

Michelle Gwinn August, 2005

Table of Contents (for the most popular topics) topic (page #s)

- 1. Getting started (3-6)
- 2. "Welcome to Manatee" page and links (7-11,21,23,26-28)
- 3. "Genome Summary" page and links (11-20)
 - -Annotation Notebook (15,37)
 - -Genome Calculations (13)
 - -Role Category Breakdown (14)
- 4. "Annotation Tools" page and links (28-38)
 - -Gene List (34-38)
 - -coordinate range (29)
 - -overlaps (30)
 - -InterEvidence (31)
- 5. Gene Curation Page (39-86)
 - -BER section (43-47)
 - -HMM section (55-57)
 - -GO section (71-75,81)
- 6. Gene Ontology (21-22,71-81)
 - -edit Gene Ontology (22)
 - -search Gene Ontology (22,76-80)
 - -Gene Ontology on the Gene Curation Page (71-75,81)
- 7. Genome Properties (23-25,57-60)
- 8. Genome Viewer (26,87-91)
- 9. TIGR role categories (35-36,38,82)
 - -Role notes (38)
 - -TIGR role entry on Gene Curation Page (82)
- 10. Edit starts (90)
- 11. Annotation Checklist (92)

What Manatee Is

- Manatee is a web-based manual annotation tool for accessing and editing annotation data
- Manatee draws information from an underlying database for its displays
- Manatee sends information entered by annotators to the underlying database for storage
- Manatee depends on TIGR's database structure (more on this later)
- Multiple users can access the same database from different computers when Manatee is run on a server (as it is at TIGR)
- Installation questions will be addressed tomorrow by Todd Creasy, our Manatee developer (for those attending the 3-day course)

Getting started with Manatee

- To log into Manatee within TIGR one must have a Sybase account and password.
- For those taking the 3-day course we have established a "training" account. This account is:
 - username = "training"
 - password = "training"
- TIGR employees have their own Sybase accounts and passwords.
- When logging into Manatee, one must enter a user name, a password, and the name of the database on which you wish to work.
- TIGR database names tend to be 3-5 letter codes:
 - during this tutorial and subsequent exercises we will be using the Shewanella oneidensis (formerly Shewanella putrefaciens) database.
 - we will be working with two versions of the Shewanella database:
 - the production database, which stores the published annotation (gsp)
 - a copy of the production database that has been set back to an unannotated state (tgsp)

Finding Manatee (working at TIGR)

GO to manatee.tigr.org and select "Prokaryotic Manatee"



Clicking on "Prokaryotic Manatee" from <u>www.manatee.org</u> takes you to the Manatee Login Page

Manatee Login							
user_name:	training						
password:	****						
database:	gsp						
Submit							

Fill in the fields with the required information.

user name and password:

3-day course: user name and password are both "training" **TIGR employees**: fill in your values

<u>database</u> = "gsp" for the tutorial portion and "tgsp" for the exercise portion (we are now in the tutorial portion).

"Welcome to Manatee"

After logging in to Manatee, you come to the "Welcome to Manatee" page.

☐ This

Here you will find several menu options and a couple search options to choose from.

I will discuss each in more detail in following slides.

NOTE: in the upper right hand corner of every Manatee page will be something like this:

Home Logged into [gsp] as mlgwinn

The "Home" link takes you back to the "Welcome to Manatee" page, from where ever you are within the Manatee tool.

This area also shows you which database you are logged into, and who is logged in. Clicking on the login name will take you back to the login page.

Welcome to Manatee							
OThis is the main	menu page for the Manatee tool. One can access genes directly (with gene's id number or name) or link to a						
	ACCESS LISTINGS						
	Annotation Tools						
	Genome Summary Gene Ontology						
	Genome Properties						
	Genome Viewer						
	 Multi Genome Annotation Tool 						
submit	© ACCESS GENE CURATION PAGE						
reset	▶ gene:						
	© SEARCH GENES BY GENE NAME						
	▶ gene name:						
	C CHANGE ORGANISM DATABASE						
	▶ database:						

The Welcome to Manatee Page "Access Gene Curation Page" option

We will look at the options in the Access Listings section in subsequent slides. First we will look at the 3 options on the bottom of this page:

Access Gene Curation Page:

This option will take you directly to a page containing gene specific information called the "Gene Curation Page" or "GCP" for short. The GCP displays most of what knowledge we have about a given protein - you will be seeing this page in much more detail later. For now just know that you can reach this page by entering either a feat_name or locus id into this box and then clicking "submit". A feat name is an internal identifier given to each gene in a genome, feat_names are not used publically. These are initially assigned by Glimmer and generally are numbered sequentially from the beginning of the DNA sequence given to Glimmer. They have the format ORF#####, where ORF stands for "open reading frame" and ##### is a 5digit zero padded number. (For more on this see the overview document.) Locus ids (loci) are assigned to proteins at the end of the annotation process. They are numbered sequentially from the origin of replication of the genome (if it can be identified). Loci are unique accessions and are used for public release and display of the proteins.

Welcome to Manatee

re

🕒 This is the main menu page for the Manatee tool. One can access genes directly (with gene's id number or name) or link to a

ACCESS LISTINGS

- Annotation Tools
- Genome Summary
- Gene Ontology
- Genome Properties
- Genome Viewer
- Multi Genome Annotation Tool

submit	© ACCESS GENE CURATION PAGE
reset	→ gene:
	© SEARCH GENES BY GENE NAME
	▶ gene name:
	• CHANGE ORGANISM DATABASE
	• database:

The Welcome to Manatee Page "Search Genes By Gene Name" option This is a keyword based search for the common names that have been given to the genes/proteins (we have a tendency to use the terms gene and protein somewhat interchangeably, however, what we are really annotating are the protein translations of the predicted genes.) Whatever keyword you enter will be treated as though it has wildcards flanking it. This means that you will get results that include names containing your keyword as an individual word and names that contain words that contain your keyword. For example, if you search with "kinase"

you could get these: "adenylate kinase" "protein kinase" "sensor histidine kinase"

as well as these: "glutamate 5-kinase" "phosphoenolpyruvate carboxykinase" "ribose-phosphate pyrophosphokinase"

The results will be in the form of a table containing additional information and links to other pages - this table format will be described later.

Welcome to Manatee

🕒 This is the main menu page for the Manatee tool. On e can access genes directly (with gene's id number orname) orlink to a

ACCESS LISTINGS

- Annotation Tools
- Genome Summary
- Gene Ontology
- Genome Properties
- Genome Viewer
- Multi Genome Annotation Tool

submit	© ACCESS GENE CURATION PAGE
reset	• gene:
	© SEARCH GENES BY GENE NAME
	gene name: keyword
	© CHANGE ORGANISM DATABASE
	→ database:

The Welcome to Manatee Page "Change Organism Database" option

To change from one database to another, one does not need to re-login, rather one need only type in the name of the database they wish to go to and click submit.

Welcome to Manatee 🕒 This is the main menu page for the Manatee tool. One can access genes directly (with gene's id number orname) or link to a ACCESS LISTINGS Annotation Tools Genome Summary Gene Ontology Genome Properties Genome Viewer Multi Genome Annotation Tool submit ACCESS GENE CURATION PAGE reset gene: SEARCH GENES BY GENE NAME gene name:

CHANGE ORGANISM DATABASE

• database: tgsp

The Welcome to Manatee Page Options under "Access Listings": "Genome Summary"

The "Annotation Tools" option is one of the most used and will be described in detail in later slides. The "Gene Ontology", "Genome Properties", and "Genome Viewer" sections are accessible here as well as elsewhere within Manatee. (There are many routes to view the various pieces of information within Manatee.) They will be described briefly as links from "Access Listings" and then in more detail as they are viewed from the Gene Curation Page (GCP) and elsewhere.

First we will look in more detail at the options under "Genome Summary"

Clicking on "Genome Summary" takes one to a new page with additional menu options (on next slide).

Welcome to Manatee 🕒 This is the main menu page for the Manatee tool. On e can access genes directly (with gene's id number or name) or link to a ACCESS LISTINGS Annotation Tools Genome Summary Gene Ontology Genome Properties Genome Viewer Multi Genome Annotation Tool submit ACCESS GENE CURATION PAGE reset > gene: SEARCH GENES BY GENE NAME gene name: CHANGE ORGANISM DATABASE database:

The "Genome Summary" page

Genome Sum	mary Page	Home Logged into [gsi] as mlgwinn							
Dhe Genome Summary information page displays specific information concering the selected genome. From here the user can view the following information : ORF counts, Role cateogory information, genes of interest, HMM and paralogous family information, membrane protein information, frameshift information, and annotation progress.									
Home	Annotation Tools	Genor	ne Summary	Gene Assignment Help					
	SUMMARY LISTS Genome Calculations Role Category Breakdown Annotation Notebook Brainet Administration	Clicking on the item in the list of options							
	 Frameshift Status Annotation Progress Report By Interpro Domain Genome Properties 	more options. Following slides will describe each of these.							
submit	 ○ ATTRIBUTES > attribute: MW or ○ EVIDENCE > evidence: HMM2 or ○ BABALOCOUS PROTEINS 	der by: 💽	descending <u></u>	e 🔽					
reset	 View by: number of family member MEMBRANE PROTEINS sort by: Y-score signal peptide cutoffs: 0.36 Y-score lower limit 5-4 S-mean lower limit transmembrane regions: lower limit constraints: genes with OMP signals genes with lipid attachmet 	ers	These are tool the data based annotation. For describe the u these.	Is that allow one to view d on various types of ollowing slides will se and output of each of					

Questions? Comments? Please feel free to send us feedback!

Links from the Genome Summary Page: "Genome Calculations"

feature name	feature count	feature type
▶ Open Reading Frame	4930	ORF
▶ ribosome binding site	4444	RBS
▶ rho-independent terminator	846	TERM
▶ Reject ORFS	169	RORF
▶ transfer RNA	101	tRNA
▶ ribosomal RNA	27	rRNA
▶ structural RNA	3	sRNA
▶ Bacteriophage	3	PHAGE

start sites	number	percent
► ATG:	3933 (3235)	79.8% (85.9%)
► GTG:	541 (350)	11.0% (9.3%)
▶ TTG:	453 (182)	9.2% (4.8%)
• OTHER:	1 (1)	0.0% (0.0%)

This page shows the various calculable and countable features of the genome. This information is newly generated each time the page is accessed so that all information is current.

13

Numbers in parentheses do not include hypothetical proteins

chromosome' Information Table								
▶ sequence id:	7974							
▶ type:	chromosome							
▶ molecule length:	4969803 bp							
► GC content:	46%							
▶ base frequencies:	(A) (C) (G) (T) 27.0% 23.0% 23.0% 27.0%							
▶ funny characters:	R Y 2 6							
► ORF count:	4758							
▶ average gene length:	895 nt							
▶ percent coding:	85.8%							
percent coding OR tRNA, rRNA, or repeat:	90.0%							
▶ skew table								

Links from the Genome Summary Page: "Role Category Breakdown"

This page shows a summary of the genes found in various broad categories based on TIGR roles and then a breakdown by TIGR sub category. Each blue role id number or "main" is a link to a table containing a list of all the genes in that category.

► ORF Summary										
Total ORFs:	4930	100.0 %								
assigned function	2521	51.1 %								
conserved hypothetical	871	17.7 %								
unknown function	378	7.7 %								
hypothetical proteins	1162	23.6 %								

▶ Role Breakdown								
role id	name	number	complete	%				
main	Unclassified	2	0	0.04%				
185	Role category not yet assigned	2	0	0.04%				
main	Amino acid biosynthesis	91	0	1.85%				
70	Aromatic amino acid family	17	0	0.34%				
71	Aspartate family	24	0	0.49%				
73	Glutamate family	21	0	0.43%				
74	Pyruvate family	13	0	0.26%				
75	Serine family	8	0	0.16%				
161	Histidine family	8	0	0.16%				
69	Other	0	0	0.00%				
main	Purines, pyrimidines, nucleosides, and nucleotides	63	0	1.28%				
123	2'-Deoxyribonucleotide metabolism	8	0	0.16%				
124	Nucleotide and nucleoside interconversions	11	0	0.22%				

14

Links from the Genome Summary Page: "Annotation Notebook" - the annotation notebook is a set of text fields associated with each TIGR role category. These are used for annotators to store information about the annotation which they feel the PIs of the project should know for purposes of writing the manuscript, generally they consist of items of particular biological interest, often involving the presence or absence of particular pathways, genes, gene order, etc. These entries are entered and edited with the "Edit Annotation Notebook" page, linked from the gene list, see page 33, 36 of this tutorial.

Shewanella oneidensis MR-1 | Annotation Notebook

Logged into [gsp] as mlgwinn

D The annotation notebook is used by annotators to note genes of interest for the project leader. The notebook is organized by role category. The project leader can view the page, and view the genes at their leisure.

Biosynthesis of cofactors, prosthetic groups, and carriers

83 : Pantothenate and coenzyme A

Appears to lack panD which I searched for with the E.coli sequence. No matches using blastp or tblastn. RTD

84 : Pyridoxine

Appears to lack serC which I searched for with the Caulobacter sequence. No high-scoring matches using blastp or tblastn. RTD

Cellular processes

96 : Detoxification

arsenate reductase operon quirky. Sequences diverge from typical prokaryotic arsB and arsC. Order of these two genes in the operon is reversed as well. The arsenate reductase repressor arsR is well conserved. There may also be a stand alone version of arsC present as well.

Protein fate

97 : Protein and peptide secretion and trafficking

general secretion pathway genes are split into two operons (gspAB) and (gspCDEFGHIJKLMN). Not sure how typical this is.

Biosynthesis of cofactors, prosthetic groups, and carriers

162 : Thiamine

Appear to lack some components of the the pathway for thiamine biosynthesis. RTD

Other links from the Genome Summary Page:

"**Project Administration**" - Clicking this link takes one to a page that displays the administrative information for the project: things like PI, grant #, etc.

"Frameshift Status" - Currently this tool is not available for people running Manatee locally outside of TIGR. For TIGR users, there is (will be) a separate section of the tutorial governing all things involving Frameshifts (this part of our SOP is currently undergoing change.) In brief, this link displays a page listing all of the genes in the genome which needed to be reviewed for the presence of a frameshift or in-frame stop codon as well as the status of each.

"Annotation Progress Report" - This links to a page that lists all of the processes that must be carried out during the annotation of a genome and provides fields in which to enter when each process was done and who did it. There is also a link to a page listing all the TIGR mainrole categories and fields for individual annotators to sign up for each category.

"by InterPro Domain" - This links to a list of genes according to membership in an InterPro domain.

"Genome Properties" - another link to this tool set, will be described in detail elsewhere.

Searches on the Genome Summary page: "Attributes"

• ATTRIBUTES	
▶ attribute: MW ▼ order by: MW ▼ descending ▼	

One can choose to view genes based on one of several "attributes" they might have. Here I have shown a selection for "MW" which stands for molecular weight. Once you choose and attribute to search by, you can then choose various ordering display options. The above choices will show the proteins in the genome according to calculated molecular weight with the heaviest ones first. (see below)



This is just the top of a very long list containing all of the proteins in the genome. One can click on any of the blue gene id links and get go to the Gene Curation Page (GCP) for the gene. (The GCP will be described in detail shortly.)

One can jump to different pages in the list by clicking on the blue numbers in the boxes above the list.

One can change the order of the list by clicking in the arrows in blue circles.

Searches on the Genome Summary page: "Evidence"



One can choose to view the genes based on one of several types of clustering evidence that has been found for them, some as the result of InterPro searches and some as a result of separate searches we perform. Here, I have selected "HMM2" (which will include both the TIGRFAM and Pfam HMM sets) and I will view the output ordered by the number of hits in the genome, the HMMs with the most hits will be listed first.

Shewanella oneidensis MR-1						Evid	lence	Displ	lay				
D													
		-											
PAG	SE	1	2	3	4	5	6	7	8	9	10	11	12
HMM	2												
count	aco	ession					H	MM na	me				
85	PF	00005	ABC	c transp	orter								
82	TIG	R01199	Heli	x-turn-l	nelix do	omain, f	is-type						
71	PF	00072	Resp	onse re	gulator	receive	er doma	in					
69	PF	02653	Bran	ched-c	hain an	nino aci	d transp	ort syst	tem / pe	ermease	compon	ent	
62	PF	00665	Integ	grase co	ore dom	ain							
62	PF	00672	HAN	IP dom	ain								
59	PF	00037	4Fe-	4Fe-4S binding domain									
57	PF	00271	Heli	case co	nserved	C-term	inal do	main					
56	PF	00872	trans	transposase, Mutator family									
56	PF	02796	Heli	Helix-turn-helix domain of resolvase									
54	PF	03466	Lys	LysR substrate binding domain									
53	PF	00486	Tran	scriptio	onal reg	ulatory	protein	, C tern	ninal				
53	PF	00990	GGE	DEF dor	nain								

By clicking on the blue numbers in each row, one will get a list of genes that hit that HMM. Clicking on the blue accession number will take one to an info page for the HMM in question.

One can reorder the list by count or accession by clicking on the blue column headers.

Numbered boxes at the top will take one to a desired page in the output.

Searches on the Genome Summary page: "Paralogous Families"

• PARALOGOUS PROTEINS

▶ view by: number of family members -

One can choose to view the genes based on membership in paralgous families, ordering either by number of family members or by family name.

Silicibacter pomeroyi DSS-3 Paralogous Families

Le The Paralogs page comes in two flavors. If launched from the Gene Summary page, it displays a summary (paralog page is launched from itself (i.e. from the paralog skim flavor) it displays an ORF list of all members for the selected pa

total families: 762 > total proteins in families: 2647

proteins	family name	description
114	fam_PF00005	ABC transporter
82	fam_PF00528	Binding-protein-dependent transport systems inner membrane component
67	fam_PF00126	transcriptional regulator, LysR family
60	fam_PF03466	LysR substrate binding domain
44	fam_TIGR01409	Tat (twin-arginine translocation) pathway signal sequence
37	fam_PF00106	oxidoreductase, short chain dehydrogenase/reductase family
34	fam_PF00072	Response regulator receiver domain
33	fam_PF00892	Integral membrane protein DUF6
31	fam_PF02653	Branched-chain amino acid transport system / permease component
30	fam_PF00392	transcriptional regulator, gntR family
29	fam_PF00583	acetyltransferase, GNAT family
27	fam_PF02518	ATPase, histidine kinase-, DNA gyrase B-, and HSP90-like domain protein
25	fam_PF06808	TRAP transporter, DctM-like membrane protein
25	fam_PF00165	transcriptional regulator, AraC family

-Paralogous families are built by first searching all of the proteins within a genome against themselves and against the HMM db. If a paralogous family matches an HMM the family will be named based on the HMM. Then further searches are done to group the proteins based on regions of sequence that did not match an HMM. Those families are given numerical names and do not have descriptions.

