**Phagehunting Program**

Phagehunting Protocols

PREPARATION · ISOLATION · PURIFICATION · AMPLIFICATION · EXTRACTION · CHARACTERIZATION · SEQUENCING · ANNOTATION · PHAMERATION · FURTHER DISCOVERY

# What's New for School Year 2014-15 in Phage Genome Annotation

Created by DJS December, 2014.

The What's New for 2015 includes changes in:
1. An additional file is required when submitting for QC
2. Notes required at submission
    a. An additional entry (ST)
    b. More specific explanations required, especially concerning the reporting of Functions
3. Guiding Principles Update
4. Shine-Dalgarno evaluation. This includes a new set of matrices for evaluating the relevance of the Shine-Dalgarno sequence
5. GeneMark
    a. Ease of Use
    b. Change of Output
6. DNA Master BLAST parameters
7. FASTA Output for ALL ORFs.
8. Starterator
9. Genome Announcements

**General Note:** Some of these changes will impact your classroom practice, so read carefully and don't hesitate to email with for further clarifcation.

1. **Additional File Required** (**Section 12**, pages 125 – 127)
   When you submit your final annotation for QC, we want you to submit two files. One of the files contains all of the info we have asked for in the past (with new additions, so read Section 12) labeled [YourPhageName]_Complete Notes.dnam5. The second file contains only reportable* functions (and notations sometimes) in the Notes section labeled [YourPhageName}_Final.dnam5.

   *Reportable function is the operative word here. You and your students will find lots of 'hints' to possible functions. Record those in your CompleteNotes file and hypothesize away as to the hows and whys a phage could use that function, but we are only going to report functions to NCBI that follow the guidelines that we have established in the annotation guide. Discernment is key!

2.  **Notes required at submission** (**Section 9.6**, pages 103 -104)
    Once again, we are messing with the notes section of your annotations.  Read
    the explanation of how to complete the notes for each gene call carefully!
    (p.104)  Improve your functional annotation by  carefully describing the
    supporting documentation for your assigned functions (**See Section 10**).

    We have a new program for you to use this year – Starterator.  See #**8** of this
    What's New for more information!  Just remember to include your Starterator
    data in the ST section of your notes.

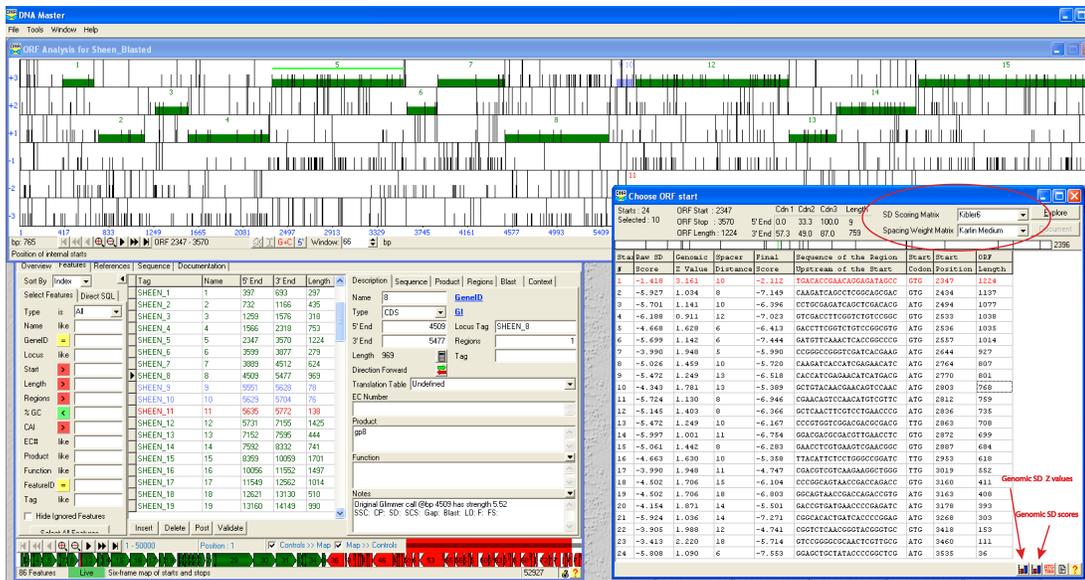3.  **Guiding Principles Update (Section 7,** p. 64-65).
    The Guiding Principles are always worth reading each year, but this year we
    have revised them a bit.  The main concept that we want to reinforce is that
    the evaluation of the Shine-Dalgarno score takes a back seat to everything else
    on the list.  One of the considerations Welkin and I would like to showcase is
    #**10a**.  This one states that the ribosome likes it when the genes have a 4 bp
    overlap.  That overlap allows the ribosome to hop to the next gene with the
    least possible effort.  If you have the option to select a start with a 4 bp
    overlap, we recommend that you do!

