

Phylogeny

Phylogeny is the description of biological relationships, usually expressed as a tree.

A statement of phylogeny among objects **assumes** homology and **depends** on classification.

Phylogenetic analysis is an investigation of the evolutionary relationships among a group of related sequences by producing a tree representation of the relationships.

In fact, phylogenetic relationships among many kinds of organisms are difficult to determine in any other way.

Simply, organisms with **high degrees of molecular similarity** are expected to be **closely related** than those that are dissimilar.

Due to the availability of molecular data, taxonomists are forced to rely on comparisons of phenotypes (how organisms looked) to infer their genotypes (the genes that gave rise to their physical appearance).

Humans, flies, mollusks: light detecting organ-eye

Protein/DNA sequences

Phylogenetic trees

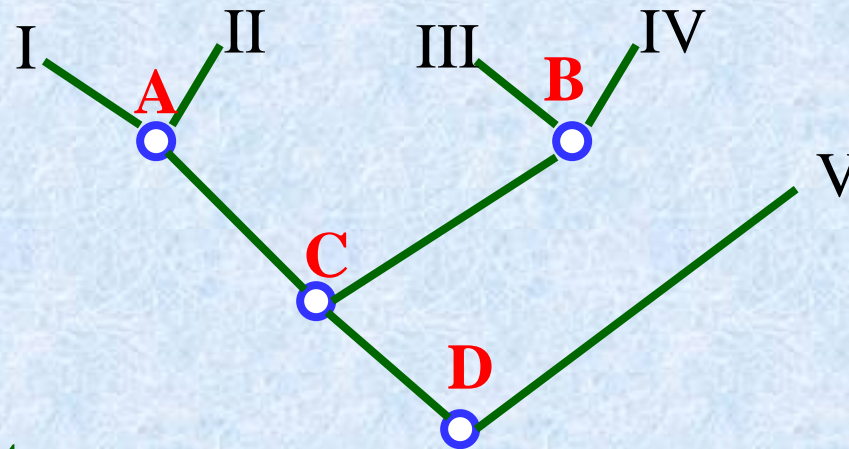
Graphical representation of the evolutionary relationship among three or more genes or organisms.

Relatedness of data; divergence times; nature of common ancestors

Nodes and branches:

Terminal nodes: at the tip of the branches (gene or organism: available data)

Internal node: common ancestor (data are not available)



Bifurcating
Multifurcating

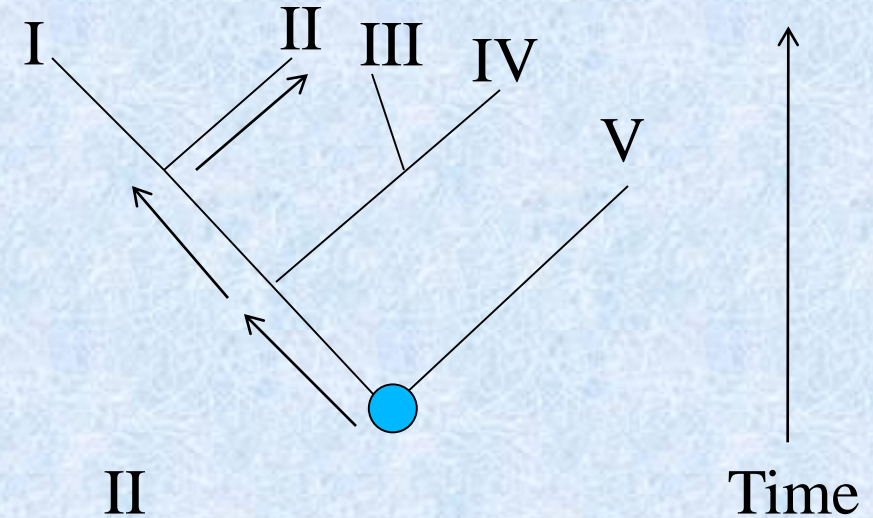
Scaled trees
Unscaled trees

Newick format

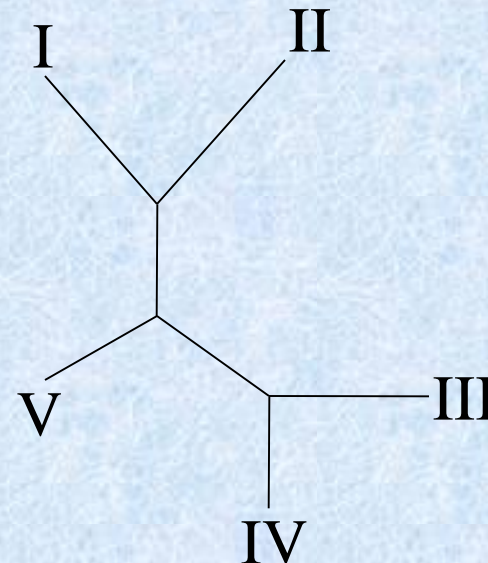
$((I, II), (III, IV)), V$

Rooted and Unrooted trees

Rooted trees: a single node is designated as a common ancestor and a unique path leads it.



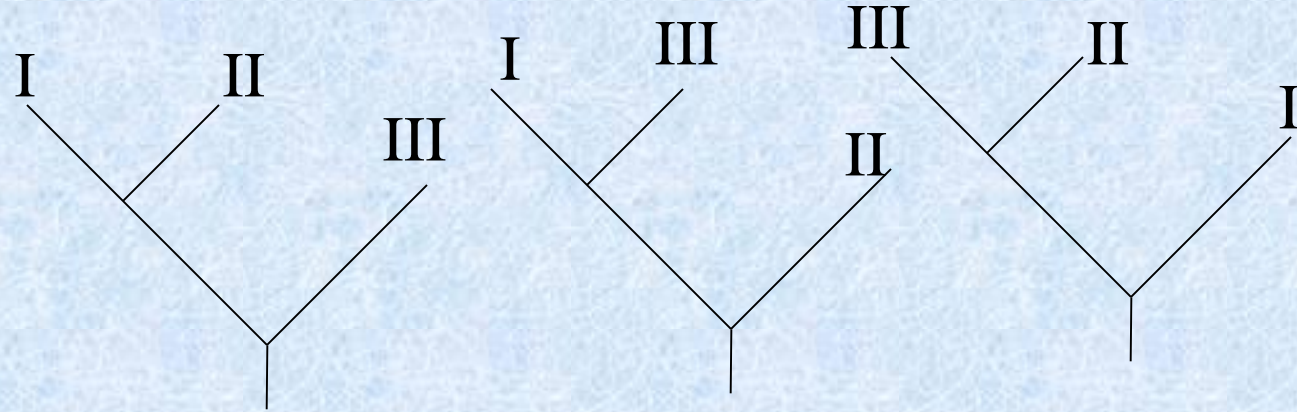
Unrooted trees: only relationship and not direction



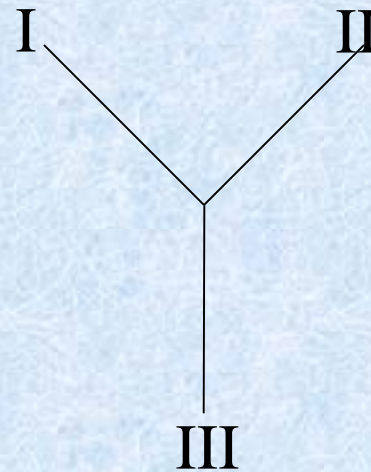
Possible rooted and unrooted trees

3 species

Rooted: 3



Unrooted: 1



$$N_R = (2n-3)!/2^{n-2}(n-2)!$$

$$N_U = (2n-5)!/2^{n-3}(n-3)!$$

2 data: Rooted: 1; Unrooted: 1

3 data: R: 3, U: 1

4 data: R: 15, U: 3

Tree construction

1. UPGMA (Unweighted Pair Group Method with Arithmetic mean)
2. Transformed Distance Method
3. Neighbor's Relation Method
4. Neighbor Joining Methods
5. Maximum Likelihood Approaches

UPGMA

Statistically based method

Requires data that can be condensed to genetic distance (distance matrix)

E.g. Species, A, B, C and D

Species	A	B	C
B	d_{AB}	-	-
C	d_{AC}	d_{BC}	-
D	d_{AD}	d_{BD}	d_{CD}

d_{AB} : the number of mismatching nucleotides (divided by total number of sites, where matches could have been found)

If A and B are related (minimum mismatches)

Form a group (AB)

Combine with other species, C and D

$$d_{(AB),C} = 1/2 (d_{AC} + d_{BC})$$

$$d_{(AB),D} = 1/2 (d_{AD} + d_{BD})$$

The species separated by the smallest distance in the new matrix can be clustered together.

For scaled branches, the distances will be averaged.

Example

A: GTGCTGCACG GCTCAGTATA GCATTTACCC TTCCATCTTC AGATCCTGAA

B: **AC**GCTGCACG GCTCAGT**GCG** **GTGCT**TACCC T**CC**ATCTTC AGATCCTGAA

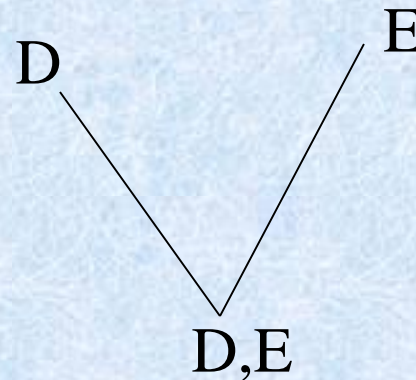
C: GTGCTGCACG GCTCGGCGCA GCATTTACCC TCCCATCTTC AGATCCTATC

D: GTATCACACG ACTCAGCGCA GCATTTGCCC TCCCGTCTTC AGATCCTAAA

E: GTATCACAT**A** **G**CTCAGCGCA GCATTTGCCC TCCCGTCTTC AGATC**TA**AAA

Pairwise comparison (distance matrix)

Species	A	B	C	D
B	9	-	-	-
C	8	11	-	-
D	12	15	10	-
E	15	18	13	5



Pairwise comparison (distance matrix)

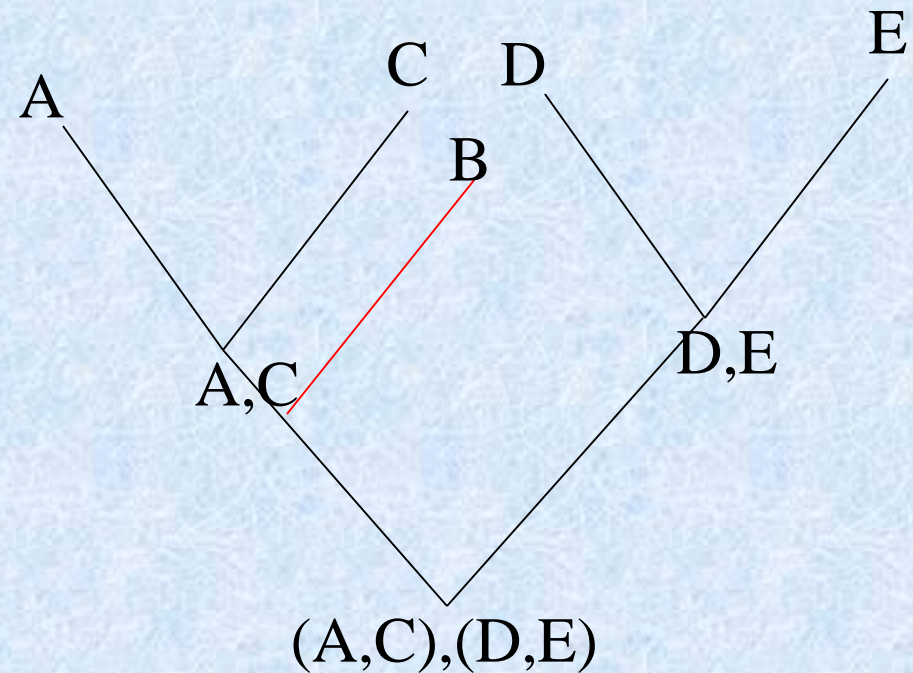
Species	A	B	C
B	9	-	-
C	8	11	-
DE	13.5	16.5	11.5

