

POY 5.0 R.C.

Program Documentation

Version 5.0.Black Sabbath Development build 2a12624f96c4

Program and Documentation

Andrés Varón
Nicholas Lucaroni
Lin Hong
Ward C. Wheeler

Documentation

Louise M. Crowley
Megan Cevasco
John S. S. Denton

Previous Version POY 4

Program and Documentation

Andrés Varón

Le Sy Vinh

Illya Bomash

Ward C. Wheeler

Documentation

Ilya Tëmkin

Megan Cevasco

Kurt M. Pickett

Julián Faivovich

Taran Grant

William Leo Smith

Andrés Varón

Jane Street Capitol, 1 New York Plaza, New York, NY, U.S.A.

Louise M. Crowley, Lin Hong, Nicholas Lucaroni, Ward C. Wheeler

Division of Invertebrate Zoology, American Museum of Natural History, New York, NY, U.S.A.

John S. S. Denton

Richard Gilder Graduate School and Department of Ichthyology, American Museum of Natural History, New York, NY, U.S.A.

Megan Cevasco

Coastal Carolina University, Department of Biology, Conway, SC, U.S.A.

Ilya Bomash

Department of Physiology And Biophysics, Weill Medical College of Cornell University, New York, NY, U.S.A.

Julián Faivovich

División Herpetología, Museo Argentino de Ciencias Naturales - CONICET, Buenos Aires, Argentina.

Taran Grant

Universidade de São Paulo, Instituto de Biociências, Departamento de Zoologia, Cidade Universitária, São Paulo, Brasil.

Kurt M. Pickett

Department of Biology, University of Vermont, Burlington, VT, U.S.A.

William Leo Smith

Department of Zoology, The Field Museum of Natural History, Chicago, IL, U.S.A.

Ilya Témkin

Northern Virginia Community College, Annandale Campus, VA, U.S.A.

Le Sy Vinh

College of Technology and Information Technology Institute, Vietnam National University, Hanoi, Vietnam.

The American Museum of Natural History

©2011 by The American Museum of Natural History,

All rights reserved. Published 2011.

Varón, A., N. Lucaroni, L. Hong, W. C. Wheeler. 2011. POY 5.0.Black Sabbath Development build 2a12624f96c4 R.C. New York, American Museum of Natural History. L. M. Crowley, M. Cevasco, J. S. S. Denton.

<http://research.amnh.org/scicomp/projects/poy.php>

Available online at <http://research.amnh.org/scicomp/projects/poy.php> and <http://code.google.com/p/poy4/> new link needed

Contents

1	What is POY5	9
1.1	The structure of POY5 documentation	9
1.2	What's new in POY5	10
2	POY5 Quick Start	11
2.1	Requirements: software and hardware	11
2.1.1	Software	11
2.1.2	Hardware	11
2.2	Obtaining and installing POY5	11
2.2.1	Installing from the binaries	11
2.2.2	Compiling from the source	13
2.3	The Graphical User Interface	14
2.3.1	POY menu bar	14
2.3.2	POY Launcher	15
2.3.3	The <i>Analyses</i> menu	17
2.3.4	The <i>View</i> menu	34
2.4	Using the Interactive Console	35
2.4.1	The interface	36
2.4.2	Starting a POY5 session using the <i>Interactive Console</i> .	36
2.4.3	Entering commands	39
2.4.4	Browsing the output	39
2.4.5	Switching between the windows	40
2.4.6	Importing data	40
2.4.7	Inspecting data	42
2.4.8	Building the initial trees	45
2.4.9	Performing a local search	46
2.4.10	Selecting trees	48
2.4.11	Visualizing the results	49
2.4.12	Interrupting a process	49

2.4.13 Reporting errors	51
2.4.14 Exiting	51
2.5 Creating and running P0Y5 scripts	51
2.6 Obtaining help	53
2.7 WWW resources	54
3 P0Y5 Commands	57
3.1 P0Y5 command structure	57
3.1.1 Brief description	57
3.1.2 Grammar specification	58
3.2 Notation	60
3.3 Command reference	62
3.3.1 build	62
3.3.2 calculate_support	65
3.3.3 clear_memory	68
3.3.4 cd	69
3.3.5 echo	70
3.3.6 exit	71
3.3.7 fuse	71
3.3.8 help	73
3.3.9 inspect	74
3.3.10 load	75
3.3.11 perturb	76
3.3.12 pwd	79
3.3.13 quit	79
3.3.14 read	80
3.3.15 rediagnose	87
3.3.16 recover	88
3.3.17 redraw	88
3.3.18 rename	89
3.3.19 report	91
3.3.20 run	102
3.3.21 save	103
3.3.22 search	104
3.3.23 select	107
3.3.24 set	111
3.3.25 store	114
3.3.26 swap	115
3.3.27 transform	120
3.3.28 use	140

CONTENTS	7
-----------------	----------

3.3.29 version	141
3.3.30 wipe	142
4 POY5 Heuristics: A Practical Guide	143
4.1 Introduction	143
4.2 Data treatment	144
4.3 Character optimization	145
4.4 Tree searching	146
4.5 Chromosome heuristics	147
4.5.1 Breakin character	147
4.5.2 Annotated chromosome character	148
4.5.3 Chromosome character	149
4.6 Transformation cost regimes	151
5 POY5 Tutorials	153
5.1 Combining search strategies	153
5.2 Searching under iterative pass	156
5.3 Bremer support	158
5.4 Jackknife support	159
5.5 Sensitivity analysis	161
5.6 Chromosome analysis: unannotated sequences	164
5.7 Chromosome analysis: annotated sequences	166
5.8 Custom alphabet and break inversion characters	167
5.9 Genome analysis: multiple chromosomes	169
Bibliography	171
General Index	178
POY 3.0 Commands Index	179

Chapter 1

What is POY5

POY5 is a flexible, multi-platform program for the phylogenetic analysis of a diversity of data types. An essential feature of POY5 is that it implements the concept of dynamic homology [51, 52] allowing optimization of unaligned sequences. POY5 offers flexibility for designing heuristic search strategies and implements an array of algorithms including multiple random addition sequence, swapping, tree fusing, tree drifting, and ratcheting. As output, POY5 generates a comprehensive character diagnosis, graphical representations of cladograms and their user-specified consensus, support values and implied alignments. In addition, POY5 can also output synteny block maps from the analysis of both chromosomal and genomic data. POY5 provides a unified approach to co-optimizing different types of data, such as morphological and molecular sequence data. In addition, POY5 can analyze entire chromosomes and genomes, taking into account large-scale genomic events (translocations, inversions, and duplications).

1.1 The structure of POY5 documentation

Chapter 2, *POY5 Quick Start*, will get you started using POY5. The first few sections are intended to provide detailed instructions on how to obtain and install POY5, introduce the user to two of the program's working environments, the *Graphical User Interface* and the *Interactive Console*. These sections also show how to initiate a POY5 session and point to the various resources to obtain further assistance. Subsequent sections build on that knowledge and give step-by-step examples on how to conduct a basic analysis and visualize the results. The following chapter, *POY5 Commands*, describes POY5 commands and their valid syntax. It also includes examples of simple

operations for every command. Chapter 4 discusses the heuristic procedures used in POY5. Their understanding helps creating building efficient search strategies. More advanced operations are described in the fifth chapter, *POY5 Tutorials*. In addition to the general index, this document contains a *POY3.0 Command Line Index*, intended to provide a link between the commands used in POY3 and the commands used in POY5.

1.2 What's new in POY5

New Features

New Commands

Major Bug fixes?

Improvements

New Compilation requirements?

Chapter 2

POY5 Quick Start

2.1 Requirements: software and hardware

2.1.1 Software

POY5 is a platform-independent, open-source program that can be compiled for many operating systems and hardware configurations, including Mac OSX, Microsoft Windows XP, and Linux. The intuitive *Graphical User Interface* of POY5 provides all the functionality for running analyses using pull-down menus and field selections, as well as creating and running POY5 scripts. Some utility programs (such as Notepad and Ghostscript for Windows,TextEdit for Mac, or Nano for Linux), can help preparing POY5 scripts and formatting data files, while others (such as Adobe Acrobat) can facilitate viewing the graphical output in PDF (Portable Document Format).

2.1.2 Hardware

POY5 runs on a variety of computers from laptops and desktops to Beowulf clusters of various sizes to symmetric multiprocessing hardware. There are no particular requirements for disk space, but XML and diagnosis report files can be large.

2.2 Obtaining and installing POY5

2.2.1 Installing from the binaries

POY5 installers for Linux, Windows XP, and Mac OSX, source code, and documentation in PDF format are available from the POY5 Download website

at the American Museum of Natural History Computational Sciences:

<http://research.amnh.org/scicomp/projects/poy.php>

The latest source code can also be obtained from POY5 Google Group website:

<http://code.google.com/p/poy4/source> new link

The following detailed step-by-step instructions will guide you through downloading and installing POY5 *binaries* for various platforms.



Windows

- Download the *poy5* folder to the desktop by selecting the *Windows* download link.
- Open the *POY_Installer.exe* from the *poy5* folder and follow the installation instructions. You will need Administrator privileges to install the application. By default, the Windows version of POY5 is installed without parallel execution support. If you have Windows XP SP2, Windows Vista or Windows 7 and possess more than one core or processor, you can take advantage of this processing power by installing the parallel components **MICH2**. Note: The developers encountered problems when using MICH2 1.4.1.p1, therefore it is advised that the user install MICH2 1.4, where no problems were encountered.



Mac OSX

- Download *poy-buildXXXX.dmg* disk image to the desktop. The complete installation of the Mac OSX version of POY5 includes MPICH2 1.4.1, which is used to communicate processes during parallel execution.
- Drag the POY5 application from the disk *poy5* and drop it into the *Applications* folder on the hard drive. Note: During the first execution in parallel you will be asked by the Firewall to unblock POY5 and MPICH, this is necessary for the successful execution of the program.



Linux

- No binaries are available for the Linux operating system. The user should compile POY5 directly from the source (see Section 2.2.2).

2.2.2 Compiling from the source

For the majority of users, downloading the binaries from the POY5 download site will suffice. However, in some cases it may be desirable—user preference for working in a command-line environment or running POY5 analyses in parallel (in the case of Linux machines or on a cluster computer)—to compile POY5 directly from the source code (see Table 2.1). If the user chooses to compile, it is advisable to check out the various configuration options that can be found in `./configure --help`.

In order to compile POY5 the following tools are required:

1. The [GNU Make 3.8](#) utility.
2. OCaml [version 3.11.2](#). or later.
3. A C compiler, for example [The GNU Compiler Collection](#).
4. [The zlib compression library](#).
5. The Linear Algebra PACKage [LAPACK](#) must be installed in order to use the likelihood option.
6. The [ncurses library](#) is necessary to compile the default interface, i.e. *ncurses* or the interactive console. If this library is not available, the *flat* interface will be compiled instead.
7. The Message Passing Interface [MICH2](#), which is used to communicate processes during parallel execution.

Download, ungzip, and untar the [POY5 source code](#). In a terminal window, change directories to the path of this uncompressed directory. In order to compile under default setting type:

```
./configure
make
make install
```

This script will compile both the *Graphical User Interface* and *Interactive Console* or *ncurses* interface. Another configuration option includes a *readline* interface. Similar to the *ncurses* interface, this allows for the use of arrow keys to modify commands and browse command history. An *html* interface is available. A *flat* interface is also available that supports the running of the

program in parallel irrespective of the operating system. This version is run through a terminal window. In order to run POY5 in parallel environments, [Message Passing Interface](#) must be invoked. More than likely, your system administrator already has one installed and should be able to provide you with the proper paths to set your config file. In order to compile this *flat* version with parallel support type:

```
./configure --enable-interface=flat --enable-mpi CC=mpicc
make
make install
```

[Note: CC=mpicc is not available for the Windows version of mpi, therefore it is not necessary to include this component in the compiling script.]

2.3 The Graphical User Interface

Two of the working environments that POY5 provides are the *Graphical User Interface* and the *Interactive Console* (also known as *nurses* interface). The *Graphical User Interface* has a user-friendly appearance like other native stand-alone applications where different functions are accessible through menus and windows. Thus, the entire analysis can be carried out by clicking on appropriate selections and, where necessary, typing specifications in designated fields. Currently, the *Graphical User Interface* is appropriate for the analysis of data with parsimony as the optimality criterion or with likelihood, with minimal model selection.

On the other hand, the *Interactive Console* requires a detailed knowledge of POY5 commands, their arguments, and the conventions of POY5 scripting. All these features are described in the *POY Commands* chapter ([3.1.1](#)).

Even though the Mac OSX version of the *Graphical User Interface* is used for screen shots throughout this chapter, the Windows and Linux(GTK+) versions contain the same items and functionality, differing only in the generic window format specific to each platform.

When POY5 is first opened, two items appear on the screen: the POY5 menu bar across the top and the *POY Launcher* window (Figure [2.1](#)). (Note that in Linux and Windows the menu bar is within the launcher window.)

2.3.1 POY menu bar

The menu bar contains the following drop-down menus:

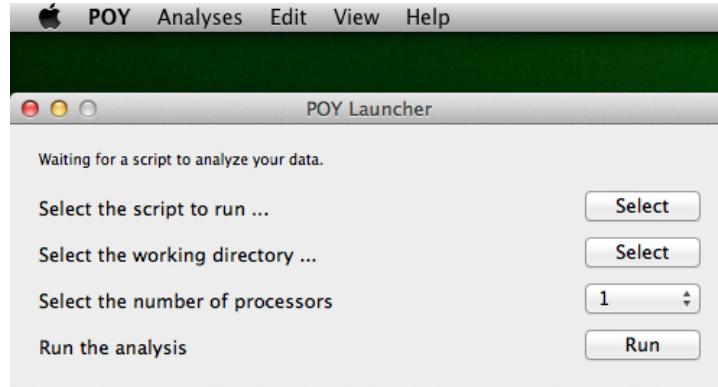


Figure 2.1: The **POY5** menu bar and the *POY Launcher* window. These items appear when **POY5** is opened.

POY (Mac OSX only) contains generic items as with other Mac OSX applications. This pull down menu allows selection of the *About POY* window (Figure 2.2) that lists the current version of **POY5**, a copy right statement, and the address of the **POY5** website. In addition, it includes a *Quit POY* tab that closes the program.

Analyses contains options for different types of tree searches, calculation of support values, tree diagnosis, and their respective outputs. Other items in this menu open the *POY Launcher* (Open Launcher) and the *Interactive Console*.

Edit contains standard tools for undoing, cutting, copying, pasting, deleting, and selecting.

View opens the *Output* window to display the results (including warning and error messages) and the current state of the analysis. This *Output* window also contains an *update* tab. It also contains the *About POY* menu item in Windows and Linux.

Help opens the **POY5 Manual** in PDF format (requires a PDF viewer).

2.3.2 POY Launcher

The *POY Launcher* is the only window that automatically opens upon starting **POY5**. This allows the user to import a previously created script,

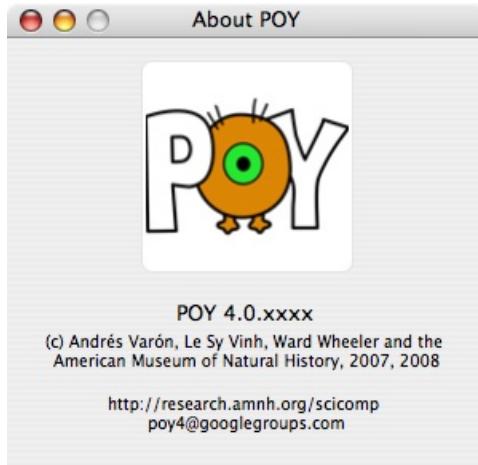


Figure 2.2: The *About POY* window.

designate a working directory, specify the number of processors, and start the analysis.

Select the script to run Allows the user to specify the location of a POY5 script.

Select the working directory The working directory is the directory that contains the input data and output files. By default, the working directory is set to be the same as the directory containing the selected POY5 script.

Select the number of processors If more than one computer processor is available, up to 8 cores can be designated for running the analysis. It is important to note that once specified, the selection is applied to *all* subsequent analyses in the current POY5 session. The selection of the number of processors is disabled for Linux platforms. Table 2.1 is a guide to the parallelization ability of POY5 depending on the operating system and the POY5 interface being used. Note that parallelization is *never* supported in interactive sessions, see Section 2.4.

Run the analysis Clicking the *Run* button starts the execution of the selected script. Once the script is initiated, the *Run* button becomes the *Cancel* button that can be used to interrupt a POY5 session.

Table 2.1: Parallelization Guide. The field 'mpi+' indicates that mpi must be enabled during compiling. A distinction is made between the interactive console or ncurses that is downloaded from the website (binaries) or that which is compiled from the source.

Operating System	GUI	Interactive Console (binaries)	Interactive Console (source)	Flat
Mac OSX	+	-	-	mpi+
Windows	+	-	-	mpi+
Linux	-	-	-	mpi+

If the *Run* button is clicked without the selected script and working directory, or the names of the scripts and working directory are entered incorrectly, POY5 issues an error message in the upper part of the *POY Launcher* window, such as **POY finished with an error**.

2.3.3 The *Analyses* menu

The *Analyses* menu is the main toolbox of the POY5 GUI interface (Figure 2.3, left). Selections are subdivided into four functional categories. The first three deal with tree searching, support calculation, and tree diagnosis; the fourth one is used for script management or interactive command execution that bypasses the menu-driven script generating. Each of the menu items is described below in the order it appears on the menu.

Most options are consistently applied through different kinds of analysis. Therefore, all options are described in detail only for the *Simple Search* analysis. The descriptions of other analyses are made with reference to the the *Simple Search* and focus on unique options.

Tree searching options

A number of different tree searching options are available through the *Graphical User Interface*. These include a *Simple Search*, *Timed Search*, *Search with Ratchet* and *Search with Perturb*.

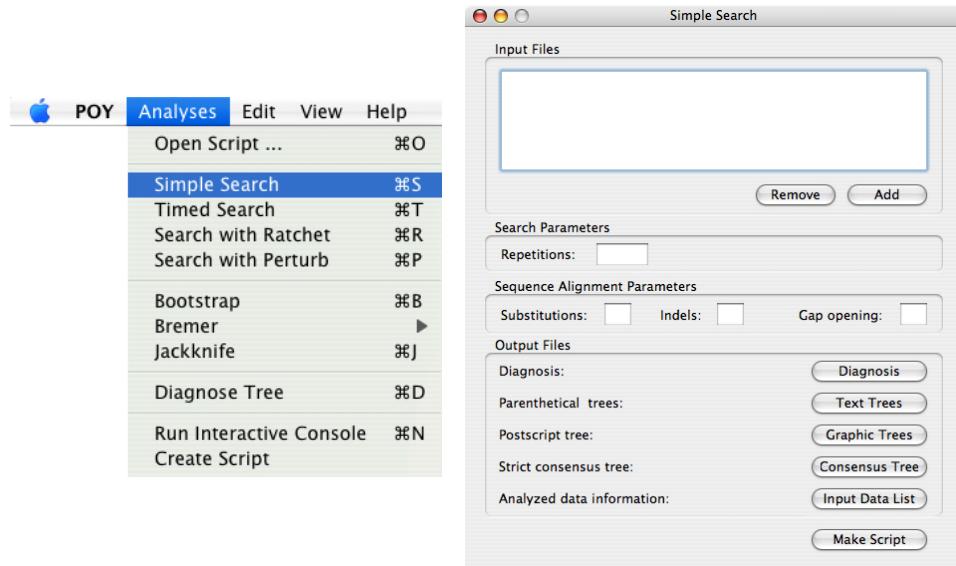


Figure 2.3: The *Simple Search* window. Selecting *Simple Search* from the *Analysis* menu (left) opens the *Simple Search* window options (right).

Simple Search

The *Simple Search* window permits the analysis of a number of different data types, including a range of molecular characters (from DNA sequences to complete genomes), custom alphabet characters, and qualitative characters, under parsimony. It is also possible to carry out a likelihood analysis of DNA sequence data under a number of different models. In the most simplest sense, a typical search involves a series of steps. First, initial trees are generated by random addition sequence from the imported character data. These trees are then subjected to branch swapping, after which trees are selected to report. The *Simple Search* window (Figure 2.3, right) provides the most common and basic options for a standard tree search in POY5 that must be selected by either clicking the appropriate buttons or by typing. Note that *all* the empty fields *must* be filled in, otherwise the default values will be used. The window is subdivided into five sections:

Input Files Contains the list of files that are to be input into POY5. These include character files in nucleotide, Hennig86, and Nexus formats, as well as tree files. (Character data in other formats can be input by specifying additional arguments in the script. See **read** (Section 3.3.14).) Note: If a Nexus file is input, a gap-opening cost can not be specified,

as gaps are treated as independent in this file format.

Search Parameters Holds one field to set the number of independent random addition replicates to be generated.

Input Parameters Holds fields to specify the optimality criterion (parsimony or likelihood), the datatype (sequence, chromosome, genome, custom alphabet, break inversion or qualitative) and whether these data types should be treated as prealigned (if possible). It is important to note, that in the *Sequence Parameters* section for a parsimony analysis below, the user may be required to specify the cost associated with a *Locus Breakpoint*, *Locus Inversion* or *Locus Indel*, depending on the data type. In these cases, *Locus* should not be taken to be the functional, biological unit in the classical sense, but only as a homologous segment of a sequence of characters.

Currently, it is only possible to perform a likelihood analysis of DNA sequences and not of the other character types. These sequences can be prealigned. A description of the *Likelihood Model Parameters* follows the *Sequence Parameters* required for a parsimony analysis below.

Sequence Parameters Holds fields to specify the substitution, indel, and gap opening costs of sequences. Enter 0 if no gap opening cost is desired. If the value of a parameter is not specified, the default values are used. The default value for both substitutions and insertion deletion events is 1 and that of gap-opening is 0. Within Sequence Parameters, an additional subsection is invoked depending on the character type (except for *Sequence* characters, where only indel, substitution and gap opening costs are involved). For more detailed explanations of the different data types see below and `transform`.

Chromosome Parameters *Chromosome* characters are multi-locus nucleotide sequences and can include nuclear chromosomes, as well as, mitochondrial and viral genomes. It is possible to submit either *Annotated* or *Unannotated* chromosomes. Within *Annotated* chromosomes, homologous regions, such as loci, are separated with the pipe symbol “|”. *Unannotated* chromosomes are entirely without delimiters. For *Unannotated* chromosomes, the Mauve Parameters, must be set by the user. These parameters (*Max. Match*, *Min. Match*, *Match Quality* and *Match Coverage*) are employed by the Mauve aligner to find regions of homology or synteny blocks between chromosomes (see the

argument `annotate` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`). Default values for these parameters are provided in the *GUI*.

Within this subsection it is necessary to specify both *Locus Indel* and rearrangement (*Locus Breakpoint* or *Locus Inversion*) costs. The cost of a *Locus Indel* is by default set to 10 plus 0.9 times the length of the locus (see the argument `locus_indel` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`). Rearrangements of homologous regions (as defined by the user in the case of *Annotated* chromosomes or as determined by the Mauve aligner as in *Unannotated* chromosomes) are then optimized using either *Locus Breakpoint* or *Locus Inversion* costs (see the arguments `locus_breakpoint` (Section 3.3.27) and `locus_inversion` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`). The default cost for both is set to 10. The user must also specify the *Median solver* for the optimization of rearrangements. The default median solver is *Caprara*, but the user can alternatively choose *BBTSP*, *Coaletsp*, *Chainlk*, *Siepel*, *Simpllk* and *Vinh* (see the argument `median_solver` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`).

The user must also specify whether the chromosome is *Circular* (true) or linear (false). Due to restrictions placed by the provided median solvers, it is not possible to submit pre-aligned data files for both *Annotated* and *Unannotated* chromosomes.

Genome Parameters *Genome* characters are multi-locus, multi-chromosomal nucleotide sequences, wherein transformations (i.e. indels, substitutions, and rearrangements) are optimized at the sequence, locus and chromosomal level. Within the genome data file, individual chromosomes are separated by the “@” symbol and the individual chromosomes remain *Unannotated*.

As with *Unannotated Chromosome* characters, homologous regions are determined using the Mauve parameters *Max. Match*, *Min. Match*, *Match Quality* and *Match Coverage*. Default values for these parameters are provided in the *GUI*. The *Locus Indel* and rearrangement (*Locus Breakpoint* or *Locus Inversion*) costs are set by the user. The cost of a *Locus Indel* is by default set to 10 plus 0.9 times the length of the locus (see the argument `locus_indel` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`).

Rearrangements of homologous regions, as determined by the Mauve aligner, are then optimized using either *Locus Breakpoint* or *Locus Inversion* costs (see the arguments `locus_breakpoint` (Section 3.3.27) and `locus_inversion` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`). The default cost for both is set to 10. A *Median solver* must also be specified for the optimization of rearrangements. The default median solver is *Caprara*, but the user can alternatively choose *BBTSP*, *Coaletp*, *Chainlk*, *Siepel*, *Simpllk* and *Vinh* (see the argument `median_solver` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`).

Two other costs must be set for the analysis of this data type—*Translocation* events and *Chromosome Indel*. The cost of the *Translocation* of a region of one chromosome to another chromosome is set by default to 10. The cost of the insertion or deletion of an entire chromosome is by default set to 10 plus 0.9 times the length of the chromosome.

As with *Chromosome* characters, it is not possible to input pre-aligned data files, due to restrictions placed by the provided median solvers.

Custom Alphabet Parameters *Custom Alphabet* characters are those that are defined by the user, and come in alphabet format. With this data type, only insertion-deletion events are allowed. *Custom Alphabet* characters can be input as prealigned. Within this subsection, the user must specify the heuristic *Level* of the median sequence calculation. *Direct Optimization* is involved in median calculation in this case. Because calculating the median states between custom alphabet strings becomes more computationally intensive (and time consuming) as the number of elements in the alphabet increases, the user should select a heuristic level of median calculation appropriate for their data. The default level is set to 2.

In addition to the data file, the user is require to upload a *Cost Matrix* that specifies the substitution and indel transformation costs for alphabet elements. By selecting the *Cost Matrix* button within this subsection, the user can upload a cost matrix that specifies these costs for their custom alphabet data. For details on the format requirements for custom alphabet data files and their associated cost matrices see the argument `custom_alphabet` (Section 3.3.14) within the command

`read`).

Break Inversion Parameters *Break Inversion* characters are an enhancement of *Custom Alphabet* characters, in which element rearrangements are calculated, as well as substitutions and indels. The median solvers provided restrict the analysis of prealigned data. The rearrangement costs for *Break Inversion* characters can be optimized using either *Locus Breakpoint* or *Locus Inversion* approaches (see the arguments `locus_breakpoint` (Section 3.3.27) and `locus_inversion` (Section 3.3.27) within the `chromosomal transformation` methods of the command `transform`). The default cost for both is set to 10. A *Median solver* must also be specified for the optimization of rearrangements. The default median solver is *Caprara*, but the user can alternatively choose *BBTSP*, *CoaletsP*, *Chainlk*, *Siepel*, *Simpllk* and *Vinh* (see the argument `median_solver` (Section 3.3.27) within the chromosomal transformation methods of the command `transform`). Linked to the median solver, the *Level* must also be specified. The *Level* specifies the heuristic level of median sequence calculation. The calculation of median states between *Break Inversion* strings becomes more computationally intensive (and time consuming) as the number of elements in the alphabet increases, the user should select a heuristic level of median calculation appropriate for their data. The default level is set to 2.

The requirements for *Break Inversion* character types are identical to those for *Custom Alphabet* characters, with respect to substitution and indel transformation costs. By selecting the Cost Matrix button within this subsection, the user can upload a cost matrix to specify these costs. The user should see the argument `custom_alphabet` (Section 3.3.14) within the command `read` for details on the format requirements for these cost matrices, which are identical in form to those for *Custom Alphabet* characters.

Qualitative Parameters *Qualitative* data are any non-sequence, pre-aligned data type (e.g. morphological, behavioral). These character types are optimized as additive, non-additive or Sankoff characters and this information should be included in the data file.

Likelihood Model Parameters Holds fields to specify the *likelihood model* of character substitution under which the analysis will be performed. Available models include JC69/Neyman, F81, K2P/K80, F84, HKY85, TN93 and GTR.

Rate Distribution specifies the nature of among-site variation. Rate variation distributions allow multipliers to be applied to separate groups of characters. These distributions can be set as to an either *Constant*, *Gamma* (for non-zero rate variation) or *Theta* (for parameterization of invariant sites).

Rate Classes it has therefore been recommended that a proportion of invariant sites be “pseudo”-estimated by increasing the number of discrete categories for the gamma distribution to account for very low rate character groups. The mean of each class is used to define the rate for the discrete category.

Gap Treatment

Cost Method

Priors The equilibrium frequencies of the characters in the stationary Markov process. The conditional probability of the data is multiplied by these frequencies to yield the likelihood score. The observed (empirical) frequencies of the characters in the data are often used as an approximation to the likelihood-optimized frequencies. This approximation is sensible for sequences evolving nearly neutrally, but this approximation breaks down for sequences under moderate or strong selection.

Output Files Designates the names and locations of files containing different kinds of results of the analysis. By default, all of these output options are generated, with default names applied. The names can then be changed in the generated script or the option can be removed entirely. As implied by their respective titles, the *tree* buttons output trees in both parenthetical (most parsimonious and consensus trees) and postscript form.

A diagnosis file provides information relating to the analysis. Information in this file includes, the cost of the tree, the rearrangement costs (in the case of the analysis of chromosome, genome and break inversion data types), as well as information about each resulting node in the tree. At each node, the user with a cost of the tree down to that node, a rearrangement cost (if applicable), the ‘children’ coming from this node

Once all the parameters are selected, click the *Make Script* button and another window—the *Script Editor*—containing the generated script, appears on screen (Figure 2.4). The script can be edited by typing in the commands

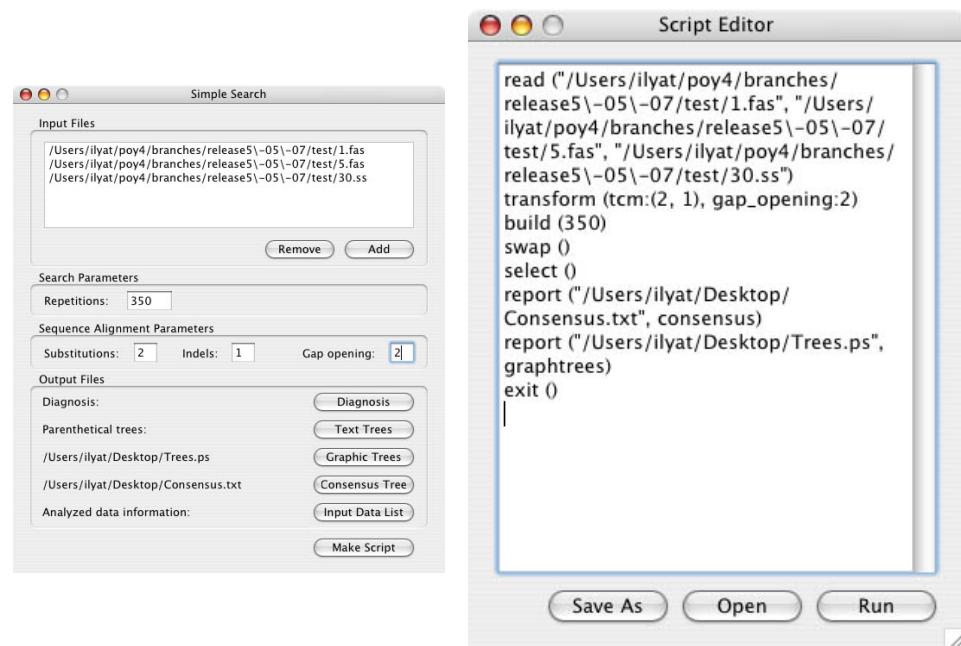


Figure 2.4: The *Simple Search* window with specified search parameters and output files (left) and the corresponding *Script Editor* window.

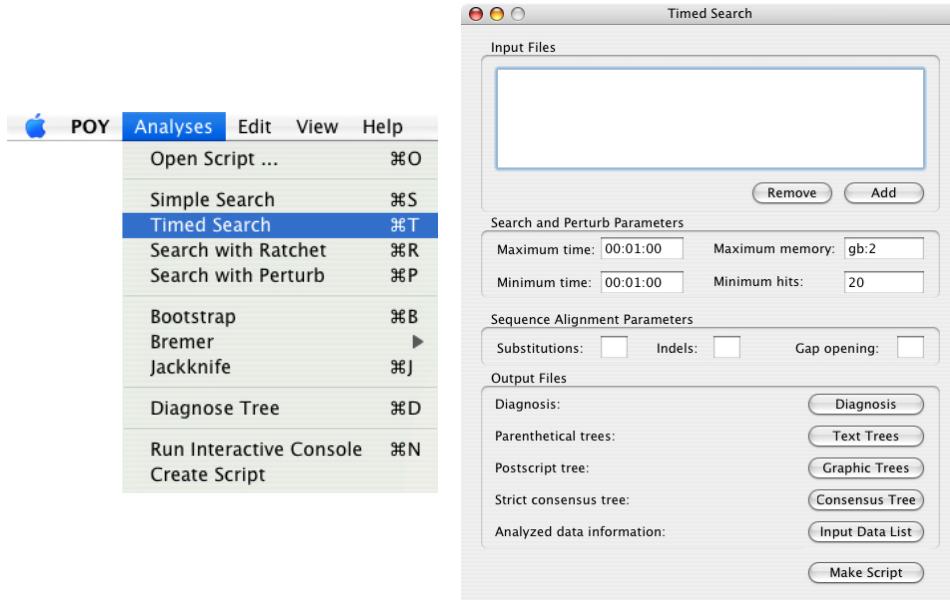


Figure 2.5: The *Timed Search* window. Selecting *Timed Search* from the *Analysis* menu (left) and viewing the *Timed Search* window options (right).

directly in the *Script Editor* window, saved (by clicking the *Save As* button), or replaced with another script (using the *Open* button). To start the analysis, click the *Run* button in the *Script Editor* window. When the *Run* button is clicked, POY5 will issue a request to save the script. Thus, not only does POY5 execute the script but it also creates the record of the type of analysis (including all user-defined specifications) that was performed.

