

Supplementary Material

Prediction of several novel families of enzymes involved in oxidative and other complex modifications of bases in nucleic acids

Lakshminarayan M Iyer¹, Mamta Tahiliani², Anjana Rao² and [L. Aravind^{1*}](mailto:aravind@mail.nih.gov)

* Address for correspondence: L. Aravind (aravind@mail.nih.gov)

¹National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA
²Department of Pathology, Harvard Medical School and Immune Disease Institute, 200 Longwood Avenue, Boston, Massachusetts 02115, USA.

Abstract

Modified bases in nucleic acids present a layer of information that directs biological function over and beyond the coding capacity of the conventional bases. While a large number of modified bases have been identified, many of the enzymes generating them still remain to be discovered. Recently, members of the 2-oxoglutarate- and iron(II)-dependent dioxygenase superfamily, which modify diverse substrates from small molecules to biopolymers, were predicted and subsequently confirmed to catalyze oxidative modification of bases in nucleic acids. Of these, two distinct families, namely the AlkB and the kinetoplastid base J binding proteins (JBP) catalyze in situ hydroxylation of bases in nucleic acids. Using sensitive computational analysis of sequences, structures and contextual information from genomic structure and protein domain architectures, we report five distinct families of 2-oxoglutarate- and iron(II)-dependent dioxygenase that we predict to be involved in nucleic acid modifications. Among the DNA-modifying families, we show that the dioxygenase domains of the kinetoplastid base J-binding proteins belong to a larger family that includes the Tet proteins, prototyped by the human oncogene Tet1, and proteins from basidiomycete fungi, chlorophyte algae, heterolobosean amoebae and bacteriophages. We present evidence that some of these proteins are likely to be involved in oxidative modification of the 5-methyl group of cytosine leading to the formation of 5-hydroxymethylcytosine. The Tet/JBP homologs from basidiomycete fungi such as *Laccaria* and *Coprinopsis* show large lineage-specific expansions and a tight linkage with genes encoding a novel and distinct family of predicted transposases, and a member of the Maelstrom-like HMG family. We propose that these fungal members are part of a mobile transposon. To the best of our knowledge, this is the first report of a eukaryotic transposable element that encodes its own DNA-modification enzyme with a potential regulatory role. Through a wider analysis of other poorly characterized DNA-modifying enzymes we also show that the phage Mu Mom-like proteins, which catalyze the N6-carbamoylmethylation of adenines, are also linked to diverse families of bacterial transposases, suggesting that DNA modification by transposable elements might have a more general presence than previously appreciated. Among the other families of 2-oxoglutarate- and iron(II)-dependent dioxygenases identified in this study, one which is found in algae, is predicted to mainly comprise of RNA-modifying enzymes and shows a striking diversity in protein domain architectures suggesting the presence of RNA modifications with possibly unique adaptive roles. The results presented here are likely to provide the means for future investigation of unexpected epigenetic modifications, such as hydroxymethyl cytosine, that could profoundly impact our understanding of gene regulation and processes such as DNA demethylation.

Contents

1. Multiple sequence alignment of novel members of the 2OGFeDO superfamily.
2. Phyletic distribution of the TET1/JBP family of proteins
- 2b. Superalignment of metazoan TET domains with intron-exon boundaries
- 2c. Multiple sequence alignment of the TET1, TET2 and TET3 proteins
- 2d. Multiple sequence alignment of the JBP1-C terminal domain
- 2e. Multiple sequence alignment of novel transposase and identical hits
- 2f. Multiple sequence alignment of the HMG domain found in the vicinity of the 2OGFeDO containing transposon
- 2g. Multiple sequence alignment of small alpha helical domain found either in the neighborhood or fused to JBP and/or the transposase
- 2h. Multiple sequence alignment of Cys cluster that is often fused to the hydroxylase domain
- 2i. Gene neighborhoods of the predicted transposase gene associated with JBP
3. Phyletic distribution, domain architectures and alignment of the algal RNA-modification associated family
4. Phyletic distribution, domain architecture and multiple sequence alignment of the fungal subfamily of AlkB proteins that are fused to SAD and R3H domains
- 5a. Phyletic distribution and multiple sequence alignment of the R3H domain-associated family of 2OGFeDO proteins
- 5b. Multiple sequence alignment of domain C-terminal to the 2OGFeDO domain, called X in Figure 2
- 5c. Multiple sequence alignment of the R3H domain fused to the 2OGFeDO domain
- 6a. Phyletic distribution and multiple sequence alignment of the DNA glycosylase associated family of 2OGFeDO domains
- 6b. Gene neighborhoods of the bacterial versions of the DNA glycosylase associated family
7. Phyletic distribution, domain architectures and multiple sequence alignment of the MOM-family of acetyltransferases
8. Table of distribution of Methylases and predicted DNA-modifying hydroxylases

