

Phylogenetics

Table of contents

1. Introduction & Aims
2. Exercise 1: *ompA* phylogeny from gene sequences
3. Exercise 2: MLST gene phylogeny using Artemis
4. Exercise 3: Phylogeny from whole genome SNPs
5. Additional Information

1. Introduction & Aims

Phylogeny has important applications in many fields of genome biology. For example, when annotating a gene in a new genome it is useful for identifying previously-annotated genes in other genomes that share a common ancestry. It is also becoming increasingly common to use phylogeny to trace the evolution and spread of bacterial diseases, and even as an epidemiological tool to help identify disease outbreaks in a clinical setting. Further analysis of genome sequences to examine recombination, molecular adaptation and the evolution of gene function, all benefit from phylogeny.

Phylogenetics is essentially about similarity, and looking at patterns of similarity between taxa to infer their relationships. Although the methods may initially appear quite daunting, the basic ideas behind phylogenetics are quite simple. We want to identify the tree that best fits our data assuming that the data evolve under a simple model. When the data are DNA, we can use our knowledge of biology to improve our chances of finding the correct tree by defining evolutionary models that make biological sense. For example, we know that transition mutations occur more frequently than transversions, so we can make our model favour trees in which this is the case with our data.

In this module we will use phylogenetics to explore the strengths and weaknesses of techniques that may be used for typing *Chlamydia trachomatis*.

The exercise will begin by using assembled sequences. Historically this is the type of data that would have been available for molecular phylogenetic analyses. We will look at the *ompA* gene, the most commonly used Chlamydia typing locus.

In the second part of the exercise, we will see how Artemis can be used to create sequence alignments from the *C. trachomatis* bcf files used in the mapping module. We will extract the sequences of genes used in a Chlamydia multilocus-typing scheme and compare the results with those from the *ompA* analysis.

Finally we will create a tree based on whole-genome SNP data and see which of the other two typing schemes it supports.

By the end of this module you should have some understanding of: 1. Multiple sequence alignment using ClustalX and/or muscle. 2. Maximum likelihood phylogenetic estimation using a nucleotide model. 3. Evaluation of support for relationships in a tree using non-parametric bootstrap replicates. 4. Interpretation of phylogenetic trees.

Back
to
top

Ex-
er-
cise
1:
ompA
phy-
logeny
from
gene
se-
quences
For
our
first
phy-
logeny
we
will
make
a
tree
us-
ing
the
se-
quences
of
the
ompA
gene
from
16
strains
of *C.*
tra-
choma-
tis.
These
se-
quences
are
pro-
vided
for
you
in a
sir
gle
file
called
ompA_assembled.mfa
in
the
Mod-

Historically,
the
most
com-
monly
used
tool
for
typ-
ing
C.
tra-
choma-
tis
iso-
lates
was
serotyp-
ing
us-
ing
the
MOMP
(ma-
jor
outer
mem-
brane
pro-
tein),
which
is en-
coded
by
the
ompA
gene.
There
are
two
bio-
vars
of *C.*
tra-
choma-
tis:
1.
the
tra-
choma
bio-
var
in-
cludes
ocu-

Based
on
MOMP
serotyp-
ing,
C.trachomatis
has
been
sub-
di-
vided
into
be-
tween
15
and
19
serotypes:
the
tra-
choma
bio-
var
in-
cludes
ocu-
lar
serotypes
A to
C
and
uro-
geni-
tal
serotypes
D
to K,
while
the
LGV
bio-
var
in-
cludes
serotypes
L1,
L2
(in
clud-
ing
L2a,
b
and
c)
and

###

Viewing
the
align-
ment
in
Seav-
iew
To
view
and
edit
align-
ments
and
pro-
duce
phy-
loge-
nies
we
will
use
a
pro-
gram
called
Seav-
iew.
Seav-
iew
is a
graph-
ical
user
in-
ter-
face
(GUI)
that
com-
bines
a
num-
ber
of
the
most
pop-
ular
align-
ment
and
phy-

First
you
should
navi-
gate
to
Mod-
ule 4
di-
rec-
tory.
“‘bash

change
to
the
work-
ing
di-
rec-
tory
cd
Mod-
ule_4_Phlogenetics

We
need
to
in-
stall
two
tools
quickly
sudo
apt-
get
in-
stall
clustalo
sudo
apt-
get
in-
stall
phyml

```
#  
The  
pass-  
word  
is  
man-  
ager  
and  
press  
'Y'  
when  
prompted  
for  
the  
in-  
stall  
#  
start  
Seav-  
iew  
seav-  
iew  
#  
load  
the  
align-  
ment  
file  
ompA_assembled.mfa  
by  
se-  
lect-  
ing  
'Open'  
from  
the  
'File'  
menu.  
««
```



If
you
have
time,
have
a
look
at
the
map-
ping
of
the
ompA
gene
in
some
of
your
sam-
ples.
The
di-
verse
na-
ture
of
the
ompA
re-
gion,
with
high
SNP
den-
sity
and
a
num-
ber
of
in-
dels,
makes
map-
ping
of
the
re-
gion
diffi-
cult
and
variant-
calling

Mul-
tiple
se-
quence
align-
ment
Be-
fore
any
phy-
loge-
netic
anal-
ysis,
we
must
make
sure
that
the
columns
in
our
data
rep-
re-
sent
ho-
mol-
o-
gous
bases.
With
gene
or
pro-
tein
se-
quence
data,
this
usu-
ally
means
align-
ing
the
nu-
cleotide
or
amino
acid
se-
quences

Seaview
al-
lows
align-
ment
us-
ing
two
pro-
grams,
clustal
and
mus-
cle.
Gen-
er-
ally
mus-
cle is
faster,
and
the
pro-
tein
align-
ments
are
of
simi-
lar
qual-
ity
to
clustal.
In
both
cases,
se-
quences
are
aligned
by
as-
sign-
ing
costs
to
par-
ticu-
lar
base
changes
and
gap
in-

On
the
'Props'
menu
choose
'View
as
pro-
teins'.



To
start
the
align-
ment,
se-
lect
'Align'
then
'Align
all'.

