

Annotation Presentation  
**Annotation Presentation**  
Week 4

**Alternative Start Codons**  
**&**  
**Novel ORFs**

You  
The Brilliant  
Student



vs.

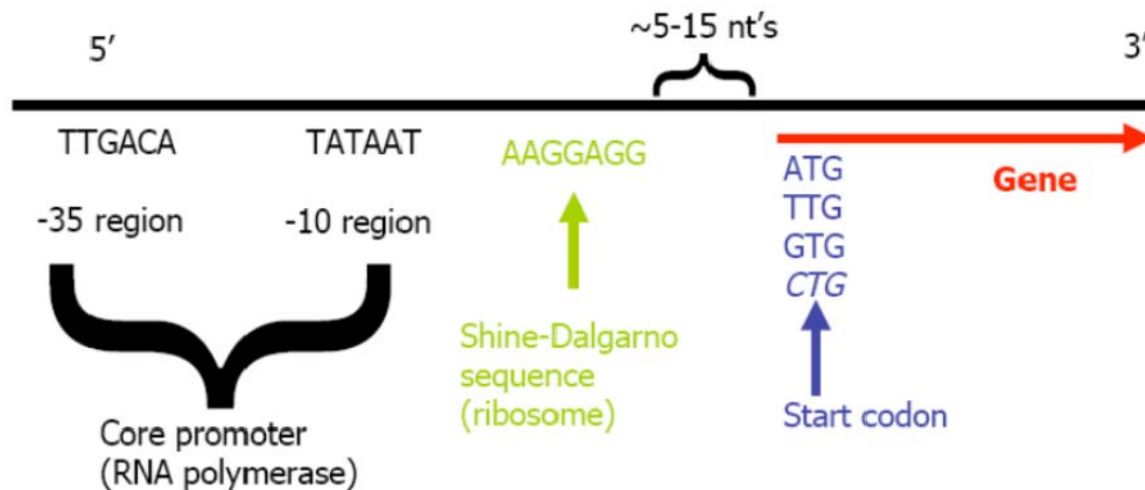


Glimmer  
The Automated  
Gene Caller

# The Automated Gene Caller (i.e., Glimmer)

What does an Automated Gene Caller do?

- Scans all nucleotides in a genome.
- Looks for “punctuation marks” that determine where genes start and stop (i.e., regulatory sequences that define where genes begin and end).



# The Brilliant Student (You)

What will you do?

Double check the work of the Gene Caller.

- ✓ Is there evidence supporting proposed start codon?
- ✓ Is there an alternative position possible for the start codon?

Examine all six reading frames.

- ✓ Is there a *novel* open reading frame (ORF) missed by Gene Caller?

Let's get started with the first task:

Verify the position for the **start codon**



# In the *imgACT* Lab Notebook...

## Basic Information Module

- [Module Instructions](#)

### IMG Gene Object ID

go to the IMG Gene Details page for the proposed gene

enter Gene Object ID (OID)

2501578154

1- Return to Gene Detail page

Quick link to Gene Details page for assigned gene (enter URL)

[http://img.jgi.doe.gov/cgi-bin/edu/main.cgi?section=GeneDetail&page=geneDetail&gene\\_oid=2501578154](http://img.jgi.doe.gov/cgi-bin/edu/main.cgi?section=GeneDetail&page=geneDetail&gene_oid=2501578154)

### DNA Coordinates

go to the **IMG Gene Details page** for the proposed gene

DNA coordinates

3948986..3951211 (+)

#### **Recall:**

You previously entered the DNA coordinates set by the automatic Gene Caller. This information was found on the Gene Detail page.

# Gene Detail Page in img/edu



INTEGRATED MICROBIAL GENOMES  
EDUCATION SITE

IMG Home | Find Genomes | Find Genes | Find Functions | Compare Genomes | Analysis Carts | MyIMG | Using IMG

Gene Search | BLAST | Phylogenetic Profilers

## Gene Detail

Loaded.

- [Gene Information](#)
- [Evidence For Function Predictions](#)
- [Sequence Search](#)
- [External Sequence Search](#)
- [IMG Sequence Search](#)
- [Homolog Display](#)

Keep the Gene Detail page open  
in separate tab while working

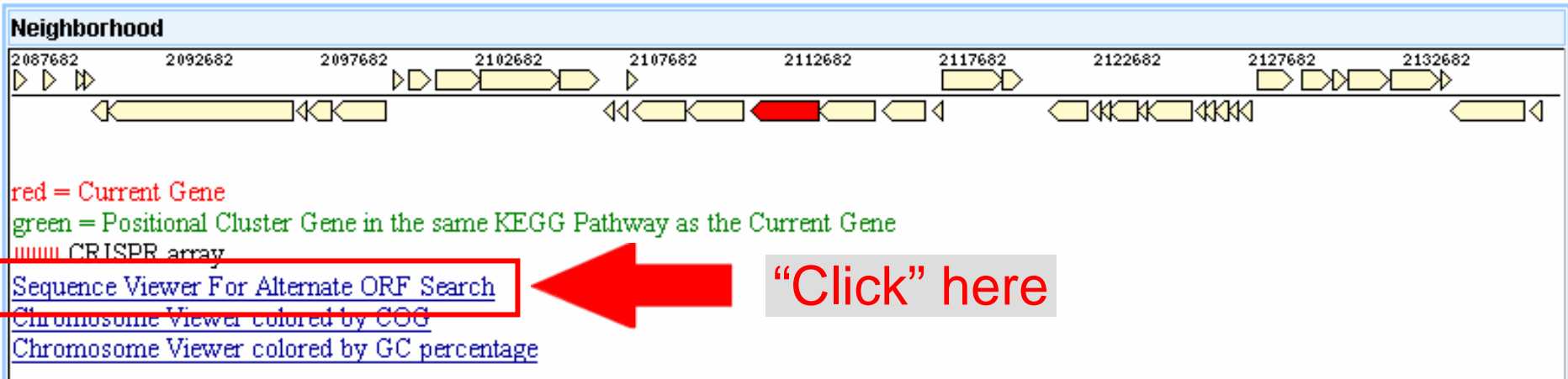
## Gene Information

Gene Information	
Gene Object ID	2500607069
Gene Symbol	
Locus Tag	PlimDRAFT_19450
Product Name	Uncharacterized anaerobic dehydrogenase, COG3383
IMG Product Source	COG3383
Genome	<a href="#">Planctomyces limnophilus DSM 3776</a>
DNA Coordinates	2111568..2113793 (-)(2226bp)
Scaffold Source	<a href="#">Planctomyces limnophilus DSM 3776 : PlimDRAFT 4083246 C168 (5423025bp)</a>
IMG ORF Type	
GC Content	0.58
External Links	
Fused Gene	No
Fusion Component	No
Protein Information	
Amino Acid Sequence Length	741aa



Scroll down

# Evidence For Function Prediction





### Sequence Viewer

Neighborhood six frame translation with putative ORF's shown below  
for gene\_oid=2500607069  
2111568..2113793(-)  
*Uncharacterized anaerobic dehydrogenase, COG3382*

Select gene neighborhood:

bp upstream.  bp downstream

Select minimum ORF size:

aa

Output Format:

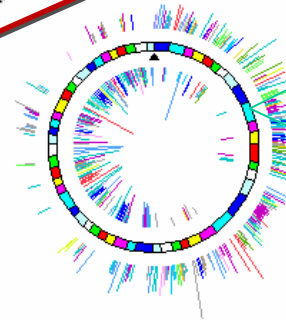
Text  Graphics

Submit

Reset

Click!

Select the region of interest:  
100 bp upstream and downstream  
of the original ORF coordinates.



+100bp

ORF

-100bp

Let's look at the  
nucleotides that  
are around  
our gene of  
interest

Use 80 amino acids as standard ORF size.  
Remember 3 nucleotides code for 1 amino acid!

The "Graphics" option will allow  
us to look at the actual  
sequence. This helps us  
complete the first task.

