

# Clustering analyses of eukaryotic proteomes

Sergei Mekhedov

NCBI

# What you will have to listen about

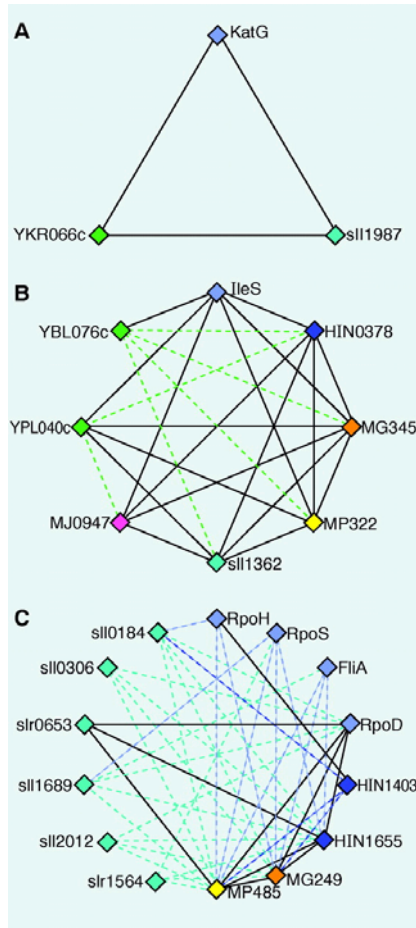
- Introduction: orthologs, paralogs, COGs, and KOGs
- Public KOG db: curated clusters for 7 species
- Orthology domains
- Relative vs. absolute KOGs
- Impossibility to continue the way KOGs were made before
- A semi-automatic approach to create curated KOGs
- Current status of KOG project for 50 eukaryotic species
- Efficient editing interface is a necessary prerequisite for accurate KOGs
- Problem of proteome updates in KOGs
- Features of new KOGnitor
- A way to improve gene predictions in available proteomes
- HMMs for KOGs – the way to save time on a long run
- Superfamilies and KOGs
- KOG applications
- Phylogenetic profile search
- What you need to be a real parasite

**Tatusov RL, Koonin EV, Lipman DJ.**

Science. 1997, 278(5338):631-637.

**A genomic perspective on protein families.**

# Examples of simple COGs based on BeTs



BeTs are best hits of a query protein in a particular proteome

solid lines – symmetrical BeTs

broken lines – asymmetrical BeTs

from Tatusov et al., 1997

**COG**

**Clusters of Orthologous Groups**

**KOG**

**Eukaryotic Clusters of Orthologous Groups**

**Public version of KOG database was created as a result of 2+ years of hard work by 20 programmers and computational biologists at NCBI. It included protein sequences from 7 eukaryotic species**

**Tatusov RL, Fedorova ND, Jackson JD, Jacobs AR, Kiryutin B, Koonin EV, Krylov DM, Mazumder R, Mekhedov SL, Nikolskaya AN, Rao BS, Smirnov S, Sverdlov AV, Vasudevan S, Wolf YI, Yin JJ, Natale DA. 2003 The COG database: an updated version includes eukaryotes. BMC Bioinformatics. 2003 4(1):41**

**Koonin EV, Fedorova ND, Jackson JD, Jacobs AR, Krylov DM, Makarova KS, Mazumder R, Mekhedov SL, Nikolskaya AN, Rao BS, Rogozin IB, Smirnov S, Sorokin AV, Sverdlov AV, Vasudevan S, Wolf YI, Yin JJ, Natale DA. 2004**

A comprehensive evolutionary classification of proteins encoded in complete eukaryotic genomes. *Genome Biol.* 2004;5(2):R7.

<http://www.ncbi.nlm.nih.gov/COG/new/>

[CE12874](#) [COGs](#) [unmask](#) [Genbank](#) [Blink](#)

7 proteins U [KOG0090](#) Signal recognition particle receptor, beta subunit (small G protein superfamily)

120 hits

		240 letters	
cel	1085	= <a href="#">CE12874</a> (240)	=
hsa	237	= <a href="#">Hs14917113</a> (271)	=
dme	231	= <a href="#">7295050_2</a> (221)	=
ath	197	= <a href="#">At5g05670</a> (260)	=
ath	184	= <a href="#">At2g18770</a> (260)	=
<hr/>			
sce	126	= <a href="#">YKL154w</a> (244)	=
cel	125	= <a href="#">CE15872</a> (184)	KOG0074
sce	122	= <a href="#">YOR094w</a> (183)	KOG0071
spo	118	= <a href="#">SPAC23H4.07c</a> (227)	=
hsa	114	= <a href="#">Hs4757774</a> (182)	KOG0074
hsa	109	= <a href="#">Hs20473689</a> (179)	KOG0070
sce	104	= <a href="#">YBR164c</a> (183)	KOG0072
hsa	104	= <a href="#">Hs4502203</a> (181)	KOG0070
dme	102	= <a href="#">7296589</a> (182)	KOG0070
hsa	101	= <a href="#">Hs6912244</a> (179)	KOG0070
cel	101	= <a href="#">CE00696</a> (181)	KOG0070
ath	101	= <a href="#">At3g22950</a> (183)	KOG0070
ath	101	= <a href="#">At2g15310</a> (205)	KOG0070
hsa	101	= <a href="#">Hs4502197_2</a> (173)	KOG0070
hsa	101	= <a href="#">Hs15208643_2</a> (145)	KOG0070
hsa	101	= <a href="#">Hs15208641_2</a> (168)	KOG0070
hsa	100	= <a href="#">Hs5031603</a> (200)	KOG0070
hsa	100	= <a href="#">Hs4502227</a> (181)	KOG0072
hsa	99	= <a href="#">Hs4502201</a> (181)	KOG0070
cel	97	= <a href="#">CE17122</a> (180)	KOG0070
sce	95	= <a href="#">YDL192w</a> (181)	KOG0070
spo	94	= <a href="#">SPBC1539.08</a> (184)	KOG0071
sce	93	= <a href="#">YDL137w</a> (181)	KOG0070
hsa	93	= <a href="#">Hs9910542</a> (198)	KOG0077
spo	92	= <a href="#">SPBC31F10.06c</a> (190)	KOG0077
hsa	92	= <a href="#">Hs4502229</a> (184)	KOG0073
spo	91	= <a href="#">SPBC4F6.18c</a> (180)	KOG0070
hsa	91	= <a href="#">Hs20472330</a> (234)	KOG0070
cel	91	= <a href="#">CE24147</a> (175)	KOG0071
ath	91	= <a href="#">At5g14670</a> (188)	KOG0070
ath	91	= <a href="#">At3g62290</a> (181)	KOG0070
ath	91	= <a href="#">At2g47170</a> (181)	KOG0070
ath	91	= <a href="#">At1g70490</a> (181)	KOG0070
ath	91	= <a href="#">At1g23490</a> (188)	KOG0070
ath	91	= <a href="#">At1g10630</a> (188)	KOG0070



Public KOG db has been edited extensively: BLAST outputs for every protein in clusters have been inspected by curators in search of errors

[CE15872](#) [COGs](#) [unmask](#) [Genbank](#) [Blink](#)

4 proteins R [KOG0074](#) GTP-binding ADP-ribosylation factor-like protein ARL3

499 hits

		184 letters	
cel	940	=	<a href="#">CE15872</a> (184) =
hsa	622	=	<a href="#">Hs4757774</a> (182) =
dme	487	=	<a href="#">7300792</a> (203) =
<hr/>			
hsa	455	-	<a href="#">Hs4502229</a> (184) KOG0073
spo	420	-	<a href="#">SPBC4F6.18c</a> (180) KOG0070
spo	420	-	<a href="#">SPAC22F3.05c</a> (186) KOG0073
sce	416	-	<a href="#">YBR164c</a> (183) KOG0072
ath	414	-	<a href="#">At2g18390</a> (185) KOG0073
dme	410	-	<a href="#">7296589</a> (182) KOG0070
cel	407	-	<a href="#">CE17122</a> (180) KOG0070
hsa	406	-	<a href="#">Hs4502203</a> (181) KOG0070
cel	406	-	<a href="#">CE00696</a> (181) KOG0070
cel	403	-	<a href="#">CE24147</a> (175) KOG0071
hsa	400	-	<a href="#">Hs4502201</a> (181) KOG0070
ath	399	-	<a href="#">At5g14670</a> (188) KOG0070
ath	399	-	<a href="#">At3g62290</a> (181) KOG0070
ath	399	-	<a href="#">At2g47170</a> (181) KOG0070
ath	399	-	<a href="#">At1g70490</a> (181) KOG0070
ath	399	-	<a href="#">At1g23490</a> (188) KOG0070
ath	399	-	<a href="#">At1g10630</a> (186) KOG0070
hsa	395	-	<a href="#">Hs4502227</a> (181) KOG0072
sce	394	-	<a href="#">YDL137w</a> (181) KOG0070
hsa	393	-	<a href="#">Hs4502211</a> (175) KOG0071
sce	391	-	<a href="#">YDL192w</a> (181) KOG0070
hsa	391	-	<a href="#">Hs4502209</a> (180) KOG0070
dme	391	-	<a href="#">7299026</a> (184) KOG0073
dme	390	-	<a href="#">7303081</a> (175) KOG0071
dme	386	-	<a href="#">7304353</a> (180) KOG0070
hsa	385	-	<a href="#">Hs4502205</a> (180) KOG0070
hsa	383	-	<a href="#">Hs6912244</a> (179) KOG0070
cel	383	-	<a href="#">CE27425</a> (184) KOG0073
dme	380	-	<a href="#">7295126</a> (179) KOG0070
spo	379	-	<a href="#">SPBC1539.08</a> (184) KOG0071
hsa	370	-	<a href="#">Hs4502197_2</a> (173) KOG0070
hsa	369	-	<a href="#">Hs20473689</a> (179) KOG0070
cel	361	-	<a href="#">CE02250</a> (180) KOG0072
cel	359	-	<a href="#">CE10492</a> (179) KOG0070
ath	356	-	<a href="#">At5g17060</a> (192) KOG0070
ath	353	-	<a href="#">At2g15310</a> (205) KOG0070
hsa	351	-	<a href="#">Hs1327624</a> (182) KOG0070

All precomputed BLAST output files were stored. With the number of proteins approaching 1,000,000 it is no more possible



## PUBLIC KOG database contains three types of clusters

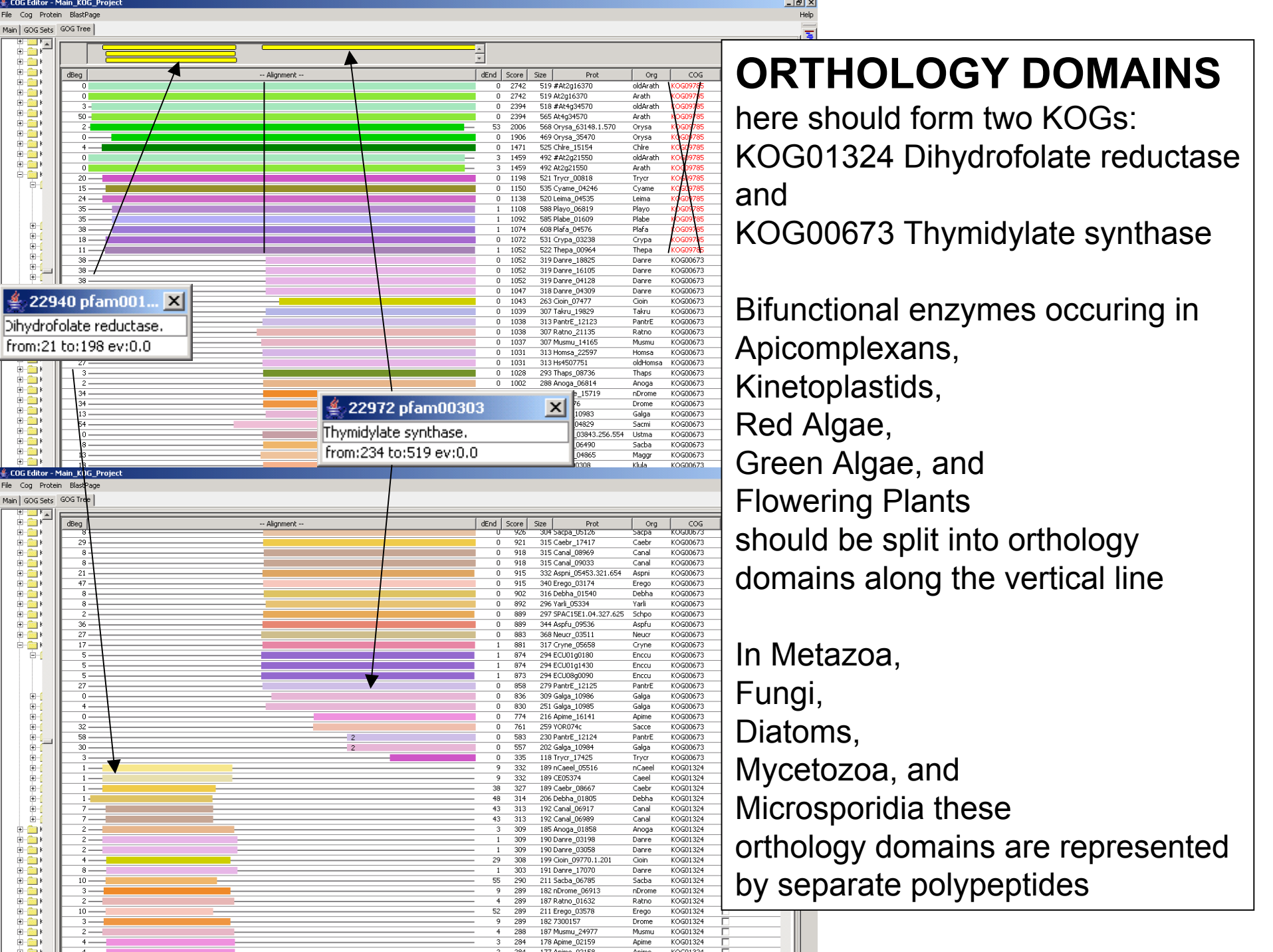
	sequences	in KOGs/LSEs		KOG	TWOGs	LSEs	total clusters
H.sapiens	38,638	26,324	68%	K 4,597	T 602	L 1,373	=6,572
D.melanogaster	13,703	10,517	77%	K 4,351	T 491	L 303	=5,145
C.elegans	20,751	17,101	82%	K 4,235	T 225	L 845	=5,305
S.pombe	5,035	4,233	84%	K 2,668	T 296	L 40	=3,004
S.cerevisiae	6,387	4,841	76%	K 2,762	T 293	L 132	=3,187
E.cuniculi	2,000	1,443	72%	K 1,073	T 6	L 29	=1,108
A.thaliana	26,406	24,154	91%	K 3,285	T 129	L 1,458	=3,872
Total	112,920	88,613	78%	4,952 KOGs	1,021 TWOGs	4,180 LSEs	

# Interface for KOG editing

- Changes are implemented on the fly
- No precomputed BLASTP results for all proteins, but only for a small selection AND only temporarily
- Real time BLASTP
- Precomputed RPS BLAST results
- Precomputed HMM search results
- Easy ways to switch to multiple alignments

KOG db deals not with orthologs but with

## **ORTHOLOGY DOMAINS**



# ORTHOLOGY DOMAINS

here should form two KOGs:  
 KOG01324 Dihydrofolate reductase  
 and  
 KOG00673 Thymidylate synthase

Bifunctional enzymes occurring in  
 Apicomplexans,  
 Kinetoplastids,  
 Red Algae,  
 Green Algae, and  
 Flowering Plants  
 should be split into orthology  
 domains along the vertical line

In Metazoa,  
 Fungi,  
 Diatoms,  
 Mycetozoa, and  
 Microsporidia these  
 orthology domains are represented  
 by separate polypeptides

**Current status of KOG db:** 50 Eukaryotic species; 46 species KOGnitorized, 1 updated;  
 712,427 protein sequences; two versions are present for 5 proteomes (+ 123,981 sequences)  
 21,816 clusters: 18,628 KOGs, 3,188 LSEs+mixed clusters

