# 15 Phylogenetic Analysis of DNA Sequence Data

SERGEI A. SUBBOTIN

*Plant Pest Diagnostic Center, California Department of Food and Agriculture, California, USA; Center of Parasitology of A.N. Severtsov Institute of Ecology and Evolution, Moscow, Russia*

## 15.1 Introduction

The goal of phylogenetics is to construct relationships that are true representations of the evolutionary history of a group of organisms or genes. The history inferred from phylogenetic analysis is usually depicted as branching in tree-like diagrams or networks. In nematology, phylogenetic studies have been applied to resolve a wide range of questions dealing with improving classifications and testing evolution processes, such as co-evolution, biogeography and many others. There are several main steps involved in a phylogenetic study:

- selection of ingroup and outgroup taxa for a study;
- selection of one or several gene fragments for a study;
- sample collection, obtaining PCR products and sequencing of gene fragments;
- visualization, editing raw sequence data and sequence assembling;
- search for sequence similarity in a public database;
- making and editing multiple alignment of sequences;
- selecting appropriate DNA model for a dataset;
- phylogenetic reconstruction using minimum evolution, maximum parsimony, maximum likelihood and Bayesian inference;
- visualization of tree files and preparation of tree for a publication; and
- sequence submission to a public database.

Molecular phylogenetic study requires particularly careful planning because it is usually relatively expensive in terms of the cost in reagents and time. The first and the most important step of any study is to define clearly the specific biological question to be answered. When the biological problem is formulated and the

literature survey pertaining to the group of interest is completed, selection of ingroup and outgroup taxa and appropriate genes should be done. These could be representatives of the same species from various locations or hosts, or different species of the same genus, or representatives of related genera or higher taxonomic categories. It is advisable to include as many samples as possible, as well as to choose several gene fragments to reduce artefactual associations between terminals.

The third stage, which includes sample collection, obtaining PCR products and sequencing, may require several weeks or months and is the most time-consuming stage in a study. Before launching a full-scale project, a pilot study with a limited sample number is recommended to determine if the selected genetic markers give sufficient resolution for phylogeny of the studied group.

The last seven stages require a computer with an Internet connection and suitable software. This chapter deals with these stages. The software and instructions for its use are discussed below. All of them are free and can be downloaded from different websites, except for PAUP* (Phylogenetic Analysis Using Parsimony*), which can be purchased from the official website. This chapter only considers data from nucleotide sequences and does not take into account amino acid sequences data.

There are many phylogenetics programs that perform similar functions (Table 15.1). The comprehensive list can be found in the websites:

- http://evolution.genetics.washington.edu/phylip/software.html;
- http://www.phylo.org;
- https://www.phylogeny.fr;
- https://isogg.org/wiki/Phylogeny_programs;
- https://en.wikipedia.org/wiki/List_of_phylogenetics_software; and
- https://molbiol-tools.ca/Phylogeny.htm.

**Table 15.1.** Phylogenetics programs.

| Software | Website |
|---|---|
| **BEAST** (**B**ayesian **E**volutionary **A**nalysis **S**ampling **T**rees) is a cross-platform program for Bayesian phylogenetic analysis of molecular sequences using MCMC | https://beast.community https://www.beast2.org |
| **BioEdit** is a **Bio**logical Sequence **Edit**or | https://bioedit.software.informer.com |
| **Clustal** is a series of widely used computer programs used in Bioinformatics for multiple sequence alignment | http://www.clustal.org |
| **GARLI** (**G**enetic **A**lgorithm for **R**apid **L**ikelihood **I**nference) is a program for inferring phylogenetic trees | https://code.google.com/archive/p/garli/ |
| **MAFFT** (**M**ultiple **A**lignment using **F**ast **F**ourier **T**ransform) is a program used to create multiple sequence alignments of amino acid or nucleotide sequences. | https://mafft.cbrc.jp/alignment/software/ |
| **MEGA** (**M**olecular **E**volutionary **G**enetic **A**nalysis) is a software for conducting statistical analysis of molecular evolution and for constructing **phylogenetic** trees. | https://www.megasoftware.net |
| **MrBayes** is a free software tool that performs Bayesian inference of phylogeny | https://nbisweden.github.io/MrBayes/ |
| **MUSCLE** (**Mu**ltiple **S**equence **C**omparison by **L**og-**E**xpectation) is a software for multiple sequence alignment of protein and nucleotide sequences. | http://www.drive5.com/muscle/ |
| **PAUP*** - **P**hylogenetic **A**nalysis **U**sing **P**arsimony (*and other methods) | http://paup.phylosolutions.com |
| **RAxML** (**R**andomized **Ax**elerated **M**aximum **L**ikelihood) is a program for sequential and parallel maximum likelihood-based inference of large phylogenetic trees. | https://cme.h-its.org/exelixis/web/software/raxml/index.html http://www.trex.uqam.ca/index.php?action=raxml&project=trex |
| **T-Coffee** (**T**ree-based **C**onsistency **O**bjective **F**unction **f**or Alignm**e**nt **E**valuation) is a multiple sequence alignment software using a progressive approach | http://www.tcoffee.org |

The reliability and practicality of the software depends on the structure and size of the data. The merits and pitfalls of various methods are the subject of often acrimonious debates in taxonomic and phylogenetic journals (Lemey *et al.*, 2009; Yang and Rannala, 2012). Generally, most software packages have been developed and maintained through the efforts of scientists in related fields and released under free software licences.

Phylogenetics is a rather complex and rapidly expanding field of research. In this chapter, some basic approaches for the analysis of nucleotide non-coding gene sequences are given and discussed with an assumption that the gene tree reflects the organism phylogeny. Many other important software with statistical testing for phylogenetic studies are not covered here, and specialized literature (e.g. Lemey *et al.*, 2009; Hall, 2017) may be recommended for researchers.

## 15.2   Visualization and Editing of Raw Sequence Data

The sequencer is a laser-based instrument that utilizes fluorescent labels to analyse the products of a sequencing reaction as they migrate through a gel. After the data are collected from a sequencing run, special software identifies and tracks the sample lanes of the gel and subsequently normalizes and integrates the data into files. Automated DNA sequencer generates two file types: (i) a four-colour chromatogram showing the results of the sequencing run; and (ii) a text file of sequence data. It is always highly recommended to check the quality of the chromatogram file of a studied sample before converting it to a sequence text file.