# After you click submit, you will see something like this:

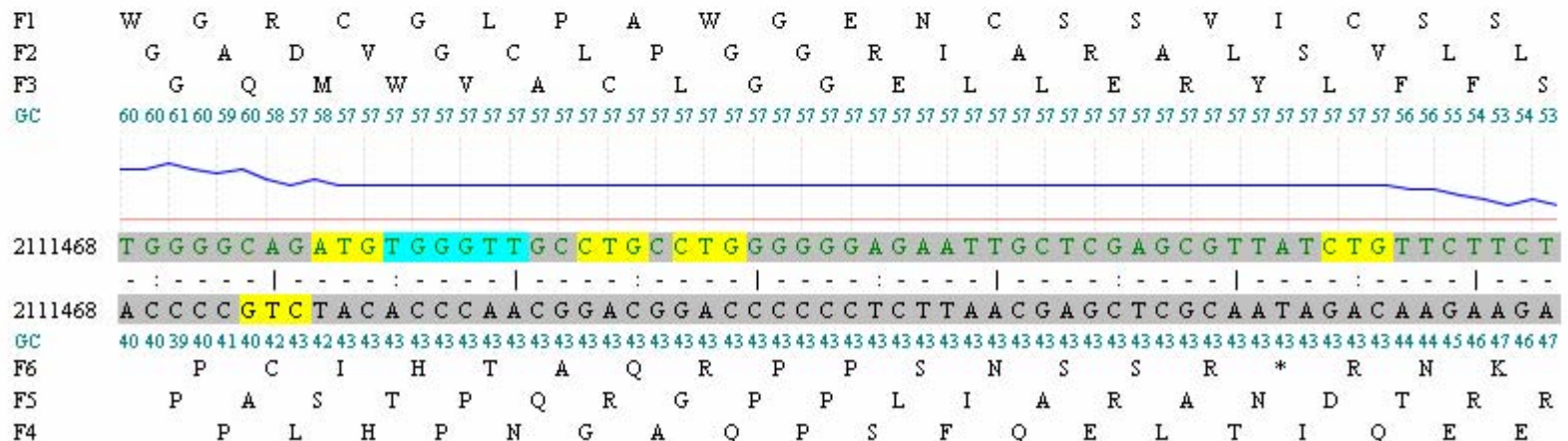


## Sequence Viewer

Loaded.

Neighborhood six frame translation with putative ORF's shown below  
for gene\_oid=2500607069  
2111568..2113793(-)  
*Uncharacterized anaerobic dehydrogenase, COG3383.*

**hint:**   indicates potential start codon region.  
  indicates possible Shine-Dalgarno region.





# WHICH STRAND IS MY GENE ON?

Example of **plus** strand gene orientation:

## DNA Coordinates

Find this on the Gene Details page for the proposed gene.

enter coordinates

5271121..5272269 (+) (1149bp)

## Evidence For Function Prediction

### Neighborhood



red = Current Gene

Example of **minus** strand gene orientation:

## DNA Coordinates

Find this on the Gene Details page for the proposed gene.

enter coordinates

2111568..2113793 (-) (2226bp)

## Evidence For Function Prediction

### Neighborhood



red = Current Gene

# HOW AM I SUPPOSED TO READ THIS??

## Deciphering Sequence Viewer in Graphics Mode



Review single letter codes for amino acids!

**img/edu** INTEGRATED MICROBIAL EDUCATION SITE

[IMG Home](#)
[Find Genomes](#)
[Find Genes](#)
[Find Functions](#)
[Compare Genomes](#)
[Analysis Carts](#)
[MyIMG](#)
[Using IMG](#)

### Sequence Viewer

Neighborhood six frame translation with putative ORF's shown for gene\_oid=2500607069  
 2111568..2113793(-)  
*Uncharacterized anaerobic dehydrogenase, COG3383.*

**hint:**    indicates potential start codon region.  
   indicates possible Shine-Dalgarno region.

Each row above (or below) designated DNA strand represents a different translational reading frame

➤ amino acids corresponding to codons in nucleotide sequence shifted by one base

AA results from forward DNA strand (F1, F2, & F3)

F1	W	G	R	C	G	L	P	A	W	G	E	N	C	S	S	V	I	C	S	S		
F2		G	A	D	V	G	C	L	P	G	G	R	I	A	R	A	L	S	V	L	L	
F3			G	Q	M	W	V	A	C	L	G	G	E	L	L	E	R	Y	L	F	F	S

AA results from reverse DNA strand (F4, F5, & F6)

2111468	T	G	G	G	C	A	G	A	T	G	T	G	G	G	T	T	G	C	T	G	C	T	G	G	G	G	A	G	A	A	T	T	G	C	T	C	G	A	G	C	G	T	T	A	T	C	T	G	T	T	C	T	T	C	T		
2111468	A	C	C	C	C	G	T	C	A	C	C	C	A	A	C	G	G	A	C	G	G	A	C	C	C	C	C	T	C	T	T	A	A	C	G	A	G	A	C	T	C	G	C	A	A	T	A	G	A	C	A	A	G	A	A	G	A
F6			P	C	I	H	T	A	Q	R	P	P	S	N	S	S	R	*	R	N	K																																				
F5			P	A	S	T	P	Q	R	G	P	P	L	I	A	R	A	N	D	T	R	R																																			
F4			P	L	H	P	N	G	A	Q	P	S	F	Q	E	L	T	I	Q	E	E																																				

# Helpful tools to keep handy while deciphering!

Insert table of genetic code and table of single letter abbreviations for amino acids

**Initiation frequency in *E. coli*: AUG > GUG > UUG > CUG**

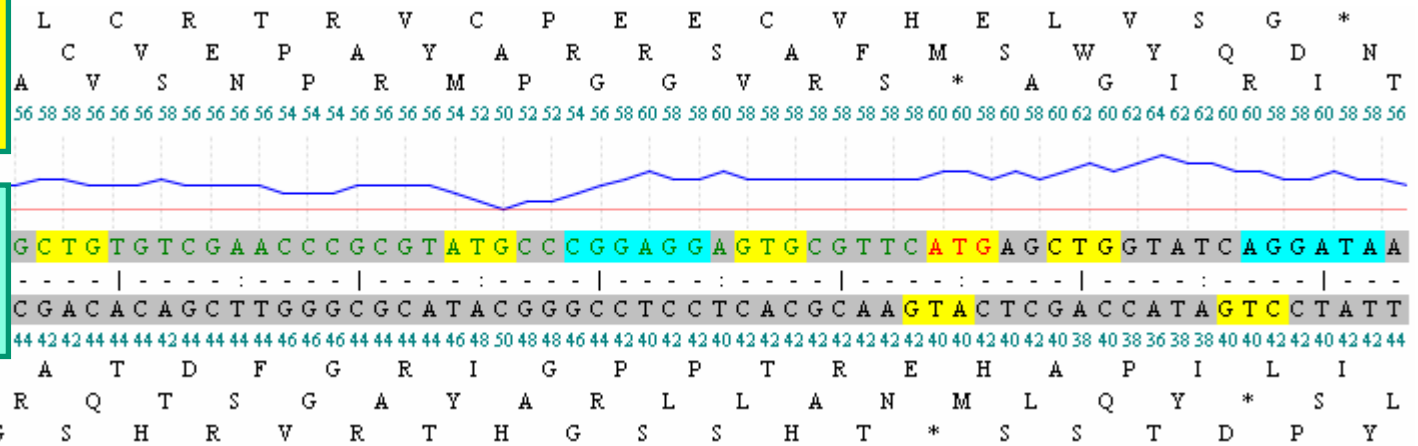
**Stop codons: UAG, UAA, UGA**



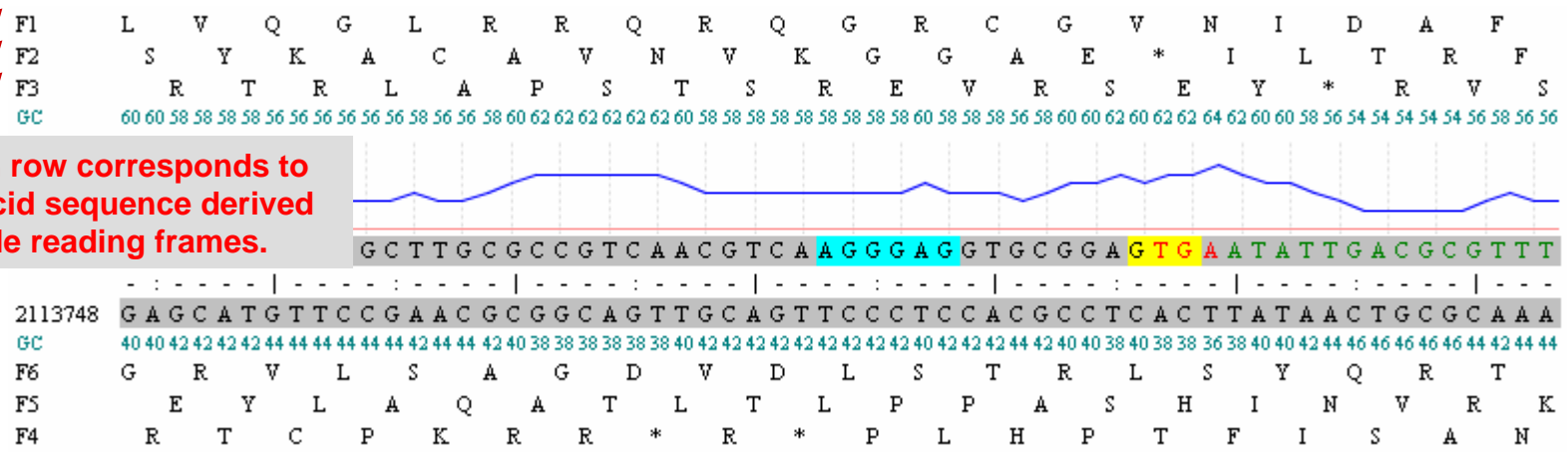
# More on... Deciphering Sequence Viewer in Graphics Mode

Yellow denotes possible alternative start codons.