sequence data

species name	Taxa	sequences	In KOGs/LSEs	%	KOGs	published	source
<i>Chlamydomonas reinhardtii</i>	Chlorophyta (green algae)	19,922	10,194	51%	4,653	No	JGI
<i>Cyanidioschyzon merolae</i>	Rhodophyta (red algae)	5,040	3,769	74%	2,735		JGI
<i>Dictyostelium discoideum</i>	Mycetozoa	13,677	9,162	67%	4,573		GB
<i>Encephalitozoon cuniculi</i>	OLD Microsporidia	2,003	1,459	73%	1,132		GB
<i>Entamoeba histolytica</i>	Entamoebidae	10,088	7,728	77%	2,930		TIGR
<i>Giardia lamblia</i>	Diplomonadida group	6,573	2,105	32%	1,242	No	GB
<i>Leishmania major</i>	Euglenozoa	8,071	7,558	93%	4,502		TIGR
<i>Trypanosoma brucei</i>	Euglenozoa	8,291					GB
<i>Trypanosoma cruzi</i>	Euglenozoa	19,666	16,844	86%	5,023		GB
<i>Thalassiosira pseudonana</i>	Bacillariophyta (diatoms)	11,406	9,311	81%	4,033		JGI
<i>Cryptosporidium parvum</i>	Apicomplexa	3,410	2,212	66%	1,780		GB
<i>Plasmodium berghei</i>	Apicomplexa	11,782	7,004	59%	3,461	No	GB
<i>Plasmodium falciparum</i>	Apicomplexa	5,279	3,365	64%	3,403		GB
<i>Plasmodium yoelii</i>	Apicomplexa	7,868					GB
<i>Theileria parva</i>	Apicomplexa	4,080					GB
<i>Arabidopsis thaliana</i>	NEW Viridiplantae	26,759	25,781	96%	6,953		GB
<i>Arabidopsis thaliana</i>	OLD Viridiplantae	27,138	25,410	94%	6,925		GB
<i>Oryza sativa</i>	Viridiplantae	57,272	38,260	67%	6,701		TIGR
<i>Aspergillus fumigatus</i>	Fungi; Ascomycota; Pezizomycotina	10,029	8,959	89%	5,232	No	GB
<i>Aspergillus nidulans</i>	Fungi; Ascomycota; Pezizomycotina	10,317	9,629	93%	5,283	No	GB
<i>Gibberella zeae</i>	Fungi; Ascomycota; Pezizomycotina	12,221	11,153	91%	5,910	No	GB
<i>Magnaporthe grisea</i>	Fungi; Ascomycota; Pezizomycotina	11,381	9,043	79%	5,397		GB
<i>Neurospora crassa</i>	Fungi; Ascomycota; Pezizomycotina	10,478	8,206	78%	5,559		GB
<i>Candida albicans</i>	Fungi; Ascomycota; Saccharomycotina	14,274	10,952	77%	4,094		GB
<i>Candida glabrata</i>	Fungi; Ascomycota; Saccharomycotina	5,216	4,992	96%	3,790	No	GB
<i>Debaryomyces hansenii</i>	Fungi; Ascomycota; Saccharomycotina	6,350	5,700	90%	4,048		GB
<i>Eremothecium gossypii</i>	Fungi; Ascomycota; Saccharomycotina	4,757	4,619	97%	3,810	No	GB
<i>Kluyveromyces lactis</i>	Fungi; Ascomycota; Saccharomycotina	5,366	5,058	94%	3,980		GB
<i>Kluyveromyces waltii</i>	Fungi; Ascomycota; Saccharomycotina	5,240	5,055	96%	3,919	No	MIT
<i>Yarrowia lipolytica</i>	Fungi; Ascomycota; Saccharomycotina	6,596	5,728	87%	4,018		GB
<i>Saccharomyces bayanus</i>	Fungi; Ascomycota; Saccharomycotina	12,040	8,488	70%	5,879		MIT
<i>Saccharomyces cerevisiae</i>	OLD Fungi; Ascomycota; Saccharomycotina	6,387	6,146	96%	4,427		GB
<i>Saccharomyces mikatae</i>	Fungi; Ascomycota; Saccharomycotina	10,350	8,481	82%	5,828		MIT
<i>Saccharomyces paradoxus</i>	Fungi; Ascomycota; Saccharomycotina	10,601	8,375	79%	6,005		MIT
<i>Schizosaccharomyces pombe</i>	OLD Fungi; Ascomycota; Schizosaccharomycetes	5,064	4,461	88%	3,360		GB
<i>Cryptococcus neoformans</i>	Fungi; Basidiomycota; Hymenomycetes	6,640	5,471	82%	3,702	No	GB
<i>Ustilago maydis</i>	Fungi; Basidiomycota; Ustilaginomycetes	6,685	5,221	78%	3,779	No	GB
<i>Anopheles gambiae</i>	Metazoa; Insecta; Diptera; Nematocera	13,895	12,405	89%	5,946		GB
<i>Apis mellifera</i>	Metazoa; Insecta; Hymenoptera	16,997	14,785	87%	4,906	No	EMBL
<i>Drosophila melanogaster</i>	NEW Metazoa; Insecta; Diptera; Brachycera	18,873					GB
<i>Drosophila melanogaster</i>	OLD Metazoa; Insecta; Diptera; Brachycera	13,936	11,552	83%	6,241		GB
<i>Caenorhabditis briggsae</i>	Metazoa; Nematoda	19,523	17,963	92%	6,938		GB
<i>Caenorhabditis elegans</i>	NEW Metazoa; Nematoda	21,193					GB
<i>Caenorhabditis elegans</i>	OLD Metazoa; Nematoda	21,246	20,461	96%	7,175		GB
<i>Ciona intestinalis</i>	Metazoa; Chordata; Ascidiacea	16,139	12,101	75%	5,684		JGI
<i>Danio rerio</i>	Vertebrata; Teleostomi	19,893	19,676	98%	4,426		GB
<i>Gallus gallus</i>	Vertebrata; Aves	28,751	26,613	93%	6,894	No	EMBL
<i>Homo sapiens</i>	NEW Vertebrata; Mammalia	27,524					GB
<i>Homo sapiens</i>	OLD Vertebrata; Mammalia	39,436	29,835	77%	9,304		GB
<i>Mus musculus</i>	Vertebrata; Mammalia	26,464	24,784	94%	9,284		GB
<i>Pan troglodytes</i>	GB Vertebrata; Mammalia	22,225					GB
<i>Pan troglodytes</i>	EMBL Vertebrata; Mammalia	38,845	36,984	95%	8,923		EMBL
<i>Rattus norvegicus</i>	Vertebrata; Mammalia	21,531					GB
<i>Takifugu rubripes</i>	Vertebrata; Teleostomi	33,357	31,879	96%	6,909		GB
<i>Tetraodon nigroviridis</i>	Vertebrata; Teleostomi	28,262	25,520	90%	7,033	No	GB

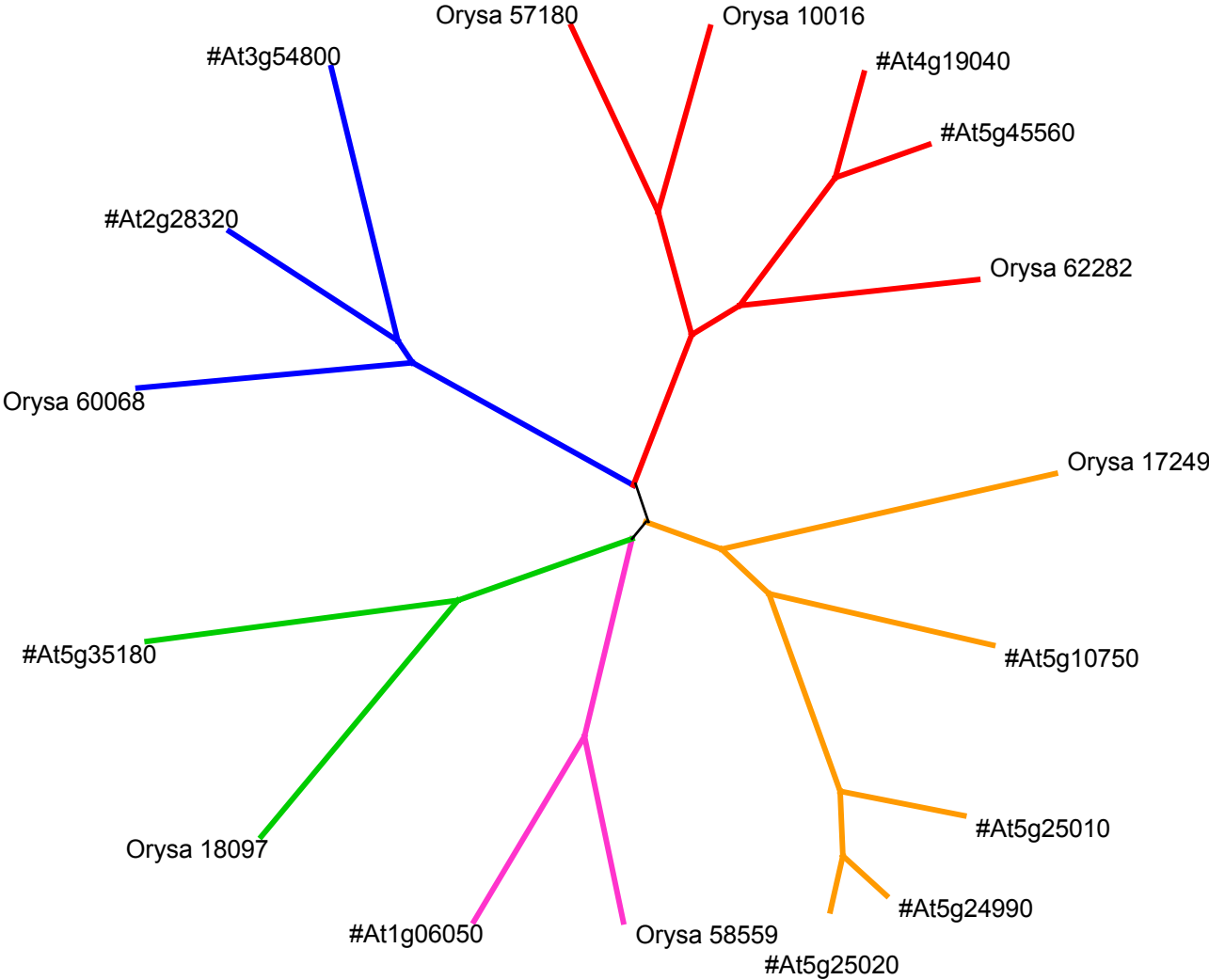
These are the results of ongoing research KOG project. These results should be considered preliminary

Absolute KOGs are not possible

KOGs are relative to the set of proteomes used in their construction

This sad truth constitutes the major difficulty in creating algorithms that would allow to add proteomes to KOGs automatically

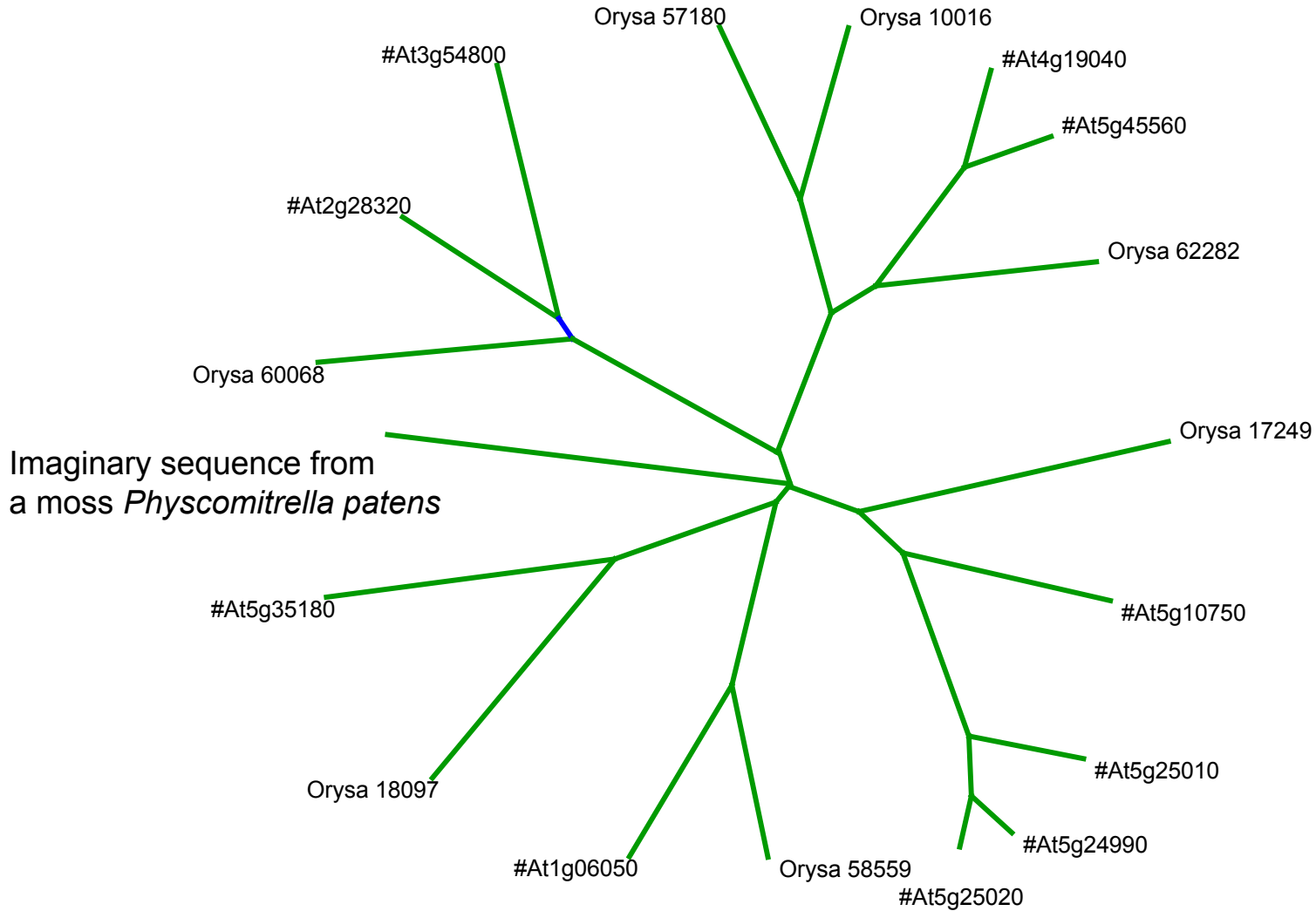
# 5 related KOGs



5 KOGs



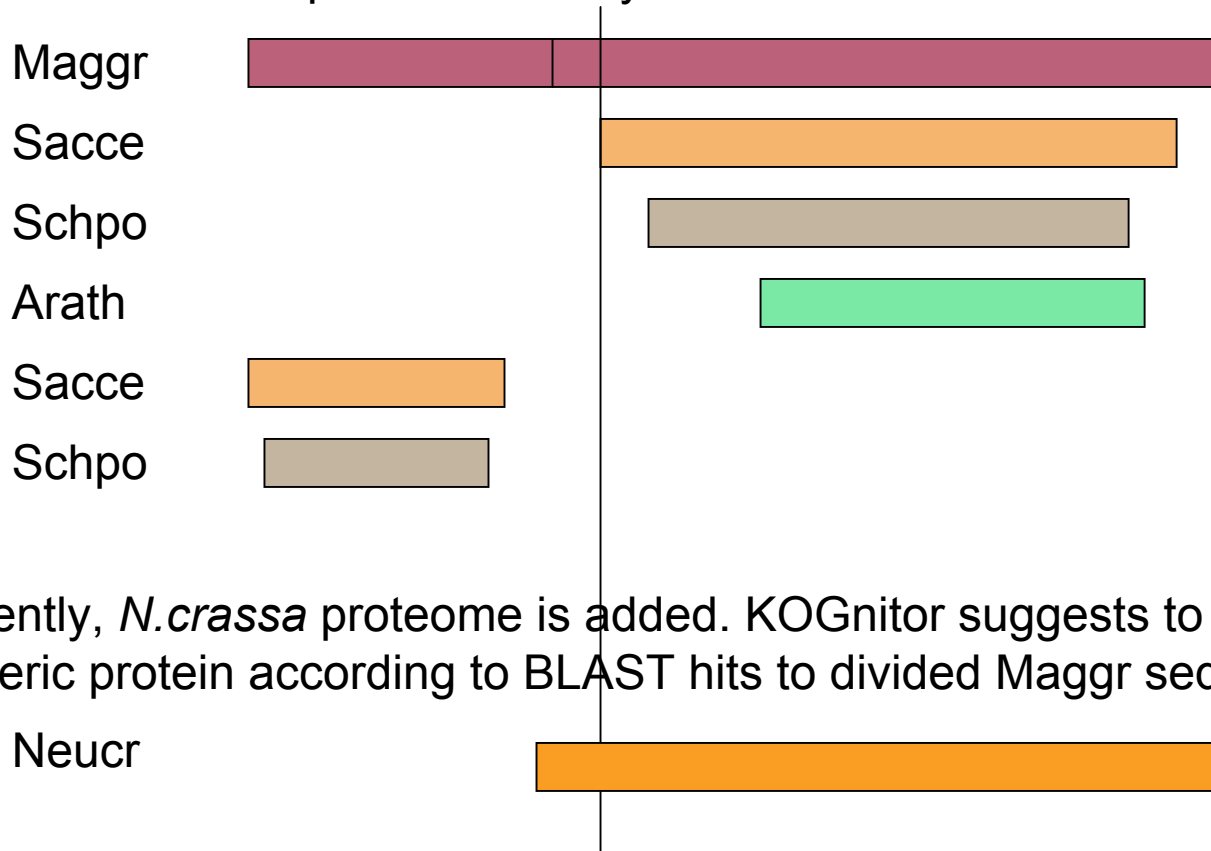
1 KOG rooted by one sequence from a new genome





# Even dividing sequences into orthology domains depends on the proteome set

A chimeric sequence from *M.grisea* was divided between two KOGs when its proteome was added to KOGs with just two fungal species. It was not possible to find division point accurately



Subsequently, *N.crassa* proteome is added. KOGnitor suggests to divide non-chimeric protein according to BLAST hits to divided Maggr sequence

Due to extensive manual curation the public KOG db project contains substantial added value if compared with results of automatic clustering using BeTs.

It would be unwise and even silly to discard these results.

One of important features of public KOGs were LSEs.

# Lineage Specific Expansions (LSE)

- [-] KOG09723 FOG: CUB dc
- [-] KOG09724 FOG: Cadher
- [-] KOG09725 Receptor-like
- [-] KOG09726 Far-red impai
- [-] KOG09727 Far-red impai
- [-] KOG09728 Far-red impai
- [-] KOG09729 Myosin bindir
- [-] KOG09730 Follistatin-like
- [-] KOG09731 Predicted pol
- [-] KOG09732 New Cog
- [-] KOG09733 Predicted hyc
- [-] KOG09734 Predicted hyc
- [-] KOG09735 Auxin respor
- [-] KOG09736 Auxin respor
- [-] KOG09737 Auxin respor
- [-] KOG09738 Auxin respor
- [-] KOG09739 Auxin respor
- [-] KOG09740 AP2 domain t
- [-] KOG09741 Ovule develo
- [-] KOG09742 Ovule develo
- [-] KOG09743 Ovule develo
- [-] KOG09744 S-adenosyl-L
- [-] KOG09745 S-adenosyl-L
- [-] KOG09746 Uncharacteri
- [-] KOG09747 Uncharacteri
- [-] KOG09748 Uncharacteri
- [-] KOG09749 Uncharacteri
- [-] KOG09750 Photosystem
- [-] KOG09751 Uncharacteri
- [-] KOG09752 Uncharacteri
- [-] KOG09753 Uncharacteri
- [-] KOG09754 Uncharacteri
- [-] KOG09755 Uncharacteri
- [-] KOG09756 Uncharacteri
  - [-] Arabidopsis thaliana
  - [-] Arabidopsis thaliana
  - [-] Oryza sativa (japoni
    - Oryza\_60068
  - [-] Thalassiosira pseud
    - Thaps\_00733
    - Thaps\_04196
    - Thaps\_08147
- [-] KOG09757 Uncharacteri
- [-] KOG09758 Uncharacteri
- [-] KOG09759 Uncharacteri
- [-] KOG09760 Uncharacteri
- [-] KOG09761 Uncharacteri
- [-] KOG09762 Uncharacteri
- [-] KOG09763 Uncharacteri
- [-] KOG09764 Uncharacteri
- [-] KOG09765 Uncharacteri
- [-] KOG09766 Uncharacteri
- [-] KOG09767 Uncharacteri
- [-] KOG09768 Uncharacteri



dBeg	dEnd	Score	Size	Prot	Org	COG	±	COG
0	2	0	4068	773 Orysa_60068	Orysa	KOG09756		
28	2	-71	2439	737 At2g28320	Arath	KOG09756		
31	2	-78	2229	733 At3g54800	Arath	KOG09756		
31	2	-126	2105	709 #At3g54800	oldArath	KOG09756		LSE00404
228	2	-653	1448	446 #At2g28320	oldArath	KOG09756		LSE00404
35	2	-110	1024	719 At5g45560	Arath	KOG09755		
35	2	-112	992	718 At4g19040	Arath	KOG09755		
35	2	-222	780	663 #At5g45560	oldArath	KOG09755		LSE00404
-19	2	11	769	778 At5g35180	Arath	KOG09754		
35	2	-190	691	679 #At4g19040	oldArath	KOG09755		LSE00404
-19	2	67	610	806 #At5g35180	oldArath	KOG09754		LSE00404
-13	2	182	608	849 Orysa_18097	Orysa	KOG09754		
88	4	-295	602	626 Orysa_62282	Orysa	KOG09755		
118	4	-177	521	685 Orysa_57180	Orysa	KOG09755		
130	4	-204	498	685 Orysa_10016	Orysa	KOG09755		
518	2	-974	462	313 #At1g06050	oldArath	KOG09752		LSE00404
518	2	-974	462	313 At1g06050	Arath	KOG09752		
503	2	-967	461	302 #At5g10750	oldArath	KOG09753		LSE00404
503	2	-967	461	302 At5g10750	Arath	KOG09753		
532	2	-952	432	350 Orysa_58559	Orysa	KOG09752		
37	2	-556	411	237 Orysa_36219	Orysa	KOG09752		
506	2	-989	401	286 #At5g25010	oldArath	KOG09753		LSE00404
506	2	-989	401	286 At5g25010	Arath	KOG09753		
506	2	-980	380	294 At5g24990	Arath	KOG09753		
506	2	-980	380	294 #At5g24990	oldArath	KOG09753		LSE00404