Timed Search

Timed Search (Figure 2.5) implements a default search strategy that effectively combines tree building with TBR branch swapping, parsimony ratchet, and tree fusing. The *Timed Search* window has the same four parameter groups described for the *Simple Search*. However, the *Search Parameters* section (called *Search and Perturb Parameters*) contains four fields specifying the search targets instead of the *Repetitions* field. These include the following:

Maximum time The maximum total execution time for the search. The time is specified as days:hours:minutes.

Minimum time The minimum total execution time for the search. The

time is specified as days:hours:minutes.

Maximum memory The maximum amount of memory allocated for the search.

Minimum hits The minimum number of times that the minimum cost must be reached before aborting the search.

This heuristic search is a powerful tool for analyzing data. The number of rounds of successive searching is limited only by the previously mentioned search targets. Therefore, when performing a *Timed Search*, it is crucial to set the maximum time such that the program has a reasonable amount of time to perform a search. Thus, it is important to have some approximation as to the length of time it would take to perform a single round of searching (e.g. build (1), followed by TBR, ratchet and fusing in the case of a parsimony analysis of DNA sequence data). Clearly, this is data and optimality criterion dependent. With this information at hand, the user can then estimate the amount of time necessary to perform a thorough search. (The user should also allow some time for the program to collate and write the results to files. If the user opted to run this analysis in parallel, this can take a considerable amount of time.)

Search with Ratchet

The parsimony ratchet is a heuristic strategy to escape the local optima during tree searching [34]. The ratchet reweights a given percentage of characters for a specified number of iterations of a search. An analysis is then performed and the resulting tree topology is evaluated using the original data matrix with all characters equally weighted to determine the length of the tree. The *Search with Ratchet* (Figure 2.6) follows the same basic steps of a simple search but includes the ratchet step after the swap. In addition to the same sequence alignment and search parameters as described for the *Simple Search* window, the *Search Parameters* section provides the following ratchet parameters fields:

Ratchet iterations The number of iterations for the parsimony ratchet.

Severity The severity parameter of the parsimony ratchet (the weight change factor for the selected characters).

Percentage The percentage of characters to be reweighted during ratcheting.

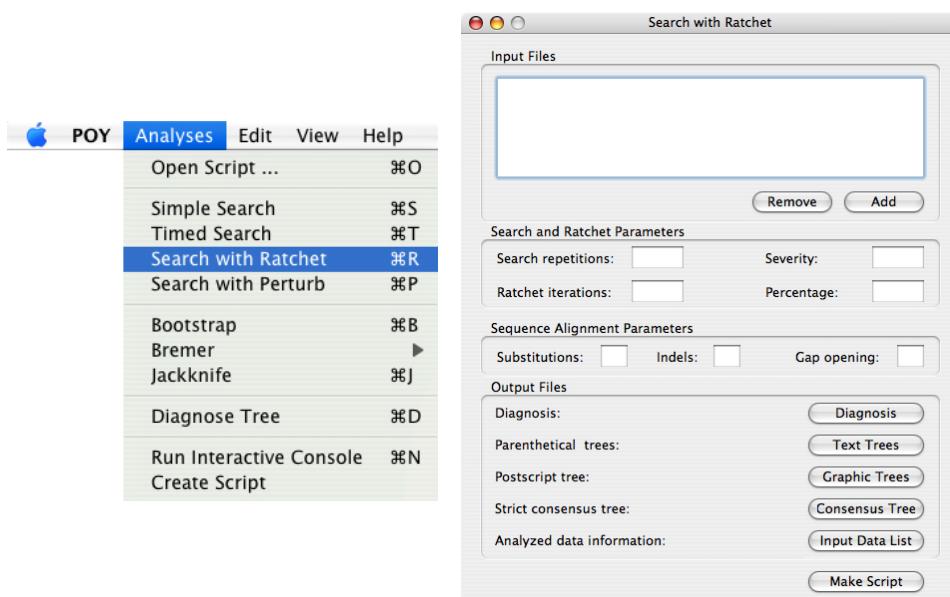


Figure 2.6: The *Search with Ratchet* window. Selecting *Search with Ratchet* from the *Analysis* menu (left) and viewing the *Search with Ratchet* window options (right).

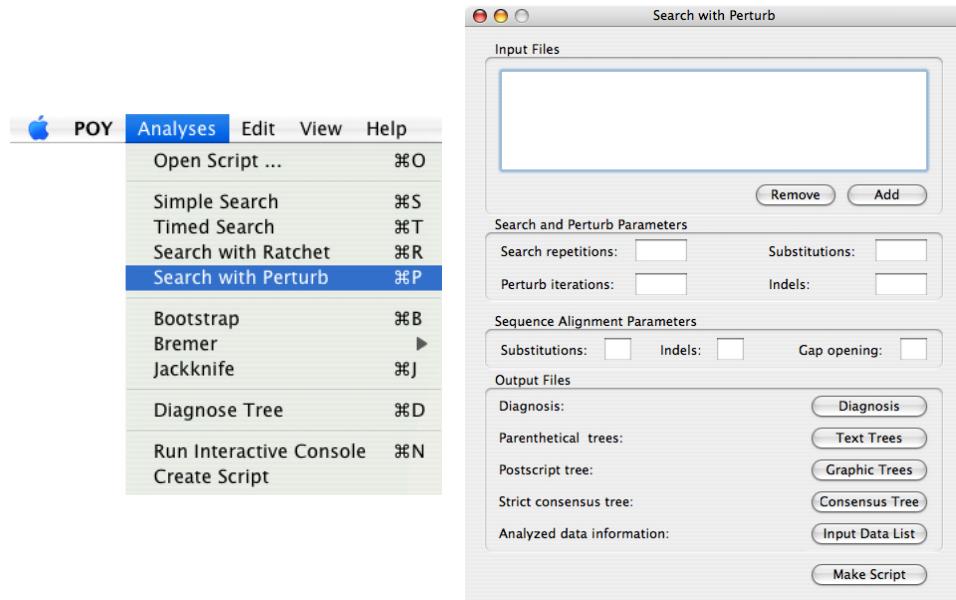


Figure 2.7: The *Search with Perturb* window. Selecting *Search with Perturb* from the *Analysis* menu (left) and viewing the *Search with Perturb* window options (right).

Search with Perturb

Search with Perturb (Figure 2.7) provides an alternative means to escape local optima by changing the transformation cost matrix of the molecular characters, a procedure similar in spirit to the parsimony ratchet. In addition to the same sequence alignment and search parameters as described for the *Simple Search* window, the *Search with Perturb* window provides three extra fields with the parameters for the transformation cost matrix perturbation as follows:

Perturb iterations Sets the number of perturb iterations to be performed.

Substitutions Specifies the cost of the perturbed substitutions.

Indels Specifies the cost of the perturbed indels.

During this heuristic search, POY5 performs a parsimony ratchet search during each iteration (default values are used, i.e. 25% probability, 2 severity).

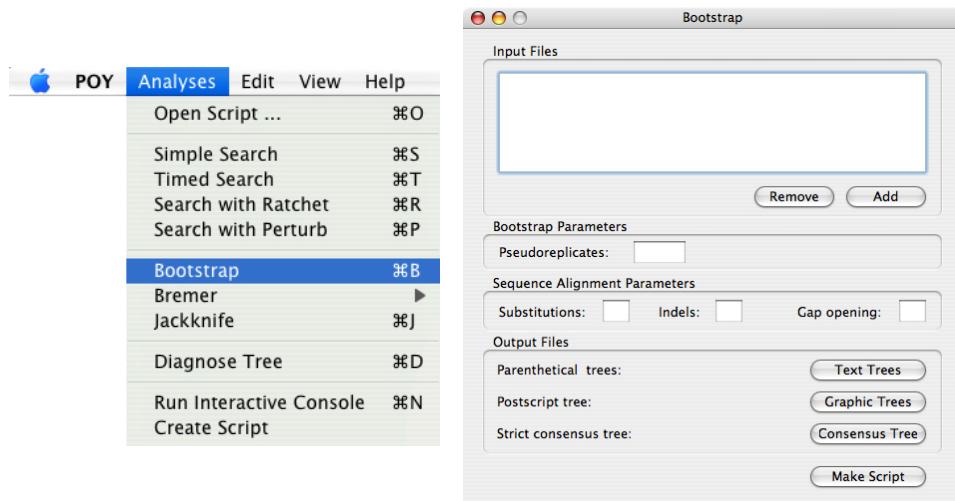


Figure 2.8: The *Bootstrap* window. Selecting *Bootstrap* from the *Analysis* menu (left) and viewing the *Bootstrap* window options (right).

Support calculation options

It is possible to calculate different support values using this interface. Two of these measures, Bootstrap and Jackknife, involve resampling techniques, while the third, Bremer support, is an optimality-based measure based on the cost of the tree.

None of the support calculation windows include functions for tree building and searching. Therefore, one of the input files must contain trees for which support values are going to be calculated.

Bootstrap

As a resampling technique, the non-parametric *Bootstrap* resamples the original data, creating a simulated dataset equal to the size of the original dataset. The characters are replaced in between each iteration. The *Bootstrap* window (Figure 2.8) specifies parameters for estimating the Bootstrap support values. In addition to the *Simple Search* window fields, it contains a *Pseudoreplicates* field, to specify the number of bootstrap pseudoreplicates.

Pseudoreplicates Specifies the number of resampling iterations.

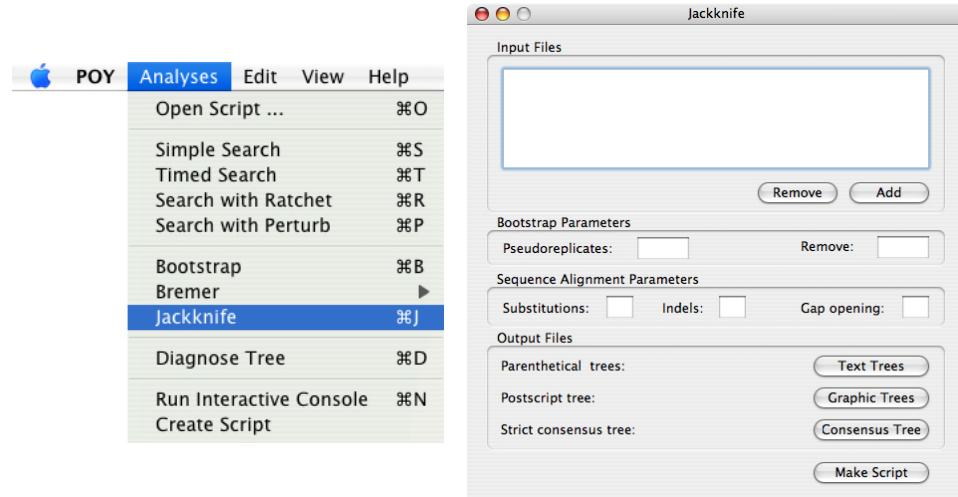


Figure 2.9: The *Jackknife* window. Selecting *Jackknife* from the *Analyses* menu (left) and viewing the *Jackknife* window options (right).

Jackknife

An alternative statistical measure of support is the *Jackknife*, wherein the original data matrix is resampled, but in this case without replacement. The *Jackknife* window (Figure 2.9) specifies parameters for estimating the Jackknife support values. In addition to the *Simple Search* window fields, *Jackknife Parameters* contains fields to specify the number of Jackknife pseudoreplicates (*Pseudoreplicates*) and the number of characters to be removed (*Remove*) during each pseudoreplicate.

Pseudoreplicates Specifies the number of resampling iterations.

Remove Specifies the percentage of characters being deleted during a pseudoreplicate.

Bremer

The *Bremer* option (Figure ??) is divided into two windows: the *Search for Bremer* window, that specifies the Bremer support [3, 24] calculation parameters, and the *Report Bremer* window to format the output of the results (Figure 2.11).

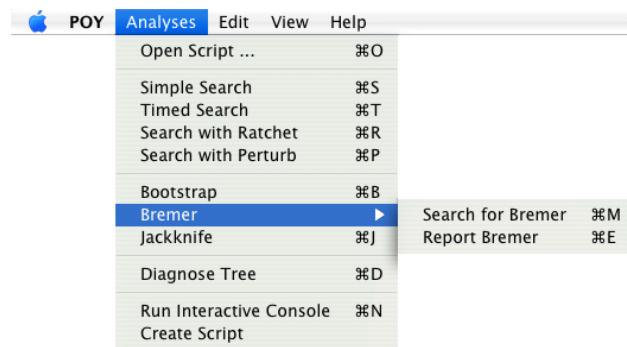


Figure 2.10: Selecting the *Bremer* windows from the *Analysis* menu.

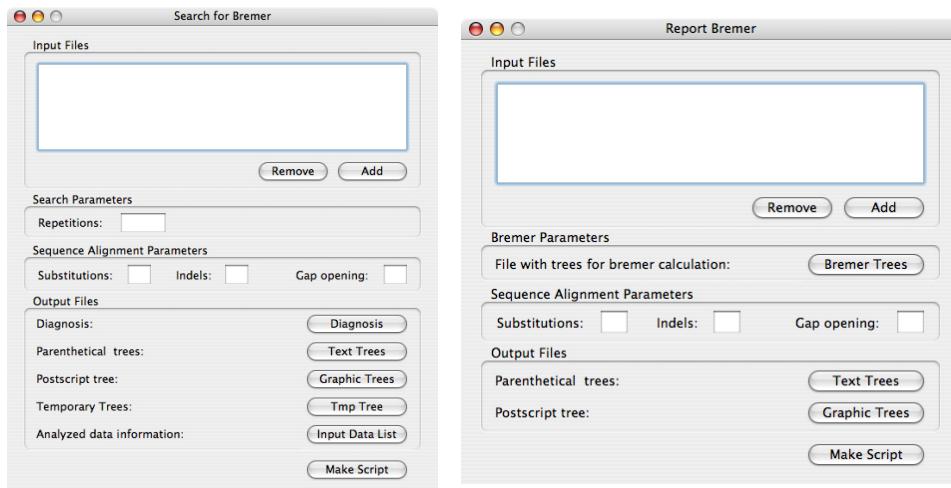


Figure 2.11: Viewing the options of the *Search for Bremer* (left) and the *Report Bremer* (right) windows.

Search for Bremer The script produced in this window collects trees visited during a search for Bremer support calculations. This search can take a long time, as the goal of this search strategy is to broadly sample variation among trees, and guarantee that all clades have Bremer support values.

In addition to the standard four sections defined for the *Simple Search* window, note that one of the output files is the *Temporary Trees* file, which contains all the information required to produce the bremer support tree results in the *Report Bremer* window. Make sure to choose a file name that does not overwrite this output.

If the search does not finish within the time frame amenable to the user the search can be interrupted and the intermediate results remain stored in the *Temporary Trees* file. As Bremer calculations are upper-bound values, terminating the search prior to completion and, thus, storing a smaller pool of visited trees may inflate support values relative to those generated by a more exhaustive search. The trees from the *Temporary Trees* file can then be reported using the *Report Bremer* window.

Report Bremer The script produced in this window takes the *Temporary Trees* file generated in the *Search for Bremer* window in the *File with trees for bremer calculation* field.

Diagnosis

Diagnose Tree

The *Diagnose Tree* window (Figure 2.12) specifies parameters for reporting a diagnosis of the input tree. This window lacks the *Search Parameters* section because the diagnosis is performed on the trees resulted from prior searches and no new trees are generated during the diagnosis procedure. It is important to add both the input tree (or trees?) in addition to the data file in order to diagnose the tree. The costs of the tree can be changed here.

Script editing and the Interactive Console

Open POY Launcher

Selecting *Open POY Script* (Figure 2.13) displays the *POY Launcher* window (Figure 2.1), the function of which is described above.

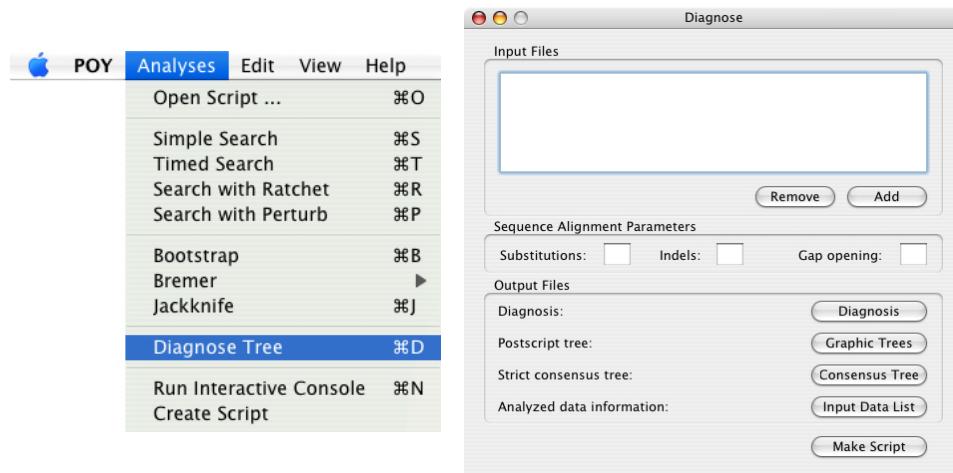


Figure 2.12: The *Diagnose* window. Selecting *Diagnose Tree* from the *Analysis* menu (left) and viewing the *Diagnose* window options (right).

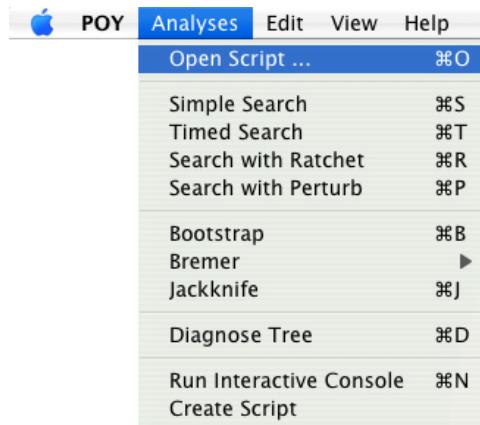


Figure 2.13: The *Open POY script* selection opens the *POY Launcher* window.

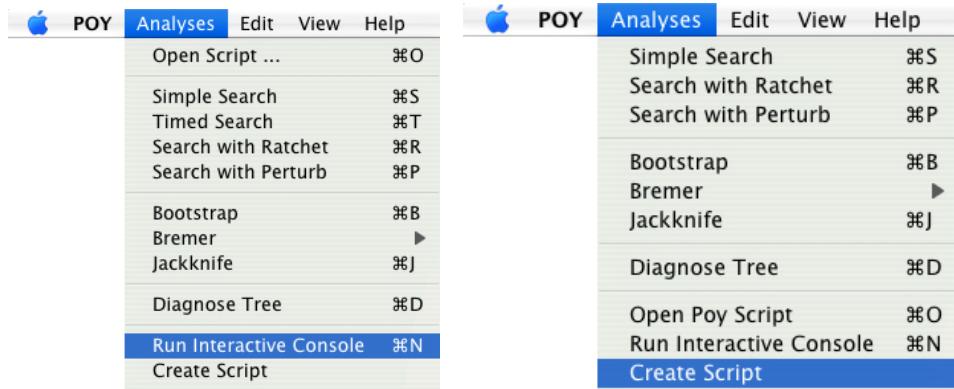


Figure 2.14: The *Run Interactive Console* selection (left) opens POY5 interactive console in a new window. The *Create Script* selection opens the *Script Editor* window (Figure 2.4).

Run Interactive Console

Selecting *Run Interactive Console* (Figure 2.14) opens the ncurses interface that enables the user to run the analysis interactively by entering POY5 commands directly via the command-line interface of the *Interactive Console*. See *Using the Interactive Console* (Section 2.4).

Create Script

The *Create Script* selection opens a blank *Script Editor* window that allows opening, creating, modifying, saving, and executing a customized script.

2.3.4 The *View* menu

The *View* menu contains the *Output* window which is subdivided into two fields: the upper *Results and Errors* and lower *Status of Search* (Figure 2.15). These fields display, respectively, the results (including warning and error messages) and the current state of the analysis. These fields are not updated automatically and in order to display the current state of the analysis the user must click the *Update* button. The *View* menu also contains the *About POY* window in Windows and Linux.

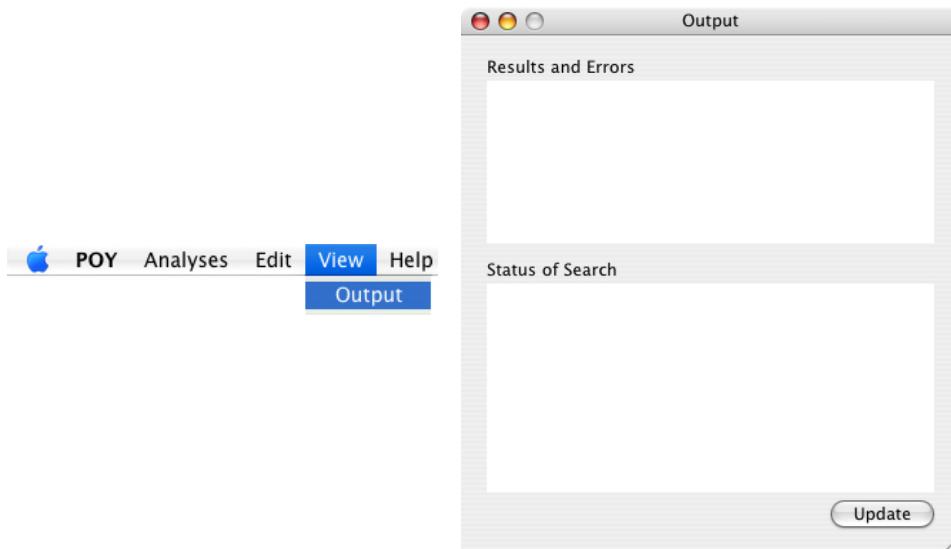


Figure 2.15: Selecting the *Output* window (left) and viewing the *Results and Errors* and *Status of Search* fields.

2.4 Using the Interactive Console

This section will help you get started using the *POY5 Interactive Console* and will prepare you for the more extensive, technical descriptions in the next chapter, *POY5 Commands*. This section will teach you how to input data files, check the data you are analyzing, generate a set of initial trees, do basic branch swapping to find a local optimum, and, finally, produce and visualize the resultant trees, their strict consensus, and generate support values in a command-line environment rather than using a *Graphical User Interface*.

For the purpose of this exercise, three data files are used available at

<http://research.amnh.org/scicomp/projects/poy.php>

- `18s.fas` and `28s.fas` contain unaligned DNA sequences (partial 18S and 28S ribosomal DNA) in FASTA format. [35]
- `morpho.ss` contains a morphological data matrix in Hennig86 format. [8]

Once POY5 has been launched and the interface (Figure 2.16) had appeared on the screen, the data can be input and the analysis can proceed. As you

follow the instructions, you are encouraged to consult the help file by using the command `help` (see Section 2.6 to learn more about POY5 commands and their arguments).

2.4.1 The interface

The *Interactive Console* provides a terminal environment with enhanced ability to display the results and the state of the analysis. It is recommended to use the console to explore and verify the data in the early steps of the analysis, and to learn the scripting language. As opposed to what? I'd imagine that a lot of people will never compile and use the flat environment. The only benefit to not using the interactive console and going with flat, is that you can run in parallel or a cluster. Using the console requires familiarity with POY5 commands, their arguments, and the conventions of POY5 scripting (which are discussed in the *POY Commands* chapter). It has four windows: *POY Output*, *Interactive Console*, *State of Stored Search*, and *Current Job* (Figure 2.16):

POY Output (Figure 2.16, upper box) displays the status of the imported data, outputs the results of the phylogenetic analyses (such as trees, character diagnoses, and implied alignments), reports errors, and displays descriptions of POY5 commands.

Interactive Console (Figure 2.16, mid-left box) is used to issue the commands interactively and to execute the commands by clicking the Return key. (See Section 3.1.1 on the description of POY5 commands.)

State of Stored Search (Figure 2.16, mid-right box) displays the time (in seconds) elapsed since the initiation of the current operation. This window also reports the number of trees currently in memory and displays the range of their costs.

Current Job (Figure 2.16, lower box) describes the currently running operation. When the operation is completed, the box is blank.

2.4.2 Starting a POY5 session using the *Interactive Console*



Windows

- Start>All Programs>POY>Interactive Console

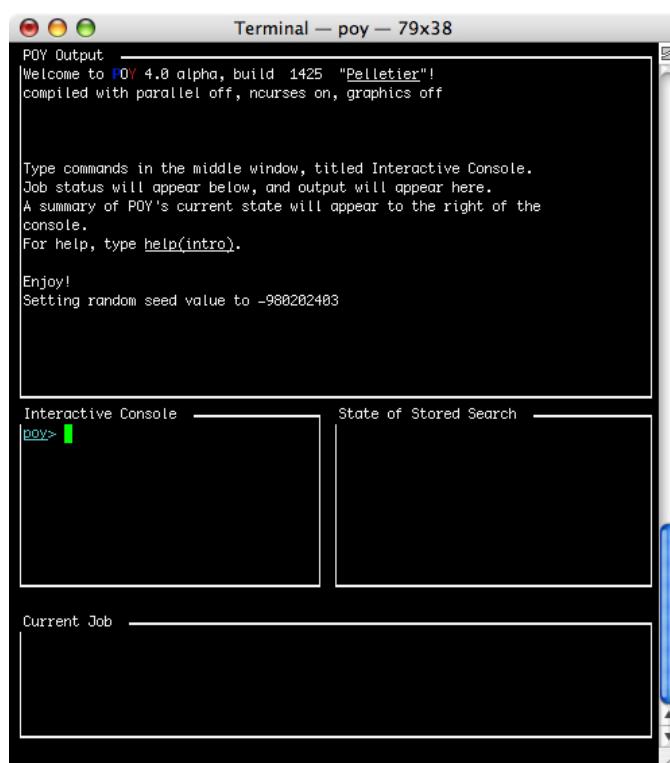


Figure 2.16: POY5 interface displayed in the Terminal window prior to analysis. Note the cursor at the POY5 prompt in the *Interactive Console* and that the *State of Stored Search* and *Current Job* windows are empty.

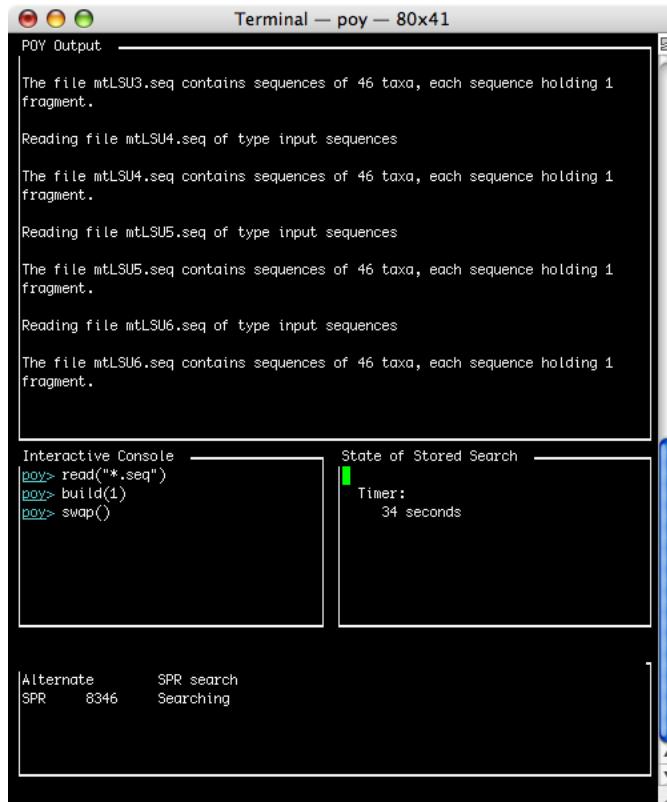


Figure 2.17: POY5 interactive console during a process. The *POY Output* window displays (by default) the information on the input data files. The *Interactive Console* lists the commands that have been consecutively executed. The *Current Job* window shows the state of the current operation and the current tree score. The *State of Stored Search* shows the time elapsed since the last command, `swap`, was initiated.



Mac OS X

- Double-click POY5 application icon to start the program.
- Select *Run Interactive Console* from the *Analyses* menu.



Linux

- Add `/opt/poy4/Resources/` to your PATH and run `ncurses_poy` from a terminal.

2.4.3 Entering commands

Once this POY5 interface is opened, the cursor appears in the *Interactive Console* portion of the window and is ready to accept commands. The interactive console does not support using the mouse and, as true for most command-line applications, the cursor can be moved using the left and right arrow keys, and the Backspace (in Windows) or Delete (in Mac) keys are used to erase individual characters to the left of the current cursor position. To eliminate the need of retying commands anew during a POY5 session, keyboard shortcuts can be used: control-P (“previous”) and control-N (“next”) will scroll through the commands previously entered during the session. In addition, the interactive console is equipped with the autocomplete feature: it involves POY5 predicting a command, an argument, or file name that the user wants to type from the first letter(s) entered. Upon typing the first letter or part of the phrase, repeatedly pressing the TAB key scrolls through the list of command, argument, and file names that begin with that letter or phrase. Autocomplete speeds up interaction with the program.

2.4.4 Browsing the output

As more output is reported in the *POY Output* window, only the most recent reports will be seen in the window. Using the Up and Down keys allows the user to scroll up and down the *POY Output* window to see the welcome line, and previously printed reports and help descriptions. Pressing Up and Down keys automatically places the cursor in the lower left corner of the *POY Output* window indicating that you are interacting with that window. Only 1000 lines are stored in the memory and the output that was reported before that will not be accessible by scrolling. The number of lines, however, can be modified by the user using the command `set()`, see `history` (Section 3.3.24). If the user desires to keep the entire output or specific items in the output, a

log can be created using the command `set()`, see `log` (Section 3.3.24)) or specific outputs can be redirected to files (see `report` (Section 3.3.19)). The user should be aware that outputting a log file can slow down the program due to I/O (input/output) delay.

2.4.5 Switching between the windows

To return to the *Interactive Console*, start typing and the cursor will automatically be placed back at the POY5 prompt. When an operation is in progress (shown in the *Current Job* window), the cursor stays in the upper left corner of the *State of Current Search* window, and switching between the *Interactive Console* and the *POY Output* window is disabled. There are no user interactions in the *Current Job* or *State of the State of Current Search*.

2.4.6 Importing data

The basic command to input data in POY5 is `read()`, which includes the list of files (in quotation marks and separated by commas) enclosed in parentheses. Suppose that we would like to simultaneously analyze morphological and molecular datasets, contained in separate data files, `morpho.ss` and `28s.fas`, respectively. We can issue a pair of `read()` commands (Figure 2.18):

```
read("morpho.ss")
read("28s.fas")
```

The syntax of `read`, like every command in POY5, contains two elements: the name of the command, in this case `read`, followed by an optional list of arguments separated by commas and enclosed in parentheses. All filenames read into POY5 should include the appropriate suffix for the file type (e.g. `.fas`, `.ss`, `.aln`, `.tre` etc.). Typically, the arguments of the command `read()` are names of data files, each being enclosed in double quotes (as shown in the example above). Even though there might be only one argument or none in some commands, parentheses (e.g. `read()`) always follow the command name. An exhaustive discussion of POY5 command structure and detailed descriptions of all commands with examples of their usage are provided in the *POY Commands* chapter (3.1.1).

In order to import data by entering the names of the files, the directory containing these files must be identified. This can be established in two ways—by using the command `cd` to redirect the path to the directory where the data is found and then reading in the data file:

The screenshot shows the POY5 Interactive Console interface. The main window displays log messages from the parser and loader. The 'POY' window at the bottom shows two consecutive `read` commands: `read ("morpho.ss")` and `read ("28s.fas")`. The 'State of Stored Search' window on the right shows 'Trees: Storing 0 trees'. A vertical bar highlights the second `read` command.

```

Reading file morpho.ss of type hennig86/None
Starting Parser
Finished Parser
The file morpho.ss defines 174 characters in 17 taxa.
Starting Loading Characters
Finished Loading Characters
Reading file 28s.fas of type input sequences
The file 28s.fas contains sequences of 17 taxa, each sequence holding 1 fragment.

Interactive Console
| poY read ("morpho.ss")
| poY read ("28s.fas")
| poY |
State of Stored Search
Trees:
Storing 0 trees

```

Figure 2.18: Importing data files using the *Interactive Console*. Two consecutive `read` commands specify both the morphological data file in Hennig86 format (`morpho.ss`), and the molecular data file in FASTA format (`28s.fas`). Note that POY5 automatically reports in the *POY Output* window the names and types of files that have been imported.

```
cd ("/Users/username/docs/poyfiles")
read ("28s.fas")
```

or by including the full path in the argument of `read`:

```
read("/Users/username/docs/poyfiles/28s.fas")
```

Most of the time users are interested in importing multiple data files to analyze an entire dataset. In this case, multiple data files can be specified as arguments for a single command. For example, importing both files, `morpho.ss` and `28s.fas`, can be written more succinctly:

```
read("morpho.ss", "28s.fas")
```

or if the full path is included in the argument of `read` as:

```
read("/Users/username/docs/poyfiles/morpho.ss",
"/Users/username/docs/poyfiles/28s.fas")
```

This is equivalent to sequentially importing each file as shown in (Figures 2.18 and 2.19).