$$(AC), B = \frac{1}{2}(AB + BC)$$

Pairwise comparison

Species	B	AC
AC	10	-
DE	16.5	12.5

$$(((A,C),B),D,E)$$



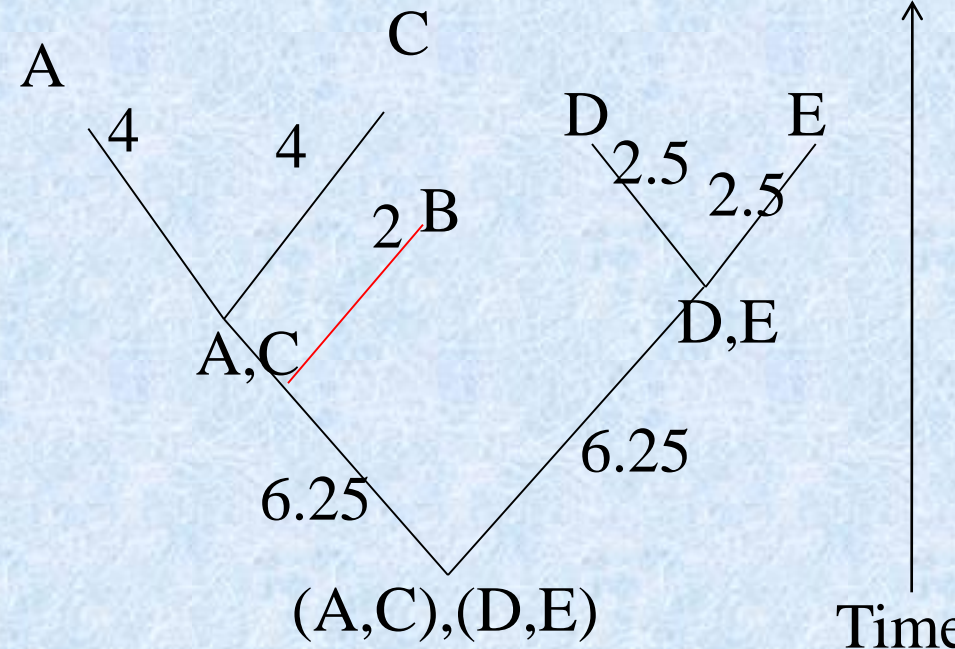
$$(((A,C),B),D,E)$$

Estimation of branch lengths

Length of the branches can also be calculated with distance matrix

Pairwise comparison (distance matrix)

Species	A	B	C	D
B	9	-	-	-
C	8	11	-	-
D	12	15	10	-
E	15	18	13	5



Neighbor's Relation Method

Another popular variant of UPGMA: tree is constructed with the smallest possible branch lengths overall.

Any unrooted tree, pairs of species that are separated from each other by just one internal node are said to be neighbors.

[UPGMA: Branch length is not additive.

E.g: $d_{AE} = 4 + 6.25 + 6.25 + 2.5 = 19$

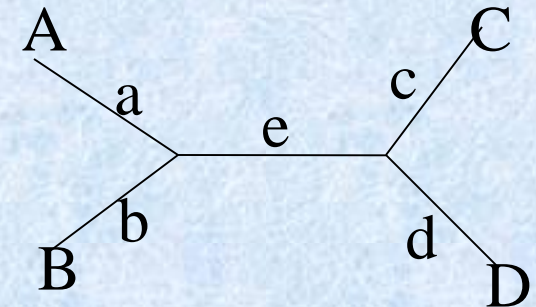
Actual case: 15]

If additivity holds good:

$$d_{AC} + d_{BD} = d_{AD} + d_{BC} = a + b + c + d + 2e = d_{AB} + d_{CD} + 2e$$

a,b,c,d: lengths of terminal branches

e: length of central branch



Four point condition:

$$d_{AB} + d_{CD} < d_{AC} + d_{BD}$$

and

$$d_{AB} + d_{CD} < d_{AD} + d_{BC}$$

Considers all possible pairwise arrangements of four species and determines the arrangement, which satisfies the four point condition.

For four species, considers all possible values,

(i) $d_{AB} + d_{CD}$; (ii) $d_{AC} + d_{BD}$ and (iii) $d_{AD} + d_{BC}$

Smallest sum with pairing is 1 and others are 0

Repeat for all possible four pairs

Ones with highest scores are grouped.

New distance matrix can be generated as was done for UPGMA

Neighbor joining method

Tree is of star-like and all species comes as a single central node regardless of the number.

Difference with other methods is the way it determines the sum of branch lengths with each reiteration of the process.

$$S_{12} = (1/(2(N-2))) \sum (d_{1k} + d_{2k}) + (1/2)d_{12} + (1/(N-2)) \sum d_{ij}$$

Any pair of species can take positions 1 and 2; k is an accepted outgroup

Simplified into $Q_{12} = (N-2)d_{12} - \sum d_{1i} - \sum d_{2i}$

All possible pairs are considered and the pairs with smallest distance is taken.

Construct new distance matrix as done with UPGMA and repeat the process.

Maximum likelihood approaches

It represent an alternative and purely statistically based method of phylogenetic reconstruction.

Probabilities are considered for every individual nucleotide substitution.

Transitions (purine to purine/ pyrimidine to pyrimidine) and transversions

Multiple substitutions occurred at one or more sites, which are not necessarily independent.

It is necessary to take into account of all these facts, which needs heavy computational power.

With current facilities, it is possible to use the method for tree construction.

Program to construct trees

- **as Windows executables** (not counting executing in a "DOS box"). Programs available as source code which is Windows-specific are listed below. (Note that compilers available on Windows systems, particularly the free Cygwin and MinGW compilers, can also be used to compile many of these generic source code). Programs run in interpreted environments such as Perl, Python, R or MATLAB can also be run under Windows if the source programs are listed above under Unix.

▪ PHYLIP	▪ DNASIS	▪ Mesquite	▪ MrModeltest	▪ MESA
▪ PAUP*	▪ MINSPNET	▪ Phyedit	▪ SymmeTREE	▪ MultiPhyl
▪ TREECON	▪ BioEdit	▪ SYN-TAX	▪ TreeJuxtaposer	▪ NimbleTree
▪ GDA	▪ ProSeq	▪ PTP	▪ Network	▪ ArboDraw
▪ SeqPup	▪ PAL	▪ DIVA	▪ Spectronet	▪ SPAGeDi
▪ MOLPHY	▪ WINCLADA	▪ TreeFitter	▪ Phylogen	▪ CBCAnalyzer
▪ GeneDoc	▪ NONA	▪ Phylo_win	▪ Phylap	▪ DualBrothers
▪ COMPONENT	▪ Phylogenetic Independence	▪ SplitsTree	▪ Dnatree	▪ PaupUp
▪ TREEMAP	▪ PEBBLE	▪ PAST	▪ IMa2	▪ Notung
▪ COMPARE	▪ HY-PHY	▪ GeneStudio Pro	▪ ProfTest	▪ SSA
▪ RAPDistance	▪ TreeExplorer	▪ Treefinder	▪ GEODIS	▪ Multidivtime
▪ TreeView	▪ Genie	▪ PPH	▪ TreeSetViz	▪ ParaFit
▪ Phylodendron	▪ Vanilla	▪ MetaPIGA	▪ TreeMe	▪ IDC
▪ POPGENE	▪ MEGA	▪ Phyltools	▪ ModelGenerator	▪ TreeMaker
▪ TFPGA	▪ TNT	▪ MSA	▪ Simplot	▪ CodonRates
▪ GeneTree	▪ GelCompar II	▪ Mgenome	▪ PHYLOGR	▪ BAli-Phy
▪ MVSP	▪ Bionumerics	▪ APE	▪ ProfDist	▪ CoMET
▪ RSTCALC	▪ TCS	▪ PHASE	▪ START2	▪ TreeDyn
▪ Genetix	▪ FORESTER	▪ PHYML	▪ IQPNNI	▪ DigTree
▪ NJplot	▪ Populations	▪ YCDMA	▪ STC	▪ Geneious
▪ unrooted	▪ T-REX	▪ NSA	▪ TreeSAAP	▪ Brownie
▪ Arlequin	▪ MrBayes	▪ BEAST	▪ Swaap	▪ Mac5
▪ DAMBE	▪ EDIBLE	▪ Clann	▪ Swaap PH	▪ BayesPhylogenies
▪ DnaSP	▪ Winboot	▪ Jevtrace	▪ TreeGraph 2	▪ BayesTraits
▪ PAML	▪ r8s	▪ MrMTgui	▪ DIVERGE	▪ MrEnt

Phylip

Phylip

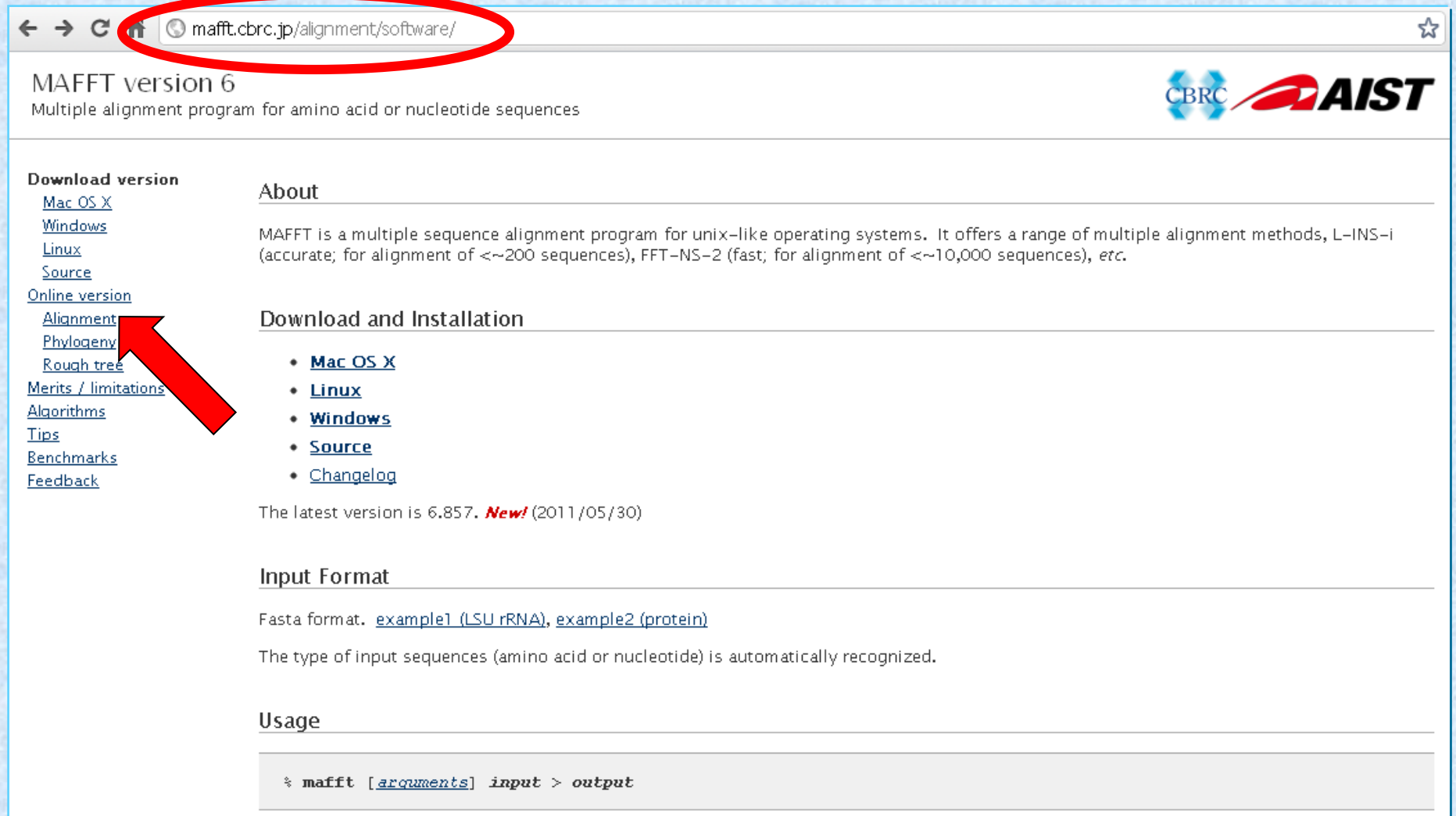
Phylip is a program to create phylogenetic tree for a given set of amino acid sequences.