There are several programs for visualization of a raw sequence data, two of which can be freely downloaded: Chromas (http://technelysium.com.au/wp/chromas/) developed by Technelysium Pty Ltd and FinchTV (http://www.geospiza.com/ftvdlinfo.html) originally designed by Geospiza. Chromas is a simple, easy-to-use viewer and editor for chromatograms from automated Sanger sequencers. Chromas has several features including automatic removal of low-quality sequence or vector sequences, copying the sequence to the clipboard in plain text, FASTA format for pasting into other applications, performing reverse and complement the sequence and chromatogram, and displaying translations in three frames along with the sequence. A chromatogram shows a sequence of peaks in four colours, each representing the base: A, green or yellow; G, black; T, red; and C, blue.

Once the sequence is obtained, the quality of sequence reading should be proofread to ensure that all ambiguous sites are correctly resolved in a chromatogram file. Good-quality sequences are characterized by well-defined peak resolution, uniform peak spacing and high signal-to-noise ratios (Fig. 15.1). If the quality of a sequence is not good and contains double or asymmetrical peaks or strong background noise, sequencing reactions should be repeated with both forward and reverse primers (Fig. 15.2). There are many reasons why a sequence reaction can fail. A good sequence generally begins approximately around 20–30 bp and lasts up to 700–800 bp. Any bad-quality sequence should be eliminated from further analysis. Sometimes, the software may miscall or miss a nucleotide or add an additional nucleotide, making proofreading of chromatograph files an important step of data verification. A single peak position within a trace may have also two peaks of different colours instead of just one. This is a common problem when sequencing a PCR product derived
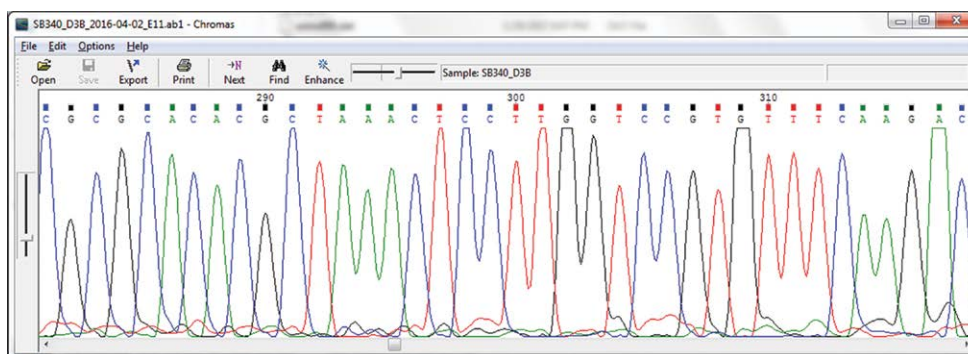


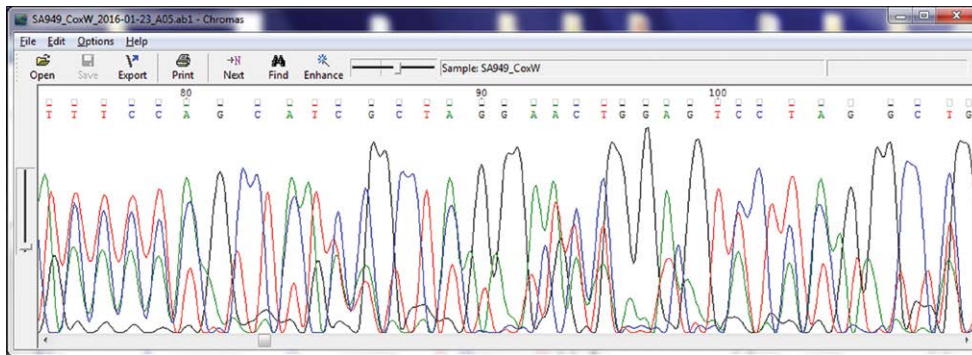**Fig. 15.1.** Chromas with good-quality sequence containing well-defined peaks.

**Fig. 15.2.** Chromas with bad-quality sequence containing double peaks and strong background.



**Fig. 15.3.** Sequence in FASTA format.

from diploid genomic DNA having polymorphic positions. In this case, it is recommended to edit the base in a chromatogram file according to the single letter nucleic acid code recommendations: R, G or A (purine); Y, T or C (pyrimidine); K, G or T (keto); M, A or C (amino); S, G or C (strong); W, A or T (weak); B, G or T or C; D, G or A or T; H, A or C or T; V, G or C or A; N, A or G or C or T (any).

Work with Chromas version 2.62 includes a few steps:

- Open a chromatogram file. Choose **File > Open**.
- If an antisense chain (3′–5′) is open, use reverse and complement option. Choose **Edit > Reverse+Complement**.
- Check quality and trim low-quality data. Choose **Edit > Trim Low Quality.**
- Convert into a FASTA format, which begins with a single-line description, followed by lines of sequence data. Choose **File > Export**. After a chromatogram file has been examined and edited, it should be exported into a FASTA format file (Fig. 15.3).

## 15.3  Sequence Assembling

If a PCR product cannot be covered by a single sequencing reaction, several reactions should be performed. The chromatograms of these reactions should be verified and then displayed in a single FASTA format file. The goal of assembling is to create a single consensus sequence covering the whole length of a studied amplicon from several partly overlapped sequences.

