# Annotation Tutorial:
# *Diaphorina citri* genome

Indian River State College, FL

# Websites you will need

Basecamp:  https://basecamp.com/

Apollo:  https://apollo.nal.usda.gov/diacit/sequences

i5k Blast or hmmer:  https://i5k.nal.usda.gov/webapp/blast/

https://i5k.nal.usda.gov/webapp/hmmer/

NCBI Blast:  https://blast.ncbi.nlm.nih.gov/Blast.cgi

MCOT database: https://citrusgreening.org/tools/blast?db_id=27

# Chose Gene or Gene family from Basecamp



This page will take you to orthologous sequences & Gene family description

## Orthologs

Test **ALL** orthologous sequences one by one in Apollo (Next slide)

# "Blat" sequences in Apollo

Menu bar: Click Tools > Search Sequence
- Enter orthologous sequence (not in FASTA)
- Indicate nucleotide or protein
- Click Search



This is one scaffold

This is another scaffold, etc.

Results in different scaffolds

Change to BLAT protein if using a protein sequence (most often).
Be sure to check "search all genome sequences" if you want anything not already on your page.
Insert the letter sequence only. Do not use the > or any spaces.

- Look at all scaffolds
- Significance that is very low and Score that is high is ideal
- *Remember* Test all orthologs given from basecamp

# Apollo Search Results

Score: A number that the BLAT results "scores" the hit. The higher is often better, but not always. There is a balance with score and Significance.

Significance: The "E-value" given by the BLAT results. The smaller the number is often better, but not always. There is a balance with Score and Significance.

ID: The scaffold location

Start or End : The number where the sequence starts or ends

Identity: How well the sequence entered matches the location of the scaffold.

Search sequence  ✕

BLAT protein ▼

Enter sequence

YVASCCKICRYDDNKSCSHGSVFRTWNFLHKRGSVTGGDYGDRTGCQPSTISPCSHHGSAPTLPSCENQK
VPKLKCHTRCTNPTYGRGFFQDKHRTTLTYWVDDNEDAIKKEILAHGPTTATFALYDDFYHYKSGVYKHT
SNAKLENYLHSGKLIGWGTENGTPYWLVINTWGPHWGDRGTVKILRGKYECAFEYLIAAGKPKN

☑ Search all genomic sequences

Search

| ID | Start | End | Score | Significance | Identity | |
|---|---|---|---|---|---|---|
| gi\|6455066... | 143862 | 143634 | 165 | 5.8e-41 | 100 | gi\|6455066... |
| 141257 | 141047 | 155 | 4.9e-38 | 100 | gi\|6455066... | 142334 |
| 142142 | 143 | 2.8e-34 | 98.44 | gi\|6455066... | 143078 | 142907 |
| 123 | 2e-28 | 100 | gi\|6455066... | 124686 | 124563 | 86 |
| 3.7e-17 | 97.56 | gi\|6455066... | 141949 | 141838 | 80 | 2.1e-15 |
| 100 | gi\|6453697... | 2 | 107 | 79 | 4.7e-15 | 100 |
| gi\|6455066... | 144400 | 144289 | 76 | 3.1e-14 | 100 | gi\|6455066... |
| 138371 | 138260 | 69 | 6.9e-12 | 89.19 | gi\|6452510... | 280 |
| 193 | 67 | 2.7e-11 | 96.55 | gi\|6453555... | 10190 | 10277 |
| 55 | 8.8e-8 | 96.55 | gi\|6455077... | 172463 | 172346 | 54 |
| 2e-7 | 64.1 | gi\|6455049... | 5657 | 5540 | 51 | 0.000001 |
| 58.97 | gi\|6455066... | 147551 | 147671 | 48 | 0.000015 | 52.5 |
| gi\|6455066... | 147725 | 147788 | 43 | 0.00038 | 85.71 | gi\|6455066... |
| 126119 | 126065 | 42 | 0.00065 | 100 | gi\|6455066... | 127962 |

This example shows the columns out of place, this may happen, The ID section starts with gi.

Choose a scaffold. By clicking on it. You can then X out of the dialog box.

# If Apollo says "No matches found"

Use i5k Blast



Indicate organism

Indicate protein
(if applicable)

Enter sequence

Results

- Chose best result(s) (Low e-value, best coverage, etc.)
- Copy sequence and Blat in Apollo.