-Output shows you the number of members in each family, the name of the family, and a description of the family (if the family is based on an HMM).

-You can view a list of the proteins in each family be clicking on the family name. You can view information about the HMM on which the family is based by clicking on the description.

Searches on the Genome Summary page: "Membrane proteins"



Output shows a table of the genes with the chosen parameters. You can reorder them using the pull-down menu and the "sort" button. The table displays all of the parameters available for each protein. Clicking on the blue gene id takes you to the Gene Curation Page (GCP) for the gene.

One can choose to view the proteins based on predicted location in a membrane. You can choose particular SignalP cutoff values, number of predicted transmembrane regions, proteins that have an OMP signal, or lipid attachment site. You can also sort the output by several different options.

Silicibacter po	omeroyi DSS-3	Membrane Protein	S				Home Logged into [gsi] as mlgwinn
D This page displ sequence, and lip	ays an ORF list bas oprotein sequence	ed on the results of fourm e. The user controls the ou	embran Itput by	e protein ar specifying	nalyses:s on e to fo	ignal ur of t	peptide sequence, membrane spanning regions, outer membrane protein hese criteria in any combination as well as the sorting key for the results.
sort by:	gene id	_	Sort				
gene id	<mark>y-score</mark> (>=0.36)	<mark>s-mean</mark> (>=0.54)	site	TM (3/5)	ОМР	LP	gene name
ORF00005	0.372	0.611	91	3			conserved hypothetical protein
ORF00119	0.557	0.739	31	4			peptide/opine/nickel uptake family ABC transporter, permease protein
ORF00160	0.546	0.723	53	3			conserved hypothetical protein
ORF00234	0.692	0.835	25	5			flagellar biosynthetic protein FliP
ORF00416	0.407	0.825	15	5			NnrU family protein
ORF00422	0.387	0.568	49	3			succinate dehydrogenase, hydrophobic membrane anchor protein
ORF00522	0.411	0.613	24	3			hypothetical protein
ORF00617	0.478	0.628	22	4			sterol desaturase, homolog
ORF00725	0.434	0.732	28	3			conserved hypothetical protein
ORF01154	0.687	0.576	90	5			membrane protein, putative
ORF01297	0.434	0.544	35	4			membrane protein, putative
ORF01321	0.654	0.724	33	4			hypothetical protein
ORF01507	0.396	0.585	41	5			oligopeptide ABC transporter, permease protein
ORF01581	0.463	0.610	29	3			hypothetical protein
ORF01589	0.370	0.547	29	3			conserved hypothetical protein
ORF01645	0.565	0.688	57	4		٥	lipoprotein, putative

The Welcome to Manatee Page Options under "Access Listings": "Gene Ontology"

This link will open a page that offers options for using the Gene Ontology (GO) system.

(For more information on the Gene Ontology system, see the Annotation Overview document, or the Gene Ontology web site, <u>www.geneontology.org</u>)

In brief, the GO offers a controlled vocabulary for the description of aspects of gene products. Currently, TIGR assigns both TIGR role categories and GO terms to all of our genes. Manatee has many built in features for the suggestion and entry of GO terms and associated information. These features will be detailed in later slides. The next slide shows a brief description of the links available here and of the "edit GO" options.

When Manatee refers to "editing" GO, we mean the creation of "TI" or TIGR terms. These are temporary terms created for use in-house at TIGR until corresponding terms are created at GO. When a need for a new term is found, we (usually Michelle) submits a request to the GO via their SourceForge tracking site that the new term be created. If a TIGR annotator needs the new term right away, they can create a TI term to use within our db. Later, when the official GO term is made, the TI term id will be replaced with the new GO term id.

Welcome to Manatee 🕒 This is the main menu page for the Manatee tool. One can access genes directly (with gene's id number or name) or link to a ACCESS LISTINGS Annotation Tools Genome Summary Gene Ontology Genome Properties Genome Viewer Multi Genome Annotation Tool submit O ACCESS GENE CURATION PAGE reset > gene: SEARCH GENES BY GENE NAME gene name: CHANGE ORGANISM DATABASE database:

The Welcome to Manatee page, links from Access Listings: "Gene Ontology"



The Welcome to Manatee Page Options under "Access Listings": "Genome Properties"

The Genome Properties system allows one to view annotation from the context of the whole genome. It predicts and/or captures information on the presence/absence of pathways, cellular structures and other features of the organism. (see overview for more details) Clicking on the Genome Properties link from the "Welcome to Manatee" page displays a table of all of the properties and their states for the organism you are working on. The state is "yes" if the property is present, "no" if the property is absent, and will have other intermediate values such as "some evidence" or "not supported" depending on the amount of evidence for a given property. Details on what is known about each property for the genome you are working on can be obtained by clicking on the blue property name. (see next slide)

Welcome to Manatee



The Welcome to Manatee page, links from Access Listings: "Genome Properties"

update genome properties										
property			state							
2-aminoethylphosphonate	degradation		- NULL	-						
curator comment	Update the status of information about a property in this organized	of or anism	prediction	S	earch for a	property in this	genome.			
assignby				360	an genome properties					
submit				Sear	ch By: Name 💌		submit reset			
Shewanella oneidensis MR-1 Genome Properties Home Logged into [gsp] as mlgwinn										
Lindate Property Search	Property									
property	state	value	assignby	prediction	experiment	curator comment	auto comment			
2-aminoethylphosphonate degradation	nonefound	0	RULES	1						
4-hydroxyphenylacetate degradation	notsupported	0.1	RULES	1						
5-aminolevulinate biosynthesis	YES		RULES	1						
5-aminolevulinate biosynthesis (glutamate pathway)	YES	1	HYBRID	1						
5-aminolevulinate biosynthesis from succinyl-CoA and glycine	NO		RULES	1			none found AND glutamate pathway present			
acetyl-CoAcarbexylase.complex	some evidence	0.25	RULES	1						

Click on the blue name of a property to learn more about the steps/requirements for the property and to see background information and references regarding the property. You can also see the genes the are involved in the property in the context of their neighbors in the genome. These pages will be shown in detail later in the tutorial but are quickly shown on the next slide.

24

Genome Property information Page (in brief, more detail will be shown later in the tutorial)

Property Definition	ı da saya da s		
property:	arginine biosynthesis from omithine, carbamoyl-p and aspartate	state:	YES
property type:	PATHWAY	value:	1
role id:	73 [update]	assignby:	RULES
GO ids:	GO:0042450 [view add update] arginine biosynthesis via omithine	date:	Apr 8 2003 11:09AM
description:	The arginine biosynthesis pathway is a three step process and a part of the urea cyc omithine carbamoyltransferases (OTCase) carries out the reaction anabolically in arg some cases, it carries out catabolic reactions. Most OTCases are homotrimers, but th organized into dodecamers built from four trimers in at least two species; the catabol Pseudomonas aeruginosa is allosterically regulated, while OTCase of the extreme th furiosus shows both allostery and thermophily.[1] The third step of the pathway yield argininosuccinate lyase, and the amino acid can be cleaved by arginase, yielding ure omithine. L-arginine also can be utilized in the creatine biosynthesis.	le.[2] The first jinine biosynth he homotrimers lic OTCase of hermophile Py ds L-arginine b ea and reconsti	enzyme, esis, but in sare rococcus y tuting

Information on the property.

Pro	Property Steps												
RE	QUI	RED		ornit hin e	carbam oyitr	ansferas	e (1)						
A	С	GC	gene id	gene name	evidence	role id	gene symbol	EC number	OP				
0	•		ORF02078	omithine carbamoyltransferase	TIGR00658	73	argF	2.1.3.3					
RE	QVII	RED		argininosuccinate synthase (2)									
A	С	GC	gene id	gene name	evidence	role id	gene symbol	EC number	OP				
0	•		ORF02077	arginin osuccinate synthase	TIGR00032	73	argG	6.3.4.5					
RE	QUI	RED		arginir	nosuccinate	e lyase (3	i)						
A	С	GC	gene id	geneid genename evidence roleid genesymbol EC number OP									
0	0		ORF02076	argininosuccinate lyase	TIGR00838	73	argH	4.3.2.1					

Information on the genes identified to be a part of the property.

25

The Welcome to Manatee page, links from Access Listings: "Genome Viewer" Genome Viewer is a tool which allows one to view the genes in context with their neighboring genes in the genome. It displays a graphic showing the 6-frame translation of a region of DNA sequence, where each horizontal bar is a different frame. Arrows representing the genes are color coded according to TIGR mainrole assignment. There are many viewing and editing options available from this page. These will be discussed in detail later in the tutorial.

	- Г	Refresh XMI	L Search	h Asmbl	Id: 7974	Database:	gsp c	latabase:	asmbl	id:	submit rese
ne to Manatee		L				L					
the main menu page for the Manatee tool. One can access genes directly (with gene's id number or name) or link i	(08	feat name	end5	end3	role id	ec num	gene syn	<u>comp</u>	<u>plete</u> <u>com</u>	name	
ACCESS LISTINGS		ORF02394	954	2504	142			com	plete prot	on/peptide sy	mporter family p
	·						prop_acc	state	prop	erty	
Annotation Tools											
Gene Ontology					1						
Genome Properties		Six Frame C	options (on clicks)	six frame	Man		Gene	Options (on	feat_name clie	ks)	Delete
→ Genome Viewer	+	View Sequence	Blast	Insert Gene	Sequence	C ORF (e Ec	Gen	end Blast Ne C ORF	C Genes	Gene C
Multi Genome Annotation Tool					Lan.	-					361
IT CACCESS CENE CUDATION DACE											
ACCESS GENE CORATION PAGE		0.0Kb	1.0Hb		2.0kb	3.0нь		4.0kb	5.0Kb	6.0kb	7.0kb
→ gene:											
C SEARCH CENES BY CENE NAME											
SEARCH GERES DT GERE NAME	·	L III A LLAD MANA	ki yang dikini	dere all is in min 1 at al	alli <mark>l a</mark> II d ila II I .					hall ddd a ca	d II la tatadiita — ul
> gene name:	1							100 11 00 10 . 111 5 1		<mark>nah da dha til a d</mark> h	
C CHANGE OPGANISM DATABASE		يت البينية الالتانية	<u></u>	0RF 023	94 1 1 1 1 1 1 1		lide de la li l		l III an that dhailt a stada		All is the second
CHARGE ORGANISH DATADASE	·	ORF	02395							ORF 02	390
→ database:					-	ORF	02393				
	-	LHL L H d I H								-	
						While Islam Charles	la collad Stiche di	II	0RF 02392	ORF 02391	
								>			
			Search								
			Coordina	te:						Search	1
			Lower Co	ordinate:		U	pper Coord	linate:		Search	
			feat_nam	e:				,		Search	1

The Welcome to Manatee page, links from Access Listings: "Multi Genome Annotation Tool" (MGAT)

The MGAT tool allows the annotation of orthologous genes from several genomes at one time. It is linked into Manatee at several points. MGAT is still undergoing development and is not currently available for public use. A separate tutorial for this tool is under construction.

Welcome to Manatee																		
Welcome to	o Manatee		welcome to MGAT									Ident Information	1					
Øthisisthem	win menu paga for the Manafes feel. One can accer	as genes desclip/juith general descender ar room of or link to a					AllVAll a	ign mana	atee	com_name			ec_num	gene_sym	role_id	mainrole	db	asmbl_id
	ACCESS LISTINGS		Access COGS By Role Cate	gories			BMA3090	ORF	704141	thiamine biosynth	esis protein ThiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gbm	772
	Genome Summary Gene Ontology Genome Properties		G Single Role Category:				GHNORE	3126 ORF	703083	thiamin biosynthe	esis protein thiC			THIC	126	Purines, pyrimidines, nucleosides, and nucleotides	ghn	1063
	Genome Viewer Multi Genome Annotation	Tool	Unclassified				AFE2099	ORF	7 <u>03389</u>	thiamine biosynth	esis protein ThiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gtf	10431
submit reset	ACCESS GENE CURAT gene:	FION PAGE		Purines, pyrimidines, nucleosides, and nuc Fatty acid and phospholipid metabolism	cleotides		<u>CC2029</u>	ORF	F06212	thiamine biosynth	esis protein ThiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gcc	12582
	SEARCH GENES BY G	ENE NAME	🖲 Main Role Category:	Central interme diary metabolism Energy metabolism	us, anu camera		NT02NM	2427 ORF	703627	thiamine biosynth serogroup a;} ^/^ protein nma0397	esis protein thic. { pirk81956k81956 [imported] - neiss	sis protein thic. {neisseria meningitidis irk81956k81956 thiamin biosynthesis imported] - neisseria meningitidis			185	Unclassified	ntnm02	1
	egene name:	DATABASE		Transport and binding proteins DNA metabolism Transcription			NMB2040	ORF	7 <u>01381</u>	thiamine biosynth	esis protein ThiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gnm	834
	• database:			Protein synthesis Protein fate		•	MXAN42	35 ORF	701813	thiamine biosynth	esis protein ThiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gmx	581
			Submit				PSPPH05	35 ORF	700554	thiamin biosynthe	esis protein ThiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gphas	340
Г							NT01PS1	591 ORF	702811	thiamin biosynthe	esis protein ThiC			thiC	126	Purines, pyrimidines, nucleosides, and nucleotides	ntps01	2
	MGAT Searc	ch information					PSPTO49	76 ORF	704688	thiamin biosynthe	esis protein ThiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gps	5676
	database :						PF5841 ORF08667 thiamin biosynthesis protein ThiC					162	Biosynthesis of cofactors, prosthetic groups, and carriers	gpf	3338			
		Biogun	thesis of col	factors prost			<u>PP4922</u>	ORF	7 <u>01030</u>	thiamin biosynthe	esis protein ThiC		thiC	thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gpp	13541
		DIUSYII	unesis of col	lactors, prosti	ieuc	: gro	GEBCOR	F0191 ORF	7 <u>00186</u>	Thiamine biosynt	hesis protein thiC			thiC	162	Biosynthesis of cofactors, prosthetic groups, and carriers	gebc	1186
				Genome A	Annota	tion												
	gene		<u>com_name</u>		SP Fams	Cogs	JFams	matrix	Tig	grFam	PFAM	start_ec						
	ORF05157	thiamin biosynth	hesis protein ThiC					X	TI	<u>GR00190</u>		No						
	ORF05164	thiH protein						X	TIC	<u>GR02351</u>		No						
	ORF03392	thiH protein, pu	tative					X				No						
	ORF01112 ApbE family protein, putative								PF02424	No				07				
	ORF05159 phosphomethylpyrimidine kinase/thiamin-phosphate pyrophosphorylase, putative								TIC	<u>GR00693</u>		No				21		

The Welcome to Manatee Page: Options under "Annotation Tools"



"Annotation Tools": "Coordinate Range"

Shewanella or	neidensis MR-I Annotation Tools		Home Logged in	to [gsp] as <u>mlgwinn</u>							
The ann_tools.	cgiscript generates the Annotator Tools webpage, whi roperties of the genome and determining the progress	ch is the entry point for accessing the Submit web made in the Annotation of the genome of interes	page for all ORFs in a genome,	as well a resource for							
Home	Annotation Tools	Genome Summary	Gene Assignr	nent Help							
	• ACCESS GENE LISTS										
	molecule: All molecules										
	all genes, ordered by role	category (excludes hypothetical	proteins)								
	man role category Unclassing single role category [role_id										
	• ACCESS GENE CURATIO	N PAGE									
_	> gene:			_							
	C ACCESS GENES BY COOP	RDINATE RANGE	-								
L	+ end5: end3:										
	OVERLAP ANALYSIS										
submit	location: //usr/local/annotation/GSP/overlap_analysis/gsp.overlap										
reset	INTEREVIDENCE REGION ANALYSIS										
	location: /usr/local/annotation/GS	P/intergenic_regions/gsp.REPC									
	C CUSTOM QUERY										
	OTHER TOOLS										
	Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search										

Input a coordinate range and you will get a list of genes whose coordinates fall anywhere in that range. List of all genes found between 10000 - 20000

A	С	gene id	locus	end5	end3	role id	gene name	gene symbol	ec
		ORF02375	SO0017	22090	18941	156	conserved hypothetical protein		
		ORF02378	SO0016	18279	18854	132	DNA-3-methyladenine glycosidase I	tag	3.2.2.20
		ORF02379	SO0015	18161	17256	137	glycyl-tRNA synthetase, alpha subunit	glyQ	6.1.1.14
		ORF02381	SO0014	17246	15180	137	glycyl-tRNA synthetase, beta subunit	glyS	6.1.1.14
		ORF02382	SO0013	14311	15111		hypothetical protein		
		ORF02383	SO0012	13791	13129	96 102,	glutathione S-transferase family protein		
		ORF02385	SO0011	10638	13055	132	DNA gyrase, B subunit	gyrB	5.99.1.3
		ORF02386	SO0010	9539	10621	132	DNA replication and repair protein RecF	recF	
		ORFA00005	SOA0024	20332	19523	154	ISSo1, transposase OrfB		
		ORFA00006	SOA0023	19154	19453	94 186,	proteic killer suppressor protein	higA	
		ORFA00007	SOA0022	18774	19079	94 186,	proteic killer active protein	higB	
		ORFA00008	SOA0021	18235	18462	154 270,	ISSo1, transposase OrfB, truncation		
		ORFA00009	SOA0020	17414	18154	154 270,	transposase family protein, truncation		
		ORFA00011	SOA0019	16733	17290	132 154,	TnSon1, resolvase		
		ORFA00012	SOA0018	16362	16739	154 156,	TnSon1, conserved hypothetical protein		
		ORFA00013	SOA0017	16075	16365	703	TnSon1, nucleotidyltransferase domain protein		
		ORFA00014	SOA0016	15911	12945	154	TnSon1, transposase		
		ORFA00015	SOA0015	12878	12732		hypothetical protein		
		ORFA00016	SOA0014	12332	12427		hypothetical protein		
		ORFA00017	SOA0013	11739	11335	132	umuD protein	umuD	3.4
		ORFA00019	SOA0012	11334	10078	132	umuC protein	umuC	

"Annotation Tools": "Overlap Analysis"

a oneidei	nsis MR-I	Annotation Tools		Home Logged into [gsp] as mlgwinn	We
ools.cgiscr ral properti	ipt gen erate: es of the ger	s the Anniotator Tools web page, v nome and determining the progre	which is the entry point for accessing the Submit v ss made in the Annotation of the genome of inter	reb page for all ORFs in a genome, as well a resource for est.	nrok
ne		Annotation Tools	Genome Summary	Gene Assignment Help	
	ACC molecul ^a al ^c m ^c si ^c ACCI ^g gene: ^c ACCI ^g gene: ^c ACCI ^c end5: ^c ACCI ^c OVER ^l location: ^c INTEI ^l location: ^c CUST ^c CUST ^c	ESS GENE LISTS le: [All molecules l genes, ordered by ro ain role category [Unc ngle role category [role ESS GENE CURATI ESS GENES BY COO end3: [KLAP ANALYSIS [/usr/local/annotation/G REVIDENCE REGIO [/usr/local/annotation/G OM QUERY		l proteins)	calls the c "hyp over on o This back the f gene and
				Shewanella oneidensis MR-1	Overlap

work on the premise that genes do not generally overlap in arvotic genomes. We look for overlapping genes predicted by mer and where we can, remove genes suspected of being false by Glimmer. Often overlap between two genes can be resolved by curation of the start site of one or both, or by the removal of a othetical protein" (one that has no similarity to anything) when it rlaps a protein with very clear similarity to other proteins. For more verlap analysis see the Annotation Overview document. display shows the pairs of overlapping genes as indicated by the ground color shifts from blue to white to blue to white. Clicking on feature id number takes you to a Gene Curation Page (GCP) for that e. Also displayed are the percent of overlap, name of the protein, notes from Glimmer regarding the protein in question.

