# A Practical Guide to Phylogenetics for Nonexperts

Damien O'Halloran[1]

[1]Department of Biological Sciences and Institute for Neuroscience, The George Washington University

Correspondence to: Damien O'Halloran at damienoh@gwu.edu

## Abstract

Many researchers, across incredibly diverse foci, are applying phylogenetics to their research question(s). However, many researchers are new to this topic and so it presents inherent problems. Here we compile a practical introduction to phylogenetics for nonexperts. We outline in a step-by-step manner, a pipeline for generating reliable phylogenies from gene sequence datasets. We begin with a user-guide for similarity search tools via online interfaces as well as local executables. Next, we explore programs for generating multiple sequence alignments followed by protocols for using software to determine best-fit models of evolution. We then outline protocols for reconstructing phylogenetic relationships via maximum likelihood and Bayesian criteria and finally describe tools for visualizing phylogenetic trees. While this is not by any means an exhaustive description of phylogenetic approaches, it does provide the reader with practical starting information on key software applications commonly utilized by phylogeneticists. The vision for this article would be that it could serve as a practical training tool for researchers embarking on phylogenetic studies and also serve as an educational resource that could be incorporated into a classroom or teaching-lab.

## Video Link

The video component of this article can be found at https://www.jove.com/video/50975/

## Introduction

In order to understand how two (or more) species evolved, it is first necessary to obtain sequence or morphological data from each sample; these data represent quantities that we can use to measure their relationship through evolutionary space. Just like when measuring linear distance, having more data available (*e.g.* miles, inches, microns) will equate to a more accurate measurement. Ergo, the accuracy with which a researcher can deduce evolutionary distance is heavily influenced by the volume of informative data available to measure relationships. Furthermore, because different samples evolve at different rates and by different mechanisms, the method that we use to measure the relationship between two taxa also directly influences the accuracy of evolutionary measurements. Therefore, because evolutionary relationships are not directly observed but instead are extrapolated from sequence or morphological data, the problem of inferring evolutionary relationships becomes one of statistics. Phylogenetics is the branch of biology concerned with applying statistical models to patterns of evolution in order to optimally reconstruct the evolutionary history between taxa. This reconstruction between taxa is referred to as the taxa's *phylogeny*.

To help bridge the gap in expertise between molecular biologists and evolutionary biologists we describe here a step by step pipeline for inferring phylogenies from a set of sequences. Firstly, we detail the steps involved in database interrogation using the Basic Local Alignment Search Tool (BLAST[1]) algorithm through the web based interface and also by using local executables; this is often the first step in obtaining a list of similar sequences to an unidentified query, although some researchers may also be interested in gathering data for a single group via web interfaces such as Phylota (http://www.phylota.net/). BLAST is an algorithm for comparing primary amino acid or nucleotide sequence data against a database of sequences to search for "hits" that resemble the query sequence. The BLAST program was designed by Stephen Altschul *et al*. at the National Institutes of Health (NIH)[1]. The BLAST server consists of a number of different programs, and here is a list of some of the most common BLAST programs:

i) *Nucleotide-nucleotide BLAST (blastn):* This program requires a DNA sequence input and returns the most similar DNA sequences from the DNA database that the user specifies (*e.g.* for a specific organism).

ii) *Protein-protein BLAST (blastp):* Here the user inputs a protein sequence and the program returns the most similar protein sequences from the protein database that the user specifies.

iii) *Position-Specific Iterative BLAST (PSI-BLAST) (blastpgp):* The user input is a protein sequence which returns a set of closely related proteins, and from this dataset a conserved profile is generated. Next a new query is generated using only these conserved "motifs" which is used to interrogate a protein database and this returns a larger group of proteins from which a new set of conserved "motifs" are extracted and then used to interrogate a protein database until an even larger set of proteins are retuned and another profile is generated and the process repeated. By including related proteins into the query in each step this program allows the user to identify sequences that are more divergent.

February 2014 |  84  | e50975 | Page 1 of 13

iv) *Nucleotide 6-frame translation-protein (blastx):* Here the user provides a nucleotide sequence input which is converted into the six-frame conceptual translation products (*i.e.* both strands) against a protein sequence database.

v) *Nucleotide 6-frame translation-nucleotide 6-frame translation (tblastx):* This program takes a DNA nucleotide sequence input and translates the input into all six-frame conceptual translation products which it compares against the six-frame translations of a nucleotide sequence database.

vi) *Protein-nucleotide 6-frame translation (tblastn):* This program uses a protein sequence input to compare against all six reading frames of a nucleotide sequence database.

Next, we describe commonly used programs for generating a Multiple Sequence Alignment (MSA) from a sequence dataset, and this is followed by a user guide to programs that determine the best-fit models of evolution for a sequence dataset. Phylogenetic reconstruction is a statistical problem, and because of this, phylogenetic methods need to incorporate a statistical framework. This statistical framework becomes an evolutionary model that incorporates sequence change within the dataset. This evolutionary model is comprised of a set of assumptions about the process of nucleotide or amino-acid substitutions, and the best-fit model for a particular dataset can be selected through statistical testing. The fit to the data of different models can be compared via likelihood ratio tests (LRTs) or information criteria to select the best-fit model within a set of possible ones. Two common information criteria are the Akaike information criterion (AIC)[2] and the Bayesian information criterion (BIC)[3]. Once an optimal alignment is generated, there are many different methods to create a phylogeny from the aligned data. There are numerous methods of inferring evolutionary relationships; broadly, they can be divided into two categories: distance-based methods and sequence-based methods. Distance-based methods compute pairwise distances from sequences, and then use these distances to obtain the tree. Sequence-based methods use the sequence alignment directly, and usually search the tree space using an optimality criterion. We outline two sequence-based methods for reconstructing phylogenetic relationships: these are PhyML[4] which implements the maximum likelihood framework, and MrBayes[5] which uses Bayesian Markov Chain Monte Carlo inference. Likelihood and Bayesian methods provide a statistical framework for phylogenetic reconstruction. By providing user information on commonly used tree-building tools, we introduce the reader to the necessary data required to infer phylogenetic relationships.