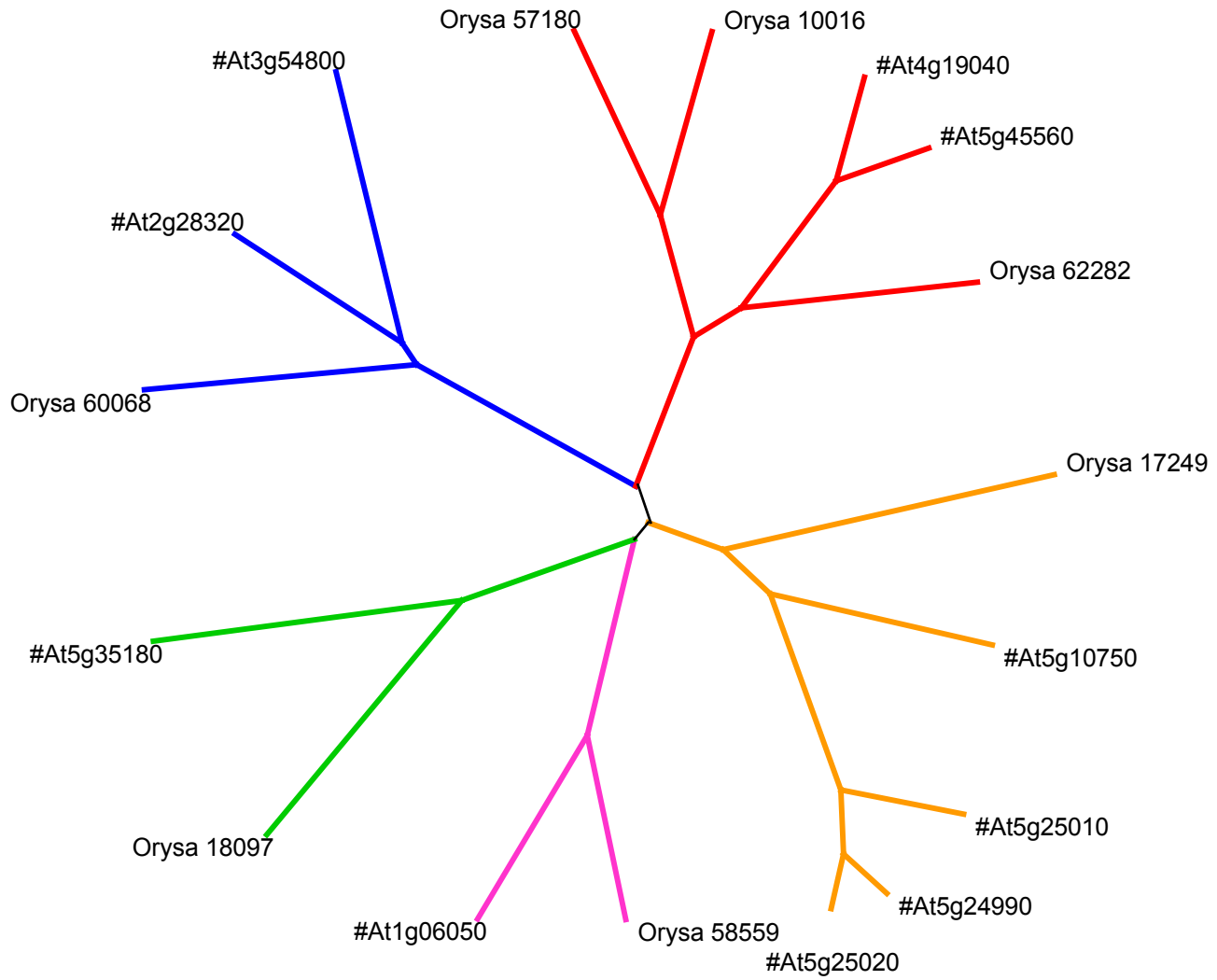
Length	679	Score	271 bits (691)	Expect	0.0
Identities	225/778 (28%)	Positives	340/778 (43%)	Gaps	163/778 (20%)

Start	End	Sequence	Length
419	419	QIAES	5
547	547	DSNCT	5
472	472	NEMER	5
593	593	IKSDE	5
532	532	FKDSE	5
652	652	-LLEE	5
587	587	SLLQRFVDCGDEFNRNRLKLIPLVPR-----ID	614
711	711	VDIGSSTVARGVVS LVLG YLNMNLVIEMAF LVQGTQEE LPEFL LGTCRLN YLDASKAV	768
615	615	VDIGSSTVA GV+ LV+G + +LV+EMAF LVQ NT EE PE L+G R++++ S A+	762
		VDIGSSTVANGVLGLVIGVITSLVVEMAFLVQANTAE EQPERLIGAVRVSHIELSSAI	

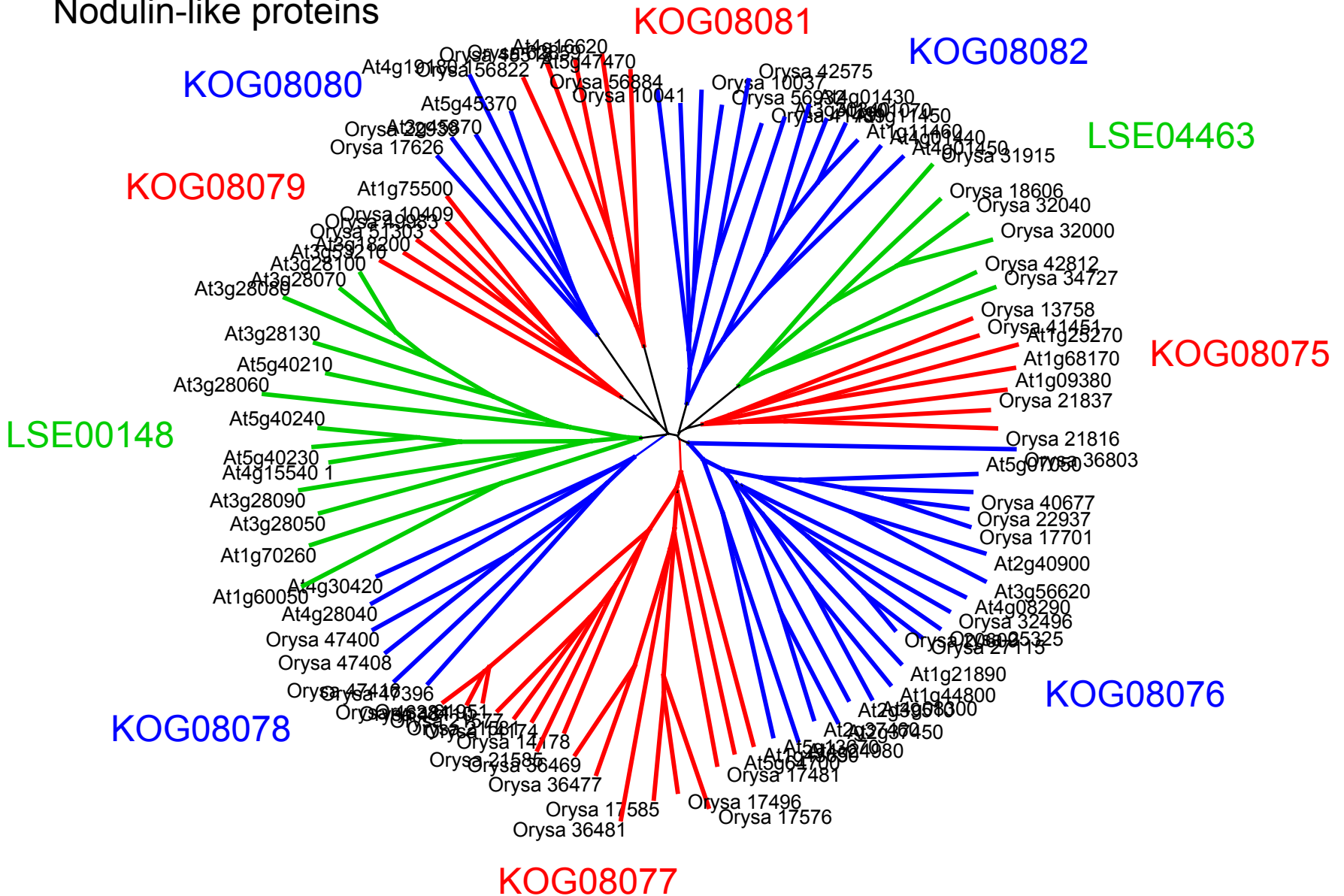
Membership in LSEs [from the public KOG database] helped enormously in editing newly created KOGs. In cases with numerous paralogs all of them were first assigned to a preexisting LSE which was subsequently split into KOGs by phylogenetic analyses.

LSE00404 →  
Uncharacterized proteins

5 KOGs



LSE00148 → 8 KOGs, 2 LSEs  
Nodulin-like proteins



With the number of protein sequences approaching 1,000,000 BLAST all-against-all does not seem to be possible

# New KOGnitor

Proteomes, one at a time are added to KOGs with KOGnitor program (Y.Wolf, unpublished) which uses BLASTP computed for every protein in the proteome against a selected set of KOGnitorized proteomes.



KOGnitor creates a list of sequences that potentially should form new KOGs

and

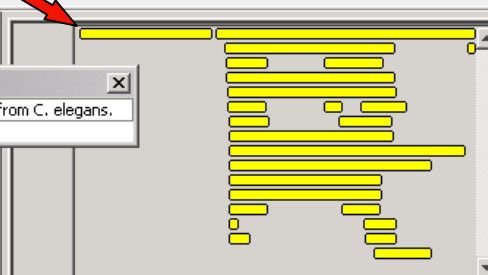
a list of sequences that are potentially chimeric and have to be split into separate orthology domains

Unfortunately, due to unidirectional hits and necessarily relaxed BLAST score cutoffs KOGnitor produces high rate of false positives. We checked completely one proteome, *Cryptosporidium parvum* and found 14% of false positives. This rate is unacceptable and we developed a plan to fix it by simultaneously applying two independent methods of sequence analysis – KOGnitor and HMM search

A plan to create curated set of KOGs semi-automatically by constructing expanding KOG\_HMM profiles:

1. Build HMM profiles for previously curated KOGs
2. Add proteomes to KOGs, one at a time with KOGnitor
3. Confirm suggested new KOGs and divisions of potentially chimeric sequences into separate orthology domains by manual curation
4. Run HMM search using KOG\_HMM profiles in KOG db
5. Compare results of KOGnitor and HMM search and manually check ONLY conflicting sequences. The great majority of sequences would not require checking.
6. Whenever necessary, add the new sequences to HMM profiles such that all legitimate members of KOGs are recognized by updated HMM profiles
7. Such an approach would result not only in curated set of KOGs, but also in creating a library of KOG\_HMM which presumably will be used more and more with accumulation of new genomes. We anticipate that BLAST searches will be largely replaced in the future by profile searches when new proteomes are KOGnitorized.

KOG00426 Ubiquitin-protein ligase  
 KOG00427 Ubiquitin conjugating enzyme  
 KOG00428 Non-canonical ubiquitin-conjugating enzyme 1  
 26125 pfam03676  
 Uncharacterised protein family (UPF0183). This family of proteins includes Lin-10 from *C. elegans*.  
 from:15 to:517 ev:9.0E-31  
 KOG00432 needs\_to\_be\_divided\_Valyl-tRNA synthetase  
 KOG00433 Isoleucyl-tRNA synthetase  
 Anopheles gambiae str. PEST | Anoga  
 Apis mellifera | Apime  
 Arabidopsis thaliana | Arath  
 Arabidopsis thaliana | oldArath  
 Aspergillus fumigatus Af293 | Aspflu  
 Aspflu\_09358  
 Aspergillus nidulans FGSC A4 | Aspni  
 Aspni\_03797.519.1518  
 Caenorhabditis briggsae | Caebr  
 Caenorhabditis elegans | Caeel  
 Caenorhabditis elegans | nCaeel  
 Candida albicans SC5314 | Canal  
 Candida glabrata CB5138 | Cangl  
 Chlamydomonas reinhardtii | Chlre  
 Ciona intestinalis | Cloin  
 Cryptococcus neoformans var. neoformans JEC21 | Cryne  
 Cyanidioschyzon merolae | Cyame



KOGnitor creates a list of potentially chimeric sequences

dBeg	-- Alignment --	dEnd	Score	Size	Prot	Org	COG	±	SC
0		98	Aspni_03797.51...	Aspni	KOG00433				
0		11	4101	1009	Aspflu_09358	Aspflu	KOG00433		
0		0	2733	516	Aspni_03797.1...	Aspni	KOG02819		
0		0	2465	997	Gibze_11395	Gibze	KOG00433		
0		5	2371	957	Neucr_07354	Neucr	KOG00433		
0		6	2261	974	Maggr_10303	Maggr	KOG00433		
35		1	1863	972	Yarli_04934	Yarli	KOG00433		
7		0	1753	980	Thaps_05451	Thaps	KOG00433		
31		2	1736	973	SPCC1885.08c	Schpo	KOG00433		
101		16	1732	1093	At5g49030	Arath	KOG00433		
16		16	1732	1008	#At5g49030	oldArath	KOG00433		
48		1	1722	1032	Cryne_04497	Cryne	KOG00433		
		0	1697	990	Klula_02672	Klula	KOG00433		
		2	1696	414	Aspflu_09357	Aspflu	KOG02819		
		8	1673	1079	Cyame_00546	Cyame	KOG00433		
		12	1653	1057	Orysa_25681	Orysa	KOG00433		
		1	1649	985	Kluwa_00429	Kluwa	KOG00433		
		3	1645	982	Cangl_04222	Cangl	KOG00433		
		0	1642	973	Canal_10752	Canal	KOG00433		
		0	1638	973	Canal_10803	Canal	KOG00433		
		0	1623	988	Erego_00480	Erego	KOG00433		
		4	1619	982	Danre_02073	Danre	KOG00433		
		1	1616	1034	Dicdi_12992	Dicdi	KOG00433		

Aspni\_03797 was predicted by KOGnitor to be chimeric. BLASTP and RPS BLAST in CDD demonstrate that this is the case. The sequence needs to be split between KOG02819 and KOG00433

Alignment		End
KPQHSAS PQGPFGRHVYVNLFGPSPYGEYIPP-TSSAYGTYVLSYPGVAFSFP LQHSA		162
KPQE + S QGP FRH+YNRLFGPSPYGEY PP S YGTYVLSYPG+AFSFP LQ+SA		
KPQE QAVSQGQPSFRHIYNRLFGPSPYGEYTPPGDQSPYGTYYVLSYPGIAFSFP LQNSA		64
WSEQCDFVALLSSSAALPATSMSIFQGPSWPEVRDKLFTQPQYPRSPALAGSKSEFLPD		222
W+EQCDFVALLSSSAALPATSMA+IFQG SWPE RDKLF++QPQYPRSPAL+GK+++ + D		
65 WAEQCDFVALLSSSAALPATSMAIFQGSWPEARDKLF SRQPQYPRSPALSGKMRDLVSD		124
EIEEFVILGAGKLEVTRRSTPSTYIRLSQTTQDLIAEFGPPDAIYRKNDRRISIHRAAG		282
EIEEF++ GAGK+EV RRS+P T I LS+TTPQDLIAEFGPPDAIYRKNDRRISIHRAAG		
125 EIEEFIVFGAGKMEVIRRSSPPTSITLSETTPQDLIAEFGPPDAIYRKNDRRISIHRAAG		184
GHGRDEMIHMSPASGRGIELHDTQSSNNSVSDSDEGISQTSALDPSSLPSECFNYFH		342
+ + +HMSFP+ GRGI++ DTDQSS+NSV+DSDS +S + +DPSSLP+ECFFNYFH		
185 NNTATD TLHMSPPGRGIDVTD TDQSSNSVTDSDSDEEVPVNNIDPSSLPTECFNYFH		244

KOG db contains comprehensive records about potentially chimeric (misassembled) protein sequences which can be used to improve gene predictions in numerous proteomes 😊.

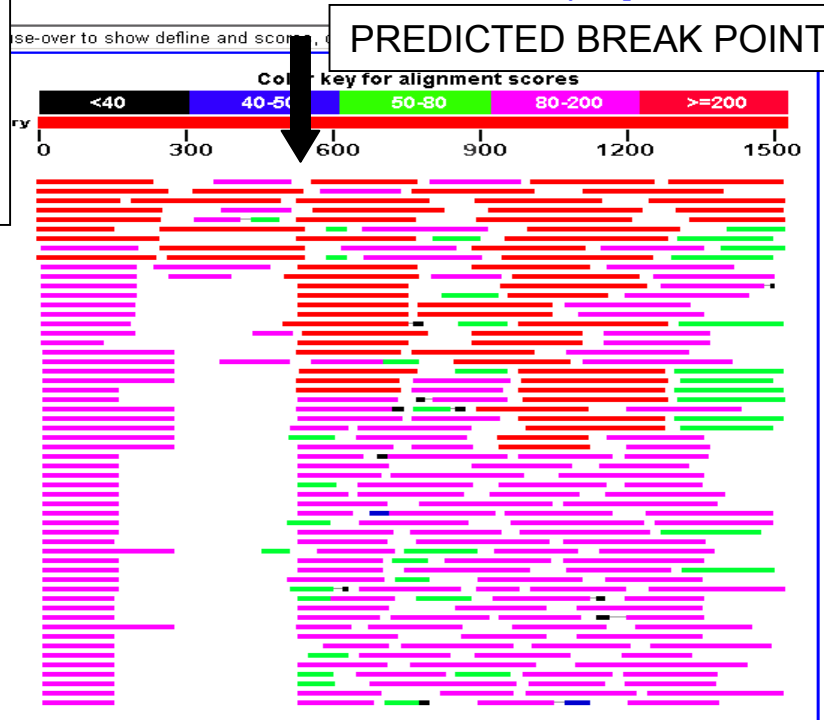
There are thousands of such sequences 😞

Is there a way to automatically confirm predicted break points?

