# RepeatAnalyzer Quickstart Guide

## Overview

The purpose of this guide is to allow you to set up and begin using the RepeatAnalyzer software package on a new computer, with either Windows, Mac OS X or Linux operating systems. This process involves installing free software from a third party source and thus requires an internet connection. A brief description of how to use the software follows, however this is not meant to be a complete source on use cases. For questions related to problems with installation or usage of this software, contact Helen Catanese at helen.catanese@wsu.edu.

## Download

*Download the source code and associated data files here.*

### Direct Download (Recommended)

On the left hand side of the page, click "Downloads." Then under the Downloads tab, select "Download repository." Be sure to unzip the directory to a convenient location.

### Using Git

If you are familiar with Git, you can also clone this repository by running > git clone https://bitbucket.org/repeatgroup/repeatanalyzer/git

## Installation

Before beginning installation, be sure you have an internet connection.

1) Go to www.continuum.io/downloads. This is the download page for the Anaconda Python distribution. Note: while Anaconda is not strictly required to run RepeatAnalyzer, it comes prepackaged with most of the libraries that RepeatAnalyzer needs, and reduces the difficulty of installing those dependencies considerably.

2) The page includes installation instructions for Windows, Mac OS X and Linux. Be sure to follow the installation instructions for your specific platform, and to choose Python version 2.7.

   a) If you are a Mac OS X user, we recommend using the graphical installer, rather than the command line installer.

3) Once the installer has downloaded, run it using the instructions provided on the page below the download link.

   a) If you are using a graphical installer (Mac OS X or Windows), select all the default options. (Confirm that both boxes under "advanced option" are checked.)

   b) If you are using a command line installer (Mac OS X or Linux) after running the installer as

directed on the page, be sure to say 'yes' when prompted to prepend Anaconda to your PATH.

4) If you are using Linux or Mac OS X, open (or restart) the terminal. If you are using Windows, once you have installed Anaconda, you will have a program called Anaconda Command Prompt, or Anaconda Prompt installed on your machine. Open that.

5) We will use the command prompt to install several additional libraries the program needs. To interact with the command prompt, simply type in a command and hit the 'enter' key. Henceforth, commands to be submitted in this way will be denoted as "> command here"

   a) If you are a Mac OS X or Linux user and at any point you see the error "conda: command not found" or "pip: command not found", it means you did not add Anaconda to your execution path. To do this add it open the \etc\paths file and add the following line [Path to Anaconda install]/anaconda/bin. To do this, you will need administrator permission.

6) In the Anaconda command prompt, enter the following commands (these are cases sensitive). If an y/n options appear, enter > y

> conda install numpy

> conda install basemap

> pip install pyshp

> pip install geopy

7) Be sure that you have downloaded the source code directory from the link above, and if necessary, unzipped the file. The files in this folder must remain together in order to function properly. Do not remove any of them from this folder (or move them to subfolders).

8) Once this is completed, if you are using Windows simply double-click the RepeatAnalyzer.py file in the source code folder to open the program. If you are using a Mac or Linux, navigate into the source folder in the terminal and run the command "python Repeatanayzer.py" to start the program. You may create a shortcut to this file to run the program from a more convenient location.

   a) If you are prompted to choose a program to open the file with, navigate to the Anaconda2 installation directory, and select the python.exe (NOT pythonw.exe) file.

      i) In rare cases on Windows the python.exe file will not be selectable as a default program. To correct this, run the following commands in the windows command prompt (cmd.exe). Be sure to run it as an administrator.

         > assoc .py=Python.File

         > ftype Python.File="[Anaconda installation path]\python.exe" "%1" %*

   b) If the window opens, but closes immediately, there is likely a problem with your installation, or a conflict with other software installed on your system (like other versions of Python). In this case, see troubleshooting.

9) When the program is run for the first time, it will warn you that the RepeatAnalyzer.dat file is missing, and ask if you wish to continue. This is normal behavior. To setup the program with the default *A. marginale* data included with the source code, just enter 'y' and then 'Anaplasma marginale' as the species name. See 'Using the Main Functions' subsection 5 for details on loading in the data for the first time.

# Basics

RepeatAnalyzer uses a primarily command line interface, which means that interacting with the main menu of the program will usually involve reading options from the screen and typing in your selection. The main menu includes 11 options (currently) which can be selected by entering the corresponding number for that option. Once an option has been selected, RepeatAnalyzer may open additional dialogs, or present other command line options.