## Protocol

## 1. Basic Local Alignment Search Tool (BLAST): Online Interface

1. Click on this link to visit the BLAST[1] web server at the National Center for Biotechnology Information (NCBI). - http://blast.ncbi.nlm.nih.gov/Blast.cgi (**Figure 1**).
2. Input a FASTA formatted text sequence (see **Figure 2** for example) into the query box.
3. Click the appropriate BLAST program and relevant database or individual species of interest to use in the search and then click "BLAST". *Note:* FASTA formatted sequence begins with a description line indicated by a ">" sign. The description must follow immediately after the ">" sign, the sequence (*i.e.* nucleotides or amino acids) follow the description on the next line. The output from the BLAST search is viewed as HTML, plain text, XML, or hit tables (Text or csv) with the default set to HTML (**Figure 3**).

## 2. Basic Local Alignment Search Tool (BLAST): Local Executables

1. Download the latest BLAST command-line BLAST executables from this link:
   ftp://ftp.ncbi.nlm.nih.gov/blast/executables/blast+/LATEST/ -
2. em>For PC users: double-click the latest blast win32.exe file and accept the license agreement and click install.
   Note: The default installation directory is C:\ncbi-blast-2.2.27+.
3. Configure the PC environment variable as follows:
   1. Click the PC "start" button, and then right click "computer",
   2. Click "Properties" and in the pop-up click on the "advanced" tab
   3. Click the "Environment Variables button" and in the new pop-up click the "new" button under the "User variables for user" section
   4. In the pop-up add the variable name "Path" and variable value "C:\ncbi-blast-2.2.27+\bin.
      *Note:* the bin directory contains the executable (*i.e.* blastp, *etc.*).
4. em>For Mac users: Open the Terminal application (to do this just open "Finder" and search "Terminal" and this will display the "terminal" icon). Into the terminal window type:
   >ftp ftp.ncbi.nih.gov
   *Note: can also type the URL used above in the example for PC*
5. To access the NCBI ftp site type "anonymous" for Name and Password, and then type:
   >cd blast/executables/LATEST
6. List the executables by typing:
   >ls
7. Get the latest version by typing the following (or whatever the latest version currently is):
   >get ncbi-blast-2.2.7-macosx.tar.gz
8. Exit the NCBI ftp server site by typing "exit".
9. Decompress the downloaded files by typing:
   >tar -xzf ncbi-blast-2.2.7-macosx.tar.gz
10. Add the location of the binaries for the blast executable to your path so that the shell can search through this directory when looking for commands by typing:
    >PATH=$PATH:new_folder_location
11. Check if this added the location to your path by typing:

>echo $PATH

12. Download a preformatted BLAST databases (which are updated daily) by clicking here:
ftp://ftp.ncbi.nlm.nih.gov/blast/db/
13. Place the database into the "db" folder.
14. em>On a PC: open a MS-DOS prompt (to do this click "start" and type "cmd" in the search bar) and change the directory to the ncbi-blast folder by typing:
C:\Users>cd ..\ [moves up one folder]
C:\>cd ncbi-blast-2.2.27+
This will change the directory to:
C:\ncbi-blast-2.2.27+>
15. Create the database using the following "makedb" command:
>makedb –in db/briggsae.fasta –dbtype prot –out db/briggsae
*Note:* In the example below (**Figure 4**) the database is named "briggsae" and is comprised of one linkage group from the organism *Caenorhabditis briggsae.*
16. Create a query protein sequence called "test" by inserting a FASTA formatted protein text sequence into the "db" folder.
17. Interrogate the database via a blastp search by typing the following command:
>blastp –query db/test.txt –db db/briggsae –out text.txt
18. em>On a Mac: download a database for local Blast searches by accessing the NCBI ftp website as per the instructions above (step 2.4) and then type:
>lcd ../databases/
19. Download the genome or sequence of interest by typing:
>get NC_[Accession #].fna
*Note:* ".fna" refers to the FASTA formatted nucleotide sequence and ".faa" refers to the FASTA formatted amino acid sequences.
20. Type "quit" to exit the ftp site.
21. Make the database by typing:
>makeblastdb -in db/mouse.faa -out mouse -dbtype prot
22. Insert a FAST formatted query sequence into the "bin" folder and interrogate the database with the following command:
> blastp -query "your query.fasta" -db "your database" -out results.txt

## 3. Generating Multiple Sequence Alignments

1. Click on these links to access commonly used Multiple Sequence Alignment (MSA) programs:
ClustalW[6] http://www.clustal.org/
Kalign[7] http://msa.sbc.su.se/cgi-bin/msa.cgi
MAFFT[8,9] http://mafft.cbrc.jp/alignment/software/
MUSCLE[10] http://www.drive5.com/muscle/
T-Coffee[11] http://www.tcoffee.org/Projects/tcoffee/
PROBCONS[12] http://toolkit.tuebingen.mpg.de/probcons
2. Click on this link - http://tcoffee.crg.cat/apps/tcoffee/do:regular - and input FASTA formatted sequence data into the query box
*Note:* A sample output from T-Coffee can be seen in **Figure 5**, similar residues are color coded.
3. Download the Clustal MSA as a command line version (ClustalW) or a graphical version (ClustalX) by clicking this link: http://www.clustal.org/clustal2/ - then click on the appropriate executable (*i.e.* win, Linux, Mac OS X).
4. Upload data as FASTA formatted sequence text and align (**Figure 6**).