A number of DNA sequence assembly programs have been developed including CAP3, which has also web-based version (http://doua.prabi.fr/software/cap3). For more advanced usage, it is recommended to install the software on your local computer. The program features include fast identification of pairs of reads
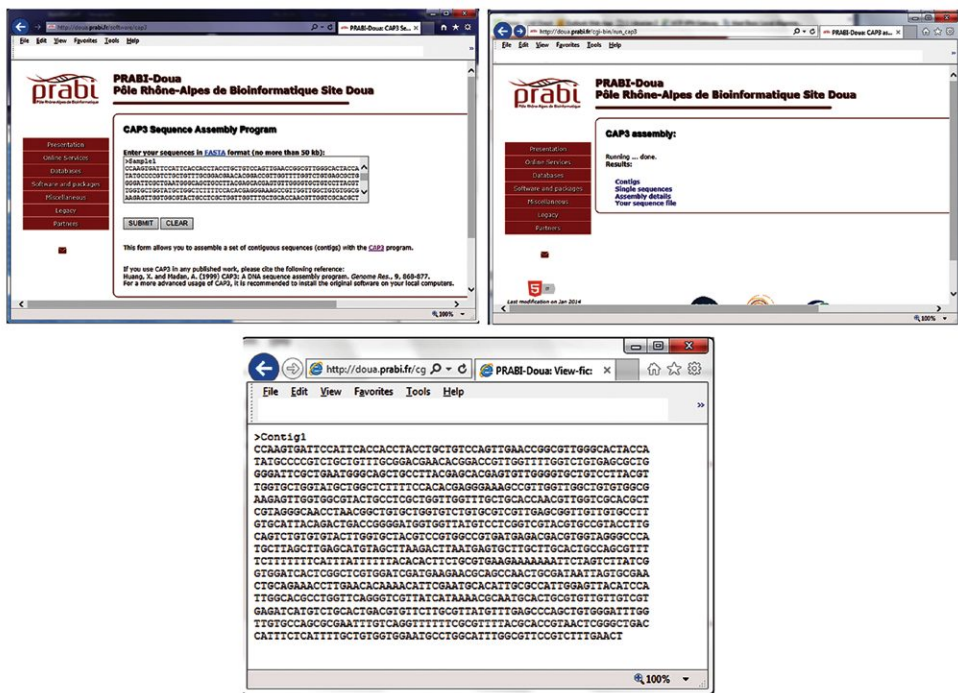
**Fig. 15.4.** CAP3 home page and output.

with an overlap, clipping of 5′ and 3′ poor regions of reads, efficient computation and evaluation of overlaps, use of forward–reverse constraints to correct errors in construction of a contig, and generation of consensus sequence for a contig (Huang and Madan, 1999).

Work with web-based CAP3 is simple and the results are self-explanatory.

- Enter your sequences in FASTA format in a window and submit the data (Fig. 15.4).
- Results are displayed as contigs, single sequence, assembling details and your sequence file. Every base in an assembly must be covered by at least two sequences of high quality. Validating sequence coverage provides a high degree of confidence in the consensus base calls.

## 15.4   Similarity Search in a Sequence Database

One of the most important steps in the study is the identification of your sequence and comparison with all known sequences collected in different databases. This procedure is called similarity search. BLAST (Basic Local Alignment Search Tool) is a powerful program for rapid searching of nucleotide and protein databases. The BLAST program finds regions of local similarity between sequences. It was developed in 1989 at the National Center for Biotechnology Information (NCBI), Maryland, USA. A BLAST query uses statistical methods to compare a DNA or protein input sequence ('query sequence') to a database of sequences ('subject sequences') and returns those sequences that have a significant level of similarity to the query sequence. The BLAST algorithm calculates similarity scores for local alignments (i.e. the most similar regions between two sequences) between the query sequence and subject sequences using specific scoring matrices, and returns a table of the best matches ('hits') from the database (Altschul, 1990; Wheeler and Bhagwat, 2007).

A BLAST search includes a few steps:

- Point your browser to the NCBI BLAST server at: http://www.ncbi.nlm.nih.gov/BLAST.
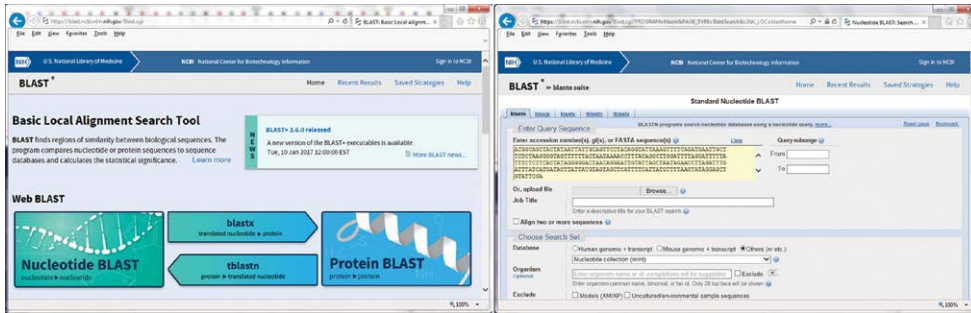- Select **Nucleotide BLAST** (Fig. 15.5).

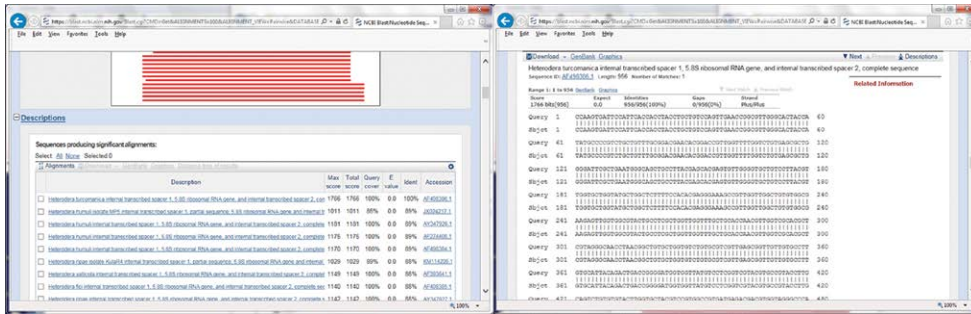**Fig. 15.5.** BLAST searching pages.



**Fig. 15.6.** BLAST output pages.

- Paste your sequence in the Query Sequence field, choose search set in the Database as **Others, Nucleotide collection. Organism group** could be also defined to limit your search to the DNA of a specific organism. Alignment parameters can be also modified (for example, increase: **Max target sequences**). The default setting uses a version of BLASTN called **megaBLAST**. Click the **BLAST** button and wait for the results.