1a. Multiple sequence alignment of novel members of the 2OGFeDO superfamily.

```

Secondary Structure      -EEEE-----saprfrfridpsplhekn-----hhhhhhh--hhhhhhh--HHHHHHHHHHHH-----EEEEEE-----
TET1_Human              SWSMYFNgcKkfr-----saprfrfridpsplhekn-----LEDNQSLATPLAPYKQYAPYAYQNOVEYENVArecl-----gs-KEGRPFSGVATC-----
TET1_Mouse              SWSMYFNgcKkfr-----semprkfrl apyvlhekg-----LEKKGELATVLAFLFKQAFYAYQNOVEYEVAgdcrl-----eg-EGRPFSGVTC-----
TET3_Human              SWSMYFNgcKkyar-----sktrpkfrlsgdnkpeeev-----LRKSFODLATEVAPLYKRLAPQAYQNOVNERIAIdcrl-----gl-KEGRFPAGVATC-----
TET3_Mouse              SWSMYFNgcKkyar-----sktrpkfrlsgdnkpeeev-----LRNSFODLATEVAPLYKRLAPQAYQNOVNEVDAIdcrl-----gl-KEGRPFSGVATC-----
TET2_Human              SWSMYFNgcKkfar-----sklprkfrl lgdgpkceek-----LESHLQNLSTLMAPTYKRLAPADAYNNQIYEHRRApecl-----gd-KEGRPFSGVATC-----
TET2_Mouse              SWSMYFNgcKkfar-----skkprkfrl hgaepkeeev-----LGSHLQNLATVIAPYKRLAPADAYNNQVFEHQApdccl-----gl-KEGRPFSGVATC-----
LOC580376_Sea urchin    SWSMYFNgcKkfar-----sktrpkfrl lsgnqgedv-----LSDRFQMATDLGPLYKRLAPEFNPNVFEESKcecl-----gk-ETGRFAQVATC-----
CG2083_Drosophila       SWSMYFNgcKkyar-----sktrpkfrlsvk-----IEDHNNIATLLAFVQVCPRSYDNQTYEHEAAdcrl-----gd-EPKPFSGVATC-----
v1g22996_Nematostella  SWSMYFNgcKkfar-----skaprkyl l-dsakeet-----LERILEGIATPIAPVYSKAAPAFANOTREERNGhecrl-----ghsAVGRPFSGVTC-----
gp2_BPCooper_109392355  RNFPGYAPRprlsv-----reacstlgrdydpq-----IYOVLESYADQFAAGLAI DPELVKRGSDQ-ASvlhdw-----rl-1-GEAKLWTSGVNV-----
FRAL2749_FaIn_111222169 RNFPGYAPRprpvc-----zegcqlglaesepg-----EHRVLERWALLERTLDIDPSIVARDAQYMTYV-dsaw-----rl-OGGRLWTSGVNI-----
GOS_4878734_Mmet_139542046 SIFQSLPrntmr-----dferfsahkksik-----YNNIFSPMNDL INIKYKLPQYERDRIVKESVwedv-----al-NKKSPLFCMNI-----
GOS_153646_Mmet_134535573 CIIGSYDrntmr-----imhnsrsvhrskaagt-----PIKAMVLAGRQSLSVIKRLPELPEYTHRESVLDrvpeqg-----rf-RCDLFTSTVNI-----
GOS_9579659_Mmet_135108850 -----vlatraqpe-----VFAGLSKVKGLMVGQVCPPEVAANQKQFVGGIhdh-----wk-KTGFPTFTVNV-----
GOS_3272744_Mmet_139110457 -----mlmccleSENLIKYMPEQYASOKKLIETtlp-----mlmccleSENLIKYMPEQYASOKKLIETtlp-----ky-RFQNLFTSSISN-----
GOS_4428895_Mmet_140212139 -----KRAMLMSCLESEKIRIKYQMPQYASOKKLIETtlp-----KRAMLMSCLESEKIRIKYQMPQYASOKKLIETtlp-----ey-RFQNLFTSSISN-----
JBP1_Tbru_6018043       GIAGYDrqtpv-----elkorktsfctyhtk-----KPAVPLPVYVSEIYVPEVPHRAADSAPDI-----LEETtlp-----ky-RFQNLFTSSISN-----
JBP2_Tbru_72391588      GIAGYDrqtpv-----elkorktsfctyhtk-----KIPORTKTKFKEVH-----KWDCIPFIEIDKQFSIHIPDRHKVOLERASLT-kd-----fq-1KNTAFTITIN-----
Tc00_1047053506605_229_Tcru_71662347 GILGYDylmnp-----krkermtfetrnwgk-----IIGCGELLQLLDQLYKENADPHYLQRRVLPPE-----ym-LFMTVFTVTVN-----