**Apollo Blat vs i5k Blast**
- Apollo Blatt is more "sensitive" than i5k Blast.
- i5k Blast will broaden the search.

# If Apollo says "No matches found"

Search your sequence on *D. citri* MCOT database

https://citrusgreening.org/tools/blast?db_id=27 **BLAST**

**Results**

⊖ **Input parameters**

| Categories | Psyllid Databases ▼ |
|---|---|
| Database | Diaphorina citri MCOT proteins ▼ db details |
| Program | blastp (protein to protein db) |
| Query | autodetect ▼ Show example |

Indicate MCOT protein DB

Indicate blastp, blastn, etc.

```
MTSESYKLFVESHPRFTTFSNYSLLKSYERPRCHFVFTLAQEHRAVAKLVGPNVSGNITF
TQSGSILLISGVVEGLKPKSTHGFHIHEKGDLSSGCASTGGHFNPYNKHHGGPTDEERHV
GDLGNIDANEHGVAAFTLSDHVASLVGPNCIIGRGVVLHSDPDDLGKGQHPDSLTTGHAG
SRIACGVIGTLDPGTDKLENSASRSSPYFTILTVFVVFLKLTN
```

Your sequence

⊖ **Results**

Basi...

**Untitled_sequence vs Diaphorina citri MCOT proteins**

| SubjectId | id% | Aln | evalue | Score | Description |
|---|---|---|---|---|---|
| MCOT10235.0.MM | 57.04 | 81/142 | 1e-48 | 163 | \| Superoxide dismutase [Cu-Zn] \| Similar to C4WTR6 \| *-*- \| PANTHER PTHR10003 \| Pfam PF00080 Length = 274 |
| MCOT13840.0.CO | 55.10 | 81/147 | 2e-46 | 157 | \| Superoxide dismutase [Cu-Zn] \| Similar to A0A023FAY3 \| ***- \| PANTHER PTHR10003 \| Pfam PF00080 Length = 238 |
| MCOT16518.0.CC | 54.05 | 80/148 | 5e-44 | 148 | \| Superoxide dismutase [Cu-Zn] \| Similar to R4V538 \| ***- \| PANTHER PTHR10003 \| Pfam PF00080 Length = 155 |

Click for sequence, and blat on Apollo

**BLAST**

# If Apollo says "No matches found"

Search your sequence in i5k blast hmmer (hidden Markov model) created from a multiple sequence alignment of the ortholog proteins

https://i5k.nal.usda.gov/webapp/hmmer/

HMMER

**Organisms**

- ☐ *Cimex lectularius*
- ☐ *Copidosoma floridanum*
- ☐ *Diachasma alloeum*
- ☑ *Diaphorina citri*
- ☐ *Drosophila biarmipes*
- ☐ *Drosophila bipectinata*
- ☐ *Drosophila elegans*
- ☐ *Drosophila eugracilis*
- ☐ *Drosophila ficusphila*
- ☐ *Drosophila kikkawai*
- ☐ *Drosophila rhopaloa*
- ☐ *Drosophila takahashii*

**Diaphorina citri**

Protein

- ☑ Protein - Diacit_RefSeq_proteins_Release_100
- ☑ Protein - Diacit_International_psyllid_consortium_proteins_v1

Tutorial
https://i5k.nal.usda.gov/webapp/hmmer/manual/

**Query Sequence / Mutliple sequence alignment**

Your sequence is detected as fasta:

```
MTSESYKLFVESHPRFTTFSNYSLLKSYERPRCHFVFTLAQEHRAVAKLVGPNVSGNITF
TQSGSILLISGVVEGLKPKSTHGFHIHEKGDLSSGCASTGGHFNPYNKHHGGPTDEERHV
GDLGNIDANEHGVAAFTLSDHVASLVGPNCIIGRGVVLHSDPDDLGKGQHPDSLTTGHAG
SRIACGVIGTLDPGTDKLENSASRSSPYFTILTVFVVFLKLTN
```

# On Apollo, once you have selected a scaffold.

Be sure that the items you want to view are chosen on the left side of the screen. For the psyllid, we start to base our models off of NCBI predictions. Be sure the following gene sets are chosen.



We will base our models off of RNAseq data, so be sure to choose some tracks under **Mapped Reads** as well.

Note: Be sure to check all RNAseq data, but it may overload the system if they are all turned on at once.

# Predicted Gene Models



*Useful tip*: You can use this as an "URL" in Apollo. As you need to look at all scaffolds (BLAT results), use this URL to organize all your potential gene models and to get back to them.

You may choose a potential gene model (Usually starting with XM) by clicking on the model and dragging it up to the User Created section and dropping it.

Zoom out if need be and locate the ends of a potential gene model.

The highlighted area shows where the section within the BLAT results of the previous step. (slides 4 & 5)

# Check Gene Models

Right click on the model in the User Created section and choose "Get sequence"



Highlight the entire section in its FASTA format and copy it.