Figures 2.18 and 2.19 also illustrate an important feature that makes POY5 different from many other phylogenetic analysis programs: every time a file is imported during a POY5 session, the input data are *added* to the current data in memory and *do not replace them*. This allows additional analytical flexibility. For example, if only morphological data are read and trees are built based on these data alone, a subsequently imported molecular character dataset will be used in conjunction with the previously imported morphological data, despite the fact that current trees in memory were generated only from morphological data (Figure 2.19):

```
read("morpho.ss")
build()
read("28s.fas")
rediagnose()
swap()
```

It must be noted that if the numbers of terminals differ among data files, *only* the data that correspond to the terminals used to generate the trees (in this case, the morphological data file) are used. The rest of the character data are ignored, unless the trees are built again with the data files containing the expanded number of terminals.

Also, because POY5 appends trees and data in memory, it is a good practice when starting a new analysis to empty the memory using the command `wipe()`.

Valid input files include nucleotide and amino acid sequence files in many formats, and morphological data in Hennig86 and Nexus formats. (For information on specific formats supported by POY5 and other types of input files see `help(read)`.)

2.4.7 Inspecting data

Once a dataset (or multiple datasets) is imported, POY5 automatically reports a brief description of contents for each loaded file in the *POY Output* (Figure 2.18). However, it may be desirable to inspect the imported data in greater detail to ensure that the format and contents of the files have been interpreted correctly. This practice helps to avoid common errors, such as inconsistently spelled terminal names, which may result in bogus results, produce error messages, and aborted jobs.

The basic command for outputting information is `report()`. One of its arguments, `data`, outputs a set of tables showing the list of terminals, the

```

Reading file morpho.ss of type hennig86/Nona
Starting Parser
Finished Parser
The file morpho.ss defines 174 characters in 17 taxa.
Starting Loading Characters
Finished Loading Characters
Starting Wagner build
Finished Wagner build
Reading file 28s.fas of type input sequences
The file 28s.fas contains sequences of 17 taxa, each sequence holding 1 fragment.

Starting Diagnosis
Finished Diagnosis
Starting Tree search
Finished Tree search

Interactive Console - State of Stored Search -
[<pre>
[m] read ("morpho.ss")
[m] build ()
[m] read ("28s.fas")
[m] rediagnose ()
[m] swap ()
[m]</pre>]
Trees:
Storing 10 trees with costs 936. to 948,
Best cost was found once

```

Figure 2.19: Building trees with morphological data only but continuing analysis using combined morphological and molecular data. This example shows how we can add data to the analysis incrementally by loading files at different points in the search. First, the morphological data are imported from `morpho.ss` file using `read()` and trees are built based on these data. Then molecular data from the `28s.fas` file are loaded into memory in addition to previously imported morphological data. Finally, subsequent analyses, `rediagnose()` and `swap()`, are conducted using the data in memory, that is the trees based on morphological data, and both morphological and molecular character sets.

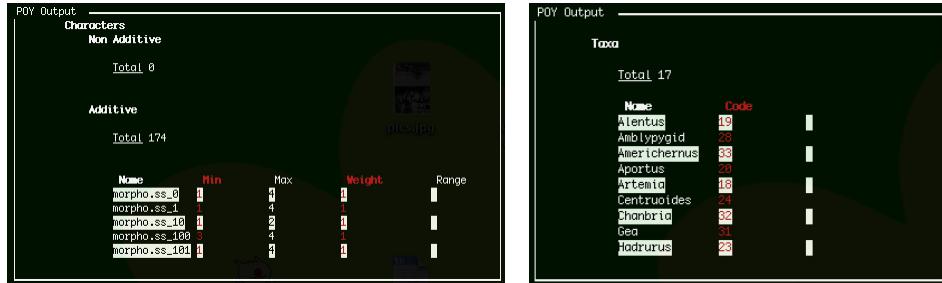


Figure 2.20: Inspecting imported data. The figure shows segments of a data report generated by `report(data)`. The left and right panels demonstrate a typical table output the character and terminal data respectively.

number and types of characters, and the lists of terminals and characters excluded from the analysis. To produce a report of the data files that were used in the previous example (`morpho.ss` and `28s.fas`), we import the data and execute `report(data)`:

```
read("morpho.ss", "28s.fas")
report(data)
```

This will generate an extensive, detailed output, partial views of which are shown in Figure 2.20. Obviously, the entire report will not be visible in the *POY Output* window. Therefore, the Up and Down arrow keys and Page Up and Page Down keys can be used to scroll. By default, *POY5* reports the results of executed commands to the *POY Output* window. However, the same output can be redirected to a file simply by adding the name of the output file in the list of argument of the command `report()` before the argument specifying the type of the requested report (in this case `data`, see the command `report` (Section 3.3.19)). For instance, if we would like to output into the file “`data_analyzed.txt`,” we would enter:

```
read("morpho.ss", "28s.fas")
report("data_analyzed.txt", data)
```

In this example, all the imported data are analyzed and, therefore, the report fields that list excluded data will appear empty. One can, however, exclude specific characters or terminals from the analysis using additional commands (see the command `report` (Section 3.3.19)).

Another useful argument of `report` is `cross_references`. This argument displays whether character data are present or absent for each terminal in each one of the imported data files. This provides a comprehensive visual

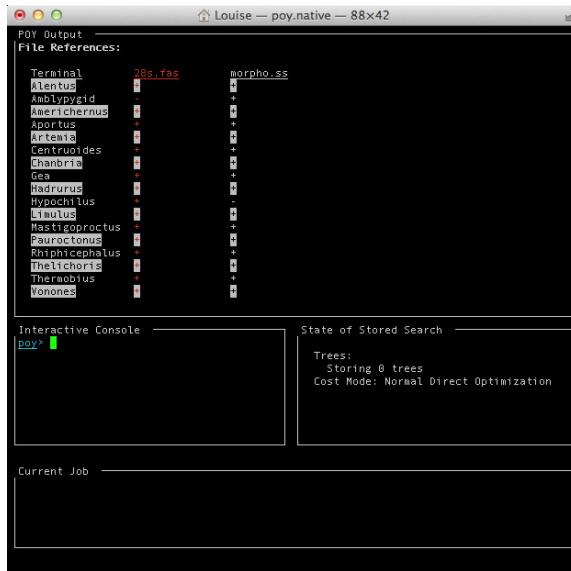


Figure 2.21: Visualizing missing data. The command `cross_references` displays a table showing whether a given terminal (in the left column) is present (“+”) or absent (“-”) in each data file. In this example, `28s.fas` is missing for Amblypygid and `morpho.ss` for Hypochilus.

overview of missing data. Building on the previous example, such output can be generated by the following sequence of commands:

```
read("morpho.ss", "28s.fas")
report("cross_refs.txt", cross_references)
```

A typical output of `cross_references` command is shown in Figure 2.21.

2.4.8 Building the initial trees

The command to build trees is `build()` (already mentioned in Section 2.4.6). After importing `morpho.ss` and `28s.fas`, executing the command `build()` without specifying any arguments (default settings) generates 10 trees by random addition sequence.

Many POY5 commands operate under default settings when executed without arguments. To learn what the default settings are for a particular

command, use either `help()` command with the command name of interest inserted in parentheses or consult the *POY Commands* chapter ([3.1.1](#)).

If the user would like to specify a number of tree building replicates different from the default value of 10, an argument `trees` followed by a colon (“`:`”) and an integer specifying the number of trees must be included in the argument list of the `build` command: `build(trees:100)`. This command has a shortcut that omits the argument `trees`. Thus, `build(trees:100)` is equivalent to `build(100)`. As defaults, the shortcuts are fully described in Section [3.1.1](#). The entire sequence of commands minimally required to import the data and build 100 trees is the following:

```
read("morpho.ss","28s.fas")
build(100)
```

As the tree building advances, the *Current Job* window displays the current status of the operation (Figure [2.22](#)). This window shows how many Wagner builds have been performed out of the total number requested, the number of terminals added in the current build, the cost of the current tree (recalculated after each terminal addition), and the estimated time for the completion of all the builds. When all the trees are generated, the *State of Stored Search* window displays the range of tree costs (the best and worst costs), the number of trees stored in memory, and the number of trees with the best cost (Figure [2.22](#)).

2.4.9 Performing a local search

Now that the trees have been generated and stored in memory, a local search can be performed to refine and improve the initial trees by examining additional topologies of potentially better cost. The command `swap()` implements an efficient strategy by performing SPR and TBR branch swapping alternately. As with other commands, the arguments of `swap()` allow the customization of the swap algorithm. In the following example, branch swapping is performed under the default settings on each of the 100 trees build in the previous step:

```
read("morpho.ss","28s.fas")
build(100)
swap()
```

Branch swapping is performed sequentially on all trees stored in memory. During swapping, the *Current Job* window reports the number of the tree



Figure 2.22: Generating Wagner trees. During the process of tree building (left panel), the *Current Job* window displays how many builds have been performed so far (57 of 100), the number of terminals added in the current build (13 of 17), a cost of a current tree recalculated after each terminal addition (362), and the estimated time (in seconds) for the completion of the operation (4 s). Because the process is not complete, the *State of Stored Search* window contains no trees. Once tree building is complete, the *State of Stored Search* window displays the best (451) and worst (472) costs, the number of trees stored in memory (100), and the number of trees with the best cost (2).

that is currently being analyzed, the method of branch swapping, the specific routine being currently performed, and the cost of the current tree (Figure 2.23). When the process is complete, the *State of Stored Search* window displays the range of tree costs (the best and worst costs), the number of trees stored in memory, and the number of trees of the best cost (Figure 2.23). Note that the local search had reduced the costs of the initial best (from 451 to 446) and narrowed the range of tree costs.

Using different combinations of `swap()` arguments allows designing a large number of search strategies of different levels of complexity. Some simple options allow the choice between SPR and TBR. More complex strategies allow keeping a specific number of best trees per single initial tree (generated during the building step). For example, the command `swap(trees:10)` will keep up to 10 best trees generated during branch swapping on a single initial tree. Consequently, if 100 trees were built initially, this command will produce up to 1,000 trees. The argument `threshold` allows the retention of suboptimal trees within a specified percent of cost difference from the current best tree. For example, `swap(trees:20, threshold:10)` will execute a swap considering trees within a ten percent cost difference of the current best tree and retain the 20 minimal length swapped trees for each initial build. Other options provide the means to sample trees as they are evaluated, timeout after certain number of seconds, transform the cost regime, and other functions in conjunction with many POY5 commands.



Figure 2.23: Performing a local search. When searching (left panel), the *Current Job* window reports the number of the tree that is currently being analyzed (73 of 100), a method of branch swapping (**Alternate**), a function being currently performed (**SPR search**), and a cost of the current tree (456). When the searching is finished (right panel), the *State of Stored Search* window displays the best (446) and worst (463) costs, the number of trees stored in memory (100), and the number of trees of the best cost cost (9) recovered from independent tree builds. Note these trees may not necessarily have unique topologies.

2.4.10 Selecting trees

Having performed the basic steps of importing character data, building initial trees, and conducting a simple local search, we obtained a set of local-optima trees in memory. Most of the time, a user would like to select only those trees that are both ptimal *and* topologically unique. The default setting of the **select()** does exactly that. Adding **select()** to our example of command sequence for the basic analysis

```
read("morpho.ss", "28s.fas")
build(100)
swap()
select()
```

selects only unique trees of best cost. The remaining trees are deleted from memory. The *State of Stored Search* window reports the number and the cost of the best tree(s) (Figure 2.24).

This command **select()**, is another multifunctional command the arguments of which are also used to select (include or exclude) specific terminals, characters, and trees.) Comparing the output reported in the *State of Stored Search* before (Figure 2.23) and after (Figure 2.24) executing **select()** shows that swapping on 9 of 100 initial trees produced the trees of best cost (446), but these trees are identical, because only one was retained when filtered using **select()**.

The screenshot shows two windows side-by-side. The left window is titled 'Interactive Console' and contains the following command history:

```

Interactive Console
000> read("morpho.ss", "28s.fos")
000> build(100)
000> swap()
000> select()
000> 

```

The right window is titled 'State of Stored Search' and displays the results of the search:

```

State of Stored Search
Trees:
  Storing 1 tree with cost 446.

```

Figure 2.24: Selecting unique best trees. Executing `select()` keeps only unique trees of best cost. The *State of Stored Search* window reports that there is only one unique tree of best cost (446).

2.4.11 Visualizing the results

There are several options for visualizing results in POY5 (see Section ??). The command `report("my_first_tree", graphtrees)` outputs a cladogram in PDF format (Figure 2.25), which can be displayed, edited, and printed using graphics software (such as Adobe Illustrator or Corel Draw). POY5 also appends the “ps” extension when generating graphic output to a file. A quick way to see the tree(s) on screen is to use the command `report(asciitrees)` that draws a cladogram in the *POY Output* window (Figure 2.25). The ascii tree(s) can also be reported to a file, if an output file name is specified within the command (`report("my_first_trees", ascitrees)`). These trees will be saved to a text file.

The command `report("my_first_trees.txt", trees)` reports the trees in memory in parenthetical notation to the file `my_first_trees` that can be imported in other programs. Other supported tree output formats include Newick and Hennig86. `report()` can also generate consensus trees in the graphical and parenthetical formats when appropriate arguments are specified (for example, `report("strict_consensus", graphconsensus)`).

2.4.12 Interrupting a process

To interrupt a process, press Control-C. By default, an error, `Error: Interrupted`, is reported in the *POY Output* window. The program does not close, however, and a new command can be entered. Interrupting the analysis cancels the execution of the last command requested by the user and restores the data and trees in memory before that last command. For example, the following two session are equivalent:

```
read("morpho.ss") <ENTER>
```

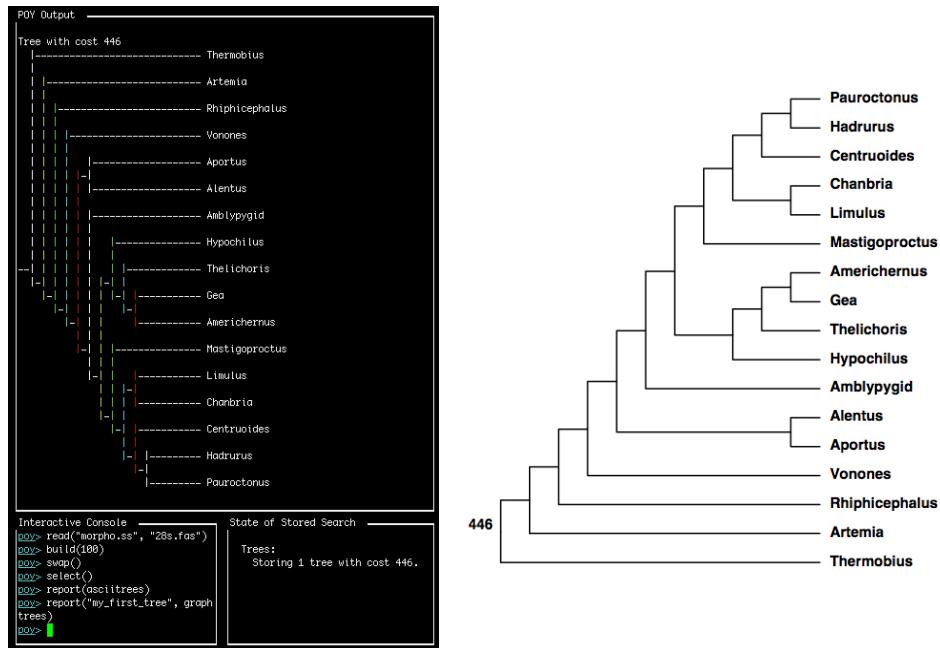


Figure 2.25: Visualizing trees. An ascii tree (left) is generated using the command `report(asciitrees)`. The same tree is reported to a file in a PDF format (right) using `report("my_first_tree", graphtrees)`. Note that both representations of trees are preceded by their costs.

and

```
read("morpho.ss") <ENTER>
read("28s.fas", "18s.fas") <CONTROL-C>
```

In both of these sessions, only the morphological dataset “morpho.ss” is read into POY5.

2.4.13 Reporting errors

If there is an error pertaining to wrong syntax (such as a typo in a command name), POY5 will indicate the location of the error by underlining the problematic part of the input with “~” in the *Interactive Console* (Figure 2.26). The description of the corresponding command, its syntax, and examples of its usage from the help file are automatically printed in the *POY Output* window. As noted above, the Up and Down keys can be used to scroll through the output and determine the source of the error. Certain types of errors are reported explicitly (Figure 2.26).

2.4.14 Exiting

To finish a POY5 session, enter the command `exit()` (Figure 2.27) or `quit()`. This will close the POY5 interface and resume the Terminal window (Mac OSX) or the Command Prompt window (Windows).

2.5 Creating and running POY5 scripts

So far, we have communicated with POY5 interactively through the *Graphical User Interface* or by executing commands from the *Interactive Console*. Another way of conducting an analysis is to run a *script*, a simple text file containing a list of commands to be performed (Figure 2.28).

Running analyses using scripts has many advantages: not only does it allow for the entire analysis to proceed from the beginning to the end at one click of a button, but it also provides means to examine the logical dependency of the commands and optimize memory consumption (see the description of `script_analysis` argument of the command `report` in the *POY Commands* chapter). Submitting jobs using scripts may produce results faster because POY5 automatically optimizes the workflow of the entire analysis by taking into account the functional relationships among various tasks and efficiently distributing the jobs and resources (such as memory and multiple processors).

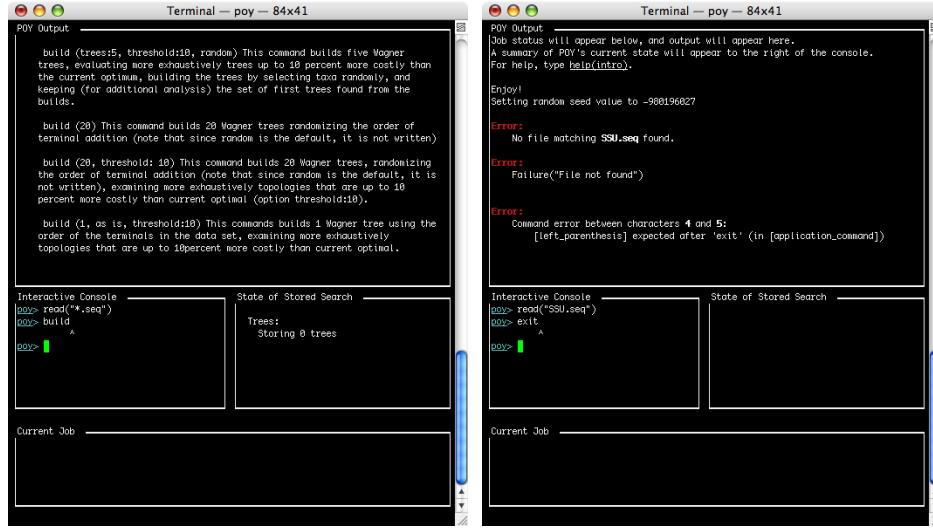


Figure 2.26: Displaying errors. POY5 displays error messages in several ways. In the example in the left panel, the command `build` was entered without parentheses, which is required for a valid POY5 command syntax; the exact place of the error is marked by “^”, in this case following the `build` commands. Examples of the proper usage of the command are automatically displayed in the *POY Output*. In other cases (right panel), error messages are explicitly reported in the *POY Output* window. The first and second error messages indicate that the data file `SSU.seq` is not present, which could have been caused either by a mistake in the name of the file, missing file, or the location of the file in a directory, other than the one specified prior to starting the POY5 session. The third error message indicates that the valid syntax of `exit` requires the parentheses following the command name (also shown by “^” in the *Interactive Console*).

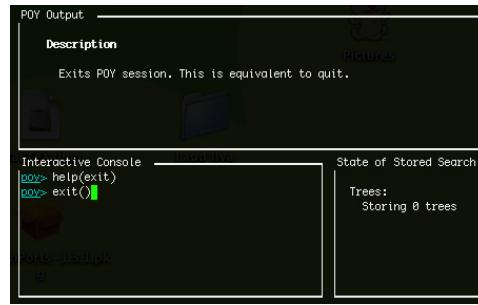


Figure 2.27: Exiting POY5

Another advantage of using scripts is that they may contain comments that are ignored by POY5 but can be helpful to describe the contents of the files and provide other annotations. The comments are enclosed in parenthesis and asterisks. For example, (* this is a comment *). Comments can be of any length and span multiple lines. Comments can also be entered interactively from the *Interactive Console*.

Obviously, using scripts requires the user to design the workflow of the process prior to conducting the analysis. POY5 scripts can be created and saved using the *Script Editor* window of the POY5 *Graphical User Interface* or any conventional text editor (such as TextPad, TextWrangler, BBEdit, Emacs, or NotePad).

POY5 scripts are extremely useful in cases when operations may take a long time to complete, eliminating the need to wait for a part of the analysis to finish in order to proceed to the next step.

There are two ways to import and run a script:

- using the *POY Launcher* in the *Graphical User Interface*;
- using the command `run()` of the *Interactive Console*; for example, `run("script.txt")`, where `script.txt` is the name of the file containing the script.

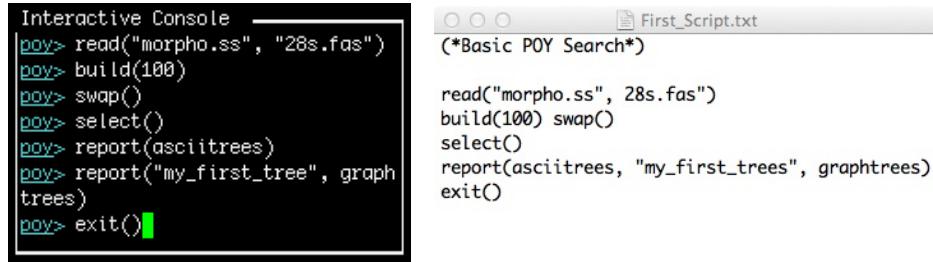
It is critical to include the command `exit()` at the end of the script. Otherwise POY5 will be waiting for further instructions to be entered after executing the script's contents.

2.6 Obtaining help

Instructions to run POY5, command descriptions, and the theory behind POY5 can be obtained from a variety of sources.

POY5.0 Program Documentation (this manual) is a comprehensive and detailed manual on all the aspects of using POY5, from installation to output and visualization of results. Included are *Quick Start*, POY5 command reference, practical guides and tutorials that make the program immediately accessible for beginners and provide in-depth information for experienced users. The documentation in PDF format can be accessed from the *Help* menu of the graphical user interface or downloaded separately from POY5 web site at

<http://research.amnh.org/scicomp/projects/poy.php>



The figure shows two side-by-side windows. The left window is titled 'Interactive Console' and contains a command-line interface with the following text:

```
poy> read("morpho.ss", "28s.fas")
poy> build(100)
poy> swap()
poy> select()
poy> report(asciitrees)
poy> report("my_first_tree", graph
trees)
poy> exit()
```

The right window is titled 'First_Script.txt' and contains a script with the following text:

```
(*Basic POY Search*)
read("morpho.ss", 28s.fas")
build(100) swap()
select()
report(asciitrees, "my_first_trees", graphtrees)
exit()
```

Figure 2.28: Using POY5 scripts. The list of commands executed interactively using the *Interactive Console* (left) and a script containing the same list of commands (right). Note, that the header of the script is a comment, enclosed in “(* *)”, that is ignored by POY5. Also note, that commands can either be listed in a row or in a column (compare `build()` and `swap()` in the console and in the script) and different arguments of the same command can either be specified separately or combined in a single argument list (compare `report()` in the console and in the script). (Both conventions are valid for interactive command submission and for scripts.)

POY interactive help can be obtained by entering `help()` at the POY5 interactive console. To obtain help on a particular command, the name of the command must be specified in the parentheses following `help()`. For example, to learn about the command `exit`, type `help(exit)`. (Figure 2.27.)

POY5 Mail Group is an Internet-based forum for discussing all issues related to POY5 and provides the best way to communicate with POY5 developers on specific issues (see *WWW resources* below). The website is located at <http://groups.google.com/group/poy5>. [new link](#)

POY Book (Wheeler et al., 2006 *Dynamic Homology and Phylogenetic Systematics: A Unified Approach Using POY* [57]) provides a review of the theory behind POY4 and by extension POY5, and contains formal descriptions of many algorithms implemented in the program and the descriptions of commands of the earlier version, POY3.

2.7 WWW resources

POY5 is an ongoing project and new versions are being continuously developed to include new procedures, improve performance, and eliminate reported bugs. Therefore, it is imperative to keep up with the program’s development

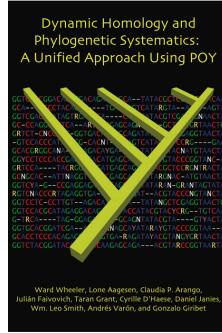


Figure 2.29: The POY Book.

and check regularly for updates. There are several Internet-based resources that offer this information.

POY5 Web Site has downloadable compressed files of POY5 binaries, source code, and documentation in PDF format. It also provides a links to the *POY Mail Group*. The website is hosted by AMNH Computational Sciences at

<http://research.amnh.org/scicomp/projects/poy.php>

POY5 source code repository contains has downloadable POY5 source code. The site is powered by Google at

<http://code.google.com/p/poy5/source> new link

POY5 Mail Group informs registered users via email of new developments, such as new versions and updates. It also provides additional resources for obtaining help and a way for reporting bugs and other problems with POY5 and its documentation. In addition, it allows users to receive and respond to each other's questions thus providing an open forum to discuss the methods and applications of POY5. The users who choose not to register, have access to the archives of the postings but will not be able to either submit or receive emails from other users and POY5 developers. The *POY5 Mail Group* is hosted by Google at

<http://groups.google.com/group/poy5> new link

Chapter 3

POY5 Commands

3.1 POY5 command structure

3.1.1 Brief description

POY5 interprets and executes *scripts* issued by the end user. These can come from *Graphical User Interface* and the command line in the *Interactive Console* of the program, or from an input file. A script is a list of *commands*, separated by any number of whitespace characters (spaces, tabs, or newlines). Each command consists of a name in lower case (`LIDENT`), followed by a list of arguments separated by commas and enclosed in parentheses. Most of the arguments are optional, in which case POY5 has default values.

In POY5, we recognize four types of command arguments: *primitive values*, *labeled values*, *commands*, and *lists of arguments*.

Primitive values can be either an integer (`INTEGER`), a real number (`FLOAT`), a string (`STRING`), or a boolean (`BOOL`).

Labeled values are lowercase identifiers (which are referred to as *label*), and an argument, separated by the colon character, “`:`”.

List of arguments are several arguments enclosed in parenthesis and separated by commas, “`,`”.

Commands are standard commands that can affect the behavior of another command when included in its list of arguments.

Thus, certain commands can function as arguments of other commands. Moreover, some commands share arguments. Although such compositive use of commands might seem complex, this structure provides much more intuitive control and greater flexibility. The fact that the same logical operation that functions in different contexts maintains the same name (typically suggestive of its function), substantially reduces the number of commands without limiting the number of operations. Using a linguistic analogy, POY5 specifies a large number of procedures by a more complex grammar (specific combinations of commands and arguments), rather than by increasing the vocabulary (the number of specific commands and arguments). For example, the command `swap` specifies the method of branch swapping. This command is used to conduct a local search on a set of trees. In addition, `swap` functions as an argument for `calculate_support` to specify the branch swapping method used in each pseudoreplicate during Jackknife or Bootstrap resampling. `swap` can also be used to set the parameters for local tree search based on perturbed (resampled or partly weighted) data as an argument of the command `perturb`. Therefore, to take the maximum advantage of POY5 functionality, it is essential to get acquainted with the grammar of POY5.

3.1.2 Grammar specification

The following is the formal specification of the valid grammar of a script in POY5:

```

script: = | command
       | command script

command: = LIDENT "(" argument list ")"

argument list: = |
               | arguments

arguments: = |
            | argument
            | argument "," arguments

argument: = | primitive
           | LIDENT
           | LIDENT ":" argument
           | command

```

```

| "(" argument list ")"
primitive: = | INTEGER
           | FLOAT
           | BOOLEAN
           | STRING

LIDENT: = [a-zA-Z_][a-zA-Z0-9_]*

INTEGER: = [0-9]++

FLOAT: = | INTEGER
         | [0-9]++ "." [0-9]++

STRING: = """ [^"]* """

```

The following examples graphically show a typical structure of valid POY5 commands formally defined above. The Figure 3.1 illustrates the syntax of the command `swap`. The name of the command, `swap`, is followed by a list of two arguments, `tbr` and `trees:2`, enclosed in parentheses and separated by a comma. Note that `trees:2` is a labeled-value argument that contains a label (`trees`) and a value (2) separated by a colon.

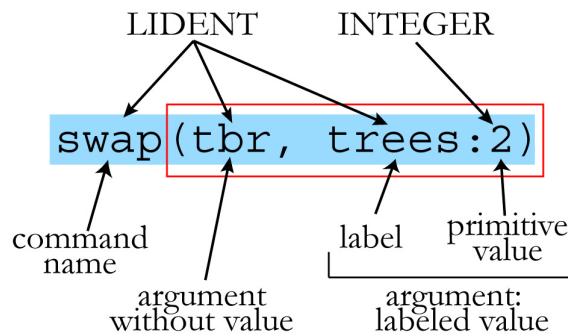


Figure 3.1: The structure of a simple POY5 command. The entire command (highlighted in blue), consists of a command name followed by a list of arguments (enclosed in red box). See text for details.

Figure 3.2 shows a more complex command structure, using the command `perturb` as an example. This is a compound command because the list of

its arguments contains another command, `swap`. This means that executing `perturb` performs a set of specified operations that contains a nested set of operations specified by `swap`. Note also, that in contrast to the first labeled-values argument `iterations`, the second labeled-values argument `ratchet` has multiple values (a float and an integer). When multiple values are specified, they must be enclosed in parentheses and separated by a comma. The third argument is a command (`swap`), therefore it is syntactically distinguished from other arguments, labeled and unlabeled alike, by having parentheses following the command name. It must be emphasized that the parentheses always follow the command name even if no arguments are specified. (If no arguments are specified, a command is executed under its default settings.)

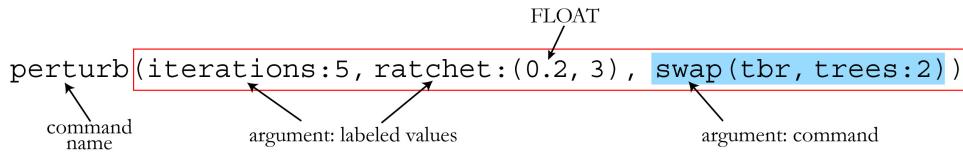


Figure 3.2: A structure of a compound POY5 command. Note that the list of arguments (enclosed in red box) includes a command (highlighted in blue). Also, note that `ratchet` accepts multiple values, a float and an integer, that are inclosed in parentheses and separated by a comma. See text for details.

3.2 Notation

Some arguments are obligatory, whereas others are not; some commands accept an empty list of arguments, but others do not; some argument labels have obligatory values, some have optional values. The POY5 commands are listed alphabetically in the next section. In the descriptions of POY5 commands below, the elements of POY5 grammar are defined in the text using the following conventions:

- A command that could be included in a POY5 script (that is can be entered in the interactive console or included in an input file) is shown in `terminal` typeface.
- Optional items are inclosed in [square brackets].
- Primitive values are shown in `UPPERCASE`.

Each command description entry contains the following sections:

- The name of the command.
- The valid syntax for the command.
- A brief description of the command's function.
- The list of descriptions of valid arguments.
- Description of default settings.
- Examples of the command's usage.
- Cross references to related commands.

NOTE

Default syntax. The default syntax for all commands is the same: it includes the command name followed by empty parentheses, e.g. `swap()`. However, within the descriptions of each command the default settings include the entire argument list for illustrative purposes only (i.e. in the case of `swap()` the entire argument list appears as `swap(trees:1, TBR, threshold:0, bfs)`).

NOTE

Command order. The effect of the order of arguments in a command depends on the context. If arguments are not logically interconnected, their order is not important. For example, the commands `build(10,randomized)` and `build(randomized,10)` are equivalent. However, executing the commands `transform(tcm:(1,1),gap-opening:4)` and `transform(gap-opening:4,tcm:(1,1))` will produce different results because `gap-opening` *modifies* the values set by `tcm`, while `tcm` *overrides* the values set by `gap-opening`.

NOTE

Output files. When an output file is specified, the file name (in double quotes and followed by a comma) must precede the argument, e.g. `report("first_trees", trees)`.

3.3 Command reference

3.3.1 build

Syntax

```
build([argument list])
```

Description

Builds Wagner trees [7]. Building multiple trees with a randomized addition of terminals allows for the evaluation of many more possible tree topologies and generates a diversity of trees for subsequent analysis. The arguments of the command `build` specify the number of trees to be generated and the order in which terminals are added during a single tree building procedure. During tree building, POY5 reports in the *Current Job* window of the ncurses interface which of the **terminal addition strategies** is currently used.