Tc00_1047053519357_10_Tcru_71421637  GIVGYDylmnp-----qkrceretftrknwss-----VDSVDSFVLMVKLIDSLFKCLIPDAYKROLNRANLR-dk-----fk-1PMTSFTVTVN-----
GOS_4675571_Mmet_139186735 GIAGYDrqtpv-----elkorktsfctyhtk-----SWNFPMHIDVSAIYKAVFPEQAAQDAVPDI-----vT-1HGSFSTLTVN-----
GOS_4428895_Mmet_140212139 GIAGYDrqtpv-----elkorktsfctyhtk-----SWNFPMHIDVSAIYKAVFPEQAAQDAVPDI-----vT-1HGSFSTLTVN-----
GOS_9294781_Mmet_135380621 GIIYFDrydrnlqngktp-----kiportktkfkevh-----KIIPORTKTKFKEVH-----KWDCIPFIEIDKQFSIHIPDRHKVOLERASLT-kd-----fq-1KNTAFTITIN-----
GOS_3208587_Mmet_139987906 SIIGYDrprip-----yrcrtafctkhd-----MYSQIPIYQSIKLFEEFLPERWQNKQNEWKTsed-----fk-1HGTFTVTVN-----
GOS_7711251_Mmet_136831790 NIPGFYEssnfs-----klperlthfrtnfd-----KYNVGLPFFQKIDSLFKCLIPDAYKROLNRANLR-dk-----fk-1PMTSFTVTVN-----
GOS_9234699_Mmet_135432669 NIPGFYEssnfs-----klperlthfrtnfd-----DYNKGLFQIQIDSLFKCLIPDAYKROLNRANLR-ph-----lk-1PMTSFTVTVN-----
GOS_464085_Mmet_144014002 NIMGYDrwdsirasKragmkp-ptrcrltsfctsrpe-----KWENVPLIQIDIAQYKRLVPKAYANORQAADS-vk-----fk-1PMTSFTVTVN-----
GOS_8006756_Mmet_136547457 NIIGYMDwtlghkymfsgvmkeikpavrsyftqnyd-----NWTPKSLVKHIDIAQYKRLAPVYKQKQAKADET-y-----fk-1KGTATTLTIN-----
GOS_8124242_Mmet_136439712 NIIGYDkrdnl-----gan-pporttaftsgvq-----KWNVVPLKIMIDIQKRLIPSNHRIQYDRANKT-d-----fv-1NGTAFSTVTVN-----
GOS_398524_Mmet_144068378 GVALRQgkalelcnlpgdgrpa-tpcifigdsamyik-----FLEWQPLFAVLGGIATVTHWPLDAGLSILSSLE-hgkldisgpenll-----dvlb-FWSNPFSTKFLI-----
GOS_137212_Mmet_134552279 GVVGFMDKsami-----yrcrtafctkhd-----NYOQGLPFFVQVDEYKLCPEYINRQKNIAGETnqg-----fv-1DPTSTFTVTVN-----
CCIG_03999_Ccin_169848807 GVACFAPawavagheksllpg-----pstplkksyik-----FVDEMEESLALVGLVSHVPELEI GCVHNMNrvadesqvkypqevh-----rvlk-TWASFTLGLSLI-----
LACBTRDAPF_136849_Lbic_170117859 GIVNLSFawfsgqghepedl-----zasidqepant-----YLQNLVLANAVGLIYAVIQDLPESGLATEKLANha-----iall-VWATFTLGLSLI-----
CCIG_09735_Ccin_169868221 GVATFSSCyyqghdgdpsps-----psaslkpppqa-----FRDMTEESAVLGGVIAITQWPLFNAGLEVLVLESLh-rtfipvdnaetld-----alln-FWSNPFSTKFLI-----
CCIG_09735_Ccin_169868221 GVALRQgkalelcnlpgdgrpa-tpcifigdsamyik-----FLEWQPLFAVLGGIATVTHWPLDAGLSILSSLE-hgkldisgpenll-----dvlb-FWSNPFSTKFLI-----
CHLREDRAPF_179132_Crei_159485382 QOAVRQppwqlkprtr-----yvnglylceeahypfg-----IEALEEMCTAVVAEAVOIAFYLYTRDFELSAWgq-----am-YGTATNLVSMI-----
CHLREDRAPF_193135_Crei_159480130 GNALIQNdqkalsk-----yavaharvrtsvnyqka-----NVDMEPEESLALVGLVSHVPELEI RFLMAYRDLVITQgk-----fm-1GATATNLVSLP-----
CHLREDRAPF_180688_Crei_159491750 ARPRVSCkydqdk-----yaaataedaate-----PERLITMTCRVTWGAVEVFFLEECRDVFLSPVgk-----vc-YDGTMTMOSVTV-----

