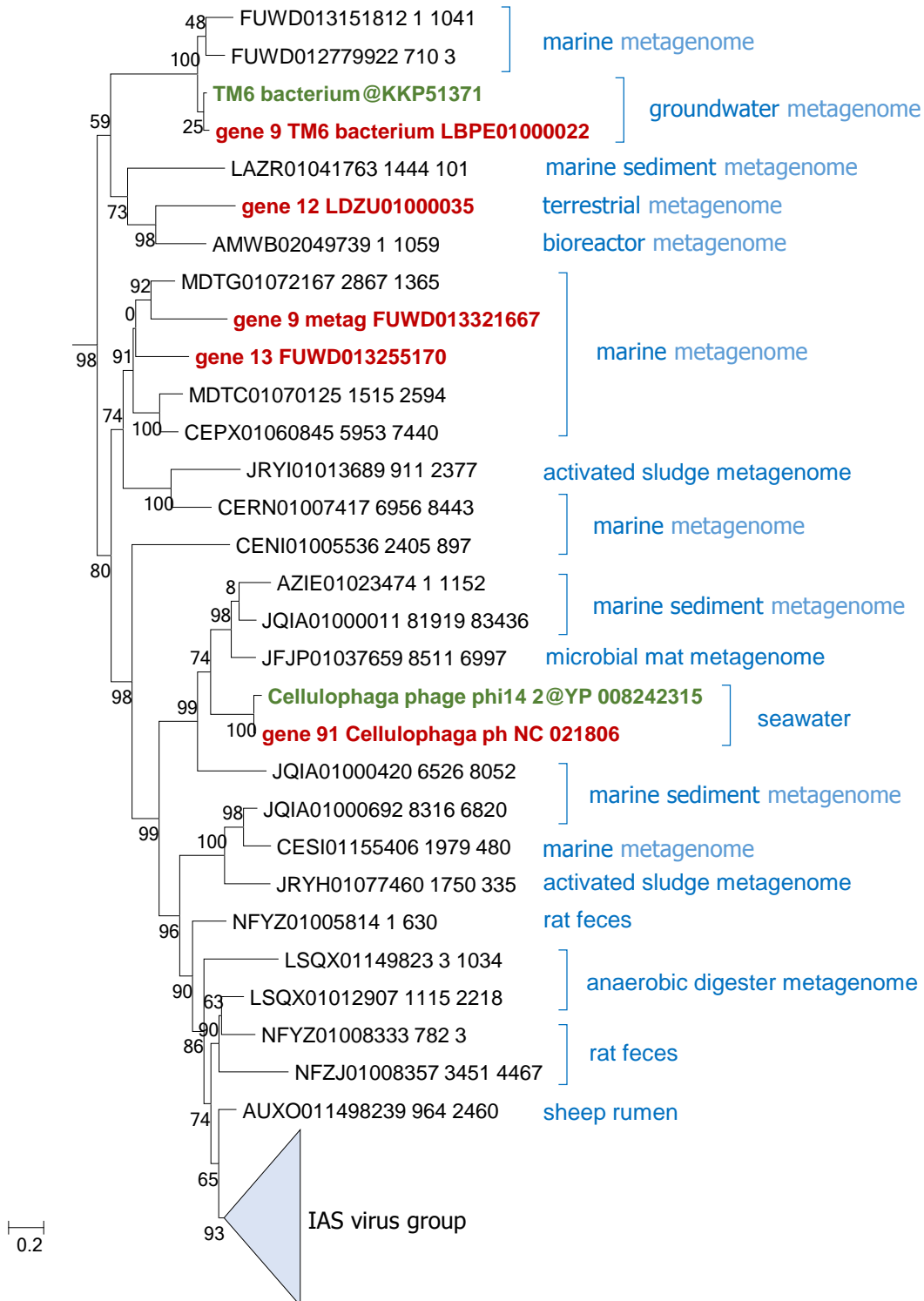
Azobacteroides phage group



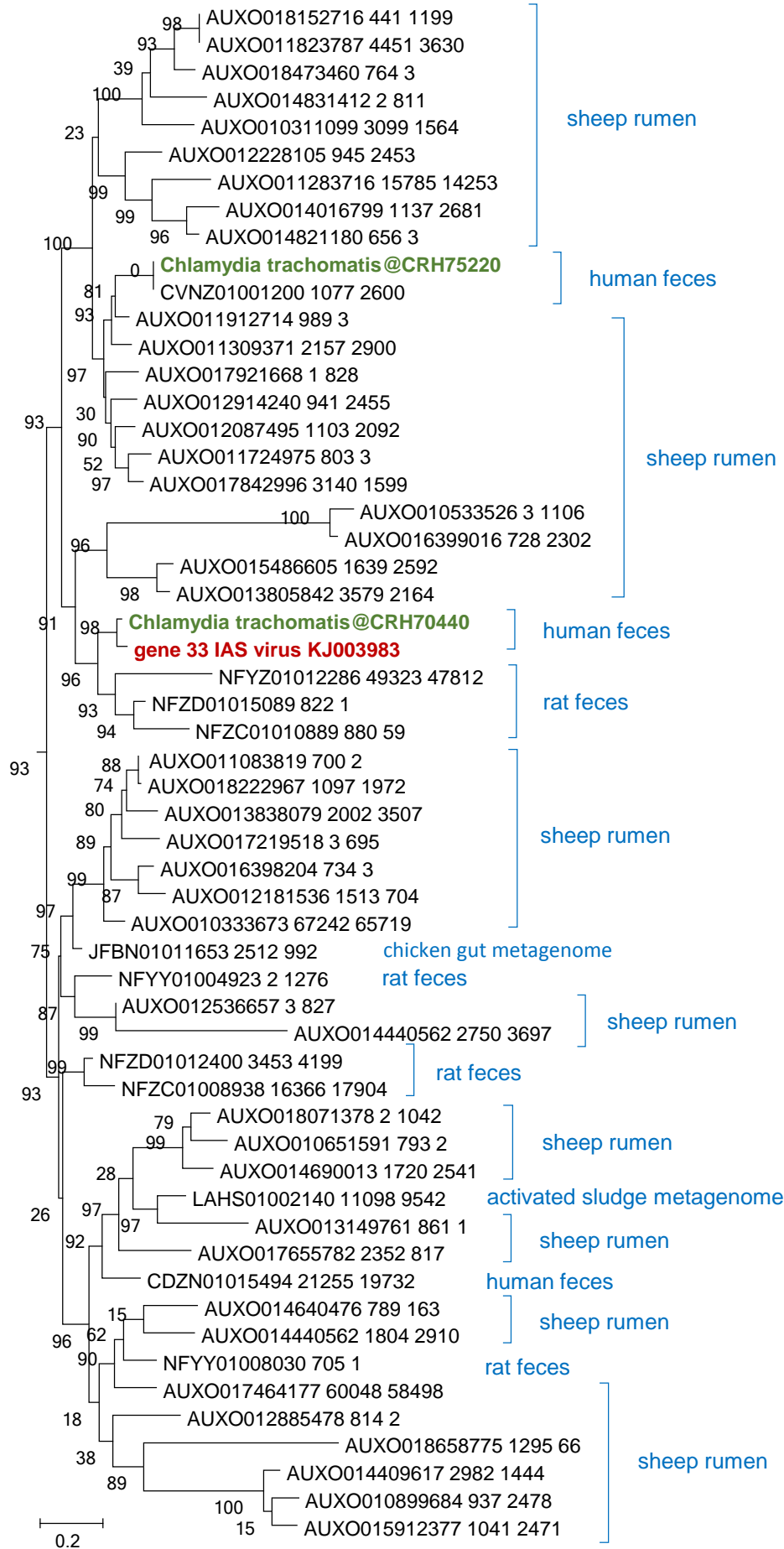
Flavobacter phage group



Cellulophaga phage/IAS virus group

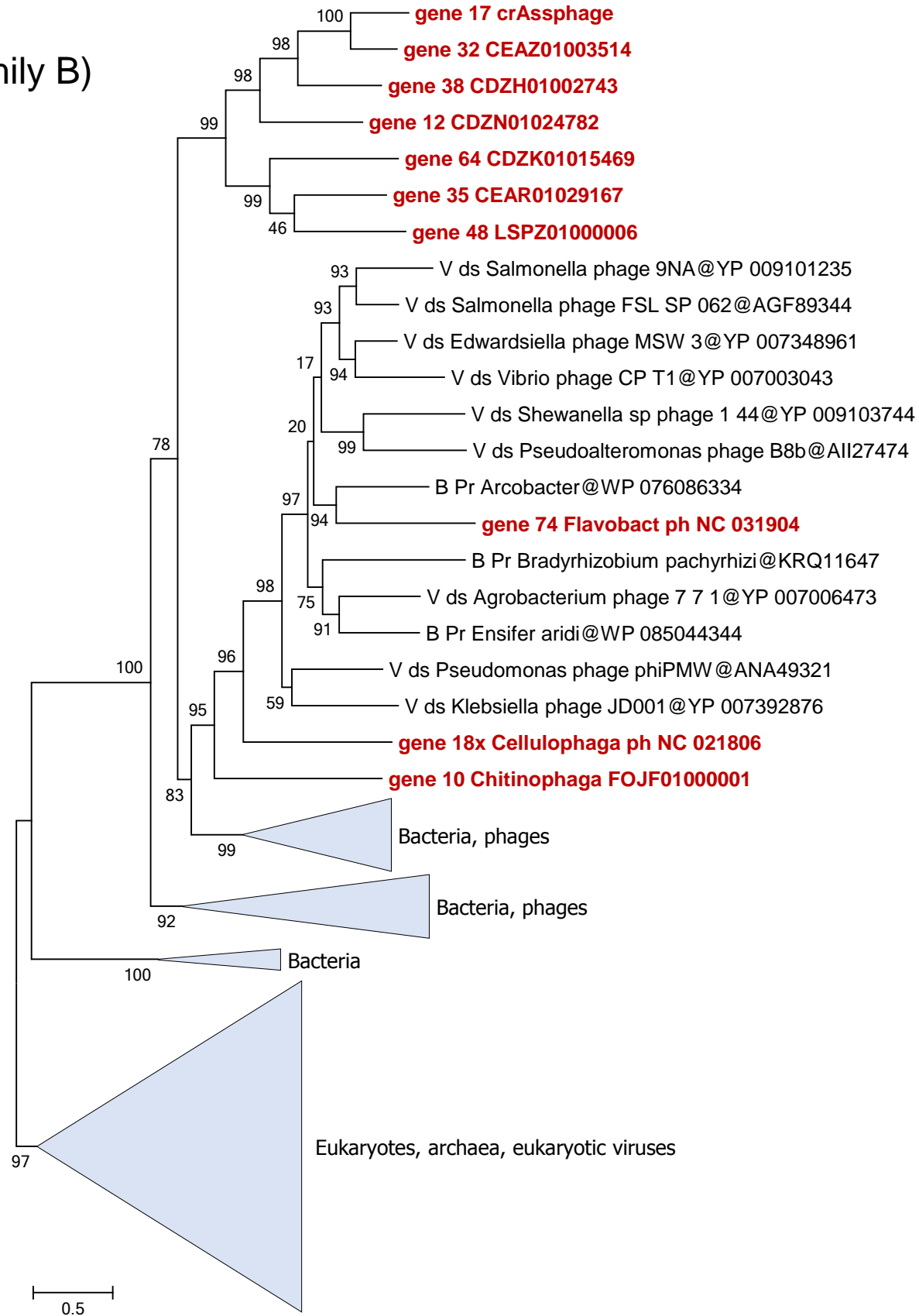


IAS virus group



Supplementary Figure 2

DNAp (family B)

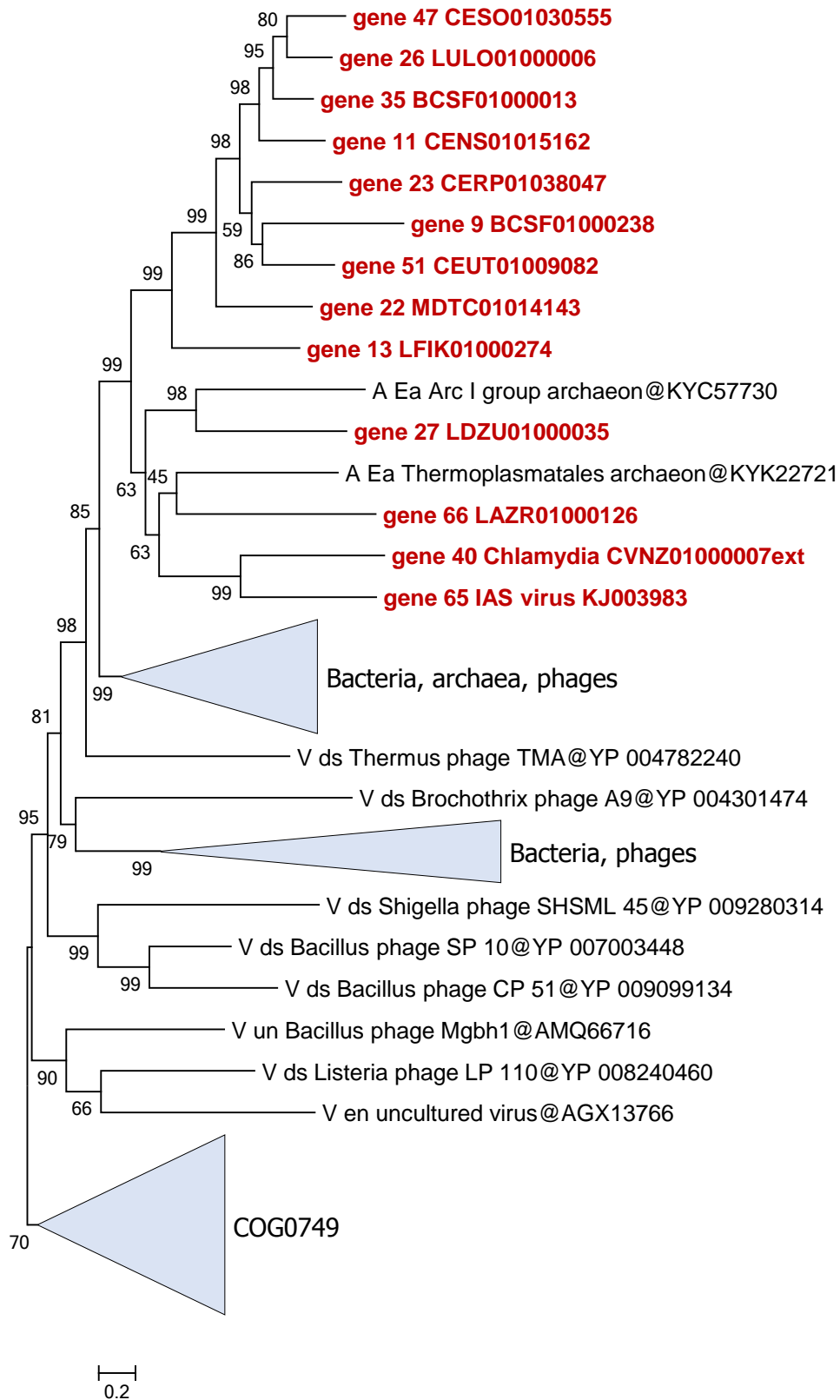


Phylogenetic trees for crAss-like family PoIA, PoIB, primase, and ligase.

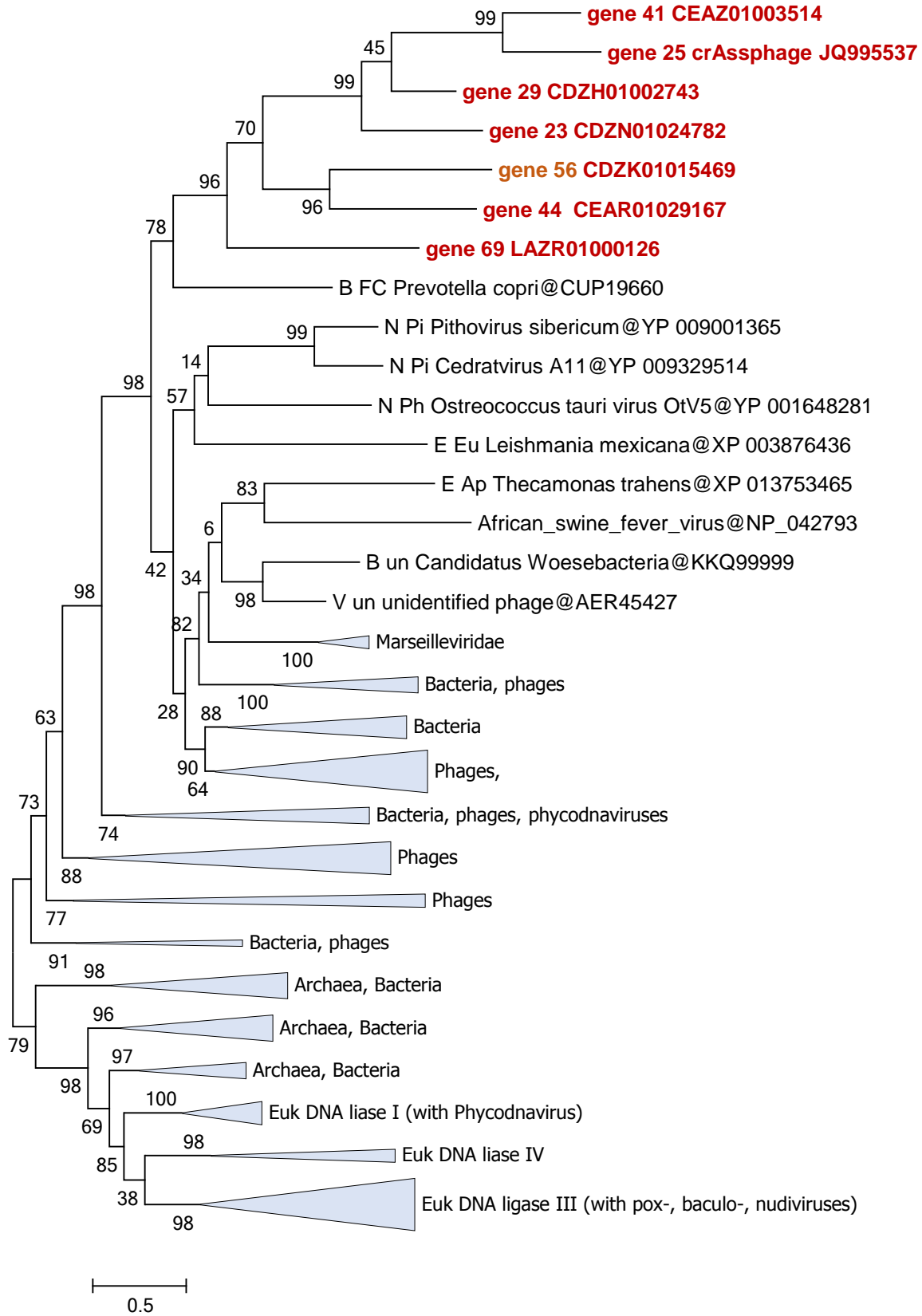
Translated wgs sequences denoted by three numbers: contig ID, orf start and orf end coordinates.

Representative sequences are shown in red, nr proteins (denoted by their source and protein ID) are shown in green. Branches corresponding to large groups are collapsed into triangles. The next 6 panels show the collapsed branches expanded. The tree was constructed using FastTree as described under Methods. Support values were obtained using 100 bootstrap replications; values greater than 50% are shown.