## 4. Determining Best-fit Models of Evolution

1. Click here to download the ProtTest[13] program:
http://darwin.uvigo.es/our-software/
2. Once ProtTest is downloaded, double-click on the ProtTest.jar file
3. Once ProtTest is launched, click on "select file" and load the sequence data (**Figure 7**).
4. Then click "start" and the program will begin (**Figure 8**).
*Note:* After completion of the run (**Figure 8**), the program will indicate the best model based on criteria *e.g.* "Best model according to AIC: WAG+I+G"

## 5. Inferring Sequence Based Phylogenies by Maximum Likelihood or Bayesian Inference

1. Downloaded PhyML[4] here:
https://code.google.com/p/phyml/
2. Launch the executable by double clicking the appropriate application (*i.e.* phyml Windows, phyml Linux, *etc*.) and the interface window will pop up (**Figure 9**).
3. Load the input sequence as a PHYLIP formatted sequence by typing:
>"file name".phy
Note: *To convert between sequence formats, use the "Readseq" web program available at -* http://iubio.bio.indiana.edu/cgi-bin/readseq.cgi.
4. Launch the program by typing "Y".
5. Download MrBayes[5] here:
http://mrbayes.sourceforge.net/download.php
6. To start the program click on the executable file and read NEXUS formatted sequence data into the program by typing:
>execute "file name".nex

7. Set the evolutionary model.
8. Select the number of generations to run by typing:
   >mcmcp ngen = 1000000 [this sets the number of generations to 1000000]
   >sump burnin =10000 [this sets the burnin to 10000]
9. Save the branch lengths in the results file by typing:
   >mcmcp savebrlens = yes
10. Run the analysis by typing:
    >mcmc
11. Summarize the trees using the "sumt" command.

# 6. Visualizing Phylogenies

1. View a list of tree viewer programs here:
   http://www.treedyn.org/overview/editors.html
2. Download the TreeView[14] program here:
   http://taxonomy.zoology.gla.ac.uk/rod/treeview.html

## Representative Results

Finding similarities to a query allows researchers to ascribe a potential identity to new sequences and also infer relationships between sequences. The file input type for BLAST[1] is FASTA formatted text sequence or GenBank accession number. FASTA formatted sequence begins with a description line indicated by a ">" sign (**Figure 2**). The description must follow immediately after the ">" sign, the sequence (*i.e.* nucleotides or amino acids) follow the description on the next line. When saving and editing sequence files, it is best to use a text editor such as "Notepad" on PC or TextWrangler (http://www.barebones.com/products/textwrangler/) for Mac. The BLAST algorithm performs "local" alignments, which searches for short stretches of sequence similarity. After the algorithm has looked up all possible "stretches" from the query sequence and maximally extended these sequences, it then assembles alignments for each query sequence pair. It is then important to understand how good these matches are, and so BLAST applies statistics to each hit which comprise an expect value (E) and a bit score. The E value gives an indication of the statistical significance of a match. The lower the E-value, the more significant the hit, for example a sequence alignment with an E-value of 0.05 means that the likelihood of this match occurring by chance alone is 5 in 100. The bit score uses a specific scoring matrix to provide an indication of how good the alignment is. The higher the bit score, the better the alignment. Similar to the online version of BLAST, there are a number of parameters that can be set via commands using the local BLAST executable. A comprehensive resource describing these commands can be found here - http://www.ncbi.nlm.nih.gov/books/NBK1762/. The output of the local search is a text file just like the output from the online BLAST interface (**Figure 4**).

A Multiple Sequence Alignment (MSA) is a sequence alignment of three or more primary sequences composed of amino acids, DNA, or RNA. ClustalW[6] released in 1994, is one of the most popular MSA tools for biologists. A user friendly online interface that provides one-stop access to several popular MSA tools can be found at the EMBL-EBI server here - http://www.ebi.ac.uk/Tools/msa. The input for each program can be FASTA formatted sequence data (see **Figure 2**) although many different formats are also accepted, and numerous mirror sites for each can be found online. Numerous parameters like gap penalties and output formats can be easily chosen. A sample output from the MSA T-Coffee can be seen in **Figure 5,** where similar residues are color coded. In some cases, the MSA tool can also be downloaded and executed locally. Clustal can be downloaded as a command line version (ClustalW) or a graphical version (ClustalX) from this website - http://www.clustal.org/clustal2/. To download, just click on the appropriate executable (*i.e.* win, Linux, Mac OS X). For Windows the program executable will download and a pop-up menu will require the user to click "Run", and then installation will begin. The program is very intuitive, sequences can be loaded from a text file containing sequences formatted as NBRF/PIR, FASTA, EMBL/Swiss-Prot, Clustal, GCC/MSF, GCG9 RSF, and GDE. Sequences are aligned by clicking "do complete alignment" from the "alignment" menu. A sample alignment of six protein sequences aligned using ClustalX can be seen in **Figure 6**. Various parameters such as font size and color can be easily modified, and editing of sequences is done by clicking on the "Edit" menu. Manually refined alignments are often superior to fully automated methods and because of this, MSA tool development is a very active area of research. Some common alignment editors can be found at the following links: Se-Al - http://tree.bio.ed.ac.uk/software/seal/; BSEdit - http://www.bsedit.org/; JalView - http://www.jalview.org/; SeaView - http://pbil.univ-lyon1.fr/software/seaview.html.

For amino-acid alignments the program ProtTest[13] is used to determine the selection of best-fit models of amino acid replacements within the data. ProtTest makes this selection by finding the model from the list of candidate models with the smallest Akaike Information Criterion (AIC), Bayesian Information Criterion (BIC) score, or Decision Theory Criterion (DT). The latest version of ProtTest (version 3.2) includes 15 different rate matrices that result in 120 different models. The user must have Java Runtime on their system to run ProtTest. Java Runtime is freely available here - http://www.java.com/en/download/chrome.jsp. Sequences are inputted as PHYLIP or NEXUS format. To convert between sequence formats, use the "Readseq" web program available at - http://iubio.bio.indiana.edu/cgi-bin/readseq.cgi. Click on "select file" and load the sequence data. Then click "start" and the program will begin. To modify the number of models selected, you can click the "models" button. Once the program begins it will display a progress bar at the bottom and list the models as they are being analyzed (**Figure 8**). The developer's manual provides a comprehensive user guide to ProtTest (can be found in the "manual" folder when the program is downloaded) and also provides descriptions of the models and parameters. More background information can be found here - https://code.google.com/p/prottest3/wiki/Background. There is also an online web interface for ProtTest which functions just like the downloaded version except that it can only handle a limited number of sequences. This web interface can be accessed by clicking here - http://darwin.uvigo.es/software/prottest2_server.html. For nucleotide datasets the program jModelTest[15] is used to examine the statistical selection of best-fit models of nucleotide substitutions by implementing the AIC, BIC, and DT criteria outlined above and also hierarchical and dynamical likelihood ration tests (hLRT and dLRT). jModelTest is optimized for Mac OS X. For the input, multiple formats are permitted. A clear step-by-step guide is available by the developers here - http://computing.bio.cam.ac.uk/local/doc/jmodeltest.pdf