4.  **How to Use the New Shine-Dalgarno Matrices (Section 8.4.2,** p. 74-75)
    This is the probably the change that will cause you to pause this year.  We
    introduced the newly developed Shine Dalgarno scoring matrices last year,
    but this year we will insist that you use them.  It is a bit daunting as you
    begin to use this, but relax, you can do it!

    Here is the skinny.  In the Choose ORF start window, set the SD Scoring
    Matrix to Kibler 6 and the Spacing Weight Matrix to Karlin Medium.  The
    new algorithms are based on the assumption that if you evaluated all of the
    starts across the genome – that includes the "real" starts, the additional
    (other) starts in an ORF with coding potential, and EVERY other possible start
    (for the frames with no coding potential) – you would get a normal
    distribution of 'scores'.  The "real" starts would actually be found outside that
    normal distribution of (ridiculous) starts because they would be remarkably
    different than the random start sites included in this data set.

    There are 2 scores to evaluate:  Raw SD score and Final score.  These scores
    are based on the math as described by D. Kibler & S. Hampson (see references
    below).  The Raw score compares the upstream sequence of a start to the SD
    sequence of E. coli – AGGAGGA.  The Final score uses that score plus the
    spacing (which is 7 – 10 bp upstream of the start).  In either case the best score
    is the largest score (which is the **least negative number).**  The second number
    to help you evaluate the score is the Z score.  You are looking for the
    Raw/Final score to be greater than 2 standard deviations from the mean.

References for the Scoring Matrices:
http://www.ics.uci.edu/~kibler/pubs/Metmbs02.pdf
http://www.ics.uci.edu/~kibler/pubs/TR03.pdf

5. **GeneMark changes** (**Section 5.3**, pages 45 – 48)
    a. The good news:  Dan has added three buttons to your sequenced phage page.  Each button will automatically generate the pdf required for 3 GeneMark predictions:
        i. GeneMark heuristic – this is the graph of coding potential that matches the output from GeneMark provided to DNA Master.  In this case, the program is trained on the sequence of the phage that you submit to it.
        ii. GeneMark smeg – this is the graph of coding potential when GeneMark is trained on the model, *M. smegmatis* mc²155.
        iii. GeneMark TB – this is the graph of coding potential when GeneMark is trained on the model, *M. tuberculosis* HR37Rv.
    b. The not so good news:  if you are like me, you really liked the landscape output of GeneMark.  It seems that the only output that is available now is available in a portrait orientation.  The really bad news is that GeneMark's servers were hacked this summer and all GeneMark programs were unavailable for a number of weeks this summer.  Their latest and best algorithms are in place now, but they did not resurrect my favorite version (2.5).
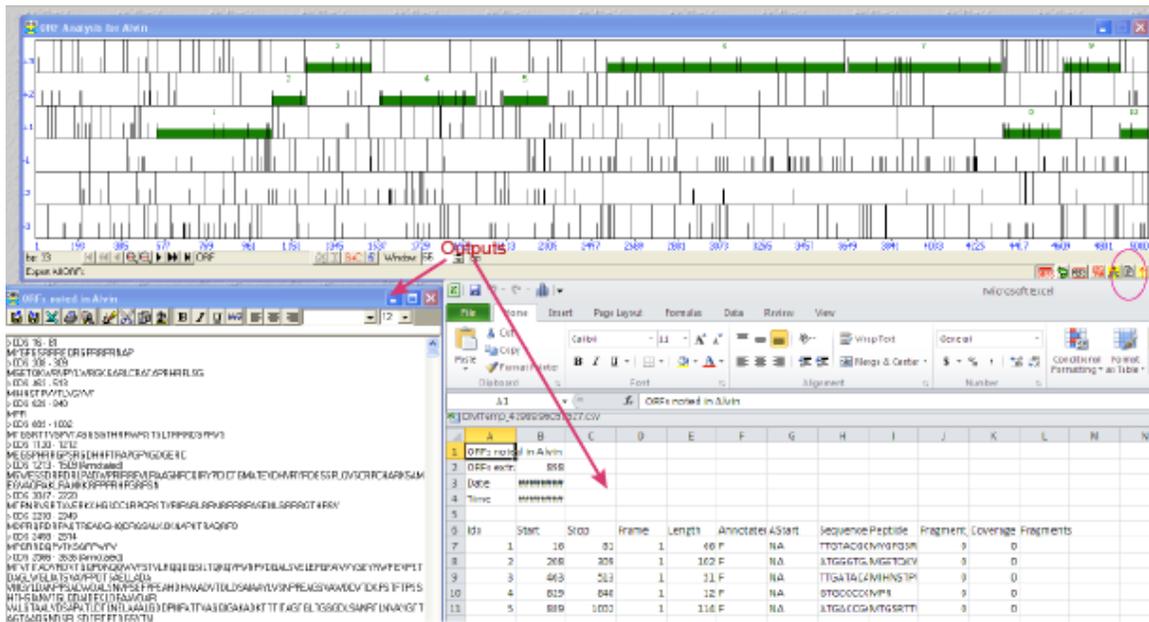
6. **DNA Master BLAST parameters**
    As long as we are on the subject of security issues on the web, NCBI changed some of their BLAST parameters too.  Because of this Dr. Lawrence changed the rates of BLAST queries and their retrievals from NCBI.  Please refer to **Section 4.5** (p. 37) for how to best BLAST an entire genome.  The changes are required to meet specifications included in NCBI's revisions to Blast batch queries.  The new settings have slowed down the queries and their retrievals.  If you exceed the expected rates of queries, your requests will not go through.