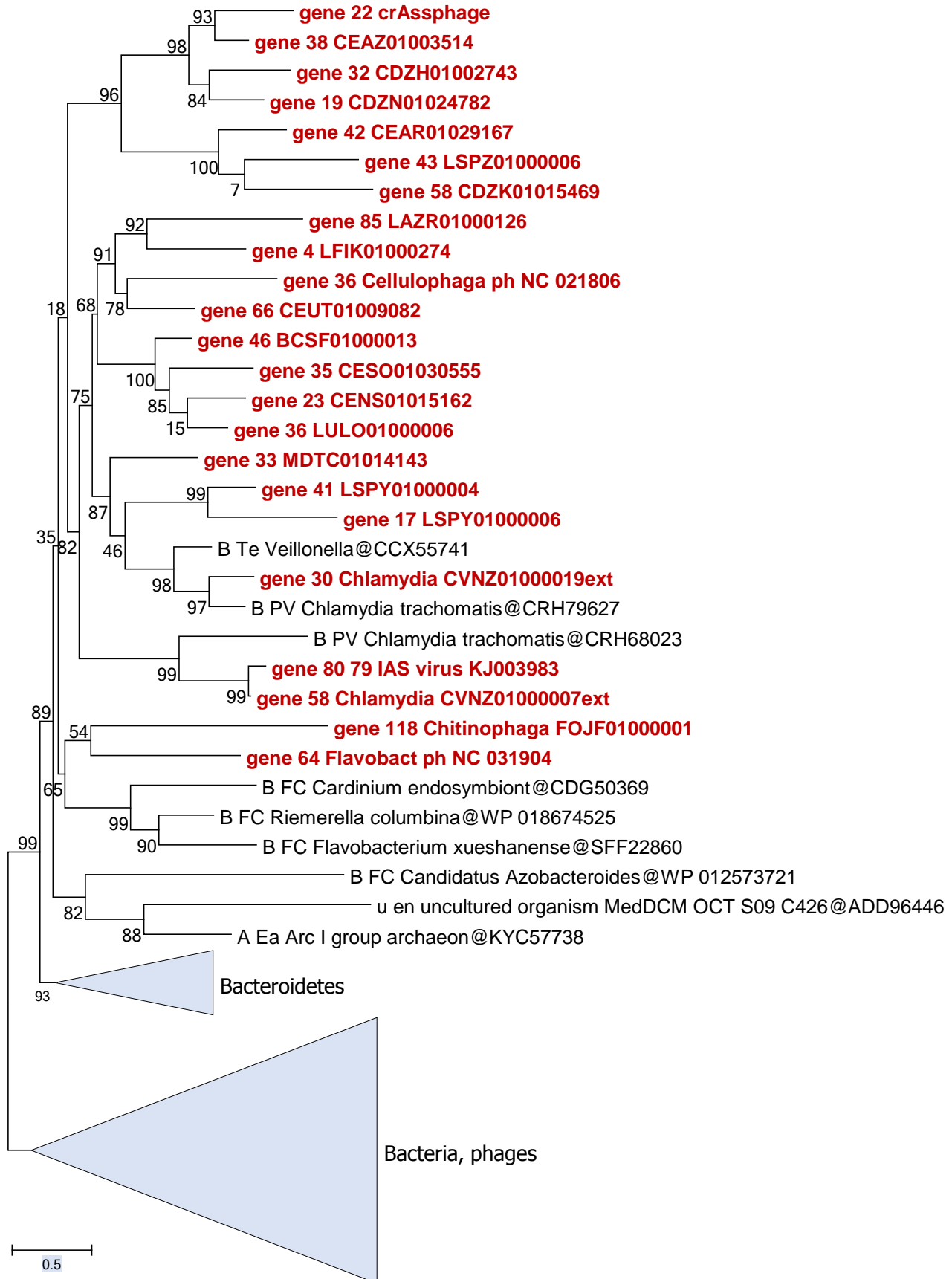
DNAp (family A)



ATP-dependent DNA ligase



Primase



Supplementary Figure 3

Multiple alignment of conserved regions of β and β' RNA polymerase subunits of the crAss-like family.

The most conserved motifs of β and β' subunits are highlighted by reverse shading. High confidence alignment is highlighted by green along the Consensus line. HHpred output demonstrating similarity with catalytic region of β' subunit is shown below.

```
gene_82_Woesebacteria_MGFQ01000035 IYKLSDGGNGLNPLWEYLFYQSLG-----GKYGP-----IDKTDVKSLYN-----
gene_22_contig0002 LQSEEDLRIVYLKEESFPVEFFGDILDNTGDSFA-----MNNLSFQRMVKNQG----
gene_52_Azobact_ph_AP017903 DSMPTDQGSFMT PAMAREME IRL-----GRWSE-----VKEDIYKQIKK-----
gene_25_842252479 EMEIFDGGGLVHIDFYRLARRLT-----SSWSK-----QDEAVYNKLLID-----
gene_45_LULO01000006 NIEEADANAFANIHFYRELMMRI-----GDWTE-----QQENLFRKSML-----
gene_30_Chitinophaga_FOJF01000001 NANEADAQSWITLPAYRELLLS-----GGRWTDAMERQYQYEMAYERLGRSKTNGFS
MODO_RS15645 ISNIADAQTFASPAFYKKWITSI-----KNGWTS-----HDEAMYNKLMDPN-----
gene_76_CEUT01009082 NVDLTDATWITPELQKQRMKGF-----GEFTP-----EVEAAYNRMM-----
gene_1_contig0005 SVNTTDAQAYITLERWKDIVKGL-----RKWNN-----KLEDSYNRLLV-----
gene_66_Cellulophaga_ph_NC_021806 DVNQTDQAWITPKRWAFILSRT-----GKWNS-----KYQSVYKILK-----
gene_34_Flavobact_ph_NC_031904 DIESTDAQEYSTATEHISILHRI-----GRISDE-----EFKVIITNKLLTEGE----
gene_64_Chlamydia_CVNZ01000007 KNTLTDGGQYRTRLSYRKMVGMMA-----GQWTE-----EMQKVYEAIEQLKADYNG
gene_49_KJ003983 GINATDQGSFRT PASMKKILEAM-----GGKWTE-----DMQEAYDRMTR-----
gene_34_775896848 DADHTDAQAFISPRMFQYIEQGL-----ARWTP-----EDDVWFEKYMS-----
gene_80_FUFK010017431 DVVKTDAQSFISIDMYRGIQQGL-----GEWGS-----RDEQAYLNETSIDPTNPN
gene_57_BCSF01000013 STVKTDAQAYITVSMYRDIQGR-----GQWTE-----RDQAYDNEMA-----
gene_107_816813760 GANKTDSMGFISLDKHRHKMEGQ-----GTWLD-----EHEAAYKNYQA-----
gene_40_MDTC01014143 KINKTDAQGFITLFEHKNLMESD-----GTWTD-----EHTDAYNKYVW-----
gene_44_935075730 TITTAADFNVITQDECIRRFKAM-----GDYDS-----FTLPSGRTLADIVADE--
gene_35_1001888980 SITTTDGISITIDVHLERYLRAS-----RGLSP-----AQEQMIKNLRE-----
gene_57_934554352 TATINDSQSIMTDIGLIKLLKAT-----GRWNP-----SEPLYQYITDLQDP---
gene_44_933950148 KSEVTDQSFITLDEFVRIYLR-----GEYDN-----YKDLIEALYDE-----
gene_10_936104480 NTTVNDQASYITFEWIRRVAGR-----GQLNK-----YIPLIERIMDR-----
gene_59_935734019 KEKVNDAQSYITFDEFIRRVYIA-----GEYNK-----YKNVIEALLS-----
gene_43_932717476 GVNANDQASYITLEEWIRRITAA-----GELDK-----YAGLIQSLTD-----
gene_46_crAssphage_JQ995537 GVNTNDAQSYITLEEWIRRITAA-----GELDK-----YAGLIQSLTD-----
CONSENSUS 0.8 -----D-----G-----
2BE5 Thermus thermophilus (beta) RDLDEEGVVRIGAEVKPGDILV.....
2BE5 Thermus thermophilus (beta') -----
Jpred4 sec.str. gene_46_crAssphage -----HHHHHHHHHH-----H-----HHHHHHHHHH-----
```

```
gene_82_Woesebacteria_MGFQ01000035 ----DVNLTSGKKILHKEAVD-----FITGDT-YNMG-----IPMRDMFK
gene_22_contig0002 ----GKIKKYGLNPKPVAFQLDS-----EGDLFFWK
gene_52_Azobact_ph_AP017903 ----GEDVGE-WGHRPLFTVMKC-----AYWGDE-MLGG-----KTGTKIAK
gene_25_842252479 ----GVEITP--EEVALFDPLKP-----QVLAQA-FEDK-----IDLRLFNK
gene_45_LULO01000006 ----GEELTK--EELTAFPTLKP-----QGFGPA-ISE-----LKQMVGLK
gene_30_Chitinophaga_FOJF01000001 YENRGDLRKHDEDLVKKGPNQSGI-----FHVIKP-IGTGVKAGAAAFADIFLDK
MODO_RS15645 ----YKIDKEVKEWLKVTGNFLAGTEKPFYGLLDSKYN-YKEGDDAFF-NNTPLMLK
gene_76_CEUT01009082 ----KSTEMQV-SKILKLTPRKS-----AARGTE-IKNG-----MNVPIREK
gene_1_contig0005 ----GKGTDS-DMLAVVAQPLKGV-----YYSLRP-VKIGNVS----LNIPTYLK
gene_66_Cellulophaga_ph_NC_021806 ----SESLDA-SEMKLAAQPLKG-----VYFGLV-----NNTPTYLK
gene_34_Flavobact_ph_NC_031904 ----KGYLTK-EELKLVFQPIKP-----VHTGTY-TNVNQD----VNRVVYIK
gene_64_Chlamydia_CVNZ01000007 KSIPSEKLAEIASMAVVFQPIKPYMF-----THEKLA-LNDKDK----VIIPVQHK
gene_49_KJ003983 ----EQRLDM-QDFYTIWNNIKPF-----VYSHES-VKIGDRV--EKVPVEHK
gene_34_775896848 ----GE-----GEWAAPFTPAYKF-----YAEQML-IDNG-----TLIADMOK
gene_80_FUFK010017431 A--GKYVDN-AGVPVVIKPIKP-----YHEELK-VRDG-----KPVLHMDK
gene_57_BCSF01000013 ----GRGFIDNEGNRPRIYPLKP-----YHEELT-VRNG-----VAELHMDK
gene_107_816813760 ----GGEFVTPDGIRPLLNPIKT-----IYDGFH-LDVDN----RVVRVVDK
gene_40_MDTC01014143 ----KGLMGDPASKKLLLNPRKT-----YYFGER-LITDSQGNQ-TITFEQIK
gene_44_935075730 ----DTPISP-SDYARIVEQLKYY-----FYKRGK-STLNNRFNTDIVFSHQDK
gene_35_1001888980 ----QKNLDP-REVHTVFEQIKTY-----AYHKS-VINDGTVA---YPKTVFK
gene_57_934554352 ----TKTFNP-TSYAKVVEQVKLFGTTRRRRGDFYHAPAVIGNVEDIFADEVSVQIK
gene_44_933950148 ----SKPIDN-YKLGELSKIQVQKN-----FYDLE-IDNDAK-----LANPIQIK
gene_10_936104480 ----SKPLRV-DDIKTFVQVQKN-----FYDYM-TYNDKIN----TYAPRQIK
gene_59_935734019 ----DKPLEE-IDFDTL-NKIQVQKN-----FYDLY-YDSARN----REVPLQIK
gene_43_932717476 ----DTPIDK-IDWTKFANKVQIQKN-----FYDLY-YDFTVG----IEVPRQVK
gene_46_crAssphage_JQ995537 ----DTPIDK-IDWTKFANKVQIQKN-----FYDLY-YDFTVG----IEVSRQVK
CONSENSUS 0.8 -----K-----
2BE5 Thermus thermophilus (beta) .....LANRHGK
2BE5 Thermus thermophilus (beta') -----
Jpred4 sec.str. gene_46_crAssphage -----HHHHHHHHHH-----EEEE-----EE-----EEEEH---
```

gene_82_Woesebacteria_MGFQ01000035
 gene_22_contig0002
 gene_52_Azobact_ph_AP017903
 gene_25_842252479
 gene_45_LULO01000006
 gene_30_Chitinophaga_FOJF01000001
 MODO_RS15645
 gene_76_CEUT01009082
 gene_1_contig0005
 gene_66_Cellulophaga_ph_NC_021806
 gene_34_Flavobact_ph_NC_031904
 gene_64_Chlamydia_CVNZ01000007
 gene_49_KJ003983
 gene_34_775896848
 gene_80_FUFK010017431
 gene_57_BCSF01000013
 gene_107_816813760
 gene_40_MDTC01014143
 gene_44_935075730
 gene_35_1001888980
 gene_57_934554352
 gene_44_933950148
 gene_10_936104480
 gene_59_935734019
 gene_43_932717476
 gene_46_crAssphage_JQ995537
 CONSENSUS 0.8
 2BE5 Thermus thermophilus (beta)
 2BE5 Thermus thermophilus (beta')
 Jpred4 sec.str. gene_46_crAssphage

TMVDHTSMLMSKEL-----GRK--IDVFATWKMFYDY-----TGRNDFAKAIR
 TSSTETITDLRNT-----SKA--YAFFGQLMDAFESKLGSE-T---KIVFASFESAMK
 TSYKTLLPD-----EYPT--LQPLVKWMDKND-----IGVLHFGSAIK
 FALFPIHPGLSKTVTDINDGNANVLDEMYADMNKHN-----LDYIVMESAVK
 LSLTPIFPQMTRI-----NGQDTALKGLLKDMYESG-----TDMVMHPTAAK
 TSAMPIYYRMS-----EGRG--LGKLYGKMSAAG-----ISYAIMESGRK
 TSFYFLFPALM-----SGTR--GNEILKAMEGQG-----VELVSFESAVK
 TAEVVLWPGL-V-----KDAG--LKSLYDAMVQEEKDYKMTGTGATLGMSVSVESAIAK
 YSQAVLFPQVAN-----ENSE--LRSIYNSMVDQG-----VDELVFDGSIK
 YSQAVLLPQLV-----AGTQ--LQSLADAMNKQD-----IGESIVLDGVK
 SSAFPILIPQLT-----AGTK--LDALRLKMELETNTGR-----FTRASFQTANK
 YAEAVLIPPELLP-----EGSK--LRDMAYWMDNN-----VDLVHSTKIVK
 NSEYMMTAFYNLIGSVLSKSPQ--LVALQEFMQDND-----IDVVFESAVK
 NSYVVLTKELA-----EGNE--LLTAMYNRMMDPNP-----IDIINTESAKK
 NSYVTVTKELSA-----EYPV--LEKLRRVMQDQD-----IHVIHEENANK
 NSYTVLTSELTA-----NYPV--ADKLRLAMLRNN-----IGVANEEGATK
 HVVFPLLREFTK-----ASPA--FEKLRAMEEQG-----IDEVNMESARK
 HSTIPLIREFTEMFKGTEEQVS--LNNLRQRMEAVGKYKGLQP----IDMANFESALK
 NSTLVIKFRMY-----KGTG--YETLYDWMKQEG-----IDSINFESGHK
 DSTIILFPQY-V-----KGTD--MERYNWMNSNK-----IDQLSFDSTHK
 DSTVVLFEFMT-----RGS--MGQLYDWMIQNN-----IDQVSPISAVK
 NAEFVLIPRFL-----GNSE--LGLLAKYMTDNN-----IGQVNFTTTEK
 NAEVLVLPRI-----EGTD--LEKVYNLMKDHG-----IDQLNTEETSK
 NAEYVLIPKLI-----KGTE--FEAIYNAMKAAG-----IDQLNVTETSK
 NAEFVLIPKLI-----KGTE--LEKVYNIMTNRG-----IHQINTVETVK
 NAEFVLIPKLI-----KGTE--LEKVYNIMTNRG-----IHQINTVETVK
 -----L-----M-----K
 GVVAKILPVE

 -----EEHHHHH-----HHHHHHHHH-----EEEEEEEE

gene_82_Woesebacteria_MGFQ01000035
 gene_22_contig0002
 gene_52_Azobact_ph_AP017903
 gene_25_842252479
 gene_45_LULO01000006
 gene_30_Chitinophaga_FOJF01000001
 MODO_RS15645
 gene_76_CEUT01009082
 gene_1_contig0005
 gene_66_Cellulophaga_ph_NC_021806
 gene_34_Flavobact_ph_NC_031904
 gene_64_Chlamydia_CVNZ01000007
 gene_49_KJ003983
 gene_34_775896848
 gene_80_FUFK010017431
 gene_57_BCSF01000013
 gene_107_816813760
 gene_40_MDTC01014143
 gene_44_935075730
 gene_35_1001888980
 gene_57_934554352
 gene_44_933950148
 gene_10_936104480
 gene_59_935734019
 gene_43_932717476
 gene_46_crAssphage_JQ995537
 CONSENSUS 0.8
 2BE5 Thermus thermophilus (beta)
 2BE5 Thermus thermophilus (beta')
 Jpred4 sec.str. gene 46 crAssphage

QTVEYYIDDEMRDS-FKSNLVWKLNRSS-----
 GSKKY-----LDSEGAYLDW-----
 SWNTD-----RAEL-FDNNGIFLG-----
 LGAKKEKNTKFKSKL-IDDVGSYQP-----
 IGSRVDLVTGNAPQF-YNEDGT-----
 VGAEF-----LNPV-YQTNGDFN-----
 TNSVT-----PTNL-VSVNQFNEYED-----
 VGAGT-----KTKI-DDDQGN-----
 EGALS-----PVSL-YSDAESTVLKT-----
 VGATT-----PNIV-TDENGDI-----
 VGATI-----KTVKPFIDINSLST-----
 VGGFG-----STDISKASDAKSLGEA-----
 VGGYG-----QVNLNYDEERFEKLRDSNGNIKVGNEPIKATDYDTFMKNLRQAIIVDE
 GFKGQ-----LFQVDQSAEGAM-----
 GARKD-----VRDA-YTDN-----
 GATGI-----VLDL-YNDS-----
 VG-----LDLSG-----
 VGASG-----IENI-SN-----
 VGGMP-----KVQL-FDISRDIAVDSKGFPILDANGKYTYTNGTKATLNIQYNEATK
 IGGTY-----PGSI-YDSEGNFIGDK-----
 VSGIT-----PVKI-HNDQGLDIEA-----
 ATTNR-----VLEF-WDAHGKFPSEK-----
 AGKAG-----VLTTL-FDEETGEVTDH-----
 AAKNR-----LIEL-WDSSTGELTDK-----
 VAQHN-----RMTL-WNNDGVLTDIA-----
 VAQHN-----RMTL-WNNDGVLTDIA-----

 E-----EEE-E-----HH-----

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDTC01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene_46_crAssphage

---VKGAIKINIYNPWNDIYDPLNIKEGKRGLDLLSDDI-----RN
---KDTYNKSFESSETDIDGV-----
---TDDLNDNYVST-----
---LQNLTSVQ-----
---INPIDS DL-----H
---PEPFSGNV-----
---IMDYNEDPFSYSPGKFR L-----NP
---IEDVTLNPIVR-----
---GEDIKFNKMKL-----
---LKSISLNPLTL-----
---IKDYNADDVNA-----NM
---LNKAYVHQ-----
KITQEEYNKAIKDIDFENSEDGDSCSYQLERAFGVLDENGEVVLDEDENLVLNEEKL
---FDGLQAREMDA-----
---LDNLKTFV-----
---LENIQSNILKK-----
---IKDINDDL SNW-----KT
---LEGLNVEI-----
RLELKGYPKGVEDF-----K
---MSKQDRELSTT-----
---LNRVDDRSILY-----
---LKR FNEDIQTK-----Y
---IQDFNNHVEDY-----K
---LQRFIVDADKY-----K
---LKEFDNNIFDN-----S
---LKKFDNNVFDN-----S

---HHH-----

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDTC01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene_46_crAssphage

DKIKVVLNYPYQDVFD-----PSRVT-----PMLQIMS YAGV-----
-QGYKYYKFNQDMLN---SLAASITEDVIPFSNIGMQGTVSKN---PNIGDI-----
--HSWDNLRWLQPM D---IHPIRERG-----LSVQLIKTLS--DDVALHPI-----
-EYGLSYFGMQMDPK---GKFGQSVT-----LGTQSTSMPLP--TNIFNNGV-----
TVIDLSYFGIQVDIN---KEFKGKVT-----LGTQSRTHIL--ANLYEGGMPADYT
-NVPFKYFGIQLETA---GTKHKQTR-----GTQLTKLAI--INLLASGIPQDIS
FIIPNDFIGRQQDTP---NKGAYEGD-----APKQAMKNIL--ANLDISSE-----
---SHLLYGKQQETK---TKGIKSTT-----FASQLKLNII--ADIDDTTL-----
---ENRYWGLQQDLT---PHTTDTQL-----EGSQVTKNILWVDPQKFSKE-----
---SNADWKLQQDLP---VKTIKPTL-----LGSQIQKNIY--SSLTDEAT-----
LVLNRDNFRIQDVPFKSDKHQDDKV-----SMGTQFFKLLF--GDGVIDKE-----
--LSYGDYRIQTNPV---EHINSSQL-----FGTQVRKLIM--ANIMDDYHYMNY
HVIPLSDYMVVQPSG---DHLTEDDL---MAIFGSQLRNILP--ADLPNDFS-----
-----SKLFMPQIIN---DKADMEAK-----MNRQIRKNIP--SMVENDVT-----
--LDSTKCLKLPQIIP---VSKKDKIT-----FSRQIRKNII--SDWLLDPN-----
---KSRNLRFPQTI---SRVEKPEI-----VFARQVRKNIL--ANLTPEDK-----
VELETRFQRIPQIVP---ENDKPVM-----MSTQARKHII--SNMLSDEV-----
--LETKNLRSPQVIE---TKTKAPLD-----GTQMAKLIL--SNILDSYN-----
QTLSHSNLYIQQQVP---SHLMDEEN-----KIGTQLQKRIL--DNLVFN D-----
-RIPITSYVIQQDIH---DHIMNAED-----ATGSQVGRITL--SGLNSEEA-----
--MRDSDFVIQQDIK---ADLLDETT-----ILGGQLVKQIM--EGLDWNNA-----
KTGWYSNLYTQQDIP---QHMDGENK-----AGLQIVKKLI--DNIGNTPE-----
EIYSYNFLYTQQETP---QHMNAENK-----AAIQIMKKIV--DNIPDTGT-----
EPYQYTYLYRQQEVP---SHLKDKEN-----KIGIQIYKLL--DNIPNDTT-----
ELFSYNYLYRQQEVP---QHMVDASN-----KAAIQIMKKML--DNLNPRTE-----
ELFSYNYLYRQQEVP---QHMVDASN-----KAAIQIMKKML--DNLNPRTE-----
-----Q-----

