# EVOLUTION AND POPULATION DIVERSITY STUDY OF INFLUENZA D VIRUS (IDV) USING BEAST ALGORITHM

A REPORT SUBMITTED TO

SAVITRIBAI PHULE PUNE UNIVERSITY

TOWARDS PARTIAL FULFILLMENT OF DEGREE

OF

MASTER OF SCIENCE (M. SC.)

IN STATISTICS

IN THE FACULTY OF SCIENCE AND TECHNOLOGY

SUBMITTED BY

Dangat Ashitosh Vilas

UNDER THE GUIDANCE OF

Dr. Mohan Kale

DEPARTMENT OF STATISTICS

AND

CENTRE FOR ADVANCED STUDIES IN

STATISTICS,

SAVITRIBAI PHULE PUNE UNIVERSITY,

PUNE-411007, INDIA. MAY, 2023

# Certificate of the Guide

This is to certify that, the following students of M.Sc. Statistics,

(1) Dangat Ashitosh Vilas

have successfully completed their project titled **Evolution and Population Diversity Study of Influenza D Virus (IDV) using BEAST Algorithm** under the guidance of **Dr. Mohan Kale** and have submitted this project report on **May,19,2023** as a part of the course ST-402, towards partial fulfillment of requirements for the degree of M.Sc. Statistics in Savitribai Phule Pune University in the academic year 2022-2023.

**Dr. Mohan Kale**
**(Project Guide)**

**Prof. T. V. Ramanathan**
**(Head of the Department)**

# Acknowledgments

We would like to express our deep and sincere gratitude to **Dr. Mohan Kale** for his valuable guidance and motivation for the successful completion and execution of the project. His profound knowledge and expertise came in handy for our project.

It gives us immense pleasure to thank **Mr. Sanket Limaye**, Bioinformatics Centre, Savitribai Phule Pune University, for providing timely inputs and suggestions required for the project.

We would also like to thank **Prof. T.V. Ramanathan, Head, Department of Statistics,** and the teaching and non-teaching staff for their invaluable assistance.

# Contents

# List of Tables

x

# List of Figures

# Chapter 1

# Introduction and Summary

## 1.1 Introduction

Influenza D Virus (IDV) was initially isolated from pigs and subsequently from bovine (cattle). Since the seroprevalence in bovine is much higher than in swine, bovine is considered to be the primary natural reservoir for IDV. The pathology and mode of transmission of IDV were investigated by experimental infection of cattle and swine, but also of guinea pigs, ferrets, camels, and horses. The virus was also transmitted to naive animals by direct contact. After experimental infection with IDV in cattle, the virus can be detected both in the upper and lower respiratory tracts and transmitted to contact animals. Previous studies have reported that IDV detection rates are higher in clinically sick cattle than in apparently healthy cattle. The recent continuous outbreaks of IDV in bovines and other hosts imply that IDV can transmit among hosts with high efficiency. The intercontinental transmission and high prevalence of IDV, especially in cattle highlight its potential threat to other agricultural animals and humans.

Bayesian analysis has been conducted to address the origin and evolutionary history of IDV. Bayesian evolutionary analysis by sampling trees (BEAST)

combines statistical techniques with molecular sequence data to reconstruct phylogenetic trees and estimate evolutionary parameters, allowing researchers to explore the genetic relationships, divergence times, and population dynamics of pathogens.

Evolutionary rate study revealed that IDV evolves with a high substitution rate. Hence, continuous monitoring of the evolution of IDV needs to be conducted in the future. As with other types of influenza viruses, IDVs will constantly evolve into more diverse lineages. The IDV vaccine research efforts in the future need to accommodate for such antigenic drift phenomena and focus on developing a universal vaccine that can effectively protect animals or humans from multiple strains or lineages of IDV. This study aims to enlighten the intricate evolutionary processes and population diversity of the IDV using the BEAST algorithm.

## 1.2   Motivation

IDV primarily affects cattle and other ruminants, leading to respiratory illnesses that can have significant implications for animal health, welfare and productivity. The cattle industry is extremely important as it stimulates the self-sufficiency of the rural economy by providing auxiliary income to farmers and in turn, initiates the economic growth of a country and consequently contributes to the Gross Domestic Product.