# Protein BLAST



Use blastp for protein BLAST

Enter the FASTA sequence into the appropriate box, choose any selection / parameters. And press BLAST at the bottom.

Indicate "Insecta". This will limited your search to insects only.

# The NCBI Results Page



Job title: CAA8EC20D7495772A225C37C77071736 (sequence:mRNA)...

RID AA3YKJER015 (Expires on 02-17 09:34 am)
Query ID lcl|Query_6881
Description CAA8EC20D7495772A225C37C77071736 (sequence:mRNA) 374 residues [gi|645505060|ref|NW_007378205.1|:108329-117924 + strand] [peptide]
Molecule type amino acid
Query Length 374

Database Name nr
Description All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environment
Program BLASTP 2.6.1+ ▶ Citation

Other reports: ▶ Search Summary [Taxonomy reports] [Distance tree of results] [Multiple alignment]

The top portion gives information on the searched sequence

## Graphic Summary

Show Conserved Domains

Putative conserved domains have been detected, click on the image below for detailed results.

The NCBI conserved domains database is used to identify the conserved domains in the orthologs and candidate genes.

Query seq.
active site
S2 subsite
Specific hits          Peptidase_C1A_CathepsinB
Superfamilies        Peptidase_C1 superfamily

Distribution of the top 101 Blast Hits on 100 subject sequences
Mouse over to see the title, click to show alignments

Color key for alignment scores
<40    40-50    50-80    80-200    >=200

Query
1    70    140    210    280    350

Query= your sequence

Results

You can see how the results align with your sequence

| | | score | score | cover | value | |
|---|---|---|---|---|---|---|
| PREDICTED: cathepsin B-like cysteine proteinase 4 [Diaphorina citri] | | 786 | 786 | 99% | 0.0 | 100% XP_0084 |
| cathepsin B [Riptortus pedestris] | | 248 | 248 | 92% | 4e-78 | 40% BAN2037 |
| PREDICTED: cathepsin B-like isoform X1 [Halyomorpha halys] | | 238 | 238 | 94% | 6e-74 | 39% XP_0142 |
| cathepsin B [Riptortus pedestris] | | 237 | 237 | 95% | 1e-73 | 39% BAN2146 |
| cathepsin B-like cysteine protease [Triatoma infestans] | | 235 | 235 | 94% | 8e-73 | 38% ABD3530 |
| PREDICTED: cathepsin B-like [Cimex lectularius] | | 234 | 234 | 91% | 2e-72 | 39% XP_0142 |
| PREDICTED: cathepsin B-like [Cimex lectularius] | | 234 | 234 | 92% | 3e-72 | 38% XP_0142 |
| hypothetical protein g.26476 [Graphocephala atropunctata] | | 234 | 234 | 94% | 3e-72 | 38% JAT2482 |
| PREDICTED: cathepsin B-like isoform X1 [Cimex lectularius] | | 233 | 233 | 94% | 6e-72 | 38% XP_0142 |
| PREDICTED: cathepsin B-like isoform X2 [Cimex lectularius] | | 232 | 232 | 94% | 7e-72 | 39% XP_0142 |
| hypothetical protein g.21611 [Homalodisca liturata] | | 232 | 232 | 94% | 1e-71 | 37% JAT0683 |
| PREDICTED: cathepsin B [Tribolium castaneum] | | 231 | 231 | 94% | 3e-71 | 37% XP_9742 |

# The NCBI Results Page cont'd

Pairwise Alignment

cathepsin B-348 precursor [Acyrthosiphon pisum]
Sequence ID: ref|NP_001119608.1| Length: 342 Number of Matches: 1
▶ See 1 more title(s)