When
the
align-
ment
pro-
cess
is
com-
plete,
Seav-
iew
will
have
in-
serted
gaps
into
the
se-
quences
so
that
ho-
mol-
o-
gous
sites
(or
at
least
ho-
mol-
o-
gous
ac-
cord-
ing
to
the
align-
ment
pro-
gram)
are
lined
up
in
columns.



If
you
in-
spect
the
align-
ment,
it
should
be
clearer
how
the
se-
quences
dif-
fer
from
one
an-
other.
Can
you
see
which
se-
quences
are
most
closely
re-
lated?

If an alignment has been problematic, requiring many gaps, it is advisable to inspect the it by eye and edit where necessary. In Seaview you can add gaps with the space bar, and remove them with the backspace (for more detailed instructions see

If
you
turn
off
pro-
tein
view
you
can
see
that
the
nu-
cleotides
are
also
now
aligned.



Phy-
logeny
esti-
ma-
tion
us-
ing
PhyML
To
esti-
mate
the
phy-
logeny,
we
will
use
a
pro-
gram
called
PhyML,
which
is in-
cluded
in
Seav-
iew.
PhyML
uses
max-
i-
mum
like-
li-
hood
(ML)
to
esti-
mate
the
tree.
We
will
use
ML
be¹⁷
cause
it is
more
ac-
cu-
rate
than



PhyML
in-
cludes
a
num-
ber
of
nu-
cleotide
sub-
sti-
tu-
tion
mod-
els.
The
strength
of
us-
ing a
Max-
i-
mum
Like-
li-
hood
method
is
that
an
ex-
plicit
model
of
nu-
cleotide
sub-
sti-
tu-
tion
is
ap-
plied,
which
can
be
more
¹⁹
logi-
cally
real-
istic
than
some
other

In-
ter-
pre-
ta-
tion
of
phy-
logeny
Once
the
run
has
fin-
ished,
click
'OK'.
The
tree
cre-
ated
by
PhyML
not
only
in-
cludes
the
topol-
ogy
of
tree
(i.e.,
the
rela-
tion-
ships
be-
tween
se-
quences)
but
also
the
branch
lengths
(i.e.,
the
amount
of
change
oc-
cur-
ring
in



If
you
are
run-
ning
be-
hind
at
this
point,
skip
now
to
part
ii of
the
exer-
cise.
You
can
come
back
to
this
part
when
you
have
time.
Oth-
er-
wise,
note
down
the
like-
li-
hood
of
the
tree.
We
will
use
it
later.

Phy-
logeny
esti-
ma-
tion
with
across
site
rate
het-
ero-
gene-
ity

Almost
all
nu-
cleotide
se-
quences
in
na-
ture
dis-
play
across
site
rate
het-
ero-
gene-
ity.
This
means
that
not
all
sites
within
the
se-
quence
evolve
at
the
same
rate;
rather
some
parts
of a
gene
evolve
faster
than
oth-
ers,
an
ac-
tive
site
of²⁵
an
en-
zyme
for
ex-
am-
ple,



Select
‘Op-
ti-
mized’

from
the
‘In-
vari-
able
sites’
box
to
al-
low
the
pro-
por-
tion
of
in-
vari-
ant
sites
to
change.

Af-
ter
se-
lect-
ing
these
pa-
ram-
e-
ters,
en-
sure
that
you
have
se-
lected
the
same
choice
of
sub-
st²⁷
tu-
tion
model
as
you
did
in

Comparing models with the likelihood ratio test. The likelihood ratio test (LRT) can be used to statistically test the difference in fit of two nested evolutionary models to the data.

'Nested'
means
that
the
more
com-
plex
model
must
in-
clude
all
of
the
pa-
ram-
eters
of
the
sim-
pler
model.

For
com-
pari-
son
of
non-
nested
mod-
els,
more
com-
plex
meth-
ods
are
avail-
able
- see
the
notes

on
model
se-
lec-
tion
at²⁹
the
end
of
this
mod-
ule.

In-

To
per-
form
a
LRT
we
must
first
cal-
cu-
late
the
like-
li-
hood
ratio
(LR)
of
our
two
mod-
els:
LR
=

$$2x(\text{neglogL1} - \text{ne-} \text{glogL2})$$

Where:

ne-
glogL1
is
the
neg-
a-
tive
log
like-
li-
hood
of
the
sim-
pler
model,
and
ne-
glogL2
is³⁰
the
neg-
a-
tive
log
like-
li-

The
LRT
statis-
tic
ap-
prox-
i-
mately
fol-
lows
a
chi-
square
dis-
tri-
bu-
tion,
so
we
can
eval-
uate
the
sig-
nifi-
cance
of
our
LR
us-
ing
chi-
square
sig-
nifi-
cance
ta-
bles
(or
cal-
cu-
late
p-
values
us-
ing
sta-
tisti-
cal
cal-
cula-
tors)

However,
as
with
any
chi-
square
sig-
nifi-
cance
test,
we
need
to
know
the
de-
grees
of
free-
dom
(df).
In a
LRT
the
df
are
the
dif-
fer-
ence
in
num-
ber
of
free
pa-
ram-
eters
be-
tween
the
two
mod-
els.

The
num-
ber
of
free
pa-
ram-
eters
for
the
mod-
els
in
Seav-
iew
are
listed
in
the
ta-
ble
of
mod-
els
on
page
6 of
this
mod-
ule.
e.g. the
GTR
model
has
8
free
pa-
ram-
e-
ters.
Both
the
gamma
pa-
ram-
eter
of
among
site
rate
vari-
a-
tion
and
the

To help you perform your LRT, we have provided a pre-formatted spread-sheet for calculating significant cancellation from your negative attitude log like-likely-hoods and the number of free parameters of the models you used.

34

Start
Li-
bre-
Of-
fice
Calc
us-
ing
the
icon
on
the
left
of
the
screen
and
open
LR_test.xls
from
the
phy-
logeny
di-
rec-
tory
of
Mod-
ule
5.
Type
the
neg-
a-
tive
log
like-
li-
hood
val-
ues
you
recorded
for
your
two
trees
³⁵
into
the
ap-
pro-
pri-
ate
boxes.
Type

QUESTION

Question:

- Do
the
addi-
tions
of
the
gamma
and
in-
vari-
ant
sites
pa-
ram-
eters
sig-
nifi-
cantly
im-
prove
the
fit of
the
model
to
the
data?