Cyan denotes candidate RBS (Shine-Dalgarno Sequences).



Recall: Each row corresponds to the amino acid sequence derived from possible reading frames.



RBS = Ribosome Binding Site = Shine-Dalgarno Sequence

➤ Always on SAME DNA STRAND AS THE START CODON



*Remember:* To initiate protein synthesis (translation), the ribosome interacts with the Shine-Dalgarno sequence in the mRNA immediately upstream of the proper start codon

Insert WebLogo of *E. coli* ribosome binding sites and start codons from <http://molbiol-tools.ca/Motifs.htm>

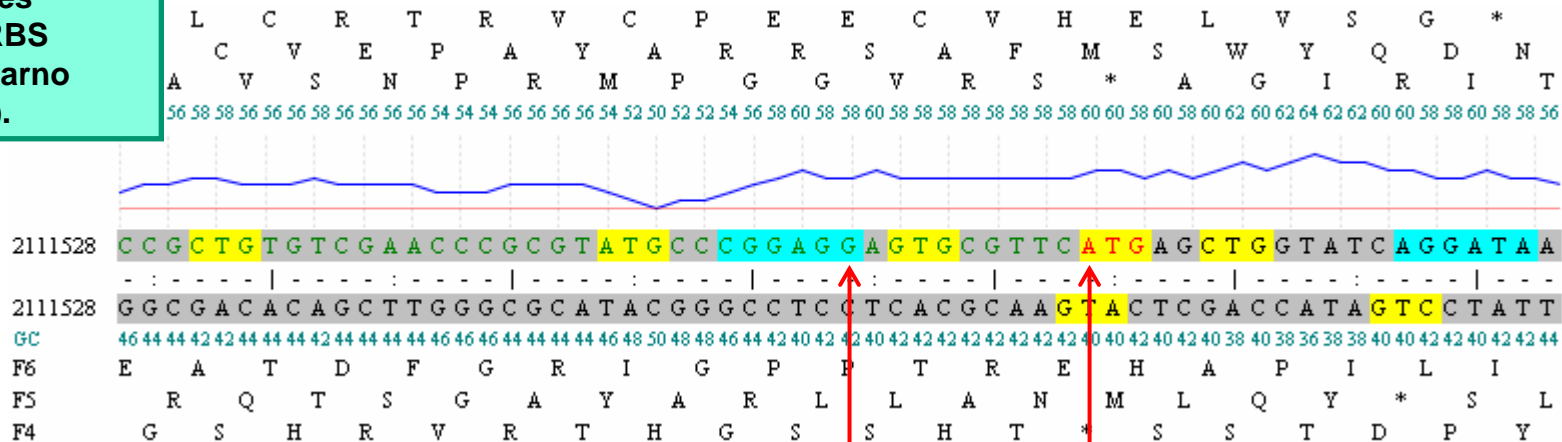
Reference: **Figure 1** in Schneider and Stephens (1990)  
Sequence logos: a new way to display consensus sequences.  
*Nucleic Acids. Research* **18**: 6097-6100.

# Decision Tree

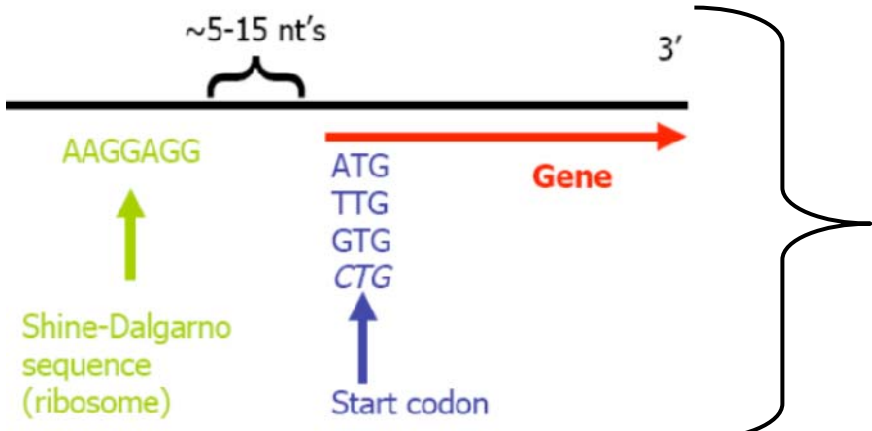
- The following slides contain a Decision Tree to aid in decision making as you proceed through the first task of this module.
- All of the steps provided in the tree will be elaborated upon throughout the presentation.
- The tree should act as an outline to help you map out your plan of **attack!**

Find your original start codon.  
 Highlighted Yellow and Red Text.  
 Does it have a Shine Dalgarno Sequence 5-15bp upstream?

Cyan denotes candidate RBS (Shine-Dalgarno Sequences).



10 bases apart



Recall: Start codons should occur about 5-15 bases downstream of the Shine-Dalgarno sequence (RBS)

Find your original start codon.

Highlighted Yellow and Red Text.

Does it have a Shine Dalgarno Sequence 5-15bp upstream?

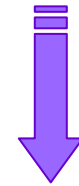
Yes

No

The DNA coordinates should have been recorded in your lab notebook as part of first module. Provide a comment in current module indicating the original coordinates are likely correct.

Task #2

# Recording results in your Lab Notebook



Scroll down

## Alternative Open Reading Frame Module

- [Module Instructions](#)

go to the IMG Gene Details page for the proposed gene

Proposed DNA coordinates

`enter in lab report`

Explanation of choice

explanation

# Recording results in your Lab Notebook

## Alternative Open Reading Frame Module

- [Module Instructions](#)

go to the IMG Gene Details page for the proposed gene

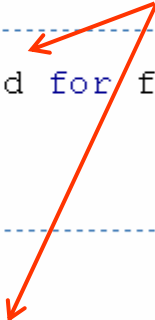
Proposed DNA coordinates

```
No change from original DNA coordinates recorded for first module:  
3948986..3951211 (+)
```

Explanation of choice

The original start codon for the ORF identified by automatic Gene Caller has Shine-Dalgarno sequence (CGGAGG) 10 bases upstream of the start codon (AUG).

1- Enter original DNA coordinates with an explanatory comment



Find your original start codon.  
**Highlighted Yellow and Red Text.**  
Does it have a Shine Dalgarno Sequence 5-15bp upstream?

Yes

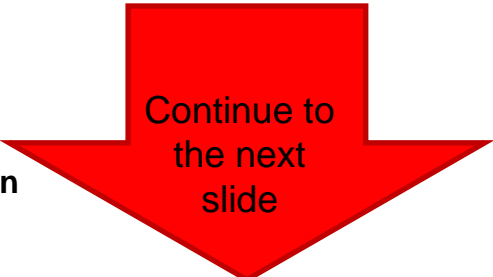
The DNA coordinates should have been recorded in your lab notebook as part of first module. Provide a comment in current module indicating the original coordinates are likely correct.