Range 1: 9 to 341 GenPept  Graphics          ▼ Next Match  ▲ Previous Match

| Score | Expect | Method | Identities | Positives | Gaps |
|-------|--------|--------|------------|-----------|------|
| 229 bits(583) | 5e-68 | Compositional matrix adjust. | 132/346(38%) | 180/346(52%) | 17/346(4%) |

```
Query  1    MIHILVFLLG---CTLVRGELYKFSDAYIDQINREANTWTAGRNFPANLSEEYLRQFLIA  57
            ++ +L+F G      VR +L    SD +ID IN    W+AGRNF  +   Y++  +
Sbjct  9    LVGLLIFSFGRVDGATVRVDLNPLSDEFIDHINSIQYYWSAGRNFHKDTPISYIKGLMGV  68

Query  58   DAKYFDQSDRPLPGDRKTYDPEYSATVPDRFDAREQWPNCGTIGHVPDTGACAAPHIFAA  117
             K    ++ P      TY+ + S  +P+ FDARE+WPNC TI  V D G+C +    F A
Sbjct  69   HEK---NAEYPKLEQLLTYN-DASTDLPETFDARERWPNCPTIREVRDQGSCGSCWAFGA  124

Query  118  VGAFSDRRCIKSKGQQNRPLSTEYVASCCKICRYDDNKSCSHGSVFRTWNFLHKRGSVTG  177
            V A SDR CI S G +N  S E + SCC  C +    C+ G    WN+   +G V+G
Sbjct  125  VEAMSDRVCIHSNGTKNFHFSAENLVSCCWTCGF----GCNGGFPGAAWNYWKTKGIVSG  180

Query  178  GDYGDRTGCQPSTISPCSHHGSAPTLPSCENQKVPKLKCHTRCTNPTYGRGFFQDKHRTT  237
            G YG   GC P  I+PC HH +      P E  K P  C  +C    Y   + QD H
Sbjct  181  GPYGSNMGCIPYEIAPCEHHVNGTRGPCKEGGKTP--TCVKKCEE-GYKVPYAQDLHHGK  237

Query  238  LTYWVDDNEDAIKKEILAHGPTTATFALYDDFYHYKSGVYKHTSNAKLENYLHSGKLIGW  297
             Y + ++ D I++EI  +GP    F +Y+DF Y++GVYKH +   L    H+ +++GW
Sbjct  238  SAYSIRNDVDQIRQEIYTNGPVEGAFTVYEDFIAYRAGVYKHVAGKALGG--HAIRILGW  295

Query  298  GTENG-TPYWLVINTWGPHWGDRGTVKILRGKYECAFEYLIAAGKP  342
            G +NG  PYWLV N+W   WG  G  KILRG  EC  E  I AG P
Sbjct  296  GVQNGEIPYWLVANSWNTDWGSDGFFKILRGSDECGIEGQINAGLP  341
```

You can see how similar your gene model is to other similar genes in related organism(s), by clicking on each result (previous slide).

Example: The gene model is 38% identical and 52% similar to gene in *A. pisum*.

# NCBI Smart BLAST

Or check gene model by using

https://blast.ncbi.nlm.nih.gov/smartblast/?LINK_LOC=BlastHomeLinkn



Smart BLAST will provide a phylogenetic tree with model organisms.

Conserved domains within the sequences

Green: subject sequences
Yellow: your sequence (query)

Sequence name that corresponds to sequences on right. (hover over to get more details)

More info. on query vs model organisms.

End of page would be BLASTp results.

# RNASeq Data

To view gene model in different colors indicating the diff. frames: (As seen) Click View, Color by CDS frame.



The colored outlined boxes represent exons, while the thin blue lines connecting them represent introns. The blue boxes without an outline represent UTRs

Look closely at the results and see how they line up with RNASeq data. You may have to zoom in close and look by sections.

Dark blue RNA seq reads indicate evidence for exons and light blue for introns.

# RNASeq Data
## More details

You can also view the RNAseq Coverage Plots.
(On the left side on Apollo (Available Tracks))



If you zoom in and see letters, this is the reference sequence, as shown here. You can view stops on the different frames by an *.

# Making Edits

## Matching Data

UTR's are dark blue and are on the 5' and 3' ends of the gene model.

Example: This is the 3' UTR of the gene model

Delete exons, or UTR (Untranslated region)

Merge exons, or gene models

Split

Duplicate

Make Intron

Set to Downstream Splice Donor

Set to Upstream Splice Donor

Set to Downstream Splice Acceptor

Set to Upstream Splice Acceptor

View History

Indicates which direction (Upstream or Downstream) the sequence is reading.
In this case = downstream (-) strand

**Important:**

**Splice sites:**
Reverse (-) strand gene model:
3' acceptor of intron, GA
5' donor of intron, TG, CG
Positive (+) strand gene model:
3' acceptor of intron, AG
5' donor of intron, GT, GC

Make sure your gene model starts with ATG (M) and ends with TAG, TGA, or TAA (stop codons).

# Fixing splice sites



**Splice sites:**
Reverse (-) strand gene model:
3' acceptor of intron, GA
5' donor of intron, TG, CG
Positive (+) strand gene model:
3' acceptor of intron, AG
5' donor of intron, GT, GC

5' donor of **intron**

3' acceptor of **intron**

Positive strand

For this example, the donor would be moved upstream.
(RNAseq evidence to do so)
Right click exon, Set to upstream splice donor.

