# Gene Modeling Tutorial using Artemis

We are working on gene structures for *Ichthyophthirius multifilliis G5*, or "Ich"—a parasite that afflicts fish. Ich is a ciliate. Ciliates use a different codon set for stop codons; TGA is the only codon encoding a 'stop' in ciliates, unlike the universal codons of TGA, TAG, TAA, which most organisms use.

## Introduction: <u>Modeling</u> <u>genes</u> <u>in</u> <u>Artemis</u>

We know from our experience working with the various kinds of evidence that some are more reliable than others. In our experience with Ich, we trust them in this order: transcripts & Training set >> TTA1_gw > paramecium > oxytricha & glimmerHMM > genezilla > augustus.

We ran Evidence modeler (EVM) after we ran genefinders, protein alignments and PASA. EVM combined the evidence using this ranking scheme, via its weights file:

| TRANSCRIPT | alignAssembly-img1_122008_pasa | 20 |
|------------|--------------------------------|----|
| PREDICTION | Train | 20 |
| PROTEIN | genewise-Paramecium_peptide_v1.fa | 7 |
| PROTEIN | genewise-oxytricha_all.fasta | 6 |
| PROTEIN | genewise-TTA1.pep | 8 |
| PREDICTION | augustus | 4 |
| PREDICTION | genezilla | 5 |
| PREDICTION | glimmerHMM | 6 |

The output of EVM is a series of gff3 files containing the evidence for each scaffold. EVM is the automated gene prediction for each locus. We will now examine these gene predictions in Artemis, and model some genes.

## SECTION I. Open and configure Artemis.



Open Artemis: Using The Artemis link in course_links, bring up the Artemis page at the Sanger Institute. Locate the link to the Java Web Start version of Artemis. Click on the link (outlined in red in the graphic), and a Welcome screen will come up.

## Setting up the options

The Artemis Welcome screen has 3 menu items: File, Options, and Windows.

- Open Options menu, check the box next to "6. Ciliate Dasycladacean and Hexamita"

- Also click on the "Set working directory" option in the Options menu, near the end of the options list. Set the working directory to your flash drive's data directory.

## Open up the Ich scaffold sequence in Artemis

Open the File menu, and click on "Open…"   A window labeled "Select a file…" will open.

- The first line reads "Look In:  and has a menu with a directory.  Make sure it is pointing to your data directory.
- The last line reads "Files of Type:"—Select "All Files".
- Find the file "genome.fasta" in your data directory.  Select it.  Click the "Open" button.

This will bring up the Artemis viewer.  Across the top are the following menus:

- File, Entries, Select, View, Goto, Edit, Create, Run, Graph, Display.

## SECTION II: Entering the evidence

We will now open up the GFF3 files containing results of our searches and alignments:
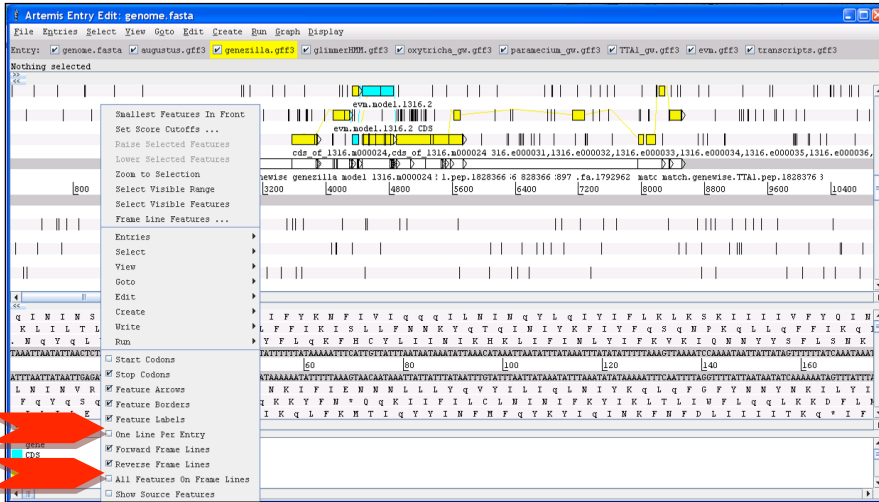
Open, in order, each of the following:

1) augustus.gff3
2) genezilla.gff3
3) glimmer.gff3
4) oxy_gw.gff3
5) para_gw.gff3
6) TTA1_gw.gff3
7) evm.gff3
8) PASA.gff3
9) Train.gff3

Open the File menu in the Artemis viewer, and choose the second line, "Read an entry…"

Change "Files of Type:" to **All Files**.

## Using Artemis

You now have all of the evidence generated by genefinders, spliced protein alignments, and PASA loaded into Artemis.

Put each kind of evidence on its own line.