PhyML is a program that estimates maximum likelihood phylogenies from alignments of nucleotide or amino acid sequences. PhyML will incorporate a large number of substitution models coupled to various options to search tree topology space (**Figure 10**). The program will

save results into two text files. The first file will contain the ML tree in Newick format which can easily be viewed using a Tree viewer (see protocol 6), and the other file will contain the statistics (filename, model, Log-likelihood scores, *etc.*) of the analysis. All parameters are very easily set by following the Menu items. More detailed descriptions of each Menu option are explained in the PhyML manual available on the PhyML download page - https://code.google.com/p/phyml/downloads/list. MrBayes[5] is a program that utilizes Bayesian MCMC inference across a number of evolutionary models to reconstruct phylogenetic relationships. The program behaves the same on all platforms and once downloaded the installer will install the executable. To start the program, simply click on the executable. There are numerous models that can be set and details of each model and their commands can be found here - http://mrbayes.sourceforge.net/wiki/index.php/Tutorial. Another help option is to type "help lset" – this will provide details on Model setting. For example "Prset aamodelpr=mixed" will permit mixed modeling or "prset aamodelpr=fixed(wag)" will set the amino acid model to the WAG model. An outgroup can be easily set by specifying the Taxon number "outgroup 30"; the program automatically lists the sequences/Taxa by number. If an outgroup is not specified the tree will be unrooted. Once the program is running (**Figure 11**) the progress can be viewed in specific intervals which can be set using the "printfreq=X" command. More details on when to stop the analysis (*i.e.* how many generations to run for) can be found in the user's manual. Clade values on a cladogram are provided in the results alongside a phylogram which is also provided in Newick format that can easily be viewed using a tree viewer (see protocol 6).

Once a phylogenetic tree is generated, the topology needs to be visualized. There are many online tools and downloadable applications used to visualize tree topologies. A partial list of popular programs can be viewed here - http://en.wikipedia.org/wiki/List_of_phylogenetic_tree_visualization_software, and a more comprehensive list can be found here - http://www.treedyn.org/overview/editors.html. TreeView[14] and TreeDyn[16] are two popular choices. Both are very user friendly and easy to become familiar with the various options. TreeView runs on Mac and Windows, using almost identical interfaces. The input can be one of several formats including NEXUS, PHYLIP, Hennig86, MEGA, and ClustalW/X. TreeView (**Figure 12**) also includes a tree editor that allows the user to move branches, reroot trees, and rearrange the appearance of the tree.
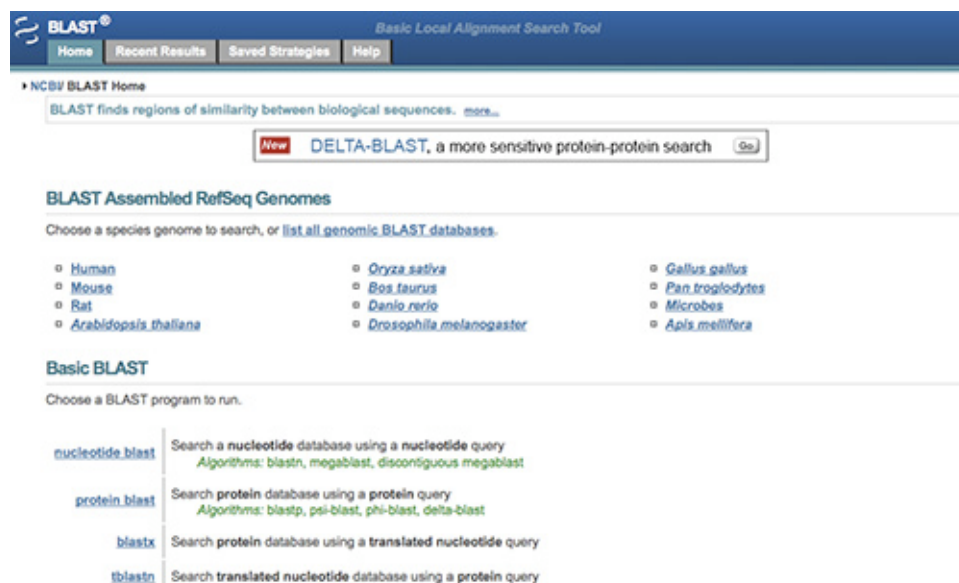


**Figure 1. NCBI BLAST web-page**. The BLAST web server contains a suite of BLAST programs and is hosted by the National Center for Biotechnology Information (NCBI). Click here to view larger image.