Zoom in

# Adding an exon



Other gene predictions can be used, not limited to, Maker, augustus, snap.

gi|645505060|ref|NW_007378205.1|:108328-117924

Maker gene predictions
maker-s772-snap-gene-0.70-mRNA-1

augustus_masked
augustus_masked-s772-abinit-gene-0.30-mRNA-1

snap_masked
snap_masked-s772-abinit-gene-0.52-mRNA-1

Control adult and nymph psyllids (Vyas 2015; Coverage Plot)

Control adult and nymph psyllids (Vyas 2015)

Click and drag exon on top of gene model and drop.

There is evidence that an exon should be added to the gene model

After

5505060|ref|NW_007378205.1|:108328-117924

s772-snap-gene-0.70-mRNA-1

ked-s772-abinit-gene-0.30-mRNA-1

2-abinit-gene-0.52-mRNA-1

as 2015; Coverage Plot)

# Isoforms

When the gene can undergo alternative splicing



RNAseq reads indicate intron and exon evidence for the same part of the model.

If there were many results containing isoforms when the gene model was initially blasted on NCBI, this is good evidence for isoforms as well.

If you have evidence for isoforms, right click your gene model in the "user-created annotations" space, and duplicate.

# Edited Gene Models



Predicted Model Conserved domain results

Conserved domain results after matching your gene model with RNASeq data

**After a change is made, check your results using the BLASTp or Smart BLAST again**

# Information editor

Information Editor (alt-click)    ✕

Select mRNA   Dcitr_cathepsin B-like 1 prot ISO 2   ▼

## gene

| | |
|---|---|
| Name | Dcitr_cathepsin B-like 1 prot ISO 2 |
| Symbol | |
| Description | |
| Created | 2015-12-17 |
| Last modified | 2016-11-04 |

### Status
◯ Approved   ◯ Delete

### DBXRefs

| DB | Accession |
|---|---|
| | |

Add   Delete

### Replaced Models

| Action | Transcript Name |
|---|---|
| | |

## mRNA

| | |
|---|---|
| Name | Dcitr_cathepsin B-like 1 prot ISO 2 |
| Symbol | |
| Description | |
| Created | 2015-12-17 |
| Last modified | 2016-01-28 |

### Status
◉ Approved   ◯ Delete

### DBXRefs

| DB | Accession |
|---|---|
| | |

Add   Delete

### Replaced Models

| Action | Transcript Name |
|---|---|
| replace | XM_008472813.1 |

Name of the gene. Do NOT use Dcitr

Symbol or abbreviation of the gene

Description of the gene, usually function. Uniprot is a good source.

Click approved once all edits are made and your gene model is complete

You will need the name of the original model that your annotation will be replacing

Example:

NCBI Predicted protein coding genes, Annotation Release 100

008471046.1

XM_008472813.1
cathepsin B-like cysteine proteinase 4 isofo

mRNA XM_008472813.1

**Primary Data**

| Name | XM_008472813.1 |
|---|---|
| Type | mRNA |
| Description | cathepsin B-like cysteine proteina |

Right click on model > View details

Fill in both sections: Gene and mRNA panels

# Information editor cont'd

Add    Delete

### Comments

Add    Delete

### Comments

NCBI model. RNA Seq data avail. BLAST match with high query and ID.

Pairwise BLAST to "Dcitr_cathepsin B-like 1 prot" shows 84% query and 43% ID

Add    Delete

Here you will enter comments about all the edits (splice sites moved, added exons, deleted exons, etc.) made to your gene model and other relevant info.

## Naming isoforms:

### gene

| Name | Dcitr- Tudor domain-containing pr |
| Symbol | TDRD7 |
| Description | involved in post-transcriptional reg |
| Created | 2015-10-02 |
| Last modified | 2015-11-05 |

Status
◉ Approved  ◯ Delete

### mRNA

| Name | ntaining protein 7 RA |
| Symbol | TDRD7 |
| Description | involved in post-transcriptional reg |
| Created | 2015-10-02 |
| Last modified | 2017-02-15 |

Status
◉ Approved  ◯ Delete

Naming convention for isoforms should have RA, RB, RC, etc. at the end for each isoform you have, respectively. (look at slide 20)

The gene section should be the same for all isoforms.

### gene

| Name | Dcitr- Tudor domain-containing pr |
| Symbol | TDRD7 |
| Description | involved in post-transcriptional reg |
| Created | 2015-10-02 |
| Last modified | 2015-11-05 |

Status
◉ Approved  ◯ Delete

### mRNA

| Name | domain-containing protein 7 RB |
| Symbol | TDRD7 |
| Description | involved in post-transcriptional reg |
| Created | 2015-11-02 |
| Last modified | 2017-02-15 |

Status
◉ Approved  ◯ Delete

But the mRNA section should be different for all isoforms, comments should contain the edits made for the particular isoform clicked on.