At any given time, RepeatAnalyzer will be set to work on a single species (*Anaplama Marginale, Anaplasma Centrale,* etc). Be sure that the program is set to the intended species when you are adding or searching data, as it will only access the repeats and strains stored for that particular species. If the species you wish to work with does not appear as on option in command 2 (change current species), you may add it with command 2, sub-option 0. Be sure to check the current species each time you open the program (if you are storing data for more than one species).

# Adding Data

## With existing data

If the RepeatAnalyzer.dat file is in the same folder as RepeatAnalyzer.py when it is executed, the stored data will be loaded. Upon opening the program you will see the main menu, showing the current species, and the number of strains and repeats stored respectively. If this is the correct species, continue as usual, otherwise, be sure to change the current species before moving on.

## Without existing data

If the RepeatAnalyzer.dat file is not in the same folder as the RepeatAnalyzer.py file when it is run, the program will first confirm that the RepeatAnalyzer.dat file is missing and ask if you wish to continue. If you enter 'y' (for yes) it will ask what species you would like to work with. At this point you will be prompted to enter the name (Genus species) of the species you will be adding data for. After you have done this the main menu of the program will appear, with the species name you entered and 0 strains, 0 repeats.

## In either case...

At the main menu, select command 5 and enter the name of the input file. AmarginaleData.txt comes packaged with the source code. If you wish to use another input file, make sure it follows the formatting outlined in the Sample File Formatting Section below. The program will then read this file into its internal storage. If there is an issue reading, the program will try to give the line number where

the issue occurred in the file, but the error may actually be on the line before or after the one given, so check that whole area and correct any formatting mistakes until it reads in with no errors. Once reading is complete, you will be prompted on whether you would like to update the geocoding for your data. If you have an internet connection, select 'y', as mapping and regional analysis functionalities will be disabled until this is completed. If you select 'n' now, geocoding can be updated later via command 11.

Keep in mind that <u>any strains or repeats with the same sequence will be stored as a single repeat or strain with multiple names</u>, so if the number of repeats is lower than you were expecting after adding a file, this is likely the reason. Check the data summary (command 10) for details on where this occurred.

Note that within the included Anaplasma marginale data file, the repeat 'tc63_3_s06' is present with no given sequence. This is not a mistake, as it was included in a publication with no sequence provided and no valid reference. It will not cause any problems in the program, only an error message when inputting the file; genotypes including it will be ignored. Keep in mind that this is the 235th repeat, as referenced in the main paper.

## Sample Input File Format

Paper:

    [citation line 1: title]

        [citation line 2:autors]

        [citation line 3: other info]

    Repeats:

        [name] : [sequence]

        {repeat for each repeat}

    Strains:

        [name (optional, but the colon is not) OR year , animalID (if using auto-naming)] :

        [space separated list of repeat names] : [location in the form country,province,county

        or country province or country]

        {repeat this for each strain. if you have one strain at multiple locations it will need to be

        listed twice, but rest assured it will only be counted once}

{Repeat this pattern for each published paper}


Unpublished:

    [author name(s), this should look like 2 of paper citation]

        [year of findings]

    Repeats:

        {this should be exactly the same as the repeats section for a paper}

    Strains:

        {likewise, this is the same as the strains section for a paper}

{Repeat this pattern for each unpublished source}

# Using the Main Functions

Each of the following subsections outlines the capabilities of one of the 11 main menu functions which RepeatAnalyzer performs. Each section includes details on any input required for the function, and any output it generates. A user can access the function by entering the associated number at the main menu.

## 1. Identify repeats

**Input:** One full or partial gene sequence either in DNA or protein form

> OR multiple gene sequences in FASTA format.

**Output:** The names of the maximal set of repeats in each sample with genotype name and publication details if any.



*Figure 1 Shows the repeat/genotype identification interface. Once a sequence or sequences have been entered into the identification window (bottom left) and the correct input type is chosen (DNA or Protein), the match window will appear with which (known) repeats occur in the sequence and all relevant information on the genotype (if it has been reported previously). All windows can be resized and/or have their contents copied as needed.*

## 2. Change current species

**Input:** The number of the species (as listed on-screen)

OR 0 followed by name of the new species.