IDV represents a relatively new and understudied member of the influenza virus family. As an emerging pathogen, it poses a potential threat to animal and potentially to human health. Understanding the virus's biology, transmission patterns, and evolutionary dynamics can aid in the development of targeted interventions, control measures, and prevention strategies to safeguard animal populations. The bovine respiratory disease complex costs the global cattle industry billions of dollars per year. As observed earlier, IDV has emerged an important role in this disease complex, thereby indicating a need for an effective and safe vaccine against IDV infection in the bovines. Currently, there are no specific vaccines or treatments for IDV in animals.

The interspecies transmission and the international appearance of IDV in worldwide animal populations represent a potential risk to global human health which clearly motivates further investigation of this novel influenza virus.

## 1.3   Objectives

- To estimate the rates of substitution (per site, per year) evolution.

- The study the time of evolution for the most recent common ancestor. IDV continuously evolves into more diverse lineages, its vaccine research efforts in the future need to accommodate for such an antigenic drift phenomenon and focus on developing a universal vaccine that can effectively protect animals or humans from multiple strains or lineages of IDV.

- To model the substitution rates for a nucleotide.

- To determine the maximum parsimony rate.

- To estimate the linkage disequilibrium.

- To conduct molecular phylogeny analysis.

- To perform biodiversity study.

## 1.4 About virus and disease:

The influenza viruses are members of the family Orthomyxoviridae. The influenza viruses A, B, C and D represent the four antigenic types of influenza viruses. In there four antigenic types, the influenza A virus is the most severe, the influenza B virus is less severe but can still cause outbreaks and the influenza C virus is usually only associated with minor symptoms.

IDV primarily spreads among animals, particularly cattle, through direct contact with infected animals or exposure to respiratory secretions. The virus can also spread indirectly through contaminated surfaces, although the exact duration of survival on surfaces is not well studied for IDV.

The respiratory secretions of infected animals, such as nasal discharge, saliva or respiratory droplets expelled during coughing or sneezing contain the virus. A close contact between infected and susceptible animals, such as being housed together in the same area or through nose-to-nose contact, facilitates the spread of IDV.

**Taxonomy of IDV is as follows:**

Realm: Riboviria

Kingdom: Orthornavirae

Phylum: Negarnaviricota

Class: Insthoviricetes

Order: Articulavirales

Family: Orthomyxoviridae

Genus: Deltainfluenzavirus

Species: Influenza D virus

## 1.5 Data Source and Description

We collect secondary data from National Center for Biotechnology Information GenBank (Reference: `https://www.ncbi.nlm.nih.gov/genbank/`). It provides a large suite of online resources for the biological information and data. GenBank sequence database is an open access annotated of all publicly available nucleotide sequences. IDV has an identical genome structure. It consists of 7 gene segments namely PB2, PB1, P3, HEF, NP, P42 and NS. The NS1 and NS2 genes has been separated from the NS segment as the independent segments. The gene sequences are further extracted, curated and aligned by a Ph. D. student of the Bioinformatics Centre, Savitribai Phule Pune University and then used for analysis.

The sequences are as follows:

| SR NO. | SEQUENCE NAME | NO. OF SEQUENCES |
|--------|---------------|------------------|
| 1 | PB2 | 175 |
| 2 | PB1 | 170 |
| 3 | P3 | 189 |
| 4 | HEF | 188 |
| 5 | NP | 183 |
| 6 | P42 | 179 |
| 7 | NS | 147 |
| 8 | NS1 | 147 |
| 9 | NS2 | 147 |

Table 1.1: Number of Sequences

# 1.6 Terminologies and Definitions

**Allele:** An alternative form of a gene (one member of a pair).

**Biodiversity**: Biodiversity is the degree of variation of life forms in a given ecosystem.

**Genetic biodiversity:** Genetic biodiversity refers to the total number of genetic characteristics in the genetic makeup of a species.

**Genome:** Complete set of genetic information found within an individual organism.

**Genotype:** An individual's genotype for a gene is the set of alleles, which the gene contains. Such as the letter Bb (B-dominant genotype and b-recessive genotype).