# An independent way to confirm predicted chimeric sequences: TBLASTN in dbEST

Chimeric Aspni\_03797  
sequence was used as  
query in TBLASTN in  
dbEST

Distribution of 1064 Blast Hits on the Query Sequence



Sequences producing significant alignments:

Query ID	Accession	Species	Score (Bits)	E Value
<a href="#">gi 70826686 gb DR710376.1 </a>	Asn_12133	Aspergillus niger pBlues...	364	4e-98
<a href="#">gi 45923932 gb CF818054.1 </a>	EST695436	Coccidioides posadasii s...	353	7e-95
<a href="#">gi 48559074 gb CO028230.1 </a>	EST806614	Coccidioides posadasii s...	335	2e-89
<a href="#">gi 25129738 gb CA581347.1 </a>	EST001022	Mycelium and yeast cells...	303	9e-80
<a href="#">gi 48516926 gb CO010037.1 </a>	EST798372	Coccidioides posadasii s...	300	7e-79

HMM profiles were built for curated versions (containing less than 50, typically a set of 12+ proteomes) of ca. 15,000 KOGs.

HMM search results in KOG db with HMM profiles as queries provide a new tool to detect conflicts between KOGnitor and HMM search and thus sort out false positives

# How the HMM search results work in the context of editing interface

HMM search results are imported into KOG db. Cutoff point is set manually for the score

- [-] KOG00391 DnaJ family DNA-dependent ATPase
- [-] KOG00392 SMF2 family DNA-dependent ATPase domain-containing protein
- [-] KOG00393 Ras-related small GTPase; Rho type
- [-] KOG00394
- [-] KOG00409 Predicted dehydrogenase
- [-] KOG00410 Predicted GTP binding protein
- [-] KOG00411 Uncharacterized membrane protein
- [-] KOG00412 Golgi transport complex COD1 protein
- [-] KOG00413 Uncharacterized conserved protein related to condensin complex
- [-] KOG00414 Chromosome condensation complex Condensin; subunit D2
- [-] KOG00426 Ubiquitin-protein ligase
- [-] KOG00427 Ubiquitin conjugating enzyme
- [-] KOG00428 Non-canonical ubiquitin conjugating enzyme 1
- [-] KOG00429 Ubiquitin-conjugating enzyme-related protein Ft1; involved in pro
- [-] KOG00430 Xanthine dehydrogenase
- [-] KOG00431 Auxilin-like protein and related proteins containing DnaJ domain
- [-] KOG00432 needs\_to\_be\_divided\_Valyl-tRNA synthetase
- [-] **KOG00433 Isoleucyl-tRNA synthetase**
- [-] KOG00434 Isoleucyl-tRNA synthetase
- [-] KOG00435 Leucyl-tRNA synthetase
- [-] KOG00436 Methionyl-tRNA synthetase
- [-] KOG00437 Leucyl-tRNA synthetase
- [-] KOG00438 Mitochondrial ribosomal protein L2
- [-] KOG00439 VAMP-associated protein involved in inositol metabolism
- [-] KOG00440 Cell cycle-associated protein Mob1-1
- [-] KOG00441 Cu2+/Zn2+ superoxide dismutase SOD1
- [-] KOG00442 Structure-specific endonuclease ERCC1-XPF; catalytic component
- [-] KOG00443 Actin regulatory proteins (gelsolin/villin family)
- [-] KOG00444 Cytoskeletal regulator Flightless-I (contains leucine-rich and gels
- [-] KOG00445 Actin regulatory protein supervillin (gelsolin/villin family)

Threshold is not yet defined

Name	COG	Sup	Score	E-Value	pf	Ing
Klula_02672	KOG00433		1.284E3	0.0	<input type="checkbox"/>	990
Cangl_04222	KOG00433		1.255E3	0.0	<input type="checkbox"/>	982
Aspni_03797	DIV		1.189E3	0.0	<input type="checkbox"/>	1518
Caebr_08776	KOG00433		1.167E3	0.0	<input type="checkbox"/>	947
Aspfu_09358			1.156E3	0.0	<input type="checkbox"/>	1009
Gibze_11395	KOG00433		1.148E3	0.0	<input type="checkbox"/>	997
Neucr_07354	KOG00433		1.146E3	0.0	<input type="checkbox"/>	957
Kluwa_00429	KOG00433		1.146E3	0.0	<input type="checkbox"/>	985
Erego_00480	KOG00433		1.132E3	0.0	<input type="checkbox"/>	988
Dicdi_12992	KOG00433		1.126E3	0.0	<input type="checkbox"/>	1034
Galga_14607	KOG00433		1.109E3	0.0	<input type="checkbox"/>	773
Yarli_04934	KOG00433		1.109E3	0.0	<input type="checkbox"/>	972
Galga_14610	KOG00433		1.076E3	0.0	<input type="checkbox"/>	767
Cryne_04497	KOG00433		1.035E3	0.0	<input type="checkbox"/>	1032
Debha_04672	KOG00433		1.016E3	0.0	<input type="checkbox"/>	988
Tetni_18015	KOG00433		974.6	0.0	<input type="checkbox"/>	933
Galga_14608	KOG00433		963.8	0.0	<input type="checkbox"/>	723
Galga_14609	KOG00433		952.4	0.0	<input type="checkbox"/>	728
Canal_10752	KOG00433		937.8	0.0	<input type="checkbox"/>	973
Canal_10803	KOG00433		934.3	0.0	<input type="checkbox"/>	973
Ratno_18798			922.8	0.0	<input type="checkbox"/>	1136
Anoga_05900	KOG00433		861.0	0.0	<input type="checkbox"/>	806
Ustma_01543	KOG00433		595.4	0.0	<input type="checkbox"/>	851
Pantr_01697			552.5	0.0	<input type="checkbox"/>	620
Galga_14611	KOG00433		308.2	0.0	<input type="checkbox"/>	470
Takru_16952	KOG00433		274.7	0.0	<input type="checkbox"/>	485
Thepa_04012			107.8	3.1E-27	<input type="checkbox"/>	1221
Trycr_07252	KOG00434		68.8	1.7E-15	<input type="checkbox"/>	1156
Enthi_08456	KOG00434		57.6	4.0E-12	<input type="checkbox"/>	1056
Cyame_04722	KOG00434		49.0	3.4E-11	<input type="checkbox"/>	1207
Canal_11674	KOG00434		48.5	3.4E-11	<input type="checkbox"/>	1088
Canal_11824	KOG00434		48.5	3.4E-11	<input type="checkbox"/>	1088
Anoga_00791	KOG00434		46.4	3.8E-11	<input type="checkbox"/>	1213
Orysa_57566	KOG00434		38.7	5.8E-11	<input type="checkbox"/>	1386
Playo_01798			37.4	6.1E-11	<input type="checkbox"/>	1294
Cangl_01703	KOG00434		36.2	6.6E-11	<input type="checkbox"/>	1072
7296537	KOG00434		35.0	7.0E-11	<input type="checkbox"/>	1081
nDrome_08034			35.0	7.0E-11	<input type="checkbox"/>	1229
nDrome_09237			35.0	7.0E-11	<input type="checkbox"/>	1229
nDrome_09754			35.0	7.0E-11	<input type="checkbox"/>	1229
Plafa_00323	KOG00434		32.2	8.1E-11	<input type="checkbox"/>	1272
Danre_05440	KOG00434		27.0	1.1E-10	<input type="checkbox"/>	1271
Danre_14019	KOG00434		27.0	1.1E-10	<input type="checkbox"/>	1271
Kluwa_00189	KOG00434		25.4	1.2E-10	<input type="checkbox"/>	1099
Ereg_02472	KOG00434		25.2	1.2E-10	<input type="checkbox"/>	1072

Property | Proteins | Prot-Cogs | **HMM Global** | HMM Local



Main | GOG Sets | GOG Tree

- [-] KOG00391 Smr 2 family DNA-dependent ATPase
- [+] KOG00392 SNF2 family DNA-dependent ATPase domain-containing protein
- [+] KOG00393 Ras-related small GTPase; Rho type
- [+] KOG00394 Ras-related GTPase
- [+] KOG00395 Ras-related GTPase
- [+] KOG00396 Uncharacterized conserved protein
- [+] KOG00397 60S ribosomal protein L11
- [+] KOG00398 Mitochondrial/chloroplast ribosomal protein L5/L7
- [+] KOG00399 Glutamate synthase
- [+] KOG00400 40S ribosomal protein S13
- [+] KOG00401 Translation initiation factor 4F; ribosome/mRNA-bridging subunit (
- [+] KOG00402 60S ribosomal protein L37
- [+] KOG00403 Neoplastic transformation suppressor Pcdcd4/MA-3; contains MA3
- [+] KOG00404 Thioredoxin reductase
- [+] KOG00405 Pyridine nucleotide-disulphide oxidoreductase
- [+] KOG00406 Glutathione S-transferase
- [+] KOG00407 40S ribosomal protein S14
- [+] KOG00408 Mitochondrial/chloroplast ribosomal protein S11
- [+] KOG00409 Predicted dehydrogenase
- [+] KOG00410 Predicted GTP binding protein
- [+] KOG00411 Uncharacterized membrane protein
- [+] KOG00412 Golgi transport complex COD1 protein
- [+] KOG00413 Uncharacterized conserved protein related to condensin complex

Apime\_11855  
 Aspni\_03797.519.1518  
 CE08373  
 Caebr\_08775  
 Chlre\_10737  
 Gibze\_09654  
 Leima\_03326  
 Plabe\_04226  
 Plabe\_09243  
 Plabe\_11680  
 Sacba\_06308  
 Takru\_19036  
 Tetni\_05919  
 Tetni\_15833  
 Thaps\_08507  
 Trycr\_17705

Get PF

Name	COG	Sup	Score	E-Value	pf	Ing
Klula_02672	KOG00433		1.284E3	0.0	<input type="checkbox"/>	990
Cangl_04222	KOG00433		1.255E3	0.0	<input type="checkbox"/>	982
Aspni_03797	Div		1.189E3	0.0	<input type="checkbox"/>	1518
aebr_08776	KOG00433		1.167E3	0.0	<input type="checkbox"/>	947
spfu_09358			1.156E3	0.0	<input type="checkbox"/>	1009
ibze_11395	KOG00433		1.148E3	0.0	<input type="checkbox"/>	997
eucl_07354	KOG00433		1.146E3	0.0	<input type="checkbox"/>	957
luwa_00429	KOG00433		1.146E3	0.0	<input type="checkbox"/>	985
rego_00480	KOG00433		1.132E3	0.0	<input type="checkbox"/>	988
icdi_12992	KOG00433		1.126E3	0.0	<input type="checkbox"/>	1034
alga_14607	KOG00433		1.109E3	0.0	<input type="checkbox"/>	773
erli_04934	KOG00433		1.109E3	0.0	<input type="checkbox"/>	972
alga_14610	KOG00433		1.076E3	0.0	<input type="checkbox"/>	767
ryne_04497	KOG00433		1.035E3	0.0	<input type="checkbox"/>	1032
Debha_04672	KOG00433		1.016E3	0.0	<input type="checkbox"/>	988
Tetni_18015	KOG00433		974.6	0.0	<input type="checkbox"/>	933
Galga_14608	KOG00433		963.8	0.0	<input type="checkbox"/>	723
Galga_14609	KOG00433		952.4	0.0	<input type="checkbox"/>	728
Canal_10752	KOG00433		937.8	0.0	<input type="checkbox"/>	973
Canal_10803	KOG00433		934.3	0.0	<input type="checkbox"/>	973
Ratno_18798			922.8	0.0	<input type="checkbox"/>	1136
Anoga_05900	KOG00433		861.0	0.0	<input type="checkbox"/>	806
Ustma_01543	KOG00433		595.4	0.0	<input type="checkbox"/>	851
Pantr_01697			552.5	0.0	<input type="checkbox"/>	620
Galga_14611	KOG00433		308.2	0.0	<input type="checkbox"/>	470
Takru_16952	KOG00433		274.7	0.0	<input type="checkbox"/>	485
Thapa_04012			107.8	3.1E-27	<input type="checkbox"/>	1221
Trycr_07252	KOG00434		68.8	1.7E-15	<input type="checkbox"/>	1156
Enthi_08456	KOG00434		57.6	4.0E-12	<input type="checkbox"/>	1056
Cyame_04722	KOG00434		49.0	3.4E-11	<input type="checkbox"/>	1207
Canal_11674	KOG00434		48.5	3.4E-11	<input type="checkbox"/>	1088
Canal_11824	KOG00434		48.5	3.4E-11	<input type="checkbox"/>	1088
Anoga_00791	KOG00434		46.4	3.8E-11	<input type="checkbox"/>	1213
Orysa_57566	KOG00434		38.7	5.8E-11	<input type="checkbox"/>	1386

Property | Proteins | Prot-Cogs | HMM Global | HMM Local

Setting the score cutoff triggers creation of a list of potential false positives which have to be checked by a human curator

- [+] KOG00425 Ubiquitin-protein ligase
- [+] KOG00426 Ubiquitin-protein ligase
- [+] KOG00427 Ubiquitin conjugating enzyme
- [+] KOG00428 Non-canonical ubiquitin conjugating enzyme 1
- [+] KOG00429 Ubiquitin-conjugating enzyme-related protein Ft1; involved in pro
- [+] KOG00430 Xanthine dehydrogenase
- [+] KOG00431 Auxilin-like protein and related proteins containing DnaJ domain
- [+] KOG00432 needs\_to\_be\_divided\_Valyl-tRNA synthetase
- [+] KOG00433 Isoleucyl-tRNA synthetase
- [+] KOG00434 Isoleucyl-tRNA synthetase
- [+] KOG00435 Leucyl-tRNA synthetase
- [+] KOG00436 Methionyl-tRNA synthetase
- [+] KOG00437 Leucyl-tRNA synthetase
- [+] KOG00438 Mitochondrial ribosomal protein L2
- [+] KOG00439 VAMP-associated protein involved in inositol metabolism
- [+] KOG00440 Cell cycle-associated protein Mob1-1
- [+] KOG00441 Cu2+/Zn2+ superoxide dismutase SOD1
- [+] KOG00442 Structure-specific endonuclease ERCC1-XPF; catalytic component
- [+] KOG00443 Actin regulatory proteins (gelsolin/villin family)
- [+] KOG00444 Cytoskeletal regulator Flightless-I (contains leucine-rich and gels
- [+] KOG00445 Actin regulatory protein supervillin (gelsolin/villin family)

- [-] KOG00391 SWI 2 family DNA-dependent ATPase
- [+] KOG00392 SNF2 family DNA-dependent ATPase domain-containing protein
- [+] KOG00393 Ras-related small GTPase; Rho type
- [+] KOG00394 Ras-related GTPase
- [+] KOG00395 Ras-related GTPase
- [+] KOG00396 Uncharacterized conserved protein
- [+] KOG00397 60S ribosomal protein L11
- [+] KOG00398 Mitochondrial/chloroplast ribosomal protein L5/L7
- [+] KOG00399 Glutamate synthase
- [+] KOG00400 40S ribosomal protein S13
- [+] KOG00401 Translation initiation factor 4F; ribosome/mRNA-bridging subunit (
- [+] KOG00402 60S ribosomal protein L37
- [+] KOG00403 Neoplastic transformation suppressor Pdc4/MA-3; contains MA3
- [+] KOG00404 Thioredoxin reductase
- [+] KOG00405 Pyridine nucleotide-disulphide oxidoreductase
- [+] KOG00406 Glutathione S-transferase
- [+] KOG00407 40S ribosomal protein S14
- [+] KOG00408 Mitochondrial/chloroplast ribosomal protein S11
- [+] KOG00409 Predicted dehydrogenase
- [+] KOG00410 Predicted GTP binding protein
- [+] KOG00411 Uncharacterized membrane protein
- [+] KOG00412 Golgi transport complex COD1 protein
- [+] KOG00413 Uncharacterized conserved protein related to condensin complex
- [+] KOG00414 Chromosome condensation complex Condensin; subunit D2
- [+] KOG00415 Predicted peptidyl leucyl-lysine transaminase

Threshold is not yet defined

Get PF

Name	COG	Sup	Score	E-Value	pf	lng
Gibze_11395	KOG00433		1.238E3	0.0	<input type="checkbox"/>	997
Galga_14609	KOG00433		1.237E3	0.0	<input type="checkbox"/>	728
Neucr_07354	KOG00433		1.236E3	0.0	<input type="checkbox"/>	957
Kluwa_00429	KOG00433		1.231E3	0.0	<input type="checkbox"/>	985
Aspni_03797	DIV		1.224E3	0.0	<input type="checkbox"/>	1518
Erego_00480	KOG00433		1.217E3	0.0	<input type="checkbox"/>	988
Dicdi_12992	KOG00433		1.18E3	0.0	<input type="checkbox"/>	1034
Yarji_04934	KOG00433		1.163E3	0.0	<input type="checkbox"/>	972
Tetri_18015	KOG00433		1.144E3	0.0	<input type="checkbox"/>	933
Ratno_18798			1.136E3	0.0	<input type="checkbox"/>	1136
Debha_04672	KOG00433		1.093E3	0.0	<input type="checkbox"/>	988
Cryne_04497	KOG00433		1.075E3	0.0	<input type="checkbox"/>	1032
Anoga_05900	KOG00433		1.069E3	0.0	<input type="checkbox"/>	806
752	KOG00433		1.006E3	0.0	<input type="checkbox"/>	973
803	KOG00433		1.004E3	0.0	<input type="checkbox"/>	973
97			961.0	0.0	<input type="checkbox"/>	620
543	KOG00433		843.8	0.0	<input type="checkbox"/>	851
511	KOG00433		821.8	0.0	<input type="checkbox"/>	470
752	KOG00433		779.9	0.0	<input type="checkbox"/>	485
Plabe_04226	KOG00433		482.4	0.0	<input type="checkbox"/>	687
Pantr_01696			413.5	0.0	<input type="checkbox"/>	249
Playo_05169			395.9	0.0	<input type="checkbox"/>	647
Chhre_10737	KOG00433		368.8	0.0	<input type="checkbox"/>	1153
Apime_11855	KOG00433		302.6	0.0	<input type="checkbox"/>	230
Thepa_04012			285.0	0.0	<input type="checkbox"/>	1221
Playo_02836			278.5	0.0	<input type="checkbox"/>	1197
Plabe_09243	KOG00433		222.0	0.0	<input type="checkbox"/>	267
Enthi_08456	KOG00434		218.4	0.0	<input type="checkbox"/>	1056
7296537	KOG00434		207.6	0.0	<input type="checkbox"/>	1081
nDrome_08034			207.6	0.0	<input type="checkbox"/>	1229
nDrome_09237			207.6	0.0	<input type="checkbox"/>	1229
nDrome_09754			207.6	0.0	<input type="checkbox"/>	1229
Canal_11674	KOG00434		206.3	0.0	<input type="checkbox"/>	1088
Canal_11824	KOG00434		206.3	0.0	<input type="checkbox"/>	1088
Anoga_00791	KOG00434		205.4	0.0	<input type="checkbox"/>	1213
Cangl_01703	KOG00434		204.3	0.0	<input type="checkbox"/>	1072
Galga_05083	KOG00434		200.2	0.0	<input type="checkbox"/>	1262
Galga_05084	KOG00434		196.0	0.0	<input type="checkbox"/>	992
Cloin_10894	KOG00434		194.4	0.0	<input type="checkbox"/>	1175
Erego_02472	KOG00434		192.6	0.0	<input type="checkbox"/>	1072
Danre_05440	KOG00434		192.6	0.0	<input type="checkbox"/>	1271
Danre_14019	KOG00434		192.6	0.0	<input type="checkbox"/>	1271
Sacmi_02885	KOG00434		191.8	0.0	<input type="checkbox"/>	1072

And now do the same with the local hmm search result...