```

Ngru1000008499 Ngru Ngru1000008499
 Ngru1000004275 Ngru Ngru1000004275
 Ngru100014031 Ngru Ngru100014031
 Ngru1000012731 Ngru Ngru1000012731
 Ngru100009725 Ngru Ngru100009725
 Ngru100002001 Ngru Ngru100002001
 Ngru1000012827 Ngru Ngru1000012827
 Ngru100013068 Ngru Ngru100013068
 Aano1000002837 Aano Aano1000002837
 Aano1000006794 Aano Aano1000006794
 Aano1000008654 Aano Aano1000008654
 Aano100001260 Aano Aano100001260
 Aano1000010511 Aano Aano1000010511
 Aano1000003820 Aano Aano1000003820
 Aano1000005385 Aano Aano1000005385
 Mpus1000003041 Mpus Mpus1000003041
 Mpus1000008227 Mpus Mpus1000008227
 Aano1000005614 Aano Aano1000005614
 PHATRDRAFT_50200 Ptri_219129852
 PHATRDRAFT_42626 Ptri_217411154
 PHATR 44207 Ptri_219113839
 Oe10g03170 Otau_116002016
 OSTLU_93380 Olic_145351887
 THAPSDRAFT_21769 Tpsa_220975450
 Aano1000005432 Aano Aano1000005432
 ALKB Ec_113638
 ALKBH3 Hsap_21040275
 Aano1000003820 Aano Aano1000003820
 CIMG_06114 Cimm_119182137
 ATEG_08110 Ater_115433188
 An09g03880 Anig_145242116
 AN3951_2 Anig_167625987
 BC1G_08942 Bfuc_154305295
 CHGG_06386 Cglo_116193337
 LACB1DRRAFT_131729 Lbic_170093061
 CNBL2390 Cneo_134117982
 SMOG_03244 Fmo_169600101
 Franean1_7129 Fsp_158318848
 GobsU_010100005758 Gobs_168699004
 BlinB01002514 Blin_62423906
 AvindRAFT_5847 Avin_67153576
 Acry_3179 Acry_14824370010095614
 An18g01550 Anig_145254402
 CCG_09674 Ccin_169866364
 CCG_09656 Ccin_169866328
 CCG_09671 Ccin_116499174
 CHL1DRRAFT_153459 Crei_159485224
 Dpui1000012176 Dpui_1000012176
 Dpui1000004224 Dpui_1000004224
 FG09449.1 Gzea_42553402
 LACB1DRRAFT_32592 Lbic_164645299
 LACB1DRRAFT_328485 Lbic_170102781
 LACB1DRRAFT_315778 Lbic_164646975
 NCU09412 Ncra_85091137
 Ngru1000002312 Ngru Ngru1000002312
 Ngru1000008142 Ngru Ngru1000008142
 Ngru1000013028 Ngru Ngru1000013028
 Pc13g0210 Pchr_211583340
 PHYPADRAFT_169264 Ppat_162672410
 Pp1a1000002318 Pp1a_Pp1a1000002318
 Ppam1000011751 Ppam_Ppam1000011751
 Psoj1000008776 Psoj_Psoj1000008776
 Smeo1000018380 Smeo_Smeo1000018380
 Smeo1000018266 Smeo_Smeo1000018266
 TSTA_016850 Tsta_218711789
 P4H Crei_159794863
 NONBRDRAFT_31538 Mbre_167519270
 OSTLU_32491 Olic_145348837
 OSTLU_17228 Olic_145351789
 Oe10g09020 Otau_116058658
 OSTLU_28402 Olic_145346166
 OSTLU_28816 Olic_145341417
 Oe10g02080 Otau_116000439
 PHATRDRAFT_47703 Ptri_219123239
 PHATRDRAFT_43074 Ptri_219110259
 PHATRDRAFT_37198 Ptri_219120466
 PHATRDRAFT_42585 Ptri_219109862
 PHATRDRAFT_49155 Ptri_219127159
 Aano1000005937 Aano Aano1000005937
 STIAU_4281 Saur_115375244
 MXAN_2402 Mxan_10109216
 Syncc9902_1289 Syn_78184862
 BL107_13860 Syn_116070736
 Syncc9605_1181 Syn_78212712
 consensus/80%
 consensus/70%
 consensus/85%
 consensus/75%

Species abbreviations: Aano : Aureococcus anophagefferens; Acry : Acidiphilium cryptum; Anid : Aspergillus nidulans; Anig : Aspergillus niger; Ater : Aspergillus terreus; Avin : Az

[Back to Contents](#)