**Figure 2. FASTA formatted sequence**. FASTA format begins with a description line indicated by a ">". The description must follow immediately after the ">" sign, the sequence (*i.e.* nucleotides or amino acids) follow the description on the next line. Click here to view larger image.
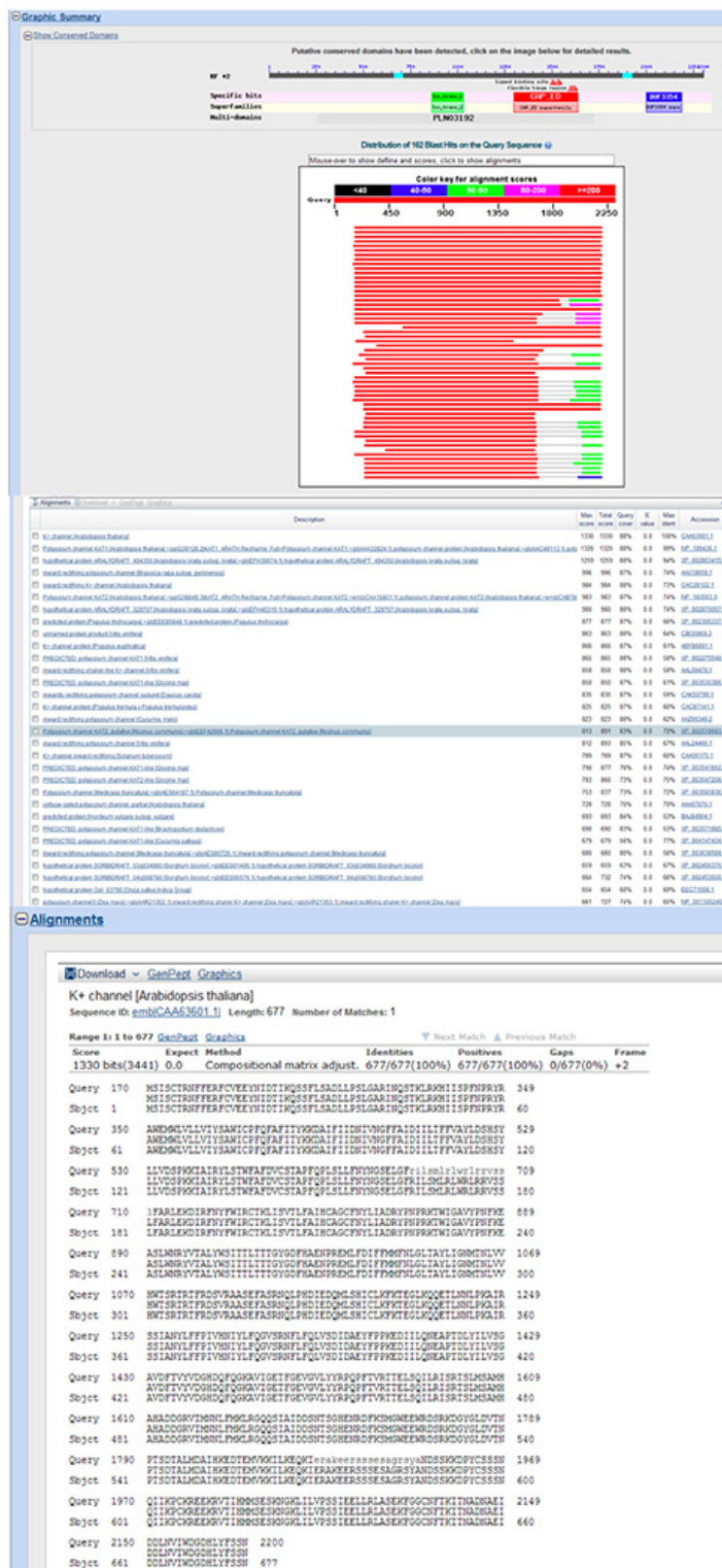
**Figure 3. HTML output from a BLAST search.** The output from the BLAST search illustrates the areas of identity within the query sequence, and also provides bit-scores, expect values and pairwise alignments with each match. Click here to view larger image.

