# Functional Annotation



May 23, 2007

Rama Maiti

# Functional Annotation Overview

- What is annotation

- Steps we take to annotate eukaryotic genes

- Software tools we use for functional annotation

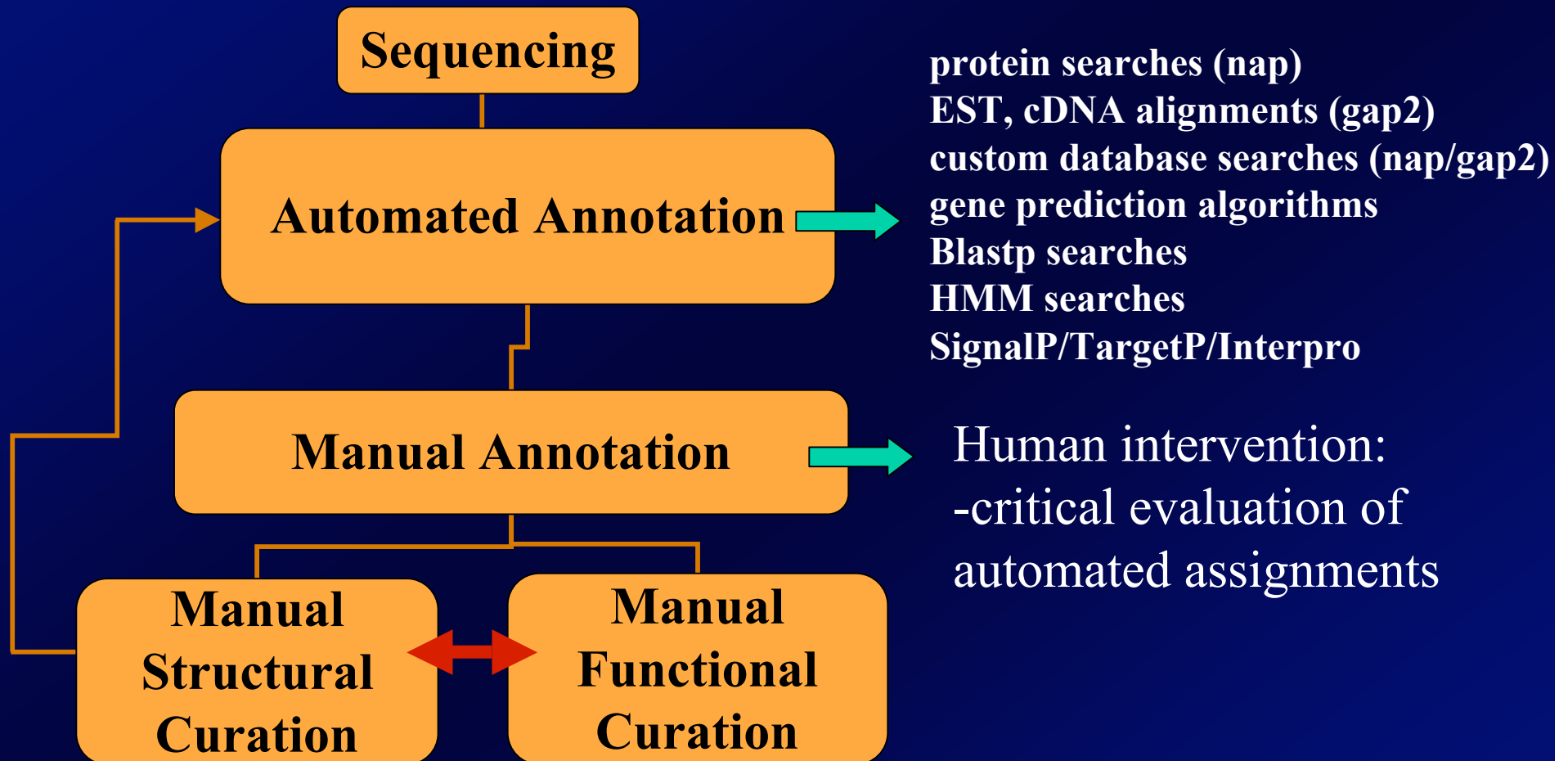- Steps we take to manually annotate or verify an automated annotation

# What are the questions?

- How did the gene get its structure and name?

- Does it really have a function assigned to it?

- Where did this information come from?

- Is it accurate? Can you rely on it?

# What is Functional Annotation?

- "To annotate" is "to make or furnish critical or explanatory notes or comments"