By default POY5 replaces the trees stored in memory with those generated in a subsequent build. For example, executing `build(10)` followed by `build(20)` will replace the 10 trees generated during the first build with 20 new trees. However, it might be desirable (for example, if computer memory were limited) to generate a large number of trees by appending trees from multiple separate builds. To keep trees from consecutive builds, a tree output file must be specified using the command `report` (Section 3.3.19) that must precede the `build` command. This will produce a file containing the trees appended from all builds. If the same file name is used for reporting trees for other analysis, the new trees are going to be appended. Alternatively, trees from different builds can be redirected to separate files if different file names are specified.

The command `build` is also used as an argument for the command `calculate_support`.

Arguments

all Turns off all **preference strategies** for adding branches and simply tries all possible addition positions for all terminals. By default only one tree is built but the number of trees can be specified by the argument `trees`.

as_is Indicates that in one of the trees to be built, the terminals are added in the order in which they appear in the imported datafiles, and all others are built using a random addition sequence.

branch_and_bound[:FLOAT] Calculates the exact solution using the Branch and Bound algorithm [20]. By default only one optimal tree is kept but the number of optimal trees to be retained can be specified by the argument **trees**. The optional float value specifies the bound (either tree cost or likelihood score).

constraint[:STRING] Builds trees using the set of constraints specified by the consensus tree input file. If no input file is provided, the constraint is calculated as the strict consensus of the trees in memory. Every tree built using this method is subjected to the same randomization as Wagner builds within each constraint.

lookahead:INTEGER The number of trees that can be kept at each build step. If the **lookahead** argument specifies a number n , and the best tree found has cost c , then the best n trees with cost at most $c + \text{threshold}$ as specified by the **threshold** (Section 3.3.1) command are held for the next build step. If no **threshold** command is specified, then it is set to 0.

nj Creates a tree using the Neighbor Joining algorithm. If more than one tree is requested, all the trees will be the same (the algorithm implementation is deterministic).

random Generates a tree at random. All possible trees have equal probability.

randomized Indicates that terminals are added in random order on every Wagner tree built. This is a default tree-building strategy.

threshold:FLOAT The numerical value specifies the extra cost over the current best tree that makes another tree acceptable for the lookahead list. This parameter is only useful if **lookahead** (Section 3.3.1) is used.

trees:INTEGER The integer value specifies the number of independent, individual Wagner tree builds. The label **trees** is optional: it is sufficient to specify only the integer value. Therefore, **build(5)** is equivalent to **build(trees:5)**. Note that **trees** is also used as an argument of the command **swap** (Section 3.3.26) but with different meaning.

The value 0 generates no trees but it *retains* all trees in memory. This is useful, for example, in the **bremer** (Section 3.3.2) support calculation,

where instead of generating new trees per each node, the searches are performed on the trees in the neighborhood of the current trees in memory.

INTEGER The integer argument specifies the number of independent, individual Wagner tree builds. This is a shortcut of the argument **trees**.

of_file:STRING Imports a tree file included in the file path of the argument.

This command is useful for importing starting trees for calculating **bremer** (Section 3.3.2) support. In other contexts the command **read** (Section 3.3.14) can be used with the same effect.

STRING This is a shortcut of the argument **of_file**.

Defaults

build(trees:10, randomized, lookahead:1, threshold:0) By default, POY5 will build 10 trees using a random addition sequence for each of them.

Examples

- **build(20)**

Builds 20 Wagner trees randomizing the order of terminal addition (note that because the argument **randomized** is specified by default, it can be omitted).

- **build(trees:20, randomized)**

A more verbose version of the previous example. By default a build is randomized, but in this case the addition sequence is explicitly set. For the total number of trees, rather than simply specifying 20, the label **trees** is used. The verbose version might be desirable to improve the readability of the script.

- **build(15, as_is)**

Builds the first Wagner tree using the order of terminals in the first imported datafile and generates the remaining 14 trees using random addition sequences.

- **build(branch_and_bound, trees:5)**

Builds trees using branch and bound method and keeps up to 5 optimal trees in memory.

3.3.2 calculate_support

Syntax

```
calculate_support([argument list])
```

Description

Calculates the requested support values. P0Y5 implements support estimation based on resampling methods (Jackknife [10] and Bootstrap [13]) and Bremer support [3, 24]. The Jackknife and Bootstrap support values are computed as frequencies of clades recovered in strict consensus trees built in each resampling iteration. The consensus trees are based on best trees recovered in each replicate with zero-length branches collapsed. All the arguments of `calculate_support` command are optional and their order is arbitrary. For examples of scripts implementing support measures see [tutorials 4.3 and 4.4](#).

The `calculate_support` command does not output support values by default. The output of support values must be requested using the command `report` (Section [3.3.19](#)). This is particularly important for Jackknife and Bootstrap support values, as these sampling techniques do not require the presence of trees in memory. Therefore, it is possible to perform the sampling for support values *before* the tree of interest has been found.

Arguments

Support calculation methods The following commands allow selecting among several methods for calculating support.

bootstrap[:INTEGER] Calculates Bootstrap support [13]. The integer value specifies the number of resampling iterations (pseudoreplicates). If the value is omitted, 5 pseudoreplicates are performed by default.

bremer Calculates Bremer support values [3, 24] for each tree in memory by performing independent constrained searches for each node. The parameters for the searches can be modified using arguments described under *Search strategy*.

jackknife[:([argument list])] Calculates Jackknife support [10] using the sampling parameters specified by the arguments. The arguments of `jackknife` are optional and their order is arbitrary. If both values are omitted, the default values of each argument is used.

NOTE

In the context of dynamic homology, the characters being sampled during pseudoreplicates are entire sequence fragments, not individual nucleotides. Consequently, the bootstrap and jackknife support values calculated for dynamic characters are not directly comparable to those calculated based on static character matrices. If it is desired to perform character sampling at the level of individual nucleotides, the dynamic characters must be transformed into static characters using `static_approx` argument of the command `transform` (Section 3.3.26) prior to executing `calculate_support`. Alternatively, an output file in the Hennig86 format can be generated based on an implied alignment using `phastwinclad` (Section 3.3.19) that can subsequently be analyzed using other programs.

It is important to remember that the local optimum for the dynamic homology characters can differ from that for the static homology characters based on the same sequence data. Therefore, it is recommended to perform an extra round of swapping on the transformed data to reach the local maximum for the static homology characters prior to calculating support values.

NOTE

The placement of the root affects calculation of Bremer support values. Therefore, it is critical to specify the root prior to executing `calculate_support`. See the description of the command `set` (Section 3.3.24) on how to specify the root.

remove:FLOAT The value of the argument `remove` specifies the percentage of characters being deleted during a pseudoreplicate. The default of `remove` is 36 percent.

resample:INTEGER The value of the argument `resample` specifies the number of resampling pseudoreplicates. The default of `resample` is 5.

Search strategy The calculation of the support values requires a local search, that is performed under the default settings unless the values of the following arguments are specified.

build For calculating Bremer support, the integer value of **build** specifies the number of independent Wagner tree builds per node. The integer value 0 (**build:0**) specifies that Bremer support values are calculated on the starting trees currently in memory, rather than on newly generated trees. The initial trees for calculating Bremer support can also be imported using the argument **of_file** of the command **build** (Section 3.3.1).

For calculating Jackknife and Bootstrap supports, it specifies the number of Wagner tree builds per pseudoreplicate. Single best trees from all pseudoreplicates are used to calculate the support values. If multiple best trees are recovered in a pseudoreplicate, **one is selected (at random? the first one?)**. If **build** is omitted from the argument list of **calculate_support**, a single random addition Wagner tree per pseudoreplicate is built by default. This is equivalent to **build(trees:1, randomized)**. See **build** (Section 3.3.1) for a detailed discussion of arguments of the command **build**.

swap Specifies the method and parameters for local tree search. If searching parameters are not specified, the search is performed under the default settings of **swap** (Section 3.3.26).

Defaults

calculate_support(bremer, build(trees:1, randomized), swap(trees:1)) By default POY5 will calculate the bremer support for each tree in memory node by node. However, if no trees are stored in memory, executing the command **calculate_support()** does not have any effect.

Examples

- **calculate_support(bremer)**
Calculates Bremer support values by performing independent searches for every node for every tree in memory. This is equivalent to executing **calculate_support()** (the default setting.)
- **calculate_support(bremer, build(trees:0), swap(trees:2))**
Calculates Bremer support values by performing swapping on each tree in memory for every node and keeping up to two best trees per search round.
- **calculate_support(bremer, build(of_file:"new_trees"), swap(tbr, trees:2))**

Calculates Bremer support values by performing TBR swapping on each tree in the file `new_trees` located in the current working directory for every node and keeping up to two best trees per search round.

- **`calculate_support(bootstrap)`**

Calculates Bootstrap support values under default settings. This command is equivalent to `calculate_support(bootstrap:5, build(trees:1, randomized), swap(trees:1))`.

- **`calculate_support(bootstrap:100, build(trees:5), swap(trees:1))`**

Calculates Bootstrap support values performing one random resampling with replacement, followed by 5 Wagner tree builds (by random addition sequence) and swapping these trees under the default settings of the command `swap`, and keeping one minimum-cost tree. The procedure is repeated 100 times.

- **`calculate_support(jackknife:(resample:1000), build(), swap(tbr, trees:5))`**

Calculates Jackknife support values randomly removing 36 percent of the characters (the default of `jackknife`), building 10 Wagner trees by random addition sequence (the default of `build`), swapping these trees using `tbr`, and keeping up to 5 minimum-cost tree in the final swap per swap (totaling up to 50 stored trees per replicate). The procedure is repeated 1000 times.

See also

- [report](#) (Section 3.3.19)
- [supports](#) (Section 3.3.19)
- [graphsupports](#) (Section 3.3.19)

3.3.3 `clear_memory`

Syntax

```
clear_memory([argument list])
```

Description

Frees unused memory. Rarely needed, this is a useful command when the resources of the computer are limited. The arguments are optional and their order is arbitrary.

Arguments

- m** Includes the alignment matrices in the freed memory.
- s** Includes the unused pool of sequences in the freed memory.

Defaults

`clear_memory()` By default POY5 clears all memory *except* for the pool of unused sequences and the matrices used for the alignments.

Examples

- `clear_memory(s)`
This command frees memory including all alignment matrices but keeping unused pool of sequences.

See also

- `wipe` (Section 3.3.30)

3.3.4 cd

Syntax

`cd(STRING)`

Description

Changes the working directory of the program. This command is useful when datafiles are contained in different directories. It also eliminates the need to navigate into the working directory before beginning a POY5 session. To display the path of the current directory, use the command `pwd` (Section 3.3.12).

Arguments

`STRING` The value specifies a path to a directory.

Examples

- `cd ("~/Users/username/docs/poyfiles")`
Changes the current directory to the directory in a Mac environment (when using a PC, the forward slashes will be replaced with backslashes.)

See also

- `pwd` (Section [3.3.12](#))

3.3.5 echo

Syntax

`echo(STRING, output class)`

Description

Prints the content of the string argument into a specified type of output. Several types of output are generated by POY5 which are specified by the “output class” of arguments (see below). If no output-class arguments are specified, the command does not generate any output.

Arguments

Output class

`error` Outputs the specified string as an error message (`stderr` in the flat interface).

`info` Outputs the specified string as an information message (`stderr` in the flat interface).

`output[:STRING]` Reports a specified string (`stdout` in the flat interface) to screen or file, if the filename string (enclosed in parentheses) is specified following `output` and separated from it by a colon, “:”.

Examples

- `echo("Building with indel cost 1", info)`
Prints to the output window in the ncurses interface and to the standard error in the flat interface the message Building with indel cost 1.

- `echo("Final trees", output:"trees.txt")`
Prints the string `Final trees` to the file `trees.txt`.
- `echo("Initial trees", output)`
Prints the string `Initial trees` to the output window in the ncurses interface, and to the standard output (`stdout` in the flat interface).

See also

- `report` (Section 3.3.19)

3.3.6 exit

Syntax

`exit()`

Description

Exits a POY5 session. This command does not accept any argument. `exit` is equivalent to the command `quit`.

NOTE

To interrupt a process without quitting a POY5 session, use Control-C. It aborts a currently running operation but keeps all the previously accumulated data in memory. It does not abort the current session permitting entering new command and continuing the session.

Examples

- `exit()`
Quits the program.

See also

- `quit` (Section 3.3.13)

3.3.7 fuse

Syntax

`fuse([argument list])`

Description

Performs Tree Fusing [16] on the trees in memory. Tree Fusing method to escape local optima by exchanging clades with identical composition of terminals between pairs of trees. Only *one* pair of trees is evaluated during a single iteration. The size of the clades being exchanged is not specified.

Arguments

keep:INTEGER Specifies the maximum number of trees to keep between iterations. By default, the number of trees retained is the same as the number of starting trees.

iterations:INTEGER Specifies the number of iterations of tree fusing to be performed. The number of iterations is effectively the number of pairwise clade exchanges. The default number of iterations is four times the number of retained trees (as specified by **keep**).

replace:argument Specifies the method for tree selection. Acceptable arguments are:

better Replaces parent trees with trees of better cost produced during a fusing iteration.

best Keeps a set of trees of the best cost regardless their origin.

The default is **best**.

swap Specifies tree swapping strategy to follow each iteration of tree fusing. No swapping is performed under default settings. See the description of the command **swap** (Section 3.3.26).

Defaults

fuse(replace:best) By default POY5 performs fusing keeping the same number of trees per iterations as the number of the starting trees. The number of iterations is four times the number starting trees. During the procedure, only the best trees are retained. No swapping is performed subsequent to tree fusing.

Examples

- **fuse(iterations:10, replace:best, keep:100, swap())**

This command executes the following sequence of operations. In the first

iteration, clades of the same composition of terminals are exchanged between two trees from the pool of the trees in memory. The cost of the resulting trees is compared to that of the trees in memory and a subset of the trees containing up to 100 trees of best cost is retained in memory. These trees are subjected to swapping under the default settings of `swap`. The entire procedure is repeated nine more times.

- **`fuse(swap(constraint))`**

This command performs tree fusing with modified settings for swapping that follows each iteration. Once a given iteration is completed, a consensus tree of the files in memory is computed and used as constraint file for subsequent rounds of swapping (see the argument `constraint` (Section 3.3.26) of the command `swap`).

See also

- `swap` (Section 3.3.7)

3.3.8 `help`

Syntax

```
help([argument])
```

Description

Reports the requested contents of the help file on screen.

Arguments

LIDENT Reports the description of the command, the name of which is specified by the LIDENT value.

STRING Reports every occurrence in the help file of the expression specified by the string value.

Defaults

`help()` By default POY5 displays the entire content of the help file on screen.

Examples

- `help(swap)`
Prints the description of the command `swap` in the *POY Output* window of the ncurses interface or to the standard error in the flat interface.
- `help("log")`
Finds every command with text containing the substring `log` and prints them in the *POY Output* window of the ncurses interface or to the standard error in the flat interface.

3.3.9 inspect

Syntax

`inspect(STRING)`

Description

Retrieves the description of a POY5 file produced by the command `save` (Section 3.3.21). If the description was not specified by the user, `inspect` reports that the description is not available. If the file is not a proper POY5 file format, a message is printed in the *POY Output* window of the ncurses interface or to the standard error of the flat interface.

POY5 files are not intended for permanent storage. They are recommended for temporary storage of a POY5 session, checkpointing the current state of the search (to avoid losing data in case the computer or the program fails), or reporting bugs. POY5 also automatically generates POY5 files in cases of terminating errors (important exceptions are out-of-memory errors).

Examples

- `inspect("initial_search.poy")`
Prints the description of the POY5 file `initial_search.poy` located in the current working directory in the *POY Output* window of the ncurses interface or to the standard error in the flat interface. If, for example, the file was saved using the command `save ("initial_search.poy", "Results of Total Analysis")`, then the output message is: `Results of Total Analysis`.

See also

- **save** (Section [3.3.21](#))
- **load** (Section [3.3.10](#))
- **cd** (Section [3.3.4](#))
- **pwd** (Section [3.3.12](#))

3.3.10 load

Syntax

```
load(STRING)
```

Description

Imports and inputs POY5 files created by the command **save**. The name of the file to be loaded is included in the string argument. All the information of the current POY5 session will be replaced with the contents of the POY5 file. If the file is not in proper POY5 file format, an error message is printed in the *POY Output* window of the ncurses interface, or the standard error in the flat interface. See the description of the command **save** (Section [3.3.21](#)) on the POY5 file and its usage.

POY5 files are not intended for permanent storage: they are recommended for temporary storage of a POY5 session, checkpointing the current state of the search (to avoid losing data in case the computer or the program fails), or reporting bugs. POY5 also automatically generates POY5 files in cases of terminating errors (an important exception is out-of-memory error).

Examples

- **load("initial_search.poy")**
Reads and imports the contents of the POY5 file `initial_search.poy`, located in the current working directory.
- **load("/Users/andres/test/initial.poy")**
Reads and imports the contents of the POY5 file `initial.poy` in the absolute path described by the argument.

See also

- `save` (Section 3.3.21)
- `inspect` (Section 3.3.9)
- `cd` (Section 3.3.4)
- `pwd` (Section 3.3.12)

3.3.11 perturb**Syntax**

```
perturb([argument list])
```

Description

Performs branch swapping on the trees currently in memory using a temporarily modified (“perturbed”) characters. Once a local optimum is found for the perturbed characters, a new round of swapping using the original (non-modified) characters is performed. Subsequently, the costs of the initial and final trees are compared and the best trees are selected. If there are n trees in memory prior to searching using `perturb`, then the n best trees are selected at the end. For example, if there are 20 trees currently in memory, 20 individual `perturb` procedures will be performed (each procedure starting with one of the 20 initial trees), and 20 final trees are produced. This command allows for movement from a local search optimum in the tree space by *perturbing* the character space (hence the name). The arguments specify the type of perturbation (`ratchet`, `resample`, and `transform`), the parameters of the subsequent search (`swap`), and the number of iterations of the `perturb` operation (`iterations`).

No new Wagner trees are generated following the perturbation of the data; the search is performed by local branch swapping (specified by `swap`). If `perturb` is executed with no trees in memory, an error message is generated. The arguments of `perturb` are optional and their order is arbitrary.

Arguments

`iterations:INTEGER` Repeats (iterates) the `perturb` procedure for the number of times specified by the integer value. The number of iterations is reported in the *Current Job* window of the ncurses interface and to the standard error in the flat interface.

ratchet[:(FLOAT, INTEGER)] Perturbs the data by implementing a variant of the parsimony ratchet [34] by reweighting characters listed in **report**(**data**). For unaligned data, **ratchet** randomly selects and reweights a fraction of sequence fragments (*not* individual nucleotides) specified by the float (decimal) value, upweighted by a factor specified by the integer value (severity). Thus, the number of sequence fragments into which the data is partitioned will impact the effectiveness of using the ratchet on dynamic character matrices. For static matrices, such as those obtained using the command **transform** (Section 3.3.27), **ratchet** randomly selects and reweights individual nucleotide positions (column vectors), as in Nixon’s original implementation [34].

Under default settings, **ratchet** selects 25 percent of characters and upweights them by a factor of 2. Unless **ratchet** is performed under default settings (that does not require the specification of the fraction of data to be reweighted and the severity value), both values must be specified in the proper order and separated by a comma. This argument is only used as an argument for **perturb**.

resample:(INTEGER, LIDENT) Resamples the data (characters or terminals) in random order with replacement. The **resample** string consists of an integer value specifying the number of items to be resampled (followed by a comma) and a lident value specifying whether characters or terminals (values **characters** and **terminals**, respectively) are to be resampled. Specifying both values is required. No default settings are available for **resample**. This command is only used as an argument of **perturb**.

swap Specifies the method of branch swapping for a local tree search based on perturbed data. If the argument **swap** is omitted, the search is performed under default settings of the command **swap** (Section 3.3.26).

transform Specifies a type of character transformation to be performed before executing a **perturb** procedure. See the command **transform** (Section 3.3.27) for the description of the methods of character type transformations and character selection.

Defaults

perturb(ratchet, swap (trees:1)) When no arguments specified, POY5 performs the ratchet procedure under default settings.

Examples

- `perturb(resample:(50,terminals), iterations:10)`
 Performs 10 successive repetitions of random resampling of 50 percents of terminals with replacement. Branch swapping is performed using alternating SPR and TBR, and keeping one minimum-cost tree (the default of `swap`).
- `perturb(iterations:20, ratchet:(0.18,3))`
 Performs 20 successive repetitions of a variant of the ratchet (see above) by randomly selecting 18 percent of the characters (sequence fragments) and upweighting them by a factor of 3. Branch swapping is performed using alternating SPR and TBR, and keeping one optimal tree (the default of `swap`).
- `perturb(iterations:1, transform (tcm:(4,3)))`
 Transforms the cost regime of all applicable characters (*i.e.* molecular sequence data) to the new cost regime specified by `transform` (cost of substitution 4 and cost of indel 3). Subsequently a single round of branch swapping is performed using alternating SPR and TBR, and keeping one optimal tree (the default of `swap`).
- `perturb(ratchet:(0.2,5), iterations:25, swap(tbr, trees:5))`
 Performs 25 successive repetitions of a variant of the ratchet (see above) by randomly selecting 20 percent of the characters (sequence fragments) and upweighting them by a factor of 5. Branch swapping is performed using TBR and keeping up to 5 optimal trees in each iteration.
- `perturb(transform(static_approx), ratchet:(0.2,5), iterations:25, swap(tbr, trees:5))`
 Transforms all applicable (*i.e.* dynamic homology sequence characters) using `transform` into static characters. Therefore, the subsequent ratchet is performed at the level of individual nucleotides (as in the original implementation), *not* sequence fragments. Thus, ratchet is performed by selecting 20 percent of the characters (individual nucleotides) and upweighting them by a factor of 5. Branch swapping is performed using TBR and keeping up to 5 optimal trees in each iteration as in the example above.

See also

- `swap` (Section 3.3.26)

- `transform` (Section 3.3.27)

3.3.12 `pwd`

Syntax

`pwd()`

Description

Prints the current working directory in the *POY Output* window of the ncurses interface and the standard error (stderr) of the flat interface. The command `pwd` does not have arguments. The default working directory is the shell's directory when POY5 started.

Examples

- `pwd()`

This command generates the following message: “The current working directory is /Users/myname/datafiles/”. The actual reported directory will vary depending on the directory of the shell when POY5 started, or if it has been changed using the command `cd()`.

See also

- `cd` (Section 3.3.4)

3.3.13 `quit`

Syntax

`quit()`

Description

Exits POY5 session. This command does not have any arguments `quit` is equivalent to the command `exit`.

NOTE

To interrupt a process without quitting a POY5 session, use Control-C. It aborts a currently running operation but keeps all the previously accumulated data in memory. It does not abort the current session permitting entering new commands and continuing the session.

Examples

- `quit()`
Quits the program.

See also

- `exit` (Section 3.3.6)

3.3.14 `read`

Syntax

```
read([argument list])
```

Description

Imports data files and tree files. Supported formats are ASN1, Clustal, FASTA, GBSeq, Genbank, Hennig86, Newick, NewSeq, Nexus, PHYLIP, POY3, TinySeq, and XML. Filenames should be enclosed in quotes and, if multiple filenames are specified, they must be separated by commas. All filenames read into POY5 should include the appropriate suffix (*e.g.* .fas, .ss, .aln). `read` automatically detects the type of the input file. `read` can use wildcard expressions (such as *) to refer to multiple files in a single step. For example, `read("biv*")` imports all data files the names of which start with `biv` or `read("*.ss")` imports all files with the extension `.ss` (given that the data files are in the current directory). Specifying filename(s) is obligatory: an empty argument string, `read()`, results in no data being read by POY5. The list of imported files and their content can be reported on screen or to a file using `report(data)`.

If a file is loaded twice, POY5 issues an error message but this will not interfere with subsequent file loading and execution of commands.

POY5 automatically reports in the *POY Output* window of the ncurses interface or to the standard error in the flat interface the names of the imported files, their file type, and a brief description of their contents. A more comprehensive report on the contents of the imported files can be requested (either on screen or to a file) using the argument `data` of the command `report` (Section 3.3.19).

NOTE

Although POY5 recognizes multiple data file formats, it does not interpret all of their contents. Instead, it will recognize and import only character data and ignore other content (such as blocks of commands, *etc.*). For certain data file formats, POY5 will interpret additional information as detailed for each file type below. It is important, however, to verify that the data was interpreted properly (using the command `report`).

The terminal names, as well as input file names, must not contain spaces, “at” or percentage symbols.

NOTE

Unlike many phylogenetic programs, POY5 does not clear the memory upon reading a second file. Instead, any subsequently read files will be added to the total data being analyzed. If a *new* taxon appears in a file, then it is assigned missing data for all previously loaded characters. If a taxon does *not* appear in a file, missing data are assigned for the characters that appear in it.

To eliminate the imported data and then to input a new data the `wipe()` command must be issued first.

NOTE

If one of the terminal names in an imported molecular file contains a space, “ ”, POY5 issues a warning. This also occurs if a taxon name appears to match a nucleotide sequence.

If one of the terminal names in an imported molecular file contains an “at” or a percentage symbols, the file will not be loaded because it may cause the program to crash when reporting results.

Arguments

Data file types To import data files, individual data file names must be included in the list of `read` arguments, enclosed in quotes, and separated by commas. If no data file types are specified, the types of the imported files are recognized automatically. To specify the data type, an additional argument explicitly denoting the data type, is included; it is followed by a colon (“：“) and the list of data file names (enclosed in parentheses), separated by commas and enclosed in quotes. This format prevents any ambiguity

in importing multiple data file types simultaneously (*i.e.* included in an argument list of a single `read`) command.

STRING Reads the file specified in the path included in the string argument.

A path can be absolute or relative to the current working directory (as printed by `pwd()`). The file type is recognized automatically.

Molecular files are assumed to contain nucleotide sequences. Valid files to read using this command are: tree files using parenthetical notation (newick, POY5 trees), Hennig86 files, Nona files, Sankoff character files as used in POY 3, FASTA files (and virtually any file generated by Genbank), and NEXUS files. Only taxon names, trees, characters, and cost regimes will be imported from each one of this files, no other commands are currently recognized.

NOTE

POY5 recognizes the characters *x* and *n* as representing any nucleotide base (a,c,t, or g). The ? symbol inserted in sequence data signifies missing data, a gap, or any nucleotide base may occur in that matrix position. For prealigned data sequence gaps are recognized by dashes.

NOTE

Continuous characters can be treated as such by assigning the lower and upper bounds of the range as polymorphic additive character states. Because additive characters are integers, such characters need to be re-scaled using the `weightfactor` of the `transform()`. Consider a continuous character `winglength`, the states of which are ranges of measurements in hundredth of a millimeter, for example 2.53-3.68 mm for a given terminal. A corresponding character state in the additive character matrix (in Hennig86 format) is [253,368]. To scale the values, a transformation is applied to the character `winglength` as follows:

```
transform((characters,names:("winglength")),
          (weightfactor:0.01)).
```

aminoacids:(STRING list) Specifies that the data listed in the string argument are amino acid sequences in FASTA format.

NOTE

Currently, IUPAC ambiguity codes for aminoacids are *not* supported and inputing files that contain aminoacid data with ambiguities results in an error message.

annotated:(STRING list) Specifies that the data listed in the string argument are chromosomal sequences with pipes (“|”) separating individual loci. This data type allows for locus-level rearrangements specified by the argument `dynamic_pam` (Section 3.3.27) of the command `transform` (Section 3.3.27). Locus homologies are determined dynamically, but based on annotated regions [47] (for a sample script using this data type see tutorial 4.7).

breakinv:(STRING, STRING, [orientation:BOOL, init3D:BOOL]) An enhancement of the data file type `custom_alphabet` allowing rearrangement events specified using `dynamic_pam()`. Syntactically, `breakinv` data type is identical to `custom_alphabet` data type. If `init3D` is set to `true`, execution will be faster, but will require very large amounts of memory.

chromosome:(STRING list) Specifies that the data in the files listed in the string argument are chromosomal sequences without predefined locus boundaries. Specifying that imported sequences are chromosome type data enables the application of parameter options that optimize chromosome-level events such as rearrangements, inversions, and large-scale insertions and deletions (including duplications). These parameter options (*e.g.* inversion cost) are specified using the command `transform` (Section 3.3.27). Unlike when using `annotated` data type, both locus-level and nucleotide-level homologies are determined dynamically [6, 46] (see tutorial 4.6). If chromosome sequences were imported as nucleotide type data, they can be converted to chromosome type data using the argument `seq_to_chrom` of `transform` (Section 3.3.27).

custom_alphabet:(STRING, STRING, [orientation:BOOL, init3D:BOOL, LIDENT]) Reads the data in the user-defined alphabet format. The first string argument is the name of a datafile that contains custom-alphabet sequences in FASTA format. The characters can be (but are not required to be) separated by spaces.

The second string argument is the name of a custom-alphabet input file that contains two parts: an alphabet itself, where the alphabet elements are separated by spaces, and a transformation cost matrix. The elements in an alphabet can be letters, digits, or both, as long as one element is not a prefix of another (“prefix-free”). For example, the following pairs of custom-alphabet elements are *not* valid because the first is a prefix of the second (which would prevent the proper parsing of an input file): AB and ABBA or 122 and 122X. The transformation cost matrix contains the rows and columns in which the positions from left to right and top to bottom correspond to the sequence of the elements as they are listed in the alphabet. An extra rightmost column and lowermost row correspond to a gap. It is important that the cost matrix must be symmetrical. An example of a valid custom alphabet input file is provided below: `customalphabet`

```
alpha beta gamma delta
0 2 1 2 5
2 0 2 1 5
1 2 0 2 5
2 1 2 0 5
5 5 5 5 0
```

In this example, the cost of transformation of `alpha` into `beta` is 2, and cost of a deletion or insertion of any of the four elements costs 5.

An example of a corresponding input file:

```
>Taxon1
alphabetagammadelta
>Taxon2
alphabetabetagammadelta
>Taxon3
alphabetabetadelta
```

The optional arguments of `custom_alphabet` include `orientation` and `init3D`, both of which require obligatory boolean values. The argument `orientation` allows the user to specify the orientation of custom-defined alphabet characters. The *tilde* symbol (“~”) preceding an alphabet character indicates the negative orientation. The options are `orientation:true` or `orientation:false`. The default option is `true`.

The argument `init3D` indicates that if program will calculate in advance the medians for all triplets of characters (a, b, c). The options are

`init3D:true` or `init3D:false`. The default option is `true`.

If the optional argument `LIDENT` is `prealigned`, the sequences are treated as aligned sequences.

The default heuristic `level` is 2. This can be changed via the command `transform`.

`custom_alphabet` can be transformed into `breakinv` using `transform`.

`genome:(STRING list)` Specifies that the data listed in the string argument are multi-chromosomal nucleotide sequences with the “@” sign separating individual chromosomes. This data type allows for the translocation of segments of sequences between chromosome

-level rearrangements specified by the argument `dynamic_pam` (Section 3.3.27) of the command `transform` (Section 3.3.27). Chromosome homologies are determined dynamically using distance threshold levels specified by the argument `chrom_hom` (Section 3.3.27) of `transform` (Section 3.3.27)(for a sample script using this data type see tutorial 4.9.

`nucleotides:(STRING list)` Specifies that the data in the list of files hold nucleotide sequences in FASTA format. The sequences can be divided in smaller fragments using a pound sign, and each fragment is visible as an individual character.

NOTE

By default, upon importing prealigned sequence data, all the gaps are removed and the sequences are treated as dynamic homology characters. To preserve the alignment the data must be imported using the `prealigned` argument of the command `read`.

`prealigned:(read argument, tcm:STRING [,gap_opening:INTEGER])`
Specifies that

the input sequences are prealigned and should be assigned the transformation cost matrix from the input file defined by the string argument. (See the argument `tcm` (Section 3.3.27) of the command `transform`.)

`prealigned:(read argument, tcm:(INTEGER, INTEGER) [,gap_opening:INTEGER])`
Specifies that

the input sequences are prealigned and should be assigned substitution and indel costs as defined by the `tcm` and `gapopening` arguments. (See the argument `tcm` (Section 3.3.27) of the command `transform`.)

Defaults

`read()` If no data files are specified, POY5 does nothing. If however, data files are listed but character type is not indicated, POY5 automatically detects data file types and interprets sequence files as nucleotides-type data.

Examples

- `read("/Users/andres/data/test.txt")`
Reads the file `test.txt` located in the path “/Users/andres/data/”.
- `read("28s.fas", "initial_trees.txt")`
Reads the file `28s.fas` and loads the trees in parenthetical notation of the file `initial_trees.txt`.
- `read("SSU*", "*.txt")`
Reads all the files the names of which start with `SSU`, and all the files with the extension `.txt`. The types of the datafiles are determined automatically.
- `read(nucleotides:(chel.fasta", "chel2.fasta"))`
Reads the files `chel.fasta` and `chel2.fasta`, containing nucleotide sequences.
- `read(aminoacids:(a.fasta", "b.fasta", "c.fasta"))`
Reads the amino acid sequence files `a.fasta`, `b.fasta`, and `c.fasta`.
- `read("hennig1.ss", "chel2.fasta", aminoacids:(a.fasta))`
Reads the Hennig86 file `hennig1.ss`, the FASTA file `chel2.fasta` containing nucleotide sequences, and the amino acid sequence file `a.fasta`.
- `read(custom_alphabet:(my_data", "alphabet"))`
Reads the first file, `my_data`, containing data in the format of a custom alphabet, which is defined in the second input file, `alphabet`. By default, the forward and reverse orientation (`orientation:true`) of custom-alphabet characters is considered and prior calculation of medians for their triplets (`init3D:true`) is performed.