File Edit Format View Help

```
BLASTP 2.2.27+


Reference: Stephen F. Altschul, Thomas L. Madden, Alejandro A.
Schaffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J.
Lipman (1997), "Gapped BLAST and PSI-BLAST: a new generation of
protein database search programs", Nucleic Acids Res. 25:3389-3402.


Reference for composition-based statistics: Alejandro A. Schaffer,
L. Aravind, Thomas L. Madden, Sergei Shavirin, John L. Spouge, Yuri
I. Wolf, Eugene V. Koonin, and Stephen F. Altschul (2001),
"Improving the accuracy of PSI-BLAST protein database searches with
composition-based statistics and other refinements", Nucleic Acids
Res. 29:2994-3005.


Database: db/briggsae.fasta
           19,404 sequences; 8,001,967 total letters


Query= gi|351058250|emb|CCD65668.1| Protein HSP-16.1 [Caenorhabditis
elegans]

Length=145
                                                        Score      E
Sequences producing significant alignments:            (Bits)   Value

  gi|268558958|ref|XP_002637470.1|  Hypothetical protein CBG19186 ...   229    2e-077
  gi|268559656|ref|XP_002637819.1|  Hypothetical protein CBG04608 ...   222    9e-075
  gi|268559628|ref|XP_002637805.1|  Hypothetical protein CBG04591 ...   221    2e-074
  gi|268558954|ref|XP_002637468.1|  Hypothetical protein CBG19184 ...   214    2e-071
  gi|268559652|ref|XP_002637817.1|  Hypothetical protein CBG04606 ...   165    2e-052
  gi|268558956|ref|XP_002637469.1|  Hypothetical protein CBG19185 ...   125    1e-036
  gi|268559650|ref|XP_002637816.1|  Hypothetical protein CBG04605 ...   122    2e-035
  gi|268558960|ref|XP_002637471.1|  Hypothetical protein CBG19187 ...   120    8e-035
  gi|268559630|ref|XP_002637806.1|  Hypothetical protein CBG04592 ...   116    3e-033
  gi|268559654|ref|XP_002637818.1|  Hypothetical protein CBG04607 ...   115    1e-032
  gi|268574778|ref|XP_002642368.1|  C. briggsae CBR-SIP-1 protein ...   107    1e-029
  gi|268579717|ref|XP_002644841.1|  C. briggsae CBR-HSP-43 protein...  70.5    6e-015
  gi|268579717|ref|XP_002644841.1|  C. briggsae CBR-HSP-43 protein...  70.5    6e-015
  gi|268554144|ref|XP_002635059.1|  C. briggsae CBR-HSP-17 protein...  60.1    3e-012
  gi|268573196|ref|XP_002641575.1|  Hypothetical protein CBG09876 ...  59.3    7e-012
  gi|268573992|ref|XP_002641973.1|  C. briggsae CBR-HSP-12.2 prote... 57.8    2e-011
  gi|268566747|ref|XP_002639803.1|  Hypothetical protein CBG02254 ...  48.9    2e-008
  gi|268570300|ref|XP_002648467.1|  Hypothetical protein CBG24755 ...  48.9    2e-008
  gi|268535548|ref|XP_002632907.1|  Hypothetical protein CBG21660 ...  43.1    2e-006
  gi|268579635|ref|XP_002644800.1|  C. briggsae CBR-HSP-25 protein...  42.4    1e-005
  gi|268579635|ref|XP_002644800.1|  C. briggsae CBR-HSP-25 protein...  42.4    1e-005
  gi|268564238|ref|XP_002647119.1|  Hypothetical protein CBG23899 ...  41.2    4e-005
  gi|268579371|ref|XP_002644668.1|  Hypothetical protein CBG14649 ...  27.7    2.4
  gi|268564057|ref|XP_002639005.1|  Hypothetical protein CBG22251 ...  26.9    3.6
  gi|268563170|ref|XP_002638772.1|  C. briggsae CBR-PTR-21 protein...  26.6    4.8
  gi|268553777|ref|XP_002634875.1|  Hypothetical protein CBG10541 ...  26.6    5.3
  gi|268562597|ref|XP_002646699.1|  Hypothetical protein CBG13076 ...  25.8    9.4


> gi|268558958|ref|XP_002637470.1|  Hypothetical protein CBG19186
[Caenorhabditis briggsae]
Length=146

 Score =  229 bits (584),  Expect = 2e-077, Method: Compositional matrix adjust.
 Identities = 112/146 (77%), Positives = 129/146 (88%), Gaps = 1/146 (1%)

Query  1    MSLYHYFRPAQRSVFGDLMRDMAQMERQFTPVCR-GSPSESSEIVNNDQKFAINLNVSQF  59
            MSLY +FRP   SV G++MRD+A+MERQ P+    +P+ +SEIVNNDQKFAINLNVSQF
Sbjct  1    MSLYPFFRPRPFSVIGEMMRDIARMERQFIPISAFEAPAAASEIVNNDQKFAINLNVSQF  60

Query  60   KPEDLKINLDGHTLSIQGEQELKTEHGYSKKSFSRVILLPEDVDVGAVASNLSEDGKLSI  119
            KPEDLKINLDG TLSIQGEQE+K EHG+SKKSFSR+ILLPEDVD+GAVASNLSEDGKLSI
Sbjct  61   KPEDLKINLDGRTLSIQGEQEVKDEHGHSKKSFSRIILLPEDVDIGAVASNLSEDGKLSI  120

Query  120  EAPKKEAIQGRSIPIQQAPVEQKTSE  145
```
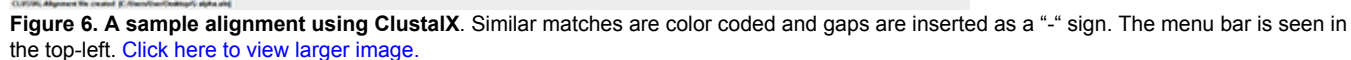
**Figure 4. A sample output from a local BLAST executable search.** The output of this search is a text file just like the output from the online BLAST interface, that include the expect value and bit score, as well as match description. Click here to view larger image.
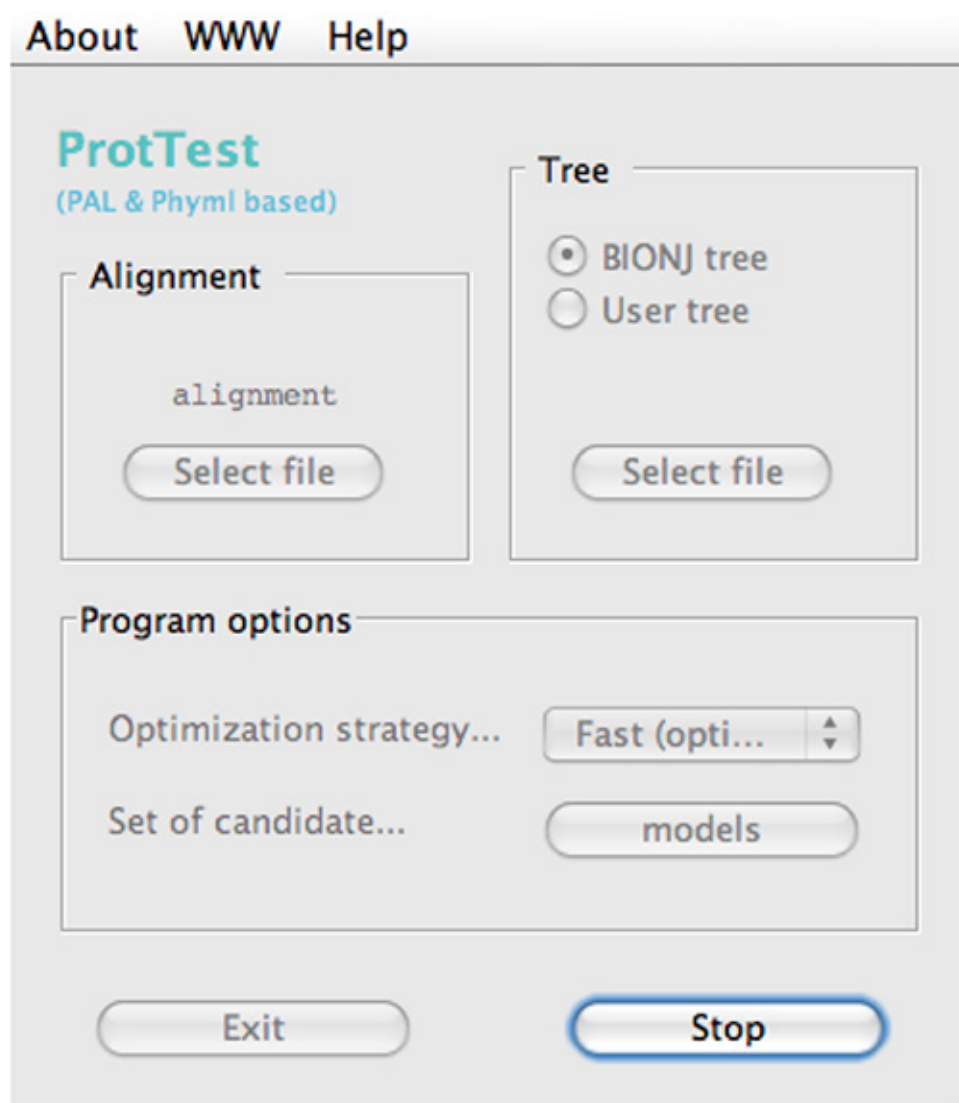
**Figure 5. Output of a MSA using T-Coffee**. The output highlights similar sites and weights the match by color. Gaps are inserted as "-" signs and the residue or nucleotide position is preserved for each taxon. Click here to view larger image.