HHH-HHHHHHHH-----HHHH-----HHHHHHHHHH--H-----

gene_82_Woesebacteria_MGFQ01000035 -----GNEVNRAGKEINRLFRTLFINAKREIDDKISEEI-----
gene_22_contig0002 -----ALSNQKMMYTMNYSSDLSQKNKIEKLLTDLL-----
gene_52_Azobact_ph_AP017903 -----INKMSERKALDME-----
gene_25_842252479 -----VSDEYANTKFSDTETWEQAI DRYHNINKELIAREASTLASELGFSY-----
gene_45_LULO01000006 GDIASWEKMTPEQKKKASSVHKDASDYLDVNEITRRSRENLLKRFNLEKVV-----
gene_30_Chitinophaga_FOJF01000001 MDAAAWNKLSETEKIQLSPSYKLISANIAILKAMADKGYQQVLDQLGISE-----
MODO_RS15645 -----EKNYTYKGNKNTAKEMFDKVTSMGNTI INKQIDKFKRDMYT-----
gene_76_CEUT01009082 -----IGNRSGVEWREEVEANEIAI SNLGREEFLTKYGI-----
gene_1_contig0005 -----IEEFHALTADKSRGRNKLRLRELGATFN-----
gene_66_Cellulophaga_ph_NC_021806 -----YTIENEAFNGSGMFQAINDTVSAMSNLSIAGLSSELGK-----
gene_34_Flavobact_ph_NC_031904 -----GFVIDGKELTGKELYQHFNKAFSDIVDSKKQELFLDLGL-----
gene_64_Chlamydia_CVNZ01000007 IGGDGTVNLGGNHGRVKNLGRNLSVFNGLIVANILESFDNFKDGISD-----
gene_49_KJ003983 -----VDVTFGKDKVTLDRDKYIELYNTLIVDQLLDSFSKVDKEFSS-----
gene_34_775896848 -----YTLPDGTTLTGEEELKQYHQTHSDIVDAQYDRLEKELGYD-----
gene_80_FUFK010017431 -----LNLEFQNLIAENIKEDTKKIEDELGLTALRNAKSL-----
gene_57_BCSF01000013 -----PLYDTFHELVSENVNEDLKSILQDELGITE-----
gene_107_816813760 GGKKGQWTQGYILGGKKIGGEEELFDLYHNVTQQLVGRSLDDLKQLGYTKEIVTETPF
gene_40_MDTC01014143 -----YNVFGGKYNGKQIKDLYNEIYAERIKRSHDALTNELGWDS-----
gene_44_935075730 -----YTIGSTVRKGTGDYSYDGSFAFEYYQMLLSANANDEMYRLLADWGAI-----
gene_35_1001888980 -----IYELNGKKLGTGNQI AVEFSNSLHSLVYDGLVRFNFKVGA-----
gene_57_934554352 -----IYELDGKNVTGKELFDEFQKTLATNIREDDAMQLLYDIGGI-----
gene_44_933950148 -----GQSLIKDFFDNFTANIQDSFKDAASRIGVSI-----
gene_10_936104480 -----IGEYKKEFFKLYVANIKDSFNLSLVKELNIPT-----
gene_59_935734019 -----GRALKQRVFRNIVSNINTSFTNACTLLNIPL-----
gene_43_932717476 -----LNDLKSQVFNNYVANIRNSFEKTC AELGIGL-----
gene_46_crAssphage_JQ995537 -----LNELKDKVFNNYVANIRNSFEKTC AELGIGL-----
CONSENSUS 0.8 -----
2BE5 Thermus thermophilus (beta) -----
2BE5 Thermus thermophilus (beta') -----
Jpred4 sec.str. gene_46_crAssphage -----HHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHHH-----

gene_82_Woesebacteria_MGFQ01000035 -----
gene_22_contig0002 -----
gene_52_Azobact_ph_AP017903 -----
gene_25_842252479 -----
gene_45_LULO01000006 -----
gene_30_Chitinophaga_FOJF01000001 -----
MODO_RS15645 -----
gene_76_CEUT01009082 -----
gene_1_contig0005 -----
gene_66_Cellulophaga_ph_NC_021806 -----
gene_34_Flavobact_ph_NC_031904 -----
gene_64_Chlamydia_CVNZ01000007 -----
gene_49_KJ003983 -----
gene_34_775896848 -----
gene_80_FUFK010017431 -----
gene_57_BCSF01000013 -----
gene_107_816813760 PNSITQETYHGSNSPDIKNFDLSLQGAHEGRGINFALTREDAKVYGTHIYEVKLN
gene_40_MDTC01014143 -----
gene_44_935075730 -----
gene_35_1001888980 -----
gene_57_934554352 -----
gene_44_933950148 -----
gene_10_936104480 -----
gene_59_935734019 -----
gene_43_932717476 -----
gene_46_crAssphage_JQ995537 -----
CONSENSUS 0.8 -----
2BE5 Thermus thermophilus (beta) -----
2BE5 Thermus thermophilus (beta') -----
Jpred4 sec.str. gene_46_crAssphage -----

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDTC01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene_46_crAssphage

-----VEPG---LKG-----EDRRIAQ-----YEKYLKRVGMD
-----SIQL-----DKSKAS-----FSSTENIVRTF
-----TPNG--QFQF-----LGS-----NRKMLDTIL-G
-----DSKG--EVSY-----KLTGGN-----IQKFADMIR-D
-----TENG--FHIA-----DKKLVAAFISR
-----DKYGVYSDRL-----N-----DAKIRNLVI-T
-----DENF--DVN-----EETLYKALI-D
-----EQSGVYRIND-----QRKLIDKVL-T
-----DSEG--KID-----KPKLYDMLE-R
-----DTNG--QIKN-----EQNF-----VKNLQDLLI-K
-----ANKLSEKLI-Q
-----AENLQRALLNK
-----KLKK-----DPSNPELRLNF-----LQRVRELIYDN
-----EDKRDAKLDH-----LQKLKERLL-S
-----LNNA--KTDT-----QRLKAKRKY-----LEKLALI-R
SKWVTTDAEAKTKTTYGIKDVDEIVVYHPDQVQIIKEEIKGLEEKEKFTNLGRVIKR
-----YQQAIQSTDL-----ETRNKAQLDF-----LLKVVDVAI-N
-----TNDG--NIKY-----TSIETDGGLRNVIGVD--LDLVLADLR-R
-----ISNG--KINI-----DNGTVVVDQR-----KVMRSVKNSLL-T
-----DENG--EIRT-----DARGAIQID-----INKLVNRFQ-E
-----DARG--NVVY-----EEGKVKID-----NNQFIALIK-D
-----NEDG--SIKL-----DANGNIEGLD-----MKLFFNKLK-R
-----DSNG--NLQF-----DEQGNILGLN-----YERLLKLAR-E
-----DDNG--HIDL-----NSNGTIKLN-----RFVFFDRFK-E
-----DDNG--HIDL-----NSNGTIKLN-----RFVFFDRFK-E
-----E-----HHH-----HHHHHHHHH-H

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDTC01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene 46 crAssphage

ALNEM--GQIGNYAE----LVLNEQ-----ISV-GIPPL-RQKLLQAYRSL
GNQPL--GGAGVYQGNIALSALDNN-----SLLNHPNV-SNLLQAQLSRM
GARPT--NATDNVLH-----AAED-----GASVLAF-----KGLGKLRSSIS
EMTKR--DIPENVQA-----SIIDL-----FSREVT-FSQGLFEKNKIETLNFNAI
EIAKR--NYPENVAAGIESLLNQP-----GDKKV--FDIIVNKGRLNLLFSL
EAEAR--DLPDNFKS---AIGIDAE-----TGDFKISLEAISNY--RAVKSIYISI
QLDPV--KMDNLIK-----LRDPNVPIATMISS-TQWILPIITKK
TLSRQ--SFPDTALI----EGLENG-----LKFDSIFQS-RKKVMNLLANE
SLKQTLGSLPENIEQGLKEVYTHIN-PVTKEHTYLKLNPLDHPFY--KRIMPVLASM
EMLDK--GSAINLLK-----SIQKN-----LPIEAMPGI-KDKLYNIVFSK
EATSR--GYSVKSLAGLKIIEKLA---AGFYEFKTP-LWSSD-SNRYESLLNSI
GVINNSRESMDNVLA----YAVDD-----GKFAMPL-FEGGL-EHDAAMAFSM
IHNNL--KYGPEVEE---ALELNED-----KTGFKIPF--NNPNL-KNKIEELLST
AARNN--KIDSNLDK---QMIIVED---ILQQRDFNVPI-SFPVY-QREYQNLISAL
QIREK--GLPDSYLG---ALNIVPD---GQFDYRFAIPL-AFPNF-QAKFEQIYFSL
TLKDK--NLPDSYLG---ALDIVPN---GEGDVRFDIPL-AFPNY-QAKFEQIFFSM
SIEER--DLPENYERALNIKTEND-----VTFELPL-AFPAY-QKRFEQIILSL
SLEER--DLPDNYLLAMKIEKLVDD-----VNAYGFEAPL-SFPPF-AKRFESILLSL
YFNET--EIDRNFNIK---ATVVVN-----GKPFIPF-YHPTI-KSRIESVLLAR
NVKSQ--NDIASISD--TVDYVMDDELGNVVDNEMMFAL-SNPPI-RTKLASYMNSL
IISDD--VDAITIRK---ALEIKED-----GLPTMPL-SYPVI-KGKLEKILASM
ELTRR--GLDSNYRK---YAEINPE-----TGLPYMPA-WTNLV-RSKIENIVNSI
ECLRQ--GLDSNILE---FFTLNEDSPYTELGRANTVMPT-YMTNM-MSKAQNVCQSM
NAARN--GADKNTLD---FLTTDEN-----GNPRYPM-YLNSI-SNKIENLVNGI
NAQQQ--GVEKALLE---FFDLDSA-----GFNNLPL-FLSNI-NSKLESANSY
NAQQQ--GVEKALLE---FFDLDSA-----GFNNLPL-FLSNI-NSKLESANSY

-----RGAPVTNPGSD
HHHH--HHHHHHHH--HHH-----E-E-----HHHHHHHHH-

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDTC01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene_46_crAssphage

ITRKAIKPFFSGLRANQVTADYFDVFEKDGIPYSLKEVERDLKM-----EITREN
LGNAGTKPRVSGAQLHMPDIGSEE-----NLKF--
GLAATVRTDLQGGAFVRTAAVGRQ-----DLRY--
VTNTVVRKMNMGNVVLQSNFGYETMSKARKQDKNVPLER-----KLKSYG
VGNNSITQKVKGENAVQVSSLGSENTVRIVRQKGVESGRN-----GLKF--
VDKLLSPKMPGGPKVQVSSALFERNNRQAVYKPAGAKTANWER-----VTDFDK
IDKQVAQVNTNKGSVVQVADIGLSSVTDIDKTEILQVGSLELAPPPIPVKFSMLSESD
LTKSTVTYSTNGANMVQISSFGLFQKPESTITALEETNKS-----SIRY--
IKSSTIDLETKGGSYIQMSGLGFREKRTVDDKTK-----GIDY--
INSAAVKLKTNGGSFIQLSNFGLDKQTADAK-----GITW--
ITNKIMKHKMPGNGFVAGSESGFRMKENLEGVDKS-----RIIY--
FKKI VNKQKIGGSAVQVSAMGIKGYEEDG-----DLKY--
FKNAIQKQKINGGNIILVSNFGLSKDL-----QTAY--
FKTNVHSVKLPGRELVQVAGPGKWKIGDEVV-----ELRH--
FNNGIFKQQQKGGKELVQIAEVGGHIESS-----ELKM--
FNNNVYTQKIKGKELVQVAEVGAYLDSQGRK-----ELRM--
FKKNVINQTINGASLVQIAELGGVQEEKGTK-----ELGF--
FKNRVLKQRFNGMSVVQVAEFGYELDSQL-----EIRQ--
ITRRVTNLKLGAVHTIQPDTFLQPAAVTLDDKGIKGTQANVQRMYLEG--QIKFSD
INKHTIRQEVAGITVPIMPNI SWTYNTIENKARLKTMSAEEMKQHYNDG--AITY--
LSKQVINKYLPGFHAPIRADIFTASNELIKHNEFYSDKELYNKTI DELVANGSITYAS
FTNRVTRQVLPGFHASQVSDVGITALSGRTDLRDLMSKVEEKHGYSLGR--KITY--
FNNAITRQKLPGFHAAQVTVNGVYSK-----RLRY--
FNREITRQRMFGWHAQAQVADFGFGGWKTTKDTATDD-----KLAY--
FNTNITRQLISGWHAQLSDFGFKVDKQTQTD-----KLOY--
FNTNITRQLISGWHAQLSDFGFKVDKQTQTD-----KLOY--
-----G-----
RPLRSLTDILSGKQGRFRQNLGKRVDSGRS-----
--HHHHHHHH--HHH-----EE--

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDTC01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene_46_crAssphage

ISQVLSKGYTRRKLHGMRKGEDGTPGEMVMPFSYMTRFGVKPSESVNDLFTIILNDG
-----Y-----KDG-----
-----MDEHG-----
KDINPETLKGKSAEEIKRLKAAA-----
-----YRKKDPNNPNS-----
LTDQEKSSVRLTSSDLKIFYTRDA-----
RNVLQSTFETEIEPNSTVYFNEE-----
-----LKDGKLPRIENG-----
-----IVDKEVLTGPKMNEDG-----
-----LVEPSDLKPPVIEKDA-----
---LDSYNGKELQGTHTSTDENG-----
-----VTDGNG-----
-----KVDKKG-----
-----LDIDPKTG-----
-----YDG-----
-----YVEEG-----
-----VRNKDG-----
-----HANG-----
DYWQSRAELNEDGTIKRDANGTP-----
-----TPAFLEAMNRKEG-----
SFIEECKRTGRSLELQAEYREDG-----
-----HKDG-----
-----HPDG-----
-----RKIGKYNDND-----
-----KKIGEVVDG-----
-----KKIGEVVDG-----

-----VIVVG-----
-----EE-----

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDT01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene_46_crAssphage

LRINIRKMDRSARMNTIAENLPKTD-----FNATPMVRQFRELEGIFKPD
-----ETVFAEVKVP-----G
-----LTEAVVNPSTLANIMPH-----GIPFSMQKQ---V
-----PTTAMEVFLPH-----QYKEFLG-----Q
-----ETLGMEVYLP-----YFKEFLGEGLVQVRPDG---V
-----PYIEVYMPH-----WFREKLLERGSKTDQ---E
GIPSLDPESGKMRISKAKVLMPP-----KDLNIGISWEQF---K
-----KLIMGDIFIPY-----SYIENIPGHEKM-----S
-----SVASAVIFLPH-----TFKKDLGINSMS----A
-----DGKNYIRPGQIFMSH-----VQIAKLVDPYAKM----D
-----IVFHKAQVFIPS---KFKNDKKELIDLFEGFNGKEGKYLTRRENG-T
-----NILYAECEIPFDISYKDSEGNVHLDNFNEYCNPDGTLKMIIESEGKK
-----RKSIDYIPCYMPA-----YKRSLIQDLLVPRKDSAGEE-----Y
-----RVKHAEIMVSE-----D
-----TSPAEVMMRR-----SDLGIAENE-----N
-----GAWHAEMLMRA-----SDLGF-----E
-----KVIGAEVALPY-----KLAQKLG-----P
-----GVYAEVALPY-----ELAAKLG-----P
I IKKNADFKLQSEYWETKADGTK-----VFHPAEIILNNWDSRFLDA-----N
FNLLNERIDKDGNYKAEVILNV-----FDREIYDNLIEQEDG-----R
-----NYHYAEVIVNP---WKIDFYKNIGTVKTTINEDGTTKD-----I
-----SQIVEILLPK-----WMVKAYNTYDNEGNL--V
-----GRYIEVLLPK-----SNFDFAKNEDGTYKEPDE---I
-----VYYAEIKLPR-----WNKELDG-----
-----TPVYYTEIKLPR-----WSSQLKG-----
-----TPVYYTEIKLPR-----WSSKLG-----

-----PQLKLHQCGLPKRMALELFKPFLLKKMEEKGIAPNVKAARRMLERQR
-----EEEE--EE-----

gene_82_Woesebacteria_MGFQ01000035
gene_22_contig0002
gene_52_Azobact_ph_AP017903
gene_25_842252479
gene_45_LULO01000006
gene_30_Chitinophaga_FOJF01000001
MODO_RS15645
gene_76_CEUT01009082
gene_1_contig0005
gene_66_Cellulophaga_ph_NC_021806
gene_34_Flavobact_ph_NC_031904
gene_64_Chlamydia_CVNZ01000007
gene_49_KJ003983
gene_34_775896848
gene_80_FUFK010017431
gene_57_BCSF01000013
gene_107_816813760
gene_40_MDT01014143
gene_44_935075730
gene_35_1001888980
gene_57_934554352
gene_44_933950148
gene_10_936104480
gene_59_935734019
gene_43_932717476
gene_46_crAssphage_JQ995537
CONSENSUS 0.8
2BE5 Thermus thermophilus (beta)
2BE5 Thermus thermophilus (beta')
Jpred4 sec.str. gene_46_crAssphage

TEVTVDLNNIISWIDNVENTLSAISTRIPSHGLQSAFLGDITFFVNDN-G--NNTYIP
FANVGD-----ILIHVRIPASGPYSDFVAKVVGYTGKAQDGRNTILTP
FEKIKSDLGGA-----VELLVHRTPLNGLQSAKCNIKDLSTEST--GQNIQLP
NVPPGKTSVK-----ALQAVGFRIPTEGLNSVEYII IKDFLPQHA--GNTIIVP
YNQRGQKVGDS-----NLLQLIGFRIPDGLHSIDFMTIKGFLPQOS--GAQVMVP
LIDYLNKSGSD-----LLTGIGFRIPQELNSVEHIKIKGFLPQEM--GDTIIVP
LLMANGKVDKE-----IFRNILAYRIPNQAISSNDSIEIVGILPPYH--SDSAIVY
TAEKLRIGKD-----VLTVVGYRIPNQLASTDALEIAGILPKEA--GDTVVVY
LEISKTVKDSR-----VLEGVAYRIPNQGFTSIDSITIGGFLPQSM--GDTVAVY
SKTLSSMIDPK-----ALRAIGYRIPNQGQSSNDPLQIVGILPEAM--GDTIVAY
LTLKEGMIDPA-----LFNNFSEFRIPSSHKSGSSIEIAGILPPEV--GDLMIVP
ISLLEHRFPGS-----TSLLAYRIPTERDYSMLNLRVKKRFSQKTA--GGTIKVP
WEVDYKIKGN-----EDLLDLVGYRIPENKYSIFKLRKIGFLPISA--GTAIMLP
IAERLGLAIGT-----TGVLYRIPNQDYSSNVPSKIVGILPRGY--SKTVIVP
IEDLIANNDR-----LTVIGYRIPNQGKNSMLPMNVVRFLEPESH--AKGIIVP
PGTKIEDVDPN-----DSRLKVIGYRIPNQGKNSMLPLKVVQFLPESH--EKAIVVP
--TNIIEGIPGE-----LRTLVGYRIPGHGKNSMIPLRIVRILPEM--GKVILVP
GDTVDSNVDQQ-----LFEIMGYRIPQGKNSMLSLKITKVLPEM--GGVIILP
GNLDLNSVPEN-----LRTMFGIRIPTEGHQSMFVAKVVGVNLNGA--SQAIIVP
FIVDMSKIDPK-----ALEMFGFRIPTEGRQSIVLFEVVGFLNTGT--SQAQFP
ITVDIDKLDVE-----ARRMIGIRIPTEGKQSMVFEVVGFLNNA--TQAIIVP
HEVTLEDLQNA-----GLDTMIGYRIPTEGKQSIAMKVVGLLDESQ--GSTIVVP
RDENGLIGLLYQLQKAKLDTIIGYRIPTEGKQSICAMKVVQFTDDAQ--GSTIVVP
--VNIEDVSED-----ARTMIGYRIPTEGKQSVIMRVVDFLPNVS--DSTVVLP
--LNIEQVPS-----LRTMIGYRIPTEGKQSICIMYVKEFLPDAY--GSTVVVP
--LNIEQVPS-----LRTMIGYRIPTEGKQSICIMYVKEFLPDAY--GSTVVVP
-----RIP-----S-----I-----P
DIKDEVWDALEEVING----KVVLLNRAIPLHRLG-----IQAFQPVLVE-GQSIQLH
----HHHHHH-----HHHHH-----EEEEEE-----EEEE

```

gene_82_Woesebacteria_MGFQ01000035    SEQTVYDDTDFDIDQRSVYFYA
gene_22_contig0002                    AAYIKASNAFDKDKMLFTFKKF
gene_52_Azobact_ph_AP017903           HIMTSLGDDNDGDKLFFAMPY
gene_25_842252479                     SEIVAKSGADEFDIDKLTLYFPN
gene_45_LULO01000006                 SELVVKAGSDEFDIDKLTLYMPS
gene_30_Chitinophaga_FOJF01000001    SEIVLKAGSDEFDIDKLNLYLKN
MODO_RS15645                          HDITNKTGSNFEIDKLYLMFPS
gene_76_CEUT01009082                 EEMTAKTGSDEFDIDKLFMLIPN
gene_1_contig0005                     NDLTAKTGADFDIDKLYIMLAE
gene_66_Cellulophaga_ph_NC_021806    TEIPTKTGSDEFDIDKMYVMLPN
gene_34_Flavobact_ph_NC_031904      KNFTKQKGLDYDIDKESAYQLN
gene_64_Chlamydia_CVNZ01000007     AQGTTIAGEDEFDVKLYFMRRE
gene_49_KJ003983                     SDIVEMSGTDEFDIDKLFMIKN
gene_34_775896848                    GNITIQTGSDEFDIDKLFALFRN
gene_80_FUFK010017431                GGITTQMGSEFDIDKMYVIQKE
gene_57_BCSF01000013                 GGITVQMGSEFDVDKMMVMQKE
gene_107_816813760                   GEITTQMGTEFDMKVFLLMPN
gene_40_MDTC01014143                 AEITTMGSEFDIDKMYLMFPE
gene_44_935075730                     EHLVTRTGWDYDIDSIYLSMKE
gene_35_1001888980                   YDLVLQTGWDEFDIDKVYAYRKN
gene_57_934554352                     QSLITRTGWDEFIDSIYAYYRN
gene_44_933950148                     DEWVLQTGADEFIDSIYGIYHT
gene_10_936104480                     DDWVAQTGSDEFIDSVYGIQYN
gene_59_935734019                     ENWVHQSGSDYDVSVDYAMTFG
gene_43_932717476                     DEWVTQTGSDEFDVSVDYGMSTK
gene_46_crAssphage_JQ995537          DEWVTQTGSDEFDVSVDYGMSTK
CONSENSUS 0.8                         -----G-DEID-----
2BE5 Thermus thermophilus (beta)     PLVCEAFNADFDGDQMAVHVPL
2BE5 Thermus thermophilus (beta')    -----EEE-----
Jpred4 sec.str. gene_46_crAssphage