- `read(annotated:(filea.txt, fileb.txt), chromosome:(filec.txt))`

Reads three files containing chromosome-type sequence data. The sequences in two files, `filea.txt` and `fileb.txt`, contain pipes (“|”) separating individual loci, whereas the sequences in the third, are without predefined boundaries.

- `read(genome:(mt_genomes, nu_genomes))`

Reads two files containing genomic (multi-chromosomal) sequence data.

- `read(prealigned:(18s.aln, tcm:(1,2)))`

Reads the prealigned data file `18s.aln` generated from the nucleotide file `18s.fasta` using the transformation costs 1 for substitutions and 2 for indels.

- `read(prealigned:(nucleotides:(*.nex), tcm:"matrix1"))`

Reads character data from all the Nexus files as prealigned data using the transformation cost matrix from the file `matrix1`.

See also

- `report` (Section 3.3.19)

3.3.15 rediagnose

Syntax

```
rediagnose()
```

Description

Performs a re-optimization of the trees currently in memory. This function is only useful for sanity checks of the consistency of the data. Its main usage is for the P0Y5 developers. This command does not have arguments.

Examples

- `rediagnose()`

See the description of the command.

3.3.16 recover

Syntax

```
recover()
```

Description

Recovers the best trees found during swapping, even if the swap was cancelled. This command functions only if the argument `recover` (Section 3.3.26) was included in a previously executed (in the current POY5 session) command `swap`. Otherwise, it has no effect.

The trees imported by `recover` are appended to those currently stored in memory.

Note that using recovered trees is not intended for temporary storage of trees. It is useful only as an intermediary operation in a given part of a POY5 session. When other commands that require clearing memory are executed (such as `build`, `calculate_support`, or another `swap`), the trees stored by `recover` can no longer be retrieved.

Examples

- `recover()`

If the command `swap` (executed earlier in the current POY5 session) contained the argument `recover`, for example, `swap(tbr,recover)`, this command will restore the best trees recovered during swapping.

See also

- `swap` (Section 3.3.26)
- `recover` (Section 3.3.26)

3.3.17 redraw

Syntax

```
redraw()
```

Description

Redraws the screen of the terminal. This command is only used in the ncurses interface, other interfaces ignore it. `redraw` clears the contents of the

Interactive Console window but retains the contents of the other windows. It does not affect the state of the search and the data currently in memory.

Examples

- `redraw()`

See the description of the command.

3.3.18 rename

Syntax

```
rename([argument list])
```

Description

Replaces the name(s) of specified item(s) (characters or terminals). This command allows for substituting taxon names and helps merging multiple datasets without modifying the original datafiles. More specifically, it can be used, for example, (1) for housekeeping purposes, when it is desirable to maintain long verbose taxon names (such as catalog or GenBank accession numbers) associated with the original datafiles but avoid reporting these names on the trees; (2) to provide a single name for a terminal in cases where the corresponding data is stored in different files under different terminal names; and (3) to simply change an outdated or invalid terminal name.

The command consists of a terminal or character identifier followed by a comma and then by either a string containing a synonymy file or a pair (or pairs) of strings containing the names of items being renamed.

NOTE

In order to change taxon names, the command `rename` must be executed *before* importing the datafiles (see command `read` (Section 3.3.14)) that contain character data the taxa to be renamed.

NOTE

Once the command `rename` is applied, subsequent commands must refer to the terminals using the new, substitute names. This is critical, for example, when importing a terminals file using the command `select` (Section 3.3.23) or specifying a root using the command `set` (Section 3.3.24).

Arguments

Identifiers The identifiers specify whether terminals or characters are being renamed. An identifier must precede the subsequent arguments.

characters Specifies that the subsequently items to be renamed are characters.

terminals Specifies that the subsequently items to be renamed are terminals.

Specifying items to be renamed These arguments allow to specify the items to be renamed either individually (by using a pair of string arguments) or in a group (by importing a *synonymy* file). The latter is useful when there are multiple items to be renamed and/or when it is desirable to substitute a single name for multiple ones.

STRING Specifies the name of the file (a *synonymy* file) that contains the list of terminals or characters to be renamed. The synonymy file has the following structure: each line contains a list of synonyms (two or more) separated by spaces. The name of the item listed first is going to be substituted for all the subsequently listed names. Consider, for example, a two-line synonymy file below:

```
alpha beta gamma
delta 1
```

When this file is imported, the items **beta** and **gamma** will be renamed as **alpha** and the item **1** will be renamed as **delta** in all subsequently imported datafiles.

(STRING, STRING) Specifies the names of individual items to be renamed. The first item is renamed as the second item: specifying ("alpha", "beta") renames the character or taxon **alpha** to **beta**. To specify multiple pairwise name substitution, several name pairs can be listed: ("alpha", "beta"), ("gamma", "delta").

NOTE

Note that when **rename** is applied by specifying pairs of synonyms in the command's argument, the substitute name is listed *second*. However, the substitute name appears *first* in a synonymy file, followed by one or more synonyms.

Examples

- `rename(terminals,"synfile")`

This command renames terminal names contained in the synonymy file `synfile` in all subsequently imported datafiles.

- `rename(terminals,("Mytilus_sp","Mytilus_edulis"))`

This command renames terminal taxon `Mytilus_sp` as `Mytilus_edulis` in all subsequently imported datafiles.

3.3.19 report

Syntax

```
report([argument list])
```

Description

Outputs the results of current analysis or loaded data in the *POY Output* window of the ncurses interface, the standard output of the flat interface, or to a file. To redirect the output to a file, the file name in quotes and followed by a comma must be included in the argument list of `report`. All arguments for `report` are optional.

Arguments

Reporting to files

`STRING` Specifies the name of the file to which all the specific types of report outputs, designated by additional arguments, are printed. If no additional arguments are specified, the data, trees, and diagnosis are reported to that file by default.

A string (text in quotes) argument is interpreted as a filename. Therefore, `"/Users/andres/text"` represents the file `text` in the directory `/Users/andres` (in Windows `C:\users\andres`). If no path is given, the path is relative to the current working directory as printed by `pwd()`. `usage`.

Terminals and characters This set of arguments reports the current status of terminals and characters from the imported data files.

compare:(BOOL, identifiers, identifiers) If the boolean argument is set to **false**, the command reports the ratios of all pairwise distances to their maximum length for the characters specified by character identifiers. If the boolean argument is set to **true**, the complement sequences for the characters specified by the second identifier are computed prior to reporting the distance.

cross_references[:identifiers[:STRING]] Reports a table with terminals being analyzed in rows, and the data files in columns. A plus sign (“+”) indicates that data for a given terminal is present in the corresponding file; a minus sign (“-”) indicates that it is not. **cross_references** is a very useful tool for visual representation of missing data.

Under default settings, cross-references are reported for all imported datafiles. To report cross-references for some of the fragments within a given file, a single character, or a subset of characters, optional arguments must be specified. A combination of a character identifier (see command **select** (Section 3.3.23)) and the file names (specified in the string value) is used to select specific datafiles to be cross-referenced. For example, if a command **cross_references:names:("file1")** is executed, the output is produced only for **file1**.

The argument **cross_references:all** generates a table that shows presence and absence of fragments contained within each file. If each datafile contains a single fragment, executing **cross_references:all** is equivalent to executing **cross_references**.

By default, the cross-reference table is printed on screen or to an output file, if specified.

data Outputs a summary of the input data. More specifically, POY5 will report the number of terminals to be analyzed, a list of included terminals with numerical identification numbers, list of synonyms (if specified), a list of excluded terminals, a number of included characters in each character-type category (*i.e.* additive, non-additive, Sankoff, and molecular) with the corresponding cost regimes, a list of excluded characters, and a list of input files. If the report is directed to a file with extension “nex” or “nexus” then the output is suitable for a nexus file (including the NEXUS header). Hennig format is produced if the report is directed to a file with extension “ss” or “hen” or “hennig”.

searchstats Outputs a summary of the results of the last search command, including the number of builds, fuse, ratchet, and the costs of the trees found.

searchstats Outputs a summary of the results of the last search command, including the number of builds, fuse, ratchet, and the costs of the trees found.

seq_stats:identifiers Outputs a summary of the sequences specified in the argument value, for all taxa. The summary includes the maximum, minimum, and average length and distance for all terminals.

terminals Reports a list and number of terminals included and excluded per input file. Use the command **select** (Section 3.3.23) for including and excluding terminals .

treestats Reports the number of trees in memory per cost.

treescosts Reports the cost of each tree separated by colons. The output contains no formatting for easy processing by any scripting language.

Trees This set of arguments outputs tree representations in parenthetical, ascii (simple text), or PDF formats. The arguments specify the types of tree outputs. They include actual trees resulting from current searches, or imported from files, their consensus trees, or trees displaying support values.

To select the root terminal in the tree representation, the command **set** (Section 3.3.24) is used.

Most analyses produce more than a single tree and it is often desirable to report only some of them. To report particular trees (for instance all optimal trees, randomly-selected trees, or all unique trees, *etc.*), first the command **select** (Section 3.3.23) must be applied to specify (select) the desired trees from all those stored in memory.

all_roots In a tree with n vertices (and therefore $n - 1$ edges), calculates the cost of the $n - 1$ rooted trees as implied by a root located in the subdivision vertex at each edge in the unrooted tree in memory.

asciitrees[:collapse[:BOOL]] Draws ascii character representations of trees stored in memory. The argument **collapse** collapses the zero length branches if the boolean value is **true** (the default); if the boolean value is **false**, the zero length branches are not collapsed.

clades Output a set of Hennig86 files. Each file, named `file.hen`, where “file” is whatever string you pass to this function contains information on each clade for one of the trees currently stored. This is similar to the utility `jack2hen` of POY3.

consensus[:INTEGER] Reports the consensus of trees in memory in parenthetical notation. If no integer value is specified, a strict consensus is calculated [39]; if integer value is specified, a majority rule consensus is computed, collapsing nodes with occurrence frequencies less than the specified integer [30]. If a value less than 51 is specified, POY5 reports an error.

graphconsensus[:INTEGER] Same as **consensus** except for consensus trees are reported in graphical format, either in the ascii format on screen or in the PDF format if redirected to a file.

graphdiagnosis Output the diagnosis in PDF format. The PDF is compressed, and contains the trees and links to see the diagnosis of each vertex in the tree.

graphsupports[:argument] This command outputs a tree with support values that have been previously calculated using the `calculate_support` (Section 3.3.2) either on screen in ascii format, or, if specified, to a file in PDF format. The argument values are the same as for `supports` (*i.e.* `bremer`, `jackknife`, and `bootstrap`).

graphtrees[:collapse[:BOOL]] If POY5 has been compiled with graphics support, it will display a window in which you can browse graphical representations of all the trees in memory. When working in this window, using “j” and “k” keys displays the previous or next tree respectively. If no graphical support available, the output is similar to that generated by the `asciitrees` argument. Pressing “q” key returns to the *Interactive Console* window. The argument `collapse` will collapse the zero length branches if true, otherwise not (default is `true`.)

supports[:argument] Outputs a newick format representation of a tree with the support values has previously been calculated using the command `calculate_support` (Section 3.3.2), either to the screen or to a file (if specified). If no argument is given, all calculated support values are printed. The arguments `bremer`, `jackknife`, and `bootstrap` specify which type of support tree to report.

To print the bremer supports of the trees in memory, using as reference trees stored in a file, `bremer` accepts an optional string argument (as in `report(supports:`

`bremer:"file.txt")`, or a list of strings enclosed in parentheses (as in `report(supports:`

`bremer:(“file1.txt”, “file2.txt”))`). The argument’s value specifies the files containing lists of trees and costs (as those generated by `visited` (Section 3.3.26)), that should be used with their annotated cost to assign the bremer support values.

To print the bremer supports of a tree that does not exist in memory (or a consensus tree) stored in a file, `bremer` accepts the value `of_file:(STRING, INTEGER, files)`, where the first argument value is the file containing the tree for which bremer supports should be computed, the second argument is the cost of the tree, and the files is that described in the previous paragraph.

If no input file is given, or if `bootstrap` or `jackknife` are requested, then the necessary information must have been calculated using `calculate-support` (Section 3.3.2).

`jackknife` and `bootstrap` accept an optional argument with two possible values: `individual`, `consensus`, or a `STRING`. `individual` reports the support value for each tree held in memory: if there are a hundred trees stored in memory, for each one, the support values for each tree are reported. `consensus` generates a “consensus” tree, with the clades that have support higher than 50 percent. `STRING` labels the branches in the input trees contained in the inputfile located in the path of the `STRING` (e.g. to assign support values to the branches of a consensus tree). The default behavior, when no `individual` or `consensus` value is provided, is `individual`.

`trees:(argument list)` Outputs the trees in memory in parenthetical notation. The argument `trees` receives an optional list of values specifying the format of the tree that has to be generated. Unless `hennig` is specified in the list of values, `trees` uses newick format in the tree output. The valid optional arguments are:

`total` Includes the total cost of a tree in square brackets after each tree.

`_cost` Include the cost in square brackets for every subtree in the tree.
(These are *not* branch lengths.)

hennig Prepends the `tread` command to the list of trees and separates them with a star; this format is suitable for Hennig86, NONA, and TNT files.

newick Outputs the trees in the Newick format, with the terminals separated with commas, and trees separated with semicolons.

nexus Outputs the trees in the NEXUS format, inside a TREE block.

NOTE

The **hennig** and **newick** arguments are mutually exclusive.

margin:INTEGER Sets the margin width of the generated trees.

nomargin Outputs the trees in a single line. This is useful for some programs (such as TreeView) that cannot read trees broken in several lines.

collapse[:BOOL] If `true`, zero length branches are collapsed (the default), but if `false`—no branches are collapsed.

If the report is directed to a file with extension “nex” or “nexus” then the output is suitable for a nexus file (trees inside a TREES block). Hennig format is produced if the report is directed to a file with extension “ss” or “hen” or “hennig”. In these two cases, all other formatting options are ignored.

Implied alignments This set of arguments outputs implied alignments [54].

fasta:identifiers The same as `implied_alignments` (Section 3.3.19) but no additional headers are added, producing a valid FASTA file. Intended for easy automation, by producing a file that other programs can read immediately.

implied_alignments[:identifiers] Outputs the implied alignments of the specified set of characters in FASTA format. The optional value of the argument specifies the characters included in the output, using the same identifiers described for the character specification in the entry for the command `select` (Section 3.3.23). If no characters are specified, then the implied alignment of all the sequence characters is generated. The output is reported on screen unless a name of an output file (in parentheses) is specified, preceding the command name and separated from it by a comma. This argument is synonymous with the argument `ia`.

`ia[:identifiers]` Synonym of `implied_alignments`.

Exporting static homology data The following commands export the static homology characters currently in memory.

`nexus` Produces a file in the NEXUS format that contains all the characters currently in memory. In order to export an implied alignment as a NEXUS file, the characters must first be transformed into static characters using the `transform` command (see the Hennig86 example in tutorial 4.2):

```
transform ((all, static_approx))
report ("report.nexus", nexus, trees:(nexus))
```

`phastwinclad` Produces a file in the Hennig86 format that contains the additive and nonadditive characters currently in memory. In order to export an implied alignment as a Hennig86 file, the characters must first be transformed into static characters using the `transform` command (see example in tutorial 4.2):

```
transform ((all, static_approx))
report ("report.ss", phastwinclad)
```

NOTE

To generate a file that contains implied alignments only for a subset of fragments, an identifier must be included in the argument list of `transform`. For example,

```
transform ((names:( "fragment_1", "fragment_2"),
static_approx))
report ("myfile.ss", phastwinclad)
```

will produce Hennig86 files only for `fragment_1` and `fragment_2`. The resulting file can be imported into other programs, such as WinClada. This is equivalent to the `phastwincladfile` command in POY3.

diagnosis This set of arguments will output the diagnosis.

diagnosis Outputs the diagnosis of each tree on screen or redirects it to a file, if specified. If the extension *.xml* is appended to the name of the output file, the diagnosis is reported in XML format, rather than in simple text format.

Other arguments

ci Calculates the ensemble consistency index (CI [9, 27]) for additive, nonadditive, and Sankoff characters. Dynamic homology characters are ignored in calculating the CI, therefore, the dynamic homology characters must be converted to static homology characters using the argument `static_approx` of the command `transform` (Section 3.3.27).

memory Reports on screen, the statistics of the garbage collector. For a precise description of each memory parameter, see the Objective Caml documentation.

ri Calculates the ensemble retention index (RI; [9]) for additive, nonadditive, and Sankoff characters. Dynamic homology characters are ignored in calculating the RI, therefore, the dynamic homology characters must be converted to static homology characters using the argument `static_approx` of the command `transform` (Section 3.3.27).

script_analysis:STRING Reports the order in which commands listed of the imported script (specified by the string argument) are going to be executed. Unlike executing individual commands interactively, when commands are submitted in a script, POY5 determines the logical interdependency of operations and processes the commands in the order that yields the same results as if they were executed sequentially. This substantially optimizes parallelization and reduces memory consumption.

The colored output in the *POY Output* window of the ncurses interface facilitates reading the output of `script_analysis`: red lines mark hard constraints that allow neither parallelization nor memory optimizations, blue lines mark constraints that allow the program to pipeline commands in parallel, and green lines mark fully parallelizable commands. When POY5 is compiled with parallel off, all the operations are sequential, therefore, each potentially parallel operation is done as sequential repetitions of the subscripts described in the output of the command, reducing memory consumption.

timer:STRING Reports the value and the user time (in seconds) elapsed between two consecutive timer reports. The string value provides a label (typically a textual description) that precedes the time report in the output produced. The first timer report displays the time elapsed since the beginning of the POY5 session. This command is useful for monitoring the execution time of specific tasks.

xslt:(STRING, STRING) Applies a user-defined xslt stylesheet to the XML output. The first string is the filename of the output, the second string is the name of the stylesheet requested to generate it.

NOTE

Extensible Stylesheet Language Transformations (XSLT) are used for the transformation of XML output into other formats. Because the XML output contains all the information regarding data and trees, using XSLT stylesheets greatly expand the capabilities of POY5 to use and display results. Examples of potential applications includes graphical display of trees with proportional branch lengths, integration of tree topologies with geographical coordinate data for spatial mapping, and generating input files for other programs.

Defaults

report(data, diagnosis, trees) By default, POY5 will print on screen the following items: the tree(s) in parenthetical notation with corresponding tree cost(s), diagnosis of each tree, and a graphical representation on the tree(s) in ascii format. This output can be re-directed to a file by specifying a file name enclosed in quotation marks, for example: `report("filename")`.

Examples

- `report("my_results")`

This commands outputs the data, trees, and diagnosis (the default) to the file `my_results`. Because no path is specified, the file is located in the current working directory.

- `report(data)`

This command displays on screen a list of included and excluded terminals, their names and codes, gene fragments, file names, and other relevant data.

- **report(treestats)**

This example displays on screen the costs of all trees in memory and the number of trees for each cost.

- **report("filename", treestats)**

This command outputs the costs of all trees in memory and the number of trees for each cost to a file `filename`.

- **report(cross_references:names("file1", "file3"))**

This command produces a table showing presence and absence of data corresponding to all terminals contained in files `file1` and `file3`. Because an output file is not specified, the table is displayed on screen.

- **report("taxa", terminals)**

This command generates a file `taxa` that contains the lists and numbers of excluded and included terminals for each of the previously imported datafiles.

- **report(trees)**

This command displays on screen the trees in memory in parenthetical notation with zero-length branches collapsed and terminals separated by spaces.

- **report(trees:(total))**

This command produces the same output as the example above but also includes the total tree cost in square brackets following each tree.

- **report("filename", trees:(collapse:false, newick))**

This command produces a file `filename` that contains all trees in Newick format with zero-length branches *not* collapsed.

- **report("filename", graphtrees)**

This command saves all trees in memory in PDF format to the file `filename.pdf`.

- **report(asciitrees, "file1", trees:(newick, nomargin), "file2", graphtrees)**

This command displays a tree in ascii format on screen and outputs to `file1` trees with zero-length branches collapsed in Newick format in a single line (using no margin, the format compatible with *TreeView*). It also writes to `file2` the graphical representation of these trees in PDF format.

- **report("hennig.ss", phastwinclad, trees:(hennig, total))**
This command outputs all the static homology characters, including their cost regime, in the file `hennig.ss`; then append to the same file the trees currently in memory using the Hennig format, including the total cost of each tree in square brackets. The generated `hennig.ss` is compatible with NONA, TNT, and Hennig86.
- **report("my_results", data, diagnosis, consensus, consensus:75, "consensus", graphconsensus)**
This command reports the requested types of outputs (*i.e.* reports on the data, diagnosis, and strict consensus and 75 percent majority-rule consensus trees in parenthetical notation) to the file `my_results`. It also outputs a strict consensus tree to the file `consensus`.
- **report(graphsupports, "bremertree", graphsupports:bremer)**
This command reports on screen all previously calculated support values placed at the nodes of ascii trees and outputs to file the `bremertree` only the tree(s) with bremer support values.
- **report(implied_alignments)**
This command reports the implied alignments for all dynamic homology characters on screen.
- **report("align_file", ia:names:(SSU, "LSU"))**
This command generates the file `align_file` that contains the implied alignments only for characters contained in datafiles SSU and LSU.
- **report("script1_analysis", script_analysis:"/users/datafiles/script1.poy")**
This command produces the file `script1_analysis` that lists the commands from the input script file `script1.poy` in the order that optimizes parallelization and memory consumption. In this example the complete path (`/users/datafiles/script1.poy`) is provided, which is not necessary if the directory containing the file `script1.poy` has already been assigned using the command `cd` (Section 3.3.4) in the same P0Y5 session.
- **report("swapping", timer:"swap end")**
This command generates the file `swapping` that contains the string `swap end` followed by the number of seconds (in decimals) elapsed since the execution of the previous `timer` argument.

- `report("new_tree_diagnosis.xml", diagnosis)`

This command reports the diagnosis to the `new_tree_diagnosis.xml` file in XML format.

See also

- `calculate_support` (Section 3.3.2)

3.3.20 run

Syntax

`run(STRING)`

Description

Runs POY5 script file or files. The filenames must be included in quotes and, if multiple files are included, they must be separated by commas. The script-containing files are executed in the order in which they are listed in the string argument. Executing scripts using `run` is useful in cases when operations take long time or many scripts need to be executed automatically, for example, when conducting sensitivity analysis[48]. There are no default settings of `run`.

NOTE

Note that if any of the scripts contain the commands `exit()` or `quit()`, POY5 will quit after executing that file. Therefore, if multiple files are submitted, only the last one must contain `exit()` or `quit()`.

Examples

- `run("script1", "script2")`

This command executes POY5 command scripts contained in the files `script1` and `script2` in the same order as they are listed in the list of arguments of `run`.

See also

- `exit` (Section 3.3.6)
- `quit` (Section 3.3.13)

3.3.21 save

Syntax

```
save(STRING [, STRING ])
```

Description

Saves the current POY5 state of the program to a file (POY5 file). The first, obligatory string argument specified the name of the POY5 file. The second, optional string argument specifies a string included in the POY5 file, that can be retrieved using the command `inspect` (Section 3.3.9).

POY5 files are not intended for permanent storage: they are recommended for temporary storing of a POY5 session by a user, checkpointing the current state of a search to avoid loss work in case the computer or the program itself fails, or to report bugs. POY5 will also automatically generate the file in many cases when a terminating error occurs (an important exception is out-of-memory errors). The format of these files might differ among different versions of POY5; consequently, these files might not be interchangeable between all the versions of the program.

Examples

- `save("alldata.poy")`

This command stores all the memory contents of the program in the file `alldata.poy` located in the current working directory, as printed by `pwd()`.

- `save("alldata.poy", "My total evidence data")`

This command performs the same operation as described in the example above, but, in addition, it includes the string `My total evidence data` with the file `alldata.poy`, which can later be retrieved using the command `inspect` (Section 3.3.9).

- `save("/Users/andres/test/alldata.poy", "My total evidence data")`

This command performs the same operation as the command described above with the important difference that the file `alldata.poy` generated in the directory `/Users/andres/test/` instead of the current working directory.

See also

- [inspect](#) (Section 3.3.9)
- [load](#) (Section 3.3.10)

3.3.22 search**Syntax**

```
search([argument list])
```

Description

search implements a default search strategy that includes tree building, swaping using TBR, perturbation using ratchet, and tree fusing. The strategy involves specifying targets for a driven search, such as maximum and minimum execution times, maximum allowed memory consumption for tree storage, minimum number of times the shortest tree is found, and an expected cost for the shortest tree. When executing **search** using parallel processing trees are exchanged upon the completion of the command (after fusing). Because the lowest cost unique trees generated are selected and stored at the end of a **search** (defined by the user with **max_time**), aggressive use of this command in a parallel environment consists of including few sequential **search** commands that will allow the processes to exchange trees and add the pool of selected best trees to subsequent iterations of the command (see the example for parallel processing).

Trees that exists in memory prior the **search** command are included in the set of trees available for the **fuse** but are not swapped.

Arguments

constraint:STRING For a complete description see **constraint** (Section 3.3.1).

hits:INTEGER Specifies the minimum number of times that the minimum cost must be reached before aborting the search. The **hits** argument is not used in parallel processing

max_time:FLOAT:FLOAT:FLOAT Maximum total execution time for the search. The time is specified as days:hours:minutes. For example, executing the search for 1.5 days can be expressed as 1:12:00 or 1.5:00:00.

memory:LIDENT:FLOAT Specifies the maximum amount of memory allocated for the stored trees during the search per processor. POY5 attempts to consume memory within the specified limit, but it may surpass it in certain operations (most notably during the ratchet). The lident value expresses the units of memory (gb for Gigabytes and mb for Megabytes), whereas the float value specifies the actual value. Keeping memory consumption within the limit is approximate and is used as a rough guide to POY5, preventing the program from overflowing the memory. Furthermore, it is important to note that when running POY5 in parallel the maximum amount of memory specified by the user is allocated to each processor being used. Under certain circumstances, however, it might be required to use more memory to avoid program failures.

min_time:FLOAT:FLOAT:FLOAT Specifies the minimum total execution time for the search. The time is specified as days:hours:minutes. This command is useful when the number of **hits** is specified but the actual cost of the tree is unknown. In this case, POY5 performs the search for at least the time specified by this argument.

target_cost:FLOAT Specifies the upper limit for the cost of the shortest tree.

visited[:STRING] For a complete description see **visited** (Section 3.3.26). Note that this argument has a tremendous execution time cost, as printing the trees becomes a bottleneck for the application.

Defaults

search(max_time:0:1:0, min_time:0:1:0, memory:gb:2) Under default parameters, the program performs a search for at most one hour using at most 2 GB of memory. If the user does not specify the value of **max_time**, the search will be aborted after one hour.

NOTE

In order to maximize computational efficiency when using `search` in parallel processing environments the `hits` argument is ignored. However, a diverse set of trees which include the current best trees found among all the processes is desirable to improve the potential of tree fusing.

POY will *only exchange trees between processes at the end of each search command*. Therefore, to guarantee that separate processes seed each other with the best trees they have found every number of hours, it is advisable to use few successive search commands when executing the program in parallel. Each search will still be run in parallel, but after each one, trees will be exchanged between processors, to initiate each successive round of search.

Examples

- `search(hits:100, target_cost:385, max_time:1:12:13)`

This command will attempt as many builds, swaps, ratchets, and tree fusings as possible within the specified time of 1 day, 12 hours, and 13 minutes, finding at least 100 hits (whichever occurs first, the time limit or the number of hits), knowing that the expected cost of the best hits is at most 385 steps.

- For Parallel Implementation of `search`

```
search(max_time:0:6:0)  
search(max_time:0:6:0)  
search(max_time:0:4:0)
```

This series of commands will attempt as many builds, swaps, ratchets, and tree fusings as possible within the specified total time of 16 hours. Trees are exchanged among processors at the end of each `search` and the best unique trees are then selected and included in the following `search` command.

See also

- `build` (Section 3.3.1)
- `swap` (Section 3.3.26)
- `transform` (Section 3.3.27)

3.3.23 select

Syntax

```
select([argument])
```

Description

Specifies a subset of terminals, characters, or trees from those currently loaded in memory to use in subsequent analysis.

Arguments

Select terminals and characters Specifies terminals and characters to use in subsequent analysis. The arguments in this group specify whether terminals or characters are being selected. *Identifiers* are used to specify which characters or terminals are being selected (see the *Character and terminal identifiers* argument group below for the description of methods for selecting specific terminals or characters).

terminals Specifies that the subsequently listed identifiers refer to *terminals* to be selected. By default, POY5 assumes that the specification refers to terminals. For example, to analyze only those terminals listed in the file `opiliones` using the character data currently loaded in memory, use the command `select(files:("opiliones"))`. This command is equivalent to `select(terminals,files:("opiliones"))`.

When the command is executed, the list of selected terminals is printed on screen. **terminals** is only valid as an argument of commands **select** and **rename** (Section 3.3.18).

NOTE

Note that once specific terminals and/or characters are selected, the excluded data cannot be restored. To be able to reconstitute the original data set or to experiment with various character and terminal selections within a given POY5 session, use the commands **store** (Section 3.3.25) and **use** (Section 3.3.28).

characters Specifies that the subsequently listed identifiers refer to *characters* to be selected.

STRING Selects terminals listed in the file specified by the string argument.

Character and terminal identifiers *Identifiers* specify which characters or terminals are analyzed. In addition to the command `select`, identifiers are used as arguments for other commands that require selection of specific terminals or characters, such as commands `report` (Section 3.3.19) and `transform` (Section 3.3.27).

all Specifies all characters or terminals.

names:(STRING list) Specifies the names of the characters or terminals.

codes:(INTEGER list) Specifies the codes of characters or terminals. The codes are unique numbers that are generated by POY5 when data files are first imported. The codes can be reported using the argument `data` (Section 3.3.19) of the command `report`. The codes are generated anew when a given data file is reloaded; therefore, they can effectively be used only within a current POY5 session.

files:(STRING list) Specifies the filename list containing lists of terminals or characters.

missing:INTEGER Selects terminals or characters to be included in the analysis based on the proportion of missing data. The integer value sets the **maximum** percentage of missing data in the analysis. Terminals or characters that have *more* missing data than that of the defined value are included in the analysis.

not missing:INTEGER Selects terminals or characters to be included from the analysis based on the proportion of missing data. The integer value sets the minimum percentage of missing data. Terminals or characters that have *less* missing data than defined by the value are included in the analysis. In effect, this selects a complement of data to the argument `missing`.

NOTE

For dynamic homology characters, the missing data refers to sequence fragments, whereas for static characters it refers to individual matrix positions. Therefore, when excluding terminals with missing data, the resulting set of selected terminals depends on the character type and might, or might not, be identical. For example, if a data file (containing sequences corresponding to a single fragment) includes a very short sequence, this sequence is not treated as missing data regardless of its length. This is because in the context of dynamic homology a fragment, rather than an individual nucleotide position, constitutes a character. On the other hand, if the same data are treated as static characters, the taxon represented by a very short sequence might be excluded if the length of the sequence exceeds the threshold defined by the value of `missing`.

static Specifies the static homology characters.

dynamic Specifies the dynamic homology characters.

not names:(STRING list) Specifies the characters or terminals other than those the names of which are listed in the string list.

not codes:(STRING list) Specifies the characters or terminals other than those the codes of which are listed in the string list.

Select trees The following arguments are used to select trees from the pool of trees currently in memory.

optimal Selects all trees of minimum cost.

best:INTEGER Selects the number of best trees specified by the integer value. Best trees are not equivalent to optimal trees because best trees can include suboptimal trees in case the value of `best` exceeds the number of optimal (minimal-cost) trees. If the number of optimal trees exceeds the value of `best`, only a subset of optimal trees (equal to the value of `best` is selected in an unspecified order).

NOTE

There is no special command in POY5 to clear trees from memory. However, selecting zero best trees using the command `select(best:0)` effectively removes all trees currently stored in memory.

within:FLOAT Selects all optimal and suboptimal trees the costs of which do not exceed the current optimal cost by the float value. For example, if the current optimal cost is 507 and the float value of `within` is 3.0, all trees with costs 507–510 are selected.

random:INTEGER Randomly selects the number of trees specified by the integer value irrespective of cost.

unique Selects only topologically unique trees (after collapsing zero-length branches) irrespective of their cost.

Defaults

select(unique, optimal) By default POY5 selects all unique trees of optimal (best) cost. The rest of the trees are removed from memory.

Examples

- **select(terminals,names:(“t1”, “t2”, “t3”, “t4”, “t5”), characters, names:(“chel.aln:0”))**

This command selects only terminals t1, t2, t3, t4, and t5 and use data only from the fragment 0 contained in the file `chel.aln`.

- **select(terminals, missing:50)**

This command excludes from subsequent analyses all the terminals that have more than 50 percent of characters missing. The list of included and excluded terminals is automatically reported on screen.

- **select(optimal)**

Selects all optimal (best cost) trees and discards suboptimal trees from memory. The pool of optimal trees might contain duplicate trees (that can be removed using `unique`).

- **select(unique, within:2.0)**

This command selects all topologically unique optimal and suboptimal

trees the cost of which does not exceed that of the best current cost by more than 2. For example, if the best current cost is 49, all unique trees that fall within the cost range 49–51 are selected.

See also

- `characters` (Section 3.3.23)
- `transform` (Section 3.3.27)

3.3.24 set

Syntax

```
set([argument list])
```

Description

Changes the settings of POY5. This command performs diverse auxiliary functions, from setting the seed of the random number generator to selecting a terminal for rooting output trees.

There is no default setting for `set` and the order of its arguments is arbitrary.