The results from a BLAST search are divided into three sections: (i) the graphic pane; (ii) a results table; and (iii) the alignments between the query and the hits. Conclusions can be based on interpretations of the BLAST results table that provides basic information about the hits together with the statistics of each hit. Information includes: Max score (highest alignment score (bit-score) between the query sequence and the database sequence segment); Total score (sum of alignment scores of all segments from the same database sequence that match the query sequence (calculated over all segments)); Query cover (the percentage of the query that aligns with the hit); E value (number of alignments expected by chance with a particular score or better); Identity (percentage of identity between the query and the hit in a nucleotide to nucleotide alignment); and Accession (GenBank sequence number). The top hits are most significant and similar to the submitted sequence (Fig. 15.6).

## 15.5  Sequence Retrieval from the Database

There are several ways to retrieve sequences to build your dataset for phylogenetic analysis:

- Sequences can be obtained from the BLAST search result page. Click on **Accession** and the hyperlink takes you to the database entry that contains this sequence. Click on **FASTA** to convert the GenBank format into a FASTA format (Fig. 15.7).
- Point your browser to http://www.ncbi.nlm.nih.gov. Select **Nucleotide** in a search field and type organism name(s) and gene name or accession number(s).
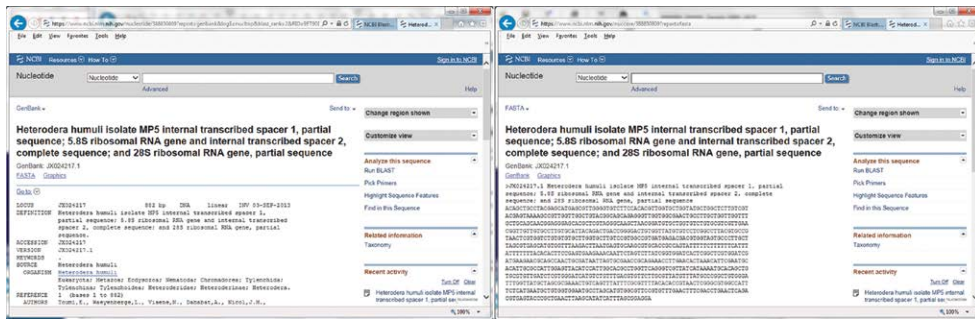
**Fig. 15.7.** GenBank entry JX024217.

## 15.5.1 Multiple alignment

The construction of alignment is the establishment of positional homology between nucleotides or amino acid bases that have descended from a common ancestral base. Errors incurred in this step can lead to an incorrect phylogeny. Multiple alignment construction is an exclusively mathematical process and is usually constructed using computer programs with particular algorithms. Most alignments are made based only on sequence information. However, aligning according to secondary structure is considered more reliable than sequence-based alignment because confidence in homology assessment is greater when comparisons are made based on structures rather than on simple characters. This approach is time consuming, requires information on secondary structure and can be done only with special alignment software.

Clustal (http://www.clustal.org) is the most popular multiple sequence alignment program and one of the most widely cited scientific publication (Larkin *et al.*, 2007). There are two main variations: ClustalW with command line interface ('W' stands for 'weighting' different parts of alignment differently) and ClustalX with a graphical user interface. Alignment can be made using web-based Clustal or software downloaded onto a computer. Clustal performs a global-multiple sequence alignment by the progressive method using a three-step process: (i) perform pairwise alignment of all the sequences by dynamic programming; it aligns each sequence against each other giving a similarity matrix; (ii) create a guide dendrogram using the similarity matrix; and (iii) start by aligning the two most similar sequences. Following the guide dendrogram, the next sequences are added in, aligning them to the existing alignment and insert gaps as necessary.

- Double click on ClustalX icon and run the program. Load sequence the sequence file. Choose **File > Load Sequences.** Under the alignment menu, choose the output format options.
- Under this menu, it is possible to change the alignment parameters (Gap Opening Penalty, Gap Extension Penalty), both for pairwise alignment and for the multiple alignment stages or run with default options.
- Choose **Alignment > Do Complete Alignment** and click **OK.** The sequence alignment is displayed in a window on the screen. The histogram below the ruler indicates the degree of similarity (Fig. 15.8).
- Clustal generates two output files, with extension 'aln' (alignment result in Clustal format) and extension 'dnd' (guide dendrogram).

It is recommended to perform phylogenetic analyses based on a series of slightly modified alignments to determine how ambiguous regions in the alignment affect the results.

In addition to Clustal, other software can be used to make alignments (Table 15.1).

## 15.5.2 Sequence alignment editing

Any automatic alignment should be visually checked and then manually edited, if necessary. GeneDoc (https://genedoc.software.informer.com) is a multiple sequence editor with a full-featured alignment visualization

including shading and structural definition features, editing and analysis tools. It has an easy-to-use point and click user interface with extensive keyboard mapping for advanced users (Nicholas *et al*., 1997).

- **Run Genedoc**. To import the file generated by Clustal, choose **File > New > Import**. In Import Type window, select **Clustal** (**ALN**) and click **>Import**. In Open window, find your file and click> **Open** and then **Done**. The alignment will be displayed (Fig. 15.9).
- To see the full sequence names in the displayed alignment: use combination **Ctrl+G** or choose **Project > Configure** and then in the field 'Max NemLen' type '30' > **OK.**
- There are many features to edit alignment. To edit sequence list, use combination **Ctrl+Q** or select **Project > Edit Sequence list**. To move nucleotide(s) or insert gap(s), use **Grab+Drag bottom** (**Ctrl+A**)**, Grab+Slide bottom** (**F5**) Insert/Delete gap in single sequence (**F6**), **Insert Gap Colum** (**F7**)**.** To replace nucleotide, use **Ctrl+U** or **Edit > Residue Edit Mode.**
- To display the alignment with only nucleotide differences, use the combination **Ctrl+G** or choose **Project > Configure**, select **Shade** and in **Residue Display Mode > Differences** and in **Difference Mode Style > Diff/ Top Sequence,** click **OK.**



**Fig. 15.8.** A multiple sequence alignment generated with ClustalX and alignment parameters.



**Fig. 15.9.** A multiple sequence alignment in GenDoc.

- To trim excess sequence in the beginning and end of alignment or remove any alignment region: use combination **Ctrl+L** or click **Edit,** select **> Select Columns** and then use the cursor to select alignment region. To delete this region: use combination **Ctrl+D,** or click **Edit,** select **> Delete All Data.**
- To prepare the alignment file for a publication, use combination **Ctrl+E** or **Edit > Select Blocks for Copy,** and use a cursor to mark block(s). When all blocks are marked, click **Edit > Copy Selected Blocks** to **> RTF file**, then **Save as**. The saved file could be opened in Word.
- To save the file use **Ctrl+S** or **File > Save.** Export the corrected GenDoc file into a FASTA format, choose **File > Export,** then type the file name, click on **Save > Done.** The result will be saved.