**Serotype:** Refers to the distinct variations within a species of bacteria or viruses or among immune cells of different individuals.

**Linkage:** Two genetic loci are said to be in linkage if the alleles at these loci segregate together more often than would be expected by chance.

**Population admixture:** It occurs when individuals from two or more previously separated populations begin interbreeding.

**Recombination:** Recombination is the process by which two Deoxyribonucleic acid (DNA) molecules share genetic information or their ancestry, producing a new combination of alleles.

**Parsimony informative sites:** These are the Sites on a chromosome containing at least two different allelic states, each represented at least twice.

**Phylogenetic tree:** It is nothing but a branching diagram or "Tree" showing the inferred evolutionary relationships among various biological species or other entities based upon similarities and differences in their genotypic or phenotypic characteristics.

**Sequence alignment:** Sequence alignment is the procedure of comparing two or more sequences by searching for a series of individual characters that

are in the same order in the sequences.

**Virus:** Small infectious agent that replicates only inside the living cells of other organisms.

**Strain:** A strain is a genetic variant or subtype of a micro-organism (e.g. virus, bacterium, fungus). For example, a "flu strain" is a certain biological form of influenza or "flu" virus.

**95% HPD Lower:** The lower bound of the highest posterior density (HPD) interval. The HPD is a credible set that contains 95% of the sampled values.

**95% HPD Upper:** The upper bound of the highest posterior density (HPD) interval. The HPD is a credible set that contains 95% of the sampled values.

**Auto-Correlation time (ACT):** The number of states in the Markov chain Monte Carlo (MCMC) chain that two samples have to be from each other for them to be uncorrelated. The ACT is estimated from the samples in the trace (excluding the burn-in).

**Effective sample size (ESS):** The ESS is the number of independent samples that the trace is equivalent to. This is essentially the chain length (excluding the burn-in) divided by the ACT.

**Maximum clade credibility (MCC) tree:** It finds the tree with the highest product of the posterior probability of all its nodes.

**Mean substitution rate per site per year:** The mean substitution rate per site per year refers to the average rate at which genetic substitutions, such as nucleotide changes, accumulate in a particular region of a genome over a specific period of time. It is commonly used in molecular clock analyses and provides an estimate of the evolutionary rate of a sequence.

**Maximum parsimony:** Maximum parsimony is a principle and method used in phylogenetic analysis to infer the most likely evolutionary tree or phylogeny from a given set of data, typically molecular sequence data such as DNA or protein sequences. The principle of maximum parsimony aims to identify the tree that requires the fewest evolutionary changes or genetic mutations to explain the observed data.

**Linkage disequilibrium (LD):** LD is a term used in genetics to describe the non-random association or correlation between alleles at different loci (positions) on a chromosome. In other words, it refers to the tendency of certain alleles at different genetic loci to occur together more often than would be expected by chance.

## 1.7    Software Used

1. **IQTREE**

   IQTREE is a fast and effective stochastic algorithm to infer phylogenetic trees by the method of maximum likelihood.

   Reference: `http://iqtree.cibiv.univie.ac.at/`

2. **TEMPEST v1.5.3**

   TEMPEST is a graphical program for looking for temporal signals in trees generated from tip-dated sequences.

   Reference: `http://tree.bio.ed.ac.uk/software/tempest/`

3. **BEAUti v1.10.4**

   BEAUti stands for Bayesian Evolutionary Analysis Utility. This program is used to import data, design the analysis, and generate the BEAST control file.

   Reference: `https://beast.community/beauti`

4. **BEAST v1.10.4**

   BEAST stands for Bayesian Evolutionary Analysis Sampling Trees. This main program takes a control file generated by BEAUti and performs the analysis.

   Reference: `https://beast.community/`

5. **LOGCOMBINER v1.10.4**

   Logcombiner will combine log files from different runs and reduce the sampling frequency. It allows you to combine log and tree files from multiple independent runs of BEAST.

   Reference: `https://beast.community/logcombiner`

6. **TRACER v1.7.2**

   Tracer is a graphical program for exploring the output of BEAST, diag-

nosing problems, and summarizing the results.