###

Phy-
logeny
esti-
ma-
tion
with
boot-
strap-
ping

Boot-
strap-
ping

is a
sta-
tisti-
cal
tech-
nique
for
adding
con-
fi-
dence
in-
ter-
vals
around
an
esti-
mate,
in
this
case,
a
tree
topol-
ogy.
Non-
parametric
boot-
strap-
ping
in-
volves
re-
peated
anal-
ysis
of
the
data
set
through
“re-

Imagine putting each site into a bag. Repli-
cate data sets are cre-
ated by ran-
domly draw-
ing sites from
the bag until
a new dataset
the same size
as the orig-
inal has been
cre-
ated. Im-
por-
tantly, after
a site has
been drawn,
it³⁸ is
re-
placed back
into the bag.
This

Trees
are
then
built
for
each
repli-
cate
data
set.
Ro-
bust
rela-
tion-
ships,
i.e. those
that
are
re-
peat-
able,
will
oc-
cur
in a
large
pro-
por-
tion
of
ran-
domised
data
sets.

Estimate
a
boot-
strapped
phy-
logeny
for
the
ompA
data
set
by
cre-
at-
ing
a
new
phy-
logeny
as
be-
fore,
with
the
addi-
tion
of
10
boot-
strap
repli-
cates.
Click
on
'Boot-
strap'
in
the
'Branch
Sup-
port'
box,
and
en-
ter
'10'
in
the
repli-
cates
box.
Pro-
cess-
ing
of

WARNING!:

boot-
strap
pro-
por-
tions
are
mea-
sures
of
ro-
bust-
ness,
or
re-
peata-
bil-
ity.

A
high
boot-
strap
value
indi-
cates
that
a
given
node
tends
to
oc-
cur
in
ev-
ery
anal-
ysis.

This
does
not
guar-
an-
tee
that
the
node
is⁴²

**cor-
rect.**

For
ex-
am-
ple,
if

You
will
no-
tice
that
PhyML
in-
cludes
an-
other
branch
sup-
port
method
called
aLRT.
aLRT
as-
sesses
whether
each
branch
on
the
tree
pro-
vides
a
sig-
nifi-
cant
like-
li-
hood
im-
prove-
ment
over
the
same
tree
with
that
branch
col-
lapsed.
We
48
will
not
use
this
method
here
or
dis-

Once
the
search
is
com-
plete,
you
can
show
the
boot-
strap
val-
ues
on
the
tree
by
tick-
ing
this
box.
Each
node
in
the
tree
now
has
an
asso-
ci-
ated
value
out
of
100,
its
boot-
strap.
Can
you
iden-
tify
any
nodes
that
are
not
ro-
bust?
Un-
for-
tu-
nately



Questions:

From
the
trees
that
you
have
pro-
duced,
which
MOMP
type
would
you
sug-
gest
the
new
vari-
ant
(NV)
strain
be-
longs
to? -
Do
the
ompA
trees
agree
with
the
sepa-
ra-
tion
of *C.*
tra-
choma-
tis
into
tra-
choma
(serotypes
A to
K)
and
LEV
(L
serotypes)
bio-
vars?

[Back to top](#)

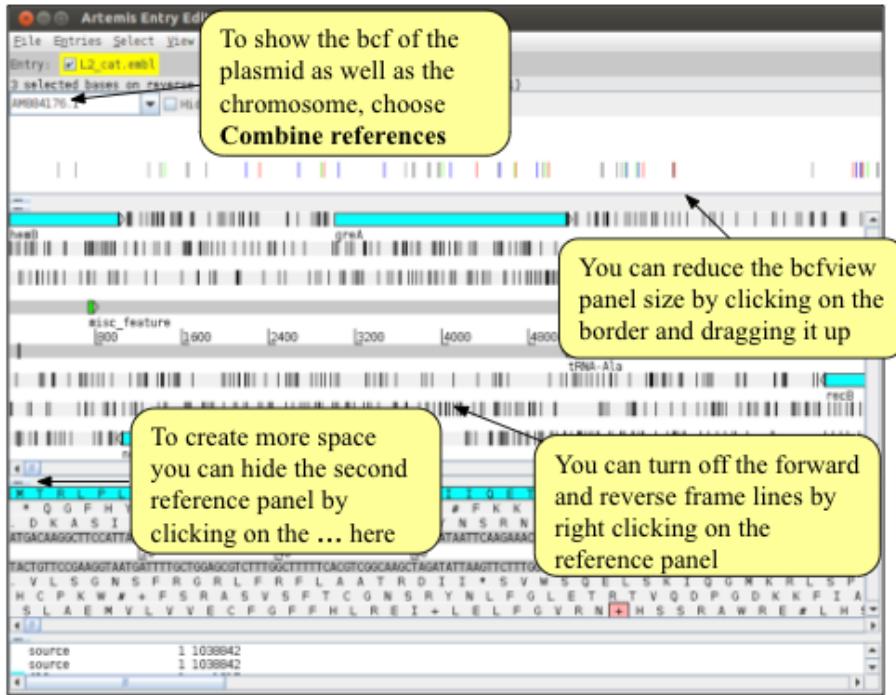
Exercise 2: MLST gene phylogeny using Artemis

A second typing method used for many bacterial species is multilocus sequence typing (MLST). MLST involves the sequencing of fragments of a number (usually 6 or 7) of housekeeping genes spread around the genome. In true MLST, each different allele for each locus is assigned a number. Each unique allelic profile defines a sequence type (ST).