It takes the multiple alignment of the sequences as input

Multiple sequence alignments can be done with ClustalW, MAFFT etc.

MAFFT is widely used to prepare the input multiple alignment file suitable for Phylip

MAFFT




The screenshot shows the MAFFT website interface. The browser's address bar at the top contains the URL mafft.cbrc.jp/alignment/software/, which is circled in red. The page title is "MAFFT version 6" with the subtitle "Multiple alignment program for amino acid or nucleotide sequences". The CBRC and AIST logos are in the top right. On the left, a sidebar lists various links: "Download version" (with sub-links for Mac OS X, Windows, Linux, and Source), "Online version" (with sub-links for Alignment, Phylogeny, and Rough tree), "Merits / limitations", "Algorithms", "Tips", "Benchmarks", and "Feedback". A large red arrow points from the "Alignment" link to the "Download and Installation" section. The "About" section describes MAFFT as a multiple sequence alignment program for Unix-like systems, offering methods like L-INS-i and FFT-NS-2. The "Download and Installation" section lists links for Mac OS X, Linux, Windows, Source, and Changelog. The "Input Format" section mentions Fasta format and provides example links. The "Usage" section shows a command-line example: `% mafft [arguments] input > output`.

← → ↻ ↗ mafft.cbrc.jp/alignment/software/ ☆

MAFFT version 6

Multiple alignment program for amino acid or nucleotide sequences



Download version

- [Mac OS X](#)
- [Windows](#)
- [Linux](#)
- [Source](#)

Online version

- [Alignment](#)
- [Phylogeny](#)
- [Rough tree](#)

- [Merits / limitations](#)
- [Algorithms](#)
- [Tips](#)
- [Benchmarks](#)
- [Feedback](#)

About

MAFFT is a multiple sequence alignment program for unix-like operating systems. It offers a range of multiple alignment methods, L-INS-i (accurate; for alignment of <~200 sequences), FFT-NS-2 (fast; for alignment of <~10,000 sequences), etc.

Download and Installation

- [Mac OS X](#)
- [Linux](#)
- [Windows](#)
- [Source](#)
- [Changelog](#)

The latest version is 6.857. **New!** (2011/05/30)

Input Format

Fasta format. [example1 \(LSU rRNA\)](#), [example2 \(protein\)](#)

The type of input sequences (amino acid or nucleotide) is automatically recognized.

Usage

```
% mafft [arguments] input > output
```


MAFFT

UPPERCASE / lowercase:

☐ Same as input

Parameters:

Scoring matrix for amino acid sequences: BLOSUM62

Scoring matrix for nucleotide sequences: 200PAM / $\kappa=2$

† Switch it to '1PAM / $\kappa=2$ ' when aligning closely related DNA sequences.

Gap opening penalty: 1.53 (1.0 – 3.0)

Offset value: 0.0 (0.0 – 1.0)

† If long gaps are not expected, set it as 0.1 or larger value.

Mafft-homologs (Collects homologs from SwissProt by BLAST and performs profile-based alignments; Protein only): [Help](#)

☐ On

☐ Show homologs (if any)

Number of homologs: 50 (5 – 200)

Threshold: $E=$ 1e-10 (1e-5 – 1e-40)

Plot **LAST** hits (DNA only):

☒ The top sequence vs the others ☐ The longest sequence vs the others

☒ Plot and alignment ☐ Plot only ☐ Alignment only

Threshold: score=39 ($E=8.4e-11$)

Submit

Reset

with <200 sequences × <1,000 nucleotides) [Help](#)

MAFFT Results

The screenshot displays the Readseq software interface. At the top, there are buttons for 'Jalview', 'Reformat', and 'Phylogenetic Tree'. The 'Reformat' button is highlighted with a red arrow. Below these buttons, the text 'to GCG, PHYLIP, MSF, NEXUS, uppercase/lowercase, etc. with Readseq' is visible. The main window shows a 'MAFFT-G-INS-i Result' with a CLUSTAL format alignment. Overlaid on this is an 'Options' dialog box. A red arrow points to the 'Output sequence format' dropdown menu, which is open and shows 'Phylip|Phylip4' selected. Other options in the dropdown include 'Pretty', 'IG|Stanford', 'GenBank|gb', 'NBRF', 'EMBL|em', 'GCG', 'DNAStrider', 'Pearson|Fast|fa', 'Phylip3.2', 'PlainRaw', 'PIR|CODATA', 'MSF', 'PAUP|NEXUS', 'Pretty', 'XML', 'Clustal', 'FlatFeat|FFF', 'GFF', and 'ACEDB'. The 'Options' dialog box also contains checkboxes for 'Remove gap symbols', 'Calculate checksum of sequences', and 'Translate bases (list as from-base to base pairs)'. The 'Select' radio button is set to 'all'. At the bottom of the dialog box, the version 'Readseq by D.G. Gilbert, 2.1.26 (18-Oct-2007)' and the URL 'http://hubio.bio.indiana.edu/soft/molbio/readseq/java/' are displayed.

10 142

```

sp|P69905| MVLSPADKTN VKAAWGKVG A HAGEYGAEAL ERMFLSFPTT KTYFPHFDLS
sp|P69907| MVLSPADKTN VKAAWGKVG A HAGEYGAEAL ERMFLSFPTT KTYFPHFDLS
sp|P06635| MVLSPADKTN VKTAWGKVG A HAGDYGAEAL ERMFLSFPTT KTYFPHFDLS
sp|P01966| MVLSAADKGN VKAAWGKVG G HAAEYGAEAL ERMFLSFPTT KTYFPHFDLS
sp|P01958| MVLSAADKTN VKAAWSKVGG HAGEYGAEAL ERMFLGFPTT KTYFPHFDLS
sp|P01959| MVLSAADKTN VKAAWSKVGG NAGEFGAEAL ERMF
sp|P01942| MVLSGEDKSN IKAAWGKIGG HGAEYGAEAL ERMF
sp|P01946| MVLSADDKTN IKNCWGKIGG HGGEYGEEAL QRMF
sp|P01965| -VLSAADKAN VKAAWGKVG G QAGAHGAEAL ERMF
sp|P60529| -VLSPADKTN IKSTWDKIGG HAGDYGGEAL DRTF

```

```

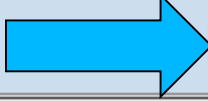
HGSAQVKGHG KKVADALTNA VAHVDDMPNA LSAL
HGSAQVKGHG KKVADALTNA VAHVDDMPNA LSAL
HGSAQVKDHG KKVADALTNA VAHVDDMPNA LSAL
HGSAQVKGHG AKVAAALTKA VEHLDDLPGA LSEL
HGSAQVKAHG KKVGDALTLA VGHLDDLPGA LSNL
HGSAQVKAHG KKVGDALTLA VGHLDDLPGA LSNL
HGSAQVKGHG KKVADALASA AGHLDDLPGA LSAL
PGSAQVKAHG KKVADALAKA ADHVEDLPGA LSTL
HGSDQVKAHG QKVADALTKA VGHLDDLPGA LSAL
PGSAQVKAHG KKVADALTTA VAHLDDLPGA LSAL

```

```

LLSHCLLVTL AAHLPAEFTP AVHASLDKFL ASVS
LLSHCLLVTL AAHLPAEFTP AVHASLDKFL ASVS
LLSHCLLVTL AAHLPAEFTP AVHASLDKFL ASVSTVLTSK YR
LLSHSLLVTL ASHLPSDFTP AVHASLDKFL ANVSTVLTSK YR
LLSHCLLSTL AVHLPNDFTP AVHASLDKFL SSVSTVLTSK YR
LLSHCLLSTL AVHLPNDFTP AVHASLDKFL STVSTVLTSK YR
LLSHCLLVTL ASHHPADFTP AVHASLDKFL ASVSTVLTSK YR
FLSHCLLVTL ACHHPGDFTP AMHASLDKFL ASVSTVLTSK YR
LLSHCLLVTL AAHHPDDFNP SVHASLDKFL ANVSTVLTSK YR
LLSHCLLVTL ACHHPTEFTP AVHASLDKFF AAVSTVLTSK YR

```



Options

Output sequence format:

Return biosequence data:
☒ Download to file
☐ View in browser

Change sequence case to
☒ No change
☐ lower
☐ UPPER

☐ Remove gap symbols:

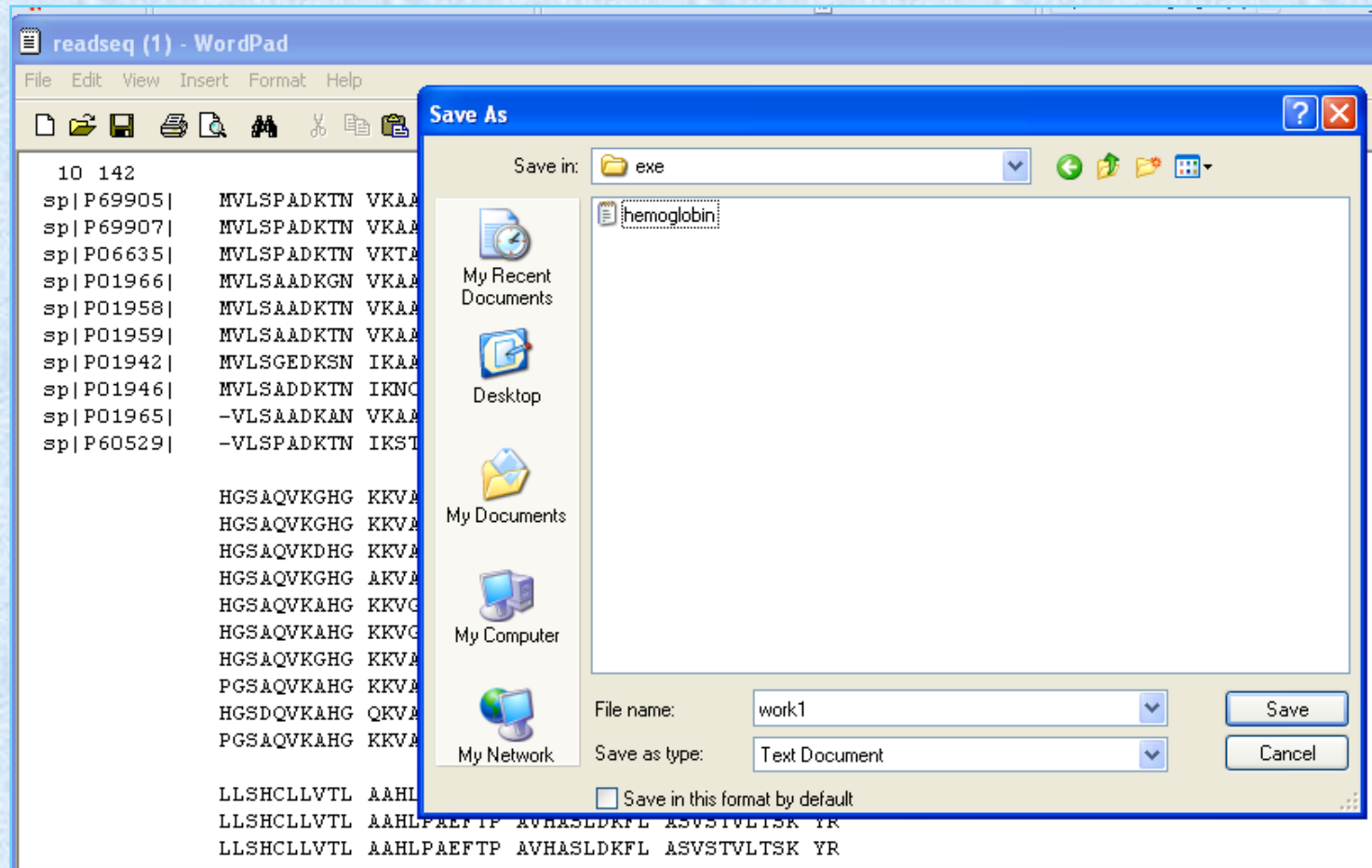
☐ Calculate checksum of sequences

Select ☒ all, or ☐ sequences by number:

☐ Translate bases (list as from-base:to-base pairs)

Saved in a temporary file “**readseq**”.