Blasting is best accomplished in Off-Peak hours (9PM – 5AM). The information you get back from an NCBI BLAST is different than a PhagesDB BLAST or Phamerator search. Do not skip one for the other. When we BLASTEd at the In Silico workshop, only 1 of 40 worked the first time (late Monday evening)!

7. **FASTA Output for ALL ORFs** (Not found in the guide…. Yet!)
If you have the opportunity to do some Mass Spec, you may want to run your data against all possible ORFs from your genome. You can use the small right hand button located in the Frames window. One button push will generate 2 outputs: 1) an Excel sheet listing all ORFs and 2) a list of FASTA sequences of all ORFs.
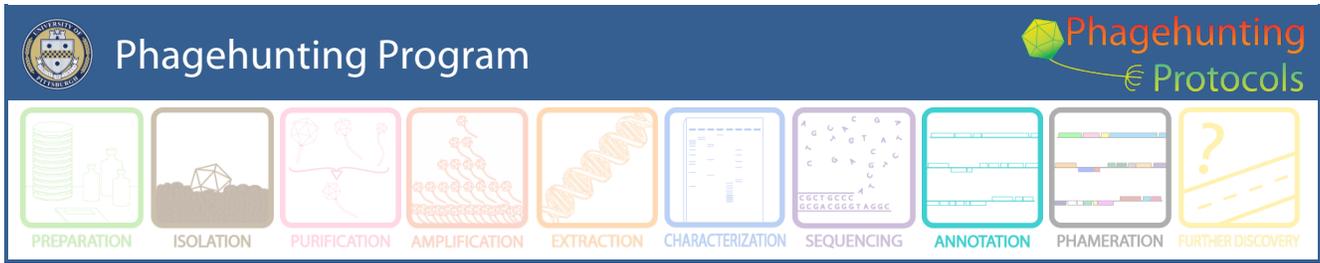


8. **Starterator (Section 6.2**, p. 59 - 61)
This is a new program written by a recent Pitt graduate, Marissa Pacey. This program allows you to compare the possible starts of all genes of the same pham. Input can include a phamerated gene, an unphamerated gene, a pham, or the whole phage genome. As you can imagine, the program will take a longer time to process the data for a whole phage genome than it will for a single gene. Starterator installation instructions are included with the other software installation instructions and on phagesDB. Use the help function of the program and the information included in **Section 6.2** (p. 59) to interpret results.

9. **Genome Announcements** are 500 word peer-reviewed articles that can be written for one or 100 genomes. We want to encourage you and your students to write a Genome Announcement for the phage(s) that you annotate this year! We have heard from the folks at Genome A and they are on board with submissions of any size. One of the benefits to you submitting your school's Genome Announcement is that you can more easily include

your students as authors.  Please know that the QC of the annotation is still a required part of the process.  Also know that we will be happy to review your announcement before submission if you want our help (and we would love to know that you are submitting it!).

Included are the "What's New" from the previous years, because they contain good recommendations!  Enjoy your bioinformatics semester!

Phagehunting Program

Phagehunting
€ Protocols

PREPARATION    ISOLATION    PURIFICATION    AMPLIFICATION    EXTRACTION    CHARACTERIZATION    SEQUENCING    ANNOTATION    PHAMERATION    FURTHER DISCOVERY

# What's New for School Year 2013-14 in Phage Genome Annotation

Created by DJS December, 2013.

**The purpose of this document is to target the most prominent changes and/or updates to DNA Master Annotation Guide.   Happy Annotating!**

## Change in final file Format (Section 12.1)

This year, the final file format requested has changed (again).  A final .dnam5 file is one that has the following properties (Figure 1).



Figure 1

1. It must be named "YourPhageName_Final.dnam5", which will help distinguish it from other versions you may have been working on.

2. **It must contain one entry and set of notes per feature**.  That means that if you have merged multiple files, you need to have evaluated the data from each source, come to a decision, and deleted erroneous or repetitive versions of each feature.  The notes for each feature should contain **everything** listed in **Section 9.6** about proper documentation of your gene calls.  You may have to delete some notes, or even rewrite some notes from scratch to meet this criterion.