Task #2

No

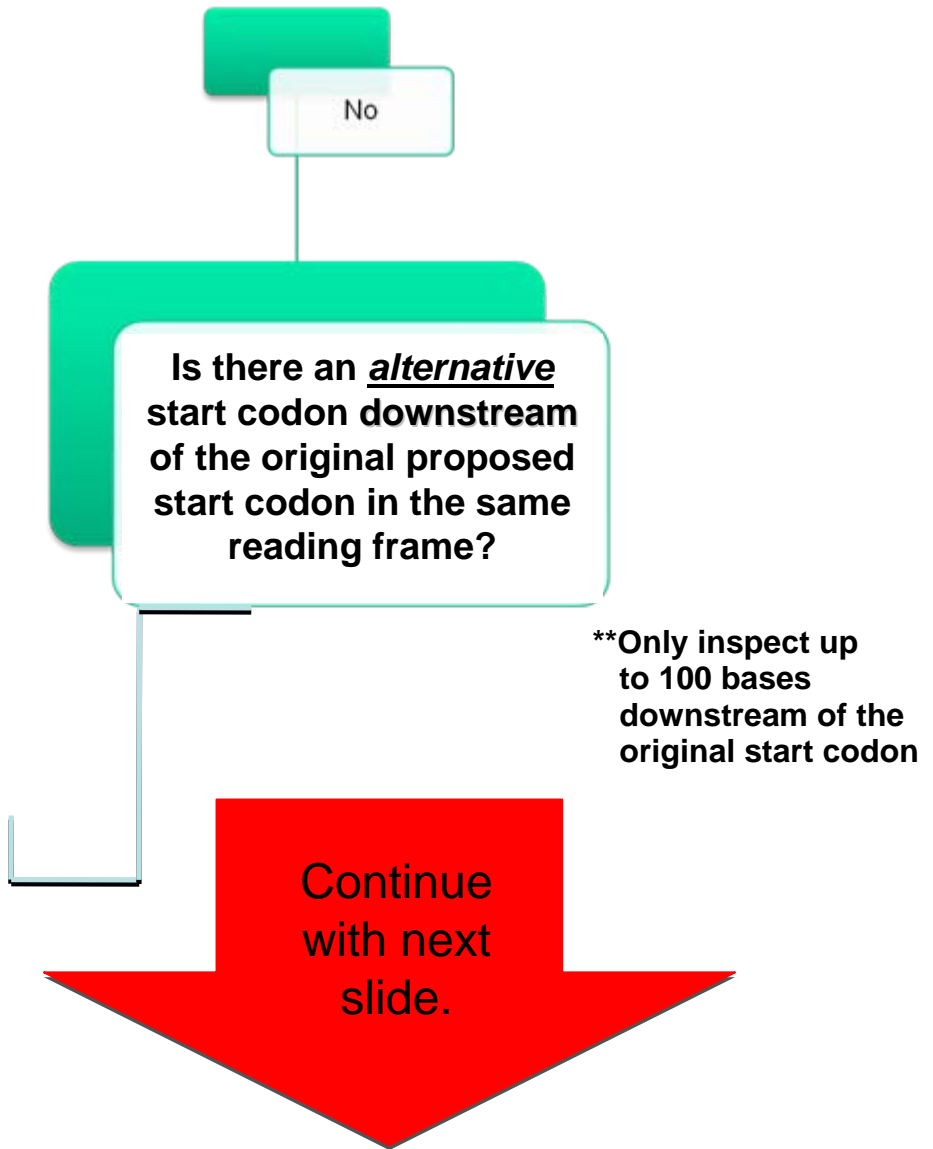
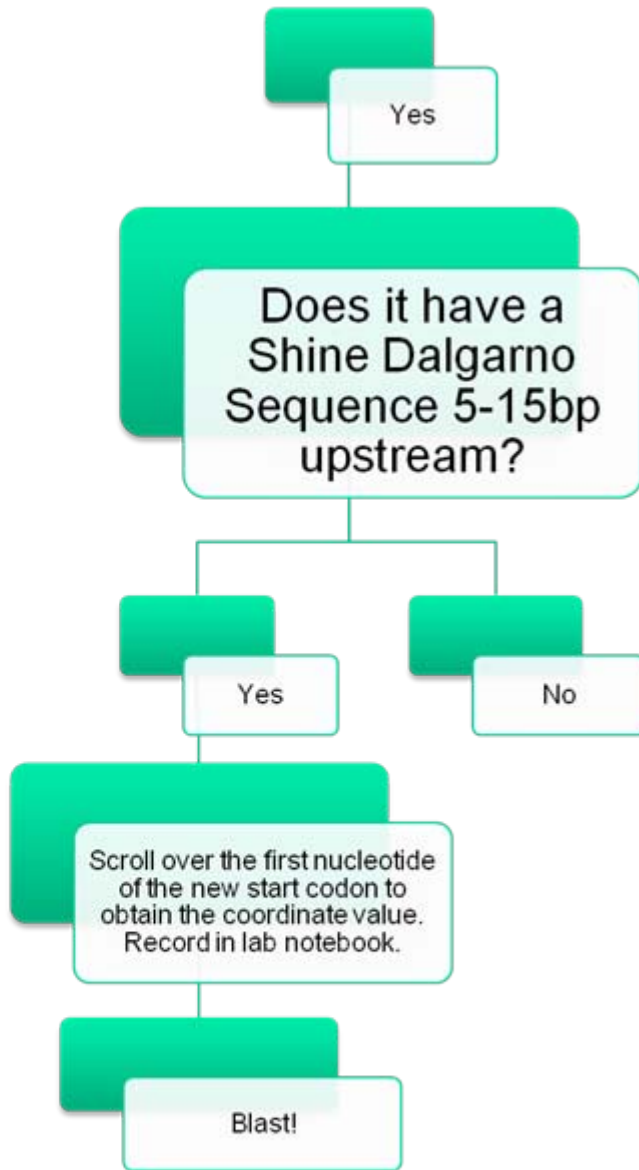
In the same reading frame as the original start codon, is there an ***alternative*** start codon upstream of the original proposed start codon?

**\*\*Only inspect up to 100 bases upstream of the original start codon**

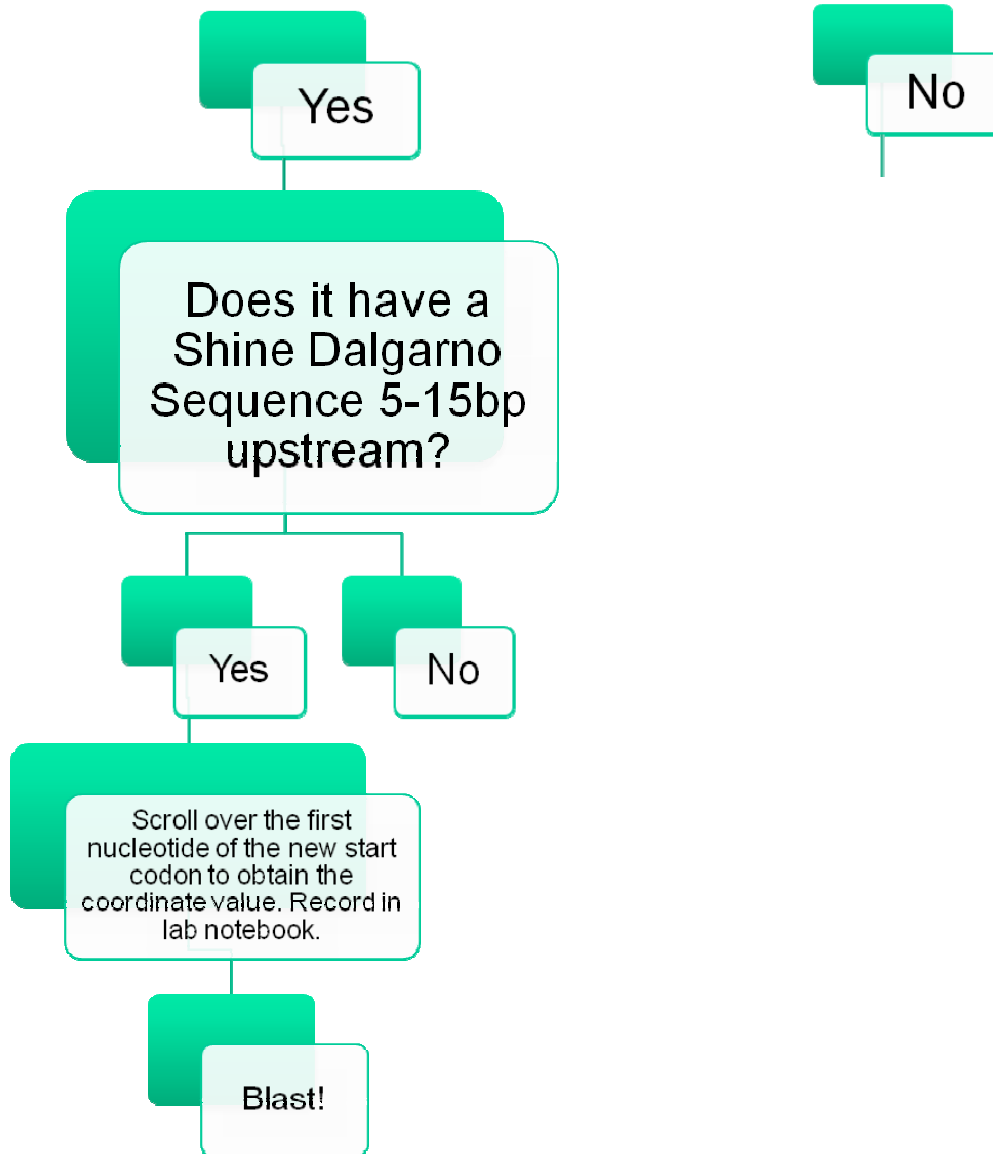












# For genes with possible alternative start codon...It's time to BLAST!

- **BLAST your results:**
  - Construct a “revised” protein sequence in FASTA format (add or subtract amino acid residues in proper reading frame to reflect new start codon position then copy/paste into lab notebook).
  - Submit as query for a BLAST search of NCBI database (Genbank, SwisProt).
- **Your results from BLAST:**
  - Compare results from original blast search with those from new blast search.
    - Determine if statistics have improved.
      - REMEMBER: higher bit score, lower e-value, higher % identity and/or longer alignment length are all good arguments that the alternative start codon is a better choice