# Gene Family Report

Once your gene or gene family is annotated, the gene report should include:

- Introduction to the gene or gene family including the pathway (use literature)

- Methods of annotation and phylogenetic analysis

- Results and Discussion (use literature)

- References cited

# Tables to include in report

## Table with gene count in:

- *Drosophila melanogaster*
- *Anopheles gambiae*
- *Tribolium castaneum*
- *Apis mellifera*
- *Nasonia vitripennis*
- *Acyrthosiphon pisum*
- *Bemicia tabaci*

## Table with blast match:

- Indicate organism and name of gene
- Query coverage
- % identity
- Bit score

Include legend for each table.
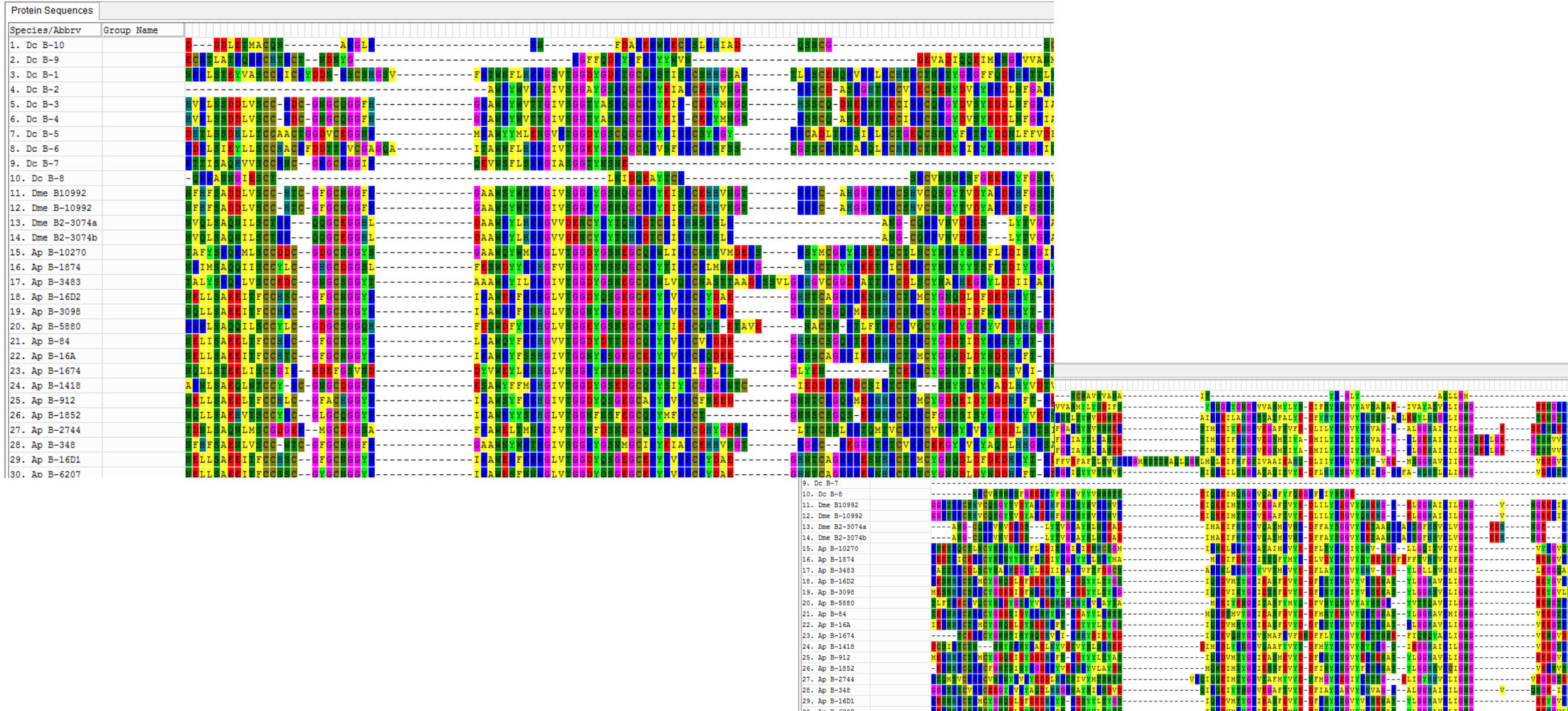
# Comparative and Phylogenetic Analysis

When your gene model(s) are completed, perform analysis in MEGA.