Shewanella d	oneidensis MR	-1 Overlap Display	Home Logged into [gsp] as mlgwinn
This page dis	plays a list of overl	apping ORFs that have been identified as candidate genes by GLIMMER(the gene find	ing program).
feature	% overlap	name	glimmer
ORFA00023	3.47	hypothetical protein	742 [-3 L= 759] [DelayedBy ORFA00024 L=45]
ORFA00022	8.77	hypothetical protein	229[-2 L= 282]
ORFA00078	6.41	ISSo8, transposase, degenerate	423[-1 L= 543]
ORFA00082	4.38	transposase, IS3 family, degenerate	254 [+3 L= 336]
ORFA00096	2.60	ISSo12, transposase	672[+1L=921]
ORFA00095	19.05	hypothetical protein	97 [-2 L= 123]
ORFA00204	3.92	conserved hypothetical protein	99[-3 L= 108]
ORFA00203	4.23	hypothetical protein	812[-1 L=1128]
ORF00040	6.46	TonB-dependent receptor C-terminal region domain protein	198[-1 L= 477][DelayedByORFA00041 L= 18]
ORF00041	14.55	hypothetical protein	197[-2 L= 210]
ORF00158	5.00	hypothetical protein	215[-3 L= 537]
ORF00157	3.33	ISSo1, transposase Orf B	693[+1L=807]
ORF00177	5.23	hypothetical protein	297[+2L= 360]
ORF00176	3.15	conserved hypothetical protein	396 [+1 L= 645]

OTHER TOOLS

Shewanei The ann_ ocating gen

submit

reset

Hor

- Data Consistency Checks
- Frameshift Reports
- Hypothetical Protein List Annotation Status
- Phage Region Viewer
- PubMed Organism Search

"Annotation Tools": "Interevidence Analysis"

Glimmer is known to sometimes miss identifying a few real genes. This is especially true for areas of the genome that have been laterally transfered. To find genes Glimmer might have missed, we run an analysis called "interevidence". This tool

gsp_7971_R3_orf1

gsp 7971 R2 orf1

gsp_7970_R8_orf1

gsp_7970_R23_orf1

2010

924

3449

30881

330

320

1345

1424

4285

30648

takes the nucleotide sequence between genes, the sequence of hypothetical proteins (those that have similarity to nothing), and any regions of proteins that have similarity to nothing, does a 6 frame translation, and then searches those translations against niaa (our in-house protein db). Any possible areas of similarity are then reviewed by annotators and missed genes are entered into the db.



Shewanella oneiden	ris MR-1	Interevide	ence Regions Display		Home	Logged into [gs	p] as <u>mlgwinn</u>				
D This page displays the output from the Interevidence Region analysis program. Areas of the genome with no genes are searched for homology in case genes were misssed by GLIMMER.											
IE gene	asmbl_ic	I IE end5	IE end3	Hit end5	Hit end3	Proposed ORF end5	Proposed ORF end3				
gsp_7971_R2833_orf1	3	383	:								
gsp_7971_R2833_orf2	3	167	:								
gsp_7971_R2833_orf3	3	383	:								
gsp_7971_R2832_orf2	293	1084	:								
gsp_7971_R2832_orf1	347	1084	:								
gsp_7971_R3_orf2	1271	1441									

ORF00003: hypothetical protein

ORFA00009: transposase family protein, truncation

INVAVPSDDDTTNEPDDFRPPCPCCGGRMIVIEVFERWRQ

370

380

360

<u> </u>	
58.2/72.8% over 54aa	Mesorhizobium loti
GPl <u>14026131</u> ldbj transposase <u>Insert characterized</u>	
<pre>gsp_7971_R2833_orf1(1.3 - 54.3 of 127 aa) GP 14026131 dbj BAB52729.1 AP003009(288 - 341 of %Match = 8.5 %Identity = 58.2 %Similarity = 72.7 Matches = 32 Mismatches = 13 Conservative Sub.s Gaps = 2 InDels = 6 Frame Shifts = 0 Primary Frame = 2 [0, 53, 0]</pre>	397) transposase {Mesorhizobium loti} = 8

gagttccaggtcgctttgttggtggcttagagggccagtctgactaccagtttt catcactccttcggtccggtttgtatcgggtcccgtatcagtgaggctgaaaactaccacacgtgagagaaaaa tgtgga ctttttqaaqacaactaqttcccaaattactaacctqactcaccqtcqtcacqqcaqtqacac-55.7 -45.7 -35.7 -15.7 -25.7 -5.7 14.3 TDFTHPVASVLASVSSRRMVVCLDLS*CGLSAA*VDLS*SVSEGGPLSEYH*AMOSTOPKACLCYRDGOOHK JT.SCD : RRAFLRHLAPVRRKRWVVYAKAPFAGPEAVLAYLSRYTHRVAISNSRLIRLDESGVTFRYKDYRRD 240 250 260 270 280 290 300 gtaccttccgtcagcacgactgtcgagtccaacgcacccccaccggtgactgaagtttttcctctgctctagtaacctgg attggatta<mark>ttcagttgtg</mark>atgttcacgggaatctttgatgacatttcgcttaaagtgcgcagatgatattgttgcactt gtccttagt<mark>ggc</mark>agg<mark>gagatttttcttcccagggccggaatggggaggttg</mark>tggct<mark>g</mark>gatgacgatt<mark>g</mark>gtt 64.3 74.3 44.3 84.3 34.3 54.3 GFLANACRRKKLALILRQLRKPQVVLASPLVKKDCLWSCPQCQLGHLQFIGLIRPQSVV

350

340

"Annotation Tools": "Other Tools" section

Shewanella on	eidensis MR-I Annotation Tools		Home Logged i	nto [gsp] as <u>mlgwinn</u>						
The ann_tools.cgi script generates the Annotator Tools webpage, which is the entry point for accessing the Submit webpage for all ORFs in a genome, as well a resource for locating general properties of the genome and determining the progress made in the Annotation of the genome of interest.										
Home	Home Annotation Tools Genome Summary Gene Assig									
	• ACCESS GENE LISTS			Data cons possible e						
	Molecule: All molecules	example,								
	all genes, ordered by role	different I								
	C main role category Uncl	consistent								
	Single role category role_	id		Frameshi						
	• ACCESS GENE CURATIO	ON PAGE		was descr						
	> gene:			basically						
	• ACCESS GENES BY COO		Hypothet							
	▶ end5: end3:			with insuff						
	OVERLAP ANALYSIS Iocation: //usr/local/annotation/GSP/overlap_analysis/gsp.overlap									
submit										
Teset	• INTEREVIDENCE REGION	any BFR (
	▶ location: //usr/local/annotation/GS									
	© CUSTOM QUERY									
		list of role								
		the work i								
	OTHER TOOLS			Phage Re						
	Data Consistency Checks Erameshift Reports	regions in								
	 Hypothetical Protein List Annotation Status Phage Region Viewer 									
	PubMed Organism Search									

Data consistency checks: Clicking this generates a list of possible errors or consistency problems in the annotation. For example, if two proteins have the same common name but different TIGR role assignments, they would be listed in the consistency check section for review.

Frameshift Reports: Similar to the "Frameshift status" link that was described earlier for the "Genome Summary" page - basically a list of genes with frameshift reports to be resolved.

Hypothetical protein list: a list of hypothetical proteins, (those with insufficient evidence to make any functional assignment) for which there is any shred of information which might lead to annotation other than "hypothetical protein", this list is generated automatically after AutoAnnotate has made its initial assignment. Those "hypothetical proteins" called by AutoAnnotate that have any BER or HMM evidence are put on this list for manual review.

Annotation status: The same page as was described from the "Genome Summary" page - lists of the steps in annotation and a list of role categories, status of completion and annotator who did the work is noted.

Phage Region Viewer: A tool that lists any identified prophage regions in the genome and the genes within them.

PubMed Organism Search: Automatically takes you to the NCBI PubMed site and gives results for a PubMed search using the organism name as keywords. Useful for finding 32 literature on the organism you are working on.

"Annotation Tools": "Access Gene Lists" section

Although all of the tools described so far in this tutorial are quite useful, the bulk of annotator time is spent in viewing and editing information that is displayed on gene lists and Gene Curation Pages that are accessed through the "Access Gene Lists" section.

This tool will create a table of genes chosen according to the options in the red box at right. As mentioned in the overview, at TIGR we organize our annotation efforts around TIGR role categories. This tool allows us to view the genes within each TIGR role category.

The first option to select in this section is which molecule you wish to annotate. Some genomes consist of just one chromosome and nothing else, while others can have multiple chromosomes or chromosome(s) and one or more plasmids. If multiple DNA molecules exist for the genome in question, the pull down menu at the top of this section will list them along with their id number. The default selection is "All molecules" as the team usually annotates all molecules at once, however, to choose just one of the molecules, simply select it from the pull-down menu.

Then choose one of the 3 options for which role categories you want to see genes from with the toggle buttons: first you can choose all role categories, second you can choose one particular main role category, and third you can choose one particular sub-role category. All of the mainrole categories are listed in the pull-down menu in the main role category selection, to choose one, simply highlight it. In order to select a particular sub-role category you must enter into the box next to "single role category" the id number of the subrole category. There is a listing of all of the TIGR role categories and their id numbers on the next two pages of this tutorial.

Once you have chosen your desired options, click submit to see a list of the genes that fit your selections.

Luddscdigdeligt generates the Annotation Tools Genome Summary Gene Assignm ome Annotation Tools Genome Summary Gene Assignm • ACCESS GENE LISTS • molecule: All molecules • • molecule: All molecules • • Gene Assignm • call genes, ordered by role category (excludes hypothetical proteins) • • • • main role category Unclassified • • • • all genes, ordered by role category (excludes hypothetical proteins) • • • • main role category Unclassified • • • • gene: • end3: • • • • end5: end3: • • • • • location: /ust/local/annotation/GSP/overlap_analysis/gsp.overlap • • • • location: /ust/local/annotation/GSP/intergenic_regions/gsp.REPC • • • • • Data Consistency Checks • Frameshift Reports • • • • • • • Data Consistency Checks • • <th>iewanella one</th> <th>cidensis MR-1 Annotation Tools</th> <th>Home Logged in</th>	iewanella one	cidensis MR-1 Annotation Tools	Home Logged in
Ome Annotation Tools Genome Summary Gene Assignm • ACCESS GENE LISTS • molecule: All molecules • • • molecule: All molecules • • • • all genes, ordered by role category (excludes hypothetical proteins) • • • main role category Unclassified • • • • single role category role_id • • • • ACCESS GENE CURATION PAGE • • • • gene: • • • • • end5: end3: • • • • oVERLAP ANALYSIS • • • • • location: /usr/local/annotation/GSP/intergenic_regions/gsp.REPC • • • CUSTOM QUERY • • • • • Data Consistency Checks • • • • • PubMed Organism Search • • • •	The ann_tools.cg ating general prop	gi script generates the Annotator Tools webpage, which is the entry point for accessing the Submit webpa perties of the genome and determining the progress made in the Annotation of the genome of interest.	ge forall ORFs in a genome, a
• ACCESS GENE LISTS • molecule: All molecules • all genes, ordered by role category (excludes hypothetical proteins) • main role category Unclassified • single role category • ACCESS GENE CURATION PAGE • gene: • ACCESS GENES BY COORDINATE RANGE • end5: end3: • OVERLAP ANALYSIS • location: / Usr/local/annotation/GSP/loverlap_analysis/gsp.overlap • INTEREVIDENCE REGION ANALYSIS • location: / Usr/local/annotation/GSP/intergenic_regions/gsp.REPC • CUSTOM QUERY OTHER TOOLS • Data Consistency Checks • Frameshift Reports • Hypothetical Protein List • Annotation Status • Phage Region Viewer • PubMed Organism Search	Home	Annotation Tools Genome Summary	Gene Assignn
main role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category insingle role category <td></td> <td>ACCESS GENE LISTS molecule: All molecules</td> <td>roteins)</td>		ACCESS GENE LISTS molecule: All molecules	roteins)
ACCESS GENE CURATION PAGE gene: ACCESS GENES BY COORDINATE RANGE end5: end3: OVERLAP ANALYSIS location: /usr/local/annotation/GSP/overlap_analysis/gsp.overlap INTEREVIDENCE REGION ANALYSIS location: /usr/local/annotation/GSP/intergenic_regions/gsp.REPC CUSTOM QUERY OTHER TOOLS > Data Consistency Checks > Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search		main role category Unclassified single role category role_id	_
> gene: • ACCESS GENES BY COORDINATE RANGE • end5: end3: • OVERLAP ANALYSIS • location: /usr/local/annotation/GSP/overlap_analysis/gsp.overlap • INTEREVIDENCE REGION ANALYSIS • location: /usr/local/annotation/GSP/intergenic_regions/gsp.REPC • CUSTOM QUERY OTHER TOOLS • Data Consistency Checks • Frameshift Reports • Hypothetical Protein List • Annotation Status • Phage Region Viewer • PubMed Organism Search		© ACCESS GENE CURATION PAGE	
 ACCESS GENES BY COORDINATE RANGE end5: end3: end3: OVERLAP ANALYSIS location: //usr/local/annotation/GSP/overlap_analysis/gsp.overlap INTEREVIDENCE REGION ANALYSIS location: //usr/local/annotation/GSP/intergenic_regions/gsp.REPC CUSTOM QUERY OTHER TOOLS Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search 		> gene:	
 end5: end3: end3: OVERLAP ANALYSIS location: /usr/local/annotation/GSP/overlap_analysis/gsp.overlap INTEREVIDENCE REGION ANALYSIS location: /usr/local/annotation/GSP/intergenic_regions/gsp.REPC CUSTOM QUERY OTHER TOOLS > Data Consistency Checks > Frameshift Reports > Hypothetical Protein List > Annotation Status > Phage Region Viewer > PubMed Organism Search 		C ACCESS GENES BY COORDINATE RANGE	
 OVERLAP ANALYSIS location: //usr/local/annotation/GSP/overlap_analysis/gsp.overlap INTEREVIDENCE REGION ANALYSIS location: //usr/local/annotation/GSP/intergenic_regions/gsp.REPC CUSTOM QUERY OTHER TOOLS > Data Consistency Checks > Frameshift Reports > Hypothetical Protein List > Annotation Status > Phage Region Viewer > PubMed Organism Search 		end5: end3:	
t , location: //usr/local/annotation/GSP/overlap_analysis/gsp.overlap Image: Interevidence region Analysis . location: //usr/local/annotation/GSP/intergenic_regions/gsp.REPC Image: Custom QUERY Image: Custom QUERY Image: Dotata Consistency Checks . Frameshift Reports . Hypothetical Protein List . Annotation Status . Phage Region Viewer . PubMed Organism Search		© OVERLAP ANALYSIS	
 INTEREVIDENCE REGION ANALYSIS location: /usr/local/annotation/GSP/intergenic_regions/gsp.REPC CUSTOM QUERY OTHER TOOLS Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search 	mit	► location: //usr/local/annotation/GSP/overlap_analysis/gsp.overlap	
 location: //usr/local/annotation/GSP/intergenic_regions/gsp.REPC CUSTOM QUERY OTHER TOOLS Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search 	et	© INTEREVIDENCE REGION ANALYSIS	
CUSTOM QUERY CUSTOM QUERY OTHER TOOLS Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search		► location: //usr/local/annotation/GSP/intergenic_regions/gsp.REPC	
OTHER TOOLS Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search		• CUSTOM QUERY	
OTHER TOOLS Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search			
 Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search 		OTHER TOOLS	
		 Data Consistency Checks Frameshift Reports Hypothetical Protein List Annotation Status Phage Region Viewer PubMed Organism Search 	
			33

Gene List: The results of your selection from the Access Listings tool are displayed in a gene list containing gene id number, locus (if available), coordinates of the gene (end5, end3), common name of the gene/protein, gene_sym, EC number, and other roles for the protein. Not all of these fields will be populated for every gene. The genes are organized by role category (if your selection included more than one.) There are many features of the gene list, and much information displayed - text describing a feature is boxed in the same color as the feature itself.