```

HHpred output

```

>pfam00623 RNA_pol_Rpb1_2 RNA polymerase Rpb1, domain 2. RNA polymerases catalyze the DNA dependent
polymerization of RNA. Prokaryotes contain a single RNA polymerase compared to three in eukaryotes (not including
mitochondrial. and chloroplast polymerases). This domain, domain 2, contains the active site. The invariant motif
-NADFDGD- binds the active site magnesium ion.
Probab=46.30 E-value=37 Score=35.72 Aligned_cols=47 Identities=26% Similarity=0.279 Sum_probs=0.0

Q ss_pred          EcCCCCCceeeEEEEEcCccCCCeEEechhHhcccCCCCchhheeEeec
Q AUX0017664536   631 RIPTEGKQSVIIMRVVDFLPNVSDSTVVLPEHWVHQSGSDYDVDSVYAMTF 681 (697)
Q Consensus       631 RIPTq~~~S~~~i~IvgFlP~~~Gd~IiVP~eiv~q~GsDFDiDKly~~~~ 681 (697)
                  |=||=-+.|+.+.++ .+.+. .+|-++.-+.+.+|||.|.|.++.|
T Consensus       105 R~PsLh~~si~~~~~i~~~~~i~~~~~c~~~naDFDGD~m~~~~~ 151 (166)
T pfam00623      105 RQPSLHRLSIMGHRV~RVLPG---KTFRLNLSVTTTPYNADFGDDEMNLHVP 151 (166)
T ss_pred         CCccccCcccccee-EEecC---CceeECcccCCcCCcCccceEeeC

```

Supplementary Note 1

Supporting evidence for remote sequence similarities used for crAssphage and IAS phage gene annotation

Query: gene_2_crAssphage_JQ995537

HHpred output:

```

>gene_2_crAssphage_JQ995537 Uncultured phage crAssphage, complete genome 249_aa||454|1203 crAssphage_JQ995537
Uncultured phage crAssphage, complete genome
MPDVKDLAQAAQGAAPESVNVVEPTTSVNPVQPTVVDTDNQSQAQVETIDDVVRRICTDGHSYVMTTIVITNIDCQERTGRNGKSYLNAFVTIASPVKGAQSMFDGTHRMGMLGAVQMPF
NQILLVMRKDKFYGRFVNYVGEAAEAGFASMYLTGVAVKVLQCFVPAGVQDHNPFTRKDNLYNVVDYDRYVYHIVGIEQPADPILVGAYNVLIKQIMEDA
RAATAAKREAKAKAASFVATAMSDDDLPF (characteristic SSB motifs and features are highlighted)

>PRK07275 single-stranded DNA-binding protein; Provisional
  Probab=42.04 E-value=23 Score =30.23 Aligned_cols=7 Identities=100% Similarity=1.659 Sum_probs=0.0

  Q ss_pred          ccCCCCC
  Q gene_2_crAssph  243 SDDDLPF   249 (249)
  Q Consensus       243 S~~D~PF   249 (249)
                    |||||
  T Consensus       156 ~~~~pf    162 (162)
  T PRK07275        156 SDDDLPF   162 (162)
  T ss_pred          ccCCCCC
  
```

Query: gene_15_crAssphage_JQ995537

CDD output:

pfam01443:Viral (Superfamily 1) RNA helicase; Helicase activity for this family has been demonstrated and NTPase activity. This helicase has multiple roles at different stages of viral RNA replication, as dissected by mutational analysis.

Pssm-ID: 307550 Cd Length: 226 Bit Score: 33.50 E-value: 0.18

```

          10          20          30          40          50          60          70          80
      .....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
3WRX_C  165 LVDGVPGCGKTKEIL-----SRVnFeedl--ilvpGRQaAEMIRRRAnas-giIVATKDNVR 218
query    46 alcGAGGTGKTfVIKyvinncwsggvigcaapthkacRVLsNsi-----gGKEvNTIQSL-F-----GFRLDvNIE 111

          90         100         110         120         130         140         150         160
      .....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
3WRX_C  219 TVDSFLMNYGKGarcQF-----KRLFIDEGLMLHtgcvNflv---emSLCDIAYVYGDtQQIPYInrv--tgfPYPAHF 287
query    112 NFDPENPAFNPV---GKdkldglKVLIIIDEASMLNaklvKYisnkckkLQIK-VIMLGSSQLPPVnektsqafLIASNT 187
  
```



```

T 5HZR_A          576 KVIQAGRFD  584 (732)
T ss_dssp          -----
T ss_pred          ccccCCCCC

```

Query: gene_18_933950148 close homolog of gene_21_crAssphage_JQ995537

HHpred output:

>TIGR00621 ssb single stranded DNA-binding protein (ssb). All proteins in this family for which functions are known are single-stranded DNA binding proteins that function in many processes including transcription, repair, replication and recombination. Members encoded between genes for ribosomal proteins S6 and S18 should be annotated as primosomal protein N (PriB). Forms in gamma-proteobacteria are much shorter and poorly recognized by this

model. Additional members of this family include phage proteins. Eukaryotic members are organellar proteins. [DNA metabolism, DNA replication, recombination, and repair]

Probab=65.26 E-value=46 Score=30.54 Aligned_cols=107 Identities=21% Similarity=0.225 Sum_probs=0.0

```

Q ss_pred          eeEehHHHHHHhhc---cCccccCCCCCccceccccC-cceeEEecccCCcCCcCccccCcceeEeccccchhc
Q gene_18_933950  213 QRQFWFRLNRYK---KGDWAFSGQGSEEGDLVFPNIVG-QGIFEEKFMLDATHFKEPSLMFDITKERIAPMDGVQSKQ 288 (347)
Q Consensus       213 pi~~WmKLLr~~kn---Kg~w~~~gs~~~gd~l~fp~FVg--GviE~~~~~ap~i~~~~l~id~~ke~i~~kd~~k~~k 288 (347)
                  .+.|=+|.....  ||++-++  .|.|.+.+. +|.-.....-.-+~+|.          ++...+++
T Consensus       54 ~v~~wg~~Ae~i~~~l~KG~~V~V-----G~L~~~~~G~~~~~I~a~~i-----~l~~~~~ 113 (164)
T TIGR00621      54 DIVIFGRLAEVAAQYLKKGSLVYV-----EGRLRTRKWEDQNGQKRKTEIIADNVQ-----LLDLLGAR--- 113 (164)
T ss_pred         EEEEEhHHHHHHHhccCCCEEEE-----EEEEEEeEEECcCCcEEEEEEEEEEEE-----EccccCcc---

```

```

Q ss_pred          cCcCccccCCCCccccccccCCCCCCCCccccCccCccccCCCCCCCCCCCC
Q gene_18_933950  289 KKAPNLAAAPGIGGIAMGAGIVNPSMPMGGFAGGVAGGFVSTESSAFAPDSEDNGGLPF 347 (347)
Q Consensus       289 ~~~~~m~g~~~a~~~np~~P~gg~~gG~a~ap~~~~~a~~~~~n~~d~lpf 347 (347)
                  +-+.+.+.|.+.+.  +|+.+.+.+.|++++.+.+.+.+.+.+.+.|.+.|.
T Consensus       114 -----d~pf 164 (164)
T TIGR00621      114 -----AQGGEFSGGGSGN---RPQGGQEQSATGKDPAQSAPQSGSADPIDGDDDEIPF 164 (164)
T ss_pred         -----ccccccccCCcC---CCCCccccCCCCccccCccccCCCCCccCCCC

```

Query: gene_19_crAssphage_JQ995537 (RecT)

HHpred output:

>TIGR01913 bet_lambda phage recombination protein Bet. This model represents the phage recombination protein Bet from a number of phage, including phage lambda. All members of this family are found in phage genomes or in putative prophage regions of bacterial genomes. [Mobile and extrachromosomal element functions, Prophage

Q Consensus 204 ~l~a~l~d~l~lt~e~i~-----l~ 278 (353)
 ..+-+.----- .++-...|.|.+.||+|.|-+++.+ .+|...- ...+ .-..+++
 T Consensus 187 ~leLf~yn-----khl~M~l~ltE~QFaqliGr~Rlyq~LP~k~k~lP-----~ll~Ds 247 (318)
 T pfam12987 187 RAALELFNMFN-----PAKELHLLQTLGNTSITEHQFCQIIGRMRLYQALPAGYQ-KRLP-----RLLLTDS 247 (318)
 T ss_pred HHHHHHHHhcC-----HHHHHHHHHhcCccchHHHHHHHHHhCCHHHH-hcCC-----ccccCh

Q ss_pred chHHHHHHHHHHHHhcCCcCCcCHHHHHHHHHhC
 Q Thu_Apr_27_21: 279 KISSRKMNVISDTYSYYFDGPGQRDILGTAWGAVNAISGYYSN 321 (353)
 Q Consensus 279 ~~~~~l~l~T~w~ynAvte~l~h 321 (353)
-.++.... +.=.....=|+|-+||-+|+.-+.
 T Consensus 248 qiN~vak~Yy~D-----nF~~~~~I~lW~yNL~T~AnKs 284 (318)
 T pfam12987 248 QINAVVKAYISD-----ENFGSNGNDITLWKLYNLLTGANKS 284 (318)
 T ss_pred HHHHHHHHhcC-----ccccCCCCcchHHHHHHHccccch

////////////////////////////////////

Query: gene_46_93454352 same family as gene_55_KJ003983 IAS virus (ThyX Thymidylate synthase)

HHpred output:

>COG1351 ThyX Thymidylate synthase ThyX [Nucleotide transport and metabolism].
 Probab=99.97 E-value=1.2e-30 Score=237.61 Aligned_cols=163 Identities=20% Similarity=0.215 Sum_probs=0.0

Q ss_pred ceeCCeEEecc--cHHHHHHHHHhCcccccC--chHHHHHHHHHhC-CCCcccCeEEEEeccCCHHHHHHHHh
 Q Mon_May_01_23: 2 NVVEPAIEFCDY--SGLRKIELIGKVCkQEEDSIKP--DSAETFCRNRLIDG-HTAIFEHYVYFNVTSIPNRIVREFVK 76 (279)
 Q Consensus 2 kii~p~ei~-----iE~agR~CYkSe~ki~--s~Fi~l~g~H~SvLEH~f~i~ 76 (279)
 ++..|-+.+.+. .+.+.|.|.|+.|-+.+ ..++|.+++++..|.|.|||++|||.|.|
 T Consensus 17 ~~~~~i~a~r~s~d~h~s~Eh~tF~I~----- 85 (273)
 T COG1351 17 KILDKGVVITIDSRGPLALIVQAARVSYPSGKLEDDGGQKDAELIRRIINEFGHESPLEHLVATFEIEG----- 85 (273)
 T ss_pred ccCCcCeEeccccchHHHHHHHHhccccccccCCcchHHHHHHHhCccChhhcceeEEEEec-----

Q ss_pred cCCeeeeccCCeHHHHHHHhCChhhHhHCCcCccCChhHHHHHhChhccccccCccccchhcccCccce
 Q Mon_May_01_23: 77 LSPYIRWSYLGNYIGFSYRVFLDIMSNRKMKAIYNDIYHPSEVNDLFYNMLLLSKEFSHLLFDDKDI TAKLEYGIDVSL 156 (279)
 Q Consensus 77 ~~~~~lv~nlR~f~e~ 156 (279)

T Consensus 86 ----- 85 (273)
 T COG1351 86 ----- 85 (273)
 T ss_pred -----

Q ss_pred eeccchHhCccccEEEEEEechHHHHHHHhCccceeeeeccccccccCCC-----CcceecCC
 Q Mon_May_01_23: 157 RIASDAEIRECAPEIYNVTYKITDDRGTHEAVRHREMSFMQESTRWCNVYAKGRLGYKYG-----RNISVIEPP 225 (279)

Q Consensus	157	~~~~~t~~~~sR~~~qlvRHR~~Sfsq~SqRYv~~~~~i~p~	225 (279)
		+ +++ + + + + - + . + + - . + .	
T Consensus	86	-----Sr~~~~Ql~RHR~aSy~e~S~RYv~~~~e~~~~~	142 (273)
T COG1351	86	-----VSRVAAHQlIRHRlASySEKSQRyVLDGDEFYVPFEPLLEDRAAFNDLCRAVDDLys	142 (273)
T ss_pred		-----cHHHHHHHHHccccCCCCceeeecCCCceecchhhcchhhhhHHHHHcccc	

Q ss_pred		CCCcHHHHHHHHHHHHHHHHHHHHCCCCCHHHHHHHhchccccEEEEEEcc	
Q Mon_May_01_23:	226	FKNEDSLEKFYDVVAGNEAIYQELTNDGEPQLARSVLPTATKSDIYISGTLD	278 (279)
Q Consensus	226	~~~~~Y~l~G~E~AR~vLP~~~~T~iv~t~Nl~	278 (279)
		. . . + . . . + . + + + . + . + . . . + + . + + + + + + + + + + +	
T Consensus	143	~~~~e~~~~~Y~l~G~E~AR~vLP~at~T~v~t~N~R	195 (273)
T COG1351	143	SKLEKVLDDLEKAYERSYELYKELLRSGIAREDARYVLPNATETRIVVTGNAR	195 (273)
T ss_pred		chhHHHHHHHHHHHHhHHHHHHHHhCCCCHHHHHHhcchhhEEEEEcHH	

////////////////////////////////////

Query: alignment based on close homologs of gene_51_crAssphage (Phage stabilization protein)

HHpred output:

```
>pfam11134 Phage_stabilize Phage stabilization protein. Members of this family are phage proteins that are
probably involved with stabilizing the condensed DNA within the capsid.
  Probab=86.29 E-value=31 Score=33.25 Aligned_cols=135 Identities=9% Similarity=0.052 Sum_probs=0.0
```

Q ss_pred		Cccccceecc--cceEEcCCCccccEEEEEEcCEEEEEcCceEEEcchheeccCCCCcccc--cceecCceeeecC	
Q gene_5_9361044	22	SLANSWRHFRA--NNYKVLskNKGNItnIIGVGTAFFVHTEHSLFYlNRDNLlKTSGNTAQLEMP-DLFEVEPIELFTSN	98 (177)
Q Consensus	22	s~in~F~n~kdl~G~I~L~L~qe~k~l~~~~~l~~~~~	98 (177)
		..+ . . + . . + . = . . + . . . + . . + . . + . . = . + . . - + + . + +	
T Consensus	161	S~l~d~s~l~fatAE~PD~iv~~~~~el~-----~lfG~TiEv~ntG~f~f~r~g~i~	228 (469)
T pfam11134	161	TDLED-ESHPDRYSAQYRAESQPDGIIGVGVRDFI-----VCFGSSSTIEYFTLTGSTTAGAALYVAQPSYMQ--	228 (469)
T ss_pred		cccc-cCccCccccceeeecCCCcEEEEECEEE-----EEEcCcEEEEEEcCCCCCcceecceEc--	

Q ss_pred		CccccCCcEEcCCcEEEECC---CEEEEECCcceechH-HHHHHhh----cCCeEEEEeCC-CEEEEEe	
Q gene_5_9361044	99	HGYGGLQHPQAWTVNSNGYFVDAD---NKRIYNFDNNHLTDLTS-DILNWMNNV---QIADAHMVTDFANARVIMCL	169 (177)
Q Consensus	99	~g~ggs~npeS~yF~D~rg~V~l~gi~IS~gm~F~g~YD~y~it~	169 (177)
		- -- . ++ . . + . + . . . + . + + + + + . + . + . . . - + . + +	
T Consensus	229	---G~a~s~s~wlg~g~V~g~IST~Ie~l~y~elsda~sy~GH~f~vl~	305 (469)
T pfam11134	229	---KGIAGTFCKCPYMDKYAIIISHPATGAPSVYLIGSGQKSPIATASIEKIIRGYTADELAAGVMESVRFDSELLIIHL	305 (469)
T ss_pred		---ccccccchhhCCEEEEEcCCCCcEEEEEcCCcEEcCHHHHHHHccccccEeeeeEECEEEEEEEc	

Q ss_pred		Cc	
-----------	--	----	--

Q gene_5_9361044 170 AY 171 (177)
Q Consensus 170 ~d 171 (177)
.+

T Consensus 306 P~ 307 (469)
T pfam11134 306 PR 307 (469)
T ss_pred CC

//

Query: gene_52_crAssphage_JQ995537 and gene_43_KJ003983 IAS virus (Tail tubular protein A)
HHpred output:

>3j4b_A Tail tubular protein A; bacteriophage, DNA ejection, tail complex, gatekeeper, viral; 12.00A
{Enterobacteria phage T7}
Probab=70.94 E-value=57 Score=29.68 Aligned_cols=70 Identities=10% Similarity=0.041 Sum_probs=0.0

Q ss_pred EEE-eCCEEEecCceEEEEEEececCCcCCCCCHHHHHHHHHHHHHHHHHhhccCCChhhhHHHHHHHHHHHH
Q gene_26_Chitin 177 FDI-RDGKIFTGFTEGLHIIFYKKEFDDSEFQLIPDNHRIIEYIKAFKLYKMYEQIFNEVADESFNQAERKYQLYKQEY 255 (291)
Q Consensus 177 y~I~~~~I~tnf~G~v~l~y~~~~d~g~plIPD~~~~eai~yi~k~l~~~~~y~~~~w~~~~A 255 (291)
|.t .+++|.++ .|.|.|.|.|.|.|.|.+.++ .+.+++|+.+.+.+.+.+.+.+.+.|.+

T Consensus 92 yd~~~~i~t~---v~v~yv~~~~d~~~~p~f~i~lA~~~~a~~~~l~~~~a~~~~a
T 3j4b_A 92 YDRTSQSDREFS----GITVNIIRLR-DYDEM----PECFRYWIVTKASRQFNRR--FFGAPEVEGVLQEEDEARRLC 159 (184)
T ss_pred EECCCCceEeCC----cEEEEEEecC-ChhhC-----CHHHHHHHHHHHHHHHHHhh--ccCHHHHHHHHHHHHHHHHH

Q ss_pred HHH
Q gene_26_Chitin 256 YDA 258 (291)
Q Consensus 256 ~~~ 258 (291)
++.