- For genomics the 'notes' are about

  - Names of the gene products
  - Functions of genes within an organism

- Elements of the functional annotation process

  - Validation of the gene structure
  - Literature search, if any is available
  - Homology / domain searches
  - Assignment of function
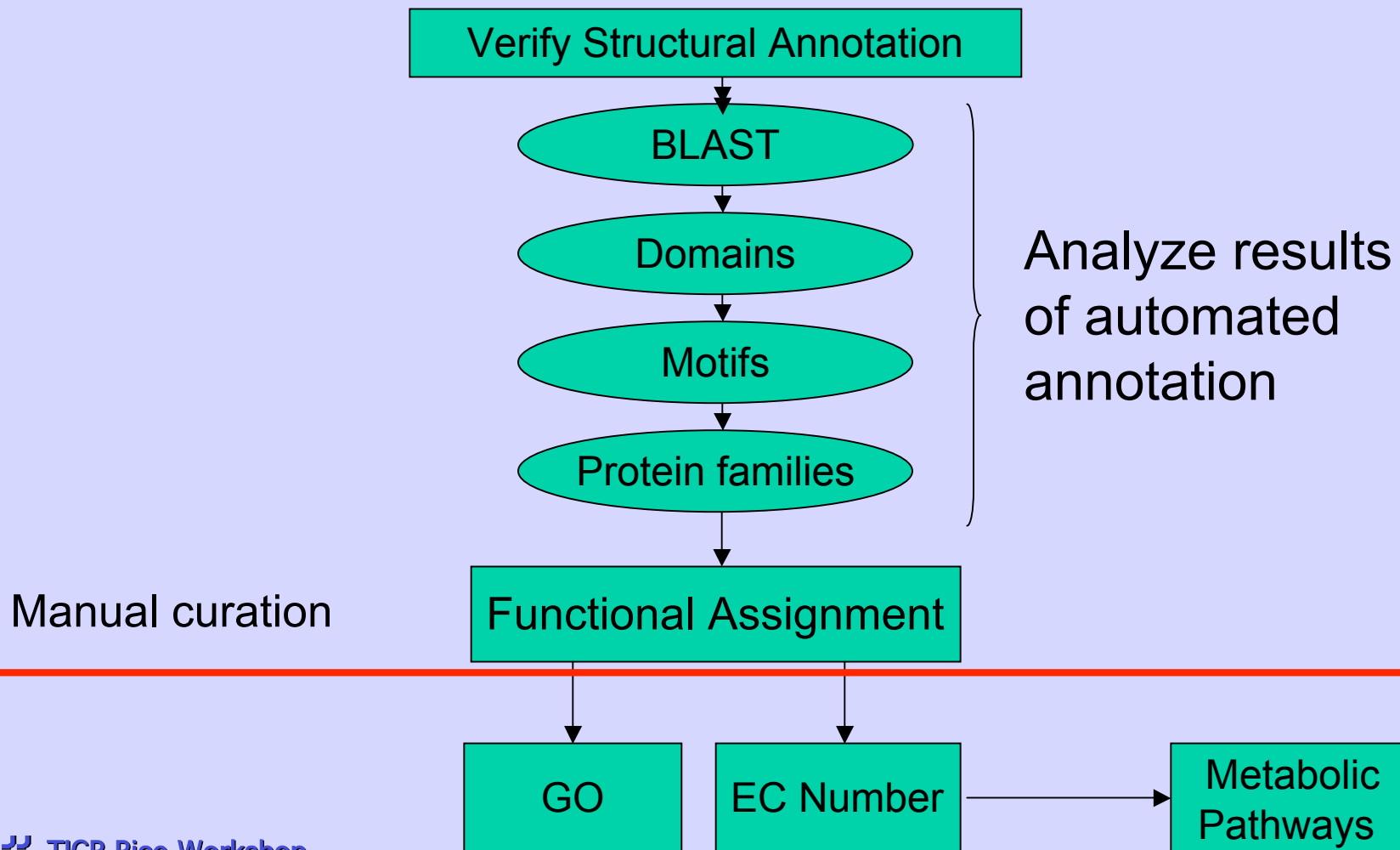  - Maintenance of data availability

# The Annotation Pipeline

**Sequencing**

**Automated Annotation**

protein searches (nap)
EST, cDNA alignments (gap2)
custom database searches (nap/gap2)
gene prediction algorithms
Blastp searches
HMM searches
SignalP/TargetP/Interpro

**Manual Annotation**

Human intervention:
-critical evaluation of
automated assignments

**Manual Structural Curation**

**Manual Functional Curation**

TIGR Rice Workshop

# Manual vs. Automated Annotation

- Automated Annotation is complicated by high volumes of data derived from different methods at different centers

- High quality annotation requires manual review and intervention.

# Steps in Functional Annotation



**Verify Structural Annotation**

BLAST

Domains

Motifs

Protein families

Analyze results of automated annotation

Manual curation

**Functional Assignment**

GO

EC Number

Metabolic Pathways

TIGR Rice Workshop

# Steps in Functional Annotation

- Analyze the gene structure (Annotation Station or preferred gene viewer)
- Name the gene product (Manatee)
  - requires analysis of the gene product
  - gene product name is primarily homology based on different bioinformatics tools
- Assign Gene Ontology terms
  - Process
  - Function
  - Component

# Homology Searching

## (Tools that are available to characterize a sequence)

- **WU BLAST** http://blast.wustl.edu/ with links to many servers

- **NCBI BLAST** http://www.ncbi.nlm.nih.gov/blast/

- **Pfam profiles** (profiles, or HMMs) http://pfam.wustl.edu/

- **TIGRFAMS** (profiles, or HMMs) http://tigrblast.tigr.org/web-hmm/

- **Prosite** (profiles & families) http://ca.expasy.org/tools/scanprosite/

- **Interpro** (families) http://www.ebi.ac.uk/InterProScan/

- **TmHMM** (transmembrane domain) http://www.cbs.dtu.dk/services/TMHMM/

- **Swiss-Prot http://au.expasy.org/sprot/**

- **SignalP** (signal peptide cleavage sites) http://www.cbs.dtu.dk/services/SignalP/

- **TargetP** (subcellular location) http://www.cbs.dtu.dk/services/TargetP/

- **PSI-BLAST** (NCBI) link at http://www.ncbi.nlm.nih.gov/BLAST/

- **Protein families and clustering**

  - **TIGR Paralogous Families** (not yet available outside of TIGR)

  - **TribeMCL** http://www.ebi.ac.uk/research/cgg/tribe/

# Manatee

- Manatee is a web-based gene evaluation and genome annotation tool.