3. All features must be validated (**Section 9.3.2**).

4. All features must be re-numbered if necessary **(Section 9.3.3)**.

5. All features must be re-BLASTed (**Section 9.3.4**).

6. Any functions are noted in the Notes fields, along with their source (**Section 9.3.3).** If it is determined that a particular gene has no function, include NKF (no known function) in the notes along with the sources for that determination.

**Note:** Our request for a different file format is because your final submission file is used to create a final version in Phamerator AND a GenBank file. These formats have differences. The file you submit is modified for both.

## Local BLASTp

There are lots of avenues to explore to determine the functions of the genes (blastp at NCBI, Phamerator, publications, HHPred to name a few). Did you know that phagesDB has a BLASTp function? You can BLASTp any protein of your choice against a protein database derived from Phamerator. The protein data at PhagesDB contains the most curated Hatfull lab data. The newest entries tend to have the most complete data. Go to the Home Page of PhagesDB and click the BLAST tab in the top Banner and choose BLASTp (Figure 2). Happy Blasting!



Figure 2

## Widen Feature Table

The default for the **Feature Table** includes **Name, 5' End** and **Length**.  **Right Click** on '**Name'** and choose "Widen Feature Table" (See Figure 3).  In this view, the Feature table includes **Tag, Name, 5'End, 3'End,** and **Length**.
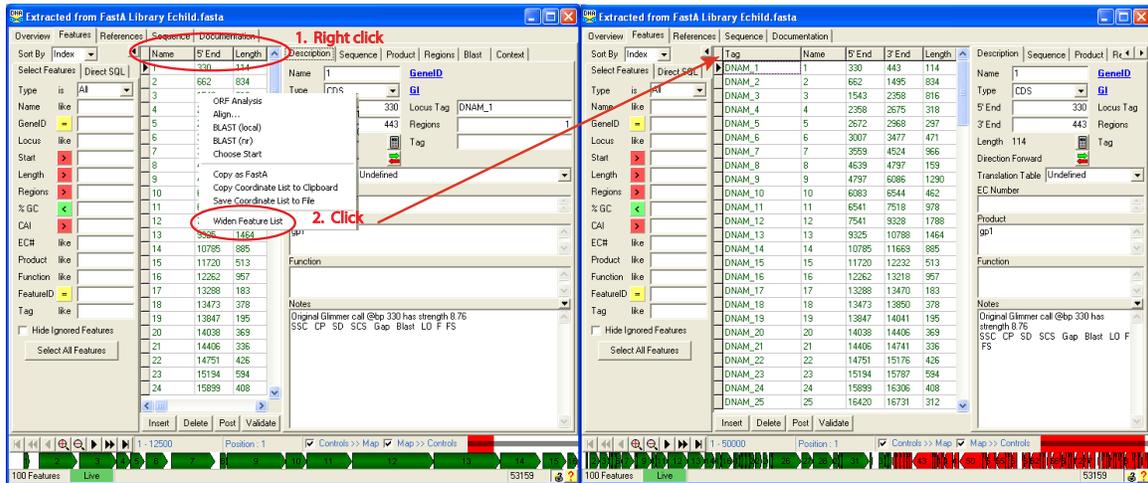


Figure 3

## Choosing Start Data

The Choose start window is dramatically new and updated!  The algorithms to score Shine Dalgarno sequences have been modified to include lots of choices.  These choices are still under evaluation.  In the meantime, you can find the "Old DNA Master" data in the SD Scoring Matrix and Spacing Weight Matrix in the upper right corner of the Choose ORF Start window.  Look for more information coming soon!
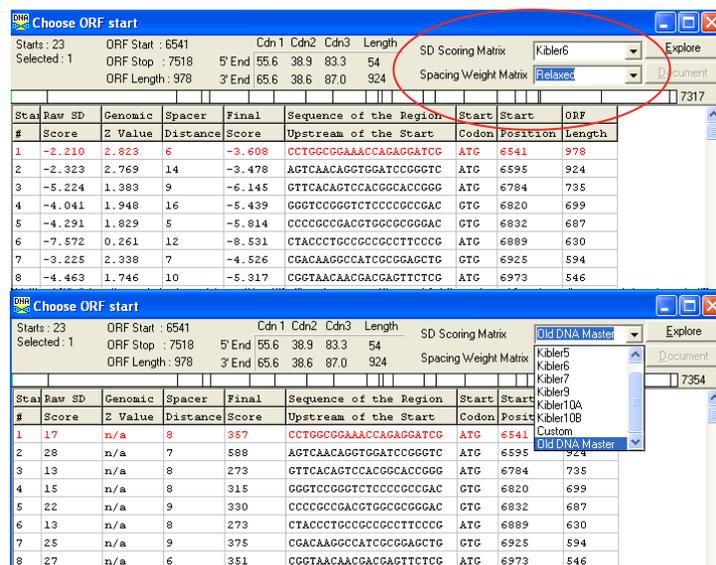


Figure 4

3

## Event Manager

Though this isn't new or different, it is often overlooked.  Go to Tools -> Event Manager.  This window gives you more information when DNA Master is not working.  Use the DNA Master Help Menu to evaluate the various events.