T Consensus 160 ~~~ 162 (184)
T 3j4b_A 160 MEY 162 (184)
T ss_pred HHH

//

Query: gene_53_Azobact_ph_AP017903, homolog of gene_85_crAssphage_JQ995537
HHpred output:

>3h4r_A Exodeoxyribonuclease 8; exonuclease, recombination, hydrolase; 2.80A {Escherichia coli}
Probab=97.52 E-value=0.00049 Score=70.62 Aligned_cols=141 Identities=15% Similarity=0.128 Sum_probs=0.0

Q ss_pred cccccchhhHHHHHHHHHHHHcccccCChHHHH-----HHHHHHHHHcCc-cCCceecCChHHHHHHHHHHHHHHHH
Q gene_55_100190 199 EATFSTETGTFTHGLMEVHSLGKESTHKDRVD-----ELMRHWDKERGK-NDNPLYIEKDEDARRMCEKITDMYEEI 270 (1664)

```

Q Consensus      199 k~tp~TeaG~H~lle~vh~l~t~wd~-----ea~w~e~g~hll~d~r~i~kit~e~i 270 (1664)
                  ..+++...|+.+|.+||...+.+.-...+.+      +++...|.++.-. ...+-.-...+.+|.+.+.+. +..
T Consensus      39  ~~~~~G~H~le~~~~~Dg~~~~g~vDlIv~dK~tg~giID~Kt~e~~~~~f~kf~n~K~vg~eht~w 117 (265)
T 3h4r_A         39  TKTKTLDLGTAFHRCRVLELEEFNSRNFIVAPEFNRRTNAGKEEEKAFLECASTGKTVITAEEGRKIELMYQSVMAL-PLG 117 (265)
T ss_dssp
T ss_pred        CCCSCCTHHHHHHHHSSHHHTCC-----CCCTTHHHHHHHHTTS-HHH
                  CCCcHHHHhHHHHHHhCcccccEEcCCCCcchhhHHHHHHHHHHhCCCccCHHHHHHHHHHHHHhC-chH

```

```

Q ss_pred        Hhchc~cEEeeeCCeEEeCCCCceeeeeeEEEEeCCCCeeEEeeccccccCCchHHHHHHHhcCcccCCCccccch
Q gene_55_100190 271  NKKY-EIVSMEAPVVLHSTEDGSPILGRCDI IALDKETGGVLDIDIKTHVENELPAEGEYNKFFQNYNLKPVNGEHTAW 349 (1664)
Q Consensus      271  ~kf~vis~E~Pl~~~~~Dg~~~~g~vDlIv~dK~tg~giID~Kt~e~~~~~f~kf~n~K~vg~eht~w 349 (1664)
                  ..-| ..+...|.+.+...+|...|+|.|.  |...|.|.  ....|.+.
T Consensus      118  ~~~~~E~~~~~G~iD~i~----~IiDyKT~----- 176 (265)
T 3h4r_A         118  QWLVESAGHAESSIYWEDPETGILCRCPDKIIP---EFHWIMDVKTTA----DIQRFKTAYYDYR----- 176 (265)
T ss_dssp        HHHHSSSCBSSCCEEEECTTTCCEEEECCSEEEG---GGTEEEEEEEES-----CHHHHHHHHHHT-----
T ss_pred        HHHhccCCceeeEEEEeCCCCeEEEEeeeec---CCeEEEEEecC---CHHHHHHHHHhC-----

```

```

Q ss_pred        hhhchHHHHHHHHH
Q gene_55_100190 350  DHANVQTMGYAAAF   363 (1664)
Q Consensus      350  d~~~~QT~aya~~~~ 363 (1664)
                  .+.|...|+++
T Consensus      177  ---~Ql~Y~~~~ 188 (265)
T 3h4r_A         177  --YHVQDAFYSDGY 188 (265)
T ss_dssp        --TTHHHHHHHHHH
T ss_pred        --hHHHHHHHHHH

```

//

Query: Alignment homologs of gene_76_crAssphage_JQ995537 (major capsid protein)
PSI-BLAST output (first round):

>ANA48974.1 major capsid protein [Pseudomonas phage PaMx41]
Length=418

Score = 198 bits (502), Expect = 4e-58
Identities = 82/488 (17%), Positives = 161/488 (33%), Gaps = 83/488 (17%)

```

Query 1      MGSMMQVYKGEYN-SGFTDENHLSNALLQPEPEELSKIITHLYGSDDSRFPLTLTEGMGN 59
             M      ++   N      ++      L + P      + ++
Sbjct 1      MSVYAGIFNTTLNPQELNMKSFAGTILRRVNGSAPLLAM----- 40

Query 60     EETIGDTEYEWKVMGRSKRPSPIVTVNSSTTNPVGSGSTFEIEFEENWFAPGHVLDVADDN 119
             +G T +      G +              T              + +E +      G + +

```



```

T 5GAI_C          53 FDVVRPVVRKLVSEMRQNPIDVLYRPKDGA----- 82 (721)
T ss_dssp        CCHHHHHHHHHHHHHSSCCECECTTC-----
T ss_pred        hhhHHHHHHHHHHhCcccEEcCCCC-----

Q ss_pred        HHHHHHHHhccchHHHHHHHHHHHHhchHHHHHHHHHHHhCceEEEEccCC-----eEEEeEe-Chhch
Q Q_6978104      165 EAFIKEFNENFIDDISAQGD LINVIDDLTDAFTIYARAYFEFVAFGACYTYRDVVG N-----QLIKRVV-SVRDA 234 (803)
Q Consensus      165 ~e~k~d~E~l~D~li~g~v~v~p~v~
      -.....+.++++.+++.+++++|++++|.+++++....      ++.++.||..-
T Consensus      83 -----d~~Ae~l~l~d~~G~G~k~d~~p~~i~~V~~P~~ 151 (721)
T 5GAI_C         83 -----RPDAADVLMGMRYRTDMRHNTAKIAVNIAVREQIEAGVGAWRLVTDYEDQSPTSNNQVIRREPIHSACSH 151 (721)
T ss_dssp        -----CGGGHHHHHHHHHTCSHHHHHHHHHHHHHTSCBEEEEEEECSSSCSBTTEEEEEEEECCTTSS
T ss_pred        -----CHHHHHHHHHHHHHHCHHHHHHHHHHHHHHCceEEEEEEecCCCCCCCCcceeEEecCCCCce

Q ss_pred        ccCCCCC--chhhhhEeeccCH---HHHHHHHHHhCCHHHHHHHHHhcccccccccccccccccecccccc
Q Q_6978104      235 FVPVNDNM--FAEDYDMFAERRMLTK---QIIDEFYEYLSEKEREALDITYQYSATTSSDRALLNWDKYMYYFGDICKS 309 (803)
Q Consensus      235 ~~~~~--yid~yvge~mti--seiid~f~LT~i~ie~~~~~
      |+.+++.. .+|+.||+..+||+ .+.+.|+.
T Consensus      152 ~~~Dp~a~d~Da~v~t~l~y~ 189 (721)
T 5GAI_C         152 VIWDSNSKLMDKSDARHCTVIHMSQNGWEDFAEKYDL----- 189 (721)
T ss_dssp        CCCCSCCCSCSSCSCSCSEEEEEEECHHHHHHHHTTC-----
T ss_pred        EEECCCCCCCCHHCeEEEEEEecccChHHHHHhCC-----

Q ss_pred        cCccccccCccccccCceEEEEEEEEEeccccceEE--ecCceEEEEeCceecC---CCcCeEEEEe--
Q Q_6978104      310 FNKDDLQHIKNTNIMARDANGLFEVWHTVWRGEIKEGILTY-SNGAFVTRIVDETYQLNP---AGGDISIEWVWR-- 382 (803)
Q Consensus      310 ~~~~~V~w~k~k~v~d~E~i~i--
      .....+.+++++...+.+++|+.+++|+.+++|+.+++|+.+++|+.+++|+.+++|+.+++|+.+++|+.+++|+.+++
T Consensus      190 ~~~~~v~v~E~w~v~d~g~g~ 269 (721)
T 5GAI_C         190 DADDIPSFQNPNDWVFPWLTDTIQIAEFYEVVEKKEAFIYQDPVTGEPVSYFKRDIKVIDDLADSGFIKIAERQIKR 269 (721)
T ss_dssp        CCCCCCSCCSCSCSSSCSEEEEEEEECCECCSSSSSCSEEBCCSSCSSTHHHHHTTCCBCCCCCECC
T ss_pred        CcccCccCcccCccccCCEEEEEEEEEEeEEEEEEeCCCCCceecchcHHHHHHHHHHCCCeEEEEEEeccc

Q ss_pred        ccceeEeeecCccceeeCcccCCCCCCCCCCEEE--EecCCCC--ccchHHhCCHHHHHHHHHHHHHHHHHHHhCC
Q Q_6978104      383 PQVYESVRIGSRATSIYPYKARPIAYNRNGKLPYNGI--AELLPGFG--RFSVVDTVIPYQVFRNIVSYHREMAIAKNKM 458 (803)
Q Consensus      383 ~~~~~k~I~r~k~p~y~r~p~p~l~v~d~p~q~n~i~k~l~a~a~
      ..|+.-+.+.+.+.+.+. .+.||+. .+...+ ..|+|.++|.|.++|.+.++.+.+.+.+.+.+.+.+.+.+.
T Consensus      270 ~~~~~g~iL~v~p~v~g~g~v~d~Q~N~s~d~s~ 342 (721)
T 5GAI_C         270 RRVYKSIIITCTAVLKDKQL-IA-----GEHIPVVFGEWGFVEDKEVYEGVVRLTKDGQRLRNMIMSFNADIVARTPK 342 (721)
T ss_dssp        CEEEEEEECSSCEEEEEEE-CS-----CSCSCCCBCSEEESSSSEEECSHHHHHHHHHHHHHHHHHTHHHHSSCC
T ss_pred        ceEEEEEEecCeEecCCCC--CC-----CCCCCEEEecceeeCCEEEEEechhccchHHHHHHHHHHHHHHHCCCC

```


Q ss_pred CCCCCcceeeeeecCCCcccCcccceeEEEEcCCccccCCceeEEEEeehCCHHHHHhhhhhhCccE
Q Q_8181294 11 IYSPTTDSFKLGYRMKNGVEADDSWSSLLSVANNNPDCAI GKDAVTIKVEELSTMQNFDEFMNVTPTMTVGTRTTGT
Q Consensus 11 ~~~~~G~S~I~~s~~n~~~~G~~~~l~~dE~Gk~~~~w~~k~~l~~G~~~~Gk
.....+|.|.|.+. ++.+.+|.++.++++||+++.+.+.|.+.++++.+.|.+.+
T Consensus 46 ~~~~~i~f~n--Gs~I~~s--~~~~rG~~~~viidE~a~~~~~l~p~l~~~~
T 200J_A 46 GIVEWN-KGS-----IELDN--GSSIGAYAS--SPDAVRGNSFAMIYIEDCAFIPNFHDSWLAIQPVISSG--RRSK
T ss_dssp CEEEC-SSE-----EETT--SEEEEEEC--SHHHHTSCSEEEESGGGSTHHHHHHHHHHHHST--TCCE
T ss_pred CeeecC-CCE-----EEeCC--CEEEEEC--CCccccCceEEEEcChhCCChHHHHHHHHhhcC--CCcE
Confidence 111000 000 01111 366666554 467889999999999999998765578988888888763 3478

Q ss_pred EEEEEcCccccchHHHHhc
Q Q_8181294 11 LMAWGTATAANMQIFEQNFY 340 (751)
Q Consensus 11 ~i~sT~n~~~~~kem~ 340 (751)
+++|.|.++. +|.|++|. |
T Consensus 46 ii~iSTP~g~-~f~l~ 324 (385)
T 200J_A 46 IIITTPNGL--NHFYDIWT 324 (385)
T ss_dssp EEEEECCSS--SHHHHHH
T ss_pred EEECCCCC--CHHHHHH
Confidence 999998875 4455554

Query: alignment for close homologs of gene_56_crAssphage_JQ995537

HHpred output:

PF12571.7 ; DUF3751 ; Phage tail-collar fibre protein

Probability: 99.18 E-value: 7.8E-19 Score: 153.26 Aligned Cols: 146 Identities: 14% Similarity: 0.202
Q ss_pred CCCCCcchHHHHHHhCCEEEEEECCHHHHHhhcccceEecceEEEE--cCCcEEEEEECCcCCEEEEEEEECCEEEEEEEcCCcCCEEEEEECcCceEEEEECcceeH
Q Q_20456 AEYSKLYITNNGQALMAKMIAGSGNIDFTKVCSSSTQY-TESQLQALTALSNIKQTTLVSKVTR--TNEVAIKIDAAYSNDLKEGYMRTLGLYAVDPDKGEILYAVCIKSNNCYMPYPNGVTVAAYLQLYTTVGADNVSLAVSPGAYATVGD 154
(600)
Q Consens a~~~~~lT~G~~~~a~a~a~~~~i~t~~~~g~g~~~~~t~l~~~~~i~~~~~i~~~~~n~~~~g~~~~eiGl~a~d~~~~eiLy~~~~~p~~~~~vsn~Vt~id~~~~t~ 154
(600) ++|.+.|++|+++++++|.++|+++|+. +.+.+.+.+.+.++++.+.+.+. .++++.+.+.+. .++|++|+|++|+++ +|++.... +|.|.+.+.+.+.+.+.++++|++|.++|+...+
T Consens ~~~~~i~t~G~~~~~l~t~~~~~g~g~~~~~t~l~~~~~i~~~~~i~~~~~n~~~~g~~~~eiGl~a~g-----l~a~~~~~v~~~~id~~~~~ 150
(150)
T PF12571 EQKYKTILTHHGERVIVEALANKIPVPLKEMAIGDNGSSITPSASQTTLVREYVRAEITDLEDPQNRHQMIAELLIPEN--VGGFIVREIGLFDEQG---GLVAVANCP--ENYKPVLEQGSGKQYRMIQVSSSDAVTSLINNIVYATRT 150
(150)
T ss_pred CccceEeCHHHHHHHHHHhCCEcCEEEEEECCECCCCCChhhcCcEEeEecceEECCCCCEEEEEEEcC--cCEEEEEEEECCECC--CCEEEEEEC--CCcCccCCEcCEEEEEEEECcCCEEEEEECcceeec

Query: gene_39_crAssphage_JQ995537

CDD output:

Thioredoxin_like super family c100388

Protein Disulfide Oxidoreductases and Other Proteins with a Thioredoxin fold; The thioredoxin (TRX)-like superfamily is a large, diverse group of proteins containing a TRX fold. Many members contain a classic TRX domain with a redox active CXXC motif. They function as protein disulfide oxidoreductases (PDOs), altering the redox state of target proteins via the reversible oxidation of their active site dithiol. The PDO members of this superfamily include the families of TRX, protein disulfide isomerase (PDI), tlpA, glutaredoxin, NrdH redoxin, and bacterial Dsb proteins (DsbA, DsbC, DsbG, DsbE, DsbDgamma). Members of the superfamily that do not function as PDOs but contain a TRX-fold domain include phosphatases, peroxiredoxins, glutathione (GSH) peroxidases, SCO proteins, GSH transferases (GST, N-terminal domain), arsenic reductases, TRX-like ferredoxins and calsequestrin, among others.

The actual alignment was detected with superfamily member TIGR00411:

```

Pssm-ID: 320948 Cd Length: 82 Bit Score: 28.31 E-value: 0.26
           10          20          30          40
.....*.....|.....*.....|.....*.....|.....*.....|.....
lcl|seqsig_MKKLI_ 44 TDTGRQEQARSAGISDAPTAYCNGDIL-RGVQSDYTIRKYLRLKL 86
Cdd:TIGR00411    39 VMENPQKAMEYGIMAVPAIVINGDVEFIGAPTKEELVEAIKKRL 82

```

////////////////////////////////////

Query: gene_41_crAssphage_JQ995537

CDD output:

TRX_family cd02947

TRX family; composed of two groups: Group I, which includes proteins that exclusively encode a TRX domain; and Group II, which are composed of fusion proteins of TRX and additional domains. Group I TRX is a small ancient protein that alter the redox state of target proteins via the reversible oxidation of an active site dithiol, present in a CXXC motif, partially exposed at the protein's surface. TRX reduces protein disulfide bonds, resulting in a disulfide bond at its active site. Oxidized TRX is converted to the active form by TRX reductase, using reducing equivalents derived from either NADPH or ferredoxins. By altering their redox state, TRX regulates the functions of at least 30 target proteins, some of which are enzymes and transcription factors. It also plays an important role in the defense against oxidative stress by directly reducing hydrogen peroxide and certain radicals, and by serving as a reductant for peroxiredoxins. At least two major types of functional TRXs have been reported in most organisms; in eukaryotes, they are located in the cytoplasm and the mitochondria. Higher plants contain more types (at least 20 TRX genes have been detected in the genome of Arabidopsis thaliana), two of which (types f and m) are located in the same compartment, the chloroplast. Also included in the alignment are TRX-like domains which show sequence homology to TRX but do not contain the redox active CXXC motif. Group II proteins, in addition to either a redox active TRX or a TRX-like domain, also contain additional domains, which may or may not possess homology to known proteins.

```

Pssm-ID: 239245 Cd Length: 93 Bit Score: 25.60 E-value: 3.0
           10          20          30          40          50          60          70          80
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*

```

lcl|seqsig_MIRID_ 7 FTRDGC DACKiAIKNITDAINEANCDIT---LNIrnTNLDDILRK-EITKFPTTVITKvdndyKRKELARLEGSFPDSYIKDIIN 87
Cdd:cd02947 17 FWAPWCGPCK-AIAPVLEELAE EYPKVkfVVDV--DENPELAE EYGVRSIPTFLFFK-----NGKEVD R VVGADPKEELEE FLE 93

Query: gene_84_crAssphage_JQ995537

CDD output:

PG_binding_3 pfam09374
Predicted Peptidoglycan domain; This family contains a potential peptidoglycan binding domain.