- Manatee displays the current annotation for prokaryotic and eukaryotic genomes.

- Manatee is an open source software available at:

  http://sourceforge.net/projects/manatee/

make high quality functional assignments using genome analyses tools. These tools consist of, but are not limited to GO classifications, blast search data, protein families.

# Verify evidence from automated annotation

- BLAST matches
- HMM
- Prosite, Interpro classifications
- Motifs
- Signal Sequence
- Target Sequence
- Transmembrane domain
- Protein families

# Functional annotation

Examine the gene structure

does it make sense with respect to the alignments?
do you need to re-curate the gene structure?

Name the gene product

Determine whether it is published,
Fully characterized?  Give it the Swiss-Prot name.
Sequenced but not characterized? Look at  the evidence.

Add comments to comment field

explain reasoning for others
add personal communication information
make comments about function or process

In many occasions after analyzing our data and make a decision about a gene function, we may need to go back and re-examine the gene structure.

# Use all possible resources...

# Example:-

A protein sequence from *Trypanosoma brucei*. Our task will be to annotate this protein sequence as fully as possible, given the tools at hand.

protein sequence:

```
>unknown_T. brucei protein_sequence
MLRRLGVRHFRRTPLLFVGGDGSIFERY
TEIDNSNERRINALKGCGMFEDEWIATE
KVHGANFGIYSIEGEKMIRYAKRSGIMP
PNEHFFGYHILIPELQRYITSIREMLCEK
QKKKLHVVLINGELFGGKYDHPSVPKT
RKTVMVAGKPRTISAVQTDSFPQYSPDL
HFYAFDIKYKETEDGDYTTLVYDEAIEL
FQRVPGLLYARAVIRGPMSKVAAFDVE
RFVTTIPPLVGMGNYPLTGNWAEGLVV
KHSRLGMAGFDPKGPTVLKFKCTAFQE
ISTDRAQGPRVDEMRNVRRDSINRAGVQ
LPDLESIVQDPIQLEASKLLLNHVCENRL
KNVLSKIGTEPFEKEEMTPDQLATLLAK
DVLKDFLKDTEPSIVNIPVLIRKDLTRYV
IFESRRLVCSQWKDILKRQSPDFSE*
```

# **Verify the gene structure**

# NCBI BLAST

NCBI BLAST tools at:
http://www.ncbi.nlm.nih.gov/blast/.

| Program | Database | Query |
|---------|----------|-------|
| BLASTN | Nucleotide | Nucleotide |
| BLASTP | Protein | Protein |
| BLASTX | Protein | Nucleotide → Protein |
| TBLASTN | Nucleotide → Protein | Protein |
| TBLASTX | Nucleotide → Protein | Nucleotide → Protein |

*Read → as "translated to"*

# BLAST: What makes a good alignment?

It depends on what you are trying to prove!

- minimum of 30% identity, better 40% & up
  - higher for short proteins
  - score is weighted for length

- full length match
  - at least 70% of both proteins

? See explanation of BLAST scores on slide 56.

# Example : run NCBI BLAST

BLASTP – protein against protein

Results:

The first hit in the BLASTP output, a 100% match, is to a genome project submission, which means that the entry is not

**Alignments**

>gi|115504417|ref|XP_001219001.1| RNA editing ligase; RNA-editing complex protein; KREL2 [Trypanosoma brucei]

gi|83642483|emb|CAJ16514.1| RNA editing ligase; RNA-editing complex protein; KREL2 [Trypanosoma brucei]

Length=416

Score = 860 bits (2222), Expect = 0.0, Method: Composition-based stats.
Identities = 416/416 (100%), Positives = 416/416 (100%), Gaps = 0/416 (0%)

```
Query  1    MLRRLGVRHFRRTPLLFVGGDGSIFERYTEIDNSNERRINALKGCGMFEDEWIATEKVHG  60
            MLRRLGVRHFRRTPLLFVGGDGSIFERYTEIDNSNERRINALKGCGMFEDEWIATEKVHG
Sbjct  1    MLRRLGVRHFRRTPLLFVGGDGSIFERYTEIDNSNERRINALKGCGMFEDEWIATEKVHG  60

Query  61   ANFGIYSIEGEKMIRYAKRSGIMPPNEHFFGYHILIPELQRYITSIREMLCEKQKKKLHV  120
            ANFGIYSIEGEKMIRYAKRSGIMPPNEHFFGYHILIPELQRYITSIREMLCEKQKKKLHV
Sbjct  61   ANFGIYSIEGEKMIRYAKRSGIMPPNEHFFGYHILIPELQRYITSIREMLCEKQKKKLHV  120
```

```
AUTHORS   ...rd,N.J., Harris,B.R., Hertz-Fowler,C.,
          ...d,C.S., Atkin,R.J., Barron,A.J.,
          ...Bringaud,F., Clark,L.N., Corton,C.H.,
          ...t,J., Fraser,A., Gruter,E., Hall,S.,
          Harper,A.D., Kay,M.P., Leech,V., Mayes,R., Price,C., Quail,M.A.,
          Rabbinowitsch,E., Reitter,C., Rutherford,K., Sasse,J., Sharp,S.,
          Shownkeen,R., MacLeod,A., Taylor,S., Tweedie,A., Turner,C.M.,
          Tait,A., Gull,K., Barrell,B. and Melville,S.E.
TITLE     The DNA sequence of chromosome I of an African trypanosome: gene
          content, chromosome organisation, recombination and polymorphism
JOURNAL   Nucleic Acids Res. 31 (16), 4864-4873 (2003)
 PUBMED   12907729
REFERENCE 2
AUTHORS   Berriman,M., Hertz-Fowler,C.V.A., Hall,N., Kerhornou,A.X.,
          Bowman,S., Quail,M., Kay,M.P., Bray-Allen,S., Lennard,N.J.,
          Clark,L.N., Harris,B.R., Melville,S., Gerrard,C., Rajandream,M.A.
          and Barrell,B.G.
TITLE     Direct Submission
JOURNAL   Submitted (20-SEP-2002) The Wellcome Trust Sanger Institute,
          Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK
REMARK    revised by [3]
REFERENCE 3  (residues 1 to 416)
AUTHORS   Hertz-Fowler,C. and Berriman,M.
TITLE     Direct Submission
```