**Figure 6. A sample alignment using ClustalX**. Similar matches are color coded and gaps are inserted as a "-" sign. The menu bar is seen in the top-left. Click here to view larger image.

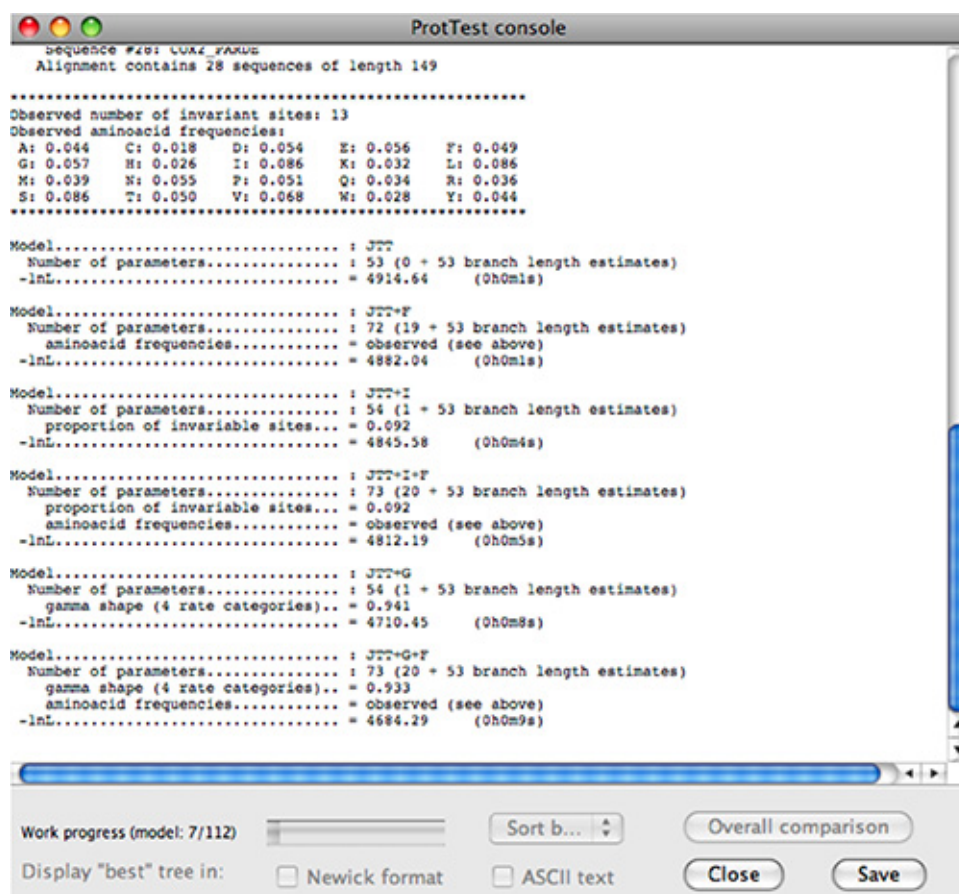**Figure 7. The ProtTest program interface.** Click here to view larger image.

**Figure 8. The ProtTest console.** ProtTest console while running an analysis. The progress bar indicates how many models have been completed, and the main window displays the log likelihood score for each model. Click here to view larger image.
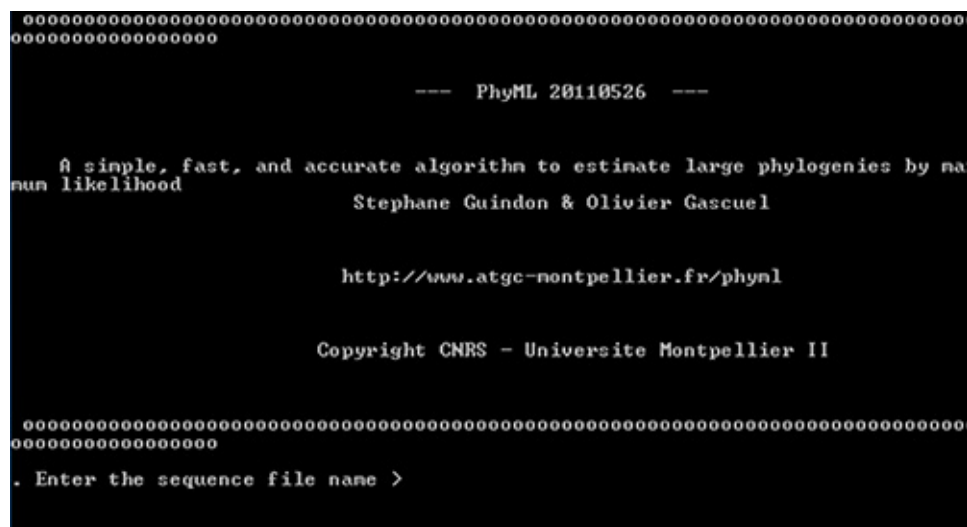


**Figure 9. The PhyML interface.** Click here to view larger image.

**Figure 10. The PhyML interface menu.** Once sequences are loaded into PhyML the first menu appears, which can be navigated by typing the letter or symbol in the square bracket. Submenus can be reached by typing the "+" sign. Click here to view larger image.