**Output:** The main menu header will change as appropriate.



*Figure 2 Shows the interface to change species or add a new species. Specifically, showing the series of commands used to add a new species to RepeatAnalyzer. The first command opens the species change menu. The second selects add a species. Finally, the third accepts the name of the species from the user. Before the species is added, RepeatAnalyzer will also ask the user to enter the command 'y' to confirm that the new name is correct.*

## 3. Search data

**Input:** Select entity type by number, repeat, strain or location. Check boxes as desired.

*Repeat*: repeat name or sequence

OR multiple of repeat names or sequences in FASTA format

*Strain*: genotype name or sequence (by repeat names)

OR multiple genotype names or sequences in FASTA format

*Location*: location name from dropdown

**Output:** A summary of information for the entity or entities selected.

See Figures 3, 4 and 5 for illustrations of each of the three search types.

*Figure 3 Shows the search by repeat interface. Once one or more repeat names or sequences (in FASTA format) is entered along with a maximum edit distance (default zero), the search result window appears with the relevant details for that repeat. The windows are scrollable and can be copied.*



*Figure 4 Shows the Search by location interface. Similar to Figure 3, search by location shows all relevant information for a given location, selected from the alphabetized list in the search window. Locations can be as broad as a country, or as narrow as a county, and broader location will include results from all narrower locations as well. (i.e. the result for Brazil includes results for any Brazilian province or county with data.) If a location does not appear on the list, then there are no genotype reports in that location.*

*Figure 5 Shows the Search by Strain interface. It can accept either a single strain name or sequence, or multiple in FASTA. For each input it returns the edit distances for all repeats in the strain, all locations it has been reported, and the associated publications.*

# 4. Map data

**Input:** Any combination of: A list of repeats by name or sequence, a list of genotypes by name or sequence, and a location. Lists are separated by semicolons. Repeats in a genotype are separated by whitespace.  Lists may be replaced by the word 'All' to include all entries, or left blank to indicate none. There is also a dropdown menu to select where the legend is placed and marker scale.

**Output:** A printout of data on the mapped entities and a world map showing where those entities occur. The map can be zoomed and it can be saved.



*Figure 6 Shows a sample mapping query and result. Each map query can be modified by a number of*

*options including region (results include only repeats/genotypes present in that region, but all locations they occur in the world), marker scale, legend location and the specific repeats/genotypes of interest. The map plot can be scaled, zoomed, panned and saved to a separate file using the options on the top left of the panel.*



*Figure 7 Shows a sample mapping query. When sequence is set to 'All' the map will show all results for the specified region. There is also an option to ignore items that appear at only a single location.*

*Figure 8 Shows a sample mapping result for the query in figure 7.*

# 5. Input data from file

**Input:** The name of a file formatted as described in the Sample Input File Format section.

**Output:** Any error in the input file will be noted. The main menu header will change as appropriate.



*Figure 9 Shows the sequence of commands to input new data for a species. The first command opens the data input menu. While the second allows the user to enter the name of the file where the data is stored. Note that if there are any errors which prevent the file from reading, the line number of the error will be shown. If any repeats are listed in strains, but never specified in the repeats section, the error shown above the third command will print, and that strain will be omitted, but the rest of the file will be read in correctly. Finally, the third command may either be entered 'y' if a stable internet connection is available, or anything else, if not. This step is required for certain program features to work correctly, and can be done later if it is skipped now (function 11).*

# 6. Regional diversity analysis

**Input:** Region, selected from a dropdown menu and checkboxes for additional plots at needed.

**Output:** A printout with all diversity scores, repeat frequencies, unique repeats and (if selected) appropriate plots.



*Figure 10 Shows the regional diversity analysis interface and a sample query. For regional diversity analysis, you can select the region of interest as well as whether or not to produce various plots (shown in Figure 11). The result has several sections including, how many genotypes each repeat in the region occurs in, a list of all repeats that occur nowhere outside the region, and a list of the various diversity scores as defined in Table 1.*

| Metric Name | Formula |
|---|---|
| GD2 | $\dfrac{Total\ \#\ Unique\ SSRs}{\#\ Genotypes}$ |
| GDM1 (Local) | $Avg\left(\dfrac{\#\ Unique\ SSRs\ in\ Genotype\ i}{Length(Genotype\ i)}\right)$ |
| GDM1 (Global) | $\dfrac{Total\ \#\ Unique\ SSRs}{Total\ Length(All\ Genotypes)}$ |
| GDM2 (Local) | $Avg\left(Deviation\left(\dfrac{Frequency(SSR\ i\ in\ Genotype\ j)}{Length(Genotype\ j)}\right)\right)$ |
| GDM2 (Global) | $Deviation\left(\dfrac{Frequency(SSR\ i)}{Total\ Length(All\ Genotypes)}\right)$ |