**Example : navigating BLAST output**

>gi|47117107|sp|P82864|TB48_TRYBB    RNA editing ligase TbMP48, mitochondrial
gi|11067028|gb|AAG27063.1|    RNA ligase MP48 [Trypanosoma brucei]
Length=416

Score =  856 bits (2212),  Expect = 0.0, Method: Composition-based stats.
Identities = 413/416 (99%), Positives = 414/416 (99%), Gaps = 0/416 (0%)

The second hit in the BLAST output, a 99% match, is to a Swiss-Prot entry.

The alignment reveals three positions with variations:

I103V (very similar, both hydrophobic) conservative

D182G (negative, hydrophilic to tiny polar) non-conservative

V364A (nonpolar, aliphatic, hydrophobic to tiny, nonpolar, aliphatic) conservative

```
Query   1    MLRRLGVRHFRRTPLLFVGGDGSIFERYTEIDNSNERRINALKGCGMFEDEWIATEKVHG    60
             MLRRLGVRHFRRTPLLFVGGDGSIFERYTEIDNSNERRINALKGCGMFEDEWIATEKVHG
Sbjct   1    MLRRLGVRHFRRTPLLFVGGDGSIFERYTEIDNSNERRINALKGCGMFEDEWIATEKVHG    60

Query   61   ANFGIYSIEGEKMIRYAKRSGIMPPNEHFFGYHILIPELQRYITSIREMLCEKQKKKLHV   120
             ANFGIYSIEGEKMIRYAKRSGIMPPNEHFFGYHILIPELQRY+TSIREMLCEKQKKKLHV
Sbjct   61   ANFGIYSIEGEKMIRYAKRSGIMPPNEHFFGYHILIPELQRYVTSIREMLCEKQKKKLHV   120

Query   121  VLINGELFGGKYDHPSVPKTRKTVMVAGKPRTISAVQTDSFPQYSPDLHFYAFDIKYKET   180
             VLINGELFGGKYDHPSVPKTRKTVMVAGKPRTISAVQTDSFPQYSPDLHFYAFDIKYKET
Sbjct   121  VLINGELFGGKYDHPSVPKTRKTVMVAGKPRTISAVQTDSFPQYSPDLHFYAFDIKYKET   180

Query   181  EDGDYTTLVYDEAIELFQRVPGLLYARAVIRGPMSKVAAFDVERFVTTIPPLVGMGNYPL   240
             E GDYTTLVYDEAIELFQRVPGLLYARAVIRGPMSKVAAFDVERFVTTIPPLVGMGNYPL
Sbjct   181  EGGDYTTLVYDEAIELFQRVPGLLYARAVIRGPMSKVAAFDVERFVTTIPPLVGMGNYPL   240

Query   241  TGNWAEGLVVKHSRLGMAGFDPKGPTVLKFKCTAFQEISTDRAQGPRVDEMRNVRRDSIN   300
             TGNWAEGLVVKHSRLGMAGFDPKGPTVLKFKCTAFQEISTDRAQGPRVDEMRNVRRDSIN
Sbjct   241  TGNWAEGLVVKHSRLGMAGFDPKGPTVLKFKCTAFQEISTDRAQGPRVDEMRNVRRDSIN   300

Query   301  RAGVQLPDLESIVQDPIQLEASKLLLNHVCENRLKNVLSKIGTEPFEKEEMTPDQLATLL   360
             RAGVQLPDLESIVQDPIQLEASKLLLNHVCENRLKNVLSKIGTEPFEKEEMTPDQLATLL
Sbjct   301  RAGVQLPDLESIVQDPIQLEASKLLLNHVCENRLKNVLSKIGTEPFEKEEMTPDQLATLL   360

Query   361  AKDVLKDFLKDTEPSIVNIPVLIRKDLTRYVIFESRRLVCSQWKDILKRQSPDFSE      416
             AKD LKDFLKDTEPSIVNIPVLIRKDLTRYVIFESRRLVCSQWKDILKRQSPDFSE
Sbjct   361  AKDALKDFLKDTEPSIVNIPVLIRKDLTRYVIFESRRLVCSQWKDILKRQSPDFSE      416
```
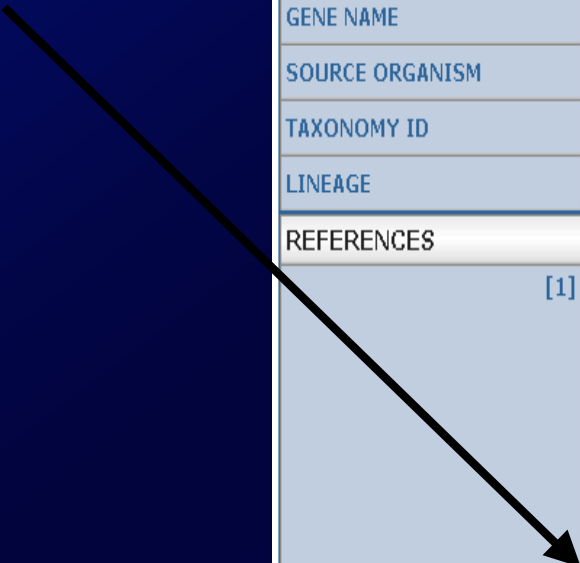
? See Glossary entry for SNP

# Swiss-Prot

Our sequence is 99% identical to the sequence of this Swiss-Prot entry.