Reference: `http://tree.bio.ed.ac.uk/software/tracer/`

7. **TREE ANNOTATOR v1.10.4**

   Tree Annotator is a post-analysis program will produce a summary tree from the output of BEAST.

   Reference: `https://beast.community/treeannotator`

8. **FIG TREE v1.4.4**

   Fig Tree is a graphical program for viewing trees, displaying summary information from TreeAnnotator.

   Reference: `http://tree.bio.ed.ac.uk/software/figtree/`

## 1.8 BEAST Software:

BEAST is a fast, flexible software architecture for Bayesian analysis of molecular sequences related by an evolutionary tree. The BEAST software package is an ambitious attempt to provide a general framework for parameter estimation and hypothesis testing of evolutionary models from molecular sequence data.

As the BEAST is a Bayesian statistical framework, it provides a role for prior knowledge in combination with the information provided by the data. The purpose behind the development of BEAST is to bring a large number of complementary evolutionary models (substitution models, insertion-deletion models, demographic models, tree shape priors, relaxed clock models, and node calibration models) into a single coherent framework for evolutionary inference. A Bayesian Markov chain Monte Carlo (MCMC) method for performing relaxed phylogenetics that is able to co-estimate phylogeny and divergence times under a new class of relaxed-clock models.

MCMC is a stochastic algorithm that produces sample-based estimates of a target distribution of choice. A Bayesian MCMC algorithm needs to evaluate the likelihood of each state in the chain and thus performance is dictated by the speed at which these likelihood evaluations can be made. BEAST attempts to minimize the time taken to evaluate a state by only recalculating the likelihood for parts of the model that have changed from the previous state. BEAST uses the Java version of the functions, thereby retaining its platform independence.

# Chapter 2

# Statistical Methods

## 2.1 Markov Chain Monte Carlo (MCMC)

Notations: $X$: genotypes of sampled individuals

$Z$: population of origin of individuals (unknown)

$P$: allele frequencies in all populations (unknown)

$Q$: admixture proportion for each individual.

$P$, Q-Dirichlet distribution (used as a prior distribution in Bayesian statistics)

MCMC is useful for obtaining the distribution $\pi(\theta)$.

In this case $\theta = (Z, P, Q)$ *and* $\pi(0) = P(Z, P, Q|X)$

The idea is to construct Markov Chain $\theta(0), \theta(1), \theta(2).....$ with stationary distribution $\pi(\theta)$ the $\theta(m)$ is approximately distributed as $\pi(\theta)$ where $m$ is sufficiently large.

This can be formalized and shown to be true provided the Markov chain satisfies certain technical conditions such as ergodicity that hold here.

Furthermore, for sufficiently large $C, 0(m), 0(m+c), 0(m+2c)$ will reasonably independent samples from $(0)$. The values of $M$ and $C$ are referred to as $M$: Burn-in period of the Markov chain. $C$: Thinning interval the value of $m$ and $C$ heavily depends on the amount of correlation between two successive states of the Markov Chain.

## 2.2 Chi-square test:

At the beginning of each run, IQTREE performs a composition chi-square test for every sequence in the alignment. The purpose is to test for homogeneity of character composition (e.g., nucleotide for DNA, amino-acid for protein sequences). A sequence is denoted as failed if its character composition significantly deviates from the average composition of the alignment. More specifically, for each sequence, compute:

$$\chi^2 = \sum_{i=1}^{n} \frac{(O_i - E_i)^2}{E_i},$$

where $k$ is the size of the alphabet
(e.g. 4 for DNA, 20 for amino acids) and the values 1 to $k$ correspond uniquely to one of the characters. $O_i$ is the character frequency in the sequence tested. $E_i$ is the overall character frequency from the entire alignment. Whether the character composition deviates significantly from the overall composition is done by testing the $\chi^2$ value using the $\chi^2$ distribution with $k-1$ degrees of freedom (df=3 for DNA or df=19 for amino acids).
This test is regarded as an explorative tool that might help to nail down problems in a dataset. One would typically not remove failing sequences by default. But if the tree shows unexpected topology the test might point in the direction of the origin of the problem.