Multiple alignment can also be automatically edited using the computer program Gblocks (http://molevol. cmima.csic.es/castresana/Gblocks.html) that eliminates poorly aligned positions and divergent regions of an alignment of DNA or protein sequences (Castresana, 2000; Talavera and Castresana, 2007). These positions may not be homologous or may have been saturated by multiple substitutions and it is convenient to eliminate them prior to phylogenetic analysis. Gblocks selects blocks in a similar way as done by hand but following a reproducible set of conditions. The selected blocks must fulfil certain requirements with respect to the lack of large segments of contiguous non-conserved positions, lack of gap positions and high conservation of flanking positions, making the final alignment more suitable for phylogenetic analysis. Gblocks is very fast in processing alignments and it is therefore highly suitable for large-scale phylogenetic analyses. Gblocks can be run on a computer and run on web: http://molevol.cmima.csic.es/castresana/Gblocks_server.html.

- Go to Gblocks web-server. Insert the alignment in a FASTA format in a window and select options.
- Click on **Get Block** bottom. Gblocks outputs files to visualize the selected blocks (Fig. 15.10).
- Click on **Resulting alignment** to see the result.

## 15.6   File Format Converting

Each phylogenetic program requires an alignment prepared in certain file format. Several popular phylogenetic programs such as PAUP*, MrBayes, Mesquite, MacClade and others use the NEXUS format widely used in bioinformatics. ForCon is a user-friendly software tool (http://bioinformatics.psb.ugent.be/webtools/ForCon/)



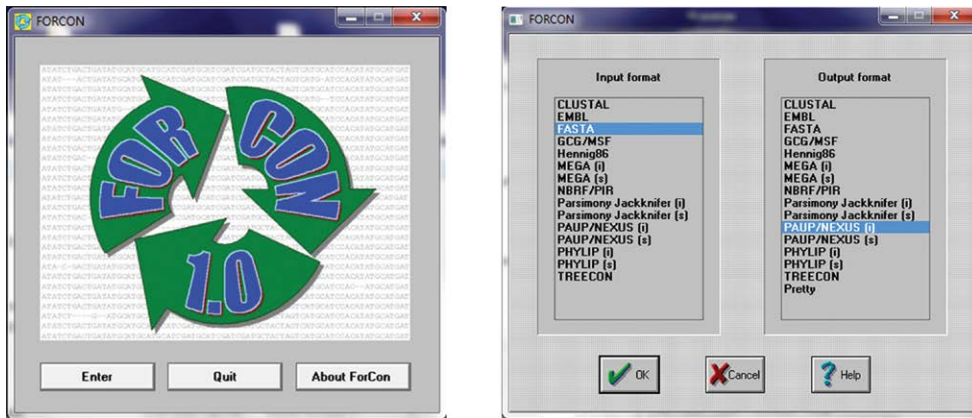**Fig. 15.10.** Gblock Server home page and output.

**Fig. 15.11.** ForCon windows.

for the conversion of nucleic acid and amino acid sequence alignments. ForCon is able to convert in both ways, i.e. reading and writing.

- Run ForCon program. Choose **Enter,** then select Input format > **FASTA** and for Output format >**PAUP/ NEXUS(i)** > OK (Fig. 15.11).
- In Open window, find your file in FASTA format > **Open > OK.** In **Save as** window, **Type** file name with file extension 'nex' > **Save**. Then choose **Select All > OK.**

The file conversion could also be done via different on-line tools, for example, http://www.ebi.ac.uk/Tools/ sfc/emboss_seqret/ and https://hcv.lanl.gov/content/sequence/FORMAT_CONVERSION/form.html.

## 15.7    Selection of Model of Sequence Evolution

When using minimum evolution (distance), maximum likelihood methods or Bayesian inference to build trees, it needs to find a model of sequence evolution that fits the DNA changes in the aligned sequences that are being used. The substitution models differ in terms of the parameters used to describe the rates at which one nucleotide replaces another during evolution. The best-fit substitution model can be selected using the Modeltest program (Posada and Crandall, 1998), MrModeltest (Nylander, 2004) or jModelTest (Darriba *et al*., 2012).

### 15.7.1    jModelTest: A tool to select the best-fit model of nucleotide substitution

jModelTest (http://code.google.com/p/jmodeltest2/) is a tool to carry out statistical selection of best-fit models of nucleotide substitution evolution that best fits the data and to use in constructing phylogenetic trees in PAUP* or MrBayes. It implements five different model selection strategies: hierarchical and dynamical likelihood ratio tests (hLRT and dLRT), Akaike and Bayesian information criteria (AIC and BIC), and a decision theory method. The jModelTest program is described by Posada (2008) and Darriba *et al*. (2012).

- Double click on jModelTest.jar to open it.
- Go **File > Load DNA alignment** and open the data set file (Fig. 15.12).
- Click on **Analysis > Compute likelihood.** A dialog box will appear that allows you to specify a number of likelihood settings, including the number of substitution schemes to be tested. Click on **Compute Likelihood** to start the analysis.
- When Likelihood score calculation is completed, click on **Analysis** and select either **Do AIC calculations** or **Do BIC calculations**. In the setting window, check a box for **Write PAUP\* block** and click on **Do AIC calculation**.
- The best-fit model is displayed with PAUP* block, which can be added to the PAUP* input file. The results can be saved with **Ctrl+S.**

**Fig. 15.12.** jModelTest windows.

The widely used General Time Reversible (GTR) model has six substitution types (**lset nst=6**), one for each pair of nucleotides and it is the most complex model. In addition to models describing the rates of change from one nucleotide to another, there are models to describe rate variation among sites in a sequence. The following are the two most commonly used models: (i) gamma distribution (G) or gamma distributed rate variation among sites proportion of invariable sites; and (ii) (I) or extent of static, unchanging sites in a dataset.