*Table 1. Metrics used to calculate genetic diversity of a geographic region.*

Length(Genotype) = the number of SSRs in in that genotype

*Figure 11 Shows the plots produced by the regional diversity analysis in Figure 10.*

# 7. Remove a species

**Input:** The ID of the species to be removed (as shown on-screen)

**Output:** None



*Figure 12 Shows the series of commands to remove a species. The first command opens the species removal menu, while the second selects the number of the species to remove.*

# 8. Remove a strain

**Input:** The repeat sequence of the strain to remove.

**Output:** The main menu header will change as appropriate.



*Figure 13 Shows the series of commands to remove a strain. The first command opens the strain removal menu. The second command allows the user to enter the repeat sequence of the strain to be removed. Finally, the third command may either be entered 'y' if the correct strain was selected.*

# 9. Generate strain names

**Input:** None

**Output:** None

Note: This function will generate any names for newly input data where the genotype name was listed as year, animalID. There is no direct input or output.

While typically this step is done when menu option 5 is run, it can only be completed after geocoding (which requires an internet connection) and so, it may optionally be skipped and run later, via this command.
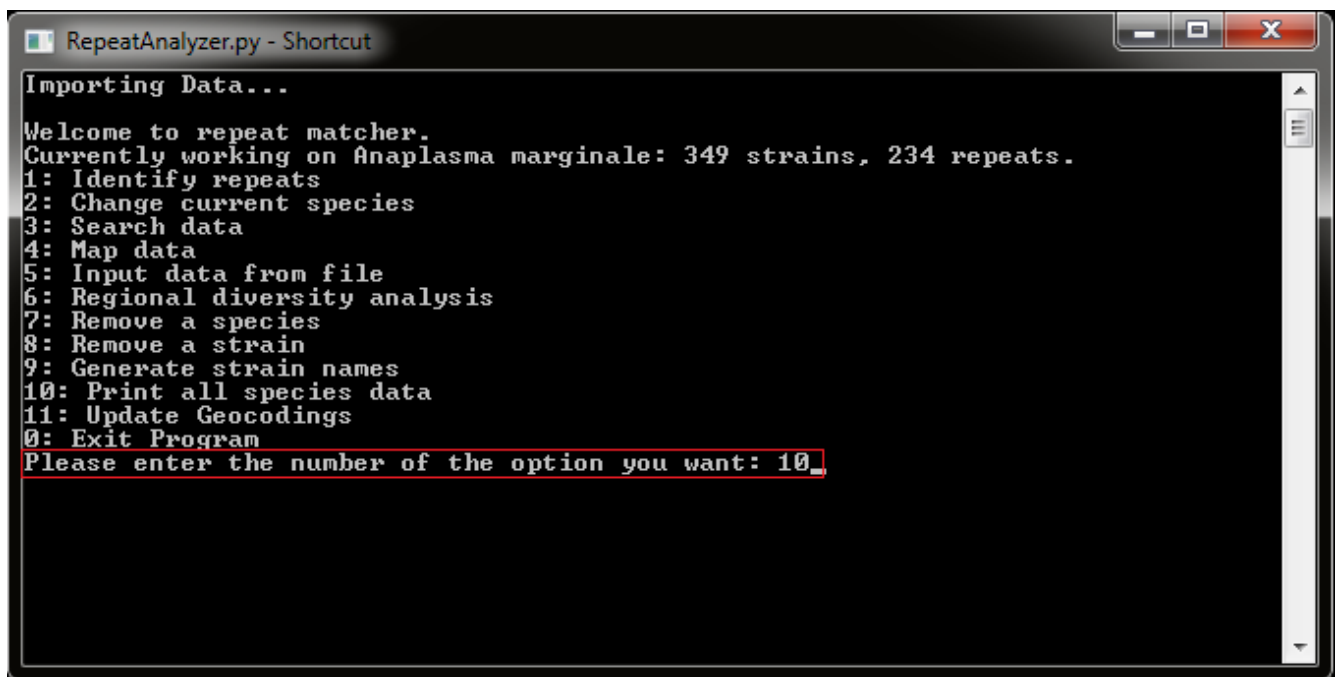
15

*Figure 14 Shows the command to automatically generate species name for species input with the correct information.*