Pssm-ID: 286462 Cd Length: 76 Bit Score: 39.38 E-value: 2.66e-05

.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MKLNK_ 66 SLNAYKSDTTISADYFVAKYKLERIRYYNDIAGK-GNNIKFLRGWIRR 112
Cdd:pfam09374 27 TLA AVKQRASAGEDALIDAICLARRAFYLRLAAK rTTNARFLRGWVNR 74

Query: gene_36_Cellulophaga_ph_NC_021806 (homolog of gene_22_crAssphage_JQ995537) primase

CDD output:

Toprim_N super family cl26789
DNA primase catalytic core, N-terminal domain;

The actual alignment was detected with superfamily member TIGR01391:

Pssm-ID: 331610 Cd Length: 415 Bit Score: 36.82 E-value: 0.02

.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
gi 526178681 14 ITSEII---LERLNPIDVYKMYIK----DEFVVRpfsSPFRSDNIPFSIYQDRRSdqILFNDFVEG-GGNCIQFVKK 84
Cdd:TIGR01391 2 IPEEFIdelKERVDIVD VISEYVKlkkkgRNYVGL---CPFHHEKTPSFSVSPEKQ---FYHCFGCGaGGDAIKFLME 73
90 100 110 120 130 140 150 160
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
gi 526178681 85 LFNSTWYQACSRIAIDFKISDDYTVDKMNTPNVNLQNTLDREVKVKSKIVNLQVKIRNYNSDDIAYWDSYGISMETLKR 164
Cdd:TIGR01391 74 IEGISFVEAVEELAKRAGIDLPEFKDQEKKEQKSKRKKLYELLELAAKFFKNQLKHTPENRAALDYLQSRGLSDETIDR 153

Toprim pfam01751
Toprim domain; This is a conserved region from DNA primase. This corresponds to the Toprim domain common to DnaG primases, topoisomerases, OLD family nucleases and RecR proteins. Both DnaG motifs IV and V are present in the alignment, the DxD (V) motif may be involved in Mg2+ binding and mutations to the conserved glutamate (IV) completely abolish DnaG type primase activity. DNA primase EC:2.7.7.6 is a nucleotidyltransferase it synthesizes the oligoribonucleotide primers required for DNA replication on the lagging strand of the replication fork; it can also


```

Query 102 SNVIDFSHNDITEYG-HINSRTFYKAIQELYTENIIRPTNKKNVYVNVHNYIFRGNINKFIQLY 164
      ++ I S++ I      I+ T++K ++EL + I T +N Y +N +Y+F G+  F++ Y
Sbjct 62 TDTILMSYDIIDMKAVKISRTTYFKGMKELVEKQFIAETMIQNYFFINPDYMFNGDRLSFFVKSY 125

```

```

>WP_074898747.1 hypothetical protein [Bacillus megaterium]
SFH67410.1 replication protein (RepL) [Bacillus megaterium]
Length=156

```

```

Score = 45.9 bits (107), Expect = 2e-06
Identities = 21/74 (28%), Positives = 38/74 (51%), Gaps = 2/74 (3%)

```

```

Query 90 IISYICDNLKYSNVIDFSHNDITEYGHINSRTFYKAIQELYTENIIRPTNKKNVYVNVHNYIFRGNINKFIQL 163
      +++YI DNL +NSN++ + +I+ I T I+ L + ++ K V ++N + I G+ NK L
Sbjct 69 VVNYILDNLDWNSNILIKTQQEIAKEAGIGFNTVNTTIKMLVDKFKLK--VKTGVIMLNPDI IAYGSHNKRAHL 140

```

//

Query: Alignment of homologs of gene_54_crAssphage

HHpred output:

```

>l0WF_A Integration Host Factor Alpha-subunit/Integration Host; protein-DNA recognition, indirect readout, IHF; 1.95A
{Escherichia coli} SCOP: a.55.1.1
Probability: 31.26 E-value: 280.0 Score: 22.05 Aligned Cols: 37 Identities: 8% Similarity: 0.191

```

```

Q ss_pred cchHHHHHHHHHHHHhCCCCCCCccccHHHHHHHHHHHHHHHHHH
Q gene_33_Azobac 1 --MTLNQIADDIILTAGYINK-EKAMTINRLQLLSWIHMYRNIIMLQ 60 (339)
Q Consensus 1 ~mmTLnqIv~I~::~::~::~s~d~::~~rqIk~I~yRallIrQ 60 (339)
      |+||..+|+..|+.... +++..+++..+...-..+...
T Consensus 1 ~::~::~~l~:::ia~::~~-----s~::v~v~::~::~~i~:: 37 (99)
T l0WF_A 1 MALTKAEMSEYLFDKLG-----LSKRDAKELVELFFEEIRRA 37 (99)
T ss_dssp -CBCHHHHHHHHHHHHC-----CCHHHHHHHHHHHHHHHHH
T ss_pred CCCCHHHHHHHHHHHhC-----CCHHHHHHHHHHHHHHHHH

```

//

Query: gene_57_crAssphage_JQ995537

HHpred output:

LTQDMVNEENTRYIIRYAFDLSGKTITMPVGCCELVFEGGIIENGTINLNKCKLTGMVGEE
Sbjct 351 LTQDMVNEENTRYIIRYAFDLSGKTITMPVGCCELVFEGGIIENGTINLNKCKLTGMVGEE 410

Phage-related protein tail component [uncultured Mediterranean phage uvMED]
Sequence ID: gi|787039080|BAR14575.1

Length: 1128Number of Matches: 1

Score Expect Identities Positives Gaps

29.2 bits(64) 5.8 12/42(29%) 20/42(47%) 0/42(0%)

Query 6 VNEENTRYIIRYAFDLSGKTITMPVGCCELVFEGGIIENGTIN 47
+N + I R ++ L G + +P + G I+ NGT N

Sbjct 275 INAQQFSQIPRRSYRLRGLKVQIPHNGIVQASGSIVYNGTFN 316

//

Query: gene_67_crAssphage

PSI-BLAST output (against viruses subset of NR database):

minor tail protein [Mycobacterium phage BrownCNA]

Sequence ID: gi|971761687|YP_009214946.1

Length: 765Number of Matches: 1

Score Expect Identities Positives Gaps

40.4 bits(93) 0.051 18/38(47%) 25/38(65%) 0/38(0%)

Query 237 LSNLTIEPIVNPKIWIGTATQYAAIAQKDNNTTYIVKS 274
++ T + N +W GTA QYAAIA K+ NT Y+VK+

Sbjct 728 VAGYTAAGLTNLNLWKGTAQAIAIATKNANTIYVVKV 765

//

Query: gene_63_crAssphage

PSI-BLAST output (against viruses subset of NR database):

Phage tail fibre adhesin Gp38 [uncultured Mediterranean phage uvMED]

Sequence ID: gi|775456205|BAQ90228.1

Length: 353Number of Matches: 1

Score Expect Identities Positives Gaps

47.0 bits(110) 2e-04 32/90(36%) 43/90(47%) 10/90(11%)

Query 13 GYNGAGMNYDGGNRRDV----NGKANAGLTLGIIGTALGAWALFGNRRSAGASILGGAG 67
G G +N +GG R D NG+ + LG GTA+G + A S G G

Sbjct 210 GGRGYSVNANGGTRGDYTTITNGNQKSGTRGLGPGSGTAIGGYGGSAGGGGANISTQAGGG 269

```

Query 68 GGGML-----GDGSTNINVFGATAGSGSGA 92
          GGGML      G+GS N +V      +G+ +GA
Sbjct 270 GGGMLITGTGGNGSANTSVGAGGSGNSAGA 299

```

Query: gene_64_crAssphage
 PSI-BLAST output (against viruses subset of NR database):

```

tail sheath protein [Vibrio phage KVP40]
Sequence ID: gi|34419589|NP_899602.1
Length: 671Number of Matches: 1

```

```

Score Expect Identities Positives Gaps
35.8 bits(81) 2.3 24/93(26%) 43/93(46%) 9/93(9%)
Query 59 DGQDLDTKSYVNSPINSLNYVAVGWSWNLSQIPYLYWGGDNTNSGAECVMFNIESMIDL 118
          +GQ + Y+ NS + +A GW Y + GG + N+GA+ MF ++ + D
Sbjct 318 NGQSIFIDEYFENSGSAYITAIAGWKTESGA--YNFGGSDANAGADDWMFGLDMLSD- 374

Query 119 EDKMPDIMKMNLCANWYGSLSRGHVTVECTAYK 151
          P+++ NL G+ + V++ T K
Sbjct 375 ----PEVLYTNLVI--AGNAAAEVSIAS TVQK 401

```

Query: gene_68_crAssphage_JQ995537
 HHpred output (against VOG database):

```

>[VOG6373]|KU160660-ALY09872.1| 21688..23658 + 656 aa|tail protein
Probab=88.77 E-value=0.4 Score=56.56 Aligned_cols=155 Identities=21% Similarity=0.308 Sum_probs=93.1

```

```

Q ss_pred          CCCCCcCHHHCCChHHHHhccccChHHhhhhhhcccCCCCcCCCCCcEEEEcCCCCCCCCcEEEEeCC-----EEEEe
Q lcl|gene_68_cr  464 LDGDGKVPASQLPSYVDDVLEGYVDETHFAEKYIEDAPVYYPTEKGIYVDISESTDYSGKTYRWSGT-----KYSVI 537 (696)
Q Consensus       464 LDasGKVPaaQLPSyVDDVLEgyYad~t~Fa~~~l~afP~~aTGEsGKIYVdld~gts~TNKtYRWSGS-----tYV~I 537 (696)
          -|~+|.||.+.|||.. .|-|.|-+... +..+..|...|+~|..| +++++|.~|-- -|.~+|
T Consensus       373 AdvtGti~Ta~LPPL--AiNe~~~vAsQ-----aaMLALTAQRGDMAIRSD-----nG~tYvLSsDsPgTLadWKEi 437 (656)
T [VOG6373]|KU16 373 ADISGTVP T S AL P PL--AVNDVFTVATQ-----AEMLALTAQRGDMAIRSD-----TGKSYALSTDSPGTLADWKEl 437 (656)
T ss_pred         hhccccccccCCcc--cccccccchH-----HHHHHHhhcCceeEecc-----CCceEEEEcCCCcchhhHHHH
Confidence        355899999999987 34443322211 22456888999999999 7888877632 34455

```

```

Q ss_pred          cCccccCcee-----eeccCCceeeechhcCCCCCceeE---EEeccCCceEEecceeE-ecccc--cccCCCC
Q lcl|gene_68_cr  538 SETLALGEVTG-----TAYDGKGGKKT T DIVNSLPKYIPSTQIK---LFRSVNGNIVIGSHHYE-FNNTT--NVYESKPF 606 (696)

```



```

Q ss_pred          cchhhcCccccccccccccccccccEeECCeEEEEEEcccccEEEcCCC
Q Q_1786342       71 KHKEKRGESLKGVIgDELtFGVNFPIVEYNDKRHLLINDQDIDYIIDGEE 120 (140)
Q Consensus       71 ~hkekr~eslkgiigdelT~vnfp~vey~krhlli~dqdidyiidgee 120 (140)
                    |.+          ++++++.++++.+.||=-+++++.
T Consensus       78 -----g~-----v~g~i~dIlAvi~ 104 (114)
T 4PJ1_X          78 -----GTK-----VVLDDKDYFLFRGDILGKYVDKL 104 (114)
T ss_dssp         -----CEE-----EETTEEEEEETTCCEEECC--
T ss_pred         -----CEE-----EEECCEEEEEEEhHHEEEECchh

```

//

Query: gene_17_IAS_virus_KJ003983
HHpred output:

>3ZK4_B DIPHOSPHONUCLEOTIDE PHOSPHATASE 1 (E.C.3.1.3.2); THREE-DOMAIN HEME-CU NITRITE REDUCTASE, ELECTRON; HET: FUC, GOL, NAG, PO4; 1.65A {LUPINUS LUTEUS}
 Probab=99.55 E-value=7.1e-15 Score=158.36 Aligned_cols=354 Identities=12% Similarity=0.002 Sum_probs=0.0
 Template_Neff=9.100

```

Q ss_pred          eeeeeeeCCCCCCEEEEEeecCC---CceecCCeEEEEecC---CCCceeeEeEeecCCCCcHHHHH-----HHHHHH
Q Q_4470287       658 TTHKCIvNGLTKGDYeyRVGRDND---PYVSEPLKFKVLAN---SDVTSFTYAXVTDQQGFNWAeyQ-----AWKKSA 725 (1147)
Q Consensus       658 ~~~~~ 725 (1147)
                    ..|. ....++.+. ..|. |.....  ..+. ....|+.+.  ....+. ....+.....  ..
T Consensus       198 ~h~l~L~p~t~Y~Y~vg~~~~~s~~~~F~t~p~~~~~f~~~gD~~~~~ 277 (571)
T 3ZK4_B          198 YIHTSFLKELWPNREYTYKLGHRLFNgtTIWSKEYHFkASPyPGQSSVQRVVI FGDMGKAeADGSNEYNnfQPGSLNtTK 277 (571)
T ss_dssp         EEEEEEEcscctTCEEEEEEEEETTSCEEECCCEEECCCTTCCSCEEEEEESCCCCCTTCcBCTTCcCTTHHHHHH
T ss_pred         eEEEEEEeCCCCCEEEEEEEecccCCceecceEEEEcCCCCCCCcEEEEEEecCCCCCCCccccCCCCcHHHHH

```

```

Q ss_pred          HHHHhCCCcCeEeecCccccCccchhhHhHhhhhHHHhcceeeCCCCccccCCCCccccCCCCccccceeecccc
Q Q_4470287       726 MMLSREEpDIQFTvNTGDITQSGNRVSEWLDYDGRQYLNNLVEMFTIGNNDLCGHNATELTNGEDATSKYSHINVLRYF 805 (1147)
Q Consensus       726 ~~~~~ 805 (1147)
                    .....+. ....+. ....|. ....+. ....+. ....+. ....+|..|.....
T Consensus       278 ~i~~~~~d~vl~GDl~y~g~w~d~f~~~~~l~P~~~~~GNHD~~~~~ 346 (571)
T 3ZK4_B          278 QIiQDLEDiDIVFHIGDLcYANGyISQWdQFTaQIEPIASTVpYMTASGNHERDWPgt-----GSFYGNLdSGG 346 (571)
T ss_dssp         HHHHTGGGCCEEEEEscscctTCTHHHHHHHHHHHHHTTSCEEEECCHHHHCCTTT-----TCSScCSTTTT
T ss_pred         HHHHhccCCEEEeccccCCcHHHHHHHHHHHHHHCCEEEeCCcCCcCCCC-----CccccCCCC

```

```

Q ss_pred          cccccCCccccCCcCCCCeEEecCceEEEEeeccccccccccccCCccccccchHhHccchHHHHHH
Q Q_4470287       806 TFELDSNTYEVewngVKYPIYSLYFNyGDFHFVSLNSEIAIASSKMYKDWSSDtyAGDRtFAEAANAQIEAWFIKDLQ 885 (1147)
Q Consensus       806 ~~~~~ 885 (1147)
                    .....+. ....+. ....+. |.+.+.+.+.+.+.+.+.  .....|..

```


The actual alignment was detected with superfamily member pfam13385:

Pssm-ID: 328935 Cd Length: 152 Bit Score: 33.51 E-value: 0.15
 10 20 30 40 50 60
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 lcl|seqsig_MGNIK_583 ITYDIFIIKVYVDGVLTAVTKETs-FPSLSDRVYFGGRIAGDKSTYLcDCNIYNFQLYDALTDFDIM 649
 Cdd:pfam13385 87 VTYDGGTLRLYVNGVLVGSSTLTGgPPSGTGGPLYIGRSPGGDDYF--NGLIDEVRIYDRALTAEEIA 152

//

**Query: Alignment is shown for CBL38474.1 CotH protein [Anaerostipes hadrus], the best hit to gene_19_IAS_virus_KJ003983
 CDD output:**

CotH protein; Members of this family include the spore coat protein H (cotH).

Pssm-ID: 331517 Cd Length: 319 Bit Score: 68.10 E-value: 1.01e-11
 10 20 30 40 50 60 70 80
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 gi 291559674 811 GTSSLQYAVKNYKIKLknpdgsKYKYSPFKNgileDTFCLKADYMESSHANNtgmaKFINDELYDTKVPPQQTNSkVRTA 890
 Cdd:pfam08757 5 GNSsREFPKksYRLKFD-----KGLERQRFPG---LRKLNLNAEFNDPSLLRN----KLAYDLFRDLGVPAPRARF-VELY 71
 90 100 110 120 130 140 150 160
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 gi 291559674 891 INGfpiqliakdsastpVYMGVFNF--NLDKGCNKsFGLDNEitggENCMSFEVSSNSDTsagafKNDTDESlrtdFEL 968
 Cdd:pfam08757 72 VNG-----EYYGLYLLveSVDKHFLLRAHGLDKD----GNLYKADDANFSLKL----SFGDPEK---YQLK 125
 170 180 190 200 210 220 230 240
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 gi 291559674 969 RYPDEDDCTSEQitekYNVLKRRLVTWVKNADETTFKNELEQYFNKEYLLKYFLQVHLFGMVDNLGKNMMLTTWDGN---I 1045
 Cdd:pfam08757 126 FEKELDEGGEED---WEDLSELIDFLNDTSEAEFEAELEKYIDVDSFLDWLAFNIIIGNTDSFSHNYLYRPDENgkWR 201
 250 260 270 280 290 300 310 320
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 gi 291559674 1046 WYPqfyDLDTQLGLDNTGYLKFYSDIDITEGVYNTSGSKLWTMV--ENVFADELSAMYKkLRtSKYRLDNILKYwYDGQV 1123
 Cdd:pfam08757 202 FIP--WDLDLAfgRDWRGIGLTLDfDEPEEGIANPEENVLFRRLldNPEFRARYIARLEELLDGVfTEERLEAK-IDALH 278
 330 340 350 360
*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
 gi 291559674 1124 AQIGELqYnkDMEAKYIKFKNDYLfmlhgRRSEhMKkWKerLLYL 1169
 Cdd:pfam08757 279 ALIAPA-LERDPQkwGGLEtATQYE---QEVEYLKDFIRQRrRYL 319

//

**Query: gene_20_IAS_virus_KJ003983
 CDD output:**

LRR_3 super family c127891 Leucine Rich Repeat;

Pssm-ID: 332712 Cd Length: 1153 Bit Score: 33.31 E-value: 1.0

```
          10          20          30          40          50          60          70          80
lcl|seqsig_MNIKI_ 130 DLTGCSKLRITIEINKCDSYKELRvsnlanletvkivscagiSSVIIINcpRLRTVNIIEYCDKLTTIQIN-NCKSLvggts 208
Cdd:PLN03210     652 DLSMATNLET LKLSDCSSIVELP-----SSIQYLN--KLEDLDMSRCENLEILPTGiNLKSL----- 706
          90          100         110         120         130         140         150         160
lcl|seqsig_MNIKI_ 209 dNYIRVANCNNINS-LDLSNNSnlkksidgcdrvkiETLKIHQTNITDVSSSTNLNDNMKLDLTEFS-----QLKT 279
Cdd:PLN03210     707 -YRLNLSGCSRLKSFpDISTN-----ISWLDLDETAIEEFPSNLRLLENLDELILCEMKseklwervQPLT 770
          170         180         190         200         210         220
lcl|seqsig_MNIKI_ 280 -FTCYYNKSVKYIAFANNQNApIPITSTFQECSNLER--IYGCVEL---SNTSYSG*YGL-FRGCSK 339
Cdd:PLN03210     771 pLMTMLSPSLTRLFLSDIPSL-VELPSSIQLNHKLEhleIENCINletlpTGINLESLESdLdLSGCSR 837
```

//

Query: gene_29_IAS_virus_KJ003983

CDD output:

>HTH_17 pfam12728 Helix-turn-helix domain; This domain is a DNA-binding helix-turn-helix domain.

Pssm-ID: 315411 Cd Length: 51 Bit Score: 33.18 E-value: 2.53e-03

```
          10          20
lcl|seqsig_MXKVI 45 LSXHQACQKLNVSRAFDFNLVREGKLP 71
Cdd:pfam12728   2  LTVEEAAELLGVSIRSTVYRLIRSGELP 28
```

//

Query: gene_42_IAS_virus_KJ003983

HHpred output:

>3CT0_A Morphogenesis protein 1; Cell wall, phi29, hydrolase, infection; HET: NAG; 1.77A {Bacteriophage phi-29}
Probab=99.47 E-value=1.8e-15 Score=138.45 Aligned_cols=136 Identities=22% Similarity=0.313 Sum_probs=0.0
Template_Neff=8.600

```
Q ss_pred           cccccChHHHHHHHHHHHHCCCCHHHHHHHHHHHHhCCCCCcC-----CCCCcceececcccHHHHHHHHhC--CCCC
Q Q_1802265         219 KINKVNPNAVRALNYFMNKGLTREQSAGLVGNLMAETGMNIRAVN-----PYSGAYGIAQWLGSRKTALFNKYG--NNPT 291 (472)
Q Consensus         219 ~~S~~knAk~I~~aLkk~G~S~a~AAGILGNi~q~ESGgNP~AvN-----~GG~A~GL~QWt~~Rf~aL~A~aG~-ni~N 291 (472)
```

```

          +.+.+.++++.|+++|+.+|+++..++|+||+||+||+||+|.+.+        .+++++||+|+|+.+.+.---++...  ++.+
T Consensus      10 ~~~~~~a~i~~~l~~~G~s~~~aagIlgn~~~ESg~~p~~~~~g~Gl~Qw~~~~~a~~~g~~~~ 89 (159)
T 3CT0_A         10 TMSEMKVNAQYILNYLSSNGWTKQAICGMLGNMQSESTINPLWQNLDEGNTSLGFGLVQWTPASNYINWANSOGLPYKN 89 (159)
T ss_dssp        CHHHHHHHHHHHHHHHHTCCHHHHHHHHHHHHHSSCTTCBGGGCTTCTTSCBTTTTBSShHHHHHHHHHHHTCCSSS
T ss_pred        ChhHhHhHHHHHHHHHHhCCCCCHHHHHHHHHHHhhcCCcCCCCccccCCCCcEEEEEEcCchHHHHHHhcccc

```



```

Q ss_pred        HHHHHHHHHhccCch-----HHHHHHhCCCHHHHHHHHHhccccCChhHHHHHHHHhCCCCcC
Q Q_1802265     292 LDQQLDIFIWHELNSSHS-----RGLRMLRQSNNPSDAAANAFGYEFSAQPLQAVRAMNAAGKNTKWK 354 (472)
Q Consensus     292 ~dQL~fainEl~y~-----s~lk~Lk~a~d~~AA~af~~~yErpgg~~~~~R~~~A~~~~~k~~ 354 (472)
                ...|+|.+++++. .... . ....+. ....+. ....+. |+.|+|+|+ .....|+|.+.+.+++
T Consensus     90 ~~~Q~~~~~e~~~~~a~~~~~yEr~~~~~R~~~A~~~~~ 159 (159)
T 3CT0_A        90 MDSELKRRIWEVNNAQWINLRDMTFKEYIKSTKTPRELAMIFLASYERPANPNQPERGDAQEYWKYKNSL 159 (159)
T ss_dssp        HHHHHHHHHHHHTCCSSCCHHHHHHTCCSCHHHHHHHCCSSCCTTTHHHHHHHHHHHHCC
T ss_pred        HHHHHHHHHhCCcCCCCcCcchHHHHHHhCCCHHHHHHHHHhCCSSCCTTCHHHHHHHHHhCCc