- [+] KOG00421 Ubiquitin-protein ligase
- [+] KOG00422 Ubiquitin-protein ligase
- [+] KOG00423 Ubiquitin-protein ligase
- [+] KOG00424 Ubiquitin-protein ligase
- [+] KOG00425 Ubiquitin-protein ligase
- [+] KOG00426 Ubiquitin-protein ligase
- [+] KOG00427 Ubiquitin conjugating enzyme
- [+] KOG00428 Non-canonical ubiquitin conjugating enzyme 1
- [+] KOG00429 Ubiquitin-conjugating enzyme-related protein Ft1; involved in pro
- [+] KOG00430 Xanthine dehydrogenase
- [+] KOG00431 Auxilin-like protein and related proteins containing DnaJ domain
- [+] KOG00432 needs\_to\_be\_divided\_Valyl-tRNA synthetase
- [+] **KOG00433 Isoleucyl-tRNA synthetase**
- [+] KOG00434 Isoleucyl-tRNA synthetase
- [+] KOG00435 Leucyl-tRNA synthetase
- [+] KOG00436 Methionyl-tRNA synthetase
- [+] KOG00437 Leucyl-tRNA synthetase
- [+] KOG00438 Mitochondrial ribosomal protein L2
- [+] KOG00439 VAMP-associated protein involved in inositol metabolism
- [+] KOG00440 Cell cycle-associated protein Mob1-1
- [+] KOG00441 Cu2+/Zn2+ superoxide dismutase SOD1
- [+] KOG00442 Structure-specific endonuclease ERCC1-XPF; catalytic component
- [+] KOG00443 Actin regulatory proteins (gelsolin/villin family)
- [+] KOG00444 Cytoskeletal regulator Flightless-1 (contains leucine-rich and gels
- [+] KOG00445 Actin regulatory protein supervillin (gelsolin/villin family)



- KOG00391 Jwi 2 family DNA-dependent ATPase
- KOG00392 SNF2 family DNA-dependent ATPase domain-containing protein
- KOG00393 Ras-related small GTPase; Rho type
- KOG00394 Ras-related GTPase
- KOG00395 Ras-related GTPase
- KOG00396 Uncharacterized conserved protein
- KOG00397 60S ribosomal protein L11
- KOG00398 Mitochondrial/chloroplast ribosomal protein L5/L7
- KOG00399 Glutamate synthase
- KOG00400 40S ribosomal protein S13
- KOG00401 Translation initiation factor 4F; ribosome/mRNA-bridging subunit (
- KOG00402 60S ribosomal protein L37
- KOG00403 Neoplastic transformation suppressor Pcdcd4/MA-3; contains MA3
- KOG00404 Thioredoxin reductase
- KOG00405 Pyridine nucleotide-disulphide oxidoreductase
- KOG00406 Glutathione S-transferase
- KOG00407 40S ribosomal protein S14
- KOG00408 Mitochondrial/chloroplast ribosomal protein S11
- KOG00409 Predicted dehydrogenase
- KOG00410 Predicted GTP binding protein
- KOG00411 Uncharacterized membrane protein
- KOG00412 Golgi transport complex COD1 protein
- KOG00413 Uncharacterized conserved protein related to condensin complex
- KOG00414 Chromosome condensation complex Condensin; subunit D2
- KOG00415 Predicted peptidyl prolyl cis-trans isomerase
- KOG00416 Ubiquitin-protein ligase
- KOG00417 Ubiquitin-protein ligase
- KOG00418 Ubiquitin-protein ligase
- KOG00419 Ubiquitin-protein ligase
- KOG00420 Ubiquitin-protein ligase
- KOG00421 Ubiquitin-protein ligase
- KOG00422 Ubiquitin-protein ligase
- KOG00423 Ubiquitin-protein ligase
- KOG00424 Ubiquitin-protein ligase
- KOG00425 Ubiquitin-protein ligase
- KOG00426 Ubiquitin-protein ligase
- KOG00427 Ubiquitin conjugating enzyme
- KOG00428 Non-canonical ubiquitin conjugating enzyme 1
- KOG00429 Ubiquitin-conjugating enzyme-related protein Ft1; involved in pro
- KOG00430 Xanthine dehydrogenase
- KOG00431 Auxilin-like protein and related proteins containing DnaJ domain
- KOG00432 needs\_to\_be\_divided\_valyl-tRNA synthetase
- KOG00433 Isoleucyl-tRNA synthetase**
- KOG00434 Isoleucyl-tRNA synthetase
- KOG00435 Leucyl-tRNA synthetase
- KOG00436 Methionyl-tRNA synthetase
- KOG00437 Leucyl-tRNA synthetase
- KOG00438 Mitochondrial ribosomal protein L2
- KOG00439 VAMP-associated protein involved in inositol metabolism
- KOG00440 Cell cycle-associated protein Mob1-1
- KOG00441 Cu2+/Zn2+ superoxide dismutase SOD1
- KOG00442 Structure-specific endonuclease ERCC1-XPF; catalytic component
- KOG00443 Actin regulatory proteins (gelsolin/villin family)
- KOG00444 Cytoskeletal regulator Flightless-1 (contains leucine-rich and gels
- KOG00445 Actin regulatory protein supervillin (gelsolin/villin family)

Aspni\_03797.519.1518  
 CE08373  
 Caebr\_08775  
 Gibze\_09654  
 Leima\_03326  
 Plabe\_11680  
 Sacba\_06308  
 Takru\_19036  
 Tetni\_05919  
 Tetni\_15833  
 Thaps\_08507  
 Trycr\_17705



Get PF

Name	COG	Sup	Score	E-Value	pf	lng
Gibze_11395	KOG00433		1.238E3	0.0	<input type="checkbox"/>	997
Galga_14609	KOG00433		1.237E3	0.0	<input type="checkbox"/>	728
Neucr_07354	KOG00433		1.236E3	0.0	<input type="checkbox"/>	957
Kluwa_00429	KOG00433		1.231E3	0.0	<input type="checkbox"/>	985
Aspni_03797	DIV		1.224E3	0.0	<input type="checkbox"/>	1518
Erego_00480	KOG00433		1.217E3	0.0	<input type="checkbox"/>	988
Dicdi_12992	KOG00433		1.18E3	0.0	<input type="checkbox"/>	1034
Yarli_04934	KOG00433		1.163E3	0.0	<input type="checkbox"/>	972
Tetni_18015	KOG00433		1.144E3	0.0	<input type="checkbox"/>	933
Ratno_18798			1.136E3	0.0	<input type="checkbox"/>	1136
Debha_04672	KOG00433		1.093E3	0.0	<input type="checkbox"/>	988
Cryne_04497	KOG00433		1.075E3	0.0	<input type="checkbox"/>	1032
Anoga_05900	KOG00433		1.069E3	0.0	<input type="checkbox"/>	806
Canal_10752	KOG00433		1.006E3	0.0	<input type="checkbox"/>	973
Canal_10803	KOG00433		1.004E3	0.0	<input type="checkbox"/>	973
Pantr_01697			961.0	0.0	<input type="checkbox"/>	620
Ustma_01543	KOG00433		843.8	0.0	<input type="checkbox"/>	851
Galga_14611	KOG00433		821.8	0.0	<input type="checkbox"/>	470
Takru_16952	KOG00433		779.9	0.0	<input type="checkbox"/>	485
Plabe_04226	KOG00433		482.4	0.0	<input type="checkbox"/>	687
Pantr_01696			413.5	0.0	<input type="checkbox"/>	249
Playo_05169			395.9	0.0	<input type="checkbox"/>	647
Chlre_10737	KOG00433		368.8	0.0	<input type="checkbox"/>	1153
Apime_11855	KOG00433		302.6	0.0	<input type="checkbox"/>	230
Thepa_04012			285.0	0.0	<input type="checkbox"/>	1221
Playo_02836			278.5	0.0	<input type="checkbox"/>	1197
Plabe_09243	KOG00433		222.0	0.0	<input type="checkbox"/>	267
Enthi_08456	KOG00434		218.4	0.0	<input type="checkbox"/>	1056
7296537	KOG00434		207.6	0.0	<input type="checkbox"/>	1081
nDrome_08034			207.6	0.0	<input type="checkbox"/>	1229
nDrome_09237			207.6	0.0	<input type="checkbox"/>	1229
nDrome_09754			207.6	0.0	<input type="checkbox"/>	1229
Canal_11674	KOG00434		206.3	0.0	<input type="checkbox"/>	1088
Canal_11824	KOG00434		206.3	0.0	<input type="checkbox"/>	1088
Anoga_00791	KOG00434		205.4	0.0	<input type="checkbox"/>	1213
Cangl_01703	KOG00434		204.3	0.0	<input type="checkbox"/>	1072
Galga_05083	KOG00434		200.2	0.0	<input type="checkbox"/>	1262
Galga_05084	KOG00434		196.0	0.0	<input type="checkbox"/>	992

Property Proteins Prot-Cogs HMM Global HMM Local

- [-] KOG00391 SWI 2 family DNA-dependent ATPase
- [-] KOG00392 SNF2 family DNA-dependent ATPase domain-containing protein
- [-] KOG00393 Ras-related small GTPase; Rho type
- [-] KOG00394 Ras-related GTPase
- [-] KOG00395 Ras-related GTPase
- [-] KOG00396 Uncharacterized conserved protein
- [-] KOG00397 60S ribosomal protein L11
- [-] KOG00398 Mitochondrial/chloroplast ribosomal protein L5/L7
- [-] KOG00399 Glutamate synthase
- [-] KOG00400 40S ribosomal protein S13
- [-] KOG00401 Translation initiation factor 4F; ribosome/mRNA-bridging subunit (
- [-] KOG00402 60S ribosomal protein L37
- [-] KOG00403 Neoplastic transformation suppressor Pcd4/MA-3; contains MA3
- [-] KOG00404 Thioredoxin reductase
- [-] KOG00405 Pyridine nucleotide-disulphide oxidoreductase
- [-] KOG00406 Glutathione S-transferase
- [-] KOG00407 40S ribosomal protein S14
- [-] KOG00408 Mitochondrial/chloroplast ribosomal protein S11
- [-] KOG00409 Predicted dehydrogenase
- [-] KOG00410 Predicted GTP binding protein
- [-] KOG00411 Uncharacterized membrane protein
- [-] KOG00412 Golgi transport complex COD1 protein

Aspni\_03797.519.1518  
 CE08373  
 Caebr\_08775  
 Gibze\_09654  
 Leima\_03326  
 Plabe\_11680  
 Sacba\_06308  
 Takru\_19036  
 Tetni\_05919  
 Tetni\_15833  
 Thaps\_08507  
 Trycr\_17705

Get PF

Name	COG	Sup	Score	E-Value	pf	lng
Gibze_11395	KOG00433		1.238E3	0.0	<input type="checkbox"/>	997
Galga_14609	KOG00433		1.237E3	0.0	<input type="checkbox"/>	728
Neucr_07354	KOG00433		1.236E3	0.0	<input type="checkbox"/>	957
Kluwa_00429	KOG00433		1.231E3	0.0	<input type="checkbox"/>	985
Aspni_03797	DIV		1.224E3	0.0	<input type="checkbox"/>	1518
Frago_00480	KOG00433		1.217E3	0.0	<input type="checkbox"/>	988
2992	KOG00433		1.18E3	0.0	<input type="checkbox"/>	1034
934	KOG00433		1.163E3	0.0	<input type="checkbox"/>	972
8015	KOG00433		1.144E3	0.0	<input type="checkbox"/>	933
18798			1.136E3	0.0	<input type="checkbox"/>	1136
04672	KOG00433		1.093E3	0.0	<input type="checkbox"/>	988
04497	KOG00433		1.075E3	0.0	<input type="checkbox"/>	1032
05900	KOG00433		1.069E3	0.0	<input type="checkbox"/>	806
10752	KOG00433		1.006E3	0.0	<input type="checkbox"/>	973
10803	KOG00433		1.004E3	0.0	<input type="checkbox"/>	973
01697			961.0	0.0	<input type="checkbox"/>	620
01543	KOG00433		843.8	0.0	<input type="checkbox"/>	851
Galga_14611	KOG00433		821.8	0.0	<input type="checkbox"/>	470
Takru_16952	KOG00433		779.9	0.0	<input type="checkbox"/>	485
Plabe_04226	KOG00433		482.4	0.0	<input checked="" type="checkbox"/>	687
Pantr_01696			413.5	0.0	<input checked="" type="checkbox"/>	249
Playo_05169			395.9	0.0	<input checked="" type="checkbox"/>	647
Chlre_10737	KOG00433		368.8	0.0	<input checked="" type="checkbox"/>	1153
Apime_11855	KOG00433		302.6	0.0	<input checked="" type="checkbox"/>	230
Thepa_04012			285.0	0.0	<input type="checkbox"/>	1221
Playo_02836			278.5	0.0	<input checked="" type="checkbox"/>	1197
Plabe_09243	KOG00433		222.0	0.0	<input checked="" type="checkbox"/>	267
Enthi_08456	KOG00434		218.4	0.0	<input type="checkbox"/>	1056
7296537	KOG00434		207.6	0.0	<input type="checkbox"/>	1081
nDrome_08034			207.6	0.0	<input type="checkbox"/>	1229
nDrome_09237			207.6	0.0	<input type="checkbox"/>	1229
nDrome_09754			207.6	0.0	<input type="checkbox"/>	1229
Canal_11674	KOG00434		206.3	0.0	<input type="checkbox"/>	1088
Canal_11824	KOG00434		206.3	0.0	<input type="checkbox"/>	1088
Anoga_00791	KOG00434		205.4	0.0	<input type="checkbox"/>	1213
Cangl_01703	KOG00434		204.3	0.0	<input type="checkbox"/>	1072
Galga_05083	KOG00434		200.2	0.0	<input type="checkbox"/>	1262
Galga_05084	KOG00434		196.0	0.0	<input type="checkbox"/>	992

Setting global and local hmm search cutoffs allows to get an automatic list of potential fragments

- [-] KOG00424 Ubiquitin-protein ligase
- [-] KOG00425 Ubiquitin-protein ligase
- [-] KOG00426 Ubiquitin-protein ligase
- [-] KOG00427 Ubiquitin conjugating enzyme
- [-] KOG00428 Non-canonical ubiquitin conjugating enzyme 1
- [-] KOG00429 Ubiquitin-conjugating enzyme-related protein Ft1; involved in pro
- [-] KOG00430 Xanthine dehydrogenase
- [-] KOG00431 Auxilin-like protein and related proteins containing DnaJ domain
- [-] KOG00432 needs\_to\_be\_divided\_Valyl-tRNA synthetase
- [-] KOG00433 Isoleucyl-tRNA synthetase
- [-] KOG00434 Isoleucyl-tRNA synthetase
- [-] KOG00435 Leucyl-tRNA synthetase
- [-] KOG00436 Methionyl-tRNA synthetase
- [-] KOG00437 Leucyl-tRNA synthetase
- [-] KOG00438 Mitochondrial ribosomal protein L2
- [-] KOG00439 VAMP-associated protein involved in inositol metabolism
- [-] KOG00440 Cell cycle-associated protein Mob1-1
- [-] KOG00441 Cu2+/Zn2+ superoxide dismutase SOD1
- [-] KOG00442 Structure-specific endonuclease ERCC1-XPB; catalytic component
- [-] KOG00443 Actin regulatory proteins (gelsolin/villin family)
- [-] KOG00444 Cytoskeletal regulator Flightless-1 (contains leucine-rich and gels
- [-] KOG00445 Actin regulatory protein supervillin (gelsolin/villin family)

Property | Proteins | Prot-Cogs | HMM Global | HMM Local

# What trouble did KOGnitor make to KOG00433?

We take the list of potential false positives and potential false negatives and compute BLASTP with output in xml format. The resulting files are loaded into KOG db temporarily to check all the sequences in question.