# Create new headings & boxes for entering amino acid sequence in FASTA format for ORF with alternate start codon

## Alternative Open Reading Frame Module

- [Module Instructions](#)

1- Add headings and box in lab notebook

go to the IMG Gene Details page for the proposed gene

## Protein Sequence for ORF with Alternate Start Codon

Construct this from Sequence Viewer on the Gene Details page for the proposed gene.

enter AA [FASTA format](#)

2- Copy/paste modified protein sequence into box

```
>OID 2500608365 Flp pilus assembly protein TadC with alternative start codon
MSNLLLEKAAPALSKALEPKSELEQSQLKIRLANAGFHSPQAPMIYLAIKTVCLVVGLVLGGGLGMYRYGTTQAGLT
TLIIAAGAGFYLPYGLAYLISKRKQAIFLQLPDVLDLLVVCVEAGLGLDAGLRRVAEELKDTAPEICGELAMCNL
QLQMGRNRRDMLHDLGVRTGVDDVKALVAIMIQADKFGSSIAQALRVQSDSMRVKRRQIAEKAQKTAVQMLFFPMV
IFIFPGIFVVVLVGPAAIKMMDQLLNKP
```

**\*\*NOTE** highlighted portion was added to the original amino acid sequence since the alternative start codon (AUG) is located 26 residues upstream of the original start codon (UUG, which encodes Leucine)

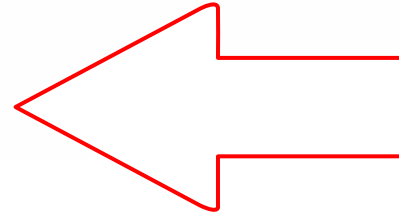
NCBI BLAST can be accessed from  
lab notebook link

## Sequence-based Similarity Data Module

- [Module Instructions](#)

### BLAST

go to <http://www.ncbi.nlm.nih.gov/blast>



## BLAST Results for ORF with original start codon (Gene Caller)

```
> [ref|NP_866981.1] [G] secretion system protein TadC [Rhodopirellula baltica SH 1]
  emb|CAD74523.1| [G] probable secretion system protein TadC [Rhodopirellula baltica
  SH 1]
```

Length=328

GENE ID: 1791454 tadC | secretion system protein TadC  
[Rhodopirellula baltica SH 1] (10 or fewer PubMed links)

Score = 263 bits (671), Expect = 1e-68, Method: Compositional matrix adjust.  
Identities = 133/227 (58%), Positives = 168/227 (74%), Gaps = 0/227 (0%)

```
Query 1 LKIRLANAGFHSPQAPMIYLAIKTVCLVVGLVGGGLGMYRYGTTQAGLTTLIIAAGAGF
      L+ +L NAGF AP+I+ I+ VC VGL+LGG G G TQ + L++ GF
Sbjct 102 LREKLINAGFRRESAPVIFKLIQLVCTGVGLMLGGVTGAVLDGLTQGMIIKLLLGMIGGF
```

✓ Compare hits from same organism

✓ Are statistics improved for proposed ORF?  
➤ If yes, enter data into notebook

## BLAST Results for ORF with alternative start codon proposed by student!

```
> [ref|NP_866981.1] [G] secretion system protein TadC [Rhodopirellula baltica SH 1]
  emb|CAD74523.1| [G] probable secretion system protein TadC [Rhodopirellula baltica
  SH 1]
```

Length=328

GENE ID: 1791454 tadC | secretion system protein TadC  
[Rhodopirellula baltica SH 1] (10 or fewer PubMed links)

Score = 271 bits (694), Expect = 3e-71, Method: Compositional matrix adjust.  
Identities = 140/253 (55%), Positives = 182/253 (71%), Gaps = 1/253 (0%)

```
Query 1 MSNLEKAAPALSKALEPKSELEQSQLKIRLANAGFHSPQAPMIYLAIKTVCLVVGLVLG 60
      ++ LEKAA ++ ++ E E +L+ +L NAGF AP+I+ I+ VC VGL+LG
Sbjct 77 LTEALEKAATPIASSVSGNEE-EMGKLREKLINAGFRRESAPVIFKLIQLVCTGVGLMLG 135
```

# Recording results in your Lab Notebook

## Alternate Start Codon for Original Open Reading Frame

Find this on the Gene Details page for the proposed gene.  
 enter coordinates  
 enter proposed DNA coordinates

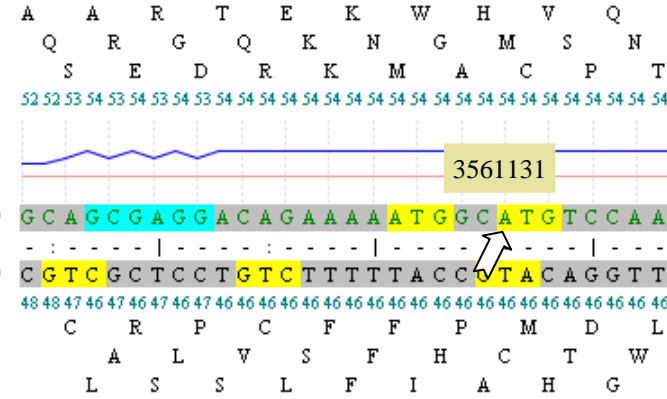
Adjust DNA coordinates  
& base pairs.

3561131..3561898 (+) (765bp)

### Reasoning

No Shine-Dalgarno sequence was found by automatic Gene Caller immediately upstream of original start codon. However, located upstream of the original start codon there is a Shine-Dalgarno sequence 14 bases upstream of an alternative start codon (ATG). Furthermore, the NCBI protein BLAST statistics for the ORF using the alternative start codon improved in comparison to the search with the ORF having the original start codon. In particular, the score was higher (271 bits compared to 263 bits) and the E-value was lower ( $3e-71$  compared to  $1e-68$ ) for the ORF with an alternative start codon. However, the percent identity was lower for the alternative compared to the original (55% compared to 58%). Also of note is that the length of the alignment did not change (328 amino acids for both). Taken together, I believe there is sufficient evidence to suggest that the alternative position for the start codon for this gene is a better choice than the original position. My conclusion is based on three results: (1) presence of Shine-Dalgarno sequence at appropriate distance from the start codon, (2) higher bit score and lower E-value for BLAST search using alternative ORF.

F1  
F2  
F3  
GC  
3561109  
3561109  
GC  
F6  
F5  
F4



1- Change section heading as shown

2- Copy/paste sub-headings & field boxes from Sequence-based Similarity Data module.

3- Fill in with your results from BLAST search using ORF with alternate start codon for gene.

### BLAST Results

Gene product name (same hit as in Module 2)

secretion system protein TadC

Organism

Rhodospirellula baltica SH 1

Length, E-Value, Score, Percent identity, Positives, and Gaps

Length=328, Score = 271 bits (694), Expect =  $3e-71$ , Method: Compositional matrix adj  
 Identities = 140/253 (55%), Positives = 182/253 (71%), Gaps = 1/253 (0%)

Alignment of the BLAST hit and the query sequence

Query	1	MSNLEKAAPALSKALEPKSELEQSQLKIRLANAGFHSPQAPMIYLAIKTVCLVVGLVLG	60
		++ LEKAA ++ ++ E E +L+ +L NAGF AP+I+ I+ VC VGL+LG	
Sbjct	77	LTEALEKAATPIASSVSGNEE-EMGKLRKLINAGFRRESAPVIFKLIQLVCTGVGLMLG	135



# Are the BLAST results always better? NO!

For example. . .