## DNA Master Genome Manager and Compare Genomes functions

DNA Master allows you to load phages genomes into your own local database, and then perform several advanced bioinformatic analyses through the "compare genomes" function. The simplest one may be the 'Map comparison' of your phage genome with other genomes of your choice.

1. **Load files**
   In order to compare files you will need to load files into Genome Manager. There are two kinds of files that you will want to load:
   a. **Current DNA Master files**
      To load a current file into your genome manager, with your .dnam5 file open, click "Genome→ Add to Database".
   b. **GenBank files**.

      To retrieve genomes already in GenBank, click "Tools→Genome Manager" (Figure 5).  The genome manager has a number of tabs, the left-most one labeled "Browse". This tab allows you to view all the files in your local database. The fourth tab from the left is labeled "Retrieve". From this tab, you can search NCBI for GenBank files and import them directly into your local database or open them as .dnam5 files.
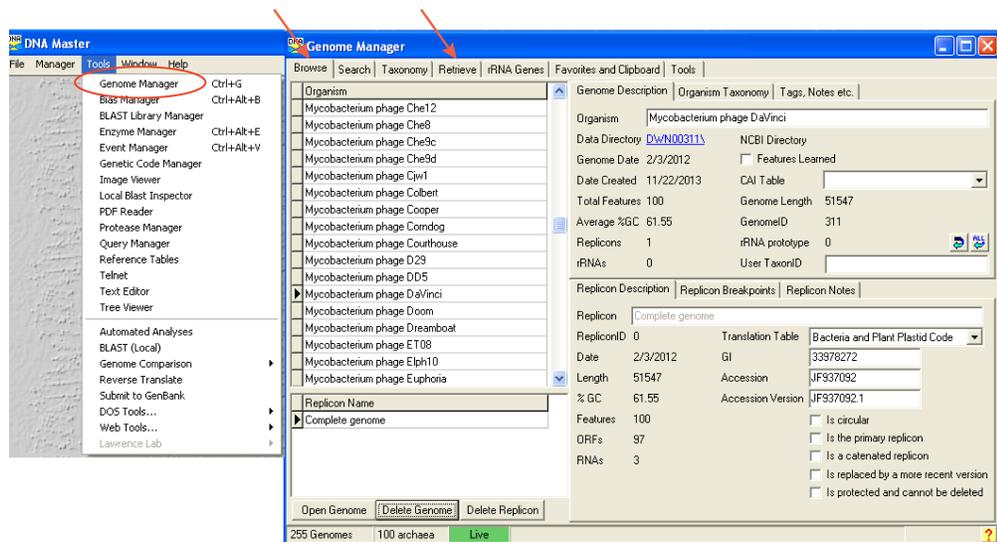


Figure 5

4

Using Figure 6 as a guide, follow this procedure to load Genome Manager with GenBank files.