Right-click and check the following:

-  One line per entry

-  All features on frame lines

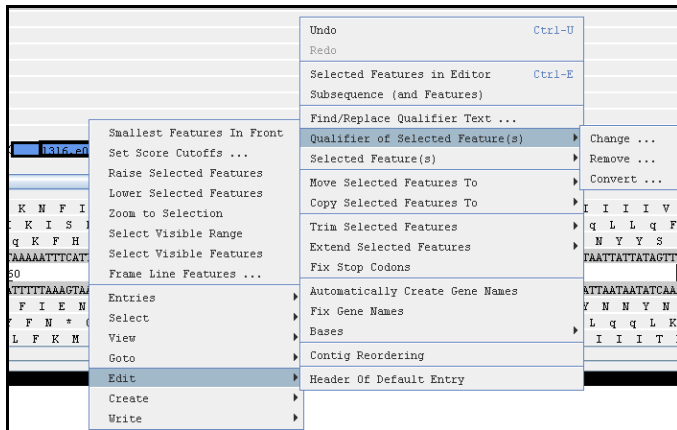As a matter of preference, you might

 click on:

-  Feature labels

## Coloring tracks (optional)

You might wish to color the tracks so you can easily distinguish them.

To choose only one track at a time, so each has a different color, you must uncheck all of the others on the Entry menu or on the Entry bar, which is immediately below the main menu. For example, uncheck everything except "augustus.gff3".
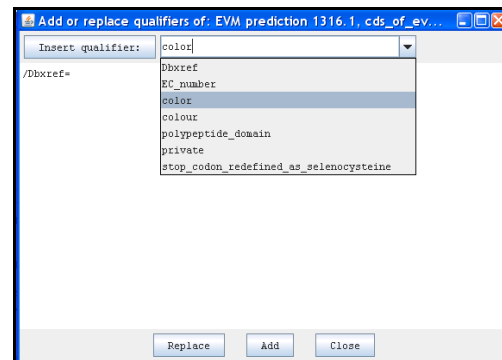
Then, on the Select menu, click on "All".
Right-clicking anywhere inside the viewer will bring up a menu.

A box will pop up. The box is shown below.

Choose 'color ' from the pull-down menu and click the "Insert qualifier" button. This inserts the line:

/color=

Now, type: 100 149 237 [the code for Cornflower Blue.}
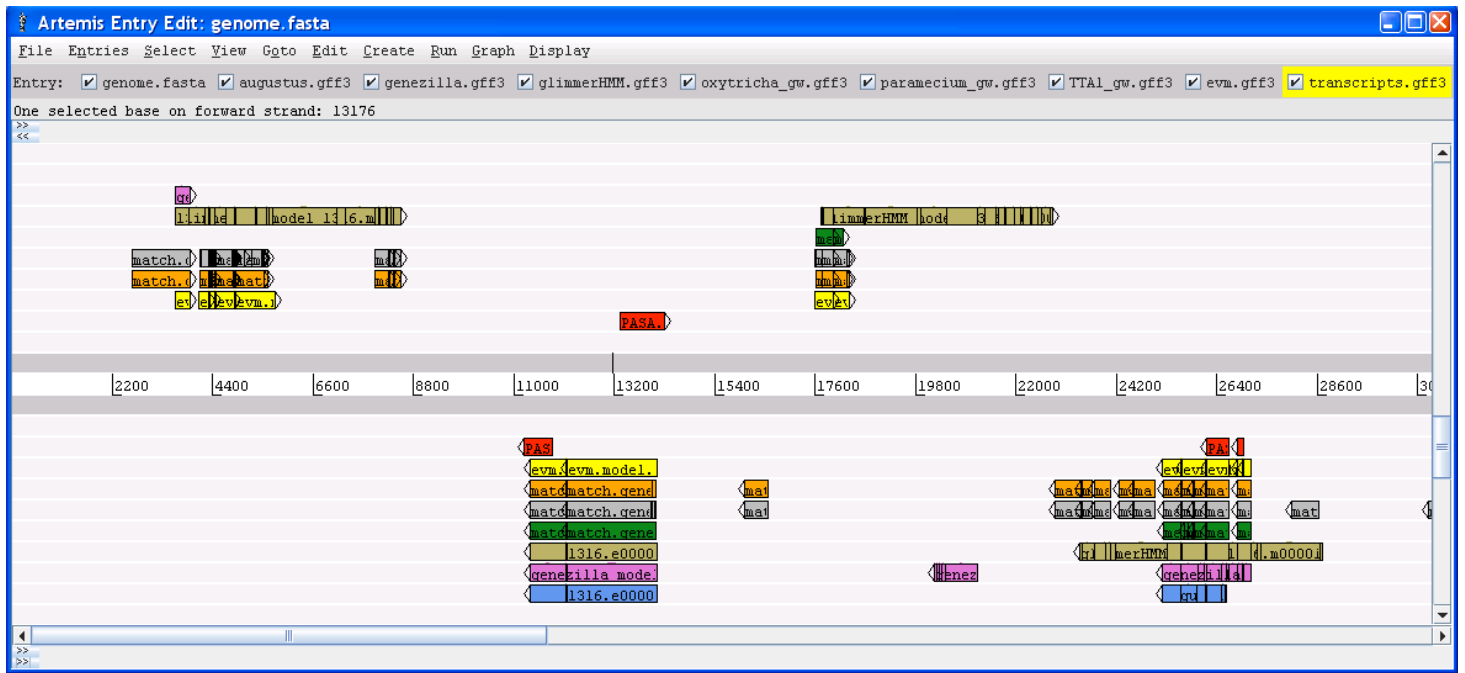
The line is now: /color= 100 149 237