			Clicking on the blue names of any mainrole category takes you to a gene list for											
Shewanella oneidensis MR-1 Gene List				inal calegory.										
	This List contains ORFs which are currently assigned to TI R microbial role categories. It is sorted by role category.													
	 All categories ►Unclassified ►Amino acid biosynthesis ►Purines, pyrimidines, nucleosides, and nucleotides ►Fatty acid and phospholipid metabolism ►Biosynthesis of cofactors, prosthetic groups, and carriers ►Central intermediary metabolism ►Energy metabolism ►Transport and binding proteins ►DNA metabolism ►Transcription ►Protein synthesis ►Protein fate ►Regulatory functions ►Signal transduction ►Cell envelope ►Cellular processes ►Mobile and extrachromosomal element functions ► Unknown function ►Hypothetical proteins ►Disrupted reading frame ►Glimmer rejects 													
Biosynthesis of cofactors, prosthetic groups, and carriers View list of Genome Properties found for this role category														
•]	Bioti	in		Role id	1: 77		1			Edit Annotation Notebool				
prop_acc st			ate	te property					is public?					
		GenProp0036		Y	ES	<u> </u>		biotin biosynthesis		YES				
А	С	gene id	locus	end5	end3			gene name		gene symbol	ec	other	roles	
	•	ORF02146 (GV)	SO0214	224657	225616	bir	h bifunc	tional protein		birA	6.3.4.15	12		
	•	ORF02395 (GV)	SO0001	774	334	mi	C prote	in		mioC		11:		
۲	•	ORF02552 (GV)	SO4626	4821460	4822251	1 bio	I protei	n		bioH				
۰	•	ORF04812 (GV)	SO2741	2856886	2858271	ade am	nosylme	ethionine8-amino-7-ox ferase	ononanoate	bioA	2.6.1.62			
۲	•	ORF04813 (GV)	SO2740	2856763	2855711	1 bio	in synti	nase		bioB	2.8.1.6			
	•	ORF04814 (GV)	SO2739	2855489	2854284	1 8-a	nino-7-o	oxononanoate synthase		bioF	2.3.1.47			
	•	ORF04816 (GV)	SO2738	2854163	2853336	6 bio	in synti	nesis protein BioC		bioC				
	•	ORF04817 (GV)	SO2737	2853281	2852586	det	niobioti	n synthase		bioD	6.3.3.3			
A green dot in the "A" column indicates this orf was given a high quality assignment by AutoAnnotate. (The only type of evidence that					rf ıt	Link to role notes for this category Click on the gene id (feat name)					ct			
will currently trigger this is an above trusted cutoff hit to an equivalog HMM.) A pink dot will appear in the "C" column once an annotator has finished annotation for the gene and marked it					vill nas it	S link to see the Gene Curation Page for each gene. Click on "GV" for Genome Viewer. The ORFs can be according to ar headers by click				s can be to any o by clickir	can be ordered 34 o any of the blue clicking on that header.			

complete

TIGR Role Categories - Page 1

Unclassified (the automated program was unable to assign a role to these)

185 Role category not yet assigned

Amino acid biosynthesis

- 70 Aromatic amino acid family
- 71 Aspartate family
- 73 Glutamate family
- 74 Pyruvate family
- 75 Serine family
- 161 Histidine family
- 69 Other

Purines, pyrimidines, nucleosides, and nucleotides

- 123 2'-Deoxyribonucleotide metabolism
- 124 Nucleotide and nucleoside interconversions
- 125 Purine ribonucleotide biosynthesis
- 126 Pyrimidine ribonucleotide biosynthesis
- 127 Salvage of nucleosides and nucleotides
- 128 Sugar-nucleotide biosynthesis and conversions
- 122 Other

Fatty acid and phospholipid metabolism

- 176 Biosynthesis
- 177 Degradation
- 121 Other

Biosynthesis of cofactors, prosthetic groups, and carriers

- 77 Biotin
- 78 Folic acid
- 79 Heme, porphyrin, and cobalamin
- 80 Lipoate
- 81 Menaquinone and ubiquinone
- 82 Molybdopterin
- 83 Pantothenate and coenzyme A
- 84 Pyridoxine
- 85 Riboflavin, FMN, and FAD
- 86 Glutathione
- 162 Thiamine
- 163 Pyridine nucleotides
- 191 Chlorophyll
- 707 Siderophores
- 76 Other

- Central intermediary metabolism
- 100 Amino sugars
- 698 One-carbon metabolism
- 103 Phosphorus compounds
- 104 Polyamine biosynthesis
- 106 Sulfur metabolism
- 179 Nitrogen fixation
- 160 Nitrogen metabolism
- 709 Electron carrier regeneration
- 102 Other

Energy metabolism

- 108 Aerobic
- 109 Amino acids and amines
- 110 Anaerobic
- 111 ATP-proton motive force interconversion
- 112 Electron transport
- 113 Entner-Doudoroff
- 114 Fermentation
- 116 Glycolysis/gluconeogenesis
- 117 Pentose phosphate pathway
- 118 Pyruvate dehydrogenase
- 119 Sugars
- 120 TCA cycle
- 159 Methanogenesis
- 105 Biosynthesis and degradation of polysaccharides
- 164 Photosynthesis
- 180 Chemoautotrophy
- 184 Other

Transport and binding proteins

- 142 Amino acids, peptides and amines
- 143 Anions
- 144 Carbohydrates, organic alcohols, and acids
- 145 Cations and iron carrying compounds
- 146 Nucleosides, purines and pyrimidines
- 182 Porins
- 147 Other
- 141 Unknown substrate

TIGR Role Categories - Page 2

DNA metabolis	sm			
132	DNA replication, recombination, and repair	Cell e	envelope	
183	Restriction/modification	91	Surface structures	
131	Degradation of DNA	89	Biosynthesis of murein sacculus and peptidoglycan	
170	Chromosome-associated proteins	90	Biosynthesis and degradation of surface polysaccarides and li	popolysaccharides
130	Other	88	Other	
Transcription				
134	Degradation of RNA	Cellu	lar processes	
135	DNA-dependent RNA polymerase	93	Cell division	
165	Transcription factors	188	Chemotaxis and motility	
166	RNA processing	701	Cell adhesion	
133	Other	702	Conjugation	
		96	Detoxification	
Protein synthes	sis	98	DNA Transformation	
137	tRNA aminoacylation	705	Sporulation and Germination	
158	Ribosomal proteins: synthesis and modification	94	Toxin production and resistance	
168	tRNA and rRNA base modification	187	Pathogenesis	
169	Translation factors	149	Adaptations to atypical conditions	
136	Other	706	Bioosynthesis of natural products	
		92	Other	
Protein fate				
97	Protein and peptide secretion and trafficking	Mobil	e and extrachromosomal element functions	
140	Protein modification and repair	186	Plasmid functions	
95	Protein folding and stabilization	152	Prophage functions	
138	Degradation of proteins, peptides, and glycopeptic	de s 54	Transposon functions	
189	Other	708	Other	
Regulatory fun	ctions	Unkn	own	
261	DNA interactions	703	Enzymes of unknown specificity	
262	RNA interactions	157	General	
263	Protein interactions			
264	Small molecule interactions	Нуро	thetical	
129	Other	156	Conserved	
		704	Domain	
Signal transdue	ction	-		
699	Two-component systems	Disru	pted reading frame	
700	PTS	270	NULL	36
710	Other	-		50
Gene list link: Edit Annotation Notebook:

Clicking on the "Edit Annotation Notebook" link on the gene list page will take you to a page where you can enter or edit annotation notes for a particular role category. It is in this text field that we store information that we think will be useful for the PI of the project in the analysis of the genome or in the preparation of the manuscript. Things such as the presence of an unexpected pathway, or the fact that a key step in another pathway is missing. Once the text is as you want it, click "submit" to store the information in the db.

Shewanella oneidensis MR-1	Annotation Notes - role id 83
----------------------------	-------------------------------

Logged into [gsp] as mlgwinn

The annotation_notebook.txt script directs the user to a web display page that contains annotators' comments about particular genes or regions of the genome that the annotators thought were unusual or interesting.

Appears to lack panD which I searched for with the E.coli sequence. No matches using blastp or tblastn. RTD

Update Reset

Gene list link: Role information page:

TIGR annotators expert in particular role categories have written "role notes" to aid new annotators and annotators unfamiliar with the category in the annotation process. These notes contain information on what genes belong in the category and what genes don't, on the pathways found in particular categories, and on the TIGR naming conventions for proteins within the category.

Any TIGR annotator can update or add text to the note field by typing it in and then clicking submit.

There is also a link to the role notes pages from the Gene Curation Page (GCP) which will be shown in the GCP section.

Shewanella oneidensis MR-1 | Role Information For Role_id 77

The role_info.cgi script is executed from the Submit web display page and directs the user to a web display page that contain Single Role Category.

Role 77 Biosynthesis of cofactors, prosthetic groups, and carriers - Biotin

Role Info:

```
Genes involved in the synthesis of biotin.
pathway from 6-carboxyhexanoyl-CoA plus L-alanine to biotin;
step
        gene
        8-amino-7-oxononanoate synthase (bioF)
TIGR00858: bioF
        adenosylmethionine-8-amino-7-oxononanoate aminotransferase
(bioA)
TIGR00508: bioA
        dethiobiotin synthetase (bioD)
TIGR00347: bioD
        biotin synthase (bioB)
TIGR00433: bioB
Other genes also involved:
BirA bifunctional protein (birA)
        acts as operon repressor, synthesizes corepressor, activates
biotin,
        and transfers activated biotin to proteins
biotin synthesis protein BioC (bioC)
        involved in an early, undefined step in biotin synthesis
biotin sulfoxide reductase (bisZ)
        changes biotin sulfoxide back to biotin, scavenging reaction
TIGR01738 bioH protein (bioH)
        in early steps of biotin biosynthesis
TIGR01204 bioW protein = 6-carboxyhexanoate--CoA ligase
        found in Bacillus and Methanoccus, involved in biotin
synthesis
        BioW plus BioF of Bacillus can replace bioC and bioH of E.
coli
                      (says PMID:2110099)
In many, but by no means all, organisms most of these genes can be
found in an operon.
mioC protein: MioC is a flavodoxin thought to function as an electron
transporter (role_id=112) and in biotin biosynthesis (role_id=77).
mioC neighbors oriC in E. coli. Early studies on mioC expression
demonstrate a dramatic effect on initiation of chromosome duplication
at oriC on minichromosomes. This role has not been demonstrated in
duplication of the wild type chromosome. Additionally, the
minichromosome is not necessarily a valid model for chromosomal
                                                                       ¥
replication Decause of this dubious association with shremosom
```

submit Update Role Note For 77

Gene Curation Page

The Gene Curation Page (GCP) is likely the most important page within Manatee, it is certainly the one that annotators spend the bulk of their time looking at and working with.

This page can be accessed within Manatee from many places:

any gene list, the "Access Gene Curation Page" option on the Genome Summary/Annotation Tools pages, Genome Viewer, and more.

The GCP is a very complex page so we will look at it in sections. I will try to organize the descriptions of each section in roughly the same order that the concepts behind each section were reviewed in the Annotation Overview.

Shewanella oneidensis MR-1 Gene Curation Page Home Logged into [gsp] as mlgwinn Image: Contract of the second second

GENE CURATION INFORMATION			ß
ORF04813 (SO2740) View BER Searches asmbl_id: 7974 Reload Page	end5/end3: 2856763 / 2855711 gene length: 1053 protein length: 350 molecular wt: 38790.13	database: gsp feat_name / locus: New Gene	
Select Display	Select Function	Refresh Searches	

GENE IDENTIFICATION			submit	hist	ory	I 🕻)
gene name:							
biotin synthase							
gene_sym: bioB							
EC number(s):		EC GO suggestions:					
2.8.1.6	100	GO:0004076 add biotin syntha	se activity (F)				
private comment:		public comment:					
Start confidence Low							
≻nt_comment		▶auto_comment					

Gene Curation Page

Gene Curation Information

This section contains basic identifying information about the gene and some search and display options.

The **feat_name** of the gene is listed at the top of the page, this number is called the "gene id" in gene lists in Manatee. The feat_name is followed in parentheses by the **locus name** (final loci are assigned to genes at the end of a project, once annotation is complete, but they may get temporary loci during the course of the project).

The **blue link** under these names is a link to a file containing the BER search results for this gene (see later slide). There is another link to this page further down the orf info page (will be seen in a later slide).

To the right of the ORF names is a box containing **coordinates, length, and molecular weight**. "end5" is the 5' coordinate for the beginning of the coding sequence, "end3" is the 3' coordinate for the end of the coding sequence.

Finally on the extreme right is a box allowing you to move to another ORF info page by typing in the feat_name or locus in the box and clicking "**new gene**". One can also change to an orf in a different genome by **changing the database** in the database box, typing in the new orf number and clicking "new gene".

If you want to reload theGCP, use the "**Reload Page**" link in this section. Do not use the browser's reload button as this can cause things to be sent to the db in error.

To generate new HMM and BER searches click "**Refresh Searches**" and enter your unix password.

Shewanella oneidensis MR-1	Gene Curation Page	Home Logged into [gsp] as mlgwinn			
D					
GENE CURATION INF	ORMATION	0			

ORF04813 (SO2740) View BER Searches asmbl_id: 7974	end5/end3: 2856763 / 2855711 gene length: 1053 protein length: 350	database: gsp feat_name / locus:
Select Display	molecular wt: 38790.13	Refresh Searches

5	GENE IDENTIFICATION			submit h	istory	[
	gene name:					
	biotin synthase					
	gene_sym: bioB					
;	EC number(s):		EC GO suggestions:			
	2.8.1.6	-log	GO:0004076 add biotin syntha	se activity (F)		
	private comment:		public comment:		_	
	Start confidence Low					
	⊁nt_comment) auto_comment			

Gene Curation Page Gene Identification

Initial information for this section comes from AutoAnnotate. The manual annotation then confirms or changes the information.

Common name: the descriptive name given to the protein

Gene sym: the gene symbol for the protein (in this case bioB) (we default to E. coli gene symbols when possible and B. subtilis for Gram + specific things)

EC#: If the protein is an enzyme, we store the Enzyme Commission number. See later slides for info on ECGO term suggestions.

private comment: a field for annotators to note information for later reference by themselves or other annotators. A good place to keep notes. **public comment**: comments meant to go out with our public accessions .

auto_comment: A link to information from the AutoAnnotate program indicating what information was used to make the preliminary annotation assignments (see next slide).

nt_comment: For non-TIGR comments. This is the place that collaborators can put comments to help the team in annotation.

Shewanella oneidensis MR-1 Gene Curation Page

GENE CURATION INFORMATION			ß
ORF04813 (SO2740) View BER Searches asmbl_id: 7974 Reload Page	end5/end3: 2856763 / 2855711 gene length: 1053 protein length: 350 molecular wt: 38790.13	database: gsp feat_name / locus: New Gene	
Select Display	Select Function	Refresh Searches	

GENE IDENTIFICATION		submit history	I
gene name:			
biotin synthase			
gene_sym: bioB			
EC number(s):		EC GO suggestions:	
2.8.1.6	Sec.	► GO:0004076 add biotin synthase activity (F)	
private comment:		public comment:	
Start confidence Low			
▶nt_comment		▶auto_comment	

Gene Curation Page - Auto Comment

Clicking on "auto_comment" pops up a text box with information on where

AutoAnnotate got the information it used for the preliminary annotation.



Gene Curation Page - BER Skim and Characterized Match

The characterized match section is where we enter the accession of a match gene whose function has been characterized in the lab (as opposed to having received its name based on sequence similarity.) This is stored as a piece of annotation evidence. This accession will pop into the go with_ev field in the proper format if you click on "Add to GO Evidence". (more on GO data later)

The BTAB SKIM section shows the top hits from the BER search file (see Annotation Overview presentation for more information on BER searches). The first column is the accession of the match protein (from various databases), the second is the percent similarity of the match, the third is the length of the match (in nucleotides), the fourth is the name of the match protein and finally, the P score from the BLAST search.

The color of the background for each entry in the skim indicates whether it is in the characterized table and at what confidence level: **green**=high confidence; **red**=automated process; **sky blue**=partial characterization; **olive**=trusted, used when multiple extremely good lines of evidence exist for function but no experiment has been done; **blue-green**=fragment/domain has been characterized; **fuzzy gray**=void, used to indicate that something that was originally thought to be characterized really is not; **gray**=omnium only

Clicking on the **blue accession number** will automatically populate the "Add accession" field in the characterized match section with that accession. Clicking on the **blue names of the proteins** in the skim will take you to a page with just the alignment to that protein.

The blue "View BER searches" link at the top of the skim section will take you to a file of all of the pairwise alignments from the BER search (see later slide). The tree icon takes you to a phylogenetic tree of the genome protein with the top hits of the skim the Belvu icon takes you to a multiple alignment of the

CHARACTERIZED M	АТСН		submit histo	ry 🗈
(aln) SP:P12996 coords: 7 / 350 score: 42 Pvalue: 7.2e-120 per id: 66.0% per sim: 79.7% [Add To GO				
			Evidence]	
Delete acces	ssion:		Add accession: [Add To GO Evidence]	
BER SKIM			sub	mit 🛙
- 🕒 Belvu	View BI	R Search	search date: Wed Oct 23 12:59:20 2002 Refresh Search	es
accession	%sim	length	description	p-value
OMNI:SO2740	100.0	349	biotin synthase {Shewanella oneidensis MR-1}	1.5e-176
SP:P36569	80.7	340	Biotin synthase (EC 2.8.1.6) (Biotin synthetase). (Serratia	2.5e-119
SP:P12996	79.7	342	Biotin synthase (EC 2.8.1.6) (Biotin synthetase). (Escherich	7.2e-120
GP:145425	79.7	342	biotin synthetase {Escherichia coli}	1.5e-119
GP:12620127	79.4	342	biotin synthase BioB {uncultured bacterium pCosHE2}	1.5e-119
OMNI:NTL03EC0855	79.4	342	biotin synthetase {Escherichia coli O157:H7 VT2-Sakai}□GPI13	5.1e-119
OMNI:NTL01YP1094	81.0	340	biotin synthase {Yersinia pestis CO92}□OMNIINTL02YP2986 biot	8.3e-119
GP:12620099	79.5	340	BioB-like protein {uncultured bacterium pCosFS1}	9.5e-118
OMNI:NTL02EC0848	79.1	342	biotin synthesis, sulfur insertion? {Escherichia coli O157:H 2	
SP:Q47862	79.2	339	Biotin synthase (EC 2.8.1.6) (Biotin synthetase). (Erwinia h	3.6e-118
SP:P12678	78.6	344	Biotin synthase (EC 2.8.1.6) (Biotin synthetase). (Salmonell	5.1e-119
OMNI:VC1112	81.8	348	biotin synthase {Vibrio cholerae El Tor N16961}□GPI9655583Ig	5.1e-119
OMNI:NTL03ST0726	78.6	344	biotin synthetase {Salmonella enterica serovar Typhi CT18}⊡G	1.1e-118
OMNI:NTL03PA00501	78.9	348	biotin synthase {Pseudomonas aeruginosa PAO1}⊡GPl9946364lgbl	7.7e-116
GP:12407614	76.8	339	biotin synthase BioB {uncultured bacterium pCosAS1}	9.1e-113
OMNI:NTL01XC0388	79.2	311	biotin synthase {Xanthomonas campestris pv. campestris ATCC3	2.8e-111
OMNI:NTL01XA0388	78.5	311	biotin synthase {Xanthomonas axonopodis pv. citri 306}□GPl21	6.6e-110
OMNI:NTL02BA0265	77.0	340	biotin synthase {Buchnera aphidicola Sg}□GPl21623185lgblAAM6	1.4e-109
OMNI:NTL01XF00065	79.4	309	biotin synthase {Xylella fastidiosa 9a5c}□GPI9104834lgblAAF8	8.4e-110
OMNI:NTL01RS0266	79.5	306	PROBABLE BIOTIN SYNTHASE PROTEIN {Ralstonia solanacearum GMI	4.7e-109
SP:P57378	77.3	342	Biotin synthase (EC 2.8.1.6) (Biotin synthetase). [Buchnera	1.1e-107
GP:15419053	79.1	328	biotin synthase {Acinetobacter calcoaceticus}	1.6e-106
OMNI:CC3521	76.2	339	biotin synthase {Caulobacter crescentus CB15}□GPI13425251lgb	3.0e-105
OMNI:NTL01BMA0776	79.8	311	BIOTIN SYNTHASE {Brucella melitensis 16M}□GPI17984969lgbIAAL	6.3e-105

Links from the Gene Curation Page - The BER alignment file

This page is accessible by clicking on the "View BER searches" link at the top of the Info page or at the top of the BTAB skim section.