## Alternative Open Reading Frame Module

- [Module Instructions](#)

go to the IMG Gene Details page for the proposed gene

## Protein Sequence for ORF with Alternate Start Codon

Construct this from Sequence Viewer on the Gene Details page for the proposed gene.

enter AA [FASTA format](#)

```
>OID 2500607070 Fe-S-cluster-containing hydrogenase components_alternative start codon
LQDLLEGFGSVSAPPSKGEQMVHRLLOKQRELSAVELFAKEHAGGTINGS
KLPDQARFYASLLPATPVGPGQYGFVDDLDRCSGCKACVTACHSLNGLD
DSETWRDVGLLIGGTETLPVMQHVTAAACHHCLEPACMTACPVNAYEKDAF
TGIVRHLDDQCFGCQYCTLACPYNVPKYHAAKGIVRKCDMCSNRLKNGEA
PACVQACPHEAISIRIVDVS RV TENAEADHFLPAAPEPYITLPTTTYRTT
RVFPRNMLPADYYSVSPQHPHWPLIVMLVLTQLSVGAFAAGSFLEEALDT
ELATAFRPIHATGALVLLGALLGASTLHLGRPLYAFRGILGFRHSWLSRE
IVAFGLFAGLAVPFFAGLCWGLPLLEVSGSPWGKLAGELLPTLSPSVACVG
VIGVFCVMIYVFTRRELWLSLERTLIRFSLTTILLGVATIWLMMWLAVGF
LSDDEWHQLAQNLTRPLARSVIILTTLKLLYDISLLRHLATFRNSPLKRS
ALLVVGPLRGFSIGRLVLGVVGGVVI PAAFAAVPIHDP IQFTTTTSAVFVG
QMWVACLGGELLERYLFFSAVSNPRMPGGVRS
```

\*\*Recall that this gene had residues deleted from the original amino acid sequence since the alternative start codon (CUG, which encodes Leucine) is located 12 residues downstream of the original start codon (AUG)

## BLAST Results for ORF with original start codon (Gene Caller)

```
>[ref|ZP_01093385.1] molybdopterin oxidoreductase, iron sulfur subunit [Blastopirellula marina DSM 3645]
gb|EA077915.1| molybdopterin oxidoreductase, iron sulfur subunit [Blastopirellula marina DSM 3645]
Length=536
```

```
Score = 374 bits (960), Expect = 2e-101, Method: Compositional matrix
Identities = 235/574 (40%), Positives = 304/574 (52%), Gaps = 61/574 (10%)
```

```
Query 28 SKGEGMVHRLLLQKQRELSAVELFAKEHAGGTINGSKLPDQARFYASLLPATPVGPGQQYG
+ G ++ LL +Q++L+AVE F++ HA T P A Y LLPA G+QY
Sbjct 13 TDGFDLIQSLLAEQDQDLTAVERFSQRHADATT-----PLMAPLYRELLPAAAPSAGEQYA
```

- ✓ Compare hits from same organism
- ✓ Are statistics improved for proposed ORF?
  - If no or not really, enter data into notebook anyway with a comment

## BLAST Results for ORF with alternative start codon proposed by student!

```
>[ref|ZP_01093385.1] molybdopterin oxidoreductase, iron sulfur subunit [Blastopirellula marina DSM 3645]
gb|EA077915.1| molybdopterin oxidoreductase, iron sulfur subunit [Blastopirellula marina DSM 3645]
Length=536
```

```
Score = 374 bits (960), Expect = 1e-101, Method: Compositional matrix adjust.
Identities = 235/574 (40%), Positives = 304/574 (52%), Gaps = 61/574 (10%)
```

```
Query 16 SKGEGMVHRLLLQKQRELSAVELFAKEHAGGTINGSKLPDQARFYASLLPATPVGPGQQYG 75
+ G ++ LL +Q++L+AVE F++ HA T P A Y LLPA G+QY
Sbjct 13 TDGFDLIQSLLAEQDQDLTAVERFSQRHADATT-----PLMAPLYRELLPAAAPSAGEQYA 67
```

# Recording results in your Lab Notebook

## Alternate Start Codon for Original Open Reading Frame

Find this on the Gene Details page for the proposed gene.  
 enter coordinates  
 enter proposed DNA coordinates

**Adjust DNA coordinates  
& base pairs.**

2113926...2115574 (-) (1746bp)

### Reasoning

No Shine-Dalgarno sequence was found by automatic Gene Caller immediately upstream of original start codon. However, located downstream of the original start codon there is a Shine-Dalgarno sequence 10 bases upstream of an alternative start codon (CTG). Although the NCBI protein BLAST statistics do not markedly improve for the ORF using the alternative start codon, they also are not significantly different (or worse). Thus, I cannot rule out the possibility that the proposed alternative start codon is a better choice than the original start codon.

### BLAST Results

Gene product name (same hit as in Module 2)

molybdopterin oxidoreductase, iron sulfur subunit

Organism

Blastopirellula marina DSM 3645

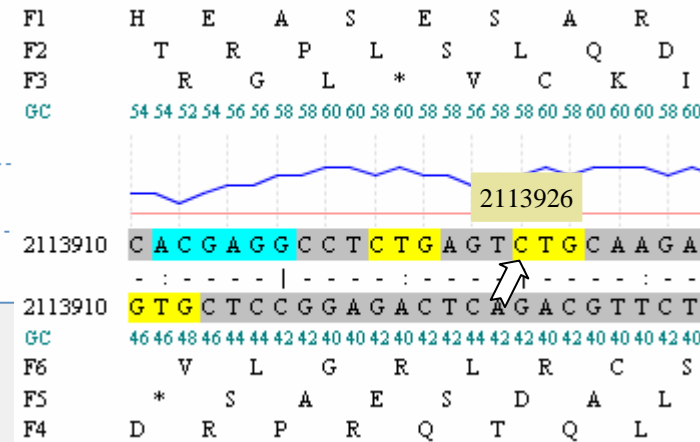
Length, E-Value, Score, Percent identity, Positives, and Gaps

Length=536, Score = 374 bits (960), Expect = 2e-101, Method: Compositional matrix adjust.  
 Identities = 235/574 (40%), Positives = 304/574 (52%), Gaps = 61/574 (10%)

Alignment of the BLAST hit and the query sequence

```

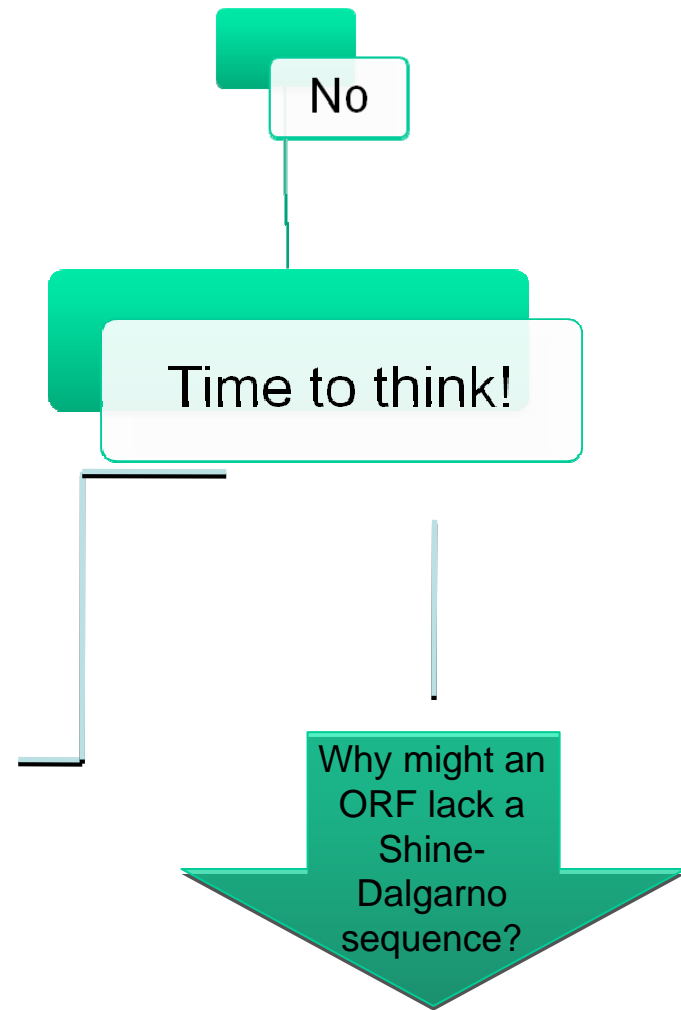
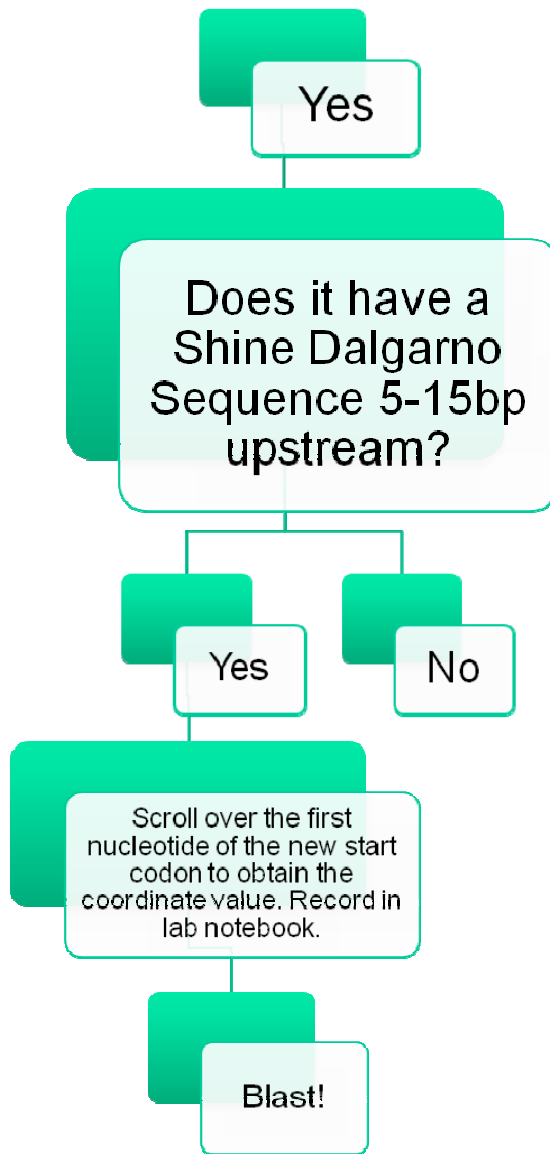
Query 16  SKGQMVHRLQLKQRELSAVELFAKEHAGGTINGSKLPDQARFYASLLPATPVGPGQQYG 75
          + G ++ LL +Q++L+AVE F++ HA T P A Y LLPA G+QY
Sbjct 13  TDGFDLIQSLLAEQQDLTAVERFSQRHADATT-----PLMAPLYRELLPAAAPSAGEQYA 67
  