Open it and save as “**work1**”



Folder: Phylip-3.69/exe/work1

Procedure to run Phylip

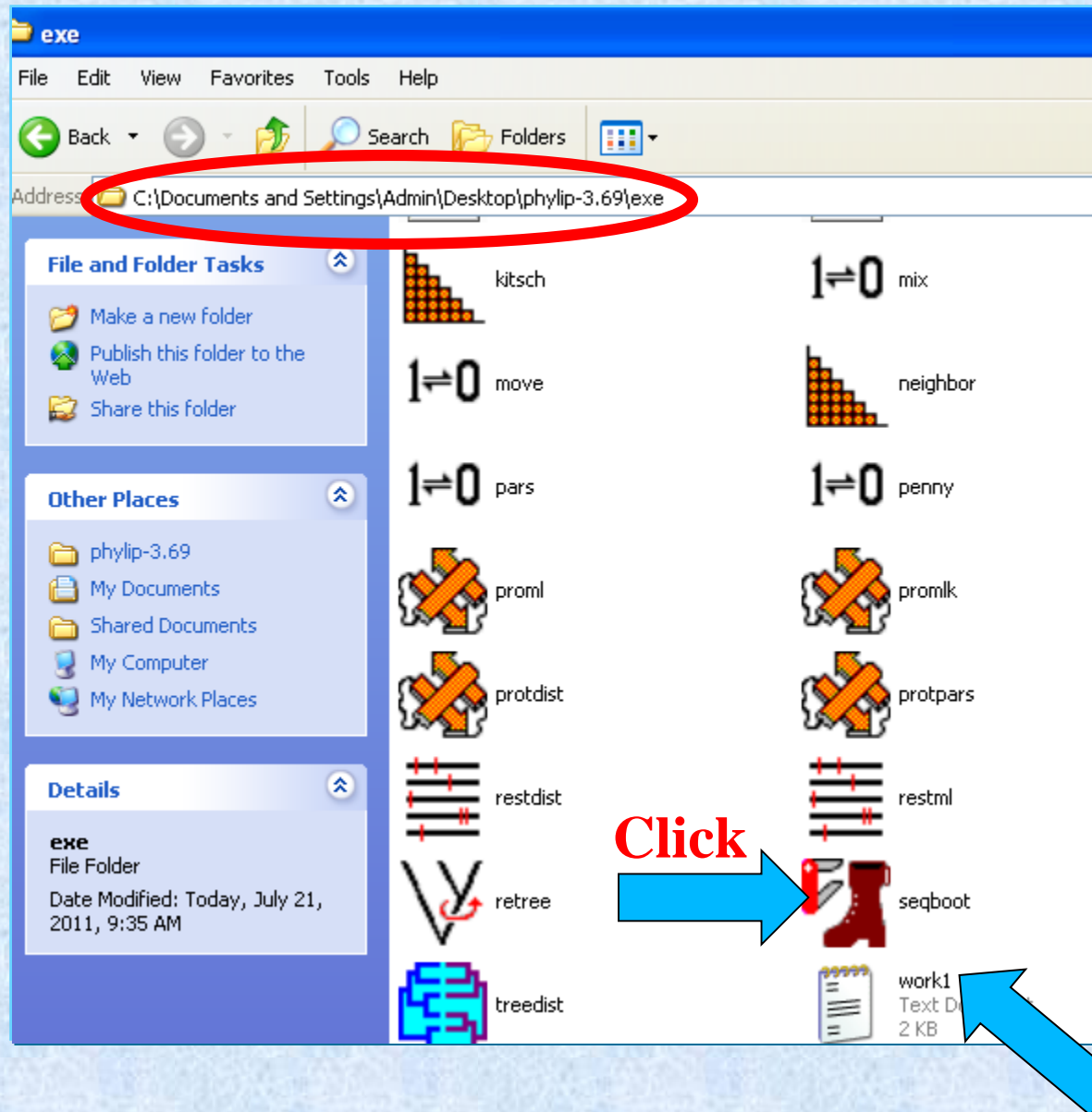
1. **Bootstrapping**: to check the confidence level

In statistics, bootstrapping is a computer-based method for **assigning measures of accuracy** to sample estimates.

Bootstrapping is the practice of estimating properties of an estimator (such as its variance) by measuring those properties when sampling from an approximating distribution.

One standard choice for an approximating distribution is the **empirical distribution** of the observed data.

This can be implemented by constructing a number of **resamples** of the observed dataset, each of which is obtained by **random sampling** with replacement from the original dataset.



seqboot.exe: can't find input file "infile"
Please enter a new file name> work1.txt

Bootstrapping algorithm, version 3.69

Settings for this run:

D Sequence, Morph, Rest., Gene Freqs? Molecular sequences
J Bootstrap, Jackknife, Permute, Rewrite? Bootstrap
% Regular or altered sampling fraction? regular
B Block size for block-bootstrapping? 1 (regular bootstrap)
R How many replicates? 100
W Read weights of characters? No
C Read categories of sites? No
S Write out data sets or just weights? Data sets
I Input sequences interleaved? Yes
0 Terminal type (IBM PC, ANSI, none)? IBM PC
1 Print out the data at start of run No
2 Print indications of progress of run Yes

Y to accept these or type the letter for one to change

completed replicate number 1
completed replicate number 2
completed replicate number 3
completed replicate number 4
completed replicate number 5
completed replicate number 6
completed replicate number 7
completed replicate number 8
completed replicate number 9
completed replicate number 10

Output written to file "outfile"

Done.

Press enter to quit.

Bootstrapping algorithm, version 3.69

Settings for this run:

D Sequence, Morph, Rest., Gene Freqs? Molecular sequences
J Bootstrap, Jackknife, Permute, Rewrite? Bootstrap
% Regular or altered sampling fraction? regular
B Block size for block-bootstrapping? 1 (regular bootstrap)
R How many replicates? 10
W Read weights of characters? No
C Read categories of sites? No
S Write out data sets or just weights? Data sets
I Input sequences interleaved? Yes
0 Terminal type (IBM PC, ANSI, none)? IBM PC
1 Print out the data at start of run No
2 Print indications of progress of run Yes

Y to accept these or type the letter for one to change

R
Number of replicates?
10

outfile contains
10 replications

Y to accept these or type the letter for one to change

y

Random number seed (must be odd)?

5

	10	142									
sp P69905	MLPPADKKT	TVVAGKGAH	HGGEAELRM	MMFFSTTK	KTYFFFDD	SHHGAAKG					
sp P69907	MLPPADKKT	TVVAGKGAH	HGGEAELRM	MMFFSTTK	KTYFFFDD	SHHGAAKG					
sp P06635	MLPPADKKT	TVVTGKGAH	HGGDGAELRM	MMFFSTTK	KTYFFFDD	SHHGAAKG					
sp P01966	MLAADKKG	GVVAGKGGH	HAAEGAELRM	MMFFSTTK	KTYFFFDD	SHHGAAKA					
sp P01958	MLAADKKT	TVVASKGGH	HGGEAELRM	MMFFGTTK	KTYFFFDD	SHHGAAKG					
sp P01959	MLAADKKT	TVVASKGGN	NGGEAELRM	MMFFGTTK	KTYFFFDD	SHHGAAKG					
sp P01942	MLGGEDKKS	SIIAGKGGH	HAAEGAELRM	MMFFSTTK	KTYFFFDD	SHHGAAKG					
sp P01946	MLAADKKT	TIINGKGGH	HGGEELRM	MMFFATTK	KTYIIIDD	SPPGAAKG					
sp P01965	-LAAADKKA	AVVAGKGGQ	QGGAGAELRM	MMFFGTTK	KTYFFFNNS	SHHGDDKGQ					
sp P60529	-LPPADKKT	TIISDKGGH	HGGDGGELRT	TTTFSTTK	KTYFFFDD	SPPGAAKG					

KKKVADDDL	AAADPSALL	HAHRRDDNN	KLLHCCCCV	TTTLAAAHL	LAAAEFFTHA
KKKVADDDL	AAADPSALL	HAHRRDDNN	KLLHCCCCV	TTTLAAAHL	LAAAEFFTHA
KKKVADDDL	AAADPSALL	HAHRRDDNN	KLLHCCCCV	TTTLAAAHL	LAAAEFFTHA
AKKVAAAAA	AAEDPSELL	HAHRRDDNN	KLLHSSSSV	TTTLAASHL	LSSSDFFTHA
KKKVGDDDL	AAGDPSNLL	HAHRRDDNN	KLLHCCCCS	TTTLAAVHL	LNNNDFFTHA
KKKVGDDDL	AAGDPSNLL	HAHRRDDNN	KLLHCCCCS	TTTLAAVHL	LNNNDFFTHA
KKKVADDDL	AAGDPSALL	HAHRRDDNN	KLLHCCCCV	TTTLAASHH	HAAADFFTHA
KKKVADDDL	AADEPSTLL	HAHRRDDNN	KFFFHCCCCV	TTTLAACHH	HGGGDDFFTHA
QKKVADDDL	AAGDPSALL	HAHRRDDNN	KLLHCCCCV	TTTLAAAHH	HDDDDFFNHA
KKKVADDDL	AAADPSALL	HAYYRRDDNN	KLLHCCCCV	TTTLAACHH	HTTTEFFTHA

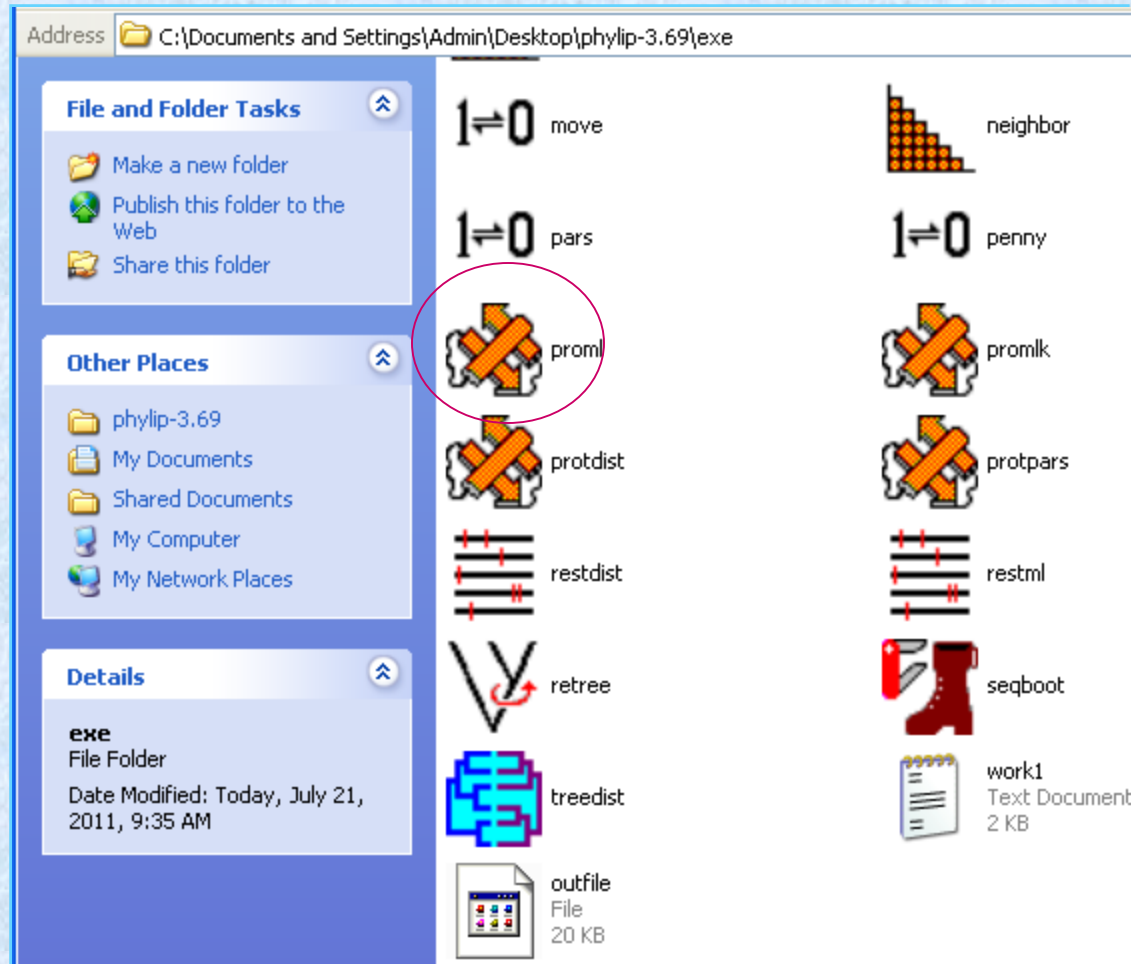
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR
ASSLLFFVVS	STLLSKYYYR	RR