Here you will find multiple pairwise alignments of the genome protein to hits found in the BER search.

In the header of each alignment will be listed the accessions and names for this protein from every database where it is found. These accessions are clickable objects and will take you to the page for the match protein in the database in question.

The background color of the header will be gold if the protein is found in the characterized table with the confidence level indicated by the color of the text for the accession found in the characterized table. (This is seen for the SP accession in this alignment.)

Names in Skim are first entry in header, not necessarily the name you want to use, check role notes for TIGR naming standards, check IUBMB EC site for official enzyme names, look in header for SwissProt as a model for the name if previous two guides are not available.

The background color in the Skim may be assigned to an entry in the header different than the one named in the Skim. Links to info pages for the match protein in the source db.



BER Alignment detail: Boxed Header

66.0/79.7% over 343aa	Escherichia coli
 SPIP12996 Biotin synthase (EC 2.8.1.6) (Biotin synthetase). 	Edit characterized
 PIRIJC2517 SYECBB biotin synthase (EC 2.8.1.6) bioB [v. 	alidated] - Escherichia coli (strain K-12) Insert characterized
 GBIAAC73862 [IGPI1786992]AE000180 biotin synthesis 	sulfur insertion? {Escherichia coli K12:} Insert characterized

-The background color of this box will be gold if the protein is in the characterized table and grey if it is not.

-The top bar lists the percent identity/similarity and the organism from which the protein comes (if available).

-The bottom section lists all of the accession numbers and names for all the instances of the match protein from the source databases (used in building NIAA for the searches.)

-The accession numbers are links to pages for the match protein in the source databases.

-A particular entry in the list will have colored text (the color corresponding to its characterized status) if that is the accession that is entered into the characterized table - this tells the annotators which link they should follow to find experimental characterization information. Only one accession for the match protein need be in the characterized table for the header to turn gold.

-There are links at the end of each line to enter the accession into the characterized table or to edit an already existing entry in the characterized table.

BER Alignment detail: alignment header

```
ORF04813( 7 - 350 of 350 aa)

SP|P12996|BIOB_ECOLI(4 - 346 of 346) Biotin synthase (EC 2.8.1.6)

%Match = 42.3

%Identity = 66.0 %Similarity = 79.7

Matches = 227 Mismatches = 69 Conservative Sub.s = 47

Gaps = 1 InDels = 3 Frame Shifts = 0

Primary Frame = 1 [343, 0, 0]
```

-It is most important to look at the range over which the alignment stretches and the percent identity

-The top line show the amino acid coordinates over which the match extends for our protein

-The second line shows the amino acid coordinates over which the match extends for the match protein, along with the name and accession of the match protein

-The last line indicates the number of amino acids in the alignment found in each forward frame for the sequence as defined by the coordinates of the gene. The primary frame is the one starting with nucleotide one of the gene. If all is well with the protein, all of the matching amino acids should be in frame 1.

-If there is a frameshift in the alignment (see overview) the phrase "Frame Shifts = #" will flash and indicate how many frameshifts there are.

BER Alignment detail: alignment of amino acids



-In these alignments the codons of the DNA sequence read down in columns with the corresponding amino acid underneath.

-The numbers refer to amino acid position. Position 1 is the first amino acid of the protein. The first nucleotide of the codon coding for amino acid 1 is nucleotide 1 of the coding sequence. Negative amino acid numbers indicate positions upstream of the predicted start of the protein.

-Vertical lines between amino acids of our protein and the match protein (bottom line) indicate exact matches, dotted lines (colons) indicate similar amino acids.

-Start sites are color coded: ATG is green, GTG is blue, TTG is red/orange

-Stop codons are represented as asterisks in the amino acid sequence. An open reading frame goes from an upstream stop codon to the stop at the end of the protein, while the gene starts at the chosen start codon.

Swiss-Prot entry - slide #1 - top of page

SwissProt is an incredibly useful database for manual annotation. All of the genes in SwissProt have been manually annotated by an experienced knowledgeable staff. In addition, along with each protein's annotation is stored additional information on references that describe the protein, cross referened databases in which the protein can be found, motifs which the protein contains, and coordinates of any known features in the protein (and much more.)

	[
	NiceProt View	of Swiss-Prot: P	<u>12996</u>	Printer-friendly view	Submit update	Quick BlastP search
	[En	try info] [Name and origin] [Refe	rences] [Comments] [Cross-reference	ces] [Keywords] [Features	[] [Sequence] [Tools]	
	Note: most headings are clicke his, even i	they don't appear as links. They link to the	<u>eser manual</u> or <u>other documents</u> .			
accession and	Entry information					
accession and	Entry name	BIOB_ECOLI				
version	Secondary accession numbers	None				
	Entered in Swiss-Prot in	Release 13, January 1990				
information	Sequence was last modified in	Release 35, November 1997				
	Annotations were last modified i	n Release 44, July 2004				
	Name and origin of the protei	n an				
name $FC#$	Protein name Synonyms	Biotin synthase		<u> </u>		
Hame, LO#	5 yhonyms	Biotin synthetase	LINK to Enzyme	Commissio	n page	
gene_symbol	Gene name	Name: bioB OrderedLocusNames: b0775	(see later slide)			
taxonomy	From	Escherichia coli [TaxID: 562]				
laxonomy	Taxonomy	Bacteria; Proteobacteria; Gamma	proteobacteria; Enterobacteriales; Enter	robacteriaceae; Escherichia.		
	Keterences	LEIC ACID				
6	MEDLINE=89066784;Pub	Med=3058702 [NCBI, ExPASy, El	31, Israel, Japan]			
references with	Otsuka A.J., Buoncristiani "The Escherichia coli biotin	M.R., Howard P.K., Flamm J., John biosynthetic enzyme sequences pred	<u>son O.;</u> licted from the nucleotide sequence of a	the bio operon ".		
links to	J. Biol. Chem. 263:19577-1	<u>9585(1988)</u> .	neted from the interestitle sequence of	une oto operone ;		
	[2] SEQUENCE FROM NUC	LEIC ACID.				
abstracts (click	"Genetic material for expres	sion of biotin synthetase enzymes.";				
	Patent number GB2216530	, 11-OCT-1989.				
on NCBI to see	STRAIN=K12 / MG1655	LEIC ACID.				
o DubMad	MEDLINE=97426617;PubMed=9278503 [NCBI, ExPASy, EBI, Israel, Japan]					via N 337 - Kielenoteiole H A
a Fubivieu	Goeden M.A., Rose D.J., Mau B., Shao Y.;					
abstract of the	"The complete genome sequence of Escherichia coli K-12."; Science 277:1453-1474(1997).					
nonori	[4] CHARACTERIZATION. PubMad=8142361 DICPU	ExDASy EDI Israal Ispar-1				
paper)	Sanyal I., Cohen G., Flint I	D.H.; ED1, ISTACI, Japall				
	"Biotin synthase: purification Biochemistry 33:3625-363	n, characterization as a [2Fe-2S] clu (1994).	ster protein, and in vitro activity of the	Escherichia coli bioB gene	product.";	
	[5] MUTAGENESIS OF CYS MEDLINE=21547100/Pub MEDLINE=21547100/Pub	TEINE RESIDUES.	PI Israal Isnan1			

Swiss-Prot entry - slide #2 - middle of page

useful functional information

links to other dbs where the protein is found or to motif clusters or protein families which this protein is a member of

Comments

- CATALYTIC ACTIVITY: Dethiobiotin + sulfur = biotin.
- COFACTOR: Binds a 4Fe-4S cluster coordinated with 3 cysteines and an exchangeable S-adenosyl-L-methionine, and a 2Fe-2S cluster coordinated with 3 cysteines and 1 arginine.
- PATHWAY: Biotin biosynthesis; last step.
- SUBUNIT: Homodimer.
- · SIMILARITY: Belongs to the biotin and lipoic acid synthetases family.

Copyright

This Swiss-Prot entry is copyright. It is produced through a collaboration between the Swiss Institute of Bioinformatics and the EMBL outstation - the European Bioinformatics Institute. There are no restrictions on its use by non-profit institutions as long as its content is in no way modified and this statement is not removed. Usage by and for commercial entities requires a license agreement (See http://www.isb-sib.ch/announce/ or send an email to http://www.isb-sib.ch/announce/

EMBL	J04423; AAA23515.1; [EMBL / GenBank / DDBJ] [CoDingSequence] A11530; CAA00965.1; [EMBL / GenBank / DDBJ] [CoDingSequence] AE000180; AAC73862.1;[EMBL / GenBank / DDBJ] [CoDingSequence]
PIR	<u>IC2517;</u> SYECBB.
PDB	1R30; 13-JAN-04.[ExPASy / RCSB / EBI]
ECO2DBASE	E038.6; 6TH EDITION.
EchoBASE	<u>EB0116;</u>
EcoGene	EG10118; bioB.
EcoCyc	EG10118; bioB.
CMR	<u>P12996;</u> b0775.
InterPro	IPR010722; BATS. IPR002684; Biotin_synth. IPR006638; Elp3/MiaB/NifB. IPR007197; Radical_SAM. Graphical view of domain structure.
Pfam	PF06968: BATS; 1. PF04055: Radical_SAM; 1. Pfam graphical view of domain structure.
SMART	SM00729; Elp3; 1.
TIGRFAMs	TIGR00433; bioB; 1.
ProDom	[Domain structure / List of seq. sharing at least 1 domain]
HOBACGEN	[Family / Alignment / Tree]
BLOCKS	P12996.
ProtoNet	<u>P12996</u> .
ProtoMap	<u>P12996</u> .
PRESAGE	<u>P12996</u> .
DIP	<u>P12996</u> .
ModBase	<u>P12996</u> .
SMR	P12996; 550A7899A2DF6082.
SWISS-2DPAGE	Get region on 2D PAGE.
UniPof	View objector of proteins with at least 500% / 000% identity

Swiss-Prot entry - slide #3 - bottom of page

keywords and sequence features with coordinates

Keyword	ls									
2Fe-2S;	2Fe-2S; 3D-structure; 4Fe-4S; Biotin biosynthesis; Complete proteome; Iron-sulfur; Transferase.									
Features	Features									
.	Feature ta	ble viewer								
Key	From To	Length	Descriptio	n						
METAL	53 53		Iron-sulfur	1 (4Fe-43).			soquence features			
METHL	57 57		Iron-sulfur Inco-sulfur	1 (4Fe-43). 1 (4Fe-48)			sequence realures			
METHL	97 97		Iron-sulfur	1 (416-43). 2 (2Fe-23).						
METAL	128 128		Iron-sulfur	2 (2Fe-23).						
METAL	188 188		Iron-sulfur	2 (2Fe-23).						
METAL	<u>260 260</u>		Iron-sulfur	2 (2Fe-23).						
CONFLICT	<u>63 63</u>		S -> T (in R	ef. <u>1</u>).						
Sequenc	e informatio)n		28648 D.		CDC(4, 550	A 7866 & 3D Ecopa (This is a shark-sum on the second sol		4	
Length: 3	940 AA	Molec	ular weight:	58648 Da		CKC04: 550.	A /899A2DF6082 [This is a checksum on the sequence]			
:	10 20) 30) 40	50	60					
THEFT	LS QVIILIIKPI	r mmartHóón	/ EKQELDEKQV	QUETLIEIKT	GHCFEDCKIC					
1	70 80	90) 100	110	120					
POSSRYKT	I SL EAERLMEVEC	i vlesarkaka	i agstrecmga	AWKNIERDM	PYLEOMVOGV					
1:	30 140) 150 I I) 160	170	180					
KAMGLEAC	T LOTISISQAO	, RIANAGLDYN	(NHNLDTSPIF	YGNIITTRTY	QERLDTLEKV					
10	20 200			220	240					
1	1 1	, 210 I I	, <u>22</u> 0	200	240					
RDAG IKVC:	SG GIVGLGETVN	C DRAGLLLQLA	NLPTPPESVP	INMLWKWKGT	PLADNDDVDA					
25	50 260	270) 280	290	300					
				1						
I DI IRTIA	WH RIPPETSTV	K ISHGREQMNE	. QIQHMCIMAG	HNSTFTGCKL	LTTPNPEEDK					
3	10 320) 330	340							
DLOLERKI	I SL NROOTAVIAG	I INECOORLEO	 ALMERTER	YNAAAT.						
and an adding	and the spectrum of the second								50	
								P12996 in FASTA format		

View of EC number info page from Swiss Institute of Bioinformatics site

NiceZyme View of EN	ZYME: EC 2.8.1.6
Official Name	
Biotin synthase.	
Alternative Name(s)	
Biotin synthetase.	
Reaction catalysed	
Dethiobiotin + sulfur <=> biotin	
Cofactor(s)	
Iron-sulfur.	
Comments	to data - it is not alemental sulfur or an iron sulfur eluster
Cross-references	to date - it is not elemental surful of an iton-surful cluster.
BRENDA	2.8.1.6
EMP/PUMA	2.8.1.6
WIT	2.8.1.6
Kyoto University LIGAND chemical database	2.8.1.6 Link to official Enzyma Commission alto
IUBMB Enzyme Nomenclature	2.8.1.6 LINK to official Enzyme Commission site
IntEnz	2.8.1.6
MEDLINE	Find literature relating to 2.8.1.6
Swiss-Prot	P54967, BIOB_ARATH; P19206, BIOB_BAC3H; P53557, BIOB_BAC3U; P57378, BIOB_BUCAI; Q8K9P1, BIOB_BUCAP; Q89AK5, BIOB_BUCBP; P12997, BIOB_CITTR; P46396, BIOB_CORGL; P12996, BIOB_ECOLI; Q47862, BIOB_ERWHE; P44987, BIOB_HAFIN; Q92JK8, BIOB_HELPJ; Q25956, BIOB_HELPY; Q58692, BIOB_METJA; P94966, BIOB_METSK; P46715, BIOB_MYCLE; Q06601, BIOB_MYCTU; P12678, BIOB_SALTY; Q59778, BIOB_SCHPO; P36569, BIOB_SERMA; P73538, BIOB_SYNY3; P32451, BIOB_YEAST;

View entry in original ENZYME format

All Swiss-Prot entries referenced in this entry, with possibility to download in different formats, align etc.

View of information page for an EC number at IUBMB site

The Enzyme Commission (EC) is part of the IUBMB and is charged with maintaining the database of enzyme classifications. In the EC system, each reaction is assigned a 4 part accession number with each part consisting of an integer, where the numbers are separated by periods. As one moves from the first number to the second to the third to the fourth the nature of the reaction becomes more specific. For example: EC2.-.- = "transferase", 2.8.-.- = "transferase, transferring sulfur-containing groups", 2.8.1.- = "sulfurtransferases", and finally 2.8.1.6 = "biotin synthase" (a specific sulfurtransferase, which is a specific class of transferases that transfer sulfur-containing groups). One can see the breakdown of all of the classes within each EC first number (they only go up to 6) by clicking on the home page for each number (see below).

IUBMB Enzyme Nomenclature

EC 2.8.1.6

Common name: biotin synthase

Reaction: dethiobiotin + sulfur = biotin

Systematic name: dethiobiotin:sulfur sulfurtransferase

Comments: an iron-sulfur enzyme. The sulfur donor has been unidentified to date - it is not elemental sulfur or an iron-sulfur cluster.

Links to other databases: BRENDA, EXPASY, KEGG, ERGO, PDB, CAS registry number: 80146-93-6 (204794-88-7, 179608-56-1, 209603-31-6, 153554-27-9, 174764-24-0 and 215108-34-2)

References:

1. Shiuan, D., Campbell, A. Transcriptional regulation and gene arrangement of Escherichia coli, Citrobacter freundii and Salmonella typhimurium biotin operons. Gene 67 (1988) 203-211. [Medline UI: 89006280]

 Zhang, S., Sanyal, I., Bulboaca, G.H., Rich, A., Flint, D.H. The gene for biotin synthase from Saccharomyces cerevisiae: cloning, sequencing, and complementation of Escherichia coli strains lacking biotin synthase. Arch. Biochem. Biophys. 309 (1994) 29-35. [Medline UI: <u>94161552</u>]

[EC 2.8.1.6 created 1999]

Return to EC 2.8.1 home page Return to EC 2.8 home page

Return to EC 2 home page Return to EC 2 home page Click here to see all the classifications within EC #2 (the transferases)

Return to IUBMB Biochemical Nomenclature home page

Links from the Gene Curation Page - Tree (may not work on laptops)



53

Links from the Gene Curation Page - BER multiple alignment (will not work on laptops)

File Edit Colour Sort Picked:									
(26×440)		-207080							
OMNIINTL01XA0388 OMNIINTL01XC0388 OMNIINTL01XC0388 OMNIINTL01RS0266 OMNIINTL03PA00501 GPI59215471emb1CAB56476.111AJ2 ORF06889 OMNIVC1112 OMNIINTL03EC0855 OMNIINTL03EC0855 OMNIINTL02EC0848 SPIP129961BI0B_ECOLI GPI1454251gb1AAA23515.111J0442 GPI126201271gb1AAG60579.11AF25 OMNIINTL03ST0726 SPIP126781BI0B_SALTY SPIQ478621BI0B_SALTY SPIQ478621BI0B_ERWHE GPI124076141gb1AAG53589.11AF24 SPIP365691BI0B_SERMA GPI126200991gb1AAG60559.111AF2	$\begin{array}{cccccccccccccccccccccccccccccccccccc$	MSVVLRHDWDRKELQALFDL PFPELLHRAASVHRAHFDPAQVQVSTLLSVKTGGCPEDCAYCP MSVVVRHDWDRKELHALFALPFPELLHRAASVHRAHFDPAEVQVSTLLSVKTGGCPEDCAYCP TPGQSPNARWSREAIEALFALPFNDLLFQAQQVHRAHFDANAVQLSTLLSIKTGGCPEDCSYCP TASVATRHDWSLAEVRALFEQPFNDLLFQAQTVHRAHFDANAVQLSTLLSIKTGGCPEDCKYCP TDACATRHDWSLAEVRALFEQPFNDLLFQAQTVHRAHFDANRVQVSTLLSIKTGACPEDCKYCP STTATLRHDWTLAEVRALFVQPFNDLLFQAQTVHRAHFDANRVQVSTLLSIKTGACPEDCKYCP MEVRHNWTVAEVKALLDKPFNDLLFEAQQVHRQHFDANRVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLDLFEAQQVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLELFEAQQIHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLELFEAQQIHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRPRWTLSQVTELFEKPLLELFEAQQIHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRRWTLSQVTELFEKPLLELFEAQQIHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRRWTLSQVTELFEKPLLELFEAQQIHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MAHRRWTLSQVTELFEKPLLELFEAQQIHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MADRIHWTVGLAQTHFFKPLELLFEAQTVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP MADRIHWTVGLAQTLFDKPLLELFEAQTVHRQHFDPRQVQVSTLLSIKTGACPEDCKYCP							

Gene Curation page - HMM hits scoring above noise

(Text describing the features of the HMM section is boxed in the same color as each feature.)