2a. Phyletic distribution of the TET/JBP family of proteins

#: Tet/JBP family; bacteriophage gp2 subfamily							
gi	gene_name	length	Species	class	GenBank_annotation		
111222169	FRAL2749	277	Frankia alni ACN14a	actinobacteria	hypothetical protein FRAL2749 [Frankia alni		
194100516	Nigel2	322	Mycobacterium phage Nigel	dsDNA viruses, no RNA stage/caudovirales	gp2 [Mycobacterium phage Nigel].		
19392355	Cooperp2	322	Mycobacterium phage Cooper	dsDNA viruses, no RNA stage/caudovirales	gp2 [Mycobacterium phage Cooper].		
Sequences from marine metagenomes that are closer to the gp2 subfamily							
139110457	GOS_3272744	133	marine metagenome		hypothetical protein GOS_3272744 [marine meta		
139186735	GOS_4675711	137	marine metagenome		hypothetical protein GOS_4675711 [marine meta		
139542046	GOS_4878734	201	marine metagenome		hypothetical protein GOS_4878734 [marine meta		
140212139	GOS_4428895	107	marine metagenome		hypothetical protein GOS_4428895 [marine meta		
142006841	GOS_3046279	235	marine metagenome		hypothetical protein GOS_3046279 [marine meta		
134535573	GOS_153646	182	marine metagenome		hypothetical protein GOS_153646 [marine metag		
135108850	GOS_9579659	181	marine metagenome		hypothetical protein GOS_9579659 [marine meta		
#:Tet/JBP family; JBP subfamily							
gi	gene_name	length	Species	class	GenBank_annotation		
6018043	JBP1	811	Crichidia fasciculata	euglenozoa>kinetoplastida	J-binding protein [Crichidia fasciculata].		
6018045	JBP1	827	Leishmania tarantolae	euglenozoa>kinetoplastida	J-binding protein [Leishmania tarantolae].		
146078722	Lin309.1540	814	Leishmania infantum JPCM5	euglenozoa>kinetoplastida	DNA J-binding protein, putative [Leishmania i		
6018041	JBP1	839	Trypanosoma brucei	euglenozoa>kinetoplastida	J-binding protein [Trypanosoma brucei].		
157865276	Linj309.1480	814	Leishmania major strain Friedlin	euglenozoa>kinetoplastida	DNA J-binding protein, putative [Leishmania m		
140720701	Tb927.7.4650	1077	Trypanosoma brucei TREU927	euglenozoa>kinetoplastida	SNF2 DNA repair protein, putative [Trypanosom		
12391588	Tb927.7.4650	1077	Trypanosoma brucei TREU927	euglenozoa>kinetoplastida	SNF2 DNA repair protein, putative [Trypanosom		
154334093	LbrM14_V2.0040	1098	Leishmania braziliensis MIMOM/BR/75/M2904	euglenozoa>kinetoplastida	J-binding protein, putative [Leishmania brazi		
146081173	Lin314.0040	1022	Leishmania infantum JPCM5	euglenozoa>kinetoplastida	J-binding protein, putative [Leishmania infan		
157866423	Lin314.0040	1022	Leishmania major strain Friedlin	euglenozoa>kinetoplastida	J-binding protein, putative [Leishmania major		
239107405	Tc00.1047053510357.10	832	Trypanosoma cruzi strain CL Brener	euglenozoa>kinetoplastida	DnaJ chaperone protein, putative [Trypanosoma		
71421637	Tc00.1047053510357.10	832	Trypanosoma cruzi strain CL Brener	euglenozoa>kinetoplastida	DnaJ chaperone protein, putative [Trypanosoma		
71653481	Tc00.1047053506753.120	831	Trypanosoma cruzi strain CL Brener	euglenozoa>kinetoplastida	DnaJ chaperone protein, putative [Trypanosoma		
14272266	Tc00.1047053508859.74	1086	Trypanosoma cruzi strain CL Brener	euglenozoa>kinetoplastida	SNF2 DNA repair protein, putative [Trypanosom		
71662347	Tc00.1047053506605.229	986	Trypanosoma cruzi strain CL Brener	euglenozoa>kinetoplastida	SNF2 DNA repair protein, putative [Trypanosom		
154332081	Tb927.7.4650	1077	Leishmania braziliensis MIMOM/BR/75/M2904	euglenozoa>kinetoplastida	DNA J-binding protein, putative [Leishmania b		
Sequences from marine metagenomes that are closer to the JBP subfamily							
139897906	GOS_3208587	293	marine metagenome		hypothetical protein GOS_3208587 [marine meta		
136831790	GOS_7711251	346	marine metagenome		hypothetical protein GOS_7711251 [marine meta		
135432669	GOS_9234698	420	marine metagenome		hypothetical protein GOS_9234698 [marine meta		
144104092	GOS_464085	112	marine metagenome		hypothetical protein GOS_464085 [marine metag		
136547457	GOS_8006756	320	marine metagenome		hypothetical protein GOS_8006756 [marine meta		
136439712	GOS_8124242	286	marine metagenome		hypothetical protein GOS_8124242 [marine meta		
144068378	GOS_398524	329	marine metagenome		hypothetical protein GOS_398524 [marine metag		
134552279	GOS_137212	298	marine metagenome		hypothetical protein GOS_137212 [marine meta		

Table with columns: accession, gene name, length, species, class, GenBank annotation. Lists various sequences and their corresponding annotations.

Table with columns: gi, gene name, length, species, class, GenBank annotation. Lists sequences and their annotations, including entries for Caps100002965 and others.

Table with columns: gi, domain architecture, gene name, length, species, class, GenBank annotation. Lists sequences and their domain architectures.

Fasta sequences of predicted JBP/TET family proteins in eukaryotes with incomplete genomes and not in GenBank or reconstructed fasta sequences. Includes a long block of FASTA format text.

Back to Contents

2b. Superalignment of metazoan TET domains with intron-exon boundaries.

Table showing superalignment of metazoan TET domains with intron-exon boundaries. Columns include domain names (TET1 exon bounds, etc.), coordinates, and sequence alignments.

Multiple sequence alignment of the TET1, TET2 and TET3 proteins. The alignment shows conserved regions across various species including Tetrahymena, Gallus, and Mus. Consensus sequences are provided at the bottom.

Back to Contents

2c. Multiple sequence alignment of the TET1, TET2 and TET3 proteins

TET-1 alignment. Intron-exon boundaries are indicated. The alignment shows conserved regions across various species including Tetrahymena, Gallus, and Mus. Consensus sequences are provided at the bottom.

TET-2 alignment. Intron-exon boundaries are indicated. The alignment shows conserved regions across various species including Tetrahymena, Gallus, and Mus. Consensus sequences are provided at the bottom.

TET-3 alignment. Intron-exon boundaries are indicated. The alignment shows conserved regions across various species including Tetrahymena, Gallus, and Mus. Consensus sequences are provided at the bottom.

Back to Contents

2d. Multiple sequence alignment of the JBP1-C terminal domain

Predicted secondary structure and multiple sequence alignment of the JBP1-C terminal domain across various species. Consensus sequences are provided at the bottom.