COG Editor - Main\_KOG\_Project

File Cog Protein BlastPage Help

Main | GOG Sets | GOG Tree

- KOG00391 Smr 2 Family DNA-dependent ATPase
- KOG00392 SNF2 Family DNA-dependent ATPase domain-containing protein
- KOG00393 Ras-related small GTPase; Rho type
- KOG00394 Ras-related GTPase
- KOG00395 Ras-related GTPase
- KOG00396 Uncharacterized conserved protein
- KOG00397 60S ribosomal protein L11
- KOG00398 Mitochondrial/chloroplast ribosomal protein L5/L7
- KOG00399 Glutamate synthase
- KOG00400 40S ribosomal protein S13
- KOG00401 Translation initiation factor 4F; ribosome/mRNA-bridging subunit (
- KOG00402 9935 COG0060
- KOG00403 Isoleucyl-tRNA synthetase [Translation, ribosomal structure and biogenesis]
- KOG00404 from:15 to:995 ev:0.0
- KOG00405 70S ribosomal protein L23
- KOG00406 Glutathione S-transferase
- KOG00407 40S ribosomal protein S14
- KOG00408 Mitochondrial/chloroplast ribosomal protein S11
- KOG00409 Predicted dehydrogenase

9935 COG0060

Isoleucyl-tRNA synthetase [Translation, ribosomal structure and biogenesis]

from:15 to:995 ev:0.0

Alignment

dBeg	-- Alignment --	dEnd	Score	Size	Prot	Org	COG	±	SC
0		0	5347	1009	Aspfu_09358	Aspfu		<input checked="" type="checkbox"/>	
0		1	4097	998	Aspni_03797.51...	Aspni	KOG00433	<input type="checkbox"/>	
0		57	2419	997	Gibze_11395	Gibze	KOG00433	<input type="checkbox"/>	
0		5	2377	957	Neucr_07354	Neucr	KOG00433	<input type="checkbox"/>	
0		4	2316	974	Maggr_10303	Maggr	KOG00433	<input type="checkbox"/>	
0		1	1925	972	Yarli_04934	Yarli	KOG00433	<input type="checkbox"/>	
0		1	1769	1032	Cryne_04497	Cryne	KOG00433	<input type="checkbox"/>	
101		16	1754	1093	At5g49030	Arath	KOG00433	<input type="checkbox"/>	
16		16	1754	1008	#At5g49030	oldArath	KOG00433	<input type="checkbox"/>	
12		0	1713	990	Klula_02672	Klula	KOG00433	<input type="checkbox"/>	
109		12	1708	1079	Cyame_00546	Cyame	KOG00433	<input type="checkbox"/>	
7		0	1706	980	Thaps_05451	Thaps	KOG00433	<input type="checkbox"/>	
74		12	1691	1057	Orysa_25681	Orysa	KOG00433	<input type="checkbox"/>	
31		2	1690	973	SPCC1885.08c	Schpo	KOG00433	<input type="checkbox"/>	
12		3	1675	982	Cangl_04222	Cangl	KOG00433	<input type="checkbox"/>	
20		0	1674	973	Canal_10752	Canal	KOG00433	<input type="checkbox"/>	
20		0	1670	973	Canal_10803	Canal	KOG00433	<input type="checkbox"/>	
30		5	1641	982	Danre_02073	Danre	KOG00433	<input type="checkbox"/>	
3		1	1630	985	Kluwa_00429	Kluwa	KOG00433	<input type="checkbox"/>	
41		1	1629	1034	Dicdi_12992	Dicdi	KOG00433	<input type="checkbox"/>	
34		3	1605	986	Galga_14606	Galga	KOG00433	<input type="checkbox"/>	
17		2	1603	996	Sacba_07540	Sacba	KOG00433	<input type="checkbox"/>	
18		3	1599	988	Debha_04672	Debha	KOG00433	<input type="checkbox"/>	
59		3	1597	1012	Musmu_13328	Musmu	KOG00433	<input type="checkbox"/>	
5		3	1595	923	PantrE_02879	PantrE	KOG00433	<input type="checkbox"/>	
5		3	1595	940	Homsa_13929	Homsa	KOG00433	<input type="checkbox"/>	
23		2	1594	1002	Sacpa_06583	Sacpa	KOG00433	<input type="checkbox"/>	
11		0	1593	988	Erego_00480	Erego	KOG00433	<input type="checkbox"/>	
23		2	1572	1002	YPL040c	Sacce	KOG00433	<input type="checkbox"/>	
23		42	1533	895	Sacmi_05887	Sacmi	KOG00433	<input type="checkbox"/>	
2		25	1505	773	Galga_14607	Galga	KOG00433	<input type="checkbox"/>	
0		0	1481	767	Galga_14610	Galga	KOG00433	<input type="checkbox"/>	
n		51	1441	851	Iistma_01543	Iistma	KOG00433	<input type="checkbox"/>	

Aspfu\_09358 was predicted by HMM search to be false negative. These BLASTP and RPS BLAST results demonstrate that it has been indeed missed and needs to be added to KOG00433

- KOG00429 Ubiquitin-conjugating enzyme-related protein Ft1; involved in pro
- KOG00430 Xanthine dehydrogenase
- KOG00431 Auxilin-like protein and related proteins containing DnaJ domain
- KOG00432 needs\_to\_be\_divided\_Valyl-tRNA synthetase
- KOG00433 Isoleucyl-tRNA synthetase
- KOG00434 Isoleucyl-tRNA synthetase
- KOG00435 Leucyl-tRNA synthetase
- KOG00436 Methionyl-tRNA synthetase
- KOG00437 Leucyl-tRNA synthetase
- KOG00438 Mitochondrial ribosomal protein L2
- KOG00439 VAMP-associated protein involved in inositol metabolism
- KOG00440 Cell cycle-associated protein Mob1-1
- KOG00441 Cu2+/Zn2+ superoxide dismutase SOD1
- KOG00442 Structure-specific endonuclease ERCC1-XPF; catalytic component
- KOG00443 Actin regulatory proteins (gelsolin/villin family)
- KOG00444 Cytoskeletal regulator Flightless-I (contains leucine-rich and gels
- KOG00445 Actin regulatory protein supervillin (gelsolin/villin family)

Alignment

Start	Alignment	End
1	MPELPRSLTRSWSTLKLPKSTFFPARVTPADQTKYLQRCDELYAQRRRPPADRPFVLH	60
1	MPELPRSLTRSWSTLKLPKSTFFPARVTPADQTKYLQRCDELYAQRRRPPADRPFVLH	60
61	DGPPYANGLHI GHALNKILKDIICRVQLAKGKRVRYVPGWDCHLPIELKALDLQKELG	120
61	DGPPYANGLHI GHALNKILKDIICRVQLAKGKRVRYVPGWDCHLPIELKALDLQKELG	120
121	NTGGSIGAAAIRRAARKLAGRTVKEQMKGFRSFGVMADWDGHWKTMDKEFEKRLGVFRE	180
121	NTGGSIGAAAIRRAARKLAGRTVKEQMKGFRSFGVMADWDGHWKTMDKEFEKRLGVFRE	180

Main | GOG Sets | GOG Tree

- [-] KOG00391 SWI 2 family DNA-dependent ATPase
- [-] KOG00392 SNF2 family DNA-dependent ATPase domain-containing protein
- [-] KOG00393 Ras-related small GTPase; Rho type
- [-] KOG00394 Ras-related GTPase
- [-] KOG00395 Ras-related GTPase
- [-] KOG00396 Uncharacterized conserved protein
- [-] KOG00397 60S ribosomal protein L11
- [-] KOG00398 Mitochondrial/chloroplast ribosomal protein L5/L7
- [-] KOG00399 Glutamate synthase
- [-] KOG00400 40S ribosomal protein S13
- [-] KOG00401 Translation initiation factor 4F; ribosome/mRNA-bridging subunit (
- [-] KOG00402 60S ribosomal protein L37
- [-] KOG00403 Neoplastic transformation suppressor Pcd4/MA-3; contains MA3
- [-] KOG00404 Thioredoxin reductase
- [-] KOG00405 Pyridine nucleotide-disulphide oxidoreductase
- [-] KOG00406 Glutathione S-transferase
- [-] KOG00407 40S ribosomal protein S14
- [-] KOG00408 Mitochondrial/chloroplast ribosomal protein S11
- [-] KOG00409 Predicted dehydrogenase

dBeg	-- Alignment --	dEnd	Score	Size	Prot	Org	COG	±	SC
0		0	618	122	Caebr_08775	Caebr	KOG00433	<input type="checkbox"/>	
0		1	497	122	nCaeel_05820	nCaeel	KOG00433	<input type="checkbox"/>	
0		1	497	122	CE08373	Caeel	KOG00433	<input type="checkbox"/>	
0		1059	304	1131	nCaeel_06559	nCaeel	KOG00433	<input type="checkbox"/>	
0		1059	304	1131	CE08372	Caeel	KOG00433	<input type="checkbox"/>	
223		137	76	456	At1g66300	Arath	KOG05893	<input type="checkbox"/>	
223		137	76	456	#At1g66300	oldArath	KOG05893	<input type="checkbox"/>	
468		305	73	859	Aspni_00278	Aspni	KOG15858	<input type="checkbox"/>	
502		436	71	982	Dicdi_06992	Dicdi	KOG03751	<input type="checkbox"/>	
25		437	70	523	Enthi_02508	Enthi	KOG00100	<input type="checkbox"/>	
413		120	68	589	Takru_24461	Takru	KOG03751	<input type="checkbox"/>	
405		120	68	581	Takru_24460	Takru	KOG03751	<input type="checkbox"/>	
		183	68	334	At5g51410	Arath	KOG00796	<input type="checkbox"/>	
		183	68	334	#At5g51410	oldArath	KOG00796	<input type="checkbox"/>	
		120	67	497	Tetni_06659	Tetni	KOG03751	<input type="checkbox"/>	
		120	66	608	Tetni_01284	Tetni	KOG03751	<input type="checkbox"/>	

Caebr\_08775, nCaeel\_05820, and CE08373 are ca. 10 times shorter than true members of KOG00433 and are likely to be either incorrectly predicted or, rather, should form a separate KOG: their C-termini are totally different from KOG00433

- [-] KOG00431 Auxilin-like protein and related proteins containing DnaJ domain
- [-] KOG00432 needs\_to\_be\_divided\_Valyl-tRNA synthetase
- [-] KOG00433 Isoleucyl-tRNA synthetase
- [-] KOG00434 Isoleucyl-tRNA synthetase
- [-] KOG00435 Leucyl-tRNA synthetase
- [-] KOG00436 Methionyl-tRNA synthetase
- [-] KOG00437 Leucyl-tRNA synthetase
- [-] KOG00438 Mitochondrial ribosomal protein L2
- [-] KOG00439 VAMP-associated protein involved in inositol metabolism
- [-] KOG00440 Cell cycle-associated protein Mob1-1
- [-] KOG00441 Cu2+/Zn2+ superoxide dismutase SOD1
- [-] KOG00442 Structure-specific endonuclease ERCC1-XPF; catalytic component
- [-] KOG00443 Actin regulatory proteins (gelsolin/villin family)
- [-] KOG00444 Cytoskeletal regulator Flightless-I (contains leucine-rich and gelsolin)
- [-] KOG00445 Actin regulatory protein supervillin (gelsolin/villin family)

Start	Alignment	End
1	MRFTDYSFEFNQEPPIHFSVTSKDKDLVSNWTRIGQVIEVIQPELLSSQLSERQRVE	60
1	M F D+SFEFN EPIHFS+TSH DKDLV FN TGRIGQVIEVIQPE+LSSQLS++QRVE	
1	MSFKDFSFEFNQEPPIHFSITSHPKDKLVHFNSTGRIGQVIEVIQPEILLSSQLSDQQRVE	60
61	YELKNLLGNSEL	72
61	YELK+LLGN EL	
61	YELKSLGNPEL	72

Info Hits

## Action after checking potential false positives list in KOG00433:

- 10 sequences are removed, 7 of which formed three new KOGs;
- 4 sequences must be added to the HMM profile

### NEW KOG

Leima\_03326 and Trycr\_17705 were false positives, were removed and formed a new KOG which also includes "XP\_823314.7.1748518 hypothetical protein Tb10.26.1000 [Trypanosoma brucei]" which is now present in the GeneBank but was absent from the downloaded T.brucei proteome

### NEW KOG

Thaps\_08507 was added to KOG00433 by a spurious unidirectional hit, it dragged Sacba\_06308 into KOG00433; both were removed and formed a new KOG

### NEW KOG

Caebr\_08775, nCaeel\_05820, and CE08373 are ca. 10 times shorter than true members of KOG00433 and are likely to be either incorrectly predicted or, rather, should form a separate KOG: their C-termini are totally different from KOG00433

### REMOVE

Tetni\_05919, Tetni\_15833, Gibze\_09654 were false positive dragged into KOG00433 by spurious BLAST hits; they were removed from KOG00433 and were made free

### ADD TO HMM

Plabe\_04226, Playo\_05169, Plabe\_09243, and Plafa\_01680 are legitimate members of KOG00433 and MUST be added to the HMM profile

Plabe\_11680, Apime\_11855, Chlre\_10737, and Takru\_19036 are fragments and are legitimate members of KOG00433

**With precomputed XML BLAST files it takes a few minutes**



# Superfamilies and KOGs

Main | GOG Sets | GOG Tree

- [-] KOG18588 New Cog
- [-] KOG18589 New Cog
- [-] KOG18590 New Cog
- [-] KOG18591 New Cog
- [-] KOG18592 New Cog
- [-] KOG18593 New Cog
- [-] KOG18594 New Cog
- [-] KOG18595 New Cog
- [-] KOG18596 New Cog
- [-] KOG18597 New Cog
- [-] KOG18598 New Cog
- [-] KOG18599 New Cog
- [-] KOG18600 New Cog

18372 KOG0577

KOG0577, Serine/threonine protein kinase [Signal transduction mechanisms]  
from:710 to:892 ev:8.0E-5

dBeg	-- Alignment --	dEnd	Score	Size	Prot	Org	COG	±	SC
0		0	5495	1066	Giala_03726	Giala			SK0001
77		144	533	nCaeel_05061	nCaeel	KOG00197			SK0001
75		143	542	Caebr_08916	Caebr	KOG00197			SK0001
284		143	522	At1g12580	Arath	KOG00032			SK0001
284		143	522	#At1g12580	oldArath	KOG00032			SK0001
83		136	558	CE27551	Caeel	KOG00197			SK0001
40		135	259	Thaps_03255	Thaps	KOG00595			SK0001
144		134	344	Enthi_08539	Enthi	KOG00032			SK0001
70		134	270	Enthi_08538	Enthi	KOG00032			SK0001
70		134	270	Enthi_08545	Enthi	KOG00032			SK0001
876		133	1083	Tetni_18164	Tetni	KOG00595			SK0001
880		133	1090	Leima_00316	Leima	KOG06353			SK0001
179		132	406	Musmu_14365	Musmu	KOG00599			SK0001
442		132	708	Erego_03523	Erego	KOG00575			SK0001

SuperCOGs

SuperCOG: SK0001: Protein kinases

- [-] SK0001: Protein kinases
  - [-] @SK0001 Free Proteins
    - [-] Anopheles gambiae str. PEST | Anoga
    - [-] Aspergillus fumigatus AF293 | AspFu
    - [-] Candida albicans SC5314 | Canal
    - [-] Chlamydomonas reinhardtii | Chlre
    - [-] Cyanidioschyzon merolae | Cyame
    - [-] Dictyostelium discoideum | Didi
    - [-] Encephalitozoon cuniculi | Enccu
    - [-] Entamoeba histolytica | Enthi
    - [-] Giardia lamblia ATCC 50803 | Giala
      - Giala\_03726
      - Giala\_04277
      - Giala\_05297
      - Giala\_02205
    - [-] Leishmania major | Leima
    - [-] Neurospora crassa | Neucr
  - [-] KOG00032 Ca2+/calmodulin-dependent protein kinase; EF-Hand protein superfamily
  - [-] KOG00192 needs revision. Tyrosine kinase specific for activated (GTP-bound) p21cdc42
  - [-] KOG00193 Serine/threonine protein kinase RAF

Create Super COG ... Delete

Giala\_03726 is a protein kinase which cannot be placed to either KOG00197 or KOG00032 with any confidence and is left as a free member of protein kinase superfamily

- [-] KOG18628 New Cog
- [-] OGI8629 Predicted protein kinase
- [-] OGI8630 Predicted protein kinase
- [-] OGI8631 Predicted protein kinase
- [-] KOG18632 Predicted protein kinase
- [-] KOG18633 Predicted protein kinase
- [-] KOG18634 Predicted protein kinase
- [-] KOG18635 Predicted protein kinase
- [-] KOG18636 Predicted protein kinase
- [-] KOG18637 Predicted protein kinase
- [-] KOG18638 Predicted protein kinase
- [-] Leishmania major | Leima
- [-] Trypanosoma brucei | Trybr
  - Trybr\_01117
  - Trybr\_03920

Start		
732	PMTLSGKQFLR---VYRQLSHRGILSFVGCCTAINNAFYVFTESPPKIKLSRVFYQDEEA	788
	P T+S + FL+ + +Q H ++ T FY+ TE L + D	
305	PGTMSPEAFLQEASIMKQCDHPNLVKLYAVCTR--EFPFYIITEYMINGSLHLHRNDGST	363
789	LQLFRSERRIRVFIRSLQAIN-FLYSNDTKLIHRDLRPNQIUWVKNSNGDPVAKIYHMG	847
	L I+ + Q N +Y + KL+HRDL + V +G PV K+ G	
364	LG-----IQALVDMAAQIANGMMLYERKLVHRDLAARNVVLVGDKISGVPVVKVADFG	416
848	FMYPLSANEPMEHNT-----IWVAPE-VVCGGIDTKSDIWSVGTILFRLL	891
	L + E T W APE CG KSD+WS G L + ++	
417	LARKLMEEDIYEARTGAKFPKIKWTAPEAATCGNFTVKSDVWSYGYLLYEIM	467

Info Hits

Most important practical application of  
KOG db would be transfer of gene  
annotation across species

Query: gi|58385127 ENSANGP00000016957 [Anopheles gambiae str. PEST]  
 Matching gi: 55240814

Show identical Best hits  
 25 BLAST hits to 17 unique species  
 0 Archaea 0 Bacteria 24 Metazoa

Animal proteins from KOG04771 are not annotated;  
 There are no hits to fungal proteins