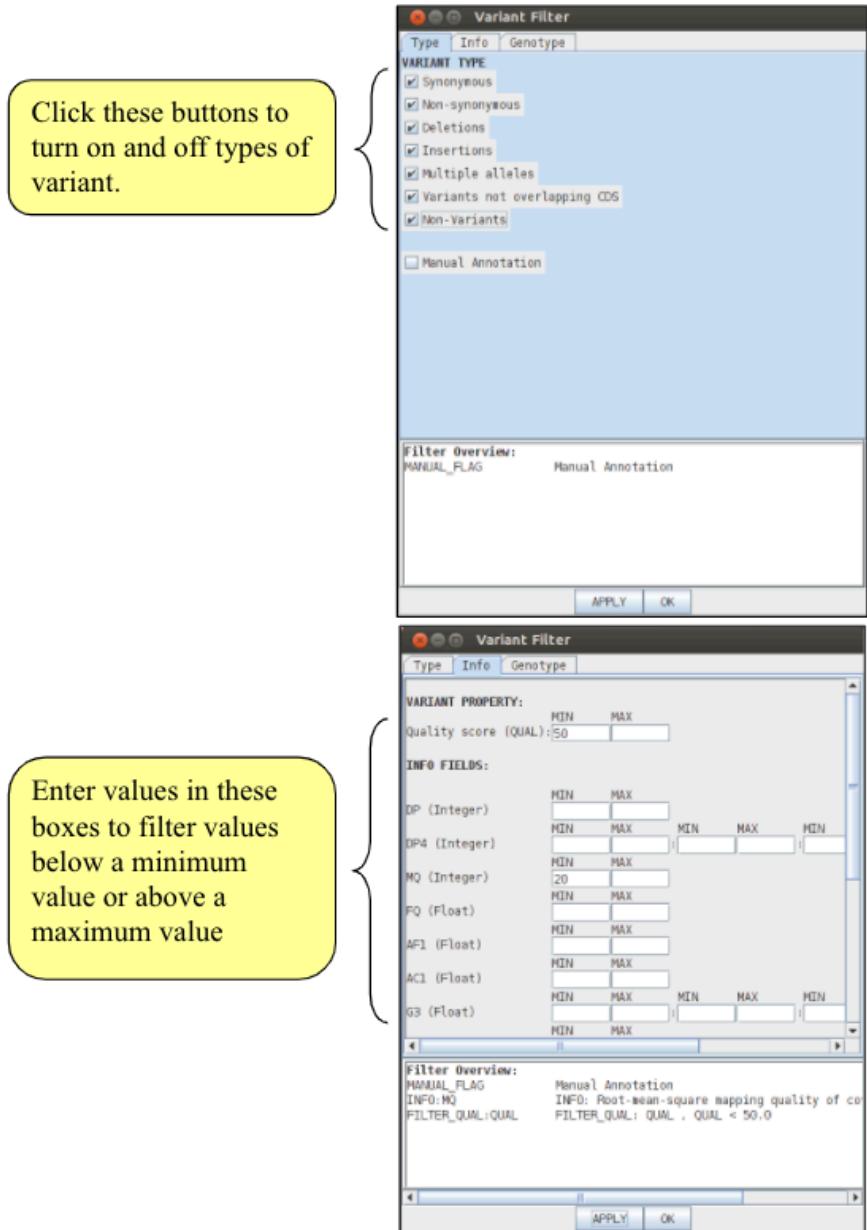
A number of MLST schemes have been devised for *C. trachomatis*, but we will use the scheme of Dean et al. (Emerg Infect Dis. 2009 Sep;15(9):1385-94.), which comprises the following seven loci.

Gene name	Locus tag in L2_cat.embl
<i>glyA</i>	CTL0691
<i>mdhC</i>	CTL0630
<i>pdhA</i>	CTL0497
<i>yhbG</i>	CTL0022
<i>pykF</i>	CTL0586
<i>lysS</i>	CTL0150
<i>leuS</i>	CTL0461

If **Artemis** is not open, start it now and open the L2_cat.fasta reference. Read in the annotation (L2_cat.embl) by '**Read Entry Into**'. Open the NV.vcf.gz file by selecting '**Read BAM / VCF**' from the '**File**' menu.



As you saw in the mapping module, bcf files contain support values for each variant. Before writing out the alignment of the MLST genes it is crucial to filter the variants so that only strongly supported variants are included in the alignment. Right click on the bcfview panel and select the Filter option. You can choose your own filters, but we would suggest something similar to the following:



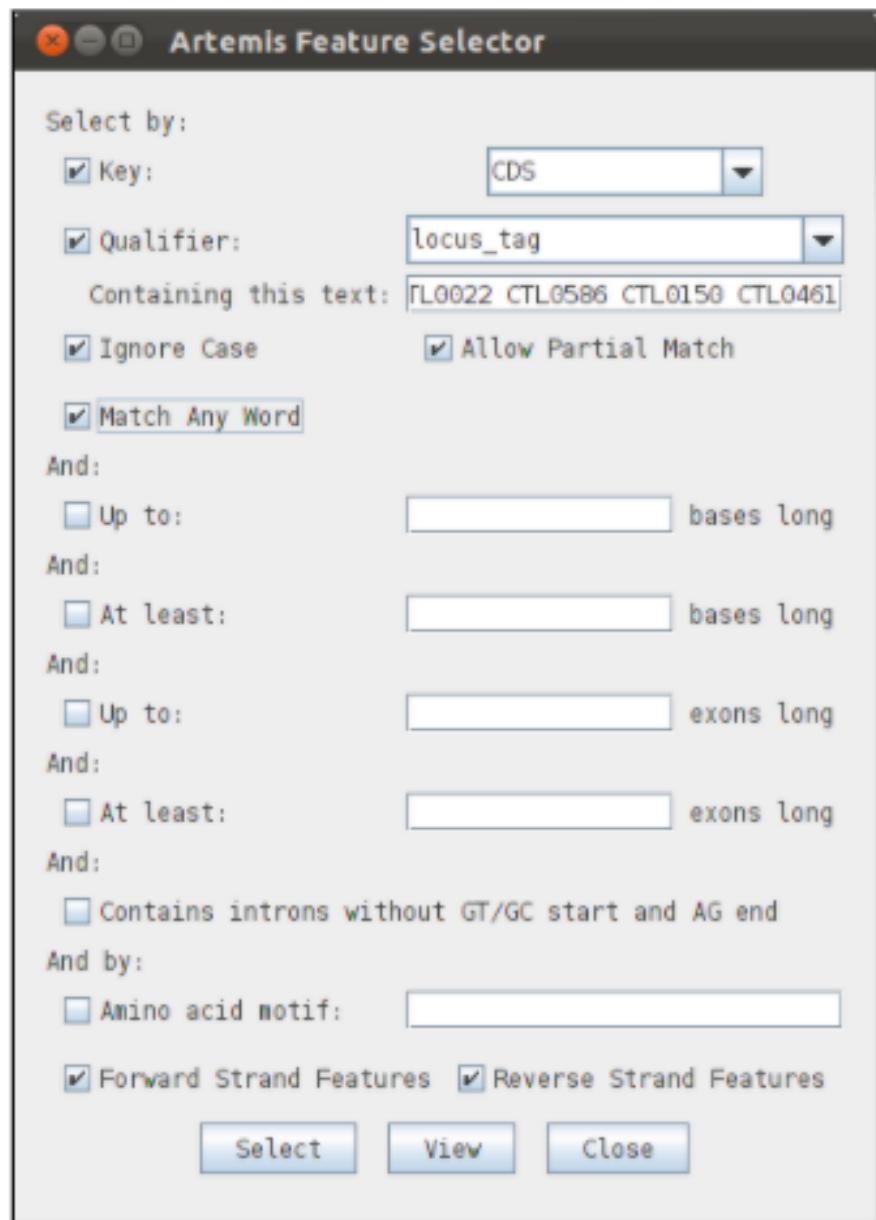
Importantly, note that we have also selected the ‘non-variants’ tick box. This option tells Artemis to also show sites which match the reference, rather than just variant sites, provided, of course, that they pass the filters. This information is necessary when saving an alignment, as it allows differentiation between regions that show no variation because they are the same as the reference from those that show no variation due to a lack of mapping. To apply the filters click