	10	142									
sp P69905	MLAAAWVGA	AGEYYGGAAL	EEERFFLLL	FYYYYFPPH	DDLLHAAQQ	QVVVKGGKKA					
sp P69907	MLAAAWVGA	AGEYYGGAAL	EEERFFLLL	FYYYYFPPH	DDLLHAAQQ	QVVVKGGKKA					
sp P06635	MLAATWVGA	AGDYYGGAAL	EEERFFLLL	FYYYYFPPH	DDLLHAAQQ	QVVVKDGKKA					
sp P01966	MLAAAWVGA	AAEYYGGAAL	EEERFFLLL	FYYYYFPPH	DDLLHAAQQ	QVVVKGGKKA					
sp P01958	MLAAAWVGA	AGEYYGGAAL	EEERFFLLL	FYYYYFPPH	DDLLHAAQQ	QVVVKAGKKG					

outfile
10 different sets

Phylogenetic tree using Maximum likelihood method

The program is **proml**



outfile obtained
from **seqboot** is the
input for **proml**

C:\Documents and Settings\Admin\Desktop\phylop-3.69\exe\proml.exe

```
proml.exe: can't find input file "infile"
Please enter a new file name> outfile

proml.exe: the file "outfile" that you wanted to
use as output file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
(please type R, A, F, or Q)
F
Please enter a new file name> work1a
```

Amino acid sequence Maximum Likelihood method, version 3.69

Settings for this run:

```
U      Search for best tree?      Yes
P      JTT, PMB or PAM probability model? Jones-Taylor-Thornton
C      One category of sites?     Yes
R      Rate variation among sites? constant rate of change
W      Sites weighted?            No
S      Speedier but rougher analysis? Yes
G      Global rearrangements?     No
J      Randomize input order of sequences? No. Use input order
O      Outgroup root?            No, use as outgroup species 1
M      Analyze multiple data sets? No
I      Input sequences interleaved? Yes
0      Terminal type (IBM PC, ANSI, none)? IBM PC
1      Print out the data at start of run No
2      Print indications of progress of run Yes
3      Print out tree              Yes
4      Write out trees onto tree file? Yes
5      Reconstruct hypothetical sequences? No
```

Y to accept these or type the letter for one to change

```
m      Multiple data sets or multiple weights? (type D or W)
```

```
d      How many data sets?
```

```
10
```

Random number seed (must be odd)?

5

Number of times to jumble?