```

>3CB7_A Lys-rich lysozyme 2 (E.C.3.2.1.17); Digestive lysozyme 2, Musca domestica; HET: ACY; 1.9A {Musca domestica}
 Probab=94.85 E-value=0.024 Score=51.02 Aligned_cols=51 Identities=24% Similarity=0.143 Sum_probs=0.0
 Template_Neff=8.300

```

Q ss_pred        HHHHHHHhCCCHHHHHHHHHHHhCCSScCC---CCcCEEEccCH
Q Q_1802265     229 RALNYFMNKGLTREQSAGLVGNLMAETGMNIRAVNP---YSGAYGIAQWLGSRK 279 (472)
Q Consensus     229 ~I~~~aLkk~G~S~a~AAGILGNiQ~ESGgNP~AvN~---GG~A~GL~QWt~Rf 279 (472)
                .|.++++.+|.+.+.+.+.+.++++|+|+|.++++. ++.++|+|+.+.+.++
T Consensus     11 ~l~~~~~g~::~l~::ia~~ES~f~~~~~g~~~~~Gl~Qi~~~~~ 65 (126)
T 3CB7_A        11 SLAREMYKLGVPKNQLARWTCTIAEHSSYNTKAVGSLNSNGSRDYGIFQINNYW 65 (126)
T ss_dssp        HHHHHHHHTTCCGGGHhHHHHHHHHHHSSBTTCCBCCCCTTSCCEETTTTEETTTT
T ss_pred        HHHHHHHHHhCCCHHHHHHHHHHHhCCcCCCCcCCCCcEeeeeEcccee

```

//

Query: gene_50_IAS_virus_KJ003983

HHpred output:

>3UQZ_A; SAM and Rossmann Fold, DNA; HET: S04; 2.7A {Streptococcus pneumoniae}
 Probab=100.00 E-value=4.4e-37 Score=323.75 Aligned_cols=268 Identities=12% Similarity=-0.037 Sum_probs=0.0
 Template_Neff=8.400

```

Q ss_pred        CCCCCCCCCChhHHHHhCCccccccccCChhHHHHHHhHHHHHHHHhccCchH-hhCCccccCCCC
Q Q_5475696     134 XXXLAILRGLDNLPISTQQFYDPAKNGTWEKGRWNKEDIKFKETIDAEFEAIKEALNSGKYDR-IVLPPIESLFHEKS 212 (837)
Q Consensus     134 ~Ls~L~GIG~~~~~L~f~~~~~a~e~~~~~e~L~ki~~~~~e~k~~~~~l~i~i~~De~ 212 (837)
                +..|.. |+|..+.++.+++... ..+..+. .....+.....+. .....+.+.++++|+

```

```

T Consensus      8 ~~~~~~gig~~~~~-----~l~~~~~l~~~~~-----it~d~ 80 (288)
T 3UQZ_A        8 IYKLLK-SGLTNQQILKLVLEYGENV-----DQELLGLADIADISGCRNPAVFMERYFQIDDAHLSKEFQKFPFSILDDC 80 (288)
T ss_dssp
T ss_pred
H H H H H H - T T C C H H H H H H H H H H G G G - - - - - T T B C H H H H H H H T T C S C H H H H H H H H H H C C H H H H H H H S S C E E E T T S T T
H H H H H h - C C C C H H H H H H H H H h h c C - - - - - C h H H H h c C H H H H h c C C C h H H H H H H h c C C H H H h H H H c C C C E e e C C C C

Q ss_pred      ccHHHhcCCCCeEEEEecCHHhc----eEEEeCCCCCHHHHHHHHHHHHHhhCCCeEEEcCCCCCHHHHHHHHHHHhC
Q Q_5475696    213 LPAKGKKGDKTISISDISKERVPKL---YEYLNKYEEVFGELEVEDKKEKNTTSVTKIISGGQTGVDITGLQVAKELD 288 (837)
Q Consensus    213 YP~~LkeI~dpP~VLY~kGnl~l----IAIVGTRk~S~YG~~~A~~ia~~LA~~Git~IVSGLA~GIDtaAH~gALeaG 288 (837)
|. + | + + + + | + + | + | + | + + . +      | | | | | + + + . + | . . . + + . + + . | + . | + + | | | + | . | | | . + | + + | + + + |
T Consensus    81 YP~~L~~~~~pP~~Lf~~G~~ll~~~~iaIvGsR~~~~~l~~~~l~~~~l~~~~~ivsG~a~GiD~~a~~~~al~~ 158 (288)
T 3UQZ_A      81 YPWLSEIYDAPVLLFYKGNLDLLKFKVAVVGSRACSKQGAKSVEKVIQGLE-NELV-IVSGLAKGIDTAAHMAALQNG 158 (288)
T ss_dssp
T ss_pred
c c H H H H h c C C C C E E E E E e c C h h h c C C c E E E e c C C C C H H H H H H H H H H H H H H - C C C E - E E c C C C C c H H H H H H H H H H C C

Q ss_pred      CCEEEEcCCCcCccchhhHHHHHHHHHCCEEEEcCCCCCCCccccHHHHHHHhccCCEEEEEcCCCCcHHHHHHHH
Q Q_5475696    289 VETGGTAPKAFLEEGIDKEDVRSYGLTEITDEEQEY TartGKKDNYTGRTDLNVKNSDGTVYFNyGNDSTGLKATRRS 368 (837)
Q Consensus    289 G~TIAVLg~Glld~yP~n~L~I~e~G~lISE~pPgt~p~~~~Fp~RNRIIAGLSdgtVVVEAg~kSGGSLiTA~ 368 (837)
+ + | + | + | + | . + . | | . . + + . + + + | + + + | + + + | + . + . + + + | + . | | + + + + + | + + + | + + + . + | + + . | + + +
T Consensus    159 g~~I~Vlp~g~l~~~~p~~~~~i~~~~lviSe~~~~~RNR~i~~lsd~vivv~~~~~G~T~~~~~ 236 (288)
T 3UQZ_A      159 GKTIAVIGTG-LDVFYPKANKRLQDYIGNDHLVLSeyGPGEQPLKFHFpARNRIIAGLCRGVIVAEAKMRSG-SLITCER 236 (288)
T ss_dssp
T ss_pred
CCEEEcSSc-TTCCSSGGGHHHHHHHHHSEEEESSCTTCCCTTHHHHHHHHHHHHCSEEEESCCTTCH-HHHHHHH
CcEEEECCC-hhhhcCHHhHHHHHHHHHCCEEEEcCCCCCchhhHHHHHHHHhCCEEEEEcCCCH-HHHHHHH

Q ss_pred      HHHhCCCEEEcCCCcChhChhccchccccCccccchHHHHHHHHHH
Q Q_5475696    369 AEEHNKPFLLNPTAKELRQWIKDNNIKTLNVAGNRGSKLAKNNNVAKTLRDAL 421 (837)
Q Consensus    369 AleqGR~VfAVPG~i~~~~s~G~N~I~ILNvAGpreS~~agI~e~a~dIle~L 421 (837)
| + + + | | | + + + | + . . . + . + | + . + . . . + + . + + . + . . . + + . + +
T Consensus    237 A~~~gk~V~~~~~g~~~Li-----~Ga~v~~~~~ 280 (288)
T 3UQZ_A      237 AMEEGRDVFAIPGSILDGLSDGCHHLI-----QEGAKLVTSgQDVLAEF 280 (288)
T ss_dssp
T ss_pred
HHHTTCEEEECSSSSTTHHHHHHH-----HTTCEECSSHHHHHHHC
HHHcCCCEEEcCCCcCCcCHHHHHHH-----HCCCEEcCHHHHHHHH

```

Query: gene_52_IAS_virus_KJ003983
CDD output:

vWFA cd00198 Von Willebrand factor type A (vWA) domain was originally found in the blood coagulation protein von Willebrand factor (vWF). Typically, the vWA domain is made up of approximately 200 amino acid residues folded into a classic a/b para-rossmann type of fold. The vWA domain, since its discovery, has drawn great interest because of its widespread occurrence and its involvement in a wide variety of important cellular functions. These include basal

membrane formation, cell migration, cell differentiation, adhesion, haemostasis, signaling, chromosomal stability, malignant transformation and in immune defenses In integrins these domains form heterodimers while in vWF it forms multimers. There are different interaction surfaces of this domain as seen by the various molecules it complexes with. Ligand binding in most cases is mediated by the presence of a metal ion dependent adhesion site termed as the MIDAS motif that is a characteristic feature of most, if not all A domains.

Pssm-ID: 238119 Cd Length: 161 Bit Score: 54.88 E-value: 1.94e-09
 10 20 30 40 50 60 70 80
 *.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MRTNL 23 LDMVIAFDTTGSMAS-SYIKAVRKHVIELIPKLFANPKLKISIVAFGdycdmnsksdfgNAYQVIDLTD--NKDELIKFV 99
Cdd:cd00198 1 ADIVFLLDVGSMGgEKLDKAKEALKALVSSLSASPPGDRVGLVTFGS-----NARVVLPLTTdtDKADLLEAI 69
 90 100 110 120 130 140 150 160
 *.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MRTNL 100 KNARNTSGGDGDefYELVIKKIVEETSW--REGSTKSVLLIADAYPHEvgysyrasisesylienNQIDWREEAKKAAAK 177
Cdd:cd00198 70 DALKKGLGGGTN--IGAALRLALELLKSAkRPNARRVIIILLTDGEPND-----GPELLAEAAARELRKL 130
 170 180
 *.....|.....*.....|.....*.....|
lcl|seqsig_MRTNL 178 GIKIDTMQCSNTHNSIWYKE LSDITNGIN 206
Cdd:cd00198 131 GITVYTIGIGDDANEDELKEIADKTTGGA 159

//

Query: gene_55_IAS_virus_KJ003983
CDD output:

Thy1 pfam02511 Thymidylate synthase complementing protein; Thymidylate synthase complementing protein (Thy1) complements the thy

Pssm-ID: 308231 Cd Length: 186 Bit Score: 66.88 E-value: 2.16e-13
 10 20 30 40 50 60 70 80
 *.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MKLIK_ 137 YTVHFI-TSRVTMDSFRTHITLSHLGESTRYCNYNKDkfdnqlTFIIPDrydieecdkytATDFYNSEDLTD EspkynWL 215
Cdd:pfam02511 52 FTFAIEgVSRVAVLRQLVRRHRIASFSQQSQRVVKLDDE-----DFVIPP-----EIAKAQSPLELLEL----YE 109
 90 100 110 120 130 140 150
 *.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....
lcl|seqsig_MKLIK_ 216 NAMINAELTYMRL LKLPQYARGVLPDVKSELISCGf kDA--WDNFFDKRCAN DAHPMAREIATAVRNKLQ 287
Cdd:pfam02511 110 EAMEEAYEAYEELLEKGVARE DARYVLPNATETRIVVTM--NarsLLHFLELRCCPRAQWEIRELAEAMLEELK 181

//

Query: gene_60_IAS_virus_KJ003983

of the ATP-dependent DNA ligase family. The catalytic core contains six conserved sequence motifs (I, III, IIIa, IV, V and VI) that define this family of related nucleotidyltransferases including eukaryotic GRP-dependent mRNA-capping enzymes. The catalytic core contains both the active site as well as many DNA-binding residues. The RNA circularization protein from archaea and bacteria contains the minimal catalytic unit, the adenylation domain, but does not contain an OB-fold domain. This family also includes the m3G-cap binding domain of snurportin, a nuclear import adaptor that binds m3G-capped spliceosomal U small nucleoproteins (snRNPs), but doesn't have enzymatic activity.

Pssm-ID: 325160 Cd Length: 325 Bit Score: 37.88 E-value: 3.66e-03

```

          10          20          30          40          50          60          70          80
lcl|seqsig_MEYQK_ 36 WECTEKIDGTNIRIYVTMEA-----GEGENPWLYGVTIKGRTNraelpsklVKKLENIFlkvdWAKvfpALTPE 104
Cdd:TIGR02307    27 WVAREKIHGTNFSIIIERDFkvtcakrtgiiLPNEFFGYHILIKNYTAS-----VKAIQDIL---ETK---AIIIV 92
          90          100         110         120         130         140         150         160
lcl|seqsig_MEYQK_ 105 DTVCIYGEYAGIQQKvgsKYIKNDVNFILFDVKFN-----DWLLKREDCEDIAkkcNVDIVSLIGYMTIPQATEFV 176
Cdd:TIGR02307    93 VSVQVFGELAGPGYQK---PVVYSKDFYAFDIKYTetsddvtlVDDYMMESFCNVP---KLKYAPLLGRGTLDELLAFD 166
          170         180         190         200         210
lcl|seqsig_MEYQK_ 177 KNGFKSRISEDKDLDD-----AEGLVLRRT-CGLRFRNGERIITKIKHCDFEKFK 224
Cdd:TIGR02307    167 VENFTTDHPALVDAGnyplegntAEGYVVKHCrPGKWLNRNGNRTIICKNSKFSEKK 223

```

Query: gene_64_IAS_virus_KJ003983

CDD output:

DUF3310 super family c113237 Protein of unknown function (DUF3310); This is a family of conserved bacteriophage proteins of unknown function.

Pssm-ID: 314594 Cd Length: 58 Bit Score: 53.05 E-value: 2.93e-11

```

          10          20          30          40          50          60          70
lcl|seqsig_MEKKD_ 8 VEHPSHYiwLKGicGIEVIDITRHM-----NFNLGNVIKIVLRSCHKSeggmsdkqKQIEDLKKARFY 70
Cdd:pfam11753    1 VNHPEHY--GAG--GIECIDVIRAQltgeefkGFCLGNAIKYLSRAGKKN-----GEEEDLKKAKWY 58

```

See <https://www.ncbi.nlm.nih.gov/pubmed/20497505> for protein experimental characterization

Query: gene_65_IAS_virus_KJ003983

CDD output:

DNA_pol_A super family c102626 Family A polymerase primarily fills DNA gaps that arise during DNA repair, recombination and replication; DNA polymerase family A, 5'-3' polymerase domain. Family A polymerase functions primarily

to fill DNA gaps that arise during DNA repair, recombination and replication. DNA-dependent DNA polymerases can be classified into six main groups based upon phylogenetic relationships with E. coli polymerase I (class A), E. coli polymerase II (class B), E. coli polymerase III (class C), euryarchaeota polymerase II (class D), human polymerase beta (class X), E. coli UmuC/DinB and eukaryotic RAP 30/Xeroderma pigmentosum variant (class Y). Family A polymerases are found primarily in organisms related to prokaryotes and include prokaryotic DNA polymerase I, mitochondrial polymerase gamma, and several bacteriophage polymerases including those from odd-numbered phage (T3, T5, and T7). Prokaryotic polymerase I (pol I) has two functional domains located on the same polypeptide; a 5'-3' polymerase and a 5'-3' exonuclease. Pol I uses its 5' nuclease activity to remove the ribonucleotide portion of newly synthesized Okazaki fragments and the DNA polymerase activity to fill in the resulting gap. The structure of these polymerases resembles in overall morphology a cupped human right hand, with fingers (which bind an incoming nucleotide and interact with the single-stranded template), palm (which harbors the catalytic amino acid residues and also binds an incoming dNTP) and thumb (which binds double-stranded DNA) subdomains.

```

Pssm-ID: 322025 Cd Length: 593 Bit Score: 196.33 E-value: 1.66e-54
      10      20      30      40      50      60      70      80
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 16 NYKVIgVEESL----KLLDPLTIVGLDTEtTGLDPWTKELKSIQLGN-----YDFQVVIDTtTINPTLyKEYLES-NRL 84
Cdd:COG0749      2 PYGTITDLAVLnawlTKLNAAANIAFDtETDGLDPHGADLVGLSVASeeeaayIPLLHGPEQLNVLAAL-KPLLEDeGIK 80
      90      100     110     120     130     140     150     160
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 85 FVGWNLKFDLkFLFRQNIvLKNVWDGYLAEKLMwLGYPPGIHSLSLKAagENyLDIEL--DKSVRGKIIYAGLTEDV--- 159
Cdd:COG0749      81 KVGQNLKYDYKVLANLGIePGVAFDtmLASYL--LNPGAGAHNLDDLA--KRYLGLETitFEDIAGKGGKQLTFADVkle 156
      170     180     190     200     210     220     230     240
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 160 --IAYSANDVKYLEKIMKlQRiELAKKGLEKAIiYEN--KFVLPLAYCEYCGIKLDADKWKAKMQKDKQRVTtaldncnk 235
Cdd:COG0749      157 kaTEYAAEDADATLRLESiLEPELLKTPVLELYEEIemPLVRVLARMERNGIKVDVQYLKELSKELGCeLAE----- 229
      250     260     270     280     290     300     310     320
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 236 wLENEpNSEYifidRqgdifnGFNLEPQvKLNWNsAKQL--IpLFKKYGVNVTTeDKVKGGTKdsIDAKSLKPQKDKCS 313
Cdd:COG0749      230 --LEEE-----IYELAGE-EFNINSPKQLgeI-LFEKLG LPPGL-KKTKTGNy-STDAEVLEKLADdHP 287
      330     340     350     360     370     380     390     400
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 314 LIPLYLEYKEAIKVTSTYGENFLKQINPVSGRIHTNYQQMGADTTRLTsGgKDKnakveyvNLLNLpANAET---RACF 389
Cdd:COG0749      288 LPKLILEYRQLAKLKSTYTDGLPKLINPDTGRIHTSFNQGTATGRlSS--SDP-----NLQNIPIRSEegrkiRKAF 358
      410     420     430     440     450     460     470     480
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 390 VAENGNKWISIDYSGQETYLMAsiANDEAIiKELTEGSgDIHSLTAyMSYHEIPRDtniKDIKKkyhDLRQDAKGIeFAI 469
Cdd:COG0749      359 VAEKGYTLISADYSQIELRiLAHLSQDEGLLRafTEGE-DIHTATA-AEVFGVPIE----EVTS---EQRRKAKAINFGL 429
      490     500     510     520     530     540     550     560
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 470 NYGGDANTISKNGKIPIEEAKKIYNAYMAGfKGLKRYQDFRRKDWFNKGyIllNPLTGhkAYIYDykelLEDKKWmatld 549
Cdd:COG0749      430 IYGMSAFGLAQLGIPRKEAKYIDRYFERYPGVKEYMERTKEEAREdGYV--ETLFGRRRYLPD----INSSN----- 497
      570     580     590     600     610     620     630     640
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MIYLV_ 550 wdyYremkiacpecetvqrVrhfFKRKSASEKQSINyPIQATGSMCLRVSMINFFeYLRsNNLLfKVLICVtpYDEINCE 629

```


Query: gene_67_IAS_virus_KJ003983

CDD output:

trimeric_dUTPase super family cl00493
Trimeric dUTP diphosphatases; Trimeric dUTP diphosphatases, or dUTPases, are the most common family of dUTPase, found in bacteria, eukaryotes, and archaea. They catalyze the hydrolysis of the dUTP-Mg complex (dUTP-Mg) into dUMP and pyrophosphate. This reaction is crucial for the preservation of chromosomal integrity as it removes dUTP and therefore reduces the cellular dUTP/dTTP ratio, and prevents dUTP from being incorporated into DNA. It also provides dUMP as the precursor for dTTP synthesis via the thymidylate synthase pathway. dUTPases are homotrimeric, except some monomeric viral dUTPases, which have been shown to mimic a trimer. Active sites are located at the subunit interface.