Arguments

Application settings Some generic application settings. Have no effect in the analyses themselves.

history:INTEGER Sets the size of the POY5 output history displayed in the *POY Output* window to the number of lines specified by the integer value. The size of the history must be greater than zero. This command has effect only in the ncurses interface. The default size of the output history is 1000 lines.

log:STRING Directs a copy of a partial output to the file specified by the string argument. The output includes the information in the *POY Output*, *Interactive Console*, and *State of Stored Search* windows of ncurses interface. Timers and current state of the search are not included in the log. If the log file already exists, POY5 will append the text to it; if the log file does not exist, then POY5 creates a new file. If the user would like to delete the contents of a pre-existing file, then

the argument `log:new:"logfile"` creates a new initially empty file named `logfile`.

`nolog` Stops outputting the log to any previously selected file. See the description of the argument `log` above.

`root:LIDENT` Specifies the terminal to root output trees. The terminal can either be indicated as a taxon name (a `STRING`, which must appear in quotes, such as `"Genus_species"`) or the code, that is automatically assigned to the taxon by POY5 at the beginning of each POY5 session (for example, `set(root:45)`). The codes can be obtained using the command `report(data)`). The terminal codes, however, are unique only within a current session.

`timer:INTEGER` Specifies the lapse of time in seconds that should have passed between reporting the total execution time of a swap and build command. If the timer is set to 0, then no time messages are generated.

Cost calculation These arguments set the tree cost estimation routines and are applied to all character types. The arguments are mutually exclusive: only the argument of `set` specified last is used.

`normal_do` Applies a standard Direct Optimization algorithm for the tree cost estimation. This is the default and fastest technique.

`exhaustive_do` Applies a standard Direct Optimization algorithm for the tree cost estimation [49, 53]. The difference with `normal_do` is that the calculation of the tree costs during a search is much more intense, always looking for the best possible alignment for every single topology (instead of a lazy and greedy strategy used by the `normal_do`).

`iterative:LIDENT` Applies the Iterative Pass optimization for the tree cost calculations. There are two forms of iterative pass: if the argument value is `exact` (the default), then a complete three dimensional alignment is computed, as described for POY 3.0 in [55]. Otherwise, if the argument value is `approximate`, then the iterations approximate the three dimensional alignment using pairwise alignments. For chromosome characters, the three dimensional alignment for three chromosomes is computed by separating into smaller three dimensional alignments of short sequences whose lengths are smaller or equal to `max_3d_len`. The default value is 200.

If the argument value is `exact`, this method improves the tree cost estimation but at the expense of a tremendous execution time. When approximate the execution time footprint is much smaller, and far less memory consumption. A heuristic strategy is to apply `iterative` at the very end of an analysis to polish the final set of trees and perform a final search.

Both arguments accept an optional integer, stating the maximum number of iterations that can be performed. If no integer is given, then the procedure iterates until no further tree cost improvement can be made.

NOTE

Due to the complexity of heuristics of the iterative pass optimization, there is no guarantee that the tree cost recovered from the search would be exactly the same as produced by the diagnosis of the same tree. However, the cost of the tree found during the search can be verified by outputting the medians from the diagnosis (see the description of the argument `diagnosis` (Section 3.3.19)) of the command `report` and determining edge costs by hand. The cost of the tree found during the search might differ from that obtained by the rediagnosing the same tree (see `rediagnose` (Section 3.3.15)), but will recover the same tree cost in subsequent rediagnoses.

Randomized routines

`seed:INTEGER` Sets the seed for the random number generator using the integer value. If unspecified, POY5 uses the system's time as seed.

NOTE

It is critical to set a seed value to insure reproducibility of the results of the analyses that require randomization routines (such as tree building).

Defaults

`set(history:1000, normal_do)` Under default settings the size of the history buffer is limited to 1,000 lines, the Direct Optimization is used for tree cost calculation, and the current time is used to specify the seed.

Examples

- `set(history:1500, seed:45, log:"mylog.txt")`
 This command increases the size of the history in the ncurses interface to 1,500 lines, sets the random number generator to 45, and initiates a log file `mylog.txt`, located in the current working directory.
- `set(root:"Mytilus_edulis")`
 This command selects terminal `Mytilus_edulis` as a root for the output trees.
- `set (iterative:exact)`
 Turns on the iterative exact algorithm in all the nucleotide sequence characters. The program will iterate on each vertex of the tree until no further tree cost improvements can be made.
- `set (iterative:approximate:2)`
 Turns on the iterative approximate algorithm in all the nucleotide sequence characters. The program will iterate either two times, or until no further tree cost improvements can be made, whichever happens first.
- `set (iterative:exact:2)`
 Same as the previous, but using the exact algorithm instead.

See also

- `report` (Section 3.3.19)

3.3.25 store

Syntax

`store(STRING)`

Description

Stores the current state of POY5 session in memory. The stored information includes character data, trees, selections, *everything*. Specifying the name of the stored state of the search (using the string argument) does *not*, however, generate a file under this name that can be examined; the name is used only to recover the stored state using the command `use`.

In combination with `use`, the command `store` is extremely useful when exploring alternative cost regimes and terminal sets within a single POY5 session.

Arguments

`STRING` Specifies the name of the stored search state of the current POY5 session.

Examples

- `store("initial_tcm")`
`transform(tcm:(1,1))`
`use("initial_tcm")`

The first command, `store`, stores the current characters and trees under the name `initial_tcm`. The second command, `transform`, changes the cost regime of molecular characters, effectively changing the data being analyzed. However, the third command, `use`, recovers the initial state stored under the name `initial_tcm`.

See also

- `use` (Section 3.3.28)
- `transform` (Section 3.3.26)

3.3.26 swap

Syntax

`swap([argument list])`

Description

`swap` is the basic local search function in POY5. This command implements a family of algorithms collectively known in systematics as branch swapping and in combinatorial optimization as hill climbing. They proceed by clipping parts of a given tree and attaching them in different positions. It can be used to perform a local search from a set of trees loaded in memory.

Swapping is performed on all trees in memory. During a search, `swap` can collect information about the visited trees and perform various kinds of checkpoints to reduce information loss in case if POY5 crashes.

`swap` is also used as an argument for other commands to specify a local search strategy in other contexts, for example, in calculating support values using the command `calculate_support` (Section 3.3.2).

All arguments of `swap` are optional and their order is arbitrary. The argument of different groups can be combined to tune the search heuristics, but the arguments within each group of are mutually exclusive. (If more than one arguments of one group is listed, only the last one is executed.)

Arguments

Neighborhood A neighborhood is a subset of topologies reachable from a given one by a given search method. The basic standard procedures for local search in phylogenetic analysis are SPR and TBR [43]. The nearest-neighbor interchanges (NNI) swapping strategy is implemented by combining the arguments `spr` and `sectorial` (see *Join method* group of arguments): `swap(spr, sectorial:1)`.

alternate Performs `spr` and `tbr` swapping iteratively until a local optimum is found. This is a specific strategy of performing `tbr`, as the trees visited by `spr` are a subset of those visited by `tbr`.

spr[:once] This argument performs `spr` swapping, starting from the current trees in memory and subsequently repeating the SPR procedure until a local optimum is found. If the optional value `once` is specified, `spr` stops once the first tree with better cost is found.

tbr[:once] This argument performs `tbr` swapping, starting from the current trees in memory and subsequently repeating the TBR procedure until a local optimum is found. If the optional value `once` is specified, `tbr` will stop once the first tree with better cost is found.

Trajectory The following arguments define the direction of the search in the defined neighborhood.

around Changes the trajectory of a search by completely exploring the neighborhood of the current tree in memory and choosing the best swap position before continuing. The default in POY5 is to choose the first one available that shows a better cost than the current best cost.

annealing:(FLOAT, FLOAT) Uses simulated annealing [26]. If the argument's value is (a, b) , POY5 accepts a tree with cost c when the

best known tree has cost d with probability $\exp(-(c-d)/t)$, where $t = a \times \exp -i/b$ and i is the number of tree evaluated in the local search.

drifting: (FLOAT, FLOAT) Uses POY5 drifting function [16]. If the argument's value is (a, b) , then POY5 always accepts a tree with better cost than the current best cost, with probability a a tree with equal cost, and with probability $1/b + d$ a tree with cost d greater than the current best cost.

Branch break order During the local search, a branch is broken and local branch swapping is performed (see *Neighborhood* group of arguments), the precise choice of which branches should be broken first can affect both the speed and the local optimum found by the program. The following arguments select among the different strategies available in POY5.

once Breaks each branch only once during a local search; that is, if a broken branch does not yield a better tree, it is never broken again, no matter how many changes occur along the search trajectory.

randomized Chooses branches uniformly at random for breakages.

distance Gives higher priority to those branches with the greatest length.

Join method After breaking a tree (using SPR or TBR), the following arguments control the selection of the positions to join the broken clades.

constraint [:INTEGER | (depth:INTEGER, file:STRING)] Sets constraints on the join locations during the search using both a tree and an optional maximum distance from the break branch. Only sets defined either in the input file, or in the strict consensus of the files in memory to consider during swapping. An integer value of **depth** specifies the maximum distance from the break branch to attempt joins. The string value for **file** specifies an input file containing a single tree that defines topological constraints. Under default settings, **constraint** will use a consensus tree from the files in memory and perform swapping with the value of **depth** set to 0 (no maximum distance is specified).

all [:INTEGER] Turns off all preference strategies to make a join, simply try all possible join positions for each pair of clades generated after a break, in a randomized order.

sectorial[:INTEGER] Join in edges at distance equal or less than the value of the argument from the broken edge, where the distance is the number of edges in the path connecting them. If no argument is given, then no distance limit is set.

Reroot order During TBR, the following options control the order of the rerooting.

bfs[:INTEGER] Reroots using breath first search [4] from the broken edge, within the arguments value distance from the root of the clade. If no value is given, there is no limit distance for the rerooting. By default, **bfs** is used with no limit distance for the rerooting.

Trajectory samples During the search, POY5 visits a large number of trees. For some applications it might be desirable to collect information about the trees examined during a search: for example, to provide backups of the state of a search (in an unlikely crash), or to examine the characteristics of the alignments. The difference from the **swap** arguments is that the user can choose any combination of trajectory samples, and that can be used during the search. None of the trajectory samples is used by default.

recover Stores the current best tree in memory that can be recovered in case of failure. If it is necessary to recover such trees after an aborted command, use **recover** (Section 3.3.16). If the program terminates normally, the stored trees are exactly those produced at the end of the **swap**. Using **recover**, however, requires twice as much memory compared to swapping without it.

timeout:INTEGER Specifies the number of seconds after which tree branch swapping is stopped. The current best tree is the result of the swap after the timeout.

timedprint:(INTEGER, STRING) **timedprint:(n, "trees.txt")** prints the current best tree in memory to the file **trees.txt**, at least every **n** seconds. However, POY5 typically underestimates the amount of time and, therefore, the samples can be slightly sparser. **timedprint** can only be used in combination with the argument **recover**.

trajectory[:STRING] **trajectory:"better.txt"** will store every new tree found with a better score during the local search in the file **better.txt**. The string is the filename where the trajectory is to

be stored, which is optional (indicated by brackets); if not added, the trees are printed in the standard output (flat interface) or the output window (ncurses interface).

`visited[:STRING] visited:"visited.txt"` will store every visited tree and its cost during the local search in the file `visited.txt`. The (optional) string is the filename where the trajectory is to be stored. If not included, the trees are printed in the standard output (flat interface) or the output window (ncurses interface).

Character transformation

`transform` Specifies a type of character transformation to be performed *prior* to swapping. See the command `transform` (Section 3.3.27) for the description of the methods of character type transformations and character selection.

Tree selection As the tree search proceeds, a tree may or may not be selected to continue the search or to return as a result. The following arguments determine under what conditions can a tree be acceptable during the search.

`threshold:FLOAT` Sets the percentage cost for suboptimal trees that are more exhaustively evaluated during the swap, meaning that trees within the threshold are subject to an extra round of swapping. For example, if the current optimal tree has cost 450, and `threshold:10` is specified, trees with cost at most 495 are swapped. `threshold` is equivalent to *slop* of P0Y3.

`trees:INTEGER` Maximum number of best trees that are retained in a search round, per tree in memory.

Defaults

`swap(trees:1, TBR, threshold:0, bfs)` By default, current trees are submitted to a round of TBR using breadth first search and one best tree per starting tree is kept.

Examples

- `swap()`

This command performs swapping under default settings.

- **swap(trees:5)**
Submits current trees to a round of SPR followed by TBR. It keeps up to 5 minimum cost trees for each starting tree.
- **swap(transform((all, static_approx)))**
Submits current trees to a round of SPR followed by TBR, using static approximations for all sequence characters.
- **swap(trees:4, transform((all, static_approx)))**
Submits current trees to a round of SPR followed by TBR, using static approximations for all characters, keeping up to 4 minimum cost trees for each starting tree.
- **swap(constraint:(depth:4))**
Calculates a consensus tree of the files in memory and uses it as constraint file, then joins at distance at most 4 from the breaking branch. This is equivalent to `swap(constraint:(4))`.
- **swap(constraint:(file:"bleh"))**
Reads the tree in file `bleh` and use it as constraint for the search. This is equivalent to `swap(constraint:(("bleh")))`.
- **swap(constraint:(file:"bleh", depth:4))**
Uses the tree in the file `bleh` as a constraint tree and joins at distance at most 4 from the breaking branch during the swap.
- **swap(recover, timedprint:(5, "timedprint.txt"))**
Saves the current best tree to file `timedprint.txt` every 5 seconds.

See also

- [transform \(Section 3.3.27\)](#)

3.3.27 transform

Syntax

`transform([argument list])`

Description

Transforms the properties of the imported characters from one type into another type. This includes changing in costs for indels and substitution,

modifying character weights, converting dynamic into static homology characters, and transforming nucleotide into chromosomal characters among other operations.

The essential arguments of the command `transform` include identifiers and methods. The methods specify what type of transformation is applied to the set of characters specified by identifiers as defined in the description of the command `select` (Section 3.3.23). Identifiers and methods are included in parentheses and separated by a comma. It is important to remember that only identifiers of *characters* (such as `names`, `codes`, among others) can be used. The parentheses separate these essential arguments from all other optional arguments that might be included in the list. Thus, if only identifiers and methods are specified, the argument list of `transform` is included in double parentheses. For example, the command `transform((all, gap-opening:1))` contains only an identifier (`all`) and a method (`gap-opening`). Minimally, only methods can be specified; in that case, the transformation is applied to all characters to which the transformation method can be applied and only a single set of parentheses is used. For instance, `transform(gap-opening:1)`, where `gap-opening` defines the transformation method.

There are no default values for `transform`, that is if no methods are specified (`transform()`), the command does nothing.

Arguments

Identifiers Identifiers specify which characters are transformed. Only identifiers of characters (*not* terminals) can be used. If identifiers are omitted, the transformation to is applied to all applicable characters. For example, `transform((all,tcm:(1,1)))` is equivalent to `transform((tcm:(1,1)))`. See the command `select` (Section 3.3.23) for detailed description of identifiers.

Methods This set of arguments specifies different transformations that can be applied to selected characters. If multiple transformation methods are applied sequentially in the same list of arguments, the effect of the methods listed earlier might be altered or canceled by methods listed after that. Thus, caution must be used in designing complex strategies with multiple character transformations. See the note on command order (Section 3.2).

`auto_static_approx` Evaluates each selected fragment and, if the number of indels appear to be low and stable between topologies, then the

character is transformed to the equivalent character using static homologies with the implied alignment [54]. This method greatly accelerates searching and is applicable only to nucleotide sequences under dynamic homology analysis.

auto_sequence_partition Evaluates each fragment and if a long region appears to have no indels, then the fragment is broken inside that region. Any number of partitions can occur along a fragment. Fragmenting long sequences greatly accelerates searching. This method is applicable only to dynamic homology characters, and requires a tree in memory.

direct_optimization Transforms the characters specified so that the initial assignment of sequences to the internal vertices of a tree use direct optimization [49]. This method is recommended for small alphabets (less than 7 elements). Otherwise fixed states is recommended. It is only applicable to dynamic homology characters.

fixed_states[:(STRING, LIDENT)] Transforms the characters specified in fixed state characters [50] where the initial assignment of sequences to the internal vertices of the tree is one of the observed sequences. If the observed sequences contain ambiguities, only those that resolve closest to another sequence are added to the set of valid states. This method is recommended for large alphabets (≥ 7 elements). It is only applicable to dynamic homology characters. The optional arguments may be specified for use with the **chromosome** character type. The STRING specifies the root name of MAUVE [6] readable files (“STRING _*i.j.alignment*” where *i* and *j* are median states) to generate graphical sequence maps between median states. If LIDENT is **ignore_polymorphism**, new potential median states will not be generated through resolution of sequence ambiguities (reducing execution time), otherwise new medians are generated as mentioned above.

gap-opening:INTEGER Sets the cost of opening a block of gaps to the specified value. Note that this cost is in addition to the standard cost of the insertion as specified by a given transformation cost matrix. The default in POY5 is not to have extension gap cost (gap-opening:0). If the gap opening cost is a , and $indel(x)$ is the cost of inserting (or deleting) a base x according to the tcm assigned to the character, the total cost of inserting (or deleting) the sequence $s[0..n]$ is $a + indel(s[0]) + indel(s[1]) + \dots + indel(s[n-1]) + indel(s[n])$. This method

is applicable only to dynamic homology characters with the nucleotide alphabet.

level:(INTEGER , LIDENT) The integer argument specifies the heuristic level in median sequence calculation. This determines the number of possible states stored at each median sequence position. For nucleotide data, all possibilities are stored. The default for amino acid sequences is 1 and for `custom_alphabet` characters it is 2. Storage and set up time increase combinatorially with level number. If the LIDENT is `first`, ties are broken (if the number of equally costly states is greater than the level number) by choosing the first median state examined; if `last`, the last state, and if `at_random` then at random.

The default choice method is `first`.

multi_static_approx Calculates the implied alignment for each tree in memory and convert them to static homology characters using the alignment's cost regime. The new character set will be the union of all those characters generated for all the trees [58]. This option is intended only for heuristic search purposes and is applicable only to dynamic homology characters.

partitioned:LIDENT Similar to `auto_sequence_partition`, with the difference that no tree is required prior the partition, and large sequence length variations at the ends of sequences are treated as missing data. If LIDENT is `clip` then a large difference in length in the sequence ends is assumed to be caused by missing data. For `clip` to take effect, at least two fragments must be found. If LIDENT is `noclip`, then sequence ends are treated like any other fragment.

prealigned Treats the sequences as prealigned and uses the cost regime according to the specified transformation cost matrix. All other cost parameters are ignored (including affine gap costs). This command requires that all the specified sequences have the same length.

search_based:(STRING, STRING) Transforms the optimization of `fixed_states` characters to `search_based` optimization [56] by adding the sequences found in the file specified by the second argument to the character specified by the first argument.

sequence_partition:INTEGER Partitions the sequences in the argument's value number of fragments of roughly the same length. This method is applicable only to dynamic homology characters.

`static_approx[:LIDENT]` Transforms the sequences to the static homology characters corresponding to their implied alignments and their transformation cost matrix [54]. The resulting characters and their number will vary depending on the characteristic of transformation cost matrix assigned to each sequence. For example, if the cost of both substitutions and indels is 1, then one non-additive character is created per each homologous position in the implied alignment. If the cost of substitutions is 1 and the cost of indels is 2, then one character is created for each homologous position, and one extra character for each homologous position with gaps. In more complex cases, a Sankoff character is created.

The LIDENT value `remove` excludes all uninformative characters information (except autapomorphies), whereas the value `keep` retains these characters. The default is `remove`. This method is applicable only to dynamic homology characters. If a non-metric transformation cost matrix is in use, this transformation will assume that the non-metricity is due to the individual insertion and deletion cost.

NOTE

The transformation of dynamic into static homology characters cannot be reverted. Therefore, caution must be taken when the transformation is applied. For example, if sequence characters have been transformed into static characters at top level using the command `transform((all, static_approx))`, all commands executed subsequently will be applied to the transformed data. However, if the transformation has been applied *within* another command (as an argument of `swap`, for instance, `swap(transform((all, static_approx)))`), the characters will be transformed only for that specific operation.

NOTE

It is important to remember that the local optimum for the dynamic homology characters can differ from that for the static homology characters based on the same sequence data. Therefore, performing additional searches on the transformed data (for example, in calculating support values based on individual nucleotides rather than on sequence fragments) can produce a discrepancy in tree costs.

trailing_insertion:STRING/(INTEGER list) The tail and prepend costs specify the cost of having an insertion of each element in the alphabet at the beginning or end of a sequence. The string is the name of a file containing the cost of a trailing insertion corresponding to each of the elements in the alphabet separated by spaces. The last element in the list is the cost of the indel of a gap (should be 0). Instead of a file, the list can be the input of the argument, in the same order, separated by commas. Synonym of the argument **ti**. This method is applicable only to dynamic homology characters.

trailing_deletion:STRING/(INTEGER list) The tail and prepend costs specify the cost of having a deletion of each element in the alphabet at the beginning or end of a sequence. The string is the name of a file containing the cost of a trailing deletion corresponding to each of the elements in the alphabet separated by spaces. The last element in the list is the cost of the indel of a gap (should be 0). Instead of a file, the list can be the input of the argument, in the same order, separated by commas. Synonym of the argument **td**. This method is applicable only to dynamic homology characters.

tcm:(INTEGER, INTEGER) Defines transformation cost matrix. The first integer value specifies substitution cost, the second integer value defines indel cost. By default, the cost of substitutions and indels are both 1. Note: previous versions of POY, had the cost of an indel set to 2 (**tcm:(1,1)**).

tcm:(STRING, INTEGER) Defines the transformation cost matrix by importing a file (specified by the string value) that contains a user defined nucleotide transformation cost matrix. This method is applicable only to dynamic homology characters. The transformation cost matrix file contains five rows and columns with values listed in the following order (left to right and top to bottom): adenine, cytosine, guanine, thymine/uracil, and indel. A similar pattern is followed for amino acids where the matrix columns and rows reflect all the amino acid names in alphabetical order (read left to right and top to bottom) with the last row and column containing a gap cost. The costs must be symmetrical (that is, the cost of the A to T substitution is equal to the cost of T to A substitution). For example:

0	2	1	2	4
2	0	2	1	4

1	2	0	2	4
2	1	2	0	4
4	4	4	4	0

The integer argument specifies the heuristic level in median sequence calculation. This determines the number of possible states stored at each median sequence position. For nucleotide data, all possibilities are stored. The default for amino acid sequences is 1 and for `custom_alphabet` characters it is 2. Storage and set up time increase combinatorially with level number.

weight:argument Changes the cost of specified characters by a constant value (weight) which is specified by either a float or an integer value. This method is applicable to any character type.

weightfactor:argument Changes the cost of specified characters by a multiplicative factor (weight factor) which is specified by either a float or an integer value. This method is applicable to any characters.

level:INTEGER Defines level value of transformation matrix. the default value of level during initialization is 3. While building tree, we need to do assignment to each node of that tree. The elements of assignment are chosen from custom-alphabet. It could be a single character or combination of n characters. Level is the maximum number of n. When we choose a combination instead of a single character, we postpone our decision to later steps of building, which could result in better tree cost

Likelihood transformation methods This set of arguments enables analysis using several variants of the maximum likelihood criterion.

likelihood:(model,rates,priors,gapmode,cost) Transforms the specified characters to either static or dynamic likelihood characters. All arguments are optional; the defaults are presented in the sub-sections below.

NOTE

Dynamic likelihood for unaligned sequence data is implemented only as most parsimonious likelihood (MPL) [1]. Both maximum average likelihood (MAL) and most parsimonious likelihood are implementable on static characters.

NOTE

Dynamic likelihood characters necessarily include gaps/indels as a state in the character alphabet. See information for `gapmode`, below.

model Defines the reversible stochastic Markov model of evolution. Models in POY, as in other standard likelihood implementations, employ globally homogeneous models [22], which estimate rates using a single matrix over the entire topology. POY requires a model of character substitution to be specified during the transformation to likelihood characters.

Parameters can be specified by following the model name with a colon, left parenthesis, and a list of numbers for the particular model, followed by a closing parenthesis, (*i.e.*, k2p:(2.0))

NOTE

All matrices are normalized so the mean rate is one. This allows branch lengths to represent the expected number of transformations averaged over the sequence length.

jc69/neymen The model of Jukes and Cantor [23] under nucleotide data, or Neyman [33] in general. A single rate class. All equilibrium frequencies are constrained to be equal. This model will work for any size alphabet.

f81 The model of Felsenstein 1981 [12]. A single rate class, proportional to the equilibrium frequencies. This model will work for any size alphabet.

k2p/k80 The model of Kimura [25]. Transitions and transversions are constrained to have separate rate classes. All equilibrium frequencies are constrained to be equal.

f84 The model of Felsenstein 1984 [11]. Rates are constrained as in k2p, but are also proportional to the equilibrium frequencies.

hky85 The model of Hasegawa [19]. Three rate classes: transversions (purine / pyrimidine), and independent classes for purine-purine and pyrimidine-pyrimidine transitions. Rates are also proportional to equilibrium frequencies.

tn93 The model of Tamura and Nei [44]. Three rate classes: transversions (purine / pyrimidine), and independent classes for purine-purine and pyrimidine-pyrimidine transitions. Rates are also proportional to equilibrium frequencies.

gtr The model of Tavaré [45]. The most general reversible, stationary model of nucleotide substitution, and the one most frequently selected by information criteria. All rates constrained to be independent. This model works for any size alphabet, but note that there are $\frac{N(N-1)}{2}$ rates, where N is the alphabet size.

file An external file containing a matrix of estimated rate values.

custom An external file in which the rate class constraints of the Q-matrix are specified. Dominant diagonal entries of this Q-matrix are constrained so as to sum rows and columns to zero. A dash ('-') may be used as a placeholder in the custom Q-matrix, and entries constrained to be equal must be specified as letters (a,b,c,...) or numbers (0,1,2,...). The matrix must be symmetric. Base frequencies are multiplied through. Note that the structure of a custom matrix can be arbitrary and need not be equal to any of the above named models. Huelsenbeck *et al* [21] note that many of the unnamed variants of the 203 possible GTR submodels outperform the named models examined in commonly-used model selection programs.

rates Specifies the distribution of rates among sites. Rate variation distributions allow multipliers to be applied to separate groups of characters. This additional parameterization frequently improves observed likelihood scores. Commonly-employed distributions for parameterizing among-site variation are the discrete gamma distribution (Γ) [60] and this distribution with an additional invariant rate class, the theta distribution (Θ) [17]. This is only useful for more than two taxa.

NOTE

Under MAL the rates are naturally averaged across the discrete distribution, while under MPL we select the best rate class for the assignment of each character. Because dynamic likelihood characters' alignment matrix would select one rate class for alignment, and this is equivalent to a constant multiplication of the branch length, dynamic MPL characters do not support rate distribution. Rate distributions are automatically taken out if applied during a transform and a warning of that affect is reported.

NOTE

It should be noted that although many researchers apply both a gamma and theta in likelihood analyses, the parameter values inferred are correlated [42]. It has therefore been recommended that a proportion of invariant sites be “pseudo”-estimated by increasing the number of discrete categories for the gamma distribution to account for very low rate character groups. The mean of each class is used to define the rate for the discrete category.

gamma[: (N, α)] Apply a discrete gamma distribution of N classes, $\gamma(\alpha, \beta)$, of rate variation across characters. The mean rate for the distribution is set to 1.0 (*i.e.* $\alpha = \beta$), thus one additional parameter is estimated. The default for N is 4. If α is not given the parameter is optimized.

theta[: $(N, \alpha, \%)$] Apply a discrete gamma distribution with an additional invariable rate class (thus $N + 1$ rate classes in total). The default for N is 4. α and $\%$ are optimized if not given.

none Apply a single rate category to all sites.

priors The equilibrium frequencies of the characters in the stationary Markov process. The conditional probability of the data is multiplied by these frequencies to yield the likelihood score. In systematics applications, the observed (empirical) frequencies of the characters in the data are often used as an approximation to the likelihood-optimized frequencies. This approximation is sensible for sequences evolving

nearly neutrally, but this approximation breaks down for sequences under moderate or strong selection.

equal Constrains each of the equilibrium frequencies to be $\frac{1}{r}$, where r is the character alphabet size.

estimate Uses observed frequencies of the characters to approximate the ML-optimized equilibrium frequencies. Under dynamic likelihood, the indel equilibrium frequency is estimated to be the minimum number of indels required to align all the sequences:

$$\pi_{gap} = \frac{\sum_{i=0}^n |S_i| - |S_i|}{n|S_i|}$$

where S_i is the longest sequence.

given: (FLOAT LIST) Uses equilibrium frequencies given as a comma-separated list specified by the user.

cost Specifies the variant of relative likelihood (*sensu* [41]) used by POY to calculate the costs at each node and score the topology.

mal Specifies the standard criterion of maximum average likelihood [12], in which the likelihood is derived from the sum over all internal vertex labelings. Because the vertex assignments are marginalized, this version of relative likelihood is not compatible with simultaneous estimation of alignment and topology (dynamic likelihood).

mpl Specifies the criterion of most parsimonious likelihood [1], in which the likelihood is derived from the set of most probable state assignments at each node. Under this criterion, vertex state assignments are not marginalized, and so this variant of relative likelihood is compatible with dynamic likelihood. However, several studies have noted that parameters under this criterion may not be identifiable (citation needed) and that the method may not be consistent [32, for example].

gapmode Defines how atomistic gaps (indels) are treated in likelihood analysis.

When implemented, the POY approach to indels on static likelihood characters follows McGuire *et al* [31], but expands the implementation outside the F84 model. POY is unique among likelihood implementations both by enabling indel parameterization for any reversible,

stationary model applied to static characters and, more crucially, for implementing dynamic likelihood for truly simultaneous alignment and topology inference. Indel-aware models are scarce, but see Rivas and Eddy [38], and Redelings and Suchard [37, 36] for alternative parameterizations, and Liu *et al* [29] for an iterative approximation of simultaneous inference under the MAL criterion.

missing Specifies that indels be treated under standard assumptions, thus making analyses in POY for a given model and alignment comparable to other implementations such as RAxML and GARLi. Under this assumption, indels are not a part of the character alphabet (“missing” data, in this context), and are given probability one.

coupled Specifies that atomistic indels are a character, with rates of nucleotide-to-indel substitution constrained to be equal to one another. For example, the Q-matrix

$$Q = \begin{bmatrix} -3\alpha - \beta & \alpha & \alpha & \alpha & \beta \\ \alpha & -3\alpha - \beta & \alpha & \alpha & \beta \\ \alpha & \alpha & -3\alpha - \beta & \alpha & \beta \\ \alpha & \alpha & \alpha & -3\alpha - \beta & \beta \\ \beta & \beta & \beta & \beta & -4\beta \end{bmatrix}$$

specifies a Jukes-Cantor matrix augmented to estimate nucleotide-to-indel substitutions as a coupled parameter separate from nucleotide-to-nucleotide substitutions. In accordance with other likelihood implementations, the diagonal entries of the Q-matrix are normalized by the trace to scale their mean rate to one.

character Specifies that atomistic indels are a character without adding an additional rate class –except under GTR where we keep the unrestricted nature of the model intact, and each indel-nucleotide rate has its own class.

Chromosomal transformation methods For chromosome and genome character types, POY5 optimizes nucleotide-, locus-, and chromosome-level variation simultaneously. The arguments in this group transform nucleotide characters into chromosomal character to allow for translocations, inversions, and indel events both at the locus-level for chromosomal data and at the chromosome-level for genomic data.

The functions to calculate breakpoint and inversion distances between two sequences of gene orders are taken from GRAPPA, Genome Rearrangements Analysis under Parsimony and other Phylogenetic Algorithms [5], available at <http://www.cs.unm.edu/~moret/GRAPPA/>.

breakinv_to_custom Transforms `breakinv` character type into `custom_alpha` characters. This transformation prevents the use of rearrangement operations.

custom_to_breakinv:([`argument list`]) Transforms `custom_alpha` characters into the `breakinv` character type to allow for rearrangement operations (translocations and inversions; duplications are not currently supported). This argument is useful, for example, when `custom_alpha` characters are used to define a sequence of individual genes and one is interested in detecting potential change in their order within a chromosome. See the command `read` (Section 3.3.14) for the description on how to load the `custom_alpha` and `breakinv` character types. The optional list of arguments includes the arguments of the `dynamic_pam` that can also be specified subsequently, as a separate step (for a sample script using this character transformation see tutorial 4.8).

seq_to_chrom:([`argument list`]) Transforms nucleotide type data into chromosome type data to allow rearrangements, inversions, and locus-level indel operations. The chromosome-specific options (*e.g.* `locus_breakpoint`, `locus_inversion`, `locus_dcj`, and `locus_indel`) can be specified by the argument `dynamic_pam`. If no `dynamic_pam` values are specified, its default values are applied.

chromosome:([`argument list`]) Specifies parameters for creating chromosome- and genome-level HTUs (medians). The arguments of `chromosome` define homologous blocks within unannotated chromosome sequences using (`min_seed_length`, `min_loci_len`, and `min_rearranged_len`); specify the cost of locus-level transformation events: (*i.e.* `locus_inversion`, `locus_dcj`, or `locus_breakpoint`, and `locus_indel`); specify the cost of chromosome-level transformation events: (*i.e.* `chrom_breakpoint`, and `chrom_indel`); take into account whether the chromosome is linear or circular (`circular`); and implement a number of heuristic procedures to accelerate computations (`median`, `swap_med`, and `med_approx`). Under default settings, the pairwise distance between two chromosome segments or two chromosomes is determined

using breakpoint rather than inversion calculations and the rest of the arguments are executed under their default settings.