Back to Contents

2e. Multiple sequence alignment of novel transposase and identical hits

Multiple sequence alignment of novel transposase and identical hits across various species. The alignment shows conserved regions and conserved motifs across different transposase families.

LOC116278 Ggal_118097475
MGC115244 Xlae_148236655
hbox_X161_1876129
LOC568530 Drecr_189524866
SMF_Hsap_160112910
LOC509419 Btau_119895718
chr1:7159298_Drecr_189534053
LOC100149527 Drecr_189535301
MRE1DRAPF_178766_Crei_159484947
CHLEDRAPF_194780_Crei_159488509
CHLEDRAPF_174339_Crei_159473346
LACB1DRAPF_307726_Lbic_170110059
LACB1DRAPF_316727_Lbic_170110668
CCin_02459_Ccin_169851012
CNA07690 Cneo_58259409
BRALDRAPF_131287_Bf1o_219494211
BRALDRAPF_95999_Bf1o_219471002
BRALDRAPF_91711_Bf1o_219462506
consensus/85%

Species abbreviations: Bf1o : Branchiostoma floridae; Btau : Bos taurus; Ccin : Coprinopsis cinerea; Clup : Canis lupus; Cneo : Cryptococcus neoformans; Crei : Chlamydomonas reinha

Identical hits in Laccaria
170103707 170103709
170094354 170120007
170114173 170114177

Back to Contents

2f. Multiple sequence alignment of the HMG domain found in the vicinity of the 20GFeD0 containing transposon

Domain boundaries
Secondary structure prediction
CCin_02072_Ccin_169851139
CCin_02773_Ccin_169842922
CCin_03995_Ccin_169848799
CCin_05495_Ccin_169864670
CCin_05591_Ccin_169862973
CCin_07037_Ccin_169850583
CCin_07254_Ccin_169868984
CCin_08112_Ccin_169860182
CCin_09453_Ccin_169867988
CCin_09747_Ccin_169868245
CCin_10129_Ccin_169848160
CCin_10218_Ccin_169857280
CCin_10454_Ccin_169868188
CCin_10664_Ccin_169852260
CCin_12596_Ccin_169850600
CCin_12707_Ccin_169864942
CCin_12859_Ccin_169869222
CCin_12891_Ccin_169869841
CCin_12944_Ccin_169866812
CCin_12950_Ccin_169862398
CCin_13090_Ccin_169862006
CCin_13128_Ccin_169862411
CCin_13284_Ccin_169869977
LACB1DRAPF_297419_Lbic_170118976
LACB1DRAPF_315026_Lbic_170108467
LACB1DRAPF_335151_Lbic_170117741
LACB1DRAPF_316530_Lbic_170117428
LACB1DRAPF_331947_Lbic_170110305
LACB1DRAPF_332251_Lbic_170111125
LACB1DRAPF_332028_Lbic_170103711
LACB1DRAPF_327601_Lbic_170100513
LACB1DRAPF_330019_Lbic_170106331
LACB1DRAPF_330403_Lbic_170107125
LACB1DRAPF_330704_Lbic_170111125
LACB1DRAPF_330207_Lbic_170117693
LACB1DRAPF_330006_Lbic_170105886
LACB1DRAPF_329656_Lbic_170105236
LACB1DRAPF_331533_Lbic_170117719
LACB1DRAPF_307132_Lbic_170120119
LACB1DRAPF_335987_Lbic_170117861
LACB1DRAPF_335987_Lbic_170117861
LACB1DRAPF_335999_Lbic_170116822
LACB1DRAPF_336286_Lbic_170120560
2c6a Hmel2128 Hmc domains
K99 Hsap_17942547
1111A 7 PDB
21e1a 49 database
1gt0D 50
1j4a 54
1aab 33
1e7jA 37
1j5nA 38
Species Abbreviations: Ccin : Coprinopsis cinerea; Lbic : Laccaria bicolor

Back to Contents

2g. Multiple sequence alignment of small alpha helical domain found either in the neighborhood or fused to JBP and/or the transposase

Predicted secondary structure
CCin_09736_Ccin_116493919
CCin_07252_Ccin_169867890
CCin_10219_Ccin_116503602
CCin_11682_Ccin_169864014
CCin_12550_Ccin_169869581
CCin_12678_Ccin_169869760
CCin_13131_Ccin_169862417
CCin_13153_Ccin_169863832
CCin_12708_Ccin_169869474
CCin_12861_Ccin_169866296
CCin_12949_Ccin_169862396
CCin_13116_Ccin_169869499
CCin_05496_Ccin_169864672
CCin_10425_Ccin_169869193
CCin_09451_Ccin_169867984
CCin_11469_Ccin_169859597
CCin_11998_Ccin_169861169
CCin_12695_Ccin_116506927
CCin_13266_Ccin_169870064
LACB1DRAPF_316487_Lbic_170117410
LACB1DRAPF_333145_Lbic_170113173
LACB1DRAPF_330096_Lbic_170117861
LACB1DRAPF_327556_Lbic_170100793
LACB1DRAPF_330405_Lbic_170106947
LACB1DRAPF_306696_Lbic_170108481
LACB1DRAPF_307133_Lbic_170107823
LACB1DRAPF_335170_Lbic_170117835
LACB1DRAPF_316718_Lbic_170117691
consensus/100%
consensus/95%
consensus/90%
consensus/85%
Species Abbreviations: Ccin : Coprinopsis cinerea; Lbic : Laccaria bicolor