3

Y to accept these or type the letter for one to change

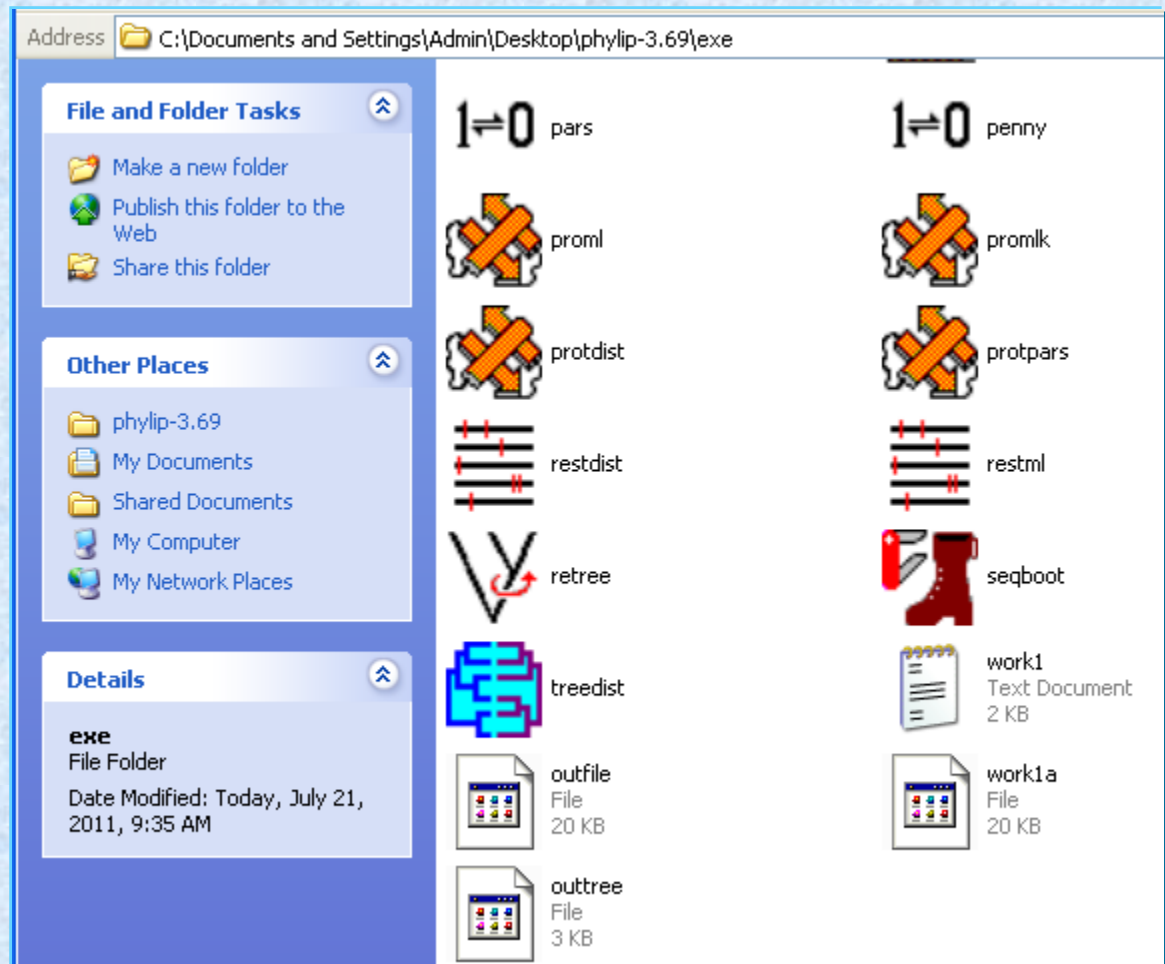
Y

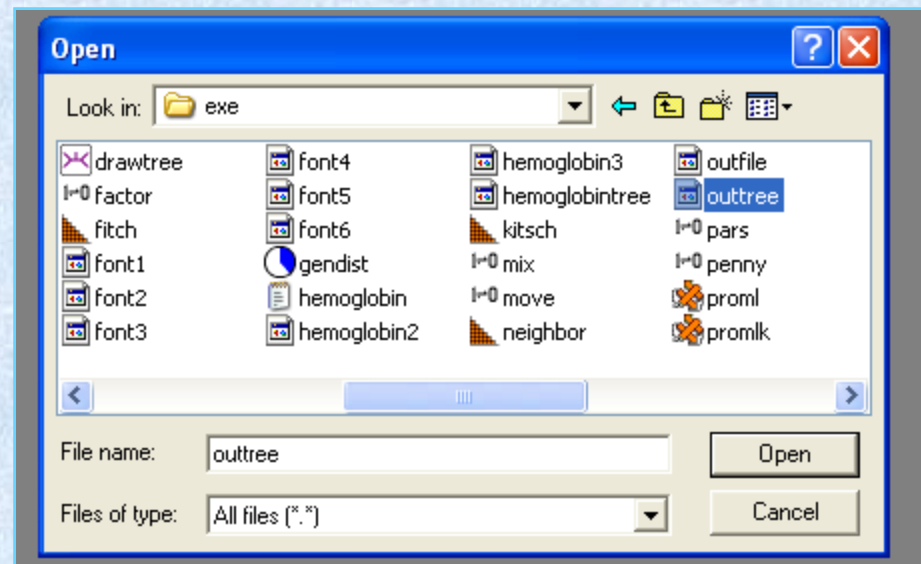
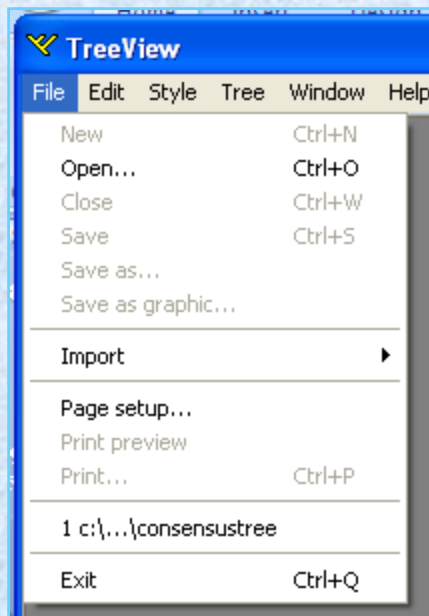
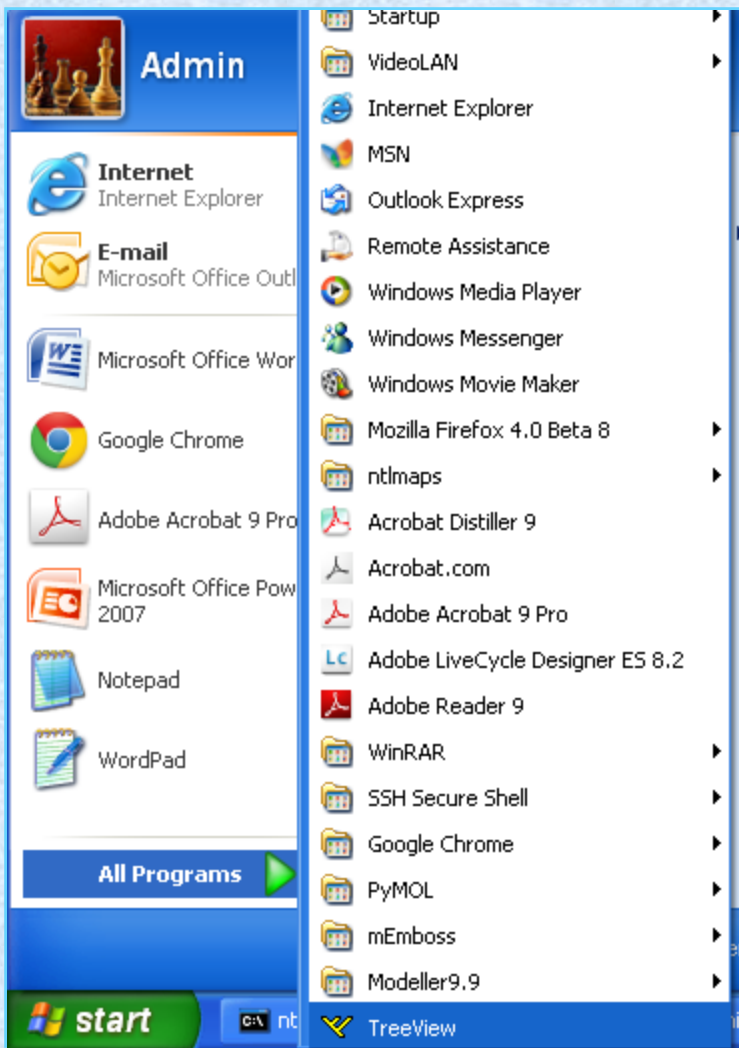
```
C:\Documents and Settings\Admin\Desktop\phylip

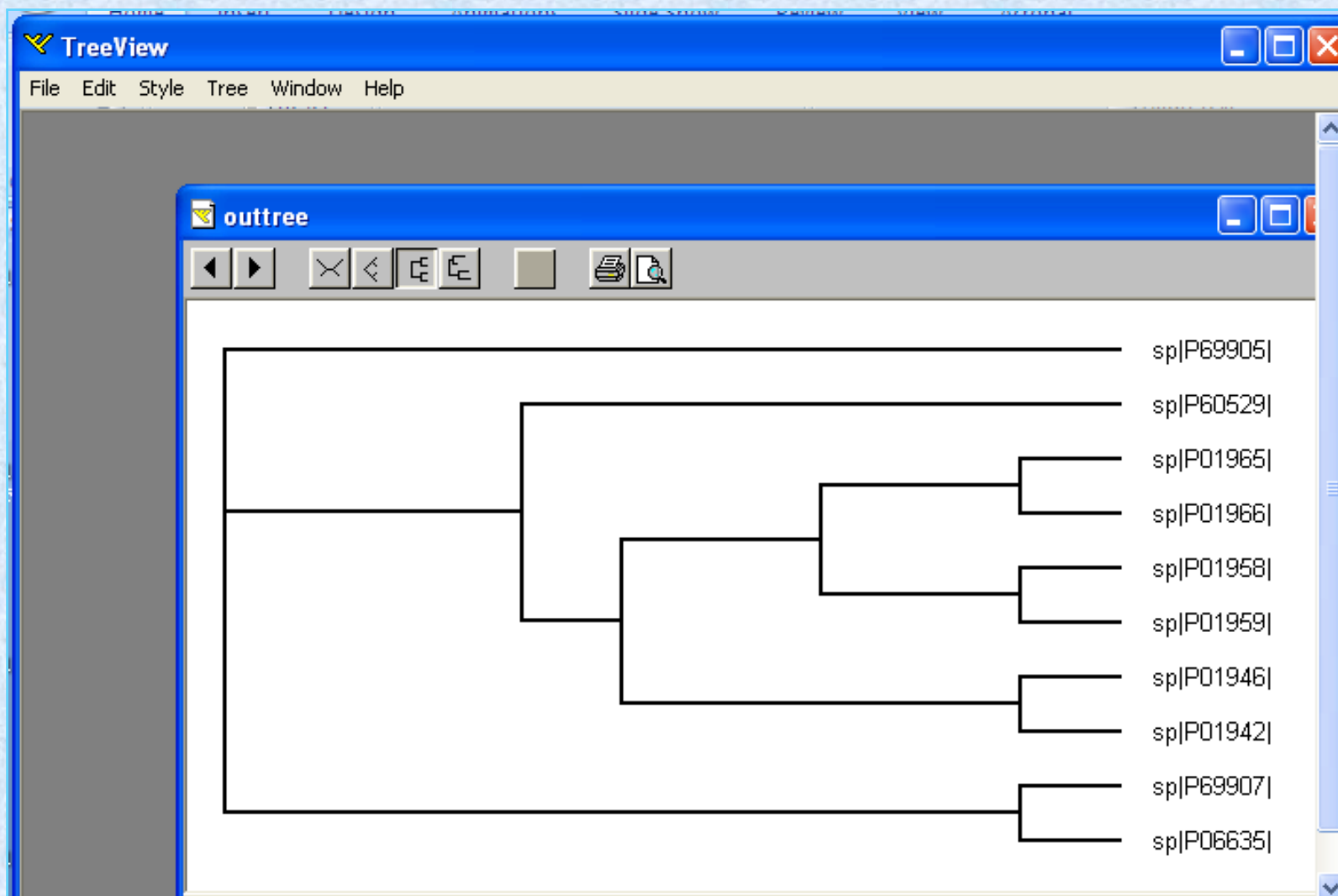
4. sp:P01959:
5. sp:P60529:
6. sp:P01946:
7. sp:P69907:
8. sp:P01958:
9. sp:P01965:
10. sp:P01942:

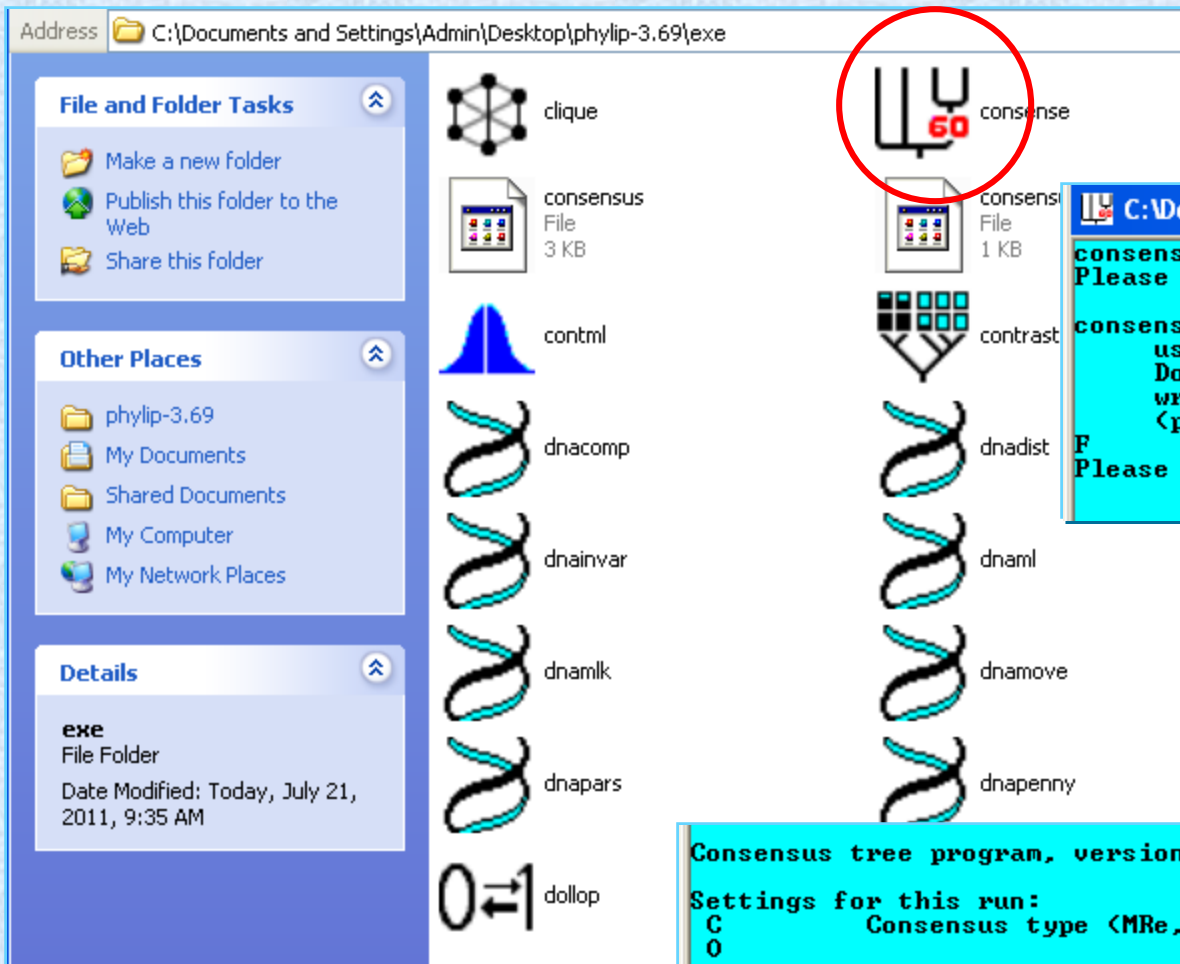
Adding species:
1. sp:P01958:
2. sp:P60529:
3. sp:P01966:
4. sp:P69907:
5. sp:P69905:
6. sp:P01959:
7. sp:P06635:
8. sp:P01965:
9. sp:P01946:
10. sp:P01942:

Output written to file "work1a"
Tree also written onto file "outtree"
Done.
Press enter to quit.
```









```
C:\Documents and Settings\Admin\Desktop\phylip-3.69\exe\consei
consense.exe: can't find input tree file "intree"
Please enter a new file name> outtree
consense.exe: the file "outfile" that you wanted to
use as output file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
<please type R, A, F, or Q>
F
Please enter a new file name> work1cons
```

```
Consensus tree program, version 3.69
Settings for this run:
C      Consensus type <MRe, strict, MR, M1>: Majority rule <extended>
0      Outgroup root: No, use as outgroup species 1

R      Trees to be treated as Rooted: No
T      Terminal type <IBM PC, ANSI, none>: IBM PC
1      Print out the sets of species: Yes
2      Print indications of progress of run: Yes
3      Print out tree: Yes
4      Write out trees onto tree file: Yes

Are these settings correct? <type Y or the letter for one to change>
y_
```

Consensus tree program, version 3.69

Settings for this run:

C Consensus type (MRe, strict, MR, ML): Majority rule (extended)
0 Outgroup root: No, use as outgroup species 1

R Trees to be treated as Rooted: No
T Terminal type (IBM PC, ANSI, none): IBM PC
1 Print out the sets of species: Yes
2 Print indications of progress of run: Yes
3 Print out tree: Yes
4 Write out trees onto tree file: Yes

Are these settings correct? (type Y or the letter for one to change)

y

consense.exe: the file "outtree" that you wanted to
use as output tree file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
(please type R, A, F, or Q)

f

Please enter a new file name> work1constree_

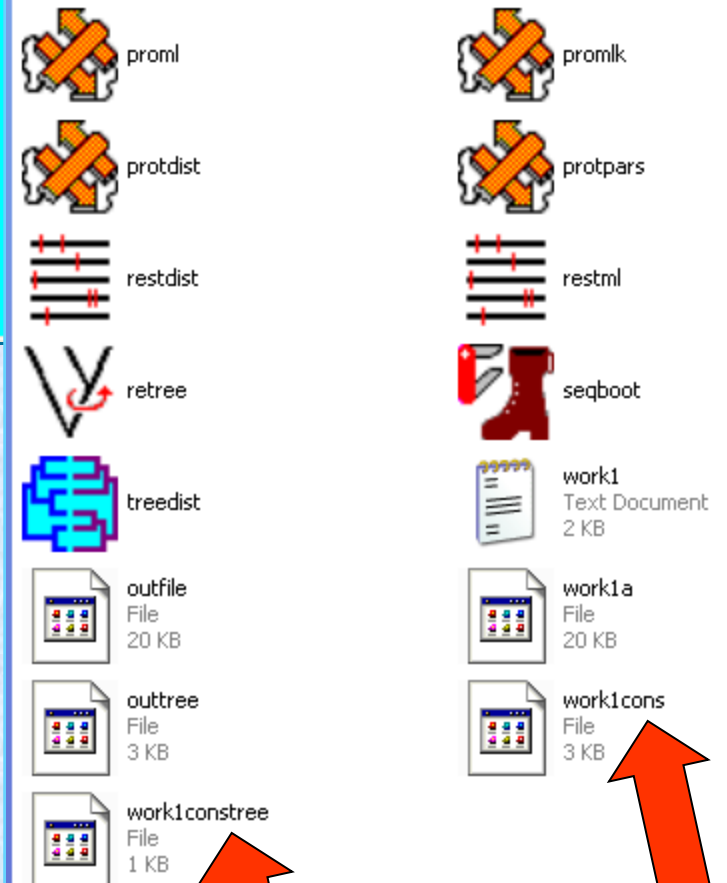
Consensus tree written to file "work1constree"