1. Construct tree with *D. citri* gene models only.

2. Construct tree with *D. citri* gene models and orthologs. Use related sequences from, but not limited to,
   - *Drosophila melanogaster*
   - *Anopheles gambiae*
   - *Tribolium castaneum*
   - *Apis mellifera*
   - *Nasonia vitripennis*
   - *Acyrthosiphon pisum*
   - *Bemicia tabaci*

Can find orthologous sequences in NCBI, Ensembl, Uniprot, etc.
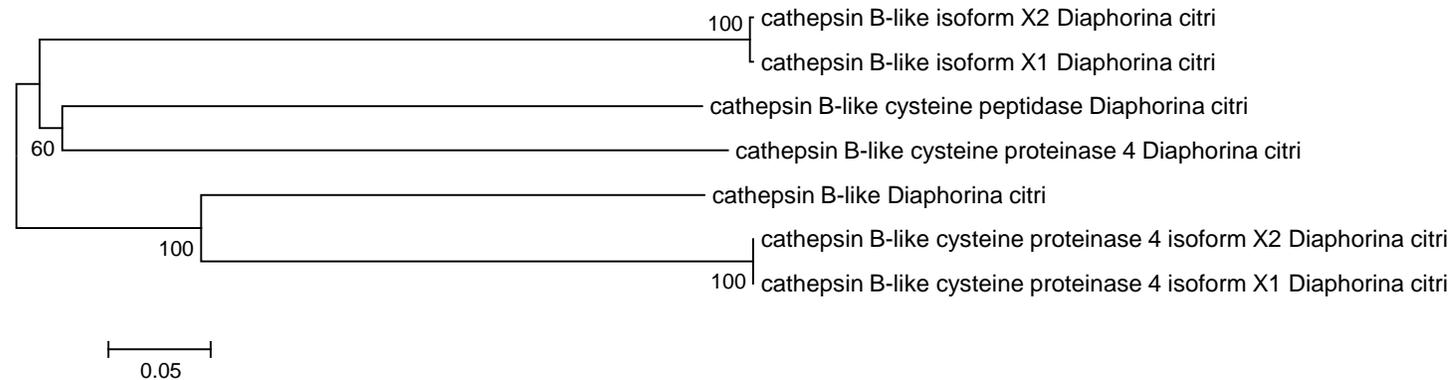
# Multiple Sequence Alignments (MSA)