Another name for this protein in the literature is 'REL2.'

## ENTRY INFORMATION

| | |
|---|---|
| ENTRY NAME | TB48_TRYBB |
| ACCESSION NUMBER | P82864 |
| Integrated into Swiss-Prot on | 2004-05-10 |
| Sequence was last modified on | 2001-03-01 (Sequence version 1) |
| Annotations were last modified on | 2006-10-31 (Entry version 24) |

## NAME AND ORIGIN OF THE PROTEIN

| | |
|---|---|
| PROTEIN NAME | RNA-editing ligase TbMP48, mitochondrial precursor |
| Synonyms | EC 6.5.1.3 <br> RNA ligase |
| GENE NAME | MP48 |
| SOURCE ORGANISM | Trypanosoma brucei brucei |
| TAXONOMY ID | 5702 [NCBI, NEWT] |
| LINEAGE | Eukaryota; Euglenozoa; Kinetoplastida; Trypanosomatidae; Trypanosoma |

## REFERENCES

| | |
|---|---|
| [1] | Panigrahi AK; Gygi SP; Ernst NL; Igo RP Jr; Palazzo SS; Schnaufer A et al. View all. <br> **Association of two novel proteins TbMP52 and TbMP48 with the Trypanosoma brucei RNA editing complex.** <br> 2001, *Mol. Cell. Biol.*, 21, 380-389. <br> *Position*: NUCLEOTIDE SEQUENCE [GENOMIC DNA], PROTEIN SEQUENCE OF 18-37; 58-72; 118-139; 143-151; 200-207; 217-224; 255-263; 302-323; 336-340; 371-384 AND 410-416, FUNCTION, AND SUBCELLULAR LOCATION.. <br> PubMed: 11134327; Medline: 20576857. |

## COMMENTS

| | |
|---|---|
| FUNCTION | Part of the RNA editing complex essential for cell variability. RNA editing in kinetoplastid mitochondria inserts and deletes uridylates at multiple sites in pre-mRNAs as directed by guide RNAs. |
| CATALYTIC ACTIVITY | ATP + (ribonucleotide)(n) + (ribonucleotide)(m) = AMP + diphosphate + (ribonucleotide)(n+m). |
| SUBCELLULAR LOCATION | Mitochondrion. |

# Swiss-Prot

Click on the NCBI hyperlink to look at this publication.

## ENTRY INFORMATION

| | |
|---|---|
| ENTRY NAME | TB48_TRYBB |
| ACCESSION NUMBER | P82864 |
| Integrated into Swiss-Prot on | 2004-05-10 |
| Sequence was last modified on | 2001-03-01 (Sequence version 1) |
| Annotations were last modified on | 2006-10-31 (Entry version 24) |

## NAME AND ORIGIN OF THE PROTEIN

| | |
|---|---|
| PROTEIN NAME | RNA-editing ligase TbMP48, mitochondrial precursor |
| Synonyms | EC 6.5.1.3 <br> RNA ligase |
| GENE NAME | MP48 |
| SOURCE ORGANISM | Trypanosoma brucei brucei |
| TAXONOMY ID | 5702 [NCBI, NEWT] |
| LINEAGE | Eukaryota; Euglenozoa; Kinetoplastida; Trypanosomatidae; Trypanosoma |

## REFERENCES

[1] Panigrahi AK; Gygi SP; Ernst NL; Igo RP Jr; Palazzo SS; Schnaufer A et al. View all.
Association of two novel proteins TbMP52 and TbMP48 with the Trypanosoma brucei RNA editing complex.
2001, Mol. Cell. Biol., 21, 380-389.
Position: NUCLEOTIDE SEQUENCE [GENOMIC DNA], PROTEIN SEQUENCE OF 18-37; 58-72; 118-139; 143-151; 200-207; 217-224; 255-263; 302-323; 336-340; 371-384 AND 410-416, FUNCTION, AND SUBCELLULAR LOCATION..
PubMed: 11134327; Medline: 20576857.

## COMMENTS

| | |
|---|---|
| FUNCTION | Part of the RNA editing complex essential for cell variability. RNA editing in kinetoplastid mitochondria inserts and deletes uridylates at multiple sites in pre-mRNAs as directed by guide RNAs. |
| CATALYTIC ACTIVITY | ATP + (ribonucleotide)(n) + (ribonucleotide)(m) = AMP + diphosphate + (ribonucleotide)(n+m). |
| SUBCELLULAR LOCATION | Mitochondrion. |

# Pubmed

- Read abstract
- If promising, read paper to be sure protein is characterized
- If characterized, it is good <u>evidence</u> for naming our sequence

1: Mol Cell Biol. 2001 Jan;21(2):380-9.