Back to Contents

2h. Multiple sequence alignment of Cys cluster that is often fused to the hydroxylase domain

```

Predicted sec str      --HHHHHHHHHHHHHHH-----HHHHHHHHHH-----HHHHHHH--
CCIG_10424_Ccin_169866911  IQIEARMDFW-----ATWV--GNIEEFKREKVNCGSPLCSV-K--CA-PLETHMCSACFYDPSRYPSNCLRLDYLRLDRLRS---QAGLSPGQVERFMPAY
CCIG_12620_Ccin_169865996  LQIEVRMDFW-----SSWL--EQNADFRGRHVNCVSGPHCTQ-G--CT-PLDDFMACEACFSNGQQYPTNCLRLTQFLRDLTLRS---RGLTSRQITDFLPVY
CCIG_13092_Ccin_169862010  FOLEADYSRF-----DRWA--LQFDKGHFTPPCGNPQCSN-G--CSGPVDEGMWCLPCWLDHGAGFERCKLVDVYLKQTLCE---BGYLTPETAPAFLEBY
CCIG_10224_Ccin_169857292  FOLQSRMEFW-----NWA--SRVCMNDPFLNCGEACHA-G--CD-ALLDHPSCVBCFQTFRGIINWCMQELIAHFLET---DCETPVEEQPFLQY
CCIG_10671_Ccin_169852274  FOLEGFELN-----EWAQALRDEDDNQKEGCCNPDGCD-TPPCD--PLDDMACGLCALQGEAYDCGFFKQYMAHLOH---DAGMPITSIPGFLEOF
CCIG_04135_Ccin_169860374  FETEADLGLN-----REWA--SLDRTGSIIRDCCNCEEGP-R--GN-PLHDYPACNSCF--NSDFGMDFCFSSHSHYHMLMS---QFYTPVHEIPLVLAEY
CCIG_13157_Ccin_169863840  FOLQRTDLN-----KAYT---AQLDHOQTCQDKPEQA-AGACE-IMEDYISCRACF---NDNEVCTLEYDYIHYLNR---DVGIPWDDLVPFLTRY
CCIG_12621_Ccin_169865998  FSVISDNGLW-----CSWT--EALDQGAIVRCSNTECQL-EG--CT-LLGDVLSCLKCFENDAIIGRCDLPLLLRRPMEF---KDLPEEKVRPFTVWF
CCIG_10672_Ccin_169852276  LEYEGSGLN-----HWT--ASRVLLEEQPTCDDPECVN-GGCD--PLRDMSCIRCFDNDKHKDTCSRRQVLSVILAIN---NPHLSEDLKCFIAHY
CCIG_12914_Ccin_169866419  MEAEOLDAN-----RRLY--POLEESHPVGVCKSDYCHDNHHECI--PSPHTMACOPCI---STGHOCSPREDFIARHVAE---AMPVTAQVAEFPVDEY
CCIG_07247_Ccin_169867880  CHDRDRFEDN--FF-----DNIG-----GPCNHEECIKFPGGCM-LTAEGLGCAPCI---KGRPCPWRKDAYIITVILG---ELEVPGKAFQFPFTKL
CCIG_09678_Ccin_169866372  YMSPLRHLF-LLLLEK--KNWG--EFLLAHRPLSRCDALACKQHGSDCS-VGDRHLGCRTCI---RNNVLCSHVTDCLLSTLL---QDAIPVAQANELISTF
CCIG_09735_Ccin_169868221  CHDRDRFEGW--FF-----DNIG-----GPCNQSQCIRFPGGCM-LTKBGLGCAPCI---KGSFPCWRKDAYI IAVILG---ELDVPYKAFEFPTKL
CCIG_11350_Ccin_169866093  VEALQDQRW-----SWT--SOWYTSQNSDESCT--CNI-PDDCV--PIEDAPGACN---VMYEGGMEELKQVREDKRTGSRCCVGVKINWDL
CCIG_04088_Ccin_169860280  FTQNRLTAA-----KAVA-DSTA--VSRSSYSANKVCDNPSCKAKPDSCH--VAPRGRGCLPCN---LALQPCSFVTDIVQRLGLE---CGLQPDDAIAFLE--
CCIG_01240_Ccin_169852992  SCVRLTDFL-----LNP--KHVK--DIVAANPLGPCSAPOCASKEVLCL--PSDVLGCTNCV---RLGLTCSHVDTCLIVFISA---KANIDTALARKFLDAL
CCIG_09972_Ccin_169869428  VMIPRIVRFY---LDP--TGWK--AFVASNPLYGCTALQCSSSDKVKCV--LSDFGCTTCV---RLGLVCSHLETCLTLALAQ---KMNIDAERARQFLRIV
CCIG_10222_Ccin_169857288  LATIQKLVDR--NL-----KEY-----CTQKCIIESFGACH--LTEBFGGPCI---ARGATCPWRKDAYI IAVIMG---ELDLSYKAFEFPTKL
CCIG_11349_Ccin_169866091  VRGLFNKEW-----NDYA---RDNPNPDRCD--PDCQ-GL-CH-PEVDFGIGCA---TRQVCSFLRDFIERRLPSL--HALPMAVNDQFYQQQ
consensus/I00%          .....C.....p.....s.....p..C..C.....C.....hh.....h...
consensus/95%          .....C.....p.....s.....p..C..C.....C.....hh.....h...
consensus/90%          h.....h.....p.h.....Cs..C.....C.....c..uC..C.....C.....hl..h.....s..h..hh...
Species Abbreviations: Ccin : Coprinopsis cinerea