A. Go to Retrieve -> Fetch By accession
B. Enter the name of the phage of interest. There are a few points about the Mycobacteriophages that you will want to keep in mind: PBI submits mycobacteriophage genomes to GenBank as Mycobacteriophage [DaVinci], when GenBank curates the file as a reference sequence, GenBank names the project Mycobacterium phage DaVinci. You made to look for both. In addition, there can be more than file for any given phage. There can be the submitted file and a the reference sequence file.. We recommend the submitted file. In addition, if the file was revised, there can be additional files.
C. When you find the files of interest, click the "keep in Mind' button.
D. That action places the Accession Number in the left box on the page.
E. Change the Fetch by Accession action to Save in local database.
F. Change windows to the Browse window.
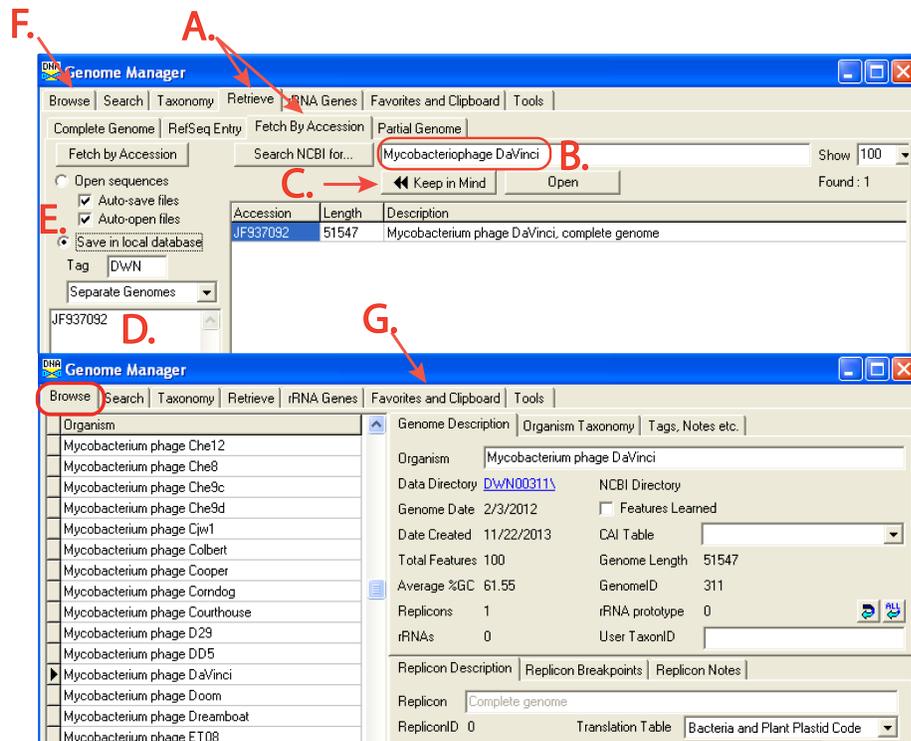G. You are ready to compare your genomes of interest. You will start by placing those genomes on a Clipboard.



Figure 6

**Note:** you can add all and any file from NCBI into your Genome Manager. Once files are available you can organize them various ways. We will use the Clipboard for this task.

## 2. Collect your genomes for comparison

To collect your genomes for comparison, select Favorites and Clipboard from the Genome Manager menu. Once there Add the genomes of interest to the Clipboard (Figure 7). Then you may close that window.
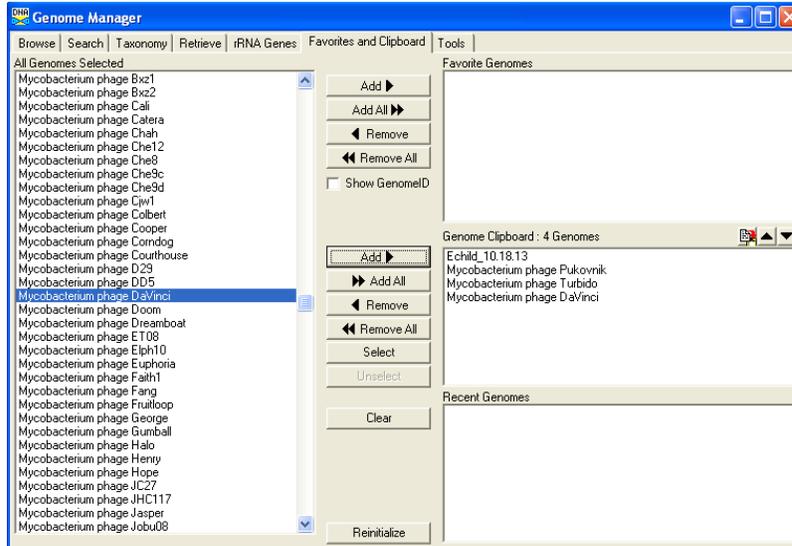


Figure 7

## 3. Compare Genomes

a. Preference Settings: To perform genome comparisons between genomes within your genome manager, you first must enter a release code into a field in the program preferences. Click "File→Preferences", and then click the far right tab labeled "Miscellaneous". At the bottom of this window, there is a field labeled "Release Code". Enter "Watson" into this field, and then click "Apply". Then click "OK", and close the window (Figure 8).
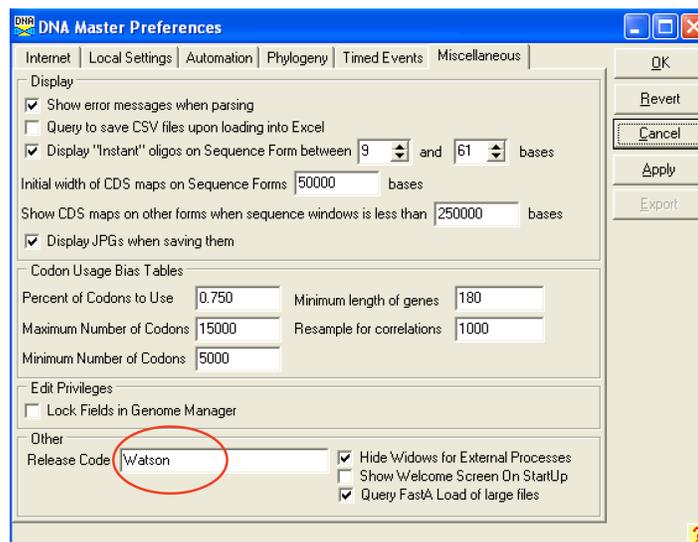