```



**1- Change section heading as shown**

**2- Copy/paste sub-headings & field boxes from Sequence-based Similarity Data module.**

**3- Fill in with your results from BLAST search using ORF with alternate start codon for gene.**



# WHY? WHY? WHY NOT?!

## Interpreting Your Negative Results

(i.e., no Shine-Dalgarno, no alternative start codon, etc.)

- Maybe there is flexibility in the amount of sequence conservation needed in the Shine-Dalgarno (S-D) that allows ribosome binding

Insert **Figure 33.9** from Garret & Grisham *Biochemistry* (2<sup>nd</sup> Ed.)  
Alignment of various Shine-Dalgarno sequences recognized by  
*E. coli* ribosomes.

# WHY? WHY? WHY NOT?!

## Interpreting Your Negative Results

(i.e., no Shine-Dalgarno, no alternative start codon, etc.)

- Consider possible mutations in DNA sequence
  - Remember “draft genome” problems?
- Consider a ribosome “skid” if your gene is part of an operon
  - Maybe your gene is in the middle or at the end of an operon. And perhaps only the first gene in the operon has a S-D upstream of the start codon and has a stop codon within 5-15 nt of the start codon for the next gene in the operon. In this case, the genes may be close enough that the ribosome does not need to completely dissociate from the RNA transcript, instead it “skids” along and begins translation of the next gene in the operon without needing to bind a second S-D.
    - **Look at the ortholog neighborhood map!**  
(Review Gene Context in HGT module)



# Recording results in your Lab Notebook

## Alternate Start Codon for Original Open Reading Frame

Find this on the Gene Details page for the proposed gene.

enter coordinates

enter proposed DNA coordinates

No change from original DNA coordinates recorded for first module:

2111568..2113793 (-) (2226bp)

Reasoning

Although the ORF identified by automatic Gene Caller does not have an obvious Shine-Dalgarno sequence upstream of the start codon, an alternative position for the start codon also was not apparent. It is possible that the S-D sequence for this gene is not well conserved or that there are point mutations in the DNA sequence that made the S-D unrecognizable (e.g., draft genome problem). Notably, upon inspection of the ortholog neighborhood, it is possible that this gene is part of an operon. As shown below, the gene (light green) preceding my gene (red) may be close enough to be transcribed together as an operon. Consistent with this hypothesis, the green gene (which encodes a Fe-S-cluster-containing hydrogenase component as ascertained by the automated Gene Caller) does indeed have a start codon preceded by a Shine-Dalgarno sequence (provided the alternative position is correct). The two genes also may have a functional relationship (e.g., hydrogenase vs. dehydrogenase) as would be expected if part of an operon. If my gene is indeed part of an operon, then it's possible that the initiation of translation occurs from a ribosome "skid" after completion of translation of the preceding gene.

*Planctomyces limnophilus* DSM 3776 : PlimDRAFT\_4083246\_C168



1- Change the heading

2- Enter original DNA coordinates with an explanation

NOW for the second task:

Hunting for Potential  
**Novel** Open Reading Frames (ORFs)



# Return to the Gene Detail Page



INTEGRATED MICROBIAL GENOMES  
EDUCATION SITE

IMG Home | Find Genomes | Find Genes | Find Functions | Compare Genomes | Analysis Carts | MyIMG | Using IMG

Gene Search | BLAST | Phylogenetic Profilers

## Gene Detail

Loaded.

- [Gene Information](#)
- [Evidence For Function Predictions](#)
- [Sequence Search](#)
- [External Sequence Search](#)
- [IMG Sequence Search](#)
- [Homolog Display](#)

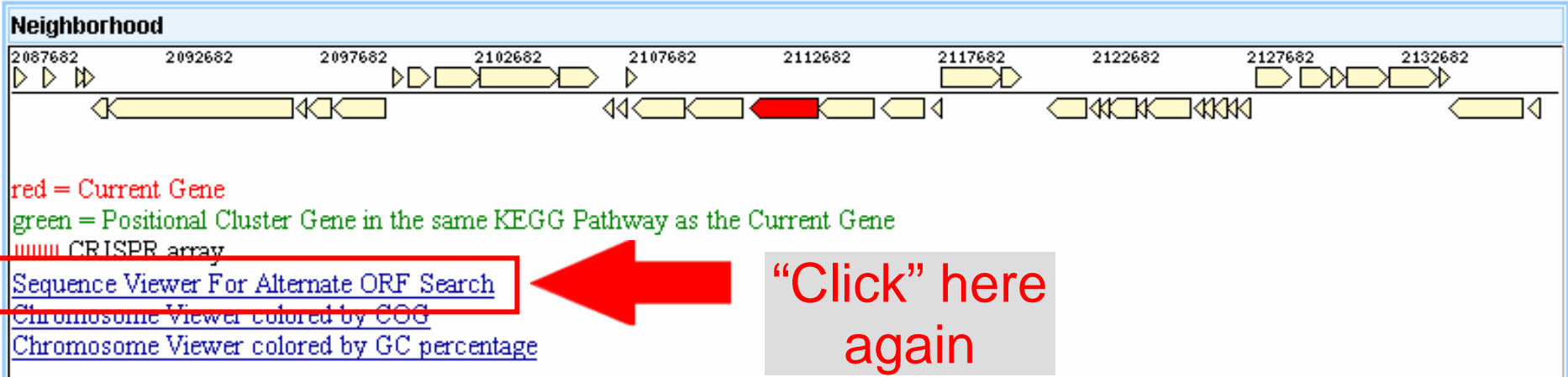
## Gene Information

Gene Information	
Gene Object ID	2500607069
Gene Symbol	
Locus Tag	PlimDRAFT_19450
Product Name	Uncharacterized anaerobic dehydrogenase, COG3383
IMG Product Source	COG3383
Genome	<a href="#">Planctomyces limnophilus DSM 3776</a>
DNA Coordinates	2111568..2113793 (-)(2226bp)
Scaffold Source	<a href="#">Planctomyces limnophilus DSM 3776 : PlimDRAFT_4083246_C168 (5423025bp)</a>
IMG ORF Type	
GC Content	0.58
External Links	
Fused Gene	No
Fusion Component	No
Protein Information	
Amino Acid Sequence Length	741aa



Scroll down

# Evidence For Function Prediction



# Change output from “Graphic” to “Text”



IMG Home

Find Genomes

Find Genes

Find Functions

Compare Genomes

Analysis Carts

MyIMG

Using IMG

## Sequence Viewer

Neighborhood six frame translation with putative ORF's shown below

for gene\_oid=2500607069

2111568..2113793(-)

*Uncharacterized anaerobic dehydrogenase, COG3383.*

Select gene neighborhood:

bp upstream.  bp downstream

Select minimum ORF size:

aa

Output Format:

Text  Graphics

Submit

Reset

Keep parameters the same as first task.

USE **TEXT** option to look for NOVEL open reading frames

Click submit!

IMG Home

Find Genomes

Find Genes

Find Functions

Compare Genomes

Analysis Carts

MyIMG

Using IMG

## Sequence Viewer

Loaded.