```

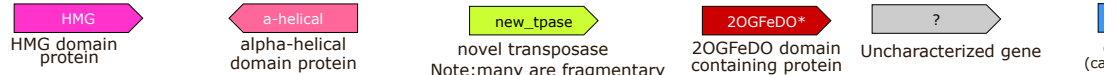
[Back to Contents](#)**2i. Gene neighborhoods of the predicted transposase gene associated with JBP**

Note the neighborhoods are rendered in the svg format and your browser should support svg to view this. Most current versions of web browsers support svg.

However, if that is not so and you see a white screen below without the genes, click [here](#) to see the neighborhood

OR if that too doesn't work, click [here](#) to access a pdf file

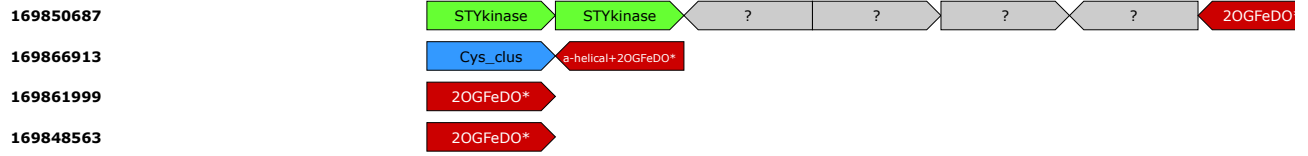
Legend for key genes:



Note: numbers on the left represent Gi numbers of the predicted transposase gene

Coprinopsis gene neighborhoods (click on the boxes to access protein sequences)

Solo 2OGFeDO like genes



2OGFeDO like genes in the vicinity of the predicted transposase gene



Back to Contents

3. Phyletic distribution, domain architecture and alignment of the algal RNA-modification associated family

Table with columns: GI, Domain architecture, gene name, Len, Organism, class, GenBank Annotat. Includes protein sequences and key features for the 20GFeDo family.

Fasta sequences of proteins not in GenBank and assigned Gags. Includes sequence headers and FASTA format content.

Back to Contents

4. Phyletic distribution, domain architecture and multiple sequence alignment of the fungal subfamily of AlkB proteins that are fused to SAD and R3H domains

Table with columns: Phyetic distribution, gene name, length, Species, class, genbank_annotat. Includes protein sequences and common domain architecture for the AlkB subfamily.

Back to Contents

5a. Phyletic distribution and multiple sequence alignment of the R3H domain-associated family of 20GFeDo proteins

Table with columns: Phyetic distribution, GI, Domain architecture, gene name, Len, Organism, class, GenBank Annotation. Includes protein sequences for the R3H domain-associated family.

Table with 4 columns: Accession ID, Gene/Protein Name, Species, and Description. Lists various sequences and their associated biological information.

Proteins that lack the 20GPeDo domain

Table listing proteins that lack the 20GPeDo domain, including accession IDs, gene names, and species.

Multiple sequence alignment of the 20GPeDo domain of the R3H domain-associated family

Multiple sequence alignment of the 20GPeDo domain of the R3H domain-associated family, showing conserved residues across various species.

GenBank accession numbers and protein names for various species including Ngru, ChRE, and others. The list includes species like Ngru1000010562, ChREDRAPT_153458, and many others, each followed by a protein name and a GenBank accession number.

Species abbreviations: Acry : Acidiphillus cryptum; Anig : Aspergillus niger; Avin : Azotobacter vinelandii; Blin : Brevibacterium linens; Ccin : Coprinopsis cinerea; Crel : Chlamy...

Fasta sequences of proteins from incomplete/unpublished genomes not in GenBank and assigned Fake gis

A large block of FASTA-formatted protein sequences. Each entry starts with a header line containing a reference ID and a protein name, followed by one or more lines of amino acid sequence in uppercase letters.

Back to Contents

6a. Phyletic distribution and multiple sequence alignment of the DNA glycosylase associated family of 20GFeD0 domains

Table with columns: gi, Domain architecture, Gene name, Len, Species, Class. Lists various protein entries and their taxonomic classifications.

Predicted secondary structure alignment showing conserved motifs and gaps across different species. Includes motifs like -HHHHH- and -GTMWPLTSSK-.

Species abbreviations: Aano : Aureococcus anophagefferens; Mbre : Monosiga brevicollis; Mxan : Myxococcus xanthus; Oluc : Ostreococcus lucimarinus; Otau : Ostreococcus tauri; Ptri

Fasta sequences of proteins from incomplete/unpublished genomes not in GenBank and assigned Fake gis

FASTA format sequences for various protein entries, starting with >gi [Ehux100018843] ref [jgi] [Emhul] [461717] est:ExtDc_igenesHEH_pg_C_70042.