on the ‘Apply’ button. Note what happens to the bcfview in the MLST gene regions. Once you are happy with your filters, close the Variant filter box.

If you have time, have a look at some of the MLST genes in Artemis. Question:
- How does their diversity within the *C. trachomatis* strains compare to that of ompA? - Can you think of a possible reason that the diversity may be different?

Rather than use fragments of these loci in a true MLST typing scheme, we will extract the sequences of the seven genes and run a phylogenetic analysis on the concatenated sequence.

There are many ways in which you could identify and select the seven MLST genes in Artemis, but one convenient method is using the ‘Feature Selector’ in the ‘Select’ menu. We want to search for CDSs with the locus tags in the MLST scheme, so in the ‘Key’ dropdown box select CDS, and in the ‘Qualifier’ box select or type locus_tag. You then now type the locus tags of the seven loci into the ‘Containing this text’ box, and make sure to tick the ‘Match Any Word’.



To find the genes click on the select button. To check the results of your search, click on the view button. This will list the features that have been selected by your search.

All features with key "CDS" with qualifier "locus_tag"

CDS	30083	30802	c	
CDS	196047	197627		
CDS	554203	556662	c	
CDS	592015	593037		
CDS	692793	694250		
CDS	747767	748747	c	
CDS	819539	821032	c	

Check that you have selected the correct seven genes. Next we will write out their sequences into a fasta file.

Artemis Entry Edit: L2_cat.fasta

File Entries Select View Info Edit Create Run Graph Display
 Entry: L2_cat.fasta L2_cat.embl
 7 selected features total bases 9717 total amino acids 3232 (CTL0022 glyA adhC pykF pdhA lacS lysS)

Add VCF ...
 VCF files

Mark new stops within CDS features

Separate out into samples
 Filter ...
 Colour By

Create
 Write
 View
 Graph
 Overview for selected features
 Show Labels
 Show details of : AM884176.1:58976 .
 Manual PASS / FAIL

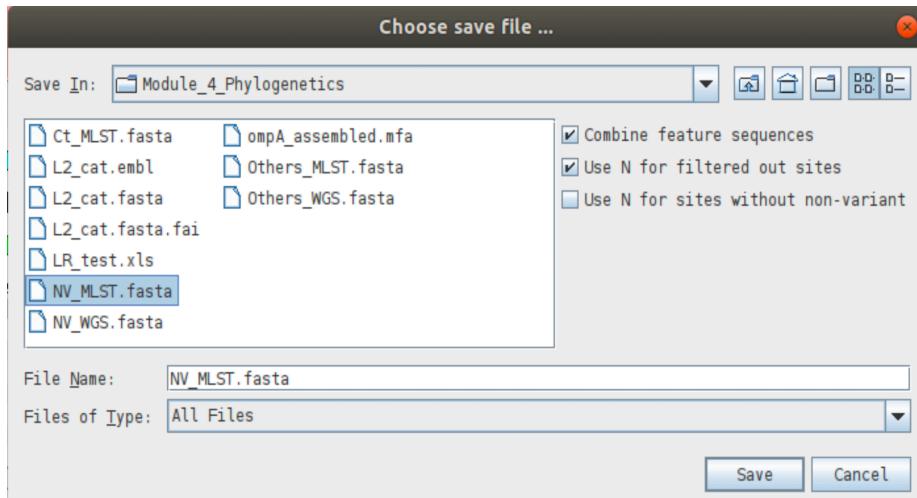
Filtered VCF
 FASTA of selected feature(s) ...
 FASTA of selected base range ...
 FASTA of variant sites only ...

Choose this option to write the sequence of one or more feature in the (e.g. CDSs) correct orientation

Choose this option to write the sequence of any selected region of the genome

We are now ready to write out the alignment of the NV MLST genes. To do this, make sure the correct genes are selected in the embl file and right click on the bcfview panel. From the menu choose '**Write**' -> '**Fasta of selected**

features’. You will be asked for a filename for the alignment file. Call it NV_MLST.fasta. Note the three options on the right hand side of the save dialogue box.



- The ‘**Combine feature sequences**’ option is useful if you have selected more than one feature in the reference when choosing to write the bcf sequences. With this option selected, the sequences of each feature will be concatenated together in one file. If you deselect this option you will save one fasta per feature selected. We need to make sure this option is selected.
- The ‘**Use N for filtered out sites**’ tells Artemis that when a site fails the chosen filters, that base should be written as an N (unknown) in the alignment. If you deselect this box any site that fails any of the filters will be saved in the alignment as the reference base. Why might it be a bad idea?
- The ‘**Use N for sites without non-variant**’ option is useful when there are non-variant sites confirming the reference sequence as this will then write out ‘N’ for each of the non-confirmed sites. If you have multiple bcf files open, a fourth option will appear.
- The ‘**Single fasta**’ option tells Artemis to save the sequences of multiple bcf files into a single fasta file. If you deselect this option, one fasta file will be saved for each individual bcf. This does not apply here.