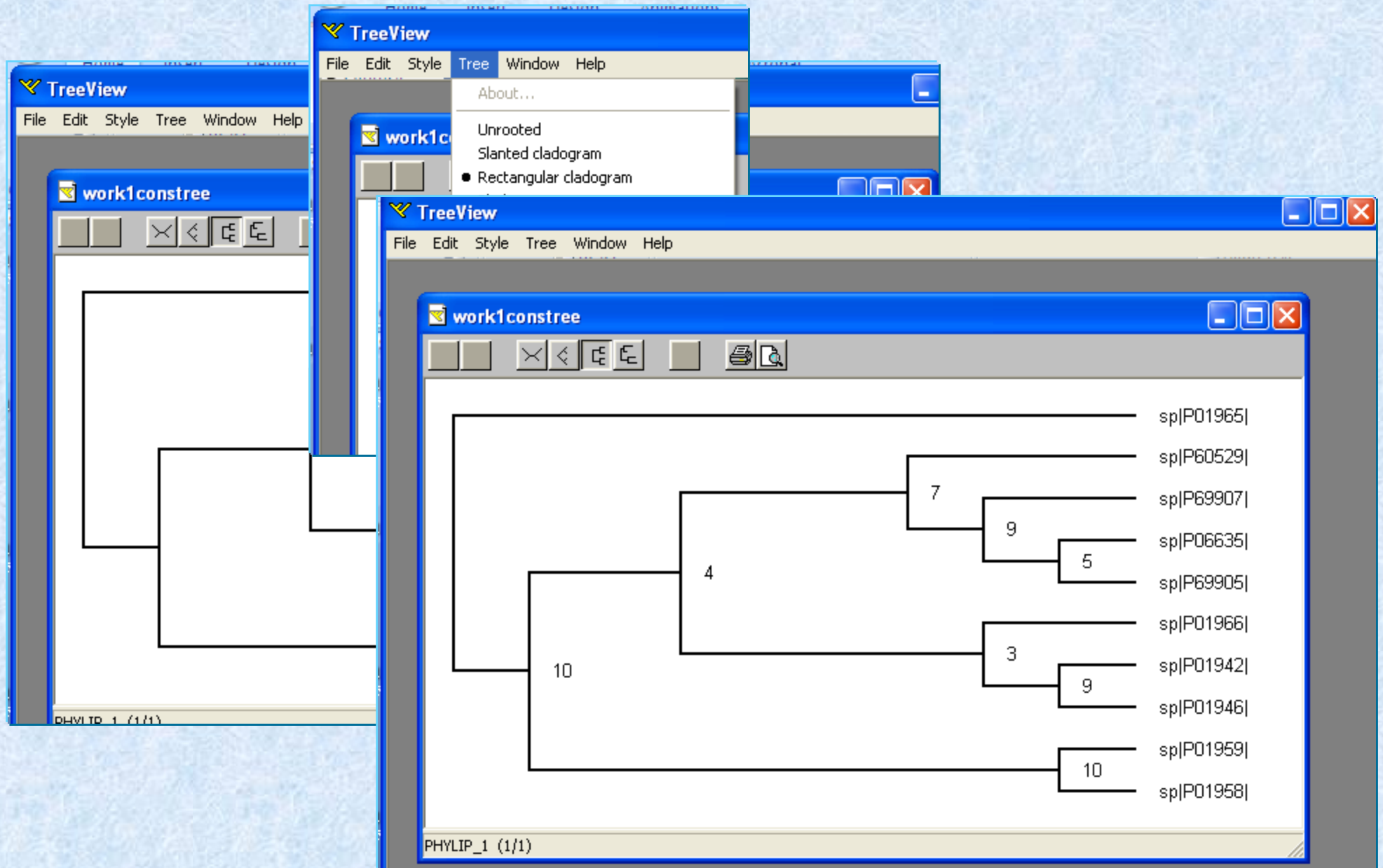
Output written to file "work1cons"

Done.

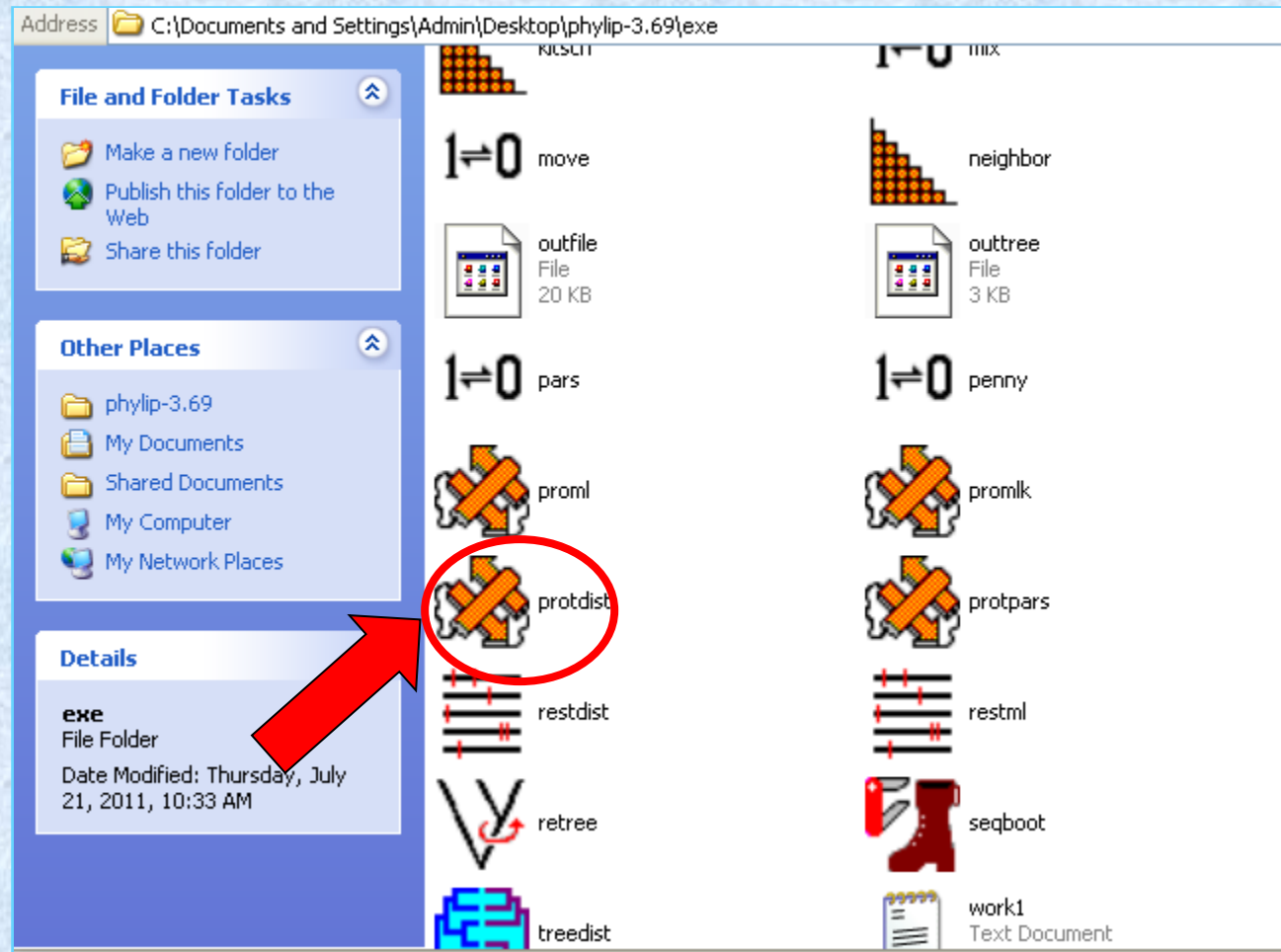
Press enter to quit.

Admin\Desktop\phylip-3.69\exe





NJ and UPGMA methods



```

protodist.exe: can't find input file "infile"
Please enter a new file name> outfile

protodist.exe: the file "outfile" that you wanted to
use as output file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
(please type R, A, F, or Q)
f
Please enter a new file name> work1-protodist

```

Protein distance algorithm, version 3.69

Settings for this run:

P	Use JTI, PMB, PAM, Kimura, categories model?	Jones-Taylor-Thornton matrix
G	Gamma distribution of rates among positions?	No
C	One category of substitution rates?	Yes
W	Use weights for positions?	No
M	Analyze multiple data sets?	No
I	Input sequences interleaved?	Yes
0	Terminal type (IBM PC, ANSI)?	IBM PC
1	Print out the data at start of run	No
2	Print indications of progress of run	Yes

Are these settings correct? (type Y or the letter for one to change)

m
Multiple data sets or multiple weights? (type D or W)

d
How many data sets?
10

Protein distance algorithm, version 3.69

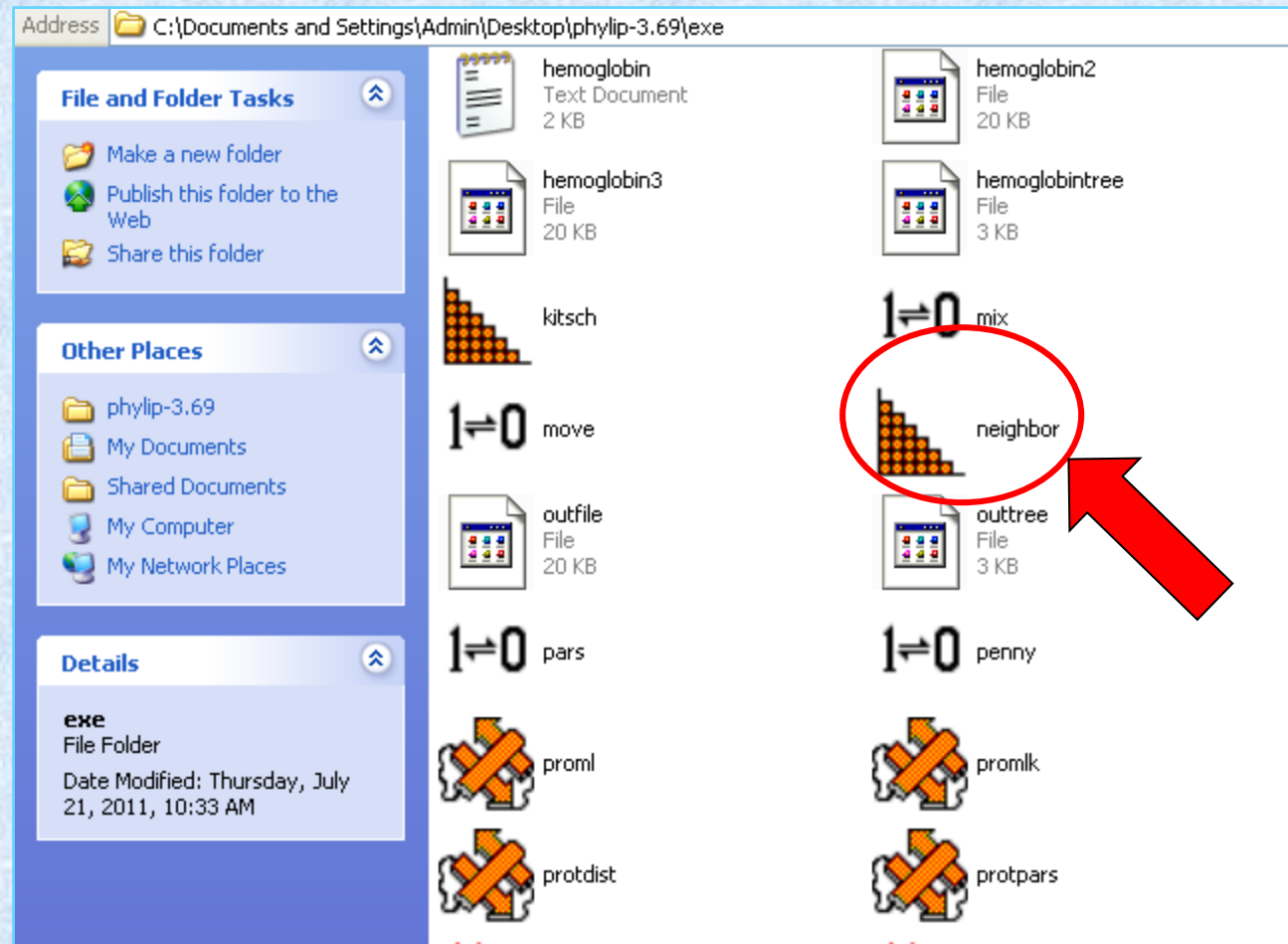
Settings for this run:

P	Use JTI, PMB, PAM, Kimura, categories model?	Dayhoff PAM matrix
G	Gamma distribution of rates among positions?	No
C	One category of substitution rates?	Yes
W	Use weights for positions?	No
M	Analyze multiple data sets?	Yes, 10 data sets
I	Input sequences interleaved?	Yes
0	Terminal type (IBM PC, ANSI)?	IBM PC
1	Print out the data at start of run	No
2	Print indications of progress of run	Yes

Are these settings correct? (type Y or the letter for one to change)

work1-protodist

sp Q9WVA2	0.000000	0.954472	1.211022	1.534523	2.384227
2.197177					
2.466581	2.303897	2.378322	2.245130	2.098663	2.066428
1.923666					
1.946335					
sp Q9Y5J9	0.954472	0.000000	1.114781	1.699316	2.007548
2.014938					
2.236423	1.873154	2.401207	2.523799	2.217026	2.237688
1.818816					
1.896587					
sp Q09783	1.211022	1.114781	0.000000	1.303653	2.771608
2.657542					
2.794046	2.191316	2.722125	2.545494	2.061901	2.150728
2.356218					
2.009404					
sp Q75DU7	1.534523	1.699316	1.303653	0.000000	2.572591
2.581259					
2.958788	2.454755	2.303727	2.134035	2.242738	1.951797
1.613806					
2.068258					
sp Q75F72	2.384227	2.007548	2.771608	2.572591	0.000000
0.263512					
0.532918	0.949154	1.037593	1.164635	1.410843	1.465146
1.892952					
2.173354					
sp Q6CJX3	2.197177	2.014938	2.657542	2.581259	0.263512
0.000000					
0.454437	0.916289	1.004302	1.030748	1.405926	1.425100
1.939355					



```
neighbor.exe: can't find input file "infile"
Please enter a new file name> work1-protdist

neighbor.exe: the file "outfile" that you wanted to
use as output file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
<please type R, A, F, or Q>
F
Please enter a new file name> work1-nj
```

Neighbor-Joining/UPGMA method version 3.69

Settings for this run:

```
N      Neighbor-joining or UPGMA tree? Neighbor-joining
O      Outgroup root? No, use as outgroup species 1
L      Lower-triangular data matrix? No
R      Upper-triangular data matrix? No
S      Subreplicates? No
J      Randomize input order of species? No. Use input order
M      Analyze multiple data sets? No
0      Terminal type (IBM PC, ANSI, none)? IBM PC
1      Print out the data at start of run No
2      Print indications of progress of run Yes
3      Print out tree Yes
4      Write out trees onto tree file? Yes
```

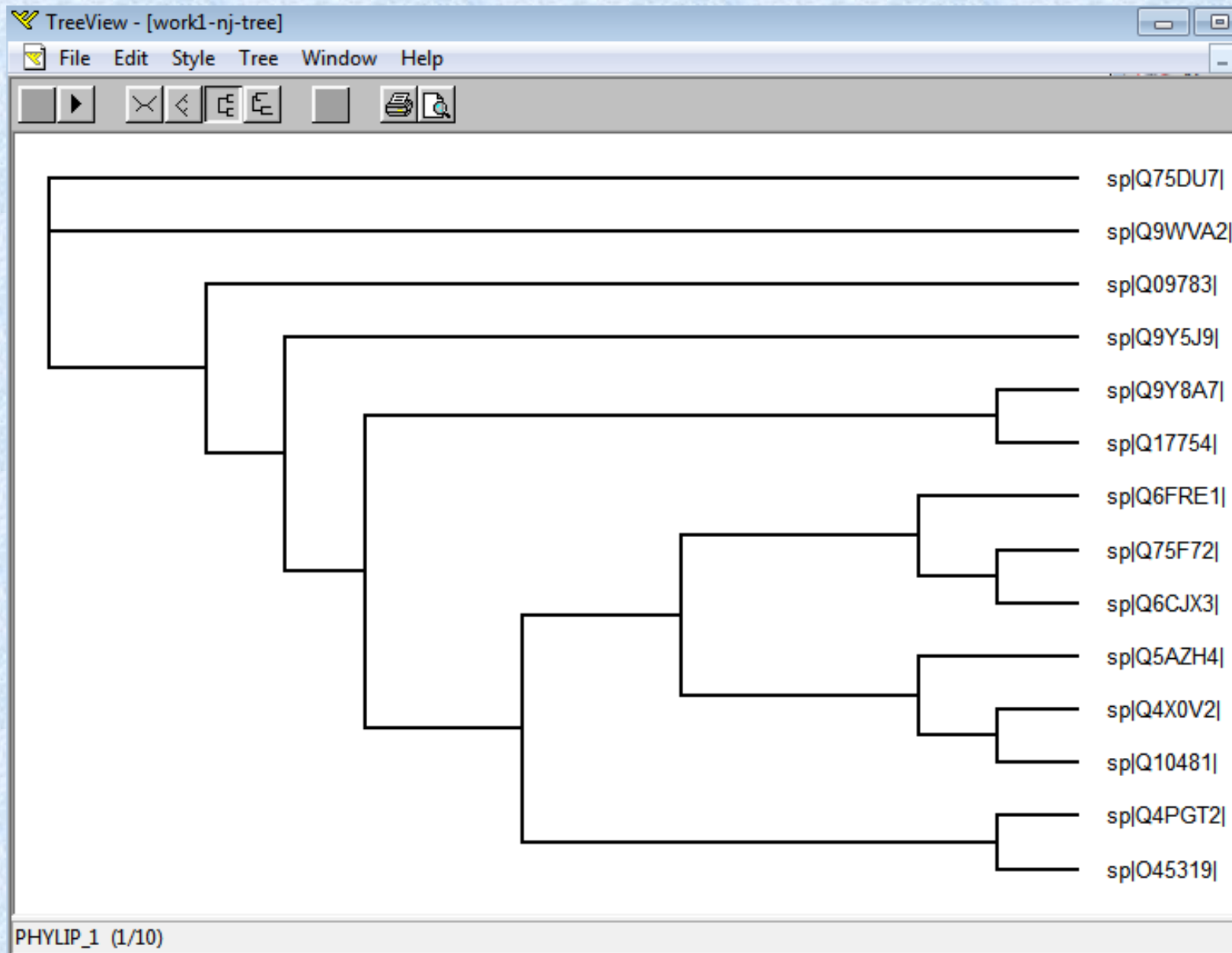
Y to accept these or type the letter for one to change

```
m      Settings for this run:
How many data sets? N      Neighbor-joining or UPGMA tree? Neighbor-joining
10      O      Outgroup root? No, use as outgroup species 1
Random number seed (must be odd)? L      Lower-triangular data matrix? No
5      R      Upper-triangular data matrix? No
S      Subreplicates? No
J      Randomize input order of species? Yes (random number seed = 5)
M      Analyze multiple data sets? Yes, 10 sets
0      Terminal type (IBM PC, ANSI, none)? IBM PC
1      Print out the data at start of run No
2      Print indications of progress of run Yes
3      Print out tree Yes
4      Write out trees onto tree file? Yes
```

Y to accept these or type the letter for one to change

```
y
neighbor.exe: the file "outtree" that you wanted to
use as output tree file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
<please type R, A, F, or Q>
f
Please enter a new file name> work1-nj-tree
```


|(sp|Q09783|:0.61703,(sp|Q9Y5J9|:0.71119,((sp|Q9Y8A7|:0.96547,
 sp|Q17754|:0.59919):0.21205,(((sp|Q6FRE1|:0.32455,(sp|Q75F72
 |:0.10986,
 sp|Q6CJX3|:0.11552):0.07268):0.29205,((sp|Q4X0V2|:0.40667,
 sp|Q10481|:0.30864):0.10478,sp|Q5AZH4
 |:0.45459):0.12561):0.17539,
 (sp|Q4PGT2|:0.67722,sp|O45319
 |:0.67765):0.02625):0.50607):0.47300):0.27258):0.09925,
 sp|Q75DU7|:0.56167,sp|Q9WVA2|:0.64156);
 ((sp|Q09783|:0.68210,sp|Q9Y5J9|:0.38881):0.03306,(((sp|Q9Y8A7
 |:0.64578,
 sp|Q17754|:0.49015):0.06214,((sp|O45319|:0.50317,sp|Q4PGT2
 |:0.42263):0.12931,
 (sp|Q5AZH4|:0.36428,(sp|Q4X0V2|:0.39579,(((sp|Q75F72|:0.12829,
 sp|Q6CJX3|:0.06919):0.12728,sp|Q6FRE1|:0.29582):0.35209,
 sp|Q10481|:0.38144):0.03146):0.07114):0.13883):0.70546):0.52085,
 sp|Q75DU7|:0.77374):0.34769,sp|Q9WVA2|:0.46575);
 ((sp|Q9Y5J9|:0.30392,sp|Q09783|:0.67455):0.07237,((sp|Q17754
 |:0.56886,
 ((sp|O45319|:0.52738,(sp|Q4PGT2|:0.49822,(sp|Q10481|:0.35123,
 (((sp|Q6CJX3|:0.08839,sp|Q75F72|:0.12819):0.10053,sp|Q6FRE1
 |:0.26834):0.28216,
 (sp|Q5AZH4|:0.33943,sp|Q4X0V2
 |:0.30149):0.12912):0.02663):0.11009):0.03799):0.40163,
 sp|Q75DU7|:0.96568):0.13386):0.11587,sp|Q9Y8A7
 |:0.62374):0.33497,sp|Q9WVA2|:0.37622);
 ((sp|Q75DU7|:0.81103,((sp|O45319|:0.60856,(((sp|Q6FRE1
 |:0.36825,
 sp|Q6CJX3|:0.11354):0.02000,sp|Q75F72|:0.15488):0.34593,
 sp|Q10481|:0.29363):0.02586,((sp|Q4PGT2|:0.56535,sp|Q5AZH4
 |:0.38322):0.10438,
 sp|Q4X0V2|:0.40429):0.07203):0.20545):0.71213,(sp|Q9Y8A7



```

consense.exe: can't find input tree file "intree"
Please enter a new file name> work1-nj-tree

consense.exe: the file "outfile" that you wanted to
use as output file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
<please type R, A, F, or Q>
f
Please enter a new file name> work1-nj-cons

```

Consensus tree program, version 3.695

Settings for this run:

```

C      Consensus type (MRe, strict, MR, ML):  Majority rule (extended)
0      Outgroup root:                        No, use as outgroup species

R      Trees to be treated as Rooted:        No
T      Terminal type (IBM PC, ANSI, none):   IBM PC
1      Print out the sets of species:        Yes
2      Print indications of progress of run:  Yes
3      Print out tree:                       Yes
4      Write out trees onto tree file:        Yes

```

Are these settings correct? <type Y or the letter for one to change>

y

```

consense.exe: the file "outtree" that you wanted to
use as output tree file already exists.
Do you want to Replace it, Append to it,
write to a new File, or Quit?
<please type R, A, F, or Q>
f

```

Please enter a new file name> work1-nj-constree

```

|((((((sp|Q9Y8A7|:10.0,sp|Q17754|:10.0):9.00,(((sp|Q4X0V2
|:10.0,sp|Q5AZH4|:10.0):6.00,
(((sp|Q6CJX3|:10.0,sp|Q75F72|:10.0):9.00,sp|Q6FRE1
|:10.0):10.0,sp|Q10481|:10.0):5.00):8.00,
sp|Q4PGT2|:10.0):6.00,sp|O45319|:10.0):10.0):9.00,sp|Q75DU7
|:10.0):7.00,sp|Q9WVA2|:10.0):7.00,
sp|Q9Y5J9|:10.0):10.0,sp|Q09783|:10.0);

```