## 15.8   Phylogenetic Analysis with PAUP*

Once the sequences are aligned, there are several methods of phylogenetic analysis that can be implemented. The main methods include minimum evolution, parsimony and likelihood methods, and Bayesian inference.

Minimum evolution (distance) calculates a measure of the distance between each pair of species and then finds a tree that predicts the observed set of distances as closely as possible. A topology showing the smallest value of the sum of all branches is chosen as an estimate of the correct tree.

Maximum parsimony (MP) is a character-based method that builds a phylogenetic tree by minimizing the total tree length. It searches for the minimum number of evolutionary steps required for a given set of data.

Maximum likelihood (ML) is a statistical method for reconstructing trees. It requires three elements: (i) sequence alignment; (ii) tree topologies; and (iii) model of character evolution. ML operates by trying to maximize the likelihood value, and the tree with the highest likelihood value is considered the best tree.

These three approaches can be performed using PAUP* software (Swofford, 2003) that is available from Sinauer Associates Inc. Publishers, Sunderland, Massachusetts, USA. PAUP* (Phylogenetic Analysis Using

Parsimony* and other methods) (http://paup.phylosolutions.com) performs phylogenetic analyses using parsimony, maximum likelihood and distance methods. The program features an extensive selection of analysis options and model choices, and accommodates DNA, RNA, protein and general data types.

The Windows versions of PAUP* require data and commands, which should be typed in NEXUS format. PAUP* is run by entering the command blocks. The PAUP* manual with command explanation is given in the PDF documentation in the PAUP* folder: http://paup.phylosolutions.com/tutorials/quick-start/ The PDF command summary can be found in the **Quick Start tutorial and PAUP* FAQ: Answers** pages.

NEXUS data files always begin with the characters #nexus but are otherwise organized into major units known as blocks. The necessary commands can be put in a PAUP* block in the original nexus file after **Assumption block** in any text editor or in PAUP* using **Edit File Open Mode**. Each command begins with a command-name and ends with a semicolon. NEXUS files can contain text comments surrounded by square brackets. The following sections present several examples of the command blocks for distance, parsimony, and likelihood analyses.

### 15.8.1 Distance method

A. Command block for reconstruction of neighbour-joining tree:

```
begin paup; [start PAUP running]
log file = NJ.log; [command starts a log file]
set criterion = distance; [command defines the optimality criterion]
outgroup 8; [command specifies that the resulting trees should be rooted to given taxon]
bootstrap nreps = 1000 search = nj; [commands specify bootstrap number and method of search]
nj; [command calculates a tree using the neighbor-joining method]
showtrees; [command to request to display one or more trees]
savetrees file = NJ.tre brlens = yes root = yes; [commands save the best rooted tree
  found during the search with branch length information in a file]
log stop; [command stops the logging of output]
end; [stop PAUP running]
```

B. Command block for calculation of nucleotide differences

```
begin paup;
log file = distance.log;
dset distance = mean; [command gives mean number of pairwise character differences,
  adjusted for missing data]
BaseFreq; [command shows base frequencies for each taxon]
showdist; [command shows output a matrix of 'distances' between taxa in a PAUP window]
log stop;
end;
```

### 15.8.2 Maximum parsimony

A. Command block for reconstruction of maximum parsimony trees

```
begin paup;
log file = MP.log;
set increase = auto; [setting automatically be increased by a number of trees equal
  to the default number 100, could be changed into 'no', if search takes long time]
set autoclose; [closes the status window automatically]
set criterion = parsimony;
outgroup 8;
bootstrap nreps = 1000 search = heuristic;
hsearch nreps = 10 addseq=random; [search for optimal trees using heuristic algorithm
  with 10 replicates using random-addition-sequence replications to be performed]
showtrees;
describetrees /apolist = yes; [command produces a depiction of the tree and set a
  list of the apomorphic characters is displayed]
```

```
savetrees file = MP.tre brlens = yes root = yes;
log stop;
end;
```

B. Command block for obtaining a strict consensus tree

```
begin paup;
log file = conMP.log;
gettrees file = MP.tre; [command to load trees into memory from a file]
outgroup 8;
contree / root = outgroup; [root a strict consensus tree]
showtrees;
contree / treefile = conMP.tre; [save strict consensus file]
log stop;
end;
```

### 15.8.3   Maximum likelihood

A. Command block for reconstruction of maximum likelihood trees

```
begin paup;
log file = ML.log;
set autoclose;
set criterion = likelihood;
Lset base=(0.1943 0.2232 0.2846) nst=6 rmat=(0.9577 3.5283 1.7964 0.5168 3.5283) rates=gamma
  shape=0.7150 ncat=4 pinvar=0; [model parameters obtained from jModelTest results]
bootstrap nreps = 100 search = faststep; [bootstrap with tree searches in each replication
  are performed using one random-sequence-addition replication and no branch swapping]
hsearch;
savetrees file = ML.tre brlens = yes root = yes;
log stop;
end;
```

- Double click on PAUP icon and run the program. Change **File Open Mode** to **Edit**. Select the executable file and click on **Open** (Fig. 15.13).
- Insert command block or blocks after assumption block in the file.



**Fig. 15.13.** PAUP* with data matrix and outputs.

Phylogenetic Analysis of DNA Sequence Data

- Open **File** and then select **Execute** 'file.extension'.
- PAUP performs analysis.

PAUP generates two files: (i) a file with extension **.log** containing information on search, resulted bootstrap tree and table and optimal tree, which can be opened in any text Editor program; and (ii) a file with extension.tre with resulting tree, which can be visualized with the TreeView.

## 15.9  Phylogenetic Analysis with MrBayes

MrBayes is a program for the Bayesian estimation of phylogeny (http://mrbayes.scs.fsu.edu/). Bayesian inference (BI) of phylogeny is based upon a quantity called the posterior probability distribution of trees, which is the probability of a tree conditioned on the observations. The conditioning is accomplished using Bayes's theorem. MrBayes uses a simulation technique called Markov chain Monte Carlo (or MCMC) to approximate the posterior probabilities of trees (Ronquist and Huelsenbeck, 2003; Ronquist *et al.*, 2012).