## 2.3 Maximum likelihood calculations based on AIC, BIC, AICc

IQTREE is a software package that implements maximum likelihood and Bayesian inference methods for phylogenetic analysis. It uses these methods to estimate the evolutionary relationships and parameters of genetic data. IQ-TREE also provides model selection tools, such as Akaike Information Criterion (AIC), Corrected Akaike Information Criterion (AICc) and Bayesian Information Criteria (BIC) to compare different models and select the most appropriate one for the data at hand. It scores each model on the basis of AIC, AICc and BIC. In BIC, a stronger penalty is applied to model complexity compared to AIC and AICc, making it more effective in selecting parsimonious models.

**1) AIC:** The AIC is a measure of the relative quality of a statistical model. Akaike proposed an information criterion, AIC, based on maximizing the expected entropy of the model. Entropy is a measure of the expected information, in this case, the Kullback - Leibler information measure. Essentially, the AIC is a penalized log-likelihood measure. Let $L$ be the likelihood function for a specific model. The AIC is,

$$AIC = -2lnL + 2k,$$

where, $k$ includes the parameters from the substitution model, such as base frequencies, substitution rates, the proportion of invariable sites, or rate variation among sites.

**2) BIC:** BIC aims to select the model that maximizes the likelihood of the data while also considering model complexity. It penalizes complex models more heavily than the AIC does. The BIC is,

$$BIC = -2lnL + k * ln(n),$$

where $k$ is the number of estimable parameters and $n$ is the sample size (for now we assume that $n$ can be approximated by the total number of characters in the alignment.)

**3) AICc:** The corrected AIC (AICc) is a modified version of the AIC that addresses the issue of small sample sizes in statistical model selection. It adjusts the AIC value to provide a more reliable estimate of model fit when the number of observations is relatively small compared to the number of parameters in the model. The AIC is given as,

$$AICc = AIC + \frac{2k(k+1)}{(n-k-1)},$$

where sample size is approximated by the total number of characters in the alignment.

**4)Weighted AIC (W-AIC):** W-AIC in the context of fitting hierarchical models using Markov Chain Monte Carlo (MCMC) methods. While traditional AIC is based on the maximized likelihood, W-AIC takes into account the full posterior distribution of model parameters.

**5)Weighted BIC (W-BIC):** The W-BIC is a variation of the BIC that incorporates additional weights to account for specific considerations in model selection. These weights are assigned to different components of the BIC formula to emphasize certain aspects of model evaluation. The standard BIC formula is given by:

$$BIC = -2(log - likelihood) + klog(n),$$

where log-likelihood is the logarithm of the likelihood function of the model, $k$ is the number of free parameters in the model and $n$ is the sample size.

**6)Weighted AICc (W-AICc):** W-AICc is a further extension of the AICc that incorporates weights to account for the possibility of model selection uncertainty. Instead of selecting a single best model based on the AICc, weighted AICc provides a weighted average of multiple models, taking into account their relative fit to the data and the complexity of each model. The weights reflect

the likelihood that each model is the best one among the set of candidate models.

## 2.4 Marginal likelihood calculations using path sampling and stepping stone sampling

Path sampling and stepping stone sampling are both methods used in Bayesian statistical analysis to estimate marginal likelihoods and to compute Bayes factors. These techniques are particularly relevant when comparing different models or hypothesis testing in the context of Bayesian model selection. Accurate estimates of marginal likelihoods can be obtained using stepping-stone sampling and path-sampling techniques.

The marginal likelihood is the probability of the observed data X under a given model $M_i$ that is averaged over all possible values of the parameters of the model $\theta_i$ with respect to the prior density on $\theta_i$

$$P(X|M_i) = \int P(X|\theta_i)P(\theta_i)dt$$

These methods require Markov chain Monte Carlo samples from a range of power posterior distributions, which represent the path between the prior and the posterior distributions.

- **Path sampling** Path sampling involves estimating the marginal likelihood for every power posterior sample. Marginal likelihood represents the probability of the observed data under a specific model, integrating over all possible values of the model's parameters. It serves as a measure of the goodness of fit of a model to the data and can be used for model

comparison.