Neighborhood six frame translation with putative ORF's shown below

for gene\_oid=2500607069

2111568..2113793(-)

*Uncharacterized anaerobic dehydrogenase, COG3383.*

**hint:** Copy and paste sequence to BLAST and InterPro scan to test ORF translation.

>2500607069\_1\_ORF2 Translation of 2500607069 in frame 1, ORF 2, threshold 80, 102aa  
ERKAGTCRLGCSAETFHGPLSTDSGKARAPLSGVHQHRPDAHRRDGIARRGGQVRHGPOA  
RRRQYAAVHGDGSGRLQGVRFRCASLHLRRLRRIGRHGVRG

>2500607069\_3\_ORF3 Translation of 2500607069 in frame 3, ORF 3, threshold 80, 140aa  
AWMRVSSPPRTVGRTTKSWREFSKRKS KDCGSSVPTRPTRGSTRTWPVKCFRGSISWLCR  
TCTTIRKRPRLTSSFQPLAGARKRGLLIPSGALVVSRRSGVRRGKRWPSTSTSFNWLLS  
IGAAESNSGAGLRPKPSLKF

>2500607069\_4\_ORF3 Translation of 2500607069 in frame 4, ORF 3, threshold 80, 123aa  
KMWKSASAFPGARLTFLIRPMRSELINVPSFSPQPAAGRTRSAIWAVSVLWYMSCTTRK  
SSRESISRARFWLIHEWAVLVQMIHSPLIFSLRIPSMISWYDQLFSVGMTLSSMLRMPAI  
FAR

>2500607069\_5\_ORF2 Translation of 2500607069 in frame 5, ORF 2, threshold 80, 426aa  
LVRPLVMHRNEDLSLSHSRSHERHRPHAASSRLDGDGPGVMGDLECSICRIDFDKRS LGI  
QLPEDRRRLRSAGLRMPLRSASPSSQONERIRLVGQFRQRAGGLKEEFCTTVGMIKTTILE  
ETPLLYGGGHPFRAGPLHPPLFFDPAVVLNAREVAGHTARTLLQNFKDFGRRPAPELLS  
AAPILSNQLKDVEVGQRFPRRTPDLLDTTNAPLGINKRPLFLAPASGWKDEVSHLGRFRI  
VVHVLHNQEIEPRKHFTGQVLVDP RVGRVGTDDPQSFDFLFENSLHDFVVRPTVLGGDDT  
LIHAQNARDLRAVIGVLEIVSAQKVGRVAEEPRTHGVALARDRIRSGSGFANVPRHQCCI  
DDRLRGTHPLVALIDPHRPPERDSLFLVDRAGESLNLVRLKTRFGCSAGQAIIVVHEAGEF  
LESRRR

NOTE: Usually the result with the longest amino acid sequence length is your original ORF.

HINT: A NOVEL ORF could be in a different reading frame from the original ORF

We are interested in all ORFs with sequence length of at least 80 aa

[IMG Home](#)[Find Genomes](#)[Find Genes](#)[Find Functions](#)[Compare Genomes](#)[Analysis Carts](#)[MyIMG](#)[Using IMG](#)

## Sequence Viewer

Loaded.

Neighborhood six frame translation with putative ORF's shown below

for gene\_oid=2500607069

2111568..2113793(-)

*Uncharacterized anaerobic dehydrogenase, COG3383.*

**hint:** Copy and paste sequence to BLAST and InterPro scan to test ORF translation.

>2500607069\_1\_ORF2 Translation of 2500607069 in frame 1, ORF 2, threshold 80, 102aa  
ERKAGTCRLGCSAETFHGPLSTDSGKARAPLSGVHQHRPDAHRRDGIARRGGQVRHGPOA  
RRRQYAAVHGDGSGRLQRGVRFRCASLHLRRLRRIGRHGVRG

>2500607069\_3\_ORF3 Translation of 2500607069 in frame 3, ORF 3, threshold 80, 140aa  
AWMRVSSPRTVGRITTKSWREFSKRKS KDCGSSVPTRPTRGSTRTWPVKCFRGSISWLCR  
TCTTIRKRPRWLTSSFQPLAGARKRGLLIPSGALVVSRRSGVRRGKRWPTSTSFNWLLS  
IGAAESNSGAGLRPKPSLKF

>2500607069\_4\_ORF3 Translation of 2500607069 in frame 4, ORF 3, threshold 80, 123aa  
KMWKSASAFPGARLTFILIRPMRSELINVPSFSPQPAAGRTRSAIWAVSVLWYMSCTTRK  
SSRESISRARFWLIHEWAVLVQMIHSPLIFSLRIPSMISWYDQLFSVGMTLSSMLRMPAI  
FAR

>2500607069\_5\_ORF2 Translation of 2500607069 in frame 5, ORF 2, threshold 80, 426aa  
LVRPLVMHRNEDLSLSHSRSHERHRPHAASSRLDGDGPGVMGDLECSICRIDFDKRS LGI  
QLPEDRRLRSAGLRMPLRSASPSSQONERIRLVGQFRQRAGGLKEEFCTTVGMIKTTILE  
ETPLL YGGGHPFRAGPLHPPLFFDPAVVLNAREVAGHTARTLLQNFKDFGRRPAPELLS  
AAPILSNQLKDVEVGQRFPRRTPDLLDTTNAPLGINKRPLFLAPASGWKDEVSHLGRFRI  
VVHVLHNQEIEPRKHFTGQVLVDP RVGRVGTDDPQSFDFLFENSLHDFVVRPTVLGGDDT  
LIHAQNARDLRAVIGVLEIVSAQKVGRVAEEPRTHGVALARDRIRSGSGFANVPRHQCCI  
DDRRLRGTHPLVALIDPHRPPERDSLFLVDRAGESLNLVRLKTRFGCSAGQAI VVHEAGEF  
LESRRR

Copy/paste the  
complete FASTA  
sequence into a  
BLAST query box

# Most translations will not produce significant hits

Sequences producing significant alignments:

		Score (Bits)	E Value
<a href="#">emb CAI59731.1 </a>	homeodomain transcription factor bW2 [Sporiso...	<a href="#">39.3</a>	0.70
<a href="#">ref YP_002179979.1 </a>	HTH-type transcriptional regulator araB [...	<a href="#">37.7</a>	2.2
<a href="#">ref YP_002200969.1 </a>	LysR-family transcriptional regulator [St...	<a href="#">36.2</a>	6.8
<a href="#">ref YP_002191839.1 </a>	lysR-family transcriptional regulator [St...	<a href="#">35.8</a>	7.8
<a href="#">ref YP_001768978.1 </a>	multi-sensor hybrid histidine kinase [Met...	<a href="#">35.4</a>	9.5

E-value too high!



. . . Or will give no hits at all

**BLAST** *Basic Local Alignment Search Tool*

Home Recent Results Saved Strategies Help

► NCBI/ BLAST/ blastp suite/ Formatting Results - YM9K9XGD011

[Edit and Resubmit](#) [Save Search Strategies](#) [►Formatting options](#) [►Download](#)

**2500607069\_1\_ORF2 Translation of 2500607069...**

**Query ID** |d|69962  
**Description** 2500607069\_1\_ORF2 Translation of 2500607069 in frame 1, ORF 2, threshold 80, 102aa  
**Molecule type** amino acid  
**Query Length** 102

**No significant similarity found. For reasons why, [click here](#)**

Other reports: [►Search Summary](#)



# What do I enter in my lab notebook?

- ✓ Create headings and boxes as indicated
- ✓ Enter “No significant hits”

## Novel ORFs

Translation of genomic region obtained from TEXT view of Sequence Viewer

BLAST Results

```
No significant hits.
```

# If your BLAST search does produce a significant hit...

A potentially valid NOVEL ORF will give you a high % identity, low e-value, & high bit score.

- ✓ Create headings and boxes as indicated in your notebook
- ✓ Fill in boxes with result from BLAST search using sequence for NOVEL ORF

## Novel ORFs

Translation of genomic region obtained from TEXT view of Sequence Viewer

### BLAST Results for significant hit(s):

Gene product name (top hit)

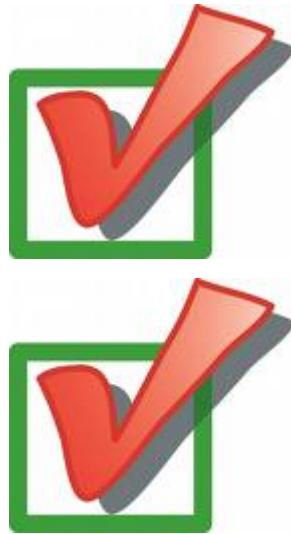
Organism

Length of alignment, Score, E-Value, Identities, Positives, Gaps

Pair-wise alignment of the database hit and the query sequence



# Module tasks complete



Are you keeping up with your annotations?

The 1<sup>st</sup> of 3 *imgACT* notebook checks will occur at the end of this week.