Pssm-ID: 294336 Cd Length: 142 Bit Score: 70.73 E-value: 5.22e-16
10 20 30 40 50 60 70 80
lcl|seqsig_MKIKV_ 35 DLRAAEDYEFEEApqasilhqkdgiktgdvkdFTKVISLGLAIQLPKGLVGRIVERSS----GVVKLNikkmgGGYIDNC 109
Cdd:TIGR00576 25 DLRAAEDVTIPP-----GERALVPTGIAIELPDGYYGRVAPRSGlalkhGVTIDN----SPGVIDAD 82
90 100 110 120 130 140 150
lcl|seqsig_MKIKV_ 110 YRGDndiwkIPV---TSMVKQTIHKGDRIQFHIELsqfatpwqklkwLFSKPKFIFVDSFNENNNR-GGFGSTGV 181
Cdd:TIGR00576 83 YRGE-----IKVilinLGKEDFTVKKGDRIAQLVVEK-----IVTEVEFEVEEELDETERGeGGFGSTGV 142

////////////////////////////////////

Query: gene_69_IAS_virus_KJ003983

CDD output:

RNR_activ_nrdG3 TIGR02826 anaerobic ribonucleoside-triphosphate reductase activating protein; Members of this family represent a set of radical SAM enzymes related to, yet architecturally different from, the activating protein for the glycine radical-containing, oxygen-sensitive ribonucleoside-triphosphate reductase (RNR) as described in model TIGR02491. Members of this family are found paired with members of a similarly divergent set of anaerobic ribonucleoside-triphosphate reductases. Identification of this protein as an RNR activating protein is partly from pairing with a candidate RNR. It is further supported by our finding that upstream of these operons are examples of a conserved regulatory element (described Rodionov and Gelfand) that is found in nearly all bacteria and that occurs specifically upstream of operons for all three classes of RNR genes. [Purines, pyrimidines, nucleosides, and nucleotides, 2'-Deoxyribonucleotide metabolism]

Pssm-ID: 274317 Cd Length: 147 Bit Score: 150.19 E-value: 2.51e-47
10 20 30 40 50 60 70 80
lcl|seqsig_MLKYV_ 3 KYVDSKVVFAEIPDEITLAINISNCPCHCEGCHSSYL-AEDIGEPLDLQHLTLNIDSNG-ITCVCIMGGDANPSEVDDI 80
Cdd:TIGR02826 1 KYVNEDIVFQVPEVSLAFYISGCPLGCPGCHSPELwHEDEGTPLTPEVLAQLLDKYRSIITCVLFLGGEWEPEALLSL 80
90 100 110 120 130
lcl|seqsig_MLKYV_ 81 AQDIKEyYPNLKVGWYSGRDYISKDIDMSNFNYIKYGHYDKDKGPLNSKTTNQVMLEI 138

Cdd:TIGR02826 81 LKYVKE-HAGLKVCLYTGREPDKPLELVQHLDYDKTGPWVETLGGDLSPTTNRQFQDI 137

//

Query: gene_70_IAS_virus_KJ003983

CDD output:

NRDD super family cl26133 Anaerobic ribonucleoside-triphosphate reductase;

Pssm-ID: 330954 Cd Length: 623 Bit Score: 612.78 E-value: 0e+00

	10	20	30	40	50	60	70	80	
lcl seqsig_MKVLK_	63	IQMEVEKILMDLAPYNVAKAYIIYRNKHEESRFIRERIDYMSNYADSDqNAASSSETDPNANVTQKNVANLDGEVYKTKN	142						
Cdd:PRK08271	1	MKVQVSEFEPGFEEELYSKLKSKYGYEMLRLEGIQIQLKVGFISEYFFAK-NAADGSKIDANANVRHKNIATYEAELMKDFF	79						
	90	100	110	120	130	140	150	160	
lcl seqsig_MKVLK_	143	RIIQRQRMKDELNVLY-PEVAKQYEIDVENHIIYPHDEASVptLKFYCQADSLYPLMTEGVGNIDGVTpSPNDLQSFSG	221						
Cdd:PRK08271	80	KLINRYLVWNKIKELFgKELADEYLRQIENHEIYVHDETS---LKPYCFSTYMTPILEDGLTKIGGES-KAPKHLSSFCG	155						
	170	180	190	200	210	220	230	240	
lcl seqsig_MKVLK_	222	QITNLTFLSSQCKGAVAFGEYFIALNYIIAEEFGDKWYEKldcvtvnshckvqrTVRDFIEKAFKQFIYGINQPAGNRS	301						
Cdd:PRK08271	156	SFINFVFAVSSQFAGAVATVEFLVYFDYFARKDYGDDYLP-----THRKEIENHLQIVVYSLNQPAARG	220						
	250	260	270	280	290	300	310	320	
lcl seqsig_MKVLK_	302	YQSPFTNVSYYDHTYFSSLFGEFYYPDGTPPEWAAVNVLQKFMFKFFNKLRTKQILTFPVETLAMVHDGKDIIDKEYKDF	381						
Cdd:PRK08271	221	YQSVFWNISYYDRNYFKAMFGGFVYPDGSTPNWEDIIALQKFFMEWFNKEREKAMLTFFPVVTAALLTDDGKCKDEDFADF	300						
	330	340	350	360	370	380	390	400	
lcl seqsig_MKVLK_	382	CAEMYAEGHSFFTYISDSADSLASCCRLRNEIENFNPSTGLTGVMTGschvITLNNRIVQDCNkayglkrnggwkeN	461						
Cdd:PRK08271	301	IAKENSKGNSFFIYISDSADSLASCCRLRNEISDNGFSYSLGAGGVNTGSINVITINLPRIAQEAR-----D	367						
	410	420	430	440	450	460	470	480	
lcl seqsig_MKVLK_	462	TSFIRDYLISILDRVYKHYAIYKTMLYEQEEKGMFAACNGGYIHMSKLYSTIGINGLNEAARFLGLKVSNNPEYIKFLQL	541						
Cdd:PRK08271	368	RDDFLEILRERVDKIHKYQLAYREIMEERIAAGMLPLYDAGFISLDKQFLTIGINGMVEAAEFMGLTVGYNEEYKDFVQE	447						
	490	500	510	520	530	540	550	560	
lcl seqsig_MKVLK_	542	ILGTIKEANKHSHhdkSRPFLFNSEVVPaesLGGKNYRWdkkDGyWVPEdenLYNSYFFDAHD-DTSVLDKMLHGRQT	620						
Cdd:PRK08271	448	VLKVIYEANEKAS--KEYGTFNTEFVPAENLGVKLAKWDREDGYGVPRQ--CYNsYSyVVEDaNTDALDKFKLHGKEL	522						
	570	580	590	600	610	620	630	640	
lcl seqsig_MKVLK_	621	AQYCDGGSACHINLEDHLSKEQYLKLIeFAVKEGTYFTFNIPNSKCDDCGYITKHPITeCPKCHSHNITWYTRVIGYLR	700						
Cdd:PRK08271	523	DKYLSGGSALHNLDERLSEEGYRKLNNIAAKTGCNYFafNVKItICNDCHHIdKRTGKRCPICGSenIDYYTRVIGYLR	602						
	650								
lcl seqsig_MKVLK_	701	PIKAFGIDRFIEAGKRVY	718						
Cdd:PRK08271	603	RVSafSKVRQkeyPRRH	620						

Query: gene_73_IAS_virus_KJ003983

CDD output:

DUF1064 super family c105706 Protein of unknown function (DUF1064); This family consists of several phage and bacterial proteins of unknown function.

Pssm-ID: 283910 Cd Length: 117 Bit Score: 33.99 E-value: 0.01

.....*.....|.....*.....|.....*.....|.....*.....|.....*.....
lcl|segsig_MPNNK_ 91 GIRYTPDFYFRY-GKLDVYIEAKGIENDVFYIKKKLFR-KFLDDKL 134
Cdd:pfam06356 61 KIKYIADFLIYHnDGLLEEVIDVKGMAKTDANIKRKLFDyKYRQVKL 106

Query: gene_75_IAS_virus_KJ003983

CDD output:

UDG_like super family c100483
Uracil-DNA glycosylases (UDG) and related enzymes; Uracil-DNA glycosylases (UDG) catalyzes the removal of uracil from DNA, which initiates the DNA base excision repair pathway. Uracil in DNA can arise as a result of mis-incorporation of dUMP residues by DNA polymerase or via deamination of cytosine. Uracil in DNA mispaired with guanine is one of the major pro-mutagenic events, causing G:C->A:T mutations. Thus, UDG is an essential enzyme for maintaining the integrity of genetic information. At least five UDG families have been characterized so far; these families share similar overall folds and common active site motifs. They demonstrate different substrate specificities, but often the function of one enzyme can be complemented by the other. Family 1 enzymes are active against uracil in both ssDNA and dsDNA, and recognize uracil explicitly in an extrahelical conformation via a combination of protein and bound-water interactions. Family 2 enzymes are mismatch specific and explicitly recognize the widowed guanine on the complementary strand, rather than the extrahelical scissile pyrimidine. This allows a broader specificity so that some Family 2 enzymes can excise uracil as well as 3, N(4)-ethenocytosine from mismatches with guanine. A Family 3 UDG from human was first characterized to remove Uracil from ssDNA, hence the name hSMUG (single-strand-selective monofunctional uracil-DNA glycosylase). However, subsequent research has shown that hSMUG1 and its rat ortholog can remove uracil and its oxidized pyrimidine derivatives from both, ssDNA and dsDNA. Enzymes in Families 4 and 5 are both thermostable. Family 4 enzymes specifically recognize uracil in a manner similar to human UDG (Family 1), rather than guanine in the complementary strand DNA, as does E. coli MUG (Family 2). These results suggest that the mechanism by which Family 4 UDGs remove uracils from DNA is similar to that of Family 1 enzyme. Although Family 5 enzymes are close relatives of Family 4, they show different substrate specificities.

Pssm-ID: 320999 Cd Length: 201 Bit Score: 134.17 E-value: 2.52e-39

.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|

TOPRIM_primases cd01029 TOPRIM_primases: The topoisomerase-primase (TORPIM) nucleotidyl transferase/hydrolase domain found in the active site regions of bacterial DnaG-type primases and their homologs. Primases synthesize RNA primers for the initiation of DNA replication. DnaG type primases are often closely associated with DNA helicases in primosome assemblies. The TOPRIM domain has two conserved motifs, one of which centers at a conserved glutamate and the other one at two conserved aspartates (DxD). This glutamate and two aspartates, cluster together to form a highly acid surface patch. The conserved glutamate may act as a general base in nucleotide polymerization by primases. The DXD motif may coordinate Mg²⁺, a cofactor required for full catalytic function. The prototypical bacterial primase. Escherichia coli DnaG is a single subunit enzyme.

```
Pssm-ID: 173779 Cd Length: 79 Bit Score: 48.81 E-value: 1.33e-08
      10    20    30    40    50    60    70    80
lcl|seqsig_LCKSK 41 NRICICSSLKDALCLWANTGIPSLAIQEGGyrMSDTAINELKRFNRIYICLDNDEAGLK---DAIQLASKTGFiNVVLP 117
Cdd:cd01029     1 DEVIIVEGYMDVLALHQAGIKNVVAAALGTA--NTEEQRLRLKRFARTVILAFDNDEAGKkaaRALELLLALGG-RVRVP 77

      ..
lcl|seqsig_LCKSK 118 QF 119
Cdd:cd01029     78 PL 79
```

////////////////////////////////////

Query: gene_83_IAS_virus_KJ003983

CDD output:

AAA_24 pfam13479 AAA domain; This AAA domain is found in a wide variety of presumed phage proteins.

```
Pssm-ID: 316040 Cd Length: 195 Bit Score: 74.22 E-value: 3.70e-16
      10    20    30    40    50    60    70    80
lcl|seqsig_MIQLP 15 NPKLMVIFGKPKSGKSSFVAVIDDNLIIIDLEDGYRALAVMKVQarTARDLQEIRDIVTKGRELHkaPYKFITIDNATRL 94
Cdd:pfam13479   1 KKLKILIYGPSGIGKTTFAKTLPKPLFIDTEGGTKSLDGD RFP--YIRSWQDFLDIIVELAAELA--DYKTIVIDTIDWA 76
      90    100   110   120   130   140   150   160
lcl|seqsig_MIQLP 95 EEMsvalaaelyratpmgasWGYMTDAKGmvIKNPKTGK PmvdpkadvrqlaNGAGWLYMRKAIRQLVDMFKNLCETLIL 174
Cdd:pfam13479   77 ERL-----CLAYVCRQNG--KGSIEDGG-----YGKGYGELAE EFRRLLDKQLQELGKNVIF 126
      170   180   190   200
lcl|seqsig_MIQLP 175 VCHVKDKQIKKNGEEM-SEM*VDLA**TGDIICGEADAVGYLYRDG 219
Cdd:pfam13479   127 TAHAKTRKDED P DGGKyTRYEPKLGKKTANLLKEWVDLVL FANYKT 172
```

////////////////////////////////////

Query: Alignment of 12 closest homologs of gene_85_IAS_virus_KJ003983

HHpred output:


```

Q ss_pred          cCcccCCCCCCCCCHHHhCCCchHHHhCCEEEEEEECCcchhhccccccCCCcchhhhhhhcCCCcCCcCCcCCcE
Q gene_186_Cellu  214 VGRIA EKSN S--ASP NAGD VFGSSFLDQLCSFNII LYDAF--KMGISQYMKVNPDRYDY LSEYFGDEDSK GKVSFETENK 289 (368)
Q Consensus       336 ~~~~~g~p~l~di~gS~I~~AD~vi~L~r~~~~~g~~~~~ 415 (519)
      .....  .|.+.+++|||.+.+.|||.|.+++..  .....  ..+  .
T Consensus       364 ~~~~~gs~~~~ad~vi~l~~~~~ 407 (444)
T 3BGW_F          364 EQRQD-----KRPMLS DLRESGQLEQDADI IEF LYRDD--YYDK ESE-----SKN-----I 407 (444)
T ss_dssp         GGSSC-----CCCCGGGCCSCSHHHHCSEEEECBGG--GTCSSCS-----STT-----E
T ss_pred         hhCCC-----CCCChHHHhccCchHHhCCEEEEEEECHH--HcCCccc-----CCc-----e

```

```

Q ss_pred          EEEEEEEcCCCCCCCCceeeeeeeeeEeccCCEEEcCCcCCc
Q gene_186_Cellu  290 IFVHLIKTRES DTPYKDI FVIEKDVSEESKARVSFKEKTPSGPKFT 335 (368)
Q Consensus       416 i~l~i~KnR~G~~~~g~v~l~~~~~fd~~~~f~e~~~~~ 461 (519)
      ..+.+.|+|.|.  + .+.+.  |+++++|.++++.+.
T Consensus       408 ~~l~i~K~R~g~~~~~f~~~~~ 444 (444)
T 3BGW_F          408 VEVIIAKHRDGPV---G-TVSLA-----FIKEYGNFVNLERRFDDR 444 (444)
T ss_dssp         EEEEECCSSSC---E-EEEE-----EETTTTEEECC-----
T ss_pred         EEEEEEEcCCCC---e-EEEE-----EeccceeEEccccccCC

```

////////////////////////////////////

Query: alignment of close homologs of gene_86_IAS_virus_KJ003983

HHpred output:

>1LJ9_A transcriptional regulator SlyA; HTH DNA binding protein, structural; 1.6A {Enterococcus faecalis} SCOP: a.4.5.28; Related PDB entries: 1LJ9_B
 Probab=92.87 E-value=0.51 Score=36.89 Aligned_cols=71 Identities=17% Similarity=0.188 Sum_probs=0.0
 Template_Neff=9.600

```

Q ss_pred          ccccCCCHHHHeehhhhhcccch-----HHHHHHHHHhCceeeecCCCC---ceEEEeHHHHHHHHHHh
Q gene_86_IAS_vi  10 XEQHNITL EFLVLYLGAKNADI-----KSISQEVIRKGLATRD LFS DN---RYIVVSNKVKDLIASII 70 (198)
Q Consensus       10 l~~~~is~e~L~l~lv~k~di-----~l~L~K~G~~~~~vt~Kf~DLf~~~~ 70 (198)
      ++.+.|+|.+.+.+.+.+.+.  ...+.|.+++|||.+.+.+.+.  ..+.+.|.+.+.
T Consensus       22 ~~~~~lt~~~~l~i~~~~~la~~~~i~~~~v~~~~l~L~g~li~~~~~d~r~~~~lT~G~----- 96 (144)
T 1LJ9_A          22 FKELSLTRGQYLVLVRCENPGIIQEKIAELIKVDRTTAARA IKRLEEQGF IYRQEDASNKIKRIYATEKGNV----- 96 (144)
T ss_dssp         TGGGTCTTTHHHHHHHHHSTTEENHHHHHHHTCCHHHHHHHHHHHHTTSEEEEECS SCTCEEEEECHHHHHH-----
T ss_pred         hHhCCChHHHHHHHHhCCCcCHHHHHHHhCCCHHHHHHHHHHHCCcEeEecCCCCceeeEEeCHHHHHH-----

```

```

Q ss_pred          cCCCCcccCCcHHHHHHHHHHHHHC
Q gene_86_IAS_vi  71 VNSDKNIVDKDEEYTRLANKLRELY 95 (198)
Q Consensus       71 ~~~~~eL~e~Y 95 (198)

```


Query: gene_89_IAS_virus_KJ003983

CDD output:

SPFH_prohibitin cd03401 Prohibitin family; SPFH (stomatin, prohibitin, flotillin, and HflK/C) superfamily; This model characterizes proteins similar to prohibitin (a lipid raft-associated integral membrane protein). Individual proteins of the SPFH (band 7) domain superfamily may cluster to form membrane microdomains which may in turn recruit multiprotein complexes. These microdomains, in addition to being stable scaffolds, may also be dynamic units with their own regulatory functions. Prohibitin is a mitochondrial inner-membrane protein which may act as a chaperone for the stabilization of mitochondrial proteins. Human prohibitin forms a hetero-oligomeric complex with Bap-37 (prohibitin 2, an SPFH domain carrying homolog). This complex may protect non-assembled membrane proteins against proteolysis by the m-AAA protease. Prohibitin and Bap-37 yeast homologs have been implicated in yeast longevity and in the maintenance of mitochondrial morphology.

```
Pssm-ID: 259799 Cd Length: 195 Bit Score: 120.31 E-value: 9.13e-34
      10      20      30      40      50      60      70      80
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MSSCG 9 VDAGCEGIKVNLYgseKGVDDVSLVTGAVWYNPFTEQVYEYPTYVQTVDYPaFTINAKDGSEFSIDPTISLKIADGKSPQ 88
Cdd:cd03401      4 VDAGEVGVVFRG---KGVKDEVLGEGLHFKIPWIIQVVIIYDVRTQPREIT-LTVLSKDGQTVNIDLSVLYRPDPEKLPE 79
      90     100     110     120     130     140     150     160
.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MSSCG 89 VFKKYRKELADVIngtLFNYVKDAFRIQLNKYTTDEIVSNRDMVEKAIEAHLKALLKENFQLE--QLTSgLKYPQSIVN 166
Cdd:cd03401      80 LYQNLGPDYEERV---LPPIVREVLKAVVAQYTAEELYTKREEVSAEIREALTERLAPFGIIVDdvLITN-IDFPDEYEK 155
      170     180     190     200
.....*.....|.....*.....|.....*.....|.....*.....|
lcl|seqsig_MSSCG 167 AVNAKNAAIQRAQKAQNELAVVKAEAEKVVAAQAEAEAN 206
Cdd:cd03401      156 AIEAKQVAEQEAERAKFELEKAEQEAERKVIEAEGEAEAQ 195
```

Supplementary Note 2

CRISPR spacers matching crAss-like family genomes

spacer host genome (genbank ID)	viral contig with protospacer	CRISPR array position	spacer start	spacer end	e-value	spacer sequence (black font)
						protospacer sequence (red font)
JIAF01000004.1	CDZH01002743	87522	79007	79036	1.28E-04	AAACTAAAGAAGAT-AAGAAAAAGCAATTA
852462_854654_14_spacer_853503_30						AAACTAAAGAAGATGAAGAATAAGCAGTTA
GG775004.1	crAssphage_JQ995537	97065	46335	46364	1.28E-04	TAT-TTTATCCAATACCTTTTTACCATTAT
67033_67769_1_spacer_67080_30						TATCTTTATCAAATACTTTTTACCATTAT
BAKP01000029.1	IAS_virus_KJ003983	99915	98754	98721	2.78E-07	AGCAGACTCATAATCAGAAGGTCATAGGTTCAAG
40807_41491_6_spacer_41202_34						AGCTGACTCATAATCAGAAGACCATAGGTTCAAG
JUJT01000004.1	Chlamydia_CVNZ01000007ext	82794	80442	80410	6.81E-08	GGA-TAGAGCGACAGCCTTCTAAGCTGTAGGTT
169511_170499_3_spacer_169659_32						GGATTAGAGCAACAGCCTTCTAAGCTGTAGGTT
JUJT01000004.1	CENS01015162	79855	11070	11101	6.81E-08	GGATAGAGCGACAGCCTTCTAAGCTGTAGGTT
169511_170499_3_spacer_169659_32						GGATAGTGCAACAGCCTTCTAAGCTGTAGGTT
KL544021.1	crAssphage_JQ995537	97065	84581	84553	1.17E-09	GTTATGAAGATAGAGGTTATCTAATGAA
123220_124110_2_spacer_123344_29						GTTATGAAGATAGAGGTTATCTAATGAA

Bacterial genomes with spacers matching crAss-like phages:

- (sheep rumen) *Prevotella* sp. HUN102 P150DRAFT_scf7180000000012_quiver.4_C, whole genome shotgun sequence
2,328,889 bp linear DNA
JIAF01000004.1 GI:607832995
- (Human oral cavity) *Prevotella scopos* JCM 17725 DNA, contig: JCM17725.contig00029, whole genome shotgun sequence
41,494 bp linear DNA
BAKP01000029.1 GI:602603948
- Parabacteroides* sp. 20_3 genomic scaffold supercont1.36, whole genome shotgun sequence
343,425 bp linear DNA
GG775004.1 GI:300829907
- (isolation_source="infected leaves") *Pectobacterium carotovorum* subsp. *carotovorum* strain BC D6 B6.scaffold4, whole genome shotgun sequence
486,112 bp linear DNA
JUJT01000004.1 GI:741149445
- (This is a reference genome for the Human Microbiome Project) *Porphyromonas* sp. 31_2 genomic scaffold acTiZ-supercont2.1, whole genome shotgun sequence
1,861,588 bp linear DNA
KL544021.1 GI:659424147