Full Text
Mol Cell B

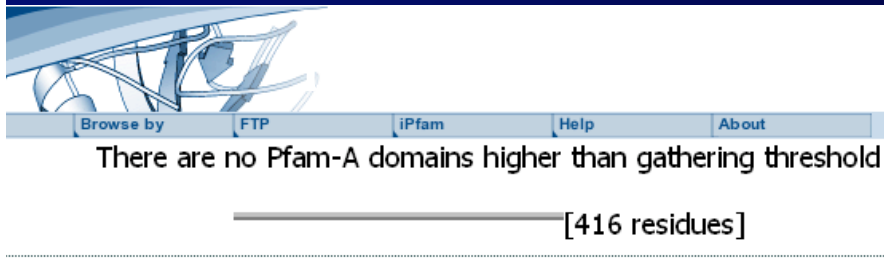**Association of two novel proteins, TbMP52 and TbMP48, with the Trypanosoma brucei RNA editing complex.**

Panigrahi AK, Gygi SP, Ernst NL, Igo RP Jr, Palazzo SS, Schnaufer A, Weston DS, Carmean N, Salavati R, Aebersold R, Stuart KD.

Seattle Biomedical Research Institute, Seattle, Washington 98109, USA.

RNA editing in kinetoplastid mitochondria inserts and deletes uridylates at multiple sites in pre-mRNAs as directed by guide RNAs. This occurs by a series of steps that are catalyzed by endoribonuclease, 3'-terminal uridylyl transferase, 3'-exouridylylase, and RNA ligase activities. A multiprotein complex that contains these activities and catalyzes deletion editing in vitro was enriched from Trypanosoma brucei mitochondria by sequential ion-exchange and gel filtration chromatography, followed by glycerol gradient sedimentation. The complex size is approximately 1,600 kDa, and the purified fraction contains 20 major polypeptides. A monoclonal antibody that was generated against the enriched complex reacts with an approximately 49-kDa protein and specifically immunoprecipitates in vitro deletion RNA editing activity. The protein recognized by the antibody was identified by mass spectrometry, and the corresponding gene, designated TbMP52, was cloned. Recombinant TbMP52 reacts with the monoclonal antibody. Another novel protein, TbMP48, which is similar to TbMP52, and its gene were also identified in the enriched complex. These results suggest that TbMP52 and TbMP48 are components of the RNA editing complex.

PMID: 11134327 [PubMed - indexed for MEDLINE]

# Domains (HMMs) TIGRFAMs search

There are no Pfam-A domains higher than gathering threshold

[416 residues]

Total score:        923.1
Trusted cutoff:     100.0
Gathering cutoff:   100.0
Noise cutoff:       -165.0

This is a very positive hit to the RNA ligase RNL2 family domain (TIGR02307).

Verify Structural Annotation
BLAST
Domains
Motifs
Protein families
Functional Assignment
GO
EC Number
Metabolic Pathways

```
hmmpfam - search a single seq against HMM database
HMMER 2.1.1 (Dec 1998)
Copyright (C) 1992-1998 Washington University School of Medicine
HMMER is freely distributed under the GNU General Public License (GPL).
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
HMM file:               ALL_LIB_bin.HMM
Sequence file:          hmmpfam-search-14395-1172255172.in
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
Query:  unknown_T.  brucei protein_sequence

Scores for sequence family classification (score includes all domains):
Model      Description                               Score    E-value   N
--------   -----------                               -----    -------  ---
TIGR02307 RNA_lig_RNL2: RNA ligase, Rnl2 family      923.1   7.7e-274   1

Parsed for domains:
Model      Domain   seq-f seq-t    hmm-f hmm-t    score  E-value
--------   ------   ----- -----    ----- -----    -----  -------
TIGR02307   1/1      25    408 ..      1   421 []   923.1 7.7e-274

Alignments of top-scoring domains:
TIGR02307: domain 1 of 1, from 25 to 408: score 923.1, E = 7.7e-274
                  *->FkkYTsleNssyrrifaeKltglglrGGEWVAlEKiHGaNFSiivee
                     F++YT+++Ns++rri+a+K++g++++   EW+A+EK+HGaNF+i+++e
  unknown_T.     25    FERYTEIDNSNERRINALKGCGMFED--EWIATEKVHGANFGIYSIE 69

                     dPNEAqDGAEkkVtfAKRtGiidPnEdGDYDFFGYhilieeytakvkAis
                     +           Ek++++AKR+Gi++PnE+    FFGYhili+e++++++i+
  unknown_T.     70 G--------EKMIRYAKRSGIMPPNEH----FFGYHILIPELQRYITSIR 107

                     dlLkekaGvikklesvivyGELaGkgyqkpvvPKsrKtvtlanKkRiISG
                     ++L+ek+  +kkl++v+++GEL+G++y++p+vPK+rKtv++a+K+R+IS
  unknown_T.    108 EMLCEKQ--KKKLHVVLINGELFGGKYDHPSVPKTRKTVMVAGKPRTIS- 154

                     vevQsdsFPQYsPDkdFyAFDIkyketGeeeddvtLvyDevlEvfervpk
                     +vQ+dsFPQYsPD++FyAFDIkyket e++d++tLvyDe++E+f+rvp+
  unknown_T.    155 -AVQTDSFPQYSPDLHFYAFDIKYKET-EDGDYTTLVYDEAIELFQRVPG 202

                     lkyAkelvRGtldEllafDNDLDSVVqvenFvtdlPaLVdlgnypLNAEA
                     l+yA++++RG++++++afD         ve+Fvt++P+LV++gnypL
  unknown_T.    203 LLYARAVIRGPMSKVAAFD---------VERFVTTIPPLVGMGNYPL---- 240
```

# Verify HMM

Total score: 923.1
Trusted cutoff: 100.0
Gathering cutoff: 100.0
Noise cutoff: -165.0

Score is well above the trusted cutoff.