Click "Replace" then "Close".

All of the augustus predictions should have turned blue.

- Choose "Edit", then
- Choose "Qualifier of Selected Feature(s), then
- Choose "Change …"

Repeat this procedure for each color change. Here are some RGB color codes:

| | | | |
|---|---|---|---|
| Forest Green | 34 | 139 | 34 |
| Dark Khaki | 189 | 183 | 107 |
| Cyan | 0 | 255 | 255 |
| Yellow | 255 | 255 | 0 |
| Goldenrod | 218 | 165 | 32 |

| | | | |
|---|---|---|---|
| Indian Red | 205 | 92 | 92 |
| Orange | 255 | 165 | 0 |
| Red | 255 | 0 | 0 |
| Orchid | 218 | 112 | 214 |
| Purple | 160 | 32 | 240 |
| Gray | 190 | 190 | 190 |
| Brown | 139 | 119 | 101 |

Your result might look something like this:



The separation and coloring of lines of evidence, while not necessary, make Artemis tracks much easier to interpret.

# Section III: Modeling genes

We have loaded lots of evidence into Artemis. The evidence included genefinder results, Genewise alignments of protein sets, transcript assembly alignments, and a training set of genes. All of these were used by Evidence modeler with the following weights:

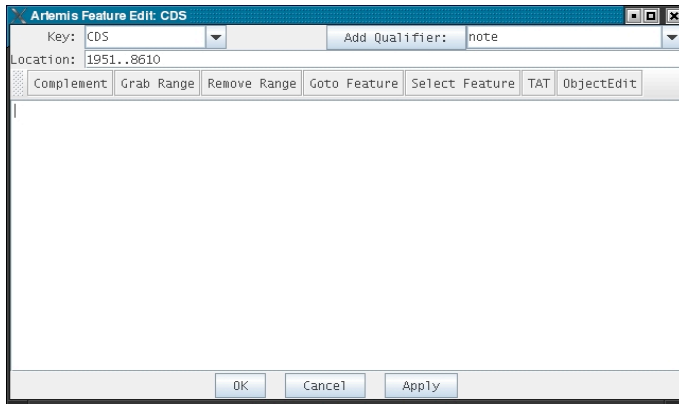| | | |
|---|---|---|
| TRANSCRIPT | alignAssembly-img1_122008_pasa | 20 |
| PREDICTION | Train | 20 |
| PROTEIN | genewise-Paramecium_peptide_v1.fa | 7 |
| PROTEIN | genewise-oxytricha_all.fasta | 6 |
| PROTEIN | genewise-TTA1.pep | 8 |
| PREDICTION | augustus | 4 |
| PREDICTION | genezilla | 5 |
| PREDICTION | glimmerHMM | 6 |

That is, the PASA transcripts and the training set had the most impact on the result, followed by Tetrahymena proteins, etc. Keep in mind that in our experience with Ich, we judge the quality of the evidence in this order: transcripts & Training set >> TTA1_gw > paramecium & glimmerHMM > oxytricha & glimmerHMM > genezilla > augustus. These are reflected in EVM. EVM is the automated gene prediction for each locus. Examine these gene predictions, and model some genes.

Model a gene

Set a default entry to save gene models to. Menu/Entries/Set Default Entry→genes.gff3. Then, save it as a GFF file: File/Save an Entry as… /GFF Format/genes.gff3. The program will ask if you want to overwrite the existing file. Say yes.