- **Stepping-stone sampling** Stepping-stone sampling estimates the likelihood ratio for a series of power posterior samples. It is a variation of path sampling that employs multiple intermediate distributions between the prior and posterior. Each intermediate distribution is obtained by progressively weighting the likelihood component of the posterior distribution. The accuracy of these two methods is very similar, but stepping-stone sampling is more computationally efficient because it requires fewer samples from the power posterior.

Although both of these methods are more accurate than the harmonic-mean estimator, they come at the cost of a substantially increased computational burden.

# Chapter 3

# Data Analysis

## 3.1 Nucleotide Substitution Models

- Jukes-Cantor (JC) Model:

  This is a simple and widely used model that assumes all nucleotide substitutions occur at equal rates. It assumes a single parameter for the overall substitution rate and does not account for the differences in mutation rates between different types of nucleotide substitutions.

- Hasegawa-Kishino-Yano (HKY) Model:

  The HKY model is a variant of the General Time Reversible (GTR) model that assumes a stationary state for nucleotide frequencies but allows for different substitution rates for transitions and transversions. It is commonly used for phylogenetic analyses and incorporates more parameters than the JC or Kimura models.

- Tamura-Nei (TN93) Model:

  This model of substitution assumes that nucleotide frequencies may vary among different sequences or lineages. In addition to nucleotide frequencies, the TN93 model considers different rates for transitions (purine-to-purine or pyrimidine-to-pyrimidine substitutions) and transversions

(purine-to-pyrimidine or vice versa substitutions). It assumes that the transition rate is different from the transversion rate, acknowledging the known bias in mutation rates observed in nucleotide substitutions.

- General Time Reversible (GTR) Model:

  This model considers six different substitution rates, one for each possible pair of nucleotides.

  $$(A \longleftrightarrow C, A \longleftrightarrow G, A \longleftrightarrow T, C \longleftrightarrow G, C \longleftrightarrow T, G \longleftrightarrow T)$$

  allowing for a more detailed representation of the evolutionary processes. It also takes into account rate heterogeneity among different sites in a sequence.

- Transitional Interchangeability Matrix (TIM) Model:

  The TIM model assumes that substitutions occur at different rates depending on the type of nucleotide or amino acid being replaced and the surrounding context. The TIM model calculates the substitution probabilities between different nucleotides or amino acids based on observed patterns of substitution in a given dataset.

## 3.2   Molecular Clocks

1. Strict Clock:

   Constant rate among branches.



Figure 3.1: Strict Clock

2. Uncorrelated Relax Clock:

   Distinct rate along each branch drawn from a chosen probability distri-
   bution.

Figure 3.2: Uncorrelated Relax Clock

3. Fixed Local Clock:

   Branches having the same rate are clustered together.



Figure 3.3: Fixed Local Clock

4. Random Local Clock:

   Branches having the same rates are distributed throughout the tree rather than being locally clustered.

0.01

Figure 3.4: Random Local Clock

## 3.3 Work flow for analysis

### 3.3.1 IQTREE

IQTREE is a fast and effective stochastic algorithm to infer phylogenetic trees by maximum likelihood.

1) Choose the alignment file and then click on submit the job.

Figure 3.5: IQTREE

2) IQTREE scores each model on the basis of AIC, w-AIC, AICc, w-AICc, BIC and w-BIC and the value corresponding to the least BIC value will be considered as the best model.

**AIC, w-AIC:** Akaike information criterion scores and weights.

**AICc, w-AICc:** Corrected AIC scores and weights.

**BIC, w-BIC:** Bayesian information criterion scores and weights.

3) Select "DOWNLOAD SELECTED JOBS" at the lower left-hand corner.

All the files required for further analysis will be hence downloaded.