Our sequence contains an RNA ligase, Rnl2 family domain, with a very strong match. Members of this TIGRfam family ligate RNA.
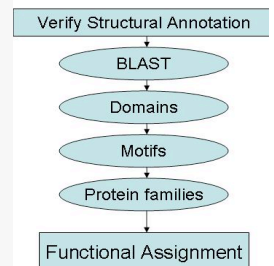
# Non-secretory protein



SignalP-HMM result:

SignalP-HMM prediction (euk models): unknown

Cleavage prob. ——
n-region prob. ——
h-region prob. ——
c-region prob. ——

MLRRLGVRHFRRTPLLFVGGDGSIFERYTEIDNSNERRINALKGCGMFEDEWIATEKVHGANFGIYSIEG

# data

>unknown
Prediction: Non-secretory protein
Signal peptide probability: 0.008
Signal anchor probability: 0.009
Max cleavage site probability: 0.006 between pos. 22 and 23

Verify Structural Annotation
BLAST
Domains
Motifs
Protein families
Functional Assignment

# TargetP

The sequence contains a mitochondrial targeting peptide, mTP.

# Transmembrane domains

There are no transmembrane domains.

## TMHMM result

[HELP](#) with output formats

```
# unknown Length: 416
# unknown Number of predicted TMHs:  0
# unknown Exp number of AAs in TMHs: 0.00491
# unknown Exp number, first 60 AAs:  0.00077
# unknown Total prob of N-in:        0.00474
unknown TMHMM2.0          outside          1    416
```

TMHMM posterior probabilities for unknown



transmembrane ——— inside ———

Verify Structural Annotation
BLAST
Domains
Motifs
Protein families
Functional Assignment

# Annotation of Example Protein

**BLAST:** A protein match at Swiss-Prot is 99% identical, with 2 conservative and one non-conservative amino acid substitutions. "RNA-editing ligase TbMP48, mitochondrial precursor" is the Swiss-Prot name for this close protein match.

This mitochondrial precursor of an RNA ligase was identified as a <u>member of a multi-protein complex that catalyzes deletion editing in vitro</u>. It was isolated from an enriched sample of Trypanosoma brucei mitochondria by sequential ion-exchange and gel filtration chromatography, followed by glycerol gradient sedimentation. The protein was not functionally characterized, but was identified as a member of an RNA-editing complex. The complex was shown to have RNA-editing function. (PMID:11134327)

**Domain:** Our sequence contains an RNA ligase, Rnl2 family <u>TIGRFAMs domain</u>, with a very strong match. Members of this TIGRfam family ligate (seal breaks in) RNA.

**Signal sequence:** none

**Targeting Sequence:** It contains a <u>mitochondrial targeting sequence</u>.

**Under the standards of this annotation project, "RNA-editing ligase TbMP48, mitochondrial precursor," is a suitable name.**

Verify Structural Annotation
→
BLAST
→
Domains
→
Motifs
→
Protein families
→
Functional Assignment

# Evidence from homology searching

**Compare sequences of unknown function to those of known function.**

Shared sequence identity <u>may</u> imply shared function:-
- Full-length match with significant identity (>30%)
- Domains and motifs
- Binding sites
- Catalytic sites


But :
- there are occurrences where one amino acid substitution changes the function of an enzyme.
- synonymous or "silent" codon substitutions may result in functional differences.
- Mutations may result in modification or deletion of function.
- all functional assignments made by similarity should be considered tentative until confirmed by experiment.

# Transitive annotation

**Beware!**

**A is like B**

**B is like C**

**C is like D**

**D is NOT like A!**

Take a conservative approach. Err on the side of missing homology rather than stretching weak data.

# Gene Ontology

GO is…

- ❖ a <u>method</u> used to structure biological knowledge using a dynamic controlled vocabulary across organisms.

- ❖ a <u>database</u> containing a shared vocabulary of descriptive terms for the description of the molecular function, biological process and cellular component of gene products.

- ❖ The Gene Ontology Consortium™ is a <u>collaboration</u> among model genome organism databases.

# Topics

- Reasons GO has been developed
- Nuts and bolts of GO
- Tools
- Searching GO
- Assigning terms
- GO Slims

# The Basics

- GO is a controlled vocabulary
- GO has three aspects, or ontologies:
  – Molecular function
  – Biological process
  – Cellular component
- The 3 aspects refer to genes and gene products

# The specificity of GO

There is a limit to how much information can be contained in the name of a protein. For example:

"translation initiation factor 2 subunit"

GO terms assigned to this tell much more:

GO:0003743    (MF)  translation initiation factor activity

GO:0005525    (MF)  GTP binding

GO:0006413    (BP)   translational initiation

GO:0005851    (CC)   eukaryotic translation initiation factor 2B complex

# The Gene Ontology is like a dictionary

**Each concept has:**

- **a name**

- **a definition**

- **an ID number**

**term**: transcription initiation

id: GO:0006352

**definition**: Processes involved in the assembly of the RNA polymerase complex at the promoter region of a DNA template resulting in the subsequent synthesis of RNA from that promoter.