We have only output the MLST gene sequences for one Chlamydia isolate. If we had opened more than one bcf file in Artemis, we could have output the sequences for all of those isolates in one go to create an alignment of MLST gene sequences. However, running Artemis with many bcf files open can be slow, especially on the USB stick. An alternative is to output the sequence for each strain into a separate file in exactly the same way as we have for NV. You can then concatenate the separate files into one alignment file using the ‘cat’

command. Remember, though that the reference will be included in each file created in Artemis.

For speed we have provided you with an alignment of the MLST genes from the 15 other *C. trachomatis* isolates you included in the *ompA* tree. These were produced in exactly the same way as the NV_MLST.fasta file you just created. We will use cat to add the reference and NV sequences and make a file containing all 17 isolates by typing:

```
cat Others_MLST.fasta NV_MLST.fasta > Ct_MLST.fasta
```

```
# once completed, open "Ct_MLST.fasta" in Seaview
```



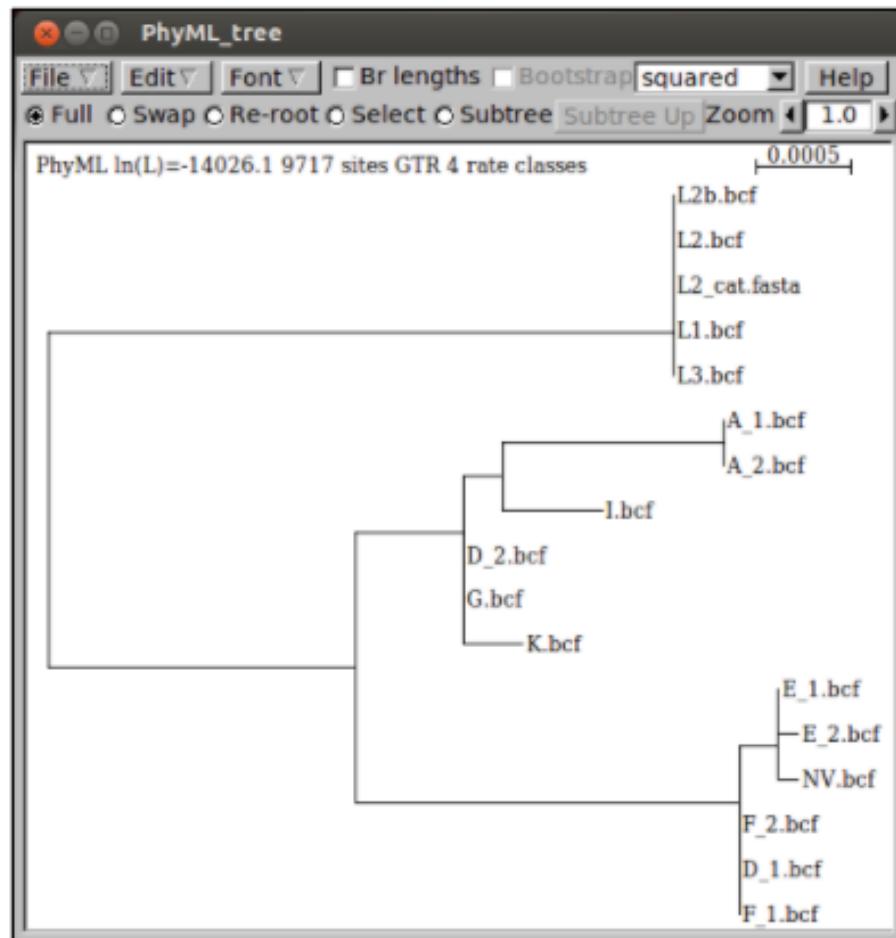
Notice that there is only one sequence for each strain. This is because Artemis has concatenated the sequences of the seven loci into one long sequence.

To check all seven genes are in the alignment you can use Seaview to translate the nucleotides into their corresponding amino acids. To do this, select ‘Props’ and then ‘View as proteins’.



To check you have seven genes in your alignment you can count the number of stop codons, which are represented as asterisks (*) in the protein view. To do this you can either scroll along the sequence or use the search box.

Construct a maximum likelihood tree of the MLST alignment as nucleotides.



Questions: - How does the tree compare with the *ompA* tree you made earlier? - Does the MLST gene tree support the splitting of *C. trachomatis* into trachoma and LGV biovars? - Do the strains cluster by serotype in the MLST tree? - What biological processes could account for these discrepancies? - Do you think *ompA* is a good gene for typing *C. trachomatis* isolates?

Back
to
top

Ex-
er-
cise
3:
Phy-
logeny
from
whole
genome
SNPs
In
the
past
few
years,
with
whole
genome
data
be-
ing
pro-
duced
for
large
num-
bers
of
bac-
te-
rial
iso-
lates,
it
has
be-
come
pos-
sible
to
use
vari-
a-
tion
from
whole
genomes
for
phy-
loge-
netic
re-
con-
struc-

—
Although
it is
pos-
sible
to
ex-
tract
whole
genome
vari-
a-
tion
from
Artemis
in
the
same
way
as
we
did
for
the
MLST
genes,
this
is
not
what
Artemis
was
de-
signed
to
do,
mak-
ing
the
pro-
cess
slow.
How-
ever,
there
are
ways
we
⁵⁸
can
ex-
tract
SNPs
from
bcf
vari-
a-

```
“‘bash
#
At
the
prompt
type
the
fol-
low-
ing
com-
mand:
sam-
tools
faidx
L2_cat.fasta
‘AM884176.1’
|
bcftools
con-
sen-
sus
../Mod-
ule_3_Mapping/NV.bcf
>
NV_WGS.fasta
```

Note:
you
are
us-
ing
the
vari-
ants
con-
tained
in
the
bcf
that
you
gen-
er-
ated
in
Mod-
ule_4
““

—
Although
this
com-
mand
line
looks
com-
plex,
it is
just
things
you've
seen
be-
fore
put
to-
gether.
The
com-
mand
is in
two
parts
sepa-
rated
by a
pipe
'|',
which
sim-
ply
tells
the
com-
mand
line
to
take
the
out-
put
from
the
first
com-
mand
afid
use
it as
the
in-
put
for
the

The
first
com-
mand
uses
sam-
tools
faidx
to
ex-
tract
a
sin-
gle
ref-
er-
ence
se-
quence
from
a
multi-
sequence
fasta
file.
In
this
case,
we
only
want
the
genome
se-
quence,
and
not
the
plas-
mid
se-
quence
that
is in
this
file.