Figure 3.6: IQTREE

### 3.3.2 TEMPEST

1. When started, TempEst displays a file selection dialog box in which you can select the contree file that you downloaded in the previous section.
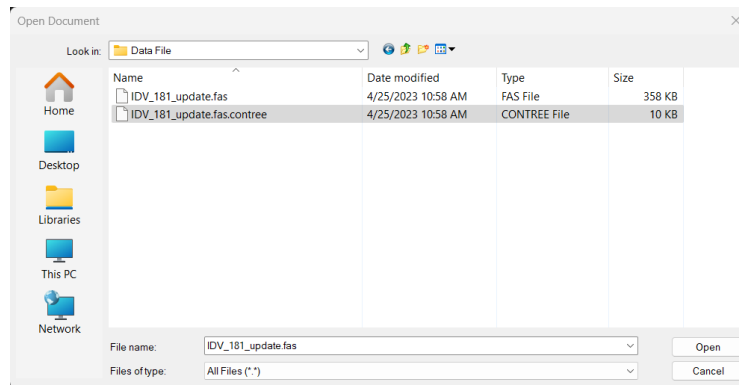


Figure 3.7: TEMPEST 1

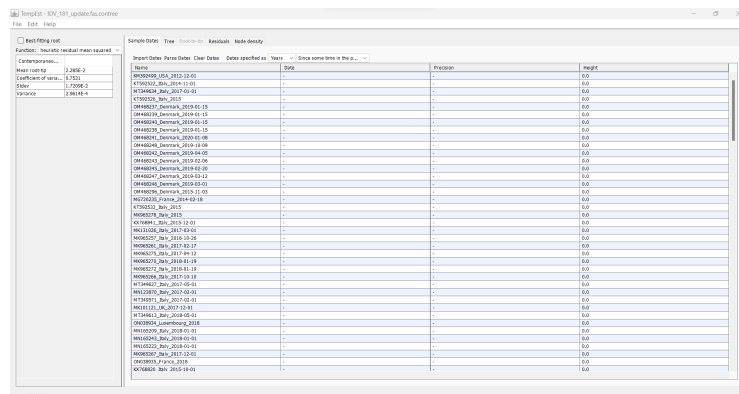2. Once the tree is loaded the main window appears:



Figure 3.8: TEMPEST 2

3. Parsing dates of sampling:

This dialog box provides a wide range of other options for extracting dates from taxon labels.

- Defined by a prefix and its order: This option finds fields that are prefixed by a particular character or string and then uses the order to specify which one. Here, the year that each virus was isolated is given at the end of the labels; hence, the order is last and separated.



Figure 3.9: TEMPEST 3

- The parsed dates are as follows:

The Height column lists the ages of the tips relative to time 0.



Figure 3.10: TEMPEST 4

4. The temporal signal and rooting

We can explore data using the tabs at the top of the window-Tree, Root-to-tip, and Residuals.

- Tree - Since the tree is constructed using a non-molecular clock model, it will be arbitrarily rooted.
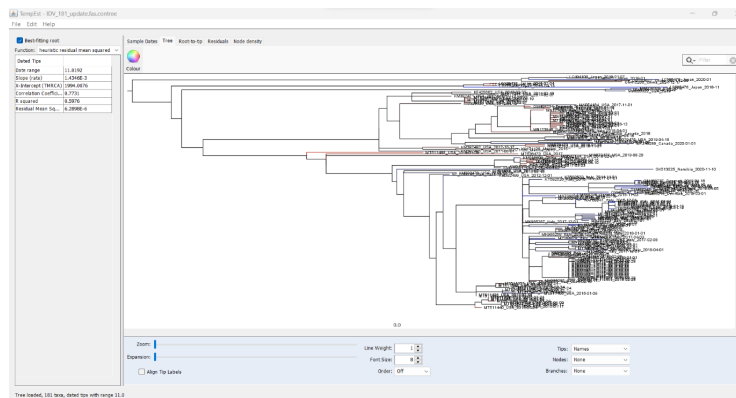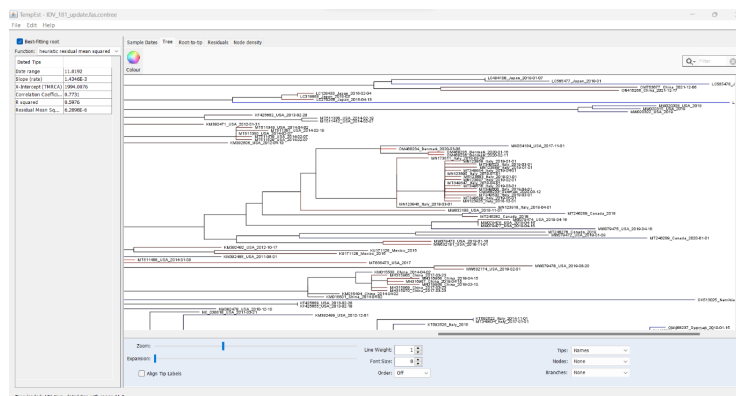


Figure 3.11: TEMPEST 5



Figure 3.12: TEMPEST 6

- Root-to-tip panel shows a plot of the divergence from the root of the tree against the time of sampling.
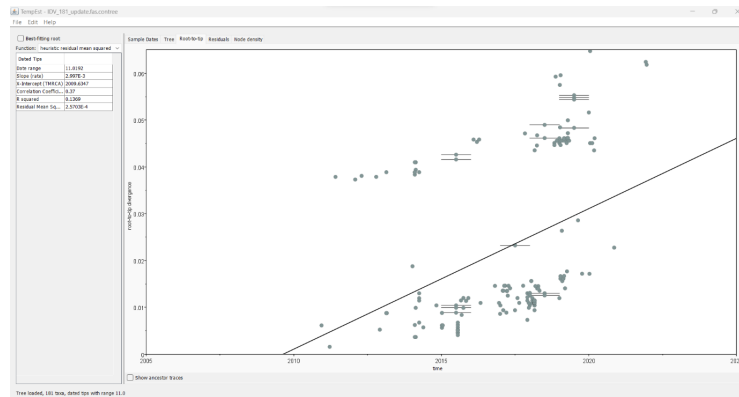
Figure 3.13: TEMPEST 7

- By clicking on the Best-fitting root, it minimizes the mean of the squares of the residuals.
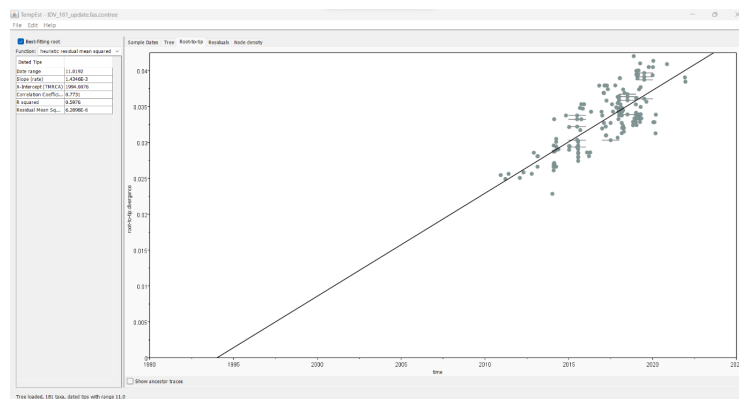


Figure 3.14: TEMPEST 8

- If one switches to the Residual panel you will see a plot of all the residuals.
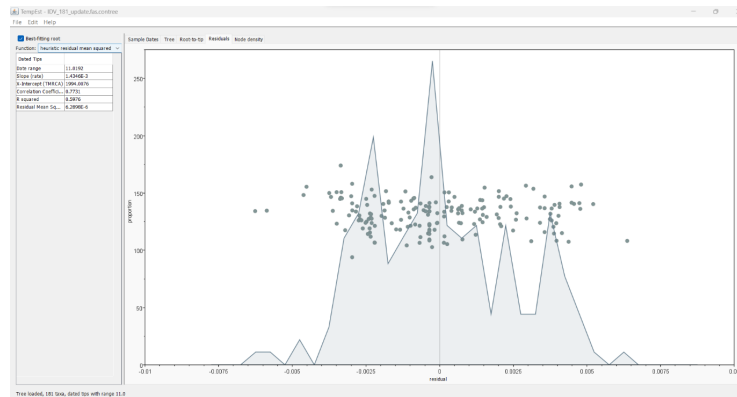