Multiple sequence alignments should be generated using MUSCLE, tcoffee, or clustal to compare the ACP gene model to the query gene set. The final model should be refined in Apollo using homology, RNAseq and proteomics evidence tracks.
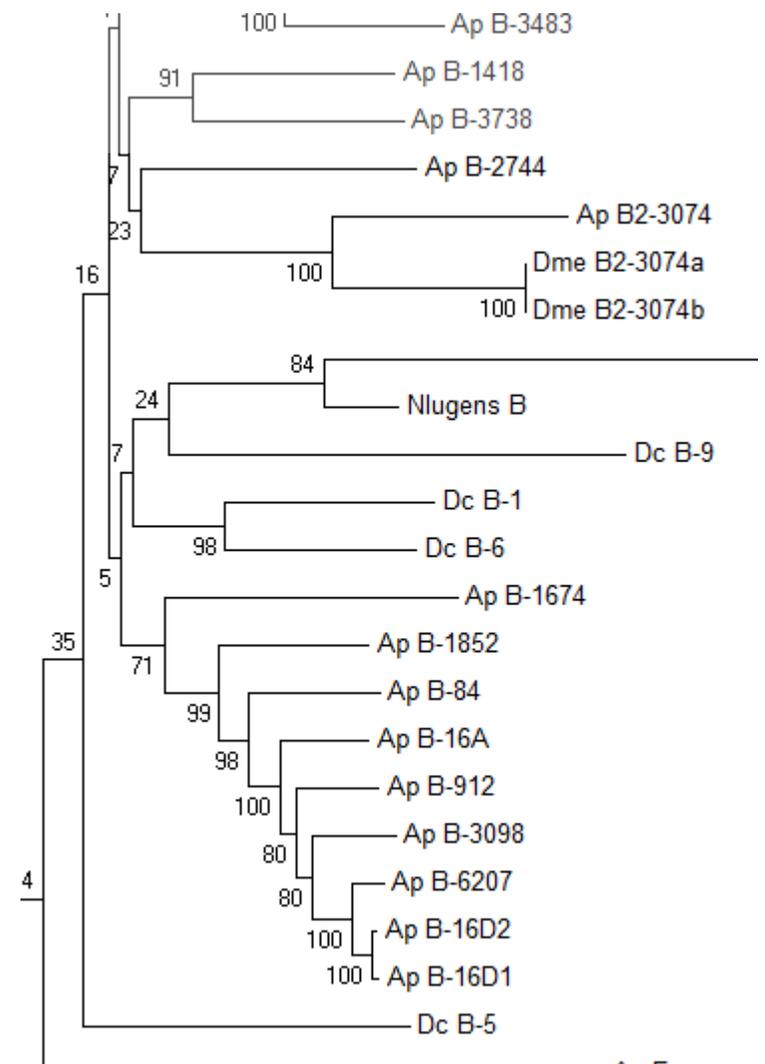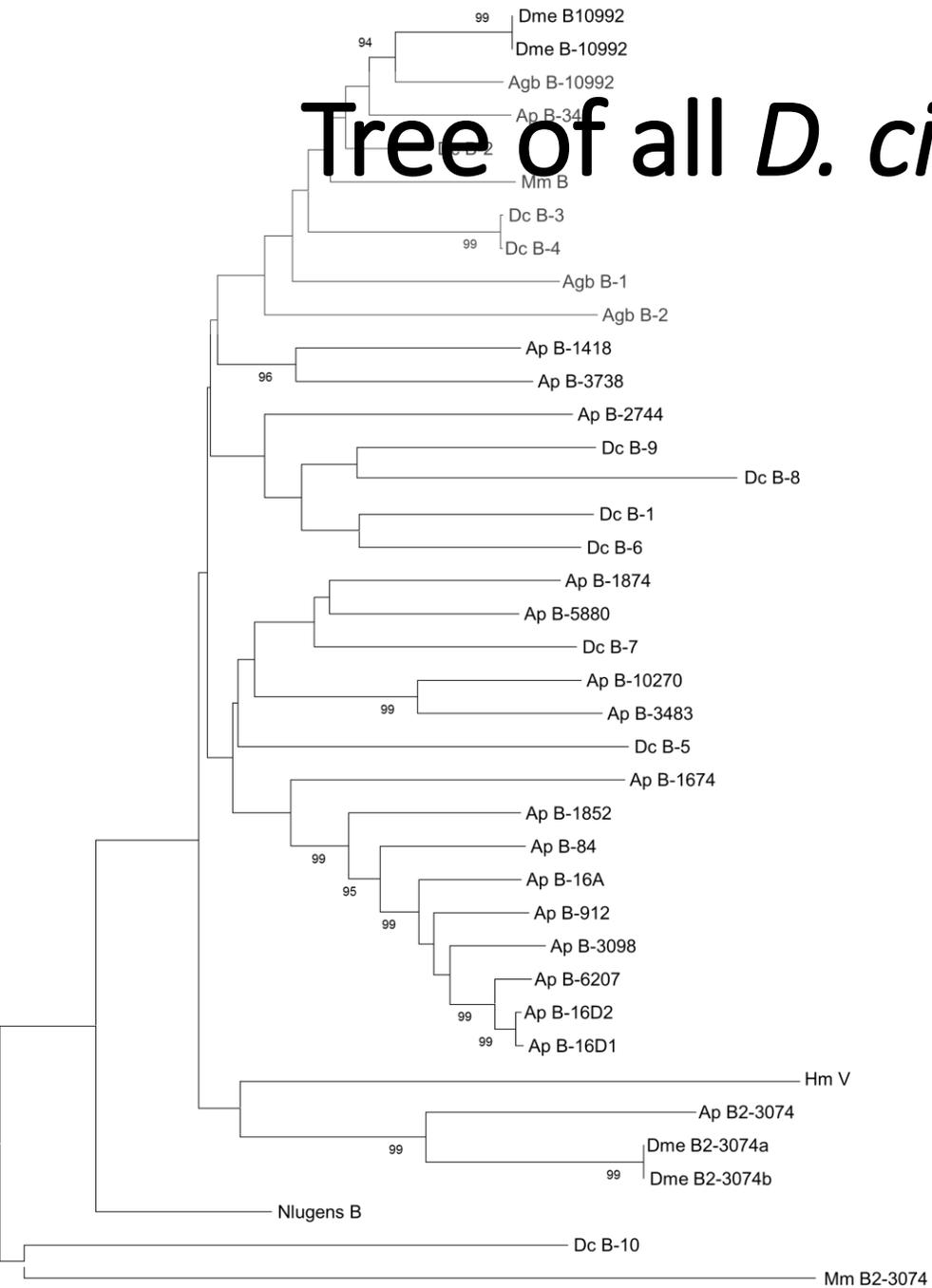
# Tree of all *D. citri* gene models

This would include all genes (paralogs) for the particular family annotated in the genome.

Example:



Tree should be included in report with legend.

# Tree of all *D. citri* gene models and other orthologs



Tree should be included in report with legend

Krystal Villalobos Ayala
Chris Cordola
Tracey Bell
Hannah Mann
Daniel DeAvila
Gabe DeAvila
Tom D'elia

Indian River State College, FL