**Figure 11. MrBayes Interface.** When MrBayes is launched the progress can be viewed in specific intervals set using the "printfreq=X" command. Although the program cannot be stopped during a run, after the specified number of generations are computed the user will be asked if they want to run more generations. Click here to view larger image.
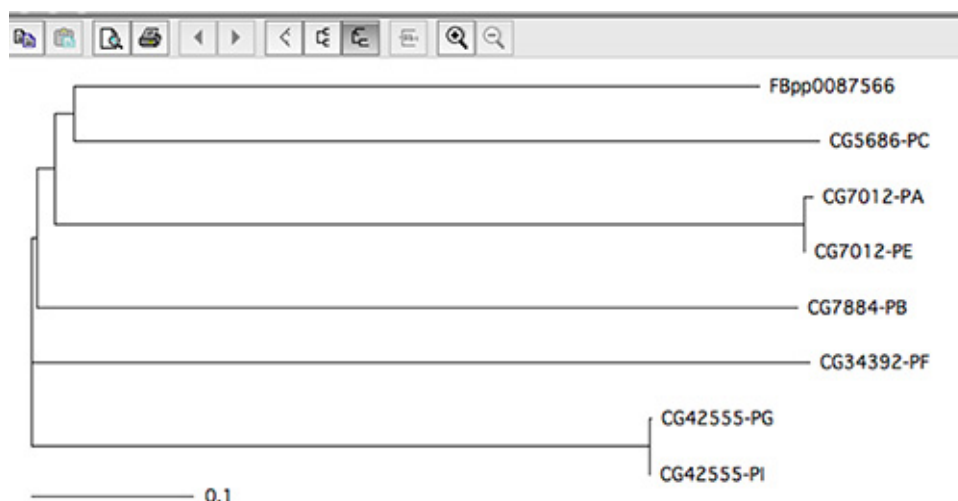
**Figure 12. The TreeView interface.** In this figure the TreeView window displays a sample tree of proteins from *Flybase* (http://flybase.org/). Files are imported by clicking the "open" option, and selecting an appropriate file type (*e.g*. Newick format). Click here to view larger image.

## Discussion

Our hope for this article is that it will serve as a starting point to guide researchers or students that are new to phylogenetics. Genome sequencing projects have become less expensive over the last few years and as a consequence the user demand for this technology is increasing, and now the production of large sequence datasets is commonplace in small labs. These datasets often provide researchers with sets of genes that require a phylogenetic framework to begin to understand their function. Furthermore, because phylogenetics is finding a home in an ever increasing number of research labs, we also intend for this article to serve as an educational device for students interested broadly in biological research. By providing user information on the "why", "how", and "where" for commonly used tree-building tools, we provide a framework for the reader to begin to familiarize themselves with these applications and how they work. However, we advise the reader to play around with all the settings within each tool in an attempt to understand how the various parameters can influence their sequence data, and to ensure compatibility between platform and software in each case. The analysis outlined above was computed using a Dell Optiplex 990 with Intel core i7 processor and a MacBook laptop with an Intel Core 2 Duo processor, however, the speed of analysis and also the specific binaries (*e.g*. 32 bit or 64 bit) will depend on the user's platform.

A challenge when compiling a user guide like this one for phylogenetics, is that the field of phylogenetics, and bioinformatics as a whole, is a rapidly expanding area of research that constantly releases new software aimed at providing better alignments, similarity predictions, or phylogenetic trees. To mitigate this problem, we tried to focus on programs that have been around for a number of years and are still popular on account of how well they work. That said, we want to point out that there are many other tools available to tackle the problems we have outlined in this article, and so encourage the reader to exploit this and incorporate multiple applications into their analyses.

## Disclosures

We have nothing to disclose.

## Acknowledgements

## References

1. Altschul SF, Carroll RJ, Lipman DJ. Weights for data related by a tree. *J. Mol. Biol*. **207**(4), 647-653 (1989).
2. Akaike H. A new look at the statistical model identification. *IEEE Trans. Automat. Contr*. **19**(6), 706-723 (1974).
3. Schwarz G. Estimating the dimension of a model. *Ann. Stat*. **6**(2), 461-464 (1978).
4. Guindon S, Gascuel O. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst. Biol*. **52**(5), 696-704 (2003).
5. Huelsenbeck JP, Ronquist F. MRBAYES: Bayesian inference of phylogenetic trees. *Bioinformatics*. **17**(8), 754-755 (2001).
6. Thompson JD, Higgins DG, Gibson TJ. CLUSTAL W: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res*. **22**(22), 4673-4680 (1994).
7. Lassmann T, Sonnhammer EL. Kalign--an accurate and fast multiple sequence alignment algorithm. *BMC Bioinformatics*. **6**, 298. doi: 10.1186/1471-2105-6-298 (2005).

8. Katoh K, Kuma K, Toh H, Miyata T. MAFFT version 5: Improvement in accuracy of multiple sequence alignment. *Nucleic Acids Res*. **33**(2), 511-518. doi: 10.1093/nar/gki198 (2005).

9. Katoh K, Misawa K, Kuma K, Miyata T. MAFFT: A novel method for rapid multiple sequence alignment based on fast fourier transform. *Nucleic Acids Res*. **30**(14), 3059-3066 (2002).

10. Edgar RC. MUSCLE: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*. **32**(5), 1792-1797. doi: 10.1093/nar/gkh340 (2004).

11. Notredame C, Higgins DG, Heringa J. T-coffee: A novel method for fast and accurate multiple sequence alignment. *J. Mol. Biol*. **302**(1), 205-217, doi: 10.1006/jmbi.2000.4042 (2000).

12. Do CB, Mahabhashyam MS, Brudno M, Batzoglou S. ProbCons: Probabilistic consistency-based multiple sequence alignment. *Genome Res*. **15**(2), 330-340, doi: 10.1101/gr.2821705 (2005).

13. Darriba D, Taboada GL, Doallo R, Posada D. ProtTest 3: Fast selection of best-fit models of protein evolution. *Bioinformatics*. **27**(8), 1164-1165, doi: 10.1093/bioinformatics/btr088; 10.1093/bioinformatics/btr088 (2011).

14. Page RD. TreeView: An application to display phylogenetic trees on personal computers. *Comput. Appl. Biosci*. **12**(4), 357-358 (1996).

15. Darriba D, Taboada GL, Doallo R, Posada D. jModelTest 2: More models, new heuristics and parallel computing. *Nat. Methods*. **9**(8), 772, doi: 10.1038/nmeth.2109; 10.1038/nmeth.2109 (2012).

16. Chevenet F, Brun C, Banuls AL, Jacq B, Christen R. TreeDyn: Towards dynamic graphics and annotations for analyses of trees. *BMC Bioinformatics*. **7**, 439, doi: 10.1186/1471-2105-7-439 (2006).