Zoom into the region 2000-9000. An easy way to do that is to click on the longest prediction in that region, and hit the "z" key on your keyboard.
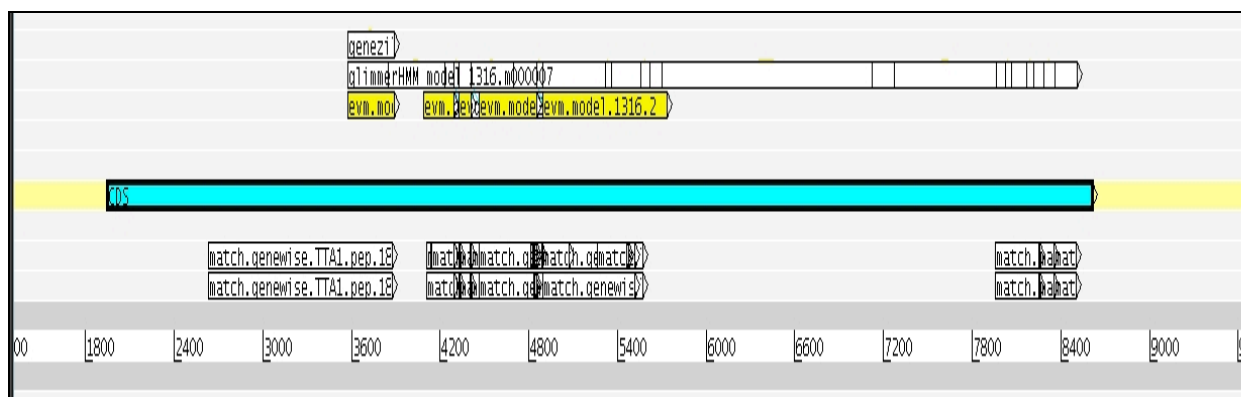
Next, drag your mouse from right to left in the positive direction above the coordinate line (above the numbers). Menu/Create/Feature from Base Range (or Control-C). A window pops up:

Click "Apply."

Note that it has a line that says "Location:". Your coordinates will differ from those shown.

Decide on the structure of the gene you are about to model. Is there an EVM model?



You might notice that EVM predicts two genes in this space. Glimmer predicted one, and genewise matches in three regions. What other evidence is there for the structure of this gene? Which is the best evidence, since there is no transcript evidence?

Let's model the longest gene predicted by EVM in this region first:
1) click on the first EVM predicted exon, with coordinates 4088 to 4286;
2) click on "Grab Range" and the Location line will change to   join(1951..8610,4088..4286)
3) Now delete the first two numbers, leaving: join(4088..4286)
4) Click on the next exon, grab range,  and the Location line will change to : join(4088..4286,4330..4405)
5) Proceed to grab each exon until the Location line reads:  join(4088..4286,4330..4405,4461..4848,4893..5726)
6) Click the "Apply" button.  Click "OK".
Is this a good gene model?  How do we know if it should include the exons that are upstream?

Check the gene model:
One way to gauge how to accurately model this gene is to use BLASTP.   If your gene is not a unique gene, you may be able to discern its correct boundaries based upon orthology. Bring up firefox.  Go to the BLASTP page at NCBI.

Click on your new gene model, and choose "View/Amino Acids of Selection as FASTA" and copy the sequence that pops up.  Paste it into the text box on the BLASTP page.

This is an almost full-length match (84%) to a *Tetrahymena thermophila* protein, which is supported by a transcript. It is probably a real gene.

There is one start site just upstream of the start codon chosen by EVM. Is this one a "better" choice? How would you check?

Save your new gene model. Click on the new CDS, File/Save Default Entry (Ctrl-S).

Zoom in to the region on the forward strand at approximately 65kb. Choose the region approximately between 64kb and 66kb on the forward strand. Drag your cursor over that region. While there is a yellow box from dragging,

Click the menu item " Create", then choose "Feature from Base Range …"

The "Artemis Feature Edit: CDS" window will pop up. Under Location: you will see the coordinates that you swiped across with your mouse.

Now, click on the PASA evidence in one of the tiers. Going back to the Feature edit window, click on "Grab Range". The coordinates of that exon are shown in the "join" statement:

join(64173..66333,64830..65583) [Your first two numbers will be different. They are the range you selected with your mouse.]

Click on the next PASA exon in this evidence tier, "grab range" and repeat for the third exon. You should now have:

join(64173..66333,64830..65583,65639..65725,65768..66000)

Let's now get rid of the approximate coordinates we chose in the beginning with our mouse.

join(64830..65583,65639..65725,65768..66000)

Click OK.

Because this gene model was based on an EST assembly, it may have untranslated regions (UTRs). Determine the appropriate start and stop codons, if possible (BLASTP is one way to do this). To trim it back to the start and stop codons, right-click the model, select Edit/Trim Selected Features/To Met. To extend to the next stop codon, select Edit/Extend Selected Features/To Previous Stop Codon or To Next Stop Codon. In some cases, it is easier to visually identify the coordinates of the stop codon and enter them in the editor. In this case, the final coordinates would be 64830... 65811.

If you wish to add UTRs, follow the instructions in the Genewise tutorial.

As time permits, continue to model genes.

Save them in the genes.gff3 file as you go.