The blue id numbers for each HMM link to an info page for that HMM.

Key information is the isology type and the "total" and "cutoff" scores.

The "Add To GO Evidence" link automatically fills the HMM information into the "with" field in the GO term entry box.

GO terms assigned to each HMM are listed under the HMM (if any). Clicking on the "Add" button here adds not only the GO term id, but also the HMM evidence.

The "Add To Annotation" link will automatically copy the annotation from the HMM to the protein.

MM						subm	it I all hmms
IGR00433: biotin synthase		gene_s	ym: bioB	ec#: 2.8.1.6	role_id: 77	,	[Add To Annotation]
Isology: equivalog Total score: 564, 1	Trusted cutoff: 300.00 Trusted cutoff2: 300.00) D (Gatheringcutoff: 30 Gatheringcutoff2: 30).00 0.00	Noise cutoff: 50.00 Noise cutoff2: 50.0) 0	Total expect: 1.5e-1 0
YiewAlignment ⊁align page	Coords Hiv 18-313 1-	MM Coords ∙350 / 350	Score 564.1	Expect 1.5e-166	Curation	[A	dd To GO Evidence]
► GO:0004076 add bio	tin synthase activity (tin biosynthesis (P)	F)					
Senome Properties state property name YES biotin biosyn	ne add C nthesis	Rules.spl GO evidence [GO]	Th de late	s section scribed on er slide			
F06968: Bio <mark>lin and Thiami</mark> i	n Synthesis associate	ed domain	ger	ie_sym: none	ec#: none ro	le_id: none	e [Add To Annota
isology, domain Total score: 181.7	Trusted cutoff: 43.70 Trusted cutoff2: 43.70	(Gathering cutoff: 25 . Gathering cutoff 2: 25	00 .00	Noise cutoff: 19.10 Noise cutoff 2: 19.10	,	Total expect: 9.8e-5
ViewAlignment ▶align page	Coords H 223-315	IMM Coords 1-115/115	Score 181.7	Expect 9.8e-52	Curation	[Ad	d To GO Evidence]
▶ No HMM-GO Suggestions To	Display.						
F04055:radical SAM doma	in protein		gene_sym: none	ec#: n	one role_id:	703	[Add To Annotation
Total score: 82.7	Trusted cutoff: 8.80 Trusted cutoff2: 8.80	G	Gatheringcutoff: 8.4 Natheringcutoff2: 8.4	ا (۱	Noise cutoff: 8.30 Noise cutoff 2: 8.30		Total expect: 6.1e-22
View Alignment ▶align page	Coords HM 50-212 1-	1M Coords 1637 163	Score 82.7	Expect 6.1e-22	Curation	[Add	1 To GO Evidence]
► GO:0003824 add cate	alytic activity (F)						

Click to see hits below noise

HMM report page - to get to this page click on an HMM accession number almost anywhere in Manatee

At the top is information about the HMM including HMM name, associated annotation (gene symbol, EC#, TIGR role, etc.) and comments from the authors.

Below is a list of all genes in the organism which hit the HMM and the scores they received. The row with the gold background is the protein of interest. Rows with a green background have scores below the trusted cutoff, rows with a purple background have scores below the noise cutoff.

hewanella oneidensis MR-1	TIGR00433 HM	M Report for ORF04813
---------------------------	--------------	-----------------------

Home | Logged into [gsp] as mlgwinn

This page displays information about a specific HMM accession as it relates to the ORF being annotated. General information about the model is presented, as well as an alignment of the model to the ORF and a list of all hits of this model to the genome. The user can follow links to more information about the model and other proteins that the model being annotated.

accession and name		TIGR00433: biotin synthase									
expanded name		biotin synth	biotin synthetase								
	gene symbol	bioB		EC number	2.8.1.6		HMM length	350			
	model type	equiva	og	trusted cutoff	300.00		noise cutoff				
	author	Loftus	BJ	created	04/20/99		last modified	09/23/03			
related accession		IPR002	684	accession type	InterPr	o assignment					
	role category	77: Biosynt	hesis of c	cofactors, prosthetic groups, and	carriers, Biotin						
	gene ontology	GO:0004076 (function): biotin synthase activity GO:0009102 (process): biotin biosynthesis									
	comment	Catalyzes the last step of the biotin biosynthesis pathway.									
	private comment										
	Edit HMM Annotation			HMM Inter Link Edit	or	All DB Hits to TIGR00433					
C	olor key										
	Protein of Interest.										
	Scores below trusted cutoff (< 300.00).									
	Scores below noise cutoff (<	50.00).									

feat_name	role_id	EC number	gene region	HMM region	score	gene name
ORF04813	77	2.8.1.6	18-313	1-350	564.1	biotin synthase
ORF03390	157		34-331	1-350	-168.2	biotin synthase family protein
ORF01034	80		76-296	1-350	-178.3	lipoic acid synthetase
ORF03392	162		62-370	1-350	-187.3	thiH protein, putative

Genome Properties - linked from the Gene Curation Page in the HMM section

If an HMM is part of a genome property, there will be a link here and an indication of the state of the property - in this case "YES" indicating that the organism has an intact biotin biosynthesis pathway. Clicking on the name of the property takes one to a property report page.

If you want to use the Genome Property as evidence for GO annotation, click the "GO" link under the "add GO evidence" section. (more on GO data later)

The "Run Rules.spl" link

						submit all hmms 🗈
TIGR00433: biotin syntl Isology: equivalog	lase	gene	e_sym: bioB	ec#: 2.8.1.6	role_id: 77	[Add To Annotation]
Total score: 564, 1	Trusted cutoff: Trusted cutoff2	: 300.00 :: 300.00	Gatheringcutoff: 3 Gatheringcutoff2: 3	00.00 100.00	Noise cutoff: 50.00 Noise cutoff2: 50.00	Total expect: 1.5e-166
View Alignment	Coords	HMM Coords	Score	Expect	Curation	[Add To GO Evidence]
⊁align page	18-313	1-350/350	564.1	1.5e-166	V	
+ GO:0009102 add	biotin biosynthe	sis (P)				

Genome Property info page (part 1): biotin biosynthesis

This has general information about the property, GO terms assigned to the property, and a place for curators to put comments regarding this property in this organism.

Property Definition								
property:	biotin biosynthesis	state:	YES					
property type:	PATHWAY	value:	1					
role id:	77	assignby:	HYBRID					
GO ids:	GO:0006355: add regulation of transcription, DNA-dependent GO:0009102: add biotin biosynthesis	date:	Mar 10 2004 3:51PM					
description:	Biotin is an essential cofactor for many carboxylation (addition of C02) reactions. This property reflects biosynthesis from pimeloyl-CoA. The source of pimeloyl-CoA may vary. BioF (EC 2.3.1.47, 8-amino-7-oxononanoate synthase, also called 7-keto-8-aminopelargonic acid synthetase) converts pimeloyl-CoA to 8-amino-7-oxononanoate. BioA (EC 2.6.1.62, adenosylmethionine-8-amino-7-oxononanoate aminotransferase) converts the product to 7,8-diaminononanoate, from which BioD (EC 6.3.3.3, dethiobiotin synthase) makes dethiobiotin. BioB (EC 2.8.1.6, biotin synthase) then makes biotin itself. Enzymes such as BioH involved in pimeloyl-CoA biosynthesis typically receive biotin-related annotations but may also appear in genomes in which biotin is not synthesized and pimeloyl-CoA is used for something else.							
auto_comment:								
curator comment:			update					

Genome Property info page (part 2): biotin biosynthesis

Prop	perty	Steps					
R	eqvi	RED		8-amino-7-oxo-nonanoate synthase (2	2]		
Α	С	GC	gene id	gene name	evidence	role id	ОР
		0	ORF04814	8-amino-7-oxononanoate synthase	GENE_CLUSTER TIGR00858	77	
R	eqvi	RED		adenosyl methionine 8-amino-7-oxononanoate tra	nsaminase (3)		
Α	С	GC	gene id	gene name	evidence	role id	ОР
		0	ORF04812	a.den.osylmethionine—8-amino-7-oxononanoate aminotransferase	TIGR00508	77	
R	eqvi	RED		dethiobiotin synthase (4)			
Α	С	GC	gene id	gene name	evidence	role id	ОР
		0	ORF04817	dethiobiotin synthase	TIGR00347	77	
R	EQVI	RED		biotin synthase (5)			
A	С	GC	gene id	gene name	evidence	role id	ОР
		O	ORF04813	biotin synthase	TIGR00433	77	
NOT	REQ	VIRED		BioC (bioC)			
Α	С	GC	gene id	gene name	evidence	role id	ОР
		0	ORF04816	biotin synthesis protein BioC	TIGR02072	77	
NOT	REQ	VIRED		bioH protein (bioH)			
A	С	GC	gene id	gene name	evidence	role id	ОР
			ORF02552	bioH protein	TIGR01738	77	
NOT	REQ	VIRED		biotin repressor (represso)			
Α	С	GC	gene id	gene name	evidence	role id	ОР
				THIS COMPONENT HAS NOT BEEN IDENTIFIED			

This section of the page shows the steps for the property, which steps are required and which steps are not, and the genes from the genome that have been identified for each step.

One can link to the GCP for each gene or to the HMM info page for the HMMs named by clicking on the gene id or HMM accession, respectively.

Genome Property info page (part 3): biotin biosynthesis

This section has reference information and a graphic showing the cluster of genes in the organism involved in the property. One can click on the arrows in the graphic to get a GCP for that gene.

Property F	References				
accession:	PMID:9847135				
title:	KEGG: Kyoto Encyclopedia of Genes and Genomes.				
authors:	Ogata H, Goto S, Sato K, Fujibuchi W, Bono H, Kanehisa M.				
alter the sec	N1-:- A-:J- D 1000 I 1.07(1):00.24				
Literature Ref	erences]			
Web Referenc	63				
title: KEGG: Bi	otin Metabolism			ORF04812	
]			
			_		

Gene Curation Page - Evidence Picture - ORF04813

All of the evidence stored for an ORF is displayed in this graphic. The black bar represents the ORF in question. Green bars represent HMMs which hit the ORF above trusted cutoff. Green HMM bars indicate above trusted score, orange indicates above noise but below trusted, red indicates below noise and is generally not shown unless an annotator has decided that the HMM should be included as evidence by toggling the curation box. The pink bar represents the characterized match to this ORF. Characterized matches are shown in different colors that at this time have no meaning. Also shown here is a secondary structure prediction (not run on all genomes). Clicking on the colored bars in the graphic opens windows with additional information on that piece of evidence. To get additional cog info, you must click on the very skinny bar all the way to the left of the cog row. The evidence picture for ORF04813 does not contain all of the possible evidence types, so later slides will show some evidence pictures from other genes.

EVIDENCE PICTURE	submit 🗎
	sec structure: Coil(-), Strand(blue), Helix(yellow) S02740 TIGR00433: biotin synthase PF06968: Biotin and Thiamin Synthesis associated doma: PF04055: radical SAM domain protein COG0502 (p-value: none) Characterized match: SP:P12996

Secondary structure prediction



Key: 📑 helix — coil 👘 strand

The biotin synthase does not have all of the evidence types that are possible, therefore, the following screen shots will show some evidence pictures from other genes displaying additional evidence types.

Following the evidence pictures will be the evidence detail pages linked to from the evidence pictures.

After all of the evidence types have been represented, the tutorial will resume with ORF04813.

Gene Curation Page - Evidence Picture (ORF03779)

Additional evidence types shown here are:

TmHMM - an HMM specific for transmembrane regions, built by the Center for Biological Sequence Analysis, Denmark

Paralogous Family membership - if a protein is a member of a paralogous family it will be represented with a blue bar, clicking on the bar takes you to a page listing all the family members. Paralogous familes are built from searching the protein set for a genome against itself. First families are built according to shared hits to HMMs, then regions not matching HMMs are searched against each other to find additional families. The families corresdponding to HMMs are given names with the HMM accession number, others are given numbers.

0 40 80 120 160 200 240 280	
sign sign mini- sec source source	ignalP:SP-HMM mHMM ec structure: Coil(-), Strand(blue), Helix(yellow) O3601 IGR02140: sulfate ABC transporter, permease protein (IGR00969: sulfate ABC transporter, permease protein F00528: Binding-protein-dependent transport systems : S00402: Binding-protein-dependent transport systems : OG0555 (p-value: none) haracterized match: SP:P16702 aralogous Domain: fam_PF00528 aralogous Domain: fam 11

NOTE: this display is from ORF03779

NOTE: this display is for ORF03779

TMHMM result

HELP with output formats

# Sequence	Length: 343			
# Sequence	Number of predict	ed TMHs: 6		
# Sequence	Exp number of AAs	in TMHs: 139.	48261	
# Sequence	Exp number, first	60 AAs: 20.9	9155	
# Sequence	Total prob of N-i	n: 0.99	734	
# Sequence	POSSIBLE N-term s	ignal sequence		
Sequence	TMHMM2.0	inside	1	8
Sequence	TMHMM2.0	TMhelix	9	31
Sequence	TMHMM2.0	outside	32	97
Sequence	TMHMM2.0	TMhelix	98	120
Sequence	TMHMM2.0	inside	121	132
Sequence	TMHMM2.0	TMhelix	133	155
Sequence	TMHMM2.0	outside	156	200
Sequence	TMHMM2.0	TMhelix	201	223
Sequence	TMHMM2.0	inside	224	267
Sequence	TMHMM2.0	TMhelix	268	290
Sequence	TMHMM2.0	outside	291	304
Sequence	TMHMM2.0	TMhelix	305	327
Sequence	TMHMM2.0	inside	328	343



plot in postscript, script for making the plot in gnuplot, data for plot

Paralogous Family display NOTE: this display is for ORF03779

Shewanella oneidensis MR-1 Paralogous Families For ORF03779

Logged into [gsp] as <u>mlgwinn</u>

This page displays information about all paralogous families contained within a specified ORF. The user is presented with three types of information; a graphic display showing the relationship of all paralogous families to the ORF, an ORF list of all ORFs that belong to a paralogous family, and a multiple sequence alignment of all ORFs that belong to a paralogous family.



Α	С	Feat_name	Common Name	Role_id(s)	<u>11</u>	PF00528	Other Fams
	•	<u>ORF02518</u> 283 aa	sulfate ABC transporter, permease protein	143	+	+	
	•	<u>ORF02519</u> 293 aa	sulfate ABC transporter, permease protein	143	+	+	
	•	<u>ORF02772</u> 226 aa	molybdenum ABC transporter, permease protein	143	+	+	
	•	<u>ORF03459</u> 245 aa	molybdenum ABC transporter, permease protein	143	+	+	
	•	<u>ORF03779</u> 289 aa	sulfate ABC transporter, permease protein	143	+	+	
	•	<u>ORF03783</u> 281 aa	sulfate ABC transporter, permease protein	143	+	+	
	•	<u>ORF00271</u> 343 aa	peptide ABC transporter, permease protein	142		+	
		ODDO0070					

Evidence picture from ORF01166

Additional evidence types shown here are signal P, lipoprotein predictions, and PROSITE hits. Signal P and PROSITE information are displayed both in the Evidence Picture and in sections of their own on the Gene Curation Page (next slide). Clicking on the bars in the graphic opens windows with additional information.

Lipoprotein predictions are based on one particular PROSITE motif, so clicking on the red lipoprotein bar will take you to the PROSITE page for the lipoprotein signature (not shown in tutorial).



NOTE: this display is for ORF01166

Gene Curation Page - PROSITE and Signal P sections on the GCP

NOTE: this display is for ORF01166

Click here to see info on PROSITE motif.

PROSITE				submit 🕒
PS01039: Bacte	rial extracellular so	olute-binding p	roteins, family 3 sig	nature.
Match sequence	e: GFDIELAKQIAK	DA		
Coords	Precision	Recall	Curation	
52/65	0.76	0.93	V	[Add To GO Evidence]
ATTRIBUTES				submit 🗈
No Frameshifts I	Detected.			

SIGNAL_P		submit 🗈
SignalP-2.0 Results: [Graphica	I Display] [Raw output for SP-HMM/NN]	
SignalP-2.0 HMM		
Prediction	No prediction generated 🗾 🗆 Curated	
Signal peptide probability	0.984	
Signal anchor probability		
Max cleavage site probability	0.340	
	Y	
		60

Click here to see output in graphical form. ⁶⁸



Tiosed by Nebe ob Millor sites. Honvia Canada Cillina Switzenand Tarwan	
The Korean ExPASy site, kr.expasy.org, is temporarily not available. Search PROSITE for Go Clear	Зу
NiceSite View of PROSITE: PDOC00798 NOTE: this display is for ORF011	66
<pre>(documentation) Bacterial extracely family 3 signature PROSITE crossreference() PSO1032 SBP BACTERIAL 3 Documentation Bacterial high affinity transport failing to external sites of the integral membrane proteins of the estimate integral membrane integral membrane proteins of the estimate integral membrane integral membr</pre>	e by d for ch).

Gene Curation Page (ORF04813) - Gene Ontology Display

Link to GO Current GO term assignments are search tool Link to GO listed in table. suggestions -Click id # to see term in tree. -Click box for GO term to be GENE ONTOLOGY deleted. Ŀ submit | go sug | search -Click "add" to add additional delete goid date evidence assigned evidence rows. (or click delete and ISS: PMID: 12368813 with TIGR_TIGRFAMS: TIGR00433 add to completely redo evidence) 07/29/04 GO:0004076 [add] [edit] (F) biotin synthase activity mlawinn -Click "edit" to edit evidence. ISS: PMID: 12368813 with TIGR_TIGRFAMS: TIGR00433 GO:0009102 [add] [edit] (P) biotin biosynthesis mlgwinn 07/29/04 -"Make ISS" (not seen in this example) can be used when the GO term and evidence assigned by AutoAnnotate are correct, clicking function component process this button marks the old association for deletion and ۲ ▼ automatically puts in the new info for insertion. add go id vith evcode reference qualifier These pull downs have commonly TIGR_CMR:annotation -ISS T used GO terms. If you choose the unknown terms from any pull-down, ISS TIGR_CMR:annotation -V ◄ the evidence will automatically fill in ISS TIGR_CMR:annotation -T ◄ (since it is always the same.) ISS TIGR_CMR:annotation -T ▼ Fill in the fields in this section to add ▼ TIGR_CMR:annotation ▼ ISS ▼ or change GO term assignments. All entries must have "ev code". "reference', and "with". 71 (more on this in a minute.....)