# GO terms

- A GO term, or ID, is attached to every function, process or component

- There are relationships between them

- Relationships are shown by a graph
  - Directed acyclic graph
  - Sometimes called a "tree"

# GO Tools

**GO tools** are available at the GO Consortium:

http://www.geneontology.org/GO.tools.shtml

<u>Developed and maintained by GO:</u>

 AmiGO - Searching through terms and annotations

 OBO-Edit - Editing and viewing the DAG

<u>Many others developed independently</u>, for:

Annotation

Gene expression/microarray data

GO Slims

# AmiGO

# The GO Browser

Search    Advanced Search    BLAST search    Browse    Help

**Filter results**

Filter by ontology

Ontology

| All |
| Biological Process |
| Cellular Component |
| Molecular Function |

Filter Gene Product Counts

Data source

| All |
| CGD |
| dictyBase |
| FlyBase |

Set filters

Remove all filters

**all : all [184843]** 🌈

⊞ ① GO:0008150 : biological_process [139437]

⊞ ① GO:0005575 : cellular_component [122434]

⊞ ① GO:0003674 : molecular_function [137219]

Graphical View
Permalink
Download as XML
Download as flat file

Last updated 2007-01-11

# the Gene Ontology

# AmiGO

**Search**    **Advanced Search**    **BLAST search**    **Browse**    **Help**

**Filter results**

Filter by ontology

Ontology

| All |
| Biological Process |
| Cellular Component |
| Molecular Function |

Filter Gene Product Counts

Data source

| All |
| CGD |
| dictyBase |
| FlyBase |

Set filters

Remove all filters

□ **all : all [184843]** 🔴

  □ ⓘ **GO:0008150 : biological_process [139437]** 🔴

    ⊞ ⓘ GO:0022610 : biological adhesion [1691]

    ⊞ ⓘ GO:0065007 : biological regulation [18316]

    ⊞ ⓘ GO:0009987 : cellular process [81676]

    ⊞ ⓘ GO:0032502 : developmental process [16502]

    ⊞ ⓘ GO:0043062 : extracellular structure organization and biogenesis [313]

    ⊞ ⓘ GO:0040007 : growth [3428]

    ⊞ ⓘ GO:0042592 : homeostatic process [1533]

    ⊞ ⓘ GO:0051179 : localization [20043]

    ⊞ ⓘ GO:0040011 : locomotion [458]

    ⊞ ⓘ GO:0051235 : maintenance of localization [200]

    ⊞ ⓘ GO:0008152 : metabolic process [53763]

    ⊞ ⓘ GO:0051704 : multi-organism process [1640]

    ⊞ ⓘ GO:0032501 : multicellular organismal process [8053]

Graphical View
Permalink
Download as XML
Download as flat file

# the Gene Ontology

# AmiGO

**Search**   **Advanced Search**   **BLAST search**   **Browse**   **Help**

**Filter results**

Filter by ontology
Ontology
| All |
| Biological Process |
| Cellular Component |
| Molecular Function |

Filter Gene Product Counts
Data source
| All |
| CGD |
| dictyBase |
| FlyBase |

Set filters

Remove all filters

Graphical View
Permalink
Download as XML
Download as flat file

⊟ **all : all [184843]** ●

⊟ ① **GO:0008150 : biological_process [139437]** ●

⊞ ① GO:0022610 : biological adhesion [1691]

⊞ ① GO:0065007 : biological regulation [18316]

⊞ ① GO:0009987 : cellular process [81676]

⊟ ① **GO:0032502 : developmental process [16502]** ●

⊞ ① GO:0009838 : abscission [6]

⊞ ① GO:0007571 : age-dependent general metabolic decline [9]

⊞ ① GO:0007568 : aging [435]

⊞ ① GO:0048856 : anatomical structure development [10162]

⊞ ① GO:0048646 : anatomical structure formation [860]

⊞ ① GO:0009653 : anatomical structure morphogenesis [5807]

⊞ ① GO:0048869 : cellular developmental process [6253]

⊞ ① GO:0016265 : death [2387]

⊞ ① GO:0048589 : developmental growth [216]

# GO information to include

Independent of interface, add:

GO ID
Evidence code
Reference
Qualifier

The date is an important part of the annotation .

In Manatee:

# Filling in the GO information

# Assigning a GO term

1) Read the literature, not just the abstract

2) Search for GO terms

3) Record the data

# GO Annotations based on similarity

- Sequence or structure
  - Similarity to GO-annotated gene products
- Domains
- EC numbers
- Pathways
- Protein families

and many more…

http://www.geneontology.org/GO.indices.shtml

# Annotating by similarity

use the evidence code 'ISS'—inferred from sequence or structural similarity.

enter the database ID of the entity used to infer similarity in the 'With' field.

# IEA: Inferred from Electronic Annotation

IEA is used when no curator has checked the annotation to verify its accuracy.

Use when an annotation:

- is based on "hits" in sequence similarity searches, if they have not been reviewed by curators

- is transferred from database records, if not reviewed by curators

- that depend directly on computation or automated transfer of annotations from a database.
  - The actual method used (BLAST search, SwissProt keyword mapping, etc.) doesn't matter.
  - If the method is match-based, a valid database ID *must* be entered in the with column.

# GO Slim

- cut-down versions of the ontologies

- useful summary of GO annotation

- versions of GO Slims available

  - Eukaryotic GO slim

  - Plant GO slim

  - Yeast GO slim

# Points to remember

- GO enables querying across annotations
- The GO Consortium website has documentation and lists available tools
- AmiGO is available online and as downloadable resource
- GO Slims summarize your annotation
- GO annotations are worth the trouble—they enhance the value of research

# MANATEE⁻

- Navigation, inspection & <u>curation</u> of gene products
  - Gene/Gene products
  - GO Assignments

- Available at:
  - http://manatee.sourceforge.net