The
sec-
ond
part
of
the
com-
mand
uses
bcftools
con-
sen-
sus
to
take
the
SNP
calls
you
gen-
er-
ated
yes-
ter-
day
in
Mod-
ule 3
and
to-
gether
with
the
ref-
er-
ence
se-
qunece
(ob-
tained
from
faidx),
to
make
a
new
fasta
se⁶³
quence
that
now
con-
tains
your
SNP

Finally
the
re-
sult
is
redi-
cted
into
a
file
called
NV_WGS.fastq.

The
se-
quence
pro-
duced
is
not
the
com-
plete
genome
of
the
NV
iso-
late,
as
map-
ping
only
al-
lows
vari-
ants
to
be
called
in re-
gions
present
in
the
ref-
er-
ence
genome,
as
where
the
amount
of
vari-
a-
tion
is
not
too
great.
In-
stead,
the
se-
quence
pro-
duced

MAKE
SURE
YOU
DO
THIS
STEP!!!!
Using
cat,
con-
cate-
nate
the
Oth-
ers_WGS.fasta
and
NV_WGS.fasta
files
into
a
sin-
gle
align-
ment
called
Ct_WGS.fasta

Try
to
open
the
align-
ment
in
Seav-
iew.
You
will
find
it is
much
larger
than
the
other
datasets
you
have
used.



The
Ns
in
the
align-
ment
rep-
re-
sent
bases
that
can-
not
be
called
from
the
map-
ping.
This
may
be
be-
cause
there
is no
map-
ping
in a
re-
gion
due
to a
true
dele-
tion,
or
be-
cause
the
map-
ping
of
that
base
fails
to
meet
one
of
the
fil-
ters
im-
posed.

Although it is possible to run a tree on such a large dataset because phylogenetic methods reduce the complexity by only analysing identical site patterns once, this may still be very slow. Instead, we will extract only those sites which contain vari-a-

```
“‘bash
#
At
the
com-
mand
line,
run
this
com-
mand:
snp-
sites
-o
Ct_WGS_SNPs.fasta
Ct_WGS.fasta
““
```

Open
the
SNP
align-
ment
in
Seav-
iew
and
make
a
tree
as
be-
fore.
Do
not
in-
clude
the
in-
vari-
ant
sites
pa-
ram-
eter,
as
this
would
not
makes
sense
– we
have
just
re-
moved
all
in-
vari-
ant
sites
from
the
dataset.

Questions:

- How
does
the
tree
of
whole
genome
SNPs
com-
pare
to
the
other
trees
you
have
made?

- What
could
cause
anal-
yses
of
dif-
fer-
ent
parts
of
the
genome
to
pro-
duce
dif-
fer-
ent
phy-
loge-
nies?

- What
does
this
tell
y72
about
ompA,
MLST
and
whole
genome
SNPs

[Back to top](#)

Additional Information

Other useful Seaview features

- File Menu:
 - **Save as...** allows you to save your alignment in many formats, including fasta, phylip (used by many phylogenetics programs), nexus (used by PAUP* and MrBayes), and many more.
 - **Save selection.** Allows you to save a region of your alignment masked with a set or an alignment of only the selected taxa.
 - **Concatenate.** Allows you to join multiple alignments together either based on the names of the taxa or the order of the taxa in the alignment
- Edit Menu:
 - In this menu there are options to delete, add, edit, reverse and reverse complement sequences in your alignment
 - You can view a dot plot of any selected pair of sequences
- Props Menu: - This menu contains options about how the sequences are shown. If you select the Allow seq. edition option, you can manually add or delete bases from the sequences
- Sites Menu:
 - This menu allows you to create a mask (set) under the alignment, which you can use to select a set of sites.
 - When sites (also taxa) are selected, any tree run will only include the selected sites and taxa.
 - You can save your selection using the save selection option in the file menu
- Search and Goto:
 - Search allows you to find a sequence of nucleotides or amino acids in a sequence.
 - Goto allows you to specify a base or amino acid number to move to in the alignment

FigTree: a more versatile tree viewer

Also included on your disk is a tree viewing program called **FigTree**. FigTree can open Newick format trees, as output by Seaview (see page 8). FigTree is more versatile than the tree viewer in Seaview, allowing you to colour branches and taxa, redraw the tree in a number of ways, collapse branches and output the results in a large number of graphics formats including eps and pdf. It is