- Open your nexus format file in any text editor program.
- Type the following command block with appropriate model (model GTR+G+I is used as an example) after the data block.
  Command block for reconstruction of BI trees:

```
begin mrbayes; [command to start MrBayes]
log start filename = BIfile.log; [command to create the log file under certain name]
outgroup 8; [command specifies that the resulting trees should be rooted to given taxon]
lset nst=6 rates=invgamma; [commands to set parameters for the model]
showmodel; [command to show model settings]mcmc ngen=1000000 printfreq=100
  samplefreq=100 nchains=4 savebrlens=yes filename=BItree; [commands to run the Markov
  chain with 1000000 generations, print the results to the screen every 100 generations,
  record the current tree and parameters values to files every 100 generations, run 4
  independent chains, save the trees with branch lengths in the tree file]
sumt filename=BItree burnin=1000 contype=halfcompat; [commands to summarize the saved
  trees, discard the first 1000 trees or 10% of total saved trees and write a majority
  rule consensus tree in BItree.con]
sump filename=BItree burnin=1000; [commands to summarize statistics for trees sampled
  during analysis]
log stop; [command to quit recording in the log file]
end;
```

  Command to set models: JC - nst = 1; K2, HKY, T3P, TN93 - nst = 2; GTR - nst = 6 and rates: Uniform - rates = equal; Gamma Distributed (G) - rates = gamma; Invariant Sites (I) - rates = propinv; G + I - rates = invgamma

- **Save** the file. Copy MrBayes in the same directory.
- Double click on MrBayes icon and run program. Type the following command: **execute** and after a space type **file name with extension**. Click > **Enter**

The program runs (Fig. 15.14). Each result line shows the generation number at the left, then a series of four numbers, each enclosed in brackets. Those numbers are the log likelihood of the trees in each of the four chain in run 1. The cold chain is enclosed in square brackets. After symbol * – chain is in run 2. The last number in each line of output to the screen is the estimated time until the run ends (Figs 15.15 and 15.16). If the average standard deviation of split frequencies is stabilized for many generations below 0.01, the run can be stopped. The results of BI analysis are saved in several files: BIfile.log, BItree.mcmc, BItree.run1.p, BItree.run2.p, BItree.run1.t, BItree.run2.t, BItree.con.

The BItree.con file can be opened in TreeView to visualize the tree.

Several other software applications can be used to reconstruct phylogenetic relationships: RAxML (https://sco.h-its.org/exelixis/web/software/raxml/) and BEAST (http://tree.bio.ed.ac.uk/software/beast) and MEGA (https://www.megasoftware.net).

**Fig. 15.14.** Output while MrBayes is running.



**Fig. 15.15.** Output of MrBayes run.

## 15.10   Visualization of Phylogenetic Trees

TreeView is a program for displaying and printing phylogenies (https://code.google.com/archive/p/ treeviewx/). TreeView provides a simple way to view the contents of a NEXUS, PHYLIP, Hennig86, Clustal, or other format tree file (Page, 1996). The PAUP* for Windows does not have a graphical interface, hence TreeView allows you to create publication quality trees from PAUP files, either directly, or by generating graphics files for editing by other programs.

- Run TreeView program. Choose **File > Open.** Then select the file in your folder. Click **Open.** The tree appears in TreeView (Fig. 15.17).
- Click on **Phylogram** icon to see the tree with length branches. Choose **Tree > Order > Select Ladderize left > OK.** Outgroup taxa appear in the bottom of the tree.
- If the tree needs to be re-rooted, select **Tree > Define outgroup** and then select Outgroup taxa. Click OK.
- The displayed tree can be saved in a graph format (.emf). Choose **File > Save as graph,** then type file name, select format. **emf** and then click **Save.**
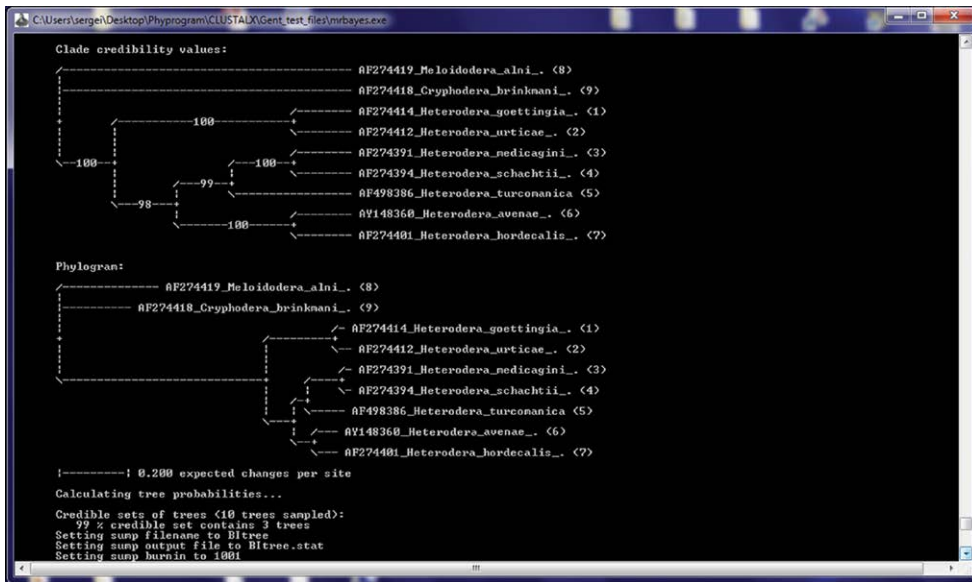
**Fig. 15.16.** Tree output of MrBayes run.



**Fig. 15.17.** Tree visualized by TreeView.

This file can be opened and edited by Adobe Illustrator, CorelDraw, Inkscape or other Graph Editors.

**FigTree** is another graphical viewer of phylogenetic trees and a program for producing publication-ready figures (http://tree.bio.ed.ac.uk/software/figtree/). In particular, it is designed to display summarized and annotated trees produced by BEAST.