Overview of steps in GO annotation:

- -Review the GO terms assigned to the gene by AutoAnnotate (if any). If they are correct and sufficient use the "Make ISS" button. (not seen here)
- -Look for any other needed GO terms in the various suggestion areas on the page: EC#s, HMMs, GO suggestions (see suggestion slide for more info)
- -If correct GO terms are unavailable on the Gene Curation Page go to the GO search pages and find the GO term you need. You get there by clicking the search link in the upper right corner of the GO section.
- -GO terms must be added in the bottom part of the Gene Ontology section. The GO term id goes in the "add go id" column
- -The ec_code column has a pull-down for choosing the ev_code you want, the default is "ISS"
- -Next is the "reference", "with", and "qualifier" columns. Additional slides following this one detail the search for and insertion of GO terms and evidence.
- -See the "overview" presentation for more info on GO

GENE	ONTO	LOGY						submit	l go sug	search	🛙
delete	lelete go id				assigned	date	evidence				
Γ	GO:000	4076 <u>(add</u>] [<u>edit]</u>	(F) biotin synthase	eactivity	mlgwinn	07/29/04	ISS: PWID: 12368813 with TIGR_TIGRFAMS:TIGR00433			
Г	GO:000	9102 <u>(add</u>] [<u>edit]</u>	(P) biotin biosynth	esis	mlgwinn	07/29/04	ISS: PMID: 12368813 with TIGR_TIGRFAMS: TIGR00433			
								1			
	function		p	rocess			comp	onent			
	_	۲		v		_			T		
		_	1	_	1				_		
add go	oid	ev code		referen	ce			vith			qualifier
		SS 🔹	TIGR_CI	AR:annotation	T			T			•
		SS 💌	TIGR_CI	R:annotation	T			V	_		
	l	SS 🔹	TIGR_CI	AR:annotation	T			T			•
		SS 🔹	TIGR_CI	R:annotation	T			T			
		SS 🔻	TIGR_C	R:annotation	T			Y			•
Gene Curation Page - GO suggestions and Auto-fill-ins

GO term suggestions and auto-fill-in buttons are located in several places on the Gene Curation Page:

-GO terms assigned to HMMs are listed under HMM hits (if any have been assigned - see the HMM slide for how these look). These are often excellent sources for GO terms. Clicking the "Add" button next to a GO term under an HMM adds both the term id and the evidence to the appropriate fields in the GO entry section. Clicking the "Add to GO evidence" button adds just the HMM accession into the "with" field in the GO entry section.

-GO terms corresponding to EC numbers are listed next to the EC box (for enzymes). Clicking the "add" button will put the GO term id into the "add go id" fields in the GO entry section.

-GO terms assigned manually to other bacterial genomes (V. cholerae, B. anthracis both a Gram + and Gram - representative), InterPro hits, Genome Properties are listed both at the bottom of the page and in a pop-up window accessed by the link in the upper right corner of the GO section. Clicking on "add" in this section puts the GO id into the "add go id" fields in the GO entry section.

-"Add to GO evidence" buttons are also available for Prosite hits, this populates the "with" field with the Prosite accession. Available when a protein has matches to Prosite. -"Add to GO evidence" is also available for the characterized match accession, this will put the accession of the characterized matching protein into the "with" field entry box. -"Add to GO evidence" is also available for Genome Properties, clicking on the "GO" link under the "add to GO evidence" column in the Genome Properties section will enter the GenProp accession in the "with" field.

See next page for screen shots.



Manatee's GO ontology and annotation search tool:

In many cases the GCP will not have a suggested GO term that meets an annotators needs. In that situation the annotator will turn to Manatee's built in GO ontology and annotation browser. There are several available functions in the search tool:

-GO term id search of ontologies - this returns a tree view of the search term in the ontology

-GO term name keyword search of ontologies - this returns a table of terms where the name, or a synonym of the name contains the keyword or where a word contains the keyword in question.

-protein name keyword search of annotations - this is a search of the annotations and returns proteins whose name match the keyword and the GO terms that were assigned to those proteins

-GO id search of annotations - search GO annotations with a GO id and see a list of proteins that have been annotated to that GO term.

-GO correlations in annotations - often a particular function term will often be assigned with a particular process term (for example: "biotin synthase" will almost always be assigned in conjunction with "biotin biosynthesis") - when one needs help finding a process or function one can search for these relationships with the correlations tool.

-EC search - uses the ec2go mapping file provided on the GO web site to look up GO terms that correspond to EC numbers

-GO BLAST - search a protein sequence against a database of proteins that have been annotated to GO, then link to the GO terms that were assigned to them. This is the only GO search tool not accessible on the GO search page - this one is found in the "Select Function" pull-down menu at the top of the Gene Curation Page.

See next page for screen shots.

Links from the Gene Curation page - GO Search Tool

Click on the "search" link in the title bar of the GO data section

Input a GO term here. This results in a GO term information page.

		In	Input an EC number and get the corresponding GO term					
sea	arch GO from GO id:	S	search GO from EC number:		search GO from keywords		search GO associations with	n gene nam
→ GO id:	submit reset	EC #:	submit reset		Exact Match			
search go_	gene_association from GO id:	5	search for GO correlations:	S	submit reset		submit Reset	
► GO id:	submit reset	GO id:	submit reset	۲ ا	Molecular Function			-
		Search	n TIGR prokaryotic data only	۲. ۲	Cellular Component			
	Search for GO te protein along with restrict the search	erms t n the h to T	hat most frequently input GO term. Che IGR prokaryotic dat	are eck ta c	e assigned to a the box to only.			
Input a genes databa	GO id and see from other ses that have beer	1	Input a search string select some or all of want to restrict your r of your input text, clic This searches the na	hei the resi ck " ame	re using the checkbo ontologies to search ults to only terms wh Exact match". of the GO term.	n i ic	es to n. If you h share all	
annota term	ited with that GO				Input text and see to genes from othe common names co	G r (O terms assigned databases whose tain the input text.	76

GO Term information page and tree view. This page is reached by clicking on GO id links or using GO id search.

Name: biotin synthase activity

Type: molecular_function

Definition: Catalysis of the reaction\: dethiobiotin + sulfur = biotin.

Comment: NONE

Synonym: NONE

Secondary ID: NONE

EC Number: 2.8.1.6

Absolute Path in GO Tree: 1 instance detected

```
+Ontology (TI:0000001)[R]
+Gene_Ontology (GO:0003673)[P]
+molecular_function (GO:0003674)[P]
+catalytic activity (GO:0003824)[I]
+transferase activity (GO:0016740)[I]
+transferase activity, transferring sulfur-contain
+sulfurtransferase activity (GO:0016783)[I]
biotin synthase activity (GO:0004076)[I]
```

View Mode: Regular

Numbers next to the terms in the tree indicate the number of genes from this organism that are annotated to that term or a child of that term - clicking on the number gives you a table of those genes and relevant info. (missing in this screen shot)

If you reach this page by clicking on a GO term on a GCP, clicking the "add" button in the tree will place that GO term in the "add" field on the GCP.

```
+Ontology (TI:0000001)[R] [add]
    +Gene Ontology (G0:0003673)[P] [add]
         +molecular function (G0:0003674)[P] [add]
               +catalytic activity (G0:0003824)[I] [add]
                    +transferase activity (GO:0016740)[I] [add]
                         +transferase activity, transferring sulfur-containing groups (G0:0016782)[I] [add]
                               +sulfurtransferase activity (G0:0016783)[I] [add]
                                      biotin synthase activity (GO:0004076)[I] [add]
                                      cysteine desulfurase activity (GO:0031071)[I] [add]
                                      3-mercaptopyruvate sulfurtransferase activity (GO:0016784)[I] [add]
                                      tRNA sulfurtransferase activity (GO:0016227)[I] [add]
                                      thiosulfate-thiol sulfurtransferase activity (G0:0050337)[I] [add]
                                      thiosulfate-dithiol sulfurtransferase activity (GO:0047362)[I] [add]
                                      thiosulfate sulfurtransferase activity (G0:0004792)[I] [add]
                               +transferase activity, transferring alkylthio groups (G0:0050497)[I] [add]
                               +CoA-transferase activity (GO:0008410)[1] [add]
                         +sulfotransferase activity (G0:0008146)[I] [add]
+transferase activity, transferring phosphorus-containing groups (G0:0016772)[I] [add]
                           pyruvyltransferase activity (G0:0046919)[I] [add]
                           CDP-alcohol phosphotransferase activity (G0:0008414)[I] [add]
                           trichothecene 3-0-acetyltransferase activity (G0:0045462)[I] [add]
                          +transferase activity, transferring alkyl or aryl (other than methyl) groups (GO:0010
                          +transferase activity, transferring glycosyl groups (GO:0016757)[I] [add]
                          +2'-phosphotransferase activity (GO:0008665)[I] [add]
                         +glucanosyltransferase activity (G0:0042123)[I] [add]
+transferase activity, transferring nitrogenous groups (G0:0016769)[I] [add]
                          +transferase activity, transferring selenium-containing groups (GO:0016785)[I] [add]
                           mannosylphosphate transferase activity (GO:0000031)[I] [add]
                          +transferase activity, transferring one-carbon groups (G0:0016741)[I] [add]
                           cobinamide phosphate guanylyltransferase activity (G0:0008820)[I] [add]
                         +transferase activity, transferring aldehyde or ketonic groups (G0:0016744)[I] [add]
lipoyltransferase activity (G0:0017118)[I] [add]
                          +transferase activity, transferring acyl groups (GO:0016746)[I] [add]
                           S-adenosylmethionine:tRNA ribosyltransferase-isomerase activity (G0:0051075)[I] [add
                           lauroyl transferase activity (GO:0008913)[I] [add]
```

```
77
```

Search results for GO term keyword: "biotin"

The first part of the table shows results from the GO term names.

The second part of the table shows results from GO term synonyms.

Note that areas of the text which matched the keyword are highlighted in red by Manatee.

Terms which are "obsolete" or "secondary" to another term will have that indicated in column one.

Click any GO term id number for a view of the term in the GO tree.

tern hits		
GO id	type	name
GO:0009374	molecular_function	biotin binding
GO:0009102	biological_process	biotin biosynthesis
GO:0042966	biological_process	biotin carboxyl carrier protein biosynthesis
GO:0004075	molecular_function	biotin carboxylase activity
GO:0009343	cellular_component	biotin carboxylase complex
GO:0042367	biological_process	biotin catabolism
GO:0006768	biological_process	biotin metabolism
GO:0004076	molecular_function	biotin synthese activity
GO:0015878	biological_process	biotin transport
GO:0015225	molecular_function	biotin transporter activity
GO:0047707	molecular_function	biotin-CoA ligase activity
GO:0004077	molecular_function	biotin-[acetyl-CoA-carboxylase] ligase activity
GO:0004078	molecular_function	biotin-[methylcrotonoyl-CoA-carboxylase] ligase activity
GO:0004079	molecular_function	biotin-[methylmalonyl-CoA-carboxytransferase] ligase activity
GO:0004080	molecular_function	biotin-[propionyl-CoA-carboxylase (ATP-hydrolyzing)] ligase activity
GO:0000106 secondary	molecular_function	biotin-apoprotein ligase activity
GO:0018271	molecular_function	biotin-protein ligase activity
GO:0047708	molecular_function	biotinidase activity
GO:0019351	biological_process	dethiobiotin biosynthesis
GO:0046450	biological_process	dethiobiotin metabolism
GO:0004141	molecular_function	dethiobiotin synthase activity
GO:0018054	biological_process	peptidyl-lysine biotinylation
GO:0009305	biological process	protein amino acid biotinylation

go_s ynonyn hits	jo_synonym hits					
GO id	type	name	synonym			
GO:0004079	molecular_function	biotin-{methylmalonyl-CoA-carboxytransferase] ligase activity	biotin-1			
GO:0018271	molecular_function	biotin-protein ligase activity	biotin-apoprotein ligase activity			
GO:0019351	biological_process	dethiobiotin biosynthesis	desthio <mark>biotin</mark> biosynthesis			
GO:0046450	biological_process	dethiobiotin metabolism	desthio <mark>biotin</mark> metabolism			

GO correlations search

Search results from query with GO:0004076 "biotin synthase activity"

Searches data set stored in our database of all associations to genes available on GO web site. First table shows percentages of occurrences of the query term with other terms. Second table shows details of all instances of query term assigned to a gene in the data set.

GO Correlations: GO:0004076 biotin synthase activity Logged into [gsp] as mlgwinn GO id correlation percentage GO name GO type 88 89 % GO:0009102 biotin biosynthesis P С 11.11 % GO:0008372 cellular_component unknown gene id gene db GO id GO type GO name gene name S0003518 SGD GO:0008372 C cellular component unknown biotin synthase S0003518 SGD biotin synthase GO:0009102 Р biotin biosynthesis BA4336 gba_TIGR biotin synthetase GO:0009102 Р biotin biosynthesis SO3925 gsp_TIGR biotin synthase family protein GO:0009102 Ρ biotin biosynthesis SO2740 Р gsp_TIGR biotin synthase GO:0009102 biotin biosynthesis CBU1007 geb TIGR GO:0009102 Р biotin biosynthesis biotin synthase VC1112 Р gvc_TIGR biotin synthase GO:0009102 biotin biosynthesis Ρ GSORF1608 ggs_TIGR GO:0009102 biotin biosynthesis biotin synthetase At2g43360 TIGR Ath1 biotin synthase (BioB) (BIO2) Ρ GO:0009102 biotin biosynthesis

Output from GO search for protein common name keyword: biotin synthase

gene id	role id	gene symbol	EC#	GO id (GO type)	ev code	GO term	gene name
SO3925	157			GO:0004076 (F)	ISS	biotin synthase activity	biotin synthase family protein
SO3925	157			GO:0009102 (P)	ISS	biotin biosynthesis	biotin synthase family protein
SO2737	77	bioD	6.3.3.3	GO:0004141 (F)	ISS	dethiobiotin synthase activity	dethiobiotin synthase
SO2737	77	bioD	6.3.3.3	GO:0009102 (P)	ISS	biotin biosynthesis	dethiobiotin synthase
SO2740	77	bioB	2.8.1.6	GO:0004076 (F)	ISS	biotin synthase activity	biotin synthase
SO2740	77	bioB	2.8.1.6	GO:0009102 (P)	ISS	biotin biosynthesis	biotin synthase
CBU1007	77	bioB	2.8.1.6	GO:0004076 (F)	ISS	biotin synthase activity	biotin synthase
CBU1007	77	bioB	2.8.1.6	GO:0009102 (P)	ISS	biotin biosynthesis	biotin synthase
VC1112	77	bioB	2.8.1.6	GO:0004076 (F)	ISS	biotin synthase activity	biotin synthase
VC1112	77	bioB	2.8.1.6	GO:0009102 (P)	ISS	biotin biosynthesis	biotin synthase
GSU1583	77	bioD	6.3.3.3	GO:0004141 (F)	ISS	dethiobiotin synthase activity	dethiobiotin synthase
GSU1583	77	bioD	6.3.3.3	GO:0009102 (P)	ISS	biotin biosynthesis	dethiobiotin synthase
P32451				GO:0005739 (C)	IDA	mitochondrion	Biotin synthase
Q84QK2				GO:0004076 (F)	IEA	biotin synthase activity	Putative biotin synthase
Q84QK2			-	GO:0005739 (C)	ISS	mitochondrion	Putative biotin synthase
Q84QK2				GO:0009102 (P)	IEA	biotin biosynthesis	Putative biotin synthase
At2g43360				GO:0004076 (F)	IGI	biotin synthase activity	biotin synthase (BioB) (BIO2)
At2g43360				GO:0009102 (P)	TAS	biotin biosynthesis	biotin synthase (BioB) (BIO2)
SPCC1235.02				GO:0004076 (F)	IEA	biotin synthase activity	biotin synthase activity
SPCC1235.02				GO:0006731 (P)	IEA	coenzyme and prosthetic group metabolism	biotin synthase activity
SPCC1235.02				GO:0006790 (P)	IEA	sulfur metabolism	biotin synthase activity
S0003518				GO:0004076 (F)	TAS	biotin synthase activity	biotin synthase
S0003518				GO:0005739 (C)	IDA	mitochondrion	biotin synthase
S0003518				GO:0009102 (P)	TAS	biotin biosynthesis	biotin synthase

Step 1. Pick an evidence

code. Most genes in bacterial genome sequencing projects will get an ev_code of "ISS". This stands for "Inferred from sequence similarity." If a gene from the sequenced organism has had experimental

characterization, then chose an appropriate experimental ev_code. All "unknown" GO terms get "ND" as ev_code. To see all ev_codes, click the "ev_code' link.

Step 2. Fill in "reference" information. For ISS terms prior to publication use "TIGR_CMR:annotation", after publication use the

Adding GO Evidence

GENE	ENE ONTOLOGY submit go sug search																
delete	elete goid a			assigned by	assign date	evidence	;										
Γ	GO:00	04076	add	edit	(F) biotin sy	F) biotin synthase activity		mlgwinn	03/29/04	ISS: PMID ISS: PMID	SS: PMID: 12368813 with Swiss-Prot: P12996 SS: PMID: 12368813 with TIGR_TIGRFAMS: TIC)96 5:TIGR0043	3		
	GO:00	09102	add	edit	(P) biotin biosynthesis		mbeanan	11/15/01	ISS: PMID ISS: PMID	PMID:12368813 with Swiss-Prot:P12996 PMID:12368813 with TIGR_TIGRFAMS:TIG		996 5:TIGR0043	3				
	function process component																
_						_							_				
I			_										<u> </u>				
add go	o i d	ev co	de		re	ference					with					qual	lifier
		ISS	•	TIGR_CN	/R:annot	ation 💌	Ī				•						•
		ISS	•	TIGR_CN	/R:annot	ation 💌	Ī				•			_			•
		ISS	•	TIGR_CM	/R:annot	ation 👤	·				•	_				_	•
		ISS	•	TIGR_CN	/R:annot	ation 💌	·				•	_				_	•
		ISS	•	TIGR_CN	/R:annot	ation 👤	Ī				•						•

PMID of the genome paper. For "unknown" terms use "GO_REF:nd". For terms with experimental evidence codes use the PMID of the paper describing the characterization.

Step 3. Fill in the "with" field. For all ISS entries you must fill in the accession of the HMM, characterized match protein, or Genome Property that led to the annotation.

Auto Fill-ins: Both the GO ids and associated evidence can be filled in automatically by clicking the "Add" buttons next to GO suggestions and the "Add to GO evidence" buttons. All info for the "unknown" terms is filled in automatically by choosing the "unknown" terms in the pull-down menus. All information for GO terms assigned to HMMs is filled in with the "Add" buttons next to GO terms under HMMs.

Qualifier should be set to "contributes_to" when annotating the function of a complex to the proteins ⁸¹ that make up the complex. (see the overview for more informatin on all of these fields)

Gene Curation Page - TIGR roles



Gene Curation Page - How to get the data into the database: The "Submit" buttons



<u>Gene Curation Page - The pull down menus</u>

If you click on the select pull down menus you will get a selection of options. Each of these when selected will generate a new page with the desired information. (Later slides show examples of some of these.)