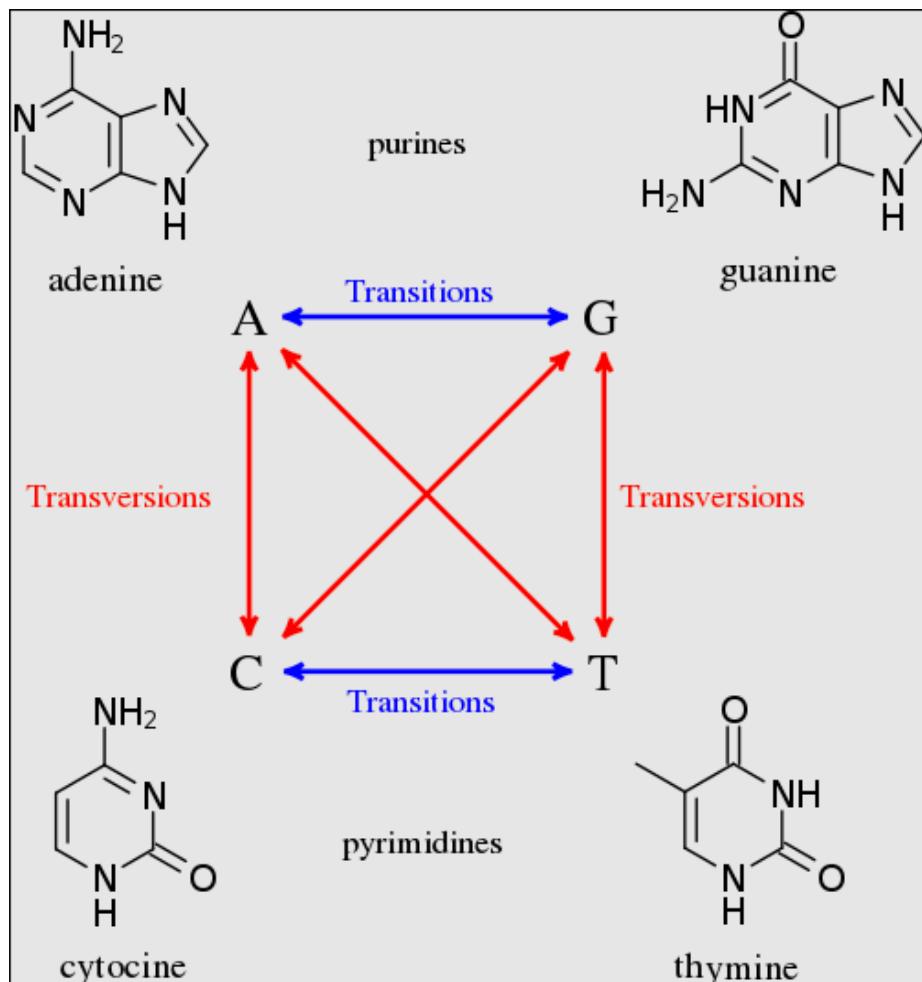
particularly useful for preparing figures for manuscripts. If you have time you may find it useful to try opening and editing your tree in FigTree.

FigTree can be opened by typing ‘`figtree`’ into the terminal.

Protein models

In this exercise you have become familiar with phylogenetic analysis of nucleotide data. Analysing sequences on the protein level is very similar, but there are some differences in the substitution models used.

When analysing nucleotide sequences our models optimise substitution rates from the data. For a JC69 model this means optimising a single rate for all changes, while for a GTR 6 substitution rates need to be optimised. The 6 nucleotide substitutions:



When analysing amino acid data we have 20 “standard” amino acids, which would require 190 substitution rates to be estimated. In most circumstances our data does not contain enough information to estimate all of these parameters, and even if it did, the calculation time would become prohibitive.

Instead we usually use a pre-made matrix of substitution rates calculated from large collections of alignments of proteins of known function. PhyML provides a range of options:

Model	Alignment source
LG	Nuclear globular proteins
WAG	
Dayhoff	
Blosum62	
MtREV	Vertebrate mitochondrial proteins
RtREV	Viral reverse transcriptase proteins
CpREV	Chloroplast proteins
DCMut	Extensnts to Dayhoff's PAM matrix
VT	Nuclear globular proteins
MtMam	Mammalian mitochondrial proteins
MtART	Arthropod mitochondrial proteins
HIVw	HIV-1 viral genes
HIVb	

The most commonly used model over the last few years has been the WAG, as it is generally applicable and in most cases gave the best results. However, the new LG model seems to be a further improvement. The gamma correction for among site rate variation, and invariant sites correction is exactly the same as with nucleotide data

Model selection

One of the most difficult and important decisions in phylogenetic analyses is which model to choose. Overly simple models are not biologically realistic, and have been shown to produce wrong trees in certain circumstances. For example, it is well known that parsimony and overly simple ML models often place long, unrelated branches together because they underestimate the number of multiple

substitutions at sites. This phenomenon, known as long branch attraction, is one of the major causes of error in phylogeny.

Given the problems associated with overly simple models, it is tempting to always use the most complex models available, which estimate more parameters from the data. Often you will find these are the best models for you to use.

AN INCREASE IN THE NUMBER OF MODEL PARAMETERS CAN ONLY INCREASE THE FIT OF THE MODEL TO THE DATA AND THEREFORE CAN ONLY IMPROVE THE LIKELIHOOD.

However, increasing the number of parameters reduces the amount of data we have to estimate the correct value for these parameters.

One way to approach choosing a model is to run a mini analysis with a range of models and compare the results. There are a number of tests that allow us to compare the likelihoods of models with different numbers of parameters. These include the likelihood ratio test, which you have used in this module, and the AIC.

You can run these tests using the programs jMODELTEST for nucleotides (a reduced version called Findmodel is available online) and PROTTEST for proteins:

- jMODELTEST: <http://darwin.uvigo.es/software/jMODELTEST.html>
- Findmodel: <http://www.hiv.lanl.gov/content/sequence/findmodel/findmodel.html>
- Prottest: <http://darwin.uvigo.es/software/prottest.html>

A final warning note: never select a model with one of these methods without thinking about the biology. If you're working on Plasmodium, would it be sensible to use a protein model produced from HIV sequences?

Alignment: some things to remember

- It is worth spending time to make a good alignment. If your alignment is wrong, your phylogeny is also likely to be wrong
- Progressive alignment is a mathematical process that is completely independent of biological reality
- Can be a very good estimate
- Can be an impossibly poor estimate
- Requires user input and skill
- Treat cautiously
- Can (usually) be improved by eye
- Often helps to have colour-coding
- Depending on the use, you should be able to make a judgement on those regions that are reliable or not
- For phylogeny reconstruction, only use those positions whose hypothesis of positional homology is strong

FINALLY... IT IS NOT USUALLY SENSIBLE TO ALIGN PROTEINS AT THE NUCLEOTIDE LEVEL: - The result might be highly-imausible and might not reflect what is known about biological processes. - It is much more sensible to translate the sequences to their corresponding amino acid sequences, align these protein sequences and then put the gaps in the DNA sequences according to where they are found in the amino acid alignment.

[Back to top](#)

License

This work is licensed under a Creative Commons Attribution 4.0 International License.