## 10. Print species data

**Input:** None

**Output:** A text printout with a summary of all data, and a spreadsheet with all repeats and their sequences for easy manipulation. These will be located in the same folder as RepeatAnalyzer.py.
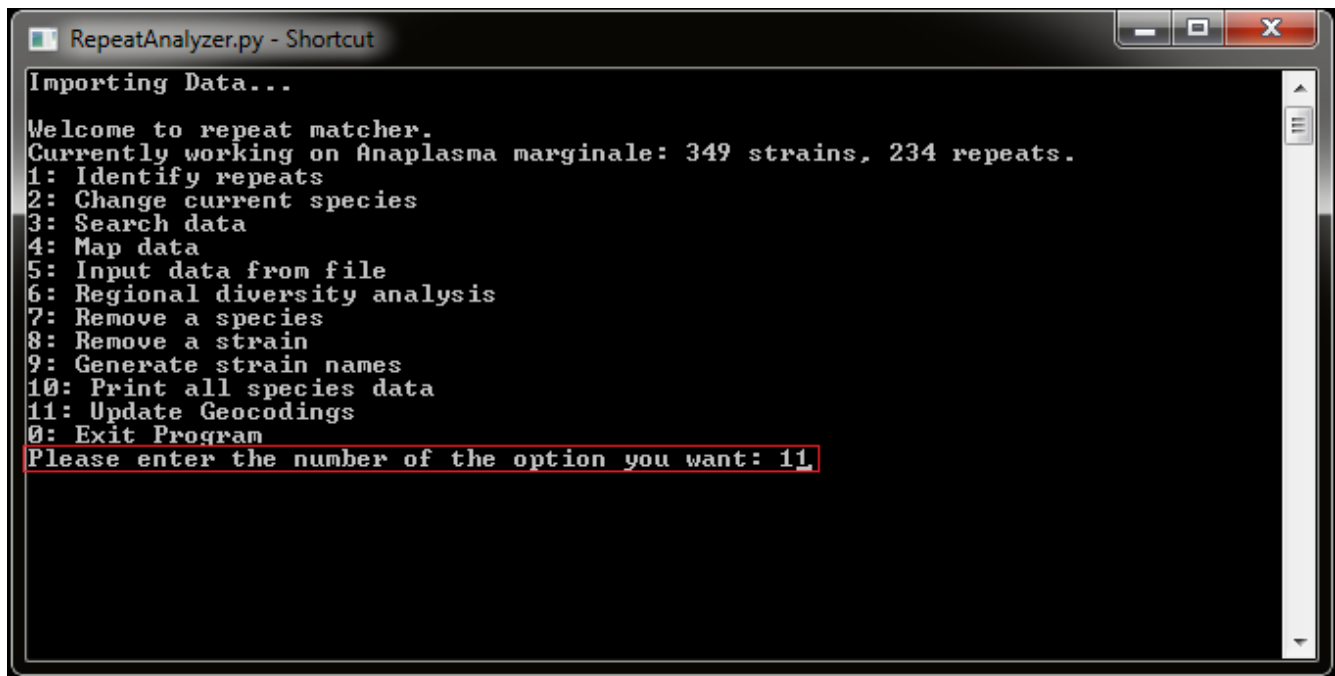


*Figure 15 Shows the command to print a summary of all species data to an external file.*

## 11. Update geocodings

**Input:** None

**Output:** Formats the names of all locations to be uniform (so that USA, U. S. A., and United States all point to the same place) and gets the coordinates for locations that have none. This must be done at least once after adding data for mapping to function properly. Note: there is a y/n option to do this after adding new data as well. <u>An internet connection is required.</u>



*Figure 16 Shows the command to update geographic codings.*

# Troubleshooting

If you have experienced problems running RepeatAnalyzer there is likely a problem with your installation, or a conflict with another program (such as an older version of Python) installed on your system. We can help resolve certain issues, but this will be considerably faster with the error messages specific to your setup. If possible, send any troubleshooting requests along with the error log acquired in the following way.

To give us more information about what your specific problem is, open a program on your computer called "spyder" (part of the Anaconda environment) and in that program open the RepeatAnalyzer.py file. Hit the 'f5' key and confirm the dialog that appears. An error message should appear on the lower right-hand side of the screen. Copy all of the text there (it may be a lot, so be sure to scroll up if needed) and send it to us.

Email information about configuration problems/conflicts to <u>helen.catanese@wsu.edu</u> and we will attempt to resolve the issue.