Keep only Cut-Off 100 Select Reset New search by GI: 58385127 Go

199 aa

	SCORE	P	ACCESSION	GI	PROTEIN DESCRIPTION
	338	20	<a href="#">NP_651884</a>	<a href="#">24651718</a>	CG11563-PA [Drosophila melanogaster]
	338	20	<a href="#">AAM29402</a>	<a href="#">21064345</a>	RE07994p [Drosophila melanogaster]
	338	20	<a href="#">EAL28137</a>	<a href="#">54638735</a>	GA11068-PA [Drosophila pseudoobscura]
	259	8	<a href="#">CAF97276</a>	<a href="#">47221358</a>	unnamed protein product [Tetraodon nigroviridis]
	241	8	<a href="#">NP_001...</a>	<a href="#">62460412</a>	hypothetical protein LOC506421 [Bos taurus]
	230	8	<a href="#">CAG33450</a>	<a href="#">48146455</a>	HSPC111 [Homo sapiens]
	229	8	<a href="#">XP_518115</a>	<a href="#">55625438</a>	PREDICTED: hypothetical protein XP_518115 [Pan troglodytes]
	229	8	<a href="#">XP_536424</a>	<a href="#">57085179</a>	PREDICTED: similar to Protein CGI-117 [Canis familiaris]
	227	8	<a href="#">NP_057475</a>	<a href="#">7705451</a>	hypothetical protein LOC51491 [Homo sapiens]
	225	8	<a href="#">BAC40235</a>	<a href="#">26353210</a>	unnamed protein product [Mus musculus]
	225	8	<a href="#">NP_848720</a>	<a href="#">30519925</a>	hypothetical protein LOC28126 [Mus musculus]
	225	8	<a href="#">XP_225176</a>	<a href="#">27682843</a>	PREDICTED: similar to DNA segment, Chr 13, Wayne State University 177, expressed [Rattus norvegicus]
	221	8	<a href="#">AAH12213</a>	<a href="#">15126561</a>	DNA segment, Chr 13, Wayne State University 177, expressed [Mus musculus]
	221	8	<a href="#">NP_001...</a>	<a href="#">50540130</a>	hypothetical protein LOC436807 [Danio rerio]
	190	8	<a href="#">AAH68216</a>	<a href="#">45872610</a>	Unknown (protein for MGC:76225) [Xenopus tropicalis]
	185	8	<a href="#">AAH73261</a>	<a href="#">49255993</a>	MGC80623 protein [Xenopus laevis]
	175	8	<a href="#">AAH32424</a>	<a href="#">34783430</a>	Unknown (protein for MGC:40298) [Homo sapiens]
	172	8	<a href="#">XP_782922</a>	<a href="#">72007783</a>	PREDICTED: similar to Protein CGI-117 [Strongylocentrotus purpuratus]
	169	7	<a href="#">CAE57069</a>	<a href="#">39578894</a>	Hypothetical protein CBG24963 [Caenorhabditis briggsae]
	168	8	<a href="#">XP_414551</a>	<a href="#">50754941</a>	PREDICTED: similar to Hypothetical protein HSPC111 [Gallus gallus]
	158	7	<a href="#">NP_492248</a>	<a href="#">17511059</a>	ZK265.6 [Caenorhabditis elegans]
	129	8	<a href="#">XP_343880</a>	<a href="#">62655980</a>	PREDICTED: similar to DNA segment, Chr 13, Wayne State University 177, expressed [Rattus norvegicus]
	109	8	<a href="#">XP_688990</a>	<a href="#">68356182</a>	PREDICTED: similar to Protein CGI-117 [Danio rerio]
	101	8	<a href="#">XP_357941</a>	<a href="#">63586355</a>	PREDICTED: similar to DNA segment, Chr 13, Wayne State University 177, expressed [Mus musculus]
	100	3	<a href="#">XP_810459</a>	<a href="#">71417059</a>	hypothetical protein, conserved [Trypanosoma cruzi]

Blink search with the human protein does not help much either

Show identical Best hits Common Tree Taxonomy Report 3D structures CDD-Search GI list

34 BLAST hits to 19 unique species [Sort by taxonomy proximity](#)

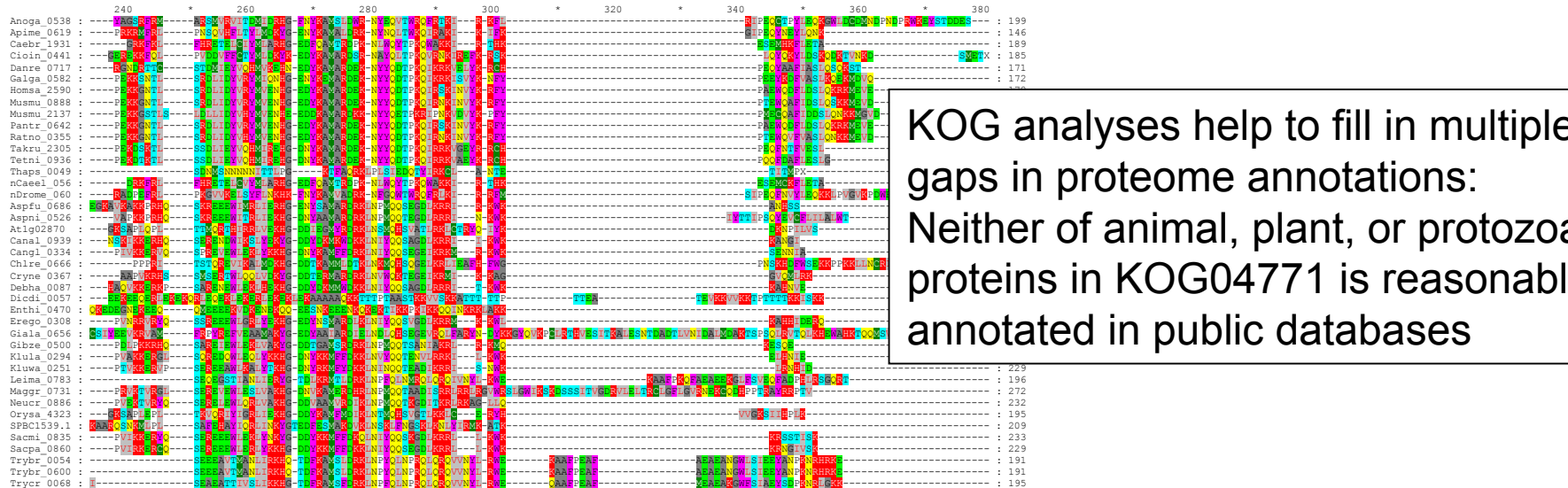
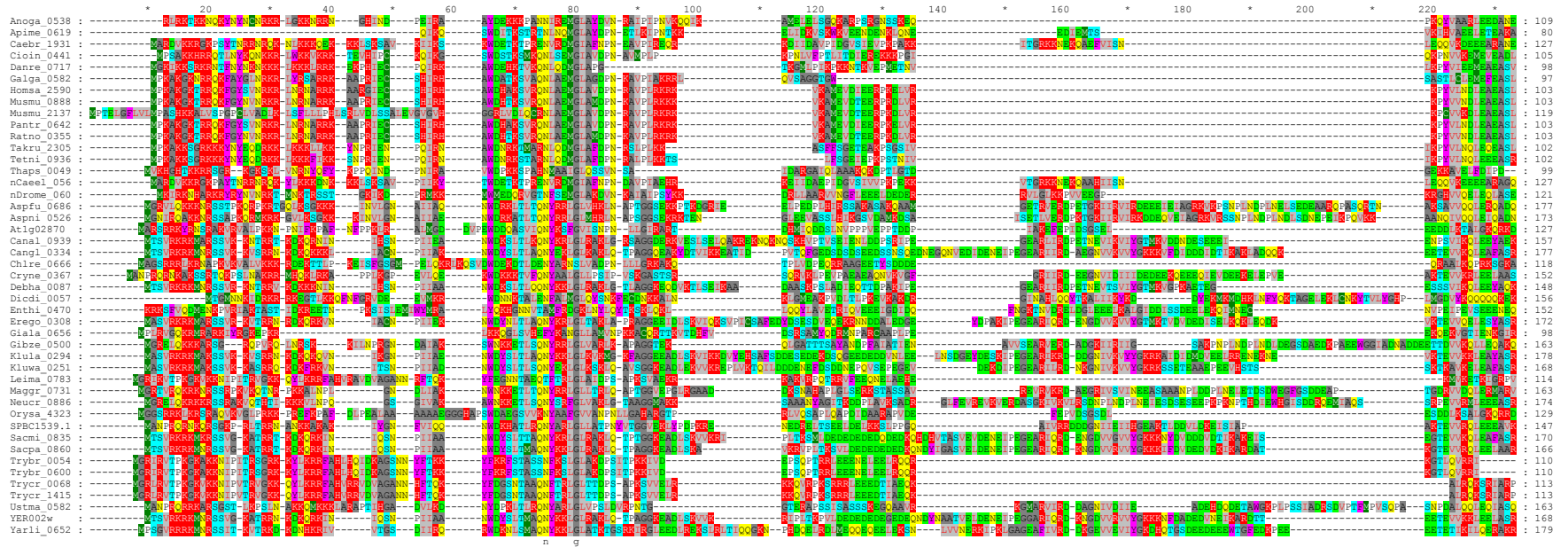
0 Archaea 0 Bacteria 32 Metazoa 0 Fungi 0 Plants 0 Viruses 2 Other Eukaryotae

Keep only  Cut-Off  Select Reset New search by GI:  Go

178 aa

	SCORE	P	ACCESSION	GI	PROTEIN DESCRIPTION
	915	28	<a href="#">XP_518115</a>	<a href="#">55625438</a>	PREDICTED: hypothetical protein XP_518115 [Pan troglodytes]
	912	30	<a href="#">CAG33450</a>	<a href="#">48146455</a>	HSPC111 [Homo sapiens]
	882	21	<a href="#">XP_536424</a>	<a href="#">57085179</a>	PREDICTED: similar to Protein CGI-117 [Canis familiaris]
	863	21	<a href="#">NP_001...</a>	<a href="#">62460412</a>	hypothetical protein LOC506421 [Bos taurus]
	857	22	<a href="#">NP_848720</a>	<a href="#">30519925</a>	hypothetical protein LOC28126 [Mus musculus]
	857	22	<a href="#">BAC40235</a>	<a href="#">26353210</a>	unnamed protein product [Mus musculus]
	848	22	<a href="#">AAH12213</a>	<a href="#">15126561</a>	DNA segment, Chr 13, Wayne State University 177, expressed [Mus musculus]
	848	22	<a href="#">XP_225176</a>	<a href="#">27682843</a>	PREDICTED: similar to DNA segment, Chr 13, Wayne State University 177, expressed [Rattus norvegicus]
	665	30	<a href="#">AAH32424</a>	<a href="#">34783430</a>	Unknown (protein for MGC:40298) [Homo sapiens]
	597	18	<a href="#">XP_414551</a>	<a href="#">50754941</a>	PREDICTED: similar to Hypothetical protein HSPC111 [Gallus gallus]
	477	17	<a href="#">AAH73261</a>	<a href="#">49255993</a>	MGC80623 protein [Xenopus laevis]
	473	22	<a href="#">XP_357941</a>	<a href="#">63586355</a>	PREDICTED: similar to DNA segment, Chr 13, Wayne State University 177, expressed [Mus musculus]
	473	17	<a href="#">AAH68216</a>	<a href="#">45872610</a>	Unknown (protein for MGC:76225) [Xenopus tropicalis]
	463	15	<a href="#">CAF97276</a>	<a href="#">47221358</a>	unnamed protein product [Tetraodon nigroviridis]
	413	15	<a href="#">NP_001...</a>	<a href="#">50540130</a>	hypothetical protein LOC436807 [Danio rerio]
	404	22	<a href="#">XP_343880</a>	<a href="#">62655980</a>	PREDICTED: similar to DNA segment, Chr 13, Wayne State University 177, expressed [Rattus norvegicus]
	340	30	<a href="#">AAF36105</a>	<a href="#">7106760</a>	HSPC185 [Homo sapiens]
	280	15	<a href="#">XP_688990</a>	<a href="#">68356182</a>	PREDICTED: similar to Protein CGI-117 [Danio rerio]
	251	9	<a href="#">XP_782922</a>	<a href="#">72007783</a>	PREDICTED: similar to Protein CGI-117 [Strongylocentrotus purpuratus]
	227	8	<a href="#">XP_313729</a>	<a href="#">58385127</a>	ENSANGP00000016957 [Anopheles gambiae str. PEST]
	197	7	<a href="#">CAES7069</a>	<a href="#">39578894</a>	Hypothetical protein CBG24963 [Caenorhabditis briggsae]
	186	7	<a href="#">NP_492248</a>	<a href="#">17511059</a>	ZK265.6 [Caenorhabditis elegans]
	163	8	<a href="#">FAL28137</a>	<a href="#">54638735</a>	GA11068-PA [Drosophila pseudoobscura]
	160	8	<a href="#">NP_651884</a>	<a href="#">24651718</a>	CG11563-PA [Drosophila melanogaster]
	160	8	<a href="#">AAM29402</a>	<a href="#">21064345</a>	RE07994p [Drosophila melanogaster]
	111	3	<a href="#">XP_810459</a>	<a href="#">71417059</a>	hypothetical protein, conserved [Trypanosoma cruzi]
	108	7	<a href="#">NP_509607</a>	<a href="#">25152153</a>	F49E2.5b [Caenorhabditis elegans]
	108	7	<a href="#">NP_001...</a>	<a href="#">71991080</a>	F49E2.5h [Caenorhabditis elegans]
	108	7	<a href="#">NP_001...</a>	<a href="#">71991088</a>	F49E2.5i [Caenorhabditis elegans]
	108	7	<a href="#">NP_001...</a>	<a href="#">71991094</a>	F49E2.5j [Caenorhabditis elegans]
	108	7	<a href="#">NP_509608</a>	<a href="#">25152156</a>	F49E2.5c [Caenorhabditis elegans]
	108	7	<a href="#">NP_509605</a>	<a href="#">25152150</a>	F49E2.5a [Caenorhabditis elegans]
	108	7	<a href="#">NP_509606</a>	<a href="#">25152159</a>	F49E2.5d [Caenorhabditis elegans]
	105	3	<a href="#">CAJ02832</a>	<a href="#">68124716</a>	hypothetical protein, conserved [Leishmania major]

# KOG04771 Nucleolar protein (NOP16) involved in 60S ribosomal subunit biogenesis



KOG analyses help to fill in multiple gaps in proteome annotations: Neither of animal, plant, or protozoan proteins in KOG04771 is reasonably annotated in public databases



# Phylogenetic pattern searches

Phylectic Pattern Search - Main\_KOG\_Project

File Main Phylet Help

#	Name	&	^
3	Arabidopsis thaliana   oldArath		
4	Aspergillus fumigatus AF293   AspFu		
5	Aspergillus nidulans FGSC A4   Aspni		
6	Caenorhabditis briggsae   Caebr		
7	Caenorhabditis elegans   Caeel		
8	Caenorhabditis elegans   nCaeel		
9	Candida albicans SC5314   Canal		
10	Candida glabrata CBS138   Cangl		
11	Chlamydomonas reinhardtii   Chlre		
12	Ciona intestinalis   Cloin		
13	Cryptococcus neoformans var. neoformans JEC21   Cryne		
14	Cryptosporidium parvum   Crypa		
15	Cyanidioschyzon merolae   Cyame		
16	Danio rerio   Danre		
17	Debaryomyces hansenii CBS767   Debha		
18	Dictyostelium discoideum   Dicdi		
19	Drosophila melanogaster   Drome		
20	Drosophila melanogaster   nDrome		
21	Encephalitozoon cuniculi   Encuc		
22	Entamoeba histolytica   Enthi		
23	Erethothecium gossypii   Erego		
24	Gallus gallus   Galga		
25	Giardia lamblia ATCC 50803   Giala		
26	Gibberella zeae PH-1   Gibze		
27	Homo sapiens   oldHomsa		
28	Homo sapiens   Homsa		
29	Kluyveromyces lactis NRRL Y-1140   Klula		
30	Kluyveromyces waltii   Kluwa		
31	Leishmania major   Leima		
32	Magnaporthe grisea 70-15   Maggr		
33	Mus musculus   Musmu		
34	Neurospora crassa   Neucr		
35	Oryza sativa (japonica cultivar-group)   Orysa		
36	Pan troglodytes   Pantr		
37	Pan troglodytes   PantrE		
38	Plasmodium berghei   Plabe		
39	Plasmodium falciparum 3D7   Plafa		
40	Plasmodium yoelii yoelii   Playo		
41	Rattus norvegicus   Ratno		
42	Saccharomyces bayanus   Sacba		
43	Saccharomyces cerevisiae   Sacce		
44	Saccharomyces mikatae   Sacmi		
45	Saccharomyces paradoxus   Sacpa		
46	Schizosaccharomyces pombe   Schpo		
47	Takifugu rubripes   Takru		
48	Tetraodon nigroviridis   Tetni		
49	Thalassiosira pseudonana CCMP1335   Thaps		
50	Theileria parva   Thepa		
51	Trypanosoma brucei   Trybr		
52	Trypanosoma cruzi   Trycr		
53	Ustilago maydis 521   Ustma		
54	Yarrowia lipolytica CLIB99   Yarl		

Get COGs

7: Caenorhabditis elegans | Caeel

COG	Pattern
KOG01923 TZ Rac1 GTPase effector FHO5	.....
KOG01925 TZ Rac1 GTPase effector FHO5	.....
KOG01972 S Uncharacterized conserved protein	.....
KOG01974 L DNA topoisomerase I-interacting protein	.....
KOG01982 L Nuclear 5'-3' exoribonuclease-interacting protein; Rai1p	.....
KOG01988 S Uncharacterized conserved protein	.....
KOG01998 T Signaling protein DOCK180	.....
KOG02008 T BTK-associated SH3-domain binding protein SAB	.....
KOG02022 YU Nuclear transport receptor LGL2 (importin beta superfamily)	.....
KOG02046 Z Calponin	.....
KOG02057 F Predicted equilibrative nucleoside transporter protein	.....
KOG02060 U Rab3 effector RIM1 and related proteins; contain PDZ and C2 domains	.....
KOG02070 F Guanine nucleotide exchange factor	.....
KOG02075 S Topoisomerase TOP1-interacting protein BTBD1	.....
KOG02083 P Na+/K+ symporter	.....
KOG02096 R WD40 repeat protein	.....
KOG02106 S Uncharacterized conserved protein; contains HELP and WD40 domains	.....
KOG02113 R Predicted RNA binding protein; contains KH domain	.....
KOG02117 S Uncharacterized conserved protein	.....
KOG02119 G Predicted bile acid beta-glucosidase	.....
KOG02122 TZ Beta-catenin-binding protein APC; contains ARM repeats	.....
KOG02124 T Glycosylphosphatidylinositol anchor synthesis protein	.....
KOG02125 T Glycosylphosphatidylinositol anchor synthesis protein	.....
KOG02135 R Proteins containing the RNA recognition motif	.....
KOG02136 K Transcriptional regulators binding to the GC-rich sequences	.....
KOG02152 D Sister chromatid cohesion protein	.....
KOG02162 A Nonsense-mediated mRNA decay protein	.....
KOG02181 K LIM domain binding protein LDB1/NLI/CLIM	.....
KOG02185 A Predicted RNA-processing protein; contains G-patch domain	.....
KOG02192 AR PolyC-binding hnRNP-K protein HRB57A/hnRNP; contains KH domain	.....
KOG02193 AR IGF-II mRNA-binding protein IMP; contains RRM and KH domains	.....
KOG02216 S Conserved coiled/coiled coil protein	.....
KOG02218 UD ER to golgi transport protein/RAD50-interacting protein 1	.....
KOG02221 U PDZ-domain interacting protein EPI64; contains TBC domain	.....
KOG02222 TR Uncharacterized conserved protein; contains TBC; SH3 and RUN dom...	.....
KOG02224 TR Uncharacterized conserved protein; contains TBC domain	.....
KOG02230 G Predicted beta-mannosidase	.....
KOG02232 T Ceramidases	.....
KOG02233 U Alpha-N-acetylglucosaminidase	.....
KOG02244 R Highly conserved protein containing a thioredoxin domain	.....