Figure 8

b. Go to Tools -> Genome Comparison-> Manual from the main menu.
c. The window that appears should look like the one in Figure 9.
d. Click Clipboard and the genomes that you selected to place on the clipboard will fill the top right field of this window.
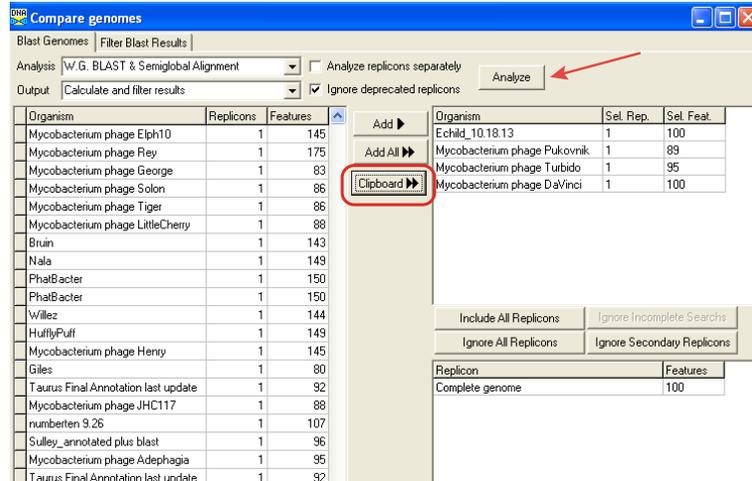


Figure 9

e. Then click Analyze.  This process will take a few minutes (It is dependent on the number of genomes that you have selected. A % completion window appears at the bottom right of the window.
f. Once complete, a new window replaces the last one with a menu as depicted in Figure 10.  Choose Map comparison from this menu. Figure 10 contains 3 separate windows.  The last one is a graphical gene comparison map, similar to the maps seen in Phamerator. Comparative analysis is based on W.G. Blast & Semiglobal Alignment.
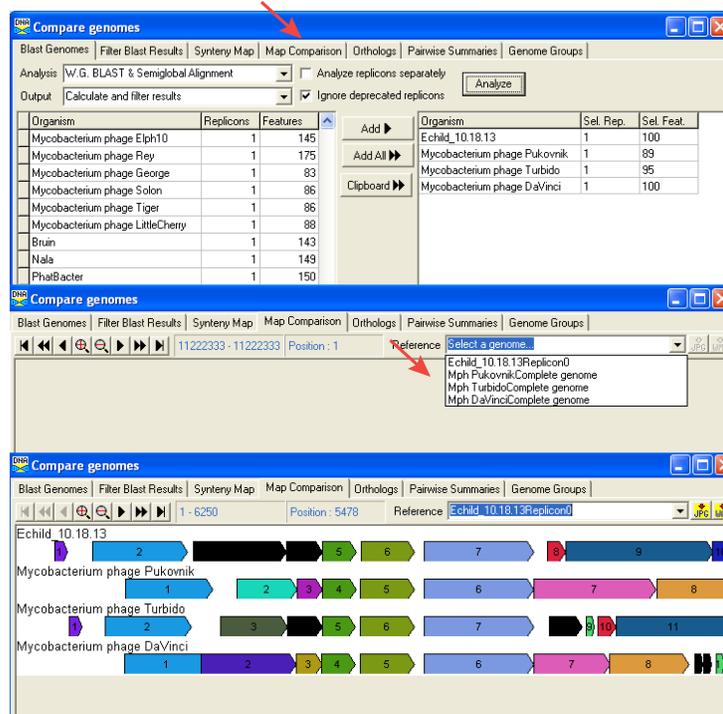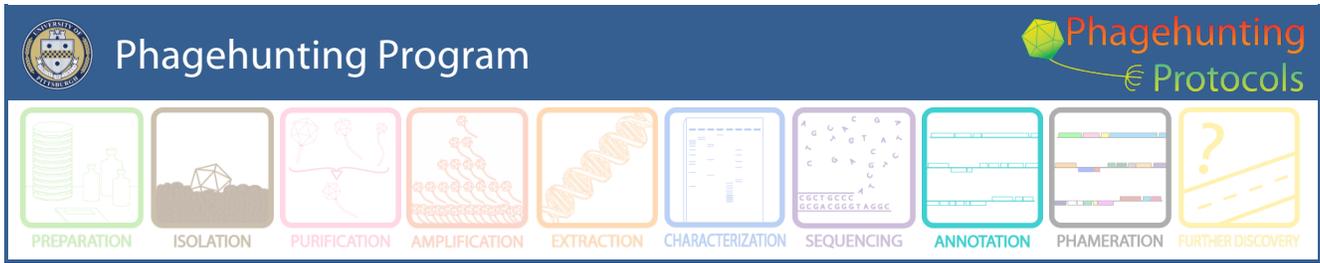
Figure 10

**Note:**  This is a quick and easy way to compare your gene calls with genomes of your choice.  It is especially helpful to be able to graphically display what is in a particular genome file.  Depending on how many genomes you've added, and how long they are, this process can take quite a bit of time (an hour). For just a few genomes, it will be relatively quick—several minutes. To read more about this, look in the DNA Master Help files for "Genome Comparison".