Shewanella oneidensis MR-1 Gene Curation Page	Home Logged into [gsp] as mlgwinn
GENE CURATION INFORMATION	
ORF04813 (SO2740) View BER Searches asmbl_id: 7974 • Reload Page	end5/end3: 2856763 / 2855711 database: gsp gene length: 1053 feat_name / locus:
Select Display ▼ Select Display Genome Region GEI View Sequences 3rd Position GC Skew sene Signal Peptide Prediction bioti Transmembrane Helix Prediction Secondary Structure Prediction	Select Function Refresh Searches Select Function Select Function Edit Start Sites Submit history Frameshift Report Submit history Translation Exception Blast Gene Against GO
gene_sym: bioB EC number(s): 2.8.1.6	EC GO suggestions: GO:0004076 add biotin synthase activity (F)
comment: Start confidence Low	pub_comment:
▶ nt_comment	▶auto_comment

Links from the Gene Curation Page - View sequence

This page shows the length (nucleotide and protein), coordinates, MW, and pl of the protein.

Also, in fasta format are the nucleotide and protein sequences.

Shewanella oneidensis | Sequence Display for O

This page displays the feat_name, nucleotide sequence and the amino acid set

length: 1050 nucleotides protein length: 350 amino acids MW: 38790.13 pl: 4.9477

Genomic sequence

>ORF04813

ATGTCGCAGTTGCAAGTTCGTCATGATTGGAAGCGGGAAGAAATCGAAGCCTTATTTGCG CTGCCGATGAATGACTTATTATTTAAAGCCCACAGTATCCACCGTGAAGAGTACGATCCT AACGAAGTGCAGATCAGCCGCTTATTGTCGATCAAAACTGGGGGCTTGTCCTGAGGATTGT AAATATTGTCCGCAGAGTGCGCGTTACGACACTGGCCTTGAAAAAGAGCGTCTCTTAGCG ATGGGCGCCGCTTGGCGTAACCCCGAAAGATAAAGATATGCCATACCTCAAGCAAATGGTG CAAGAGGTGAAAGCCCTCGGCATGGAAACCTGTATGACCTTAGGGATGTTAAGTGCCGAG CAAGCCAATGAGTTGGCCGAAGCAGGCCTTGACTATTACAACCACAATTTAGATACCTCG CCTGAATACTACGGCGATGTGATCACCACCCGTACCTATCAAAACCGCTTAGATACCTTA AGCCATGTGCGCGCATCGGGCATGAAAGTTTGCTCTGGCGGCATTGTCGGCATGGGCGAG AAGGCTACTGACAGAGCCGGTTTATTACAACAACTGGCTAATTTACCCCCAGCATCCGGAT TCTGTGCCGATCAATATGTTAGTCAAAGTAGCGGGTACCCCCTTTGAAAAACTTGATGAT TTAGATCCACTCGAGTTTGTCCGAACCATCGCCGTGGCGCGTATTTTAATGCCACTGTCG CGGGTGCGTTTATCCGCAGGCCGTGAAAATATGAGCGATGAACTGCAGGCCATGTGTTTC TTT6C666C6C6AACTC6ATTTTTAC66CT6TAA6TTACT6ACCAC6CCCAACCCC6AA GAAAGTGATGATATGGGGTTGTTCCGTCGCCTGGGTTTACGCCCTGAGCAGGCGCGCAGCC **GCCTCTATTGATGATGAGCAAGCGGTATTAGCTAAAGCTGCGGCTTATCAAGATAAAGCT** TCAGCTCAGTTTTATGATGCGGCGGCACTA

CDS

>ORF04813

ATGTCGCAGTTGCAAGTTCGTCATGATTGGAAGCGGGAAGAAATCGAAGCCTTATTTGCG CTGCCGATGAATGACTTATTATTTAAAGCCCACAGTATCCACCGTGAAGAGTACGATCCT AACGAAGTGCAGATCAGCCGCTTATTGTCGATCAAAACTGGGGCCTTGTCCTGAGGATTGT AAATATTGTCCGCAGAGTGCGCGTTACGACACTGGCCTTGAAAAAGAGCGTCTCTTAGCG ATGGGCGCCGCTTGGCGTAACCCGAAAGATAAAGATATGCCATACCTCAAGCAAATGGTG CAAGAGGTGAAAGCCCTCGGCATGGAAACCTGTATGACCTTAGGGATGTTAAGTGCCGAG CAAGCCAATGAGTTGGCCGAAGCAGGCCTTGACTATTACAACCACAATTTAGATACCTCG CCTGAATACTACGGCGATGTGATCACCCACCCGTACCTATCAAAACCGCTTAGATACCTTA AGCCATGTGCGCGCATCGGGCATGAAAGTTTGCTCTGGCGGCATTGTCGGCATGGGCGAG AAGGCTACTGACAGAGCCGGTTTATTACAACAACTGGCTAATTTACCCCCAGCATCCGGAT TCTGTGCCGATCAATATGTTAGTCAAAGTAGCGGGTACCCCCTTTGAAAAACTTGATGAT TTAGATCCACTCGAGTTTGTCCGAACCATCGCCGTGGCGCGTATTTTAATGCCACTGTCG CGGGTGCGTTTATCCGCAGGCCGTGAAAATATGAGCGATGAACTGCAGGCCATGTGTTTC TTTGCGGGCGCGAACTCGATTTTTACGGCTGTAAGTTACTGACCACGCCCAACCCCGAA GAAAGTGATGATATGGGGTTGTTCCGTCGCCTGGGTTTACGCCCTGAGCAGGGCGCAGCC GCCTCTATTGATGATGAGCAAGCGGTATTAGCTAAAGCTGCGGCTTATCAAGATAAAGCT TCAGCTCAGTTTTATGATGCGGCGGCACTA

Protein

>ORF04813

MSQLQVRHDWKREEIEALFALFMNDLLFKAHSIHREEYDPNEVQISRLLSIKTGACFEDC KYCPQSARYDTGLEKERLLAMETVLTEARSAKAAGASRFCMGAAWRNPKDKDMPYLKQMV QEVKALGHETCMTLGHLSAEQANELAEAGLDYYNHNLDTSPEYYGDVITTRTYQNRLDTL SHVRASGHKVCSGGIVGMGEKATDRAGLLQQLANLPQHPDSVPINHLVKVAGTPFEKLDD LDPLEFVRTIAVARILMPLSRVRLSAGRENMSDELQAMCFFAGANSIFYGCKLLTTPNPE ESDDMGLFRRLGLRPEQGAAASIDDEQAVLAKAAAYQDKASAQFYDAAAL

Links from the Gene Curation Page - Third position GC skew

In organisms whose DNA has a high GC content it can sometimes be helpful to look at third position GC skew to help resolve overlaps.

Due to the nature of the genetic code, the third position is the least constrained of a codon and therefore will be able to reflect the higher GC content of the overall genome. Therefore one should see a markedly higher GC content in the third position of the correct frame.





ORF Management in Manatee: Genome Viewer

Refresh XML Search Asmbl Id: 7974	Database: gsp		database:	asmbl_id:	submit reset
feat name end5 end3 rol	e id ec num ger	ne sym complete	com name		
Six Frame Options (on six frame clicks)		Gene	Options (on feat_na	ame clicks)	
View Sequence · Blast · Insert Gene	View Sequence	Annotate ORF .	Edit Start C Blast	ORF C Merge Genes	○ Delete Gene ○
		· · ·			
· · · · · · · · · · · · · · · · · · ·	î			· · · · ·	î
0.0kb 1.0kb 2.0kb	3.0kb	4.0kb	5.0kb	6.0kb 7.0kb	8.0kb
THE A DAY MUTCHES . IN A SUBJECT WAS A SUBJECT AND	lle likitaa mula mul	ւ Մերուս Մենսին հենունես։ Սու	. I. Shill also have been the state	n and all to be following a subsc	
Laukadan taal alanti tatu bi tatu an in bi		<u> 11 14 10 1 11 4</u> 11 1 10 11 11 1 1	h II. N. Li Latin di An	<u> 18 a 18 an - Int Dat Databilit II 1</u>	
ORF 02394					0RF 02389
005 02395			-4	05602390	
			ببالطبابات التعاقياك ابت		
<u>1111 - I I al al al 1100 - I I al 1100 - II I I al 1000 - II I I I I al 1000 - II I I I I I I I I I I I I I I I I</u>	0RF 02393		ul II.I. III.I.I.I.		
ULUE: (1.101001) 10110 (1.1100) (1.1100)	II. III kole kitaan takaka kaanaadhadd	de la la la caractería de	0RF 02392 0RF 0	2391 ORF08000	

Clicking on the "Genome Viewer" option on the "Welcome to Manatee" Page, selecting "GV" next to a gene id in a gene list, or selecting "Genome Region" in the "Select Display" pull-down on the GCP will take you to our Genome Viewer tool. Here you can view the genes from the whole genome in relation to each other, edit their starts, merge them, insert new genes, and delete genes. Mousing over the genes fills in the information boxes near the top of the display with coordinates, com_name, etc.

Search		
Coordinate:		Search
Lower Coordinate:	Upper Coordinate:	Search
feat_name:		Search

To get to a specific region of the genome, enter coordinates or feat_name in the search section at the bottom of every Genome Viewer page. One can also use the pulldown menu on the Gene Curation Page or the "GV" link on a gene list page.

ORF Management in Manatee - 6-frame analysis

	Refresh XMI	Search	Asmbl Id: 7974	Database: gsp		database:	asmb	l_id:	submit	reset
	feat name	end5 2861863	end3 role i 2858294	<u>d ec num ger</u>	ne sym complete	com name				
ſ	Six Frame	Options (on s	ix frame clicks)		Gen	e Options (on fea	at_name clicks	5)		
l	View Sequence	ce 🖲 Blast 🤇	ି Insert Gene ି	View Sequence C	Annotate ORF @	Edit Start C E	Blast ORF C	Merge Genes	C Delete G	iene O
	51mb	2.852mb	2.853mb	2.854mb	2.855mb	2.856mb	2.857mb	2.858mb	2.859	Mb
	ORF 04820						ORF 0	4812		
							n hulla lidandal a u b			
	واللحية المسالية	ala di kasil i								
						na na tali na ra 1000 r			na na hini na hili ni n	
	111 <u>111111111111111</u>					<u> , </u>	ابليا ار وابيريناي			
	1 10-0-1 10-00 1000 1-		Allik III williad di U. alaan I. Uada da da	والمراقبة المتعاولة ومعاقدا والمراجع	dhar dhi bi dhe bi bir	0RF04813				

To analyze regions in the 6-frame translation (options boxed in pink), click on the button for the activity you wish and then click in the open reading frame.

"View sequence" gives you the nucleotide and amino acid sequence of the ORF.

"Blast" Blasts the ORF.

"Insert Gene" inserts the ORF. (see later slide for more on this)

ORF Management in Manatee: gene adustments

Refresh XML Search Asmbl Id: 7974	Database: gsp	database: asmbl_	id:submitreset
feat_name end5 end3 role 2861863 2858294	id ec num gene sym complete	com_name	
Six Frame Options (on six frame clicks) View Sequence Insert Gene	Gene (View Sequence C Annotate ORF •	Dptions (on feat_name clicks) Edit Start	erge Genes C Delete Gene C
51mb 2,852mb 2,853mb	1 2.855mb 2.855mb 2.8	56mb 2.857mb	1 2.858mb 2.859mb
ORF 04520 III III ORF 04520 III III IIII IIII IIII IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII			
d <u>L − 4 − 1 − 11 − 11 − − 1 − − 11 − 10 − 1 − 1</u>			· · · · · · · · · · · · · · · · · · ·

To make adjustments to existing genes in the database, click on the option you want to do, then click on the arrow for the gene of interest.

New pages will pop up with information specific for your request. (see later slides)

ORF Management in Manatee: start edits

This page can be reached from the "Gene Options" menu on the Genome Viewer page or from the Gene Curation Page "Select Function pull-down by selecting "Edit Start Site

Purple text represents the gene of interest. Blue text represents other genes in the region. It is important nc to introduce overlap with other genes when changing a start site. Often editing a start site will remove overlap between two genes. Occasionally an annotator may want to extend a start into an upstream gene and will find tha the upstream gene in question is a small hypothetical with no homology to anything. In such a case the annotato should consider deleting the short hypothetical, since it becomes likely that it is not a real gene.

To edit a start, click on the start you want in the 6-frame representation. The new coordinate for the selected start site will appear in the "New End5 box. To save the change to the database, click "Submit".

		Region of 2856886 to 2858271	
on		Green nucleotides are ribosome binding sites	
ite		Describer indexedues are noosonic ontaing sites.	
		Purple bold amino acids are amino acids of your query gene.	
		Blue bold amino acids are amino acids of one of genes in the region (other than query gene).	
		Click on a new start to change the start site, then click submit to enter change to the database.	
•		East name: ODE04912	
		Peat_name: OPF04012	
nc			
es		5' End: 2856886 3' End: 28582/1 New End5: 285886	
00		Submit	
lar		Start Edits	
~~	2856725	TTCCCGCTTCCAATCATGACGAACTTGCAACTGCGACATTGAACACCCCTTTTATTTTGT	Nucleotide
an		<u>F P L P I M T N L O L R H * T P F Y F C</u>	Frame 1
rt		<u>S R F Q S * R T C N C D I E H P F I F V</u>	Frame2
uι		PASNHDELATATLNTLLFLY	Frame3
tha		<u>K G S G I M Y F K C S R C Q Y G K * K Q</u>	Frame4
u ie		<u>R G A E L * S S A V A V N F V R K N K</u>	Frame5
		<u>E R K W D H R V O L O S M S C G K I K T</u>	Frame6
	2856785	ATTTTACCTTGGCTAGGATAACCTCAGCCCTTAAACTGCCAACCAGTGATACAG	Nucleotide
y tc			Frame1
oto		FILS " L N L S F " I V N A N V " I K	Frame3
alu			Frame4
		Y K V K A L I V E A R L S D V G V L S V	Frame5
			Frame6
v	2856845	GTTTACCACTGATTAATTTTCAATCAACGCTGTGAGCTTTTATGCGCAATTTACTCGATT	Nucleotide
,		V X H * L I F N Q R C E L L C A I X S I	Frame1
		FTTD * FSINAYSFYAQFTRE	Frame2
		L P L I N F Q S T L * A F M R N L L D F	Frame3
		T * W Q N I K L * R Q S S K H A I * E I	Frame4
1 I		<u>PKGSILK*DVSHAKIRLKSS</u>	Frame5
J	2055005	N Y Y S * N E I L A T L K * A C N Y R N	Frame6
	2856905	TTGACTTTGATAGCGCCCATATTTGGCACCCTTATACCTCCATGACTCGTGCACTTCCTG	Nucleotide
			Frame1
d		DFDSAHIWHPYTSMTRALPV	Frame3
		K V K I A G M N P V R I G G H S T C K R	Frame4
ICD		KSKSLAWIQCG*VEMVRASG	Frame5
		Q S Q Y R G Y K A G K Y R W S E H V E Q	Frame6

ORF Management in Manatee: other gene adjustments

Insert

Insert g	Insert gene for the following coordinates?					
end5: (693924	end3: 693500				
YES						

In all cases you will be asked to confirm your request before it is carried out.

Merge

The gene shown will be preserved (for one coordinate and annotation info). The one you enter will be used to determine the new extended coordinates and then deleted. Be sure you want this entered gene deleted!

Gene: ORF00755			
Locus: TP0640	methyl-accepting chemotaxis protein		
gene_sym: mcp2-3	ec#:	asmbl_id: 6333	
end5: 699711	end3: 701552		
Enter name of gene to merge with	gene id	GO	

Delete

Are You Sure You Want To Delete This Gene?

Gene: ORF00755			
Locus: TP0640	methyl-accepting chemotaxis protein		
gene_sym: mcp2-3	ec#:	asmbl_i	d: 6333
end5: 699711	end3: 701552	YES	

Annotation Checklist

- Look for HMM hits
 - evaluate what the HMMs are telling you exact function? family membership? domain?
- Look at BER results
 - looking for proteins in the skim which are characterized (colored backgrounds)
 - many proteins are characterized but not marked so in our tables may need to check proteins with white backgrounds to see if they are characterized
 - color coding does not indicate quality of match only that the match protein has been experimentally characterized
 - evaluate the alignment what percent ID over what length? active sites? binding sites?
 - fill in characterized match accession number (by clicking on the accession in left column)
- Look at TMHMM, SignalP, Prosite, region, etc.
- Use multiple alignment (belvu link) and tree(tree icon link) as needed to differentiate function.
- Decide what you think the protein should be named
- Fill in appropriate fields for common name, gene symbol, EC#, comment.
- Decide what GO terms you need
 - find them on the page (HMMs, EC number, GO suggestions) or with the GO search tool
 - change/remove any IEA GO annotations
 - add GO evidence from HMMs, BER, Genome Properties, Prosite, etc.
- Review TIGR role and change as needed
- Check start site
 - look in BER and at the BER generated multiple alignment (belvu link)
 - adjust if necessary using "edit start" function in pull down or in the Genome Viewer section
 - check start site box when finished curation
- Check "complete", click "submit" and your done!

External Manatee's limitations: things available in limited capacity

- Refresh Searches button will not work, but you can submit sequences for re-searching, we hope to set up an automated pipeline for this, but the system is not in place yet. Currently there is an HMM search page on the CMR that can also be used.
- SignalP in the pull-down can work if you install it locally. Or you can go to the CBS site to run it on the fly: <u>http://www.cbs.dtu.dk/services</u>
- BER tree view should work if you have Java on your machine, but may be tricky
- BER multiple alignment you can view them with belvu if you are running on Linux (what we do), or you can try other multiple alignment tools, a possibility is:
 - ftp://ftp-igbmc.u-strasbg.fr/pub/ClustalX/

External Manatee's limitations: things not available

- consistency checks
- all frameshift scripts
- translation exceptions
- intergenic region analysis
- overlap analysis
- annotation status
- hypothetical protein list

Acknowledgements

Heading up the effort: Owen White Jeremy Peterson

Prokaryotic Annotation: Bill Nelson (Team leader) Bob Dodson Bob Deboy Scott Durkin Sean Daugherty Ramana Madupu Lauren Brinkac Steven Sullivan M.J. Rosovitz Sagar Kothari Susmita Shrivastava CMR: Tanja Davidsen (Team leader) Nikhat Zafar Qi Yang

HMMs: Dan Haft Jeremy Selengut

Bioinformatics Engineers: Todd Creasy Liwei Zhou Sam Angiuoli Charles Lu Anup Mahurkar

And the many other TIGR employees present and past who have contributed to the development of these tools and to building the annotation protocols we use. Thanks also go to the funding agencies that support our work including NIH, NSF, and DOE.