17425 cogs

Transfer Pattern Save

Conn: 2

Among eukaryotes 17,425 KOGs are present only in the crown group



# Only 157 KOGs are represented in all KOGnitorized proteomes

KOG00010 OR Ubiquitin-like protein  
KOG00032 T Ca<sup>2+</sup>/calmodulin-dependent protein kinase; EF-Hand protein superfamily  
KOG00034 T Ca<sup>2+</sup>/calmodulin-dependent protein phosphatase (calcineurin subunit B); EF-Hand superfamily protein  
KOG00050 AD mRNA splicing protein CDC5 (Myb superfamily)  
KOG00070 U GTP-binding ADP-ribosylation factor Arf1  
KOG00077 U Vesicle coat complex COPII; GTPase subunit SAR1  
KOG00084 TU GTPase Rab1/YPT1; small G protein superfamily; and related GTP-binding proteins  
KOG00101 O Molecular chaperones HSP70/HSC70; HSP70 superfamily  
KOG00102 O Molecular chaperones mortalin/PBP74/GRP75; HSP70 superfamily  
KOG00131 A Splicing factor 3b; subunit 4  
KOG00166 U Karyopherin (importin) alpha  
KOG00173 O 20S proteasome; regulatory subunit beta type PSMB7/PSMB10/PUP1  
KOG00175 O 20S proteasome; regulatory subunit beta type PSMB5/PSMB8/PRE2  
KOG00188 J Alanyl-tRNA synthetase  
KOG00190 O Protein disulfide isomerase (prolyl 4-hydroxylase beta subunit)  
KOG00191 O Thioredoxin/protein disulfide isomerase  
KOG00206 R P-type ATPase  
KOG00211 T Protein phosphatase 2A regulatory subunit A and related proteins  
KOG00213 A Splicing factor 3b; subunit 1  
KOG00219 L Mismatch repair ATPase MSH2 (MutS family)  
KOG00232 C Vacuolar H<sup>+</sup>-ATPase V0 sector; subunits c/c'  
KOG00242 Z Kinesin-like protein  
KOG00264 B Nucleosome remodeling factor; subunit CAF1/NURF55/MSI1  
KOG00270 S WD40 repeat-containing protein  
KOG00284 A Polyadenylation factor I complex; subunit PFS2  
KOG00304 A mRNA deadenylase subunit  
KOG00307 U Vesicle coat complex COPII; subunit SEC31  
KOG00328 J Predicted ATP-dependent RNA helicase FAL1; involved in rRNA maturation; DE  
KOG00331 A ATP-dependent RNA helicase  
KOG00335 A ATP-dependent RNA helicase  
KOG00355 B DNA topoisomerase type II  
KOG00358 O Chaperonin complex component; TCP-1 delta subunit (CCT4)  
KOG00360 O Chaperonin complex component; TCP-1 alpha subunit (CCT1)  
KOG00361 O Chaperonin complex component; TCP-1 eta subunit (CCT7)  
KOG00362 O Chaperonin complex component; TCP-1 theta subunit (CCT8)  
KOG00366 O Protein geranylgeranyltransferase type II; beta subunit  
KOG00371 T Serine/threonine protein phosphatase 2A; catalytic subunit  
KOG00374 TR Serine/threonine specific protein phosphatase PP1; catalytic subunit  
KOG00397 J 60S ribosomal protein L11  
KOG00417 O Ubiquitin-protein ligase  
KOG00425 O Ubiquitin-protein ligase  
KOG00446 UR Vacuolar sorting protein VPS1; dynamin; and related proteins  
KOG00469 J Elongation factor 2  
KOG00481 L DNA replication licensing factor; MCM5 component  
KOG00531 T Protein phosphatase 1; regulatory subunit; and related proteins  
KOG00580 D Serine/threonine protein kinase  
KOG00583 T Serine/threonine protein kinase  
KOG00594 R Protein kinase PCTAIRE and related kinases  
KOG00657 G Glyceraldehyde 3-phosphate dehydrogenase  
KOG00676 Z Actin and related proteins  
KOG00706 T Predicted GTPase-activating protein  
KOG00712 O Molecular chaperone (DnaK superfamily)

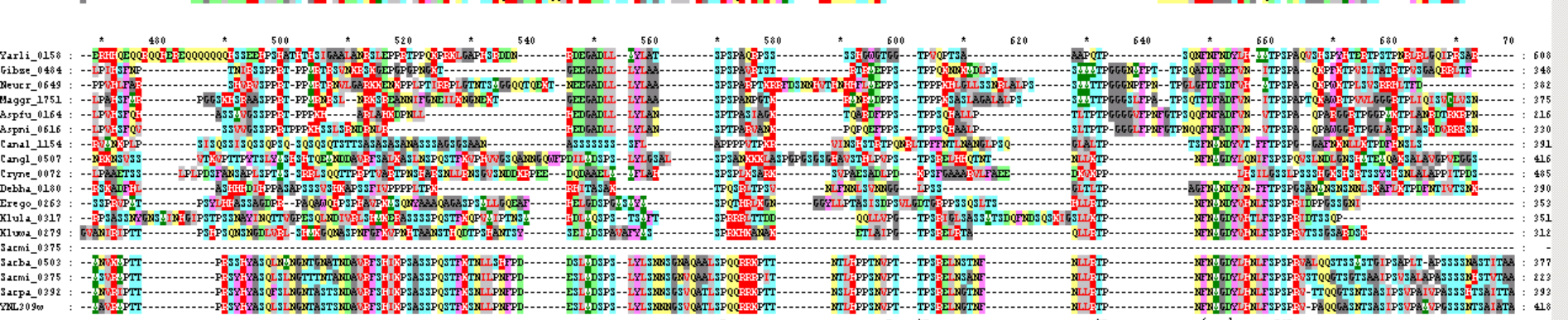
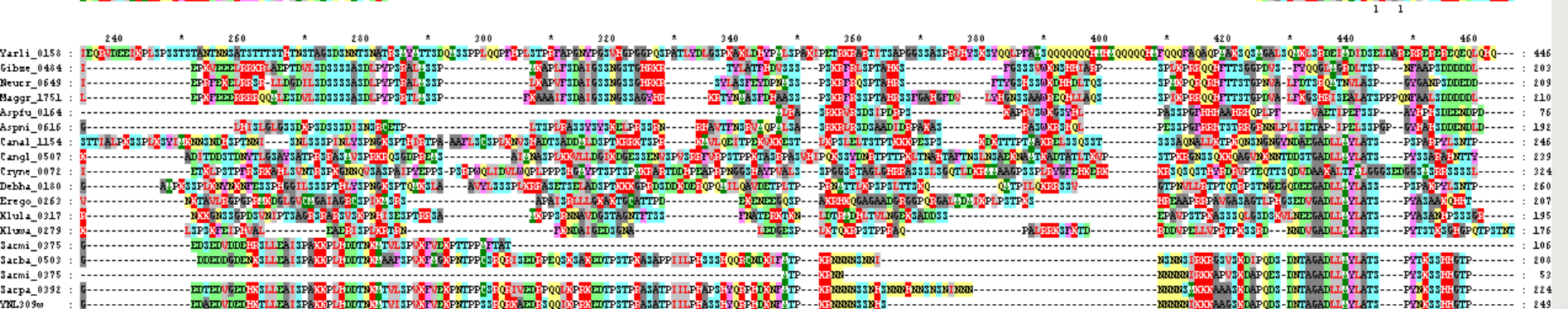
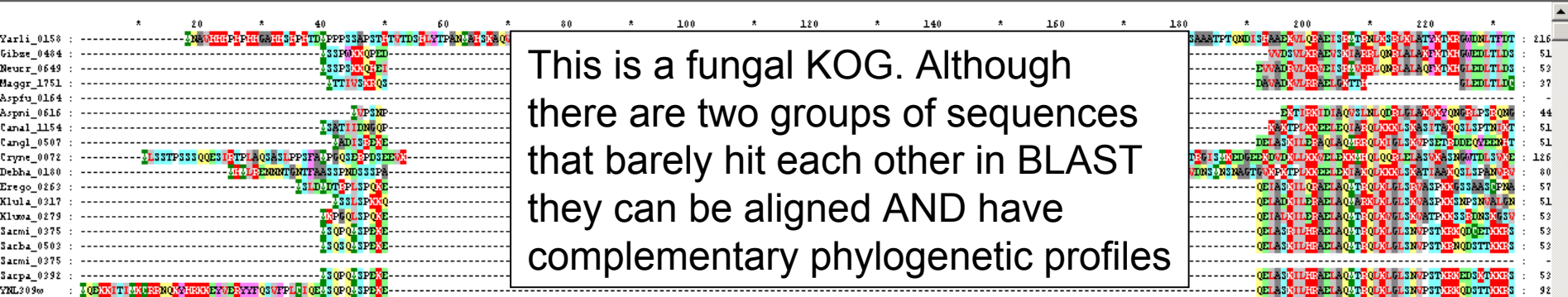
Because it is really hard to believe that some organisms do not have conventional ribosomal proteins we have to assume that analysis of phylogenetic profiles is the way to detect potentially missed genes

*Giardia lamblia*, *Plasmodium falciparum*, *Plasmodium bergeri*,  
*Cryptosporidium parvum*, *Encephalitozoon cuniculi*, *Leishmanija major*,  
and *Trypanosoma cruzi* share only 413 KOGs:

Functional categories in KOG db

102 Translation  
64 Posttranslational modification, protein turnover, chaperones  
43 RNA processing and modification  
31 Replication  
30 General function prediction only  
26 Transcription  
24 Intracellular trafficking, secretion  
15 Signal transduction mechanisms  
12 Energy production and conversion  
11 Unknown function  
11 Chromatin structure and dynamics  
9 Cytoskeleton  
7 Cell cycle control, cell division, chromosome partitioning  
6 Lipid transport and metabolism  
5 Amino acid transport and metabolism  
5 Nucleotide transport and metabolism  
3 Cell wall/membrane/envelope biogenesis  
2 Coenzyme transport and metabolism  
1 Inorganic ion transport and metabolism  
1 Secondary metabolites biosynthesis, transport and catabolism

Detection of complementary phylogenetic profiles in related KOGs will allow to make more aggressive hypotheses about orthology relationships and thus predict functions for greater number of proteins



Navigation toolbar with icons for file operations, editing, and viewing. Includes labels like 'C', 'Q', 'P', 'E', 'S', 'H', 'I', 'L', 'D', 'G', 'M', 'U'.

Sequence alignment view for the first block. Y-axis labels include Yarl1\_0158, Gibbe\_0484, Neucr\_0649, Maggr\_1751, AspFu\_0164, Aspni\_0616, Canal\_1154, Cangl\_0507, Cryne\_0072, Debba\_0180, Erego\_0263, Klula\_0217, Klwxa\_0279, Sacmi\_0375, Sacba\_0503, Sacmi\_0375, Sacmi\_0392, YNL309e. A large white box with black text reads 'One KOG. No kidding.' in the center of the alignment.

Sequence alignment view for the second block. Y-axis labels are the same as the first block. The alignment continues with various amino acid sequences and gap characters.

Sequence alignment view for the third block. Y-axis labels are the same as the first block. The alignment continues with various amino acid sequences and gap characters.

Sequence alignment view for the fourth block. Y-axis labels are the same as the first block. The alignment continues with various amino acid sequences and gap characters.

# Current status of KOG db: 50 Eukaryotic species; 46 species KOGnitorized, 1 updated;

712,427 protein sequences; two versions are present for 5 proteomes (+ 123,981 sequences)

21,816 clusters: 18,628 KOGs, 3,188 LSEs+mixed clusters

sequence data

species name	Taxa	sequences	In KOGs/LSEs	%	KOGs	published	source
<i>Chlamydomonas reinhardtii</i>	Chlorophyta (green algae)	19,922	10,194	51%	4,653	No	JGI
<i>Cyanidioschyzon merolae</i>	Rhodophyta (red algae)	5,040	3,769	74%	2,735		JGI
<i>Dictyostelium discoideum</i>	Mycetozoa	13,677	9,162	67%	4,573		GB
<i>Encephalitozoon cuniculi</i>	OLD Microsporidia	2,003	1,459	73%	1,132		GB
<i>Entamoeba histolytica</i>	Entamoebidae	10,088	7,728	77%	2,930		TIGR
<i>Giardia lamblia</i>	Diplomonadida group	6,573	2,105	32%	1,242	No	GB
<i>Leishmania major</i>	Euglenozoa	8,071	7,558	93%	4,502		TIGR
<i>Trypanosoma brucei</i>	Euglenozoa	8,291					GB
<i>Trypanosoma cruzi</i>	Euglenozoa	19,666	16,844	86%	5,023		GB
<i>Thalassiosira pseudonana</i>	Bacillariophyta (diatoms)	11,406	9,311	81%	4,033		JGI
<i>Cryptosporidium parvum</i>	Apicomplexa	3,410	2,212	66%	1,780		GB
<i>Plasmodium berghei</i>	Apicomplexa	11,782	7,004	59%	3,461	No	GB
<i>Plasmodium falciparum</i>	Apicomplexa	5,279	3,365	64%	3,403		GB
<i>Plasmodium yoelii</i>	Apicomplexa	7,868					GB
<i>Theileria parva</i>	Apicomplexa	4,080					GB
<i>Arabidopsis thaliana</i>	NEW Viridiplantae	26,759	25,781	96%	6,953		GB
<i>Arabidopsis thaliana</i>	OLD Viridiplantae	27,138	25,410	94%	6,925		GB
<i>Oryza sativa</i>	Viridiplantae	57,272	38,260	67%	6,701		TIGR
<i>Aspergillus fumigatus</i>	Fungi; Ascomycota; Pezizomycotina	10,029	8,959	89%	5,232	No	GB
<i>Aspergillus nidulans</i>	Fungi; Ascomycota; Pezizomycotina	10,317	9,629	93%	5,283	No	GB
<i>Gibberella zeae</i>	Fungi; Ascomycota; Pezizomycotina	12,221	11,153	91%	5,910	No	GB
<i>Magnaporthe grisea</i>	Fungi; Ascomycota; Pezizomycotina	11,381	9,043	79%	5,397		GB
<i>Neurospora crassa</i>	Fungi; Ascomycota; Pezizomycotina	10,478	8,206	78%	5,559		GB
<i>Candida albicans</i>	Fungi; Ascomycota; Saccharomycotina	14,274	10,952	77%	4,094		GB
<i>Candida glabrata</i>	Fungi; Ascomycota; Saccharomycotina	5,216	4,992	96%	3,790	No	GB
<i>Debaryomyces hansenii</i>	Fungi; Ascomycota; Saccharomycotina	6,350	5,700	90%	4,048		GB
<i>Eremothecium gossypii</i>	Fungi; Ascomycota; Saccharomycotina	4,757	4,619	97%	3,810	No	GB
<i>Kluyveromyces lactis</i>	Fungi; Ascomycota; Saccharomycotina	5,366	5,058	94%	3,980		GB
<i>Kluyveromyces waltii</i>	Fungi; Ascomycota; Saccharomycotina	5,240	5,055	96%	3,919	No	MIT
<i>Yarrowia lipolytica</i>	Fungi; Ascomycota; Saccharomycotina	6,596	5,728	87%	4,018		GB
<i>Saccharomyces bayanus</i>	Fungi; Ascomycota; Saccharomycotina	12,040	8,488	70%	5,879		MIT
<i>Saccharomyces cerevisiae</i>	OLD Fungi; Ascomycota; Saccharomycotina	6,387	6,146	96%	4,427		GB
<i>Saccharomyces mikatae</i>	Fungi; Ascomycota; Saccharomycotina	10,350	8,481	82%	5,828		MIT
<i>Saccharomyces paradoxus</i>	Fungi; Ascomycota; Saccharomycotina	10,601	8,375	79%	6,005		MIT
<i>Schizosaccharomyces pombe</i>	OLD Fungi; Ascomycota; Schizosaccharomycetes	5,064	4,461	88%	3,360		GB
<i>Cryptococcus neoformans</i>	Fungi; Basidiomycota; Hymenomycetes	6,640	5,471	82%	3,702	No	GB
<i>Ustilago maydis</i>	Fungi; Basidiomycota; Ustilaginomycetes	6,685	5,221	78%	3,779	No	GB
<i>Anopheles gambiae</i>	Metazoa; Insecta; Diptera; Nematocera	13,895	12,405	89%	5,946		GB
<i>Apis mellifera</i>	Metazoa; Insecta; Hymenoptera	16,997	14,785	87%	4,906	No	EMBL
<i>Drosophila melanogaster</i>	NEW Metazoa; Insecta; Diptera; Brachycera	18,873					GB
<i>Drosophila melanogaster</i>	OLD Metazoa; Insecta; Diptera; Brachycera	13,936	11,552	83%	6,241		GB
<i>Caenorhabditis briggsae</i>	Metazoa; Nematoda	19,523	17,963	92%	6,938		GB
<i>Caenorhabditis elegans</i>	NEW Metazoa; Nematoda	21,193					GB
<i>Caenorhabditis elegans</i>	OLD Metazoa; Nematoda	21,246	20,461	96%	7,175		GB
<i>Ciona intestinalis</i>	Metazoa; Chordata; Ascidiacea	16,139	12,101	75%	5,684		JGI
<i>Danio rerio</i>	Vertebrata; Teleostomi	19,893	19,676	98%	4,426		GB
<i>Gallus gallus</i>	Vertebrata; Aves	28,751	26,613	93%	6,894	No	EMBL
<i>Homo sapiens</i>	NEW Vertebrata; Mammalia	27,524					GB
<i>Homo sapiens</i>	OLD Vertebrata; Mammalia	39,436	29,835	77%	9,304		GB
<i>Mus musculus</i>	Vertebrata; Mammalia	26,464	24,784	94%	9,284		GB
<i>Pan troglodytes</i>	GB Vertebrata; Mammalia	22,225					GB
<i>Pan troglodytes</i>	EMBL Vertebrata; Mammalia	38,845	36,984	95%	8,923		EMBL
<i>Rattus norvegicus</i>	Vertebrata; Mammalia	21,531					GB
<i>Takifugu rubripes</i>	Vertebrata; Teleostomi	33,357	31,879	96%	6,909		GB
<i>Tetraodon nigroviridis</i>	Vertebrata; Teleostomi	28,262	25,520	90%	7,033	No	GB

# Some conclusions

- Clusters of orthology domains (KOGs) will form a platform for comparative genomics and evolutionary studies
- At present KOG db cannot be updated completely automatically
- A new version of eukaryotic clusters of orthology domains is being developed; it will contain ca. 25,000 clusters for 50 proteomes and HMM profiles for all clusters
- An approach to add proteomes semi-automatically is developed which is based on picking conflicts between BLASTP (KOGnitor) and HMM searches
- KOG db will serve as a tool for annotation of newly sequenced genomes, for fixing errors and filling in gaps in existing annotations, for detecting incorrect gene models and genes missed by gene predictions

# **Acknowledgements**

**Sergei Smirnov, NCBI**

**Yuri Wolf , NCBI**

**Alex Sorokin , NCBI**

**Eugene Koonin , NCBI**





**Database connection error. Please try again later**