Original order of gene	g_1	g_2	g_3	g_4	g_5
Inversion	g_1	g_2	$-g_5$	$-g_4$	$-g_3$
Translocation	g_3	g_4	g_5	g_1	g_2
Translocation and inversion	$-g_5$	$-g_4$	$-g_3$	g_1	g_2

Figure 3.3: Examples of gene rearrangements: inversions and translocations.

g_1	g_2	g_3	g_4	g_5
g_1	g_2	$-g_5$	$-g_4$	$-g_3$

Figure 3.4: Rearrangement calculations between chromosomal or genomic data of six genes g_1, \dots, g_6 , where the rearrangement events are detected as either two breakpoints $(g_2, g_3), (g_5, g_6)$ or a single inversion (g_1, g_2, g_3) .

annotate: (LIDENT,FLOAT,FLOAT,FLOAT,FLOAT) Specifies the annotation method for unannotated chromosome sequences. If LIDENT is **mauve** the MAUVE algorithm is used [6] with float parameters quality of lcb, lcb coverage, minimum length of lcb (percent), maximum length of lcb (percentage). Percent calculations are based on the shorter of two sequences compared.

med_approx: BOOL Approximates chromosome medians using a fixed-states approach. This is most useful to accelerating tree building and searching operations for large chromosomal data sets. The boolean value **true** applies the fixed-states optimization. The default value is **false**.

locus_breakpoint: INTEGER Breakpoints are chromosomal regions affected by rearrangements. The *Locus Breakpoint* calculates the breakpoint distance [2] between two pairs of chromosomes given the cost for rearrangement specified by an integer value. The breakpoint distance calculation considers a chromosome or genome $G = (x_1, \dots, x_n)$ of n gene, wherein each gene appears exactly once

and its orientation is either positive or negative. Gene orders are altered by gene rearrangement operations: gene inversion, gene translocation, gene inversion and translocation (see Figure 3.3). The breakpoint distance takes into account rearrangements but not inversions. Given G and G' , a pair of genes (g_i, g_j) is a breakpoint if (g_i, g_j) occur consecutively in G but neither (g_i, g_j) nor $(-g_j, -g_i)$ occur consecutively in G' [40]. The breakpoint distance between G and G' is the number of breakpoints between them. Figure 3.4 shows two breakpoints between G and G' . The breakpoint can be calculated easily in linear time. This argument *cannot* be used in conjunction with `locus_inversion` or `locus_dcj`. The default value of `locus_breakpoint` is 10.

locus_dcj:INTEGER Calculates the double cut and join distance [59] between two chromosome segments given the cost for each event specified by the integer value. The inversion distance takes in consideration rearrangements and inversions. The dcj distance is normally smaller than the inversion and breakpoint distances. This argument *cannot* be used in conjunction with `locus_breakpoint` or `locus_inversion`.

locus_inversion:INTEGER Calculates the inversion distance [18] between two chromosome segments given the cost for inversion specified by the integer value. The inversion distance takes in consideration rearrangements and inversions. Given G and G' , the inversion distance between them is the number of inversions to convert chromosome or genome G into G' [18]. Figure 3.4 shows one inversion between G and G' . The inversion can be calculated in linear time. The breakpoint distance is normally larger than inversion distance and dcj distance. This argument *cannot* be used in conjunction with `locus_breakpoint` or `locus_dcj`.

locus_indel:(INTEGER, FLOAT) Specifies the cost for insertion/deletion of a chromosome segment. The integer value sets the gap opening cost (o), whereas the float value sets the gap extension cost (e). The indel cost for a fragment of length l is specified by the following formula: $o + l \times e$. The default values are $o = 10, e = 1.0$.

min_seed_length:INTEGER Specifies the minimum length of identical (invariant, completely conserved) contiguous sequence fragments during comparison between two chromosomes. The integer value of `min_seed_length` is the number of nucleotides. Correct identification of such fragments facilitates detecting chromosome

rearrangement events and accelerates other operations (such as tree building and swapping). However, if `min_seed_length` value is set too low (detecting many short fragments) or too high (such that no identical fragments are detected) the time required for subsequent searching procedures may significantly increase. The optimal `min_seed_length` value depends on the specifics of a given dataset (see Figure 3.5). The default value of `min_seed_length` is 9.

`min_loci_len:INTEGER` Creates a pairwise alignment between two chromosomes to detect conserved areas (“blocks”). However, only blocks of lengths (in number of nucleotides) greater or equal to the specified `min_loci_len` value are considered as hypothetically homologous blocks and used as anchors to divide chromosomes into fragments. Thus, increasing the value of `min_loci_len` decreases the chance of inferring small-size rearrangements (see Figure 3.5). The default value is 100.

`min_rearrangement_len:INTEGER` Two seeds are said to be *non-rearranged*, if their distance is less than the predefined threshold value set for `min_rearrangement_len`. In other words, it is unlikely that rearrangement operations can occur between two non-rearranged seeds if they are connected (see Figure 3.5). The default value is 100

`max_3d_len:INTEGER` The three dimensional alignment for three chromosomes is computed by separating into smaller three dimensional alignments of short sequences whose lengths are smaller or equal to `max_3d_len`. The default value is 200

`chrom_breakpoint:INTEGER` Calculates the breakpoint distance [2] between two sequences of multiple chromosomes given the cost for rearrangement specified by an integer value. The breakpoint distance takes into account locus rearrangements between non-homologous chromosomes (translocations) but not inversions. For further discussion on how breakpoint distance is calculated see the argument `locus_breakpoint`. The default value of `chrom_breakpoint` is 100.

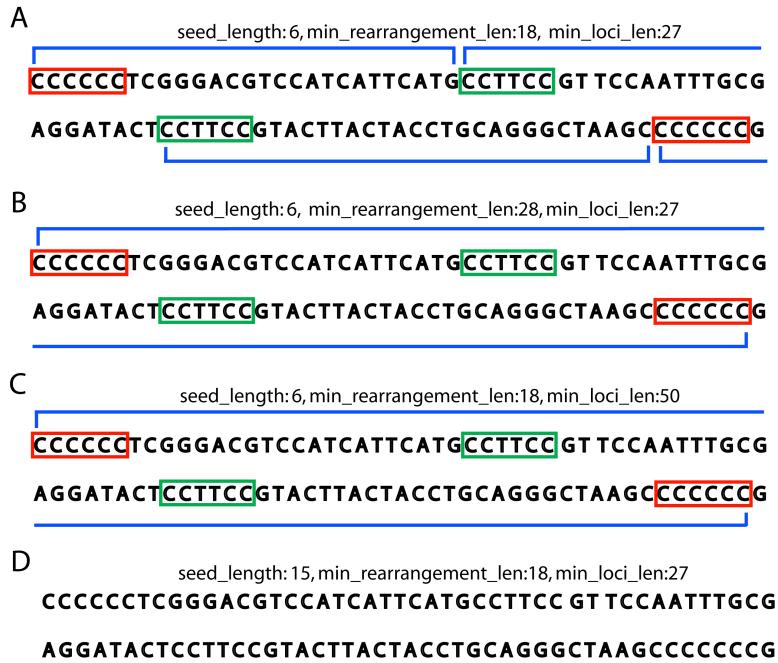


Figure 3.5: The effect of dynamic parameter arguments (`seed_length`, `min_rearrangement_len`, and `min_loci_len`) on determining homologous blocks in unannotated chromosomal sequences. *Case A* allows for rearrangements between the short homologous blocks within the blue brackets constructed upon six nucleotide seeds (red and green), *Case B* restricts rearrangement between seeds that are less than 28 nucleotides apart. *Case C* restricts rearrangement by requiring homologous blocks at least 50 nucleotides in length. *Case D* restricts recognition of homology between sequences by setting seed length of conserved seeds at 15.

NOTE

Note that the arguments `locus_breakpoint` and `chrom_breakpoint` cannot be used simultaneously with the arguments `locus_inversion` and `chrom_inversion` as they designate *alternative* methods of calculating distance between two chromosomes. If both arguments are specified, the latter will be executed. The order of other arguments of `dynamic_pam` is arbitrary.

chrom_hom:(FLOAT) Specifies the lower limit of distance between two chromosomes beyond which the chromosomes are not considered to be homologous. The default value of `chrom_hom` is 0.75.

chrom_indel:(INTEGER, FLOAT) Specifies the cost for insertion and deletion of a chromosome in analysis of multiple chromosomes. The integer value sets gap opening cost (o), whereas the float value sets gap extension cost (e). The indel cost for a fragment of length l is specified by the following formula: $o + l \times e$. The default values are $o = 10, e = 1.0$.

circular:BOOL Specifies if chromosome is circular (boolean value `true`) or linear (boolean value `false`). The default value of `circular` is `false` (linear chromosome).

median:INTEGER Specifies the number alternative locus and chromosome rearrangements of the best cost selected (randomly) for each HTU (hypothetical taxon unit) or median. Limiting the number of rearrangements stored in memory (smaller value of `median`) is heuristic strategy to accelerate calculations at the expense of thoroughness of the search. By default, only 1 rearrangement is retained (the first one found). If more than one rearrangement is specified, the selected number of rearrangements is selected in random order from the pool of all generated rearrangements.

median_solver:INTEGER Specifies the median solver. User can choose from default(`caprara`),`vinh`,`siepel`,`bbtsp`,`coaletsp`,`chainedlk` and `simplelk`. The default median solver is `caprara` median solver. Note: Concorde package is required for `chainedlk` and `simplelk` median solver. Download Concorde package from <http://www.tsp.gatech.edu//concorde/downloads/downlo>

max_kept_wag:INTEGER Defines the maximum number of Wagner-based possible ancestral sequence alignments kept to create the

next set of alignments during the pairwise alignment with rearrangement process. The default value of `max_kept_wag` is 1, however, at every step in the pairwise alignment with rearrangement process, the original order (1...n) is always considered as a potential solution.

`swap_med:INTEGER` Specifies the maximum number of swapping iterations to search for best pairwise alignment of two chromosomes taking into account locus-level rearrangement events. Limiting the number of swapping iterations accelerates the search at the expense of thoroughness. The default value is 1.

Defaults

`transform()` If no arguments are given, this command does nothing.

Examples

- `transform((all, tcm:(1,1)))`

Applies the transformation cost matrix (1,1) to all characters, meaning that substitutions and gaps receive the same weight.

- `transform((all, tcm:"molmatrix"))`

Applies the character transformation matrix "molmatrix" to all characters.

- `transform((all, tcm:(1,1)))`

This command is equivalent to `transform((dynamic, tcm:(1,1)))`.

- `transform(tcm:(1,1), gap-opening:1)`

Applies the transformation cost matrix and the gap opening cost to all characters. In this example the cost for substitutions is 1, the gap opening cost is 2 (1 set by `gap-opening` + 1 set by `tcm`), and the gap extension cost is 1 (set by `tcm`).

- `transform(tcm:(2,2), ti:(1,1,1,1,0), td:(1,1,1,1,0))`

Assigns to all characters the symmetric transformation cost matrix with cost 2 for every indel and substitution, but for those insertions and deletions at the ends of the sequences, the cost assigned will only be 1.

- **transform((static, weightfactor:2))**

This command reweights all the static homology characters by a multiplicative factor of 2, while keeping the weighting scheme that has been specified before.

- **transform((static, weight:4.2))**

Applies the same weight (a float value 4.2) to all static homology characters.

- **transform((dynamic, weight:4))**

Applies the same weight (an integer value 4) to all dynamic homology characters.

- **transform((all, tcm:(1,1)), (names:("gen1", "gen2"), static_approx), (names:("gen3"), tcm:"molmatrix"))**

Applies the substitution and indel costs 1 to all characters, then applies static approximation using that tcm to characters in files **gen1** and **gen2**, and for file **gen3**, it invokes a different transformation cost matrix, contained in the file **molmatrix**. Beware that the file name should be exactly as it was reported with **report(data)**, which differs from the actual file name (**report (data)** reports files as fileX:N).

- **transform((all, tcm:(1,1)), (names:("gen1:3", "gen2:10", "gen3:1", "gen4:5"), static_approx), (names:("gen5", "gen6"), tcm:"Molmatrix1"))**

Applies **tcm** (1,1) to all characters, then applies static approximation to the sequence data contained in files **gen1**, **gen2**, **gen3**, and **gen4** according to this transformation cost matrix, and applies the custom transformation cost matrix contained in the file **Molmatrix1** to the sequence data contained in files **gen5** and **gen6**.

- **transform(names:("mycustom.fas"),(tcm:("mycustom.mat",2)))**

This examples transforms the cost matrix used to optimize **custom_alphabet** character “**mycustom.fas**” to “**mycustom.mat**” and uses heuristic level 2.

- **transform(fixed_states)**

Transformes all sequence characters into fixed states characters.

- **transform(level:2)**

Transforms all heuristic median calculations to level 2 and breaking ties **first**

- **transform((names:(“gen1”, “gen4”), fixed_states))**
Transformed only specified sequence characters (gen1 and gen4) into fixed states characters.
- **transform(custom_to_breakinv:(circular:true))**
In this example, all custom_alphabet data is transformed into the breakinv data type and is treated as a circular chromosome.
- **transform(search_based:((“mycustom_char1.fas”, “mycustom_extra1.fas”), (“mycustom_char2.fas”, “mycustom_extra2.fas”)))**
This command adds the sequences found in “mycustom_extra1.fas” and “mycustom_extra2.fas” to **fixed_states** characters “mycustom_char1.fas” and “mycustom_char2.fas” respectively.
- **transform(search_based:“mycustom_extra.fas”)**
Adds sequences in “mycustom_extra.fas” to all **fixed_states** characters.
- **transform(seq_to_chrom:(locus_indelel:(50, 1.0), min_seed_length:15))**
All applicable (*i.e.* sequence) data are transformed into chromosome data with the minimum length of identical contiguous sequence fragments which form the seeds of homologous blocks set at 15 nucleotides and the locus-level gap opening cost is set at 50 with a gap extension cost at 1.0.
- **transform(fixed_states:(“shrimp”, ignore_polymerism))**
This example shows the optional arguments for a MAUVE-based fixed-states analysis. Here, MAUVE genome alignment files will be generated with the names “shrimp-*i.j.alignment” where *i* and *j* are median states. Sequence ambiguities will not be resolved to generate additional medians beyond those determined by the data.*
- **transform(chromosome:(annotate:(mauve,35.0,0.35,0.01,0.10)))**
This command specifies the use of MAUVE annotation with 35.0 for quality of lcb, 35% coverage of all sequences by lcb’s, 1% min length of lcb (100 for length 10,000 sequence), and 10% max length (1000 for length 10,000 sequence).

3.3.28 use

Syntax

use(STRING)

Description

Restores from memory the state of a POY5 session (that includes character data, selections, trees, all other data and specifications) that had previously been saved during the session using the command `store` (Section 3.3.25). The recalled session replaces the current session. The string argument specifies the name of the stored state.

In combination with `store` (Section 3.3.25), the command `use` is very useful for exploring alternative cost regimes and terminal sets within a single POY5 session.

Examples

- `store("initial_tcm")`
`transform(tcm:(1,1))`
`use("initial_tcm")`

The first command, `store`, stores the current characters and trees under the name `initial_tcm`. The second command, `transform`, changes the cost regime of molecular characters, effectively changing the data being analyzed. However, the third command, `use`, recovers the initial state stored under the name `initial_tcm`.

See also

- `store` (Section 3.3.25)
- `transform` (Section 3.3.26)

3.3.29 version

Syntax

`version()`

Description

Reports the POY5 version number in the output window of the ncurses interface, or to the standard error in the flat interface.

Examples

- `version ()`

3.3.30 wipe

Syntax

wipe()

Description

Eliminates the data stored in memory (all character data, trees, *etc.*).

Examples

- `wipe ()`

|||||| .working

Chapter 4

POY5 Heuristics: A Practical Guide

4.1 Introduction

As the level of phylogenetic analysis increases—from individual loci to chromosomes to genomes containing multiple chromosomes—so does computational complexity. In POY5, a significant increase in computational time results from combining in single process cladogram searching with co-optimization of nucleotide pairwise alignments, rearrangements of loci within a chromosome, and rearrangements of chromosome fragments within the genome . As a result, a phylogenetic analysis involves a set of nested computationally “hard” (NP-complete) problems that makes finding the exact solution impossible. In addition, the increasing sequence length heterogeneity (at the levels of nucleotides, loci, and chromosomes) and the ever-growing sizes of datasets further contribute to computational complexity making it impossible to obtain an exact solution in a reasonable time.

To circumvent the problem computational intractability, and, hence, the speed of the analyses, POY5 employs a battery of approximate, or heuristic, methods that function at different levels of analysis. As with all heuristic procedures, a tradeoff is involved: a substantial decrease in execution time comes at a price of obtaining possibly less accurate and less precise results (however, the extent of the tradeoff is difficult to evaluate in the analyses of real large datasets). Therefore, it becomes important to understand the combined effect of different heuristic methods, so that the chosen search strategy balances the computational time with a “reasonable” accuracy of the result.

Here we provide general guidelines for using different heuristic methods, explore their combined effect, and suggest the choice of parameters that can be explored to provide the best result for specific cases. Real datasets differ greatly in size and complexity, so that no single optimal strategy can be suggested. These guidelines, however, should enable the investigator to design an efficient strategy that will tailor to the peculiarities of a given dataset.

In addition to heuristic methods, this chapter attempts to assist with the selection of transformation cost regimes. Alternative cost regimes can significantly affect the outcome of the analysis, that becomes particularly apparent in dealing with large, genome-level datasets, where multiple cost regimes are used simultaneously to specify transformations at different levels of analysis. Most difficulties stem from selecting the most reasonable combination of parameters that affect optimization of DNA sequence data at the levels of nucleotides, loci, and chromosomes.

4.2 Data treatment

Direct optimization (see *Character optimization* section below) involves comparing all potential nucleotide homologies between two sequences. Consequently, the time it takes is proportional to the product of the lengths of the sequences compared. This procedure can be time consuming for long and greatly differing in length DNA fragments. In cases where unambiguous (such as long completely conserved regions) sequence fragments can be identified, partitioning the long sequences into smaller fragments delimited by these regions can significantly reduce computational time. Such economy is reached because nucleotide homologies are not examined over the separate partitions. This strategy assumes that the fragments are mutually exclusive and are putatively homologous across terminals.

At the level of nucleotides, individual fragments in a locus can be separated by the pound symbols (“#”) or contained as individual files (that is, treated as partitions). When “#” are used, their number must be the same across homologous sequences. Alternatively, the argument of `auto-sequence_partition` of the command `transform` (Section 3.3.26). At the chromosome level, individual loci can be separated by pipes (“|”).

Level of analysis	Heuristic	Implementation
Nucleotides	Fragment sequences	Manually separating fragments or use <code>transform(auto_sequence_partition)</code>
Locus	Fragment chromosome	Manually insert pipes separating loci
Chromosomes	NA	NA

4.3 Character optimization

Minimizing overall cladogram cost is an NP hard problem dependent on the lowest cost assignment of HTU sequences. POY implements direct optimization (DO; [49]) and fixed-states optimization (FSO; [50]) heuristics to determine the set of HTU sequences comprising the internal nodes of each cladogram constructed. Direct Optimization decomposes the problem into a series of two-node comparisons, calculating locally optimal solutions, which generates the total cladogram cost. An advantage of direct optimization is that it allows for the exploration of a large diversity of putative homologies and selects the scheme that yields the most optimal solution. This is useful in analyzing sequences of different length, where site-to-site homologies are uncertain. Because the procedure is based on a greedy algorithm, it requires multiple iterations (independent initial cladogram builds) and extensive tree searches to reach a potentially global minimum. In contrast, fixed-states optimization does not calculate HTU sequences but rather optimizes those observed in terminal taxa. These internal node sequences then are diagnosed using dynamic programming based on a matrix of edit costs between sequences. In the fixed-states implementation cladogram optimization is independent of sequence lengths, and as the number of sequences increase so to does the pool from which the HTU sequences are drawn, thereby improving cladogram cost estimation. Because of these properties fixed-states optimization is recommended as an initial approximation strategy for large data sets of variable length sequences.

Level of analysis	Heuristic	Implementation
Nucleotides	DO	Default strategy
Nucleotides	FSO	<code>transform(fixedstates)</code>
Loci	FSO	<code>transform(dynamic_-pam:(approx))</code>
Chromosomes	NA	NA

Further approximations and economies can be achieved by varying parameters of commands, such as selecting a limited subset of trees for subsequent analysis limiting the number of replicates, and examining intermediary results from an interrupted analysis.

4.4 Tree searching

The heuristic approaches to cladogram searching include random addition of taxa, branch swapping (TBR and SPR), simulated annealing (the ratchet and tree-drifting), and genetical algorithms (tree fusing). These techniques, frequently used in combination, allow a more efficient exploring of tree space and provide the means of finding more globally optimal solutions. These methods are widely used in phylogenetics [14, 57], although POY5 implements additional modifications of these procedures.

Typical search strategy in POY5 involves consecutive application of tree search algorithms that begin with generating multiple, randomly selected starting points [Random Addition Sequences (RAS) or Wagner trees]. The importance of multiple starting trees cannot be overemphasized and a successful search shall maximize the number of RAS. However, making a tree search more exhaustive by increasing the number of starting trees comes at a price of longer computation time. Therefore, it is advised here to estimate the amount of time it takes to complete a single replicate and takes this information in consideration when designing a more exhaustive strategy. The number of replicates used by POY5 practitioners for datasets of moderate size (70-100 terminals) ranges from 100 to 250. Here are some examples of search strategies:

RAS+SPR/TBR+Ratchet The strategy is for a thorough search for a data set of 100 or fewer taxa. A diversity of starting points is generated by multiple RAS, each refined by a local search (TBR or a combination of SPR and TBR, the latter is an efficient default strategy in POY5). Ratcheting is used to examine tree space that potentially has not been explored by the local searches.

RAS+SPR/TBR+Ratchet+Tree Fusing Adding tree fusing step allows for combining the best sectors of cladograms that can potentially yield a tree of shorter length. Empirical studies showed that adding tree fusing after replicate rounds enhances the results only when dealing with data sets with more than 50 taxa.

RAS+SPR/TBR+Ratchet+Tree Drifting+Tree Fusing Tree Drifting can be used in place of or in addition to the Ratchet.

Input Trees+SPR/TBR+Ratchet+Tree Drifting+Tree Fusing For more exhaustive searches, the best trees obtained from the initial searches using the strategies outlined above, can be used as input trees for subsequent analyses. In doing so, the RAS step can be omitted because searching starts with trees approximating the globally optimal tree(s).

The aggressiveness of searches can be adjusted by varying parameters of the branch swapping, ratchet, tree fusing, and tree drifting commands.

Further economies can be reached by using a combination of different character optimization methods. For example, initial searches can be conducted under faster static approximation (that converts sequence data into static homology characters; see *Character optimization* section), whereas the final refinement can be performed using direct optimization.

4.5 Chromosome heuristics

Analysis of chromosomal data requires heuristic procedures to estimate rearrangement events in addition to nucleotide transformations. Chromosomal data are divided into four different classes, namely breakinv, annotated, chromosome and genome. In the following we discuss the complexity of each class.

4.5.1 Breakinv character

Breakinv character is the most simple form of chromosome data. Each breakinv character (chromosome) is presented by a sequence of general alphabet characters each codes for one gene. The transformation cost matrix among characters is calculated in advance and provided by users. During the tree search, *pairwise alignments with rearrangements* (PAR) between two breakinv characters is constructed. The PAR generalizes the ordinary pairwise alignment by allowing rearrangements of character order. Since an exact solution for the PAR problem is likely intractable, we developed a heuristic approach that is a compromise between computational expense and alignment quality [47]. The method is comprised of two phases: first, it creates an initial PAR using stepwise addition strategy; second, it improves the initial PAR by pairwise position swapping techniques. The second

Table 4.1: The influence of *swap_med* to running time and tree cost on a dataset containing 22 taxa with approximate 20 genes

<i>swap_med</i>	tree cost	time (seconds)
0	954	28
1	948	52
2	871	77
4	882	97
8	852	102

step is repeated several iterations until either no improvement is found or the number of swap iterations exceeds a user-defined maximum number, *swap_med parameter*. The runtime complexity of the approaches is $O(n^4 \times swap_med)$ where n is the number of genes.

Table 4.1 shows that the increase of *swap_med* results in increases the runtime. However, it does not guarantee the improvement of tree cost. The *swap_med* default is one.

4.5.2 Annotated chromosome character

The annotated character type is a more general presentation of chromosome data than the breakinv character. Each annotated character consists of a sequence of loci/genes separated by pipes (“ | ”). This data type allows for locus-level rearrangements as well as nucleotide transformations. Locus homologies are determined dynamically, but based on annotated regions [47]. Given that two annotated characters each have m genes, the pairwise alignment method first calculates pairwise distances among genes are calculated and then applies the algorithm to reconstruct the PAR. The runtime complexity of the algorithm is $O(m^2 \times l^2 + n^4 \times swap_med)$ where l is the average length of genes. Note that annotated character does not require characters to have the same number of genes. Although a large of number equally optimal PARs could be constructed between two annotated characters, only a user-defined maximum number of PARs, *median*, are kept during the tree search. Our experience is that the increase of *median* does not usually result in the improvement of tree cost. The default value of *median* is one.

It takes approximately 4 minutes to construct a Wagner tree of 10 taxa each contains 8 genes of length approximate 300 nucleotides. We conducted a SPR search on the constructed Wanger tree to examine the runtime and

tree improvement. SPR search takes approximately 13 minutes and reduces the tree cost about one percent.

4.5.3 Chromosome character

Unannotated chromosomal sequences are the most general presentation of chromosome data where each chromosome consists of a long sequence of nucleotides. To analyze this character type we developed an approach to construct a *comprehensive chromosome pairwise alignment* that fulfills four conditions: (1) all putative homologies among loci are determined automatically, (2) each locus is either aligned with only one putatively homologous locus or considered as a locus indel, (3) loci are allowed to rearrange (4) the total cost to transform one genome into another genome (i.e. nucleotide transformation costs, locus indel costs, and locus order rearrangement costs) is minimized. To this end, the approach consists of two phases. First, reliable homologies between two chromosomes are detected automatically. Second, conserved areas serve as anchors to divide each chromosome into a sequence of separated loci. To construct comprehensive chromosome pairwise alignments the same method used for annotated chromosomes is applied [46]. Note that only a user-defined maximum number of PARs and medians are considered during the tree search.

To find reliable homologies between two chromosomes, we apply a three-step algorithm. First, identical segments, called *seeds*, with lengths greater or equal to user-defined *seed_length* between two genomes are detected using suffix tree structure. Detected seeds whose distance is not greater than *rearranged_len* are connected to construct larger conserved areas, called *blocks*. Blocks whose lengths are greater than the user-defined significant block length threshold *sig_block_len* are considered as reliable homologies.

A discussion of the influence of *seed_len*, *sig_block_len* and *rearranged_len* follows:

The best default value of *seed_len* is t . The higher the *seed_len* value, the fewer seeds are detected, that, in turn influences the number of blocks recognized. Conversely, if the value of *seed_length* is low, an increased number of seeds and, consequently, a greater number of short blocks are detected.

The default value of *sig_block_len* is 100. If the value of *sig_block_len* is low, small-size rearrangements are allowed; whereas if the value of *sig_block_len* is high only large-size rearrangements can be detected.

The *rearranged_len* parameter sets a threshold value under which homologous blocks separated by non-homologous regions can be considered as a single block. The default for this parameter is 100. Therefore, if two inferred

Table 4.2: The influence of *seed_length* to the runtime and tree cost of 11 corona viruses 27-32kb in length

<i>seed_length</i>	tree cost	runtime (minutes)
7	109085	16
9 (default)	91099	12
11	91524	13

Table 4.3: The default and suggested values of different parameters for chromosome characters

name	default value	cost	suggested value
seed_length	9		5-15
sig_block_len	100		60-150
rearranged_len	100		50-1000
breakpoint	10		10-50
inversion	none		15-35
approx	false		large data sets
median	1		1-2
swap_med	1		1-2
locus_indel	opening 10, extension 1		opening 10, extension 1

homologous blocks are separated by less than 100 nucleotides they will be treated as a single block in calculating of rearrangement events.

Thus, the combination of parameters **seed_length**, **sig_block_len**, and **rearranged_len** significantly influence the estimation of inferred rearrangements.

To examine the running time, we collect 11 corona viruses 27-13kb in length. The program takes about 12 minutes to reconstruct a Wagner tree, and around one hour for SPR swapping.

Table 4.3 summarizes the default and suggested values of different parameters for chromosome characters.

4.6 Transformation cost regimes

In analyses at the level of nucleotides, there are three general approaches to selecting transformation cost regimes most commonly used by P0Y5 practitioners.

Equal costs This approach assigns the same cost to all substitutions and indels, and does not take into account gap extension cost. For rationale for using this cost regime see Frost et al. [15]

Homology maximization This approach, developed by De Laet [28], assigns costs 2, 3, and 1 to transformations, gap opening, and gap extension respectively.

Parameter sensitivity analysis This method, suggested by Wheeler [48], explores the effect of varying transformation costs by comparing results of analyses conducted under different cost regimes. Partition incongruence can subsequently be computed for each cladogram and the parameter set that minimizes incongruence is selected as optimal.

More specifically, it depends on relative costs of nucleotide- and locus-level transformations. Nucleotide-level transformations are specified by tcm argument, the locus-level rearrangements are specified by locus_breakpoint or inversion costs. If locus_level rearrangement costs are extremely high, the rearrangements are not going to be counted. On the other hand, if their cost is very low (equal or slightly above that of the nucleotide-level rearrangements), rearrangements can be frequent (depending on the seed_block_len and seed_length settings).

When DNA sequence data is combined with morphological data, the cost for morphological character transformations is customarily set to be the same as for substitutions.

Chapter 5

POY5 Tutorials

These tutorials are intended to provide guidance for more sophisticated applications of POY5 that involve multiple steps and a combination of different commands. Each tutorial contains a POY5 script that is followed by detailed commentaries explaining the rationale behind each step of the analysis. Although these analyses can be conducted interactively using the *Interactive Console* or running separate sequential analyses using the *Graphical User Interphace*, the most practical way to do this is to use POY5 scripts (see *POY4 Quick Start* for more information on POY5 scripts).

It is important to remember that the numerical values for most command arguments will differ substantially depending on type, complexity, and size of the data. Therefore, the values used here should not be taken as optimal parameters.

The tutorials use sample datasets that are provided with POY5 installation but can also be downloaded from the POY5 site at

<http://research.amnh.org/scicomp/projects/poy.php>

The minimal required items to run the tutorial analyses are the POY5 application and the sample datafiles. Running these analyses requires some familiarity with POY5 interface and command structure that can be found in the preceding chapters.

5.1 Combining search strategies

The following script implements a strategy for a thorough search. This is accomplished by generating a large number of independent initial trees by random addition sequence and combining different search strategies that aim

at thoroughly exploring local tree space and escape the effect of composite optima by effectively traversing the tree space. In addition, this script shows how to output the status of the search to a log file and calculate the duration of the search.

```
(* search using all data *)
read("9.fas","31.ss", aminoacids:(“41.aa”))
(* We select as root the taxon with name t1. If we wanted
the taxon Locusta_migratoria, we would write:
root:"Locusta_migratoria" *)
set(seed:1,log:"all_data_search.log",root:"t1")
report(timer:"search start")
transform(tcm:(1,2),gap_opening:1)
build(250)
swap(threshold:5.0)
select()
perturb(transform(static_approx),iterations:15,ratchet:(0.2,3))
select()
fuse(iterations:200,swap())
select()
report("all_trees",trees:(total),"constree",graphconsensus,
"diagnosis",diagnosis)
report(timer:"search end")
set(nolog)
exit()
```

- (* search using all data *) This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- `read("9.fas","31.ss", aminoacids:(“41.aa”))` This command imports all the nucleotide sequence datafiles (all files with the extension `.seq`), a morphological datafile `morph.ss` in Hennig86 format, and an aminoacid datafile `myosin.aa`.
- `set(seed:1,log:"all_data_search.log",root:"t1")` The `set` command specifies conditions prior to tree searching. The `seed` is used to ensure that the subsequent randomization procedures (such as tree building and selecting) are reproducible. Specifying the log produces a file, `all_data_search.log` that provides an additional means to monitor the process of the search. The outgroup (`taxon1`) is designated by

the `root`, so that all the reported trees have the desired polarity. By default, the analysis is performed using direct optimization.

- `report(timer:"search start")` In combination with `report(timer:"search end")`, this command reports the amount of time that the execution of commands enclosed by `timer` takes. In this case, it reports how long it takes for the entire search to finish. Using `timer` is useful for planning a complex search strategy for large datasets that can take a long time to complete: it is instructive, for example, to know how long a search would last with a single replicate (one starting tree) before starting a search with multiple replicates.
- `transform(tcm:(1,2),gap_opening:1)` This command sets the transformation cost matrix for molecular data to be used in calculating the cost of the tree. Note, that in addition to the substitution and indel costs, the `transform` specifies the cost for gap opening.
- `build(250)` This command begins tree-building step of the search that generates 250 random-addition trees. A large number of independent starting point insures that a large portion of tree space have been examined.
- `swap(threshold:5.0)` `swap` specifies that each of the 250 trees is subjected to alternating SPR and TBR branch swapping routine (the default of `POY5`). In addition to the most optimal trees, all the suboptimal trees found within 5% of the best cost are thoroughly evaluated. This step ensures that the local searches settled on the local optima.
- `select()` Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- `perturb(transform(static_approx),iterations:15,ratchet:(0.2,3))` This command subjects the resulting trees to 15 rounds of ratchet, re-weighting 20% of characters by a factor of 2. During ratcheting, the dynamic homology characters are transformed into static homology characters, so that the fraction of nucleotides (rather than of sequence fragments) is being re-weighted. This step, that begins at multiple local maxima, is intended to further traverse the tree space in search of a global optimum.
- `fuse(iterations:200,swap())` In this step, up to 200 swappings of subtrees identical in terminal composition but different in topology, are

performed between pairs of best trees recovered in the previous step. This is another strategy for further exploration of tree space. Each resulting tree is further refined by local branch swapping under the default parameters of `swap`.