PREPARATION    ISOLATION    PURIFICATION    AMPLIFICATION    EXTRACTION    CHARACTERIZATION    SEQUENCING    ANNOTATION    PHAMERATION    FURTHER DISCOVERY

# What's New and Newly Recommended in the Annotation Guide

Created by djs November, 2012.  Last updated on 1.22.2013,  Revised 4.19.2103  Revised 10.22.2013 and still included because it is just good advice!

**Phamerator News:**
The user manual for Phamerator is located in Help menu.  Check it for updates!

**We recommend that you use the Notes Template:**  (This is not new, but we are still recommending it.)
This template can be added to the preferences setting and be added to the **Notes** window of all auto-annotated gene predictions
The suggested template is **SSC:  CP:  SD:  SCS:  Gap: Blast:  LO: F:  FS:**.

- **SSC:**  Start/stop coordinates.  (This may seem redundant because there are "Start" and "Stop" fields that already contain this information, but it serves as a double-check that all changes you made are actually contained in the final file.)

- **CP:**  Whether or not your start includes all the coding potential identified by GeneMark.

- **SD:**  Whether or not the start has the best SD score of all this ORF's possible starts.

- **SCS** (Start choice source):  Whether or not the gene was called by Glimmer and GeneMark, and if the start was called by same.

- **Gap** (or overlap):  Any significant gap or overlap with preceding gene (in basepairs).

- **Blast:**  The best BLAST match, and the alignment of the gene start with that BLAST match. (For example, "Matches KBG gp32, Query 1 to Subject 1", or "Aligns with Thibault gp45 q3:s45".)

- **LO** (Longest ORF):  Whether or not the coordinates you have chosen yield the longest possible gene for that ORF.

- **F** (Function):  Gene Function

- **FS** (Function source):  source for the function (see **Section 10**).  If the function assignment comes from a Hatfull-approved map in the Appendix, please also enter it into the field labeled "Function" directly above the "Notes" field.  Otherwise, only enter the putative functional assignment in the Notes.

- Anything else you think is important.  In particular if you made a different choice than previous annotators have made in published genomes, and feel very strongly about your choice, this is the place to let us know.   **Example:**  If your gene start does

not match the published starts of similar genes in GenBank, an explanation of why not.  ("Published Thibault gp45 start not present in my sequence" or "Thibault start caused a 200 bp overlap with upstream gene")

Caution:  Do not add hard returns in the template of the **Notes** window because it takes up too much space and is a formatting problem.

**GeneMark options:**
1.  It can be difficult for students to understand the data imported into DNA Master has the same value as the GeneMark *M.Tb* or *M. smegmatis* graphical data. You can get the same data (in a somewhat different graphic output – portrait vs. landscape representation) at http://exon.gatech.edu/heuristic_hmm2.cgi  where you can run the GeneMark heuristic model on the web.

2.  There is an updated version of GeneMark (2.8) for bacterial models.  It can be found at http://exon.gatech.edu/gmhmm2_prok.cgi.  This is relatively new and untested.  Preliminary data suggests it will evaluate the genomes identically to GeneMark 2.5.

3.  When using the web-based GeneMark against a model organism, the model organism you pick is based on the assumption that the phage and host have a relationship.  That relationship may not be helpful when calling genes.  I recommend running both the *M. smegmatis* and *M.Tb* and compare.  One may provide better information than the other.  If you do this, please send us a note with your evaluation.

**Provisional Cluster Assignments:**
No cluster assignment should be made on phagesdb.org before the genome sequence has been completed and evaluated.  If you assign a cluster based on other data, it will be considered a provisional assignment.

**Sequence Changes:**
If you retrieve your sequence from phagesdb.org, your phage sequence will be 'finished' and in the correct orientation.  If you are working on additional sequences or find that the sequence does need changed, DNA Master can make those changes without losing your database information for that genome. (Which means you won't lose the work you have done to this point.)  You can reverse-complement the genome or add and subtract bases.  However, changes made on the screen do not become part of the binary database file UNTIL you post them to that file.  Posting in the sequencing window is accomplished by clicking on the **Raw** button.