**Inkscape** (https://inkscape.org) is professional quality vector graphics software that runs on Windows, Mac and Linux. It is used for creating a wide variety of graphics and can be used for final preparation of a phylogenetic tree for a publication. The program can be freely downloaded from the website.

**Fig. 15.18.** Inkscape with a tree.

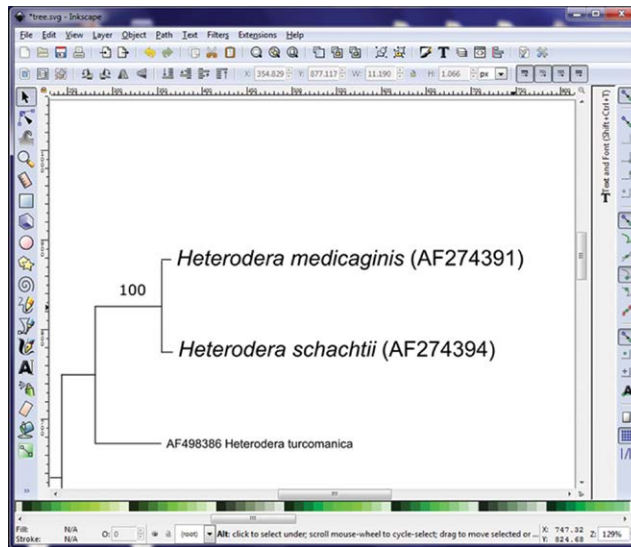The phylogenetic tree should contain clear terminal labels with species names and GenBank accession numbers and numbers with bootstrap or posterior probability values placed on appropriate branches (Fig. 15.18).

## 15.11 Sequence Submission

The final and important step of a phylogenetic study is the submission of new sequences in public databases. Once sequences are submitted and accession numbers are assigned, these numbers must be included in a publication and published tree. New sequences should be submitted in one of the public database: GenBank, a comprehensive public database of nucleotide sequences and supporting bibliographical and biological annotation built and distributed by NCBI (Benson *et al.*, 2010); EMBL (the European Molecular Biology Laboratory Nucleotide Sequence Database in Europe); or DDBJ (DNA Data Bank of Japan). Daily data exchange within these databases ensures worldwide coverage.

There are some options for submitting data to GenBank:

- **BankIt** (https://www.ncbi.nlm.nih.gov/WebSub/), a WWW-based submission tool with wizards to guide the submission process; or
- **Submission Portal** (https://submit.ncbi.nlm.nih.gov), a unified system for multiple submission types. Currently only ribosomal RNA (rRNA), rRNA-ITS, Influenza or Norovirus sequences can be submitted with the GenBank component of this tool.

## 15.12 References

Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J. (1990) Basic local alignment search tool. *Journal of Molecular Biology* 215, 403–410. DOI: 10.1016/S0022-2836(05)80360-2

Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J. and Sayers, E.W. (2010) GenBank. *Nucleic Acids Research* 38, D46–D51. DOI: 10.1093/nar/gkp1024

Castresana, J. (2000) Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Molecular Biology and Evolution* 17, 540–552. DOI: 10.1093/oxfordjournals.molbev.a026334

Darriba, D., Taboada, G.L., Doallo, R. and Posada, D. (2012) jModelTest 2: more models, new heuristics and parallel computing. *Nature Methods* 9, 772. DOI: 10.1038/nmeth.2109

Hall, B.G. (2017) *Phylogenetic Trees Made Easy*. 5th edn. Oxford University Press, Oxford, UK.

Huang, X. and Madan, A. (1999) CAP3: A DNA sequence assembly program. *Genome Research* 9, 868–877. DOI: 10.1101/gr.9.9.868

Larkin, M.A., Blackshields, G., Brown, N.P. *et al.* (2007) Clustal W and Clustal X version 2.0. *Bioinformatics* 23, 2947–2948. DOI: 10.1093/bioinformatics/btm404

Lemey, P. Salemi, M. and Vandmme, A.-M. (eds) (2009) *The Phylogenetic Handbook: A Practical Approach to Phylogenetic Analysis and Hypothesis Testing*. 2nd edn. Cambridge University Press, Cambridge, UK. DOI: 10.1111/j.1541-0420.2010.01388.x

Nicholas, K.B., Nicholas Jr, H.B. and Deerfield II, D.W. (1997) GeneDoc: analysis and visualization of genetic variation. *EMBNEW News* 4, 14.

Nylander, J.A.A. (2004) *MrModeltest Version 2. Program distributed by the author*. Evolutionary Biology Centre, Uppsala University, Uppsala, Sweden.

Page, R.D.M. (1996) TREEVIEW: An application to display phylogenetic trees on personal computers. *Computer Applications in the Biosciences* 12, 357–358.

Posada, D. (2008) jModelTest: Phylogenetic model averaging. *Molecular Biology and Evolution* 25, 1253–1256. DOI: 10.1093/molbev/msn083

Posada, D. and Crandall, K.A. (1998) Modeltest: testing the model of DNA substitution. *Bioinformatics* 14, 817–818. DOI: 10.1093/bioinformatics/14.9.817

Ronquist, F. and Huelsenbeck. J.P. (2003) MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics* 19, 1572–1574. DOI: 10.1093/bioinformatics/btg180

Ronquist, F., Teslenko, M., van der Mark, P. *et al.* (2012) MrBayes3.2: Efficient Bayesian phylogenetic inference and model choice across a large model space. *Systematic Biology* 61, 532–542. DOI: 10.1093/sysbio/sys029

Swofford, D.L. (2003) PAUP*: phylogenetic analysis using parsimony (*and other methods), version 4.0b 10. Sinauer Associates, Sunderland, Massachusetts, USA.

Talavera, G. and Castresana, J. (2007) Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments. *Systematic Biology* 56, 564–577. DOI: 10.1080/10635150701472164

Wheeler, D. and Bhagwat, M (2007) Chapter 9. BLAST QuickStart In: Bergman N.H. (ed.). *Methods in Molecular Biology, Vol. 395: Comparative Genomics*. Humana Press Inc., Totowa, New Jersey, USA, pp. 149–175.

Yang, Z. and Rannala, B. (2012) Molecular phylogenetics: principles and practice. *Nature Review: Genetics* 13, 303–314. DOI: 10.1038/nrg3186