- `select()` Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- `report("all_trees",trees:(total),"constree",
graphconsensus,"diagnosis",diagnosis)` This command produces a series of outputs of the results of the search. It includes a file containing best trees in parenthetical notation and their costs (`all_trees`), a graphical representation (in PDF format) of the strict consensus (`constree`), and the diagnoses for all best trees (`diagnosis`).
- `report(timer:"search end")` This command stops timing the duration of search, initiated by the command `report(timer:"search start")`.
- `set(nolog)` This command stops reporting any output to the log file, `all_data_search.log`.
- `exit()` This command ends the POY5 session.

5.2 Searching under iterative pass

The following script implements a strategy for a thorough search under iterative pass optimization. The iterative pass optimization is a very time consuming procedure that makes it impractical to conduct under this kind of optimization (save for very small datasets that can be analyzed within reasonable time). The iterative pass, however, can be used for the most advanced stages of the analysis for the final refinement, when a potential global optimum has been reached through searches under other kinds of optimization (such as direct optimization). Therefore, this tutorial begins with importing an existing tree (rather than performing tree building from scratch) and subjecting it to local branch swapping under iterative pass.

```
(* search using all data under ip *)
read("9.fas","31.ss",aminoacids:(41.aa))
read("inter_tree.tre")
```

```
transform(tcm:(1,2),gap_opening:1)
set(iterative:approximate:2)
swap()
select()
report("all_trees",trees:(total),"constree",graphconsensus,
"diagnosis",diagnosis)
exit()
```

- (* search using all data under ip *) This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- **read("9.fas","31.ss",aminoacids:("41.aa"))** This command imports all the nucleotide sequence datafiles (all files with the extension **.seq**), a morphological datafile **morph.ss** in Hennig86 format, and an aminoacid datafile **myosin.aa**.
- **read("inter_tree.tre")** This command imports a tree file, **inter-tree.tre**, that contains the most optimal tree from prior analyses.
- **transform(tcm:(1,2),gap_opening:1)** This command sets the transformation cost matrix for molecular data to be used in calculating the cost of the tree. Note, that in addition to the substitution and indel costs, the **transform** specifies the cost for gap opening.
- **set(iterative:approximate:2)** This command sets the optimization procedure to iterative pass such that approximated three dimensional alignments generated using pairwise alignments will be considered. The program will iterate either two times, or until no further tress cost improvements can be made.
- **swap(around)** This commands specifies that the the imported tree is subjected to alternating SPR and TBR branch swapping routine (the default of P0Y5) following the trajectory of search that completely evaluates the neighborhood of the tree (by using **around**).
- **select()** Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- **report("all_trees",trees:(total),"constree", graphconsensus,"diagnosis",diagnosis)** This command produces

a series of outputs of the results of the search. It includes a file containing best trees in parenthetical notation and their costs (`all-trees`), a graphical representation (in PostScript format) of the strict consensus (`constree`), and the diagnoses for all best trees (`diagnosis`).

- `exit()` This command ends the POY5 session.

5.3 Bremer support

This tutorial builds on the previous tutorials to illustrate Bremer support calculation on trees constructed using dynamic homology characters

```
(* Bremer support part 1: generating trees *)
read("18s.fas",28s.fas")
set(root:"Americhernus")
build(200)
swap(all,visited:"tmp.trees", timeout:3600)
select()
report("my.tree",trees)
exit()

(* Bremer support part 2: Bremer calculations *)
read("18s.fas",28s.fas","my.tree")
report("support_tree.pdf",graphsupports:bremer:"tmp.trees")
exit()
```

- (* Bremer support part1: generating trees *) This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- `read("18s.fas","28s.fas")` This command imports the nucleotide sequence files `18s.fas`, `28s.fas`.
- `set(root:"Americhernus")` The `set` command specifies conditions prior to tree searching. The outgroup (`Americhernus`) is designated by the `root`, so that all the reported trees have the desired polarity.
- `build(200)` This command initializes tree-building and generates 200 random-addition trees.
- `swap(all,visited:"tmp.trees", timeout:3600)` The `swap` command specifies that each of the trees be subjected to an alternating SPR

and TBR branch swapping routine (the default of POY5). The `all` argument turns off all swap heuristics. The `visited:"tmp.trees"` argument stores every visited tree in the file specified. Although the visited tree file is compressed to accommodate the large number of trees it will accumulate, the argument `timeout` can be used to limit the number of seconds allowed for swapping also limiting the size of the file. Alternately the `swap` command can be performed as a separate analysis and terminated at the users discretion to maximize the number of trees generated.

- `select()` Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- `report("my.tree",trees)` This command will save the swapped tree, `my.tree` to a file.
- `exit()` This commands ends the POY5 session.
- `(* Bremer support part 2: Bremer calculations *)` A comment indicating the intent of the commands which follow.
- `read("18s.fas","28s.fas","my.tree")` This command imports the nucleotide sequence files `18s.fas`, `28s.fas` and the tree file, `my.tree` for which the support values will be generated. It is important to only read the selected "`my.tree`" file rather than the expansive "`tmp.trees`" file which will be used in bremer calculations.
- `report("support_tree.pdf",graphsupports:bremer:"tmp.trees")` The `report` command in combination with a file name and the `graphsupports` generates a pdf file designated by the name `support-tree.pdf` with bremer values for the selected trees held in `tmp.trees`. It is strongly recommended that this more exhaustive approach is used for calculating Bremer supports rather than simply using the `graphsupports` default s.
- `exit()` This commands ends the POY5 session.

5.4 Jackknife support

This tutorial illustrates calculating Jackknife support values for trees constructed with static homology characters. Although it is possible to calculate

both Jackknife and Bootstrap support values for trees constructed using dynamic homology characters, it is not recommended because resampling of dynamic characters occurs at the fragment (rather than nucleotide) level. Alternately dynamic homology characters can be converted to static characters using the transform argument `static_approx`. Character transformation, however, may result in a discrepancy between tree costs generated using dynamic homology characters and those generated using static homology characters. Therefore after such a data transformation is performed, an extra round of swapping is recommended to insure that the local minimum tree length is reached for the static homology characters prior to calculating support values support values.

```
(* Jackknife support for static homology trees *)
read(prealigned:("28s.aln",tcm:(1,2)))
set( root:"Americhernus")
build()
swap()
select()
calculate_support(jackknife:(remove:0.50,resample:1000),
build(5),swap(tbr,trees:3))
report("jacktree",graphsupports)
exit()
```

- (* Jackknife support for static homology trees *) This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- `read(prealigned:("28s.aln",tcm:(1,2)))` This command imports the prealigned nucleotide sequence file `28s.aln`, and treats the characters as static with the prescribed transformation cost matrix.
- `set(root:"Americhernus")` The `set` command specifies conditions prior to tree searching. The outgroup (`Americhernus`) is designated by the `root`, so that all the reported trees have the desired polarity.
- `build()` This commands begins the tree-building step of the search that generates by default 10 random-addition trees. It is essential that trees are either specified from a file or that trees are built and loaded in memory before attempting to calculate support values.

- **swap()** The `swap` command specifies that each of the trees be subjected to an alternating SPR and TBR branch swapping routine (the default of `POY5`).
- **select()** Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- **calculate_support(jackknife,(remove:0.50,resample:1000))** The `calculate_support` command generates support values as specified by the `jackknife` argument for each tree held in memory. During each pseudoreplicate half of the characters will be deleted as specified in the argument `remove:0.50`.
- **report("jacktree",graphsupports:jackknife)** The `report` command in combination with a file name and the `graphsupports` generates a pdf file with jackknife values designated by the name specified (*i.e.* `jacktree`).
- **exit()** This command ends the `POY5` session.

5.5 Sensitivity analysis

This tutorial demonstrates how data for parameter sensitivity analysis is generated. Sensitivity analysis [48] is a method of exploring the effect of relative costs of substitutions (transitions and transversions) and indels (insertions and deletions), either with or without taking gap extension cost into account. The approach consists of multiple iterations of the same search strategy under different parameters, (*i.e.* combinations of substitution and indel costs. Obviously, such analysis might become time consuming and certain methods are shown here how to achieve the results in reasonable time. This tutorial also shows the utility of the command `store` and how transformation cost matrixes are imported and used.

`POY5` does not comprehensively display the results of the sensitivity analysis or implements the methods to select a parameter set that produces the optimal cladogram, but the output of a `POY5` analysis (such as the one presented here) generates all the necessary data for these additional steps.

For the sake of simplicity, this script contains commands for generating the data under just two parameter sets. Using a larger number of parameter sets can easily be achieved by replicating the repeated parts of the script and substituting the names of input cost matrixes.

```
(* sensitivity analysis *)
read("9.fas")
set(root:"t1")
store("original_data")
transform(tcm:"111.txt")
build(100)
swap(timeout:3600)
select()
report("111.tre",trees:(total) , "111con.tre",consensus,
"111con.pdf",graphconsensus)
use("original_data")
transform(tcm:"112.txt")
build(100)
swap(timeout:3600)
select()
report("112.tre",trees:(total), "112con.tre",consensus,
"112con.pdf",graphconsensus)
exit()
```

- **(* sensitivity analysis *)** This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- **read("9.fas")** This command imports all dynamic homology nucleotide data.
- **set(root:"t1")** The outgroup (**taxon1**) is designated by the **root**, so that all the reported trees have the desired polarity.
- **store("original_data")** This commands stores the current state of analysis in memory in a temporary file, **original_data**.
- **transform(tcm:"111.txt")** This command applies a transformation cost matrix from the file **111.txt** to for subsequent tree searching.
- **build(100)** This commands begins tree-building step of the search that generates 250 random-addition trees. A large number of independent starting point insures that thee large portion of tree space have been examined.
- **swap(timeout:3600)** **swap** specifies that each of the 100 trees build in the previous step is subjected to alternating SPR and TBR branch

swapping routine (the default of P0Y5). The argument `timeout` specifies that 3600 seconds are allocated for swapping and the search is going to be stopped after reaching this limit. Because sensitivity analysis consists of multiple independent searches, it can take a tremendous amount of time to complete each one of them. In this example, `timeout` is used to prevent the searches from running too long. Using `timeout` is optional and can obviously produce suboptimal results due to insufficient time allocated to searching. A reasonable timeout value can be experimentally obtained by the analysis under one cost regime and monitoring time it takes to complete the search (using the argument `timer` of the command `set`). The advantage of using `timeout` is saving time in cases where a local optimum is quickly reached and the search is trapped in its neighborhood.

- `select()` Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- `report("111.tre",trees:(total) , "111con.tre",consensus, "111con.pdf", graphconsensus)` This command produces a file containing best tree(s) in parenthetical notation and their costs (`111.tre`), a file containing the strict consensus in parenthetical notation (`111con.tre`), and a graphical representation (in PDF format) of the strict consensus (`111con.pdf`).
- `use("original_data")` This command restored the original (non-transformed) data from the temporary file `original_data` generated by `store`.
- `transform(tcm:"112.txt")` This command applies a different transformation cost matrix from the file `112.txt` to for another round of tree searching under this new cost regime.
- `build(100)` This command begins tree-building step of the search that generates 100 random-addition trees. A large number of independent starting point insures that the large portion of tree space have been examined.
- `swap(timeout:3600)` `swap` specifies that each of the 100 trees build in the previous step is subjected to alternating SPR and TBR branch swapping routine (the default of P0Y5) to be interrupted after 3600

seconds (see the description in the previous iteration of the command above).

- **select()** Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- **report("112.tre",trees:(total),"112con.tre",consensus,"112con.pdf", graphconsensus)** This command produces a set of the same kinds of outputs as generated during the first search (see above) but under a new cost regime.
- **exit()** This command ends the POY5 session.

5.6 Chromosome analysis: unannotated sequences

This tutorial illustrates the analysis of chromosome-level transformations using unannotated sequences, i.e., contiguous strings of sequences without prior identification of independent regions. Prior to attempting an analysis of unannotated chromosomes it is necessary to enable the "long sequences" option when compiling the POY4 program.

```
(* Chromosome analysis of unannotated sequences *)
read(chromosome:("mit5.txt"))
transform((all,dynamic_pam:(locus_breakpoint:20,locus_indel:
(10,1.5),circular:true,median:2,seed_length:15,rearranged_len:
45,sig_block_len:50,median:2,swap_med:1)))
build()
swap()
select()
report("chrom",diagnosis)
report("consensustree",graphconsensus)
exit()
```

- **(* Chromosome analysis of unannotated sequences *)** This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- **read(chromosome:("mit5.txt"))** This command imports the unannotated chromosomal sequence file `mit5.txt`. The argument `chromosome` specifies the characters as unannotated chromosomes.

- `transform((all,dynamic_pam:(locus_breakpoint:20,locus_indel:(10,1.5),circular:true,seed_length:15,rearranged_len:50,sig_block_len:50,median:2,swap_med:1)))` The `transform` followed by the argument `dynamic_pam` specifies the conditions to be applied when calculating chromosome-level HTUs (medians). The argument `locus_breakpoint:20` applies a breakpoint distance between chromosome loci with the integer value determining the rearrangement cost. The argument `locus_indel:10,1.5` specifies the indel costs for the chromosomal segments, whereby the integer 10 sets the gap opening cost and the float 1.5 sets the gap extension cost. As the type of chromosomal sequences being analyzed are of mitochondrial origin, the argument `circular:true` treats each chromosome sequence as a continuous rather than linear. The argument `seed_length:15` sets the minimum length of identical continuous fragments (seeds) at 15. As seeds are the foundation for larger homologous blocks setting the seed length to an integer appropriate for the data is critical to optimizing the efficiency with which the program correctly identifies chromosomal fragments and detects rearrangements. The `rearranged_len` argument sets the lower limit for number of nucleotides between two seeds such that each is considered independent of the other. Independent seeds belong to separate homologous blocks such that rearrangement events between blocks can be detected. The argument `sig_block_len` provides the integer value determining the minimum number of nucleotides constituting an homologous block. In this example, because the data are mitochondrial containing relatively short homologous tRNA sequences, both the `rearranged_len` and the `sig_block_len` were set to values below the defaults for these arguments. The `median` specifies the number of best cost locus-rearrangements which will be considered for each HTU (median), while the `swap_med` argument specifies the maximum number of swapping iterations performed in searching for the best pairwise alignment between two chromosomes. Because values for the `median` and `swap_med` arguments set above the default (1) will significantly increase the calculation time, the default values are recommended for larger chromosomal data sets.
- `build()` This command begins the tree-building step of the search that generates by default 10 random-addition trees. It is highly recommended that a greater number of Wagner builds be implemented when analyzing data for purposes other than this demonstration.

- **swap()** The **swap** command specifies that each of the trees be subjected to an alternating SPR and TBR branch swapping routine (the default of POY5).
- **select()** Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- **report("chrom",diagnosis)** The **report** command in combination with a file name and the **diagnosis** outputs the optimal median states and edge values to a specified file (**chrom**).
- **report("consensustree",graphconsensus)** The **report** command in combination with a file name and the **graphconsensus** generates a pdf strict consensus file of the trees generated (**consensustree**).
- **exit()** This commands ends the POY5 session.

5.7 Chromosome analysis: annotated sequences

This tutorial illustrates the analysis of chromosome-level transformations using annotated sequences, i.e., contiguous strings of sequences with prior identification of independent regions delineated by pipes "|".

```
(* Chromosome analysis of annotated sequences *)
read(annotated:("aninv2"))
transform((all,dynamic_pam:(inversion:20,locus_indel:(10,
1.5),circular:false,median:1,swap_med:1)))
build()
swap()
select()
report("Annotated",diagnosis)
report("consensustree",graphconsensus)
exit()
```

- (* Chromosome analysis of annotated sequences *) This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- **read(annotated:("aninv2"))** This command imports the annotated chromosomal sequence file **aninv2**. The argument **annotated** specifies the characters.

- `transform((all,dynamic_pam:(inversion:20,locus_indel:(10,1.5),median:1,swap_med:1)))` The `transform` followed by the argument `dynamic_pam` specifies the conditions to be applied when calculating chromosome-level HTUs (medians). The argument `inversion:20` applies an inversion distance between chromosome loci with the integer value determining the rearrangement cost. The argument `locus_indel:10,1.5` specifies the indel costs for the chromosomal segments, whereby the integer 10 sets the gap opening cost and the float 1.5 sets the gap extension cost. The default values are applied to the arguments `circular median` and `swap_med` arguments to minimize the time require for these nested search options. To more exhaustively perform these calculations trees generated from initial builds can be imported to the program and reevaluated with values greater than 1 designated for the `median` and `swap_med` arguments.
- `build()` This commands begins the tree-building step of the search that generates by default 10 random-addition trees. It is highly recommended that a greater number of Wagner builds be implemented when analyzing data for purposes other than this demonstration.
- `swap()` The `swap` command specifies that each of the trees be subjected to an alternating SPR and TBR branch swapping routine (the default of POY5).
- `select()` Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- `report("Annotated",diagnosis)` The `report` command in combination with a file name and the `diagnosis` outputs the optimal median states and edge values to a specified file (`Annotated`).
- `exit()` This commands ends the POY5 session.

5.8 Custom alphabet and break inversion characters

This tutorial illustrates the analysis of the break inversion character type. Break inversion characters are generated by transforming user-defined `custom_alphabet` characters. For example, observations of developmental stages could be represented in a corresponding array such that for each terminal

taxon there is a sequence of observed developmental stages which are represented by a user-defined alphabet. To allow rearrangement as well as indel events to be considered among alphabet elements, requires transforming the `custom_alphabet` sequences to `breakinv` characters.

```
(* Custom Alphabet to Breakinv characters *)
read(custom_alphabet:("ca1.fas","m1.fas"))
transform((all,seq_to_breakinv:()))
transform((all,dynamic_pam:(locus_breakpoint:20,
locus_indel:(10,1.5),median:1,swap_med:1)))
build()
swap()
select()
report("breakinv",diagnosis)
report("consensustree",graphconsensus)
exit()
```

- `(* Custom Alphabet to Breakinv characters *)` This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- `read(custom_alphabet:("ca1.fas","m1.fas"))` This command imports the user-defined `custom_alphabet` character file `ca1.fas` and the accompanying transformation matrix `m1.fas`.
- `transform((all,seq_to_breakinv:()))` This command transforms `custom_alphabet` characters to `breakinv` characters which allow for rearrangement operations.
- `transform((all,dynamic_pam:(locus_breakpoint:20,locus_
indel:(10,1.5),median:1,swap_med:1)))` The `transform` followed by the argument `dynamic_pam` specifies the conditions to be applied when calculating medians. The argument `locus_breakpoint:20` applies a breakpoint distance calculation where the integer value specifies the rearrangement cost of `breakinv` elements. The argument `locus_-indel:10,1.5` specifies the indel costs for each `breakinv` element, whereby the integer 10 sets the gap opening cost and the float 1.5 sets the gap extension cost. The default values are applied to the `median` and `swap_med` arguments to minimize the time require for these nested search options. To more exhaustively perform these calculations trees generated from initial builds can be imported to the program and

reevaluated with values greater than 1 designated for the `median` and `swap_med` arguments

- `build()` This command begins the tree-building step of the search that generates by default 10 random-addition trees. It is highly recommended that a greater number of Wagner builds be implemented when analyzing data for purposes other than this demonstration.
- `swap()` The `swap` command specifies that each of the trees be subjected to an alternating SPR and TBR branch swapping routine (the default of P0Y5).
- `select()` Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.
- `report ("breakinv",diagnosis)` The `report` command in combination with a file name and the `diagnosis` outputs the optimal median states and edge values to a specified file (`breakinv`).
- `exit()` This command ends the P0Y5 session.

5.9 Genome analysis: multiple chromosomes

This tutorial illustrates the analysis of genome-level transformations using data from multiple chromosomes. Prior to attempting an analysis of unannotated chromosomes it is necessary to enable the "long sequences" option when compiling the P0Y4 program.

```
(* Genome analysis of multiple chromosomes *)
read (genome:"gen5bp")
transform((all,dynamic_pam:(chrom_breakpoint:80,chrom_indel:
(15,2.5),inversion:20,locus_indel:(10,1.5),median:1,
swap_med:1)))
build()
swap()
select()
report("genome",diagnosis)
report("genconsensus",graphconsensus)
exit()
```

- (* Genome analysis of multiple chromosomes*) This first line of the script is a comment. While comments are optional and do not affect the analyses, they provide are useful for housekeeping purposes.
- `read(genome:("gen5bp"))` This command imports the genomic sequence file `mit5.txt`. The argument `genome` specifies the characters as data consisting of multiple chromomsomes.
- `transform((all,dynamic_pam:(chrom_breakpoint:80,chrom_indel:(15, 2.5),locus_breakpoint:20,locus_indel:(10,1.5),median:1,swap_med:1)))` The `transform` followed by the argument `dynamic_pam` specifies the conditions to be applied when calculating genome-level HTUs (medians). The argument `chrom_breakpoint:80` applies a breakpoint distance between chromosomes with the integer value determining the rearrangement cost. The argument `chrom_-indel:15,1.5` specifies the indel costs for each entire chromosome, whereby the integer sets the gap opening cost and the float sets the gap extension cost. The argument `inversion:20` applies an inversion distance between chromosome loci with the integer value determining the rearrangement cost. The argument `locus_indel:10,1.5` specifies the indel costs for the chromosomal segments, whereby the integer 10 sets the gap opening cost and the float 1.5 sets the gap extension cost. The default values are applied to the `median` and `swap_med` arguments to minimize the time require for these nested search options. To more exhaustively perform these calculations trees generated from initial builds can be imported to the program and reevaluated with values greater than 1 designated for the `median` and `swap_med` arguments
- `build()` This commands begins the tree-building step of the search that generates by default 10 random-addition trees. It is highly recommended that a greater number of Wagner builds be implemented when analyzing data for purposes other than this demonstration.
- `swap()` The `swap` command specifies that each of the trees be subjected to an alternating SPR and TBR branch swapping routine (the default of POY5).
- `select()` Upon completion of branch swapping, this command retains only optimal and topologically unique trees; all other trees are discarded from memory.

- **report("genome",diagnosis)** The `report` command in combination with a file name and the `diagnosis` outputs the optimal median states and edge values to a specified file (`genome`).
- **report("genconsens",graphconsensus)** The `report` command in combination with a file name and the `graphconsensus` generates a pdf strict consensus file of the trees generated (`genconsens`).
- **exit()** This commands ends the POY5 session.

Bibliography

- [1] D. Barry and J. Hartigan. Statistical analysis of hominid molecular evolution. *Stat. Sci.*, 2:191–210, 1987.
- [2] M. Blanchette, G. Bourque, and D. Sankoff. *Genome Informatics*, chapter Breakpoint phylogenies, pages 25–34. Universal Academy Press, Tokyo, 1997. S. Miyano and T. Takagi—eds.
- [3] K. Bremer. The limits of amino acid sequence data in angiosperm phylogenetic reconstruction. *Evolution*, 42:795–803, 1988.
- [4] T. H. Cormen, C. E. Leiserson, R. L. Rivest, and C. Stein. *Introduction to Algorithms*. The MIT Press, 2nd edition, 2001.
- [5] Bader D. A, B. M. E. Moret, T. Warnow, S. K. Wyman, M. Yan, J. Tang, A. C. Siepel, and A. Caprara. Grappa, version 2.0. <http://www.cs.unm.edu/~moret/grappa>. Technical report, University of New Mexico, 2002.
- [6] A. C.E. Darling, B. Mau, F. R. Blattner, and N. T. Perna. Mauve: Multiple alignment of conserved genomic sequence with rearrangements. *Genome Research*, 14:1394–1403, 2004.
- [7] J. S. Farris. A method for computing Wagner trees. *Systematic Zoology*, 19:83–92, 1970.
- [8] J. S. Farris. Hennig86, 1988.
- [9] J. S. Farris. The retention index and the rescaled consistency index. *Cladistics*, 5:417–419, 1989.
- [10] J. S. Farris, V. A. Albert, M. Källersjö, Lipscomb, and A. G. Kluge. Parsimony jackknifing outperforms neighbor-joining. *Cladistics*, 12(2):99–124, 1996.

- [11] J. Felsenstein. *PHYLIP*, 1980.
- [12] J. Felsenstein. Evolutionary trees from dna sequences: a maximum likelihood approach. *J. Mol. Evol.*, 17:368–376, 1981.
- [13] J. Felsenstein. Confidence limits on phylogenies: An approach using the bootstrap. *Evolution*, 39(4):783–791, 1985.
- [14] J. Felsenstein. *Inferring Phylogenies*. Sinauer, Sunderland, MA, 2004.
- [15] D. R. Frost, M. T. Rodrigues, T. Grant, and T. A. Titus. Phylogenetics of lizard genus *Tropidurus* (squamata: Tropiduridae: Tropidurinae): direct optimization, descriptive efficiency, and sensitivity analysis of congruence between molecular data and morphology. *Molecular Phylogenetics and Evolution*, 21(3):352–371, 2001.
- [16] P. Goloboff. Analyzing large data sets in reasonable times: solutions for composite optima. *Cladistics*, 15(4):415–428, 1999.
- [17] X. Gu, Y. X. Fu, and W. H. Li. Maximum likelihood estimation of the heterogeneity of substitution rate among nucleotide sites. *Mol. Biol. Evol.*, 12:546–557, 1995.
- [18] S. Hanenhalli and P. A. Pevzner. Transforming a cabbage into a turnip (polynomial algorithm for sorting signed permutations by reversals). In *Proceedings of the 27th Annual ACM-SIAM Symposium on the Theory of Computing*, pages 178–189, 1995.
- [19] M. Hasegawa, H. Kishino, and T. Yano. A new molecular clock of mitochondrial dna and the evolution of hominoids. *Proc. Japan Acad.*, 60:95–98, 1984.
- [20] M. D. Hendy and D. Penny. Branch and bound algorithms to determine minimal evolutionary trees. *Mathematical Biosciences*, 60:133–142, 1982.
- [21] J.P. Huelsenbeck, B. Larget, and M.E. Alfaro. Bayesian phylogenetic model selection using reversible jump markov chain monte carlo. *Molecular biology and evolution*, 21(6):1123–1133, 2004.
- [22] V. Jayaswal, L.S. Jermiin, and J. Robinson. Estimation of phylogeny using a general markov model. *Evolutionary Bioinformatics Online*, 1:62, 2005.

- [23] T. H. Jukes and C. R. Cantor. Evolution of protein molecules. In N. H. Munro, editor, *Mammalian Protein Metabolism*, pages 21–132. Academic Press, New York, 1969.
- [24] M. Källersjö, J. S. Farris, A. G. Kluge, and C. Bult. Skewness and permutation. *Cladistics*, 8:275–287, 1992.
- [25] M. Kimura. A simple method for estimating evolutionary rate of base substitution through comparative studies of nucleotide sequences. *J. Mol. Evol.*, 16:111–120, 1980.
- [26] S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by simulated annealing. *Science*, 220(4598):671–680, May 1983.
- [27] A. G. Kluge and J. S. Farris. Quantitative phyletics and the evolution of anurans. *Systematic Zoology*, 30:1–32, 1969.
- [28] J. De Laet. Parsimony and the problem of inapplicables in sequence data. In V. A. Albert, editor, *Parsimony, Phylogeny, and Genomics*, page 81116. Oxford University Press, 2005.
- [29] K. Liu, S. Raghavan, S. Nelesen, C.R. Linder, and T. Warnow. Rapid and accurate large-scale coestimation of sequence alignments and phylogenetic trees. *Science*, 324(5934):1561, 2009.
- [30] T. Margush and F. R. McMorris. Consensus n -trees. *Bulletin of Mathematical Biology*, 43:239–244, 1981.
- [31] G. McGuire, M.C. Denham, and D.J. Balding. Models of sequence evolution for dna sequences containing gaps. *Molecular Biology and Evolution*, 18(4):481–490, 2001.
- [32] E. Mossel, S. Roch, and M. Steel. Shrinkage effect in ancestral maximum likelihood. *Computational Biology and Bioinformatics, IEEE/ACM Transactions on*, 6(1):126–133, 2009.
- [33] J. Neyman. Molecular studies in evolution: a source of novel statistical problems. In S. S. Gupta and J. Yackel, editors, *Statistical Decision Theory and Related Topics*, pages 1–27, 1971.
- [34] K. C. Nixon. The parsimony ratchet, a new method for rapid parsimony analysis. *Cladistics*, 15(4):407–414, 1999.
- [35] W. R. Pearson and D. J. Lipman. Improved tools for biological sequence comparison. *PNAS*, 85:2444–2448, 1988.

- [36] B. Redelings and M. Suchard. Incorporating indel information into phylogeny estimation for rapidly emerging pathogens. *BMC evolutionary biology*, 7(1):40, 2007.
- [37] B.D. Redelings and M.A. Suchard. Joint bayesian estimation of alignment and phylogeny. *Systematic Biology*, 54(3):401–418, 2005.
- [38] E. Rivas and S.R. Eddy. Probabilistic phylogenetic inference with insertions and deletions. *PLoS computational biology*, 4(9):e1000172, 2008.
- [39] F. J. Rohlf. Consensus indices for comparing classifications. *Mathematical Biosciences*, 59:131–144, 1982.
- [40] D. Sankoff and M. Blanchette. The median problem for breakpoints in comparative genomics. *Computing and Combinatorics 3rd Annual Int. Conf. COCOON 97*, 1276:251–263, 1997.
- [41] M. Steel and D. Penny. Parsimony, likelihood, and the role of models in molecular phylogenetics. *Molecular Biology and Evolution*, 17(6):839–850, 2000.
- [42] J. Sullivan, D. L. Scofford, and G. J. P. Naylor. The effect of taxon-sampling on estimating rate heterogeneity parameters on maximum-likelihood models. *Mol. Biol. Evol.*, 52:7649–664, 1999.
- [43] D. L. Swofford and G. J. Olsen. Phylogeny reconstruction. In D. Hillis and C. Moritz, editors, *Molecular Systematics*, chapter 11, pages 411–501. Sinauer Ass. Inc., Sunderland, Massachusetts, USA, 1990.
- [44] H. Tamura and M. Nei. Estimation of the number of nucleotide substitutions in the control region of mitochondrial dna in humans and chimpanzees. *Mol. Biol. Evol.*, 10:512–526, 1993.
- [45] S. Tavaré. Some probabilistic and statistical problems on the analysis of dna sequences. *Lec. Math. Life Sci.*, 17:57–86, 1986.
- [46] L. S. Vinh, A. Varon, D. Janies, and W. C. Wheeler. Towards phylogenomic reconstruction. In *Proceedings of the International Conference on Bioinformatics and Computational Biology*, pages 98–104, Las Vegas, Nevada, USA, 2007. CSREA Press.
- [47] L. S. Vinh, A. Varon, and W. C. Wheeler. Pairwise alignment with rearrangements. *Genome informatics*, 17(2):141–151, 2006.

- [48] W. C. Wheeler. Sequence alignment, parameter sensitivity, and the phylogenetic analysis of molecular data. *Systematic Biology*, 44(3):321–331, 1995.
- [49] W. C. Wheeler. Optimization alignment: The end of multiple sequence alignment in phylogenetics? *Cladistics*, 12(1):1–9, 1996.
- [50] W. C. Wheeler. Fixed character states and the optimization of molecular sequence data. *Cladistics*, 15(4):379–385, 1999.
- [51] W. C. Wheeler. Homology and DNA sequence data. In G.P. Wagner, editor, *The Character Concept in Evolutionary Biology*, pages 303–318. Academic Press, New York, 2001.
- [52] W. C. Wheeler. Homology and the optimization of DNA sequence data. *Cladistics*, 17:S3–S11, 2001.
- [53] W. C. Wheeler. *Optimization Alignment: down, up, error, and improvements*. Techniques in Molecular Systematics and Evolution. Birkhäuser, Basel, Boston, Berlin, 2002.
- [54] W. C. Wheeler. Implied alignment. *Cladistics*, 19:261–268, 2003.
- [55] W. C. Wheeler. Iterative pass optimization. *Cladistics*, 19:254–260, 2003.
- [56] W. C. Wheeler. Search-based character optimization. *Cladistics*, 19:348–355, 2003.
- [57] W. C. Wheeler, L. Aagesen, C. P. Arango, J. Faivoich, T. Grant, C. D’Haese, D. Janies, W. L. Smith, A. Varón, and G. Giribet. *Dynamic Homology and Systematics: A Unified Approach*. American Museum of Natural History, 2006.
- [58] W. C. Wheeler, J. Gatesy, and R. DeSalle. Elision: A method for accommodating multiple molecular sequence alignments with alignment-ambiguous sites. *Molecular Phylogenetics and Evolution*, 4(1):1–9, 1995.
- [59] S. Yancopoulos, O. Attie, and R. Friedberg. Efficient sorting of genomic permutations by translocation, inversion and block interchange. *Bioinformatics*, 21:3340–3346, 2005.
- [60] Z. Yang. Maximum likelihood phylogenetic estimation from dna sequences with variable rates over sites: Approximate methods. *J. Mol. Evol.*, 39:105–111, 1994.

General Index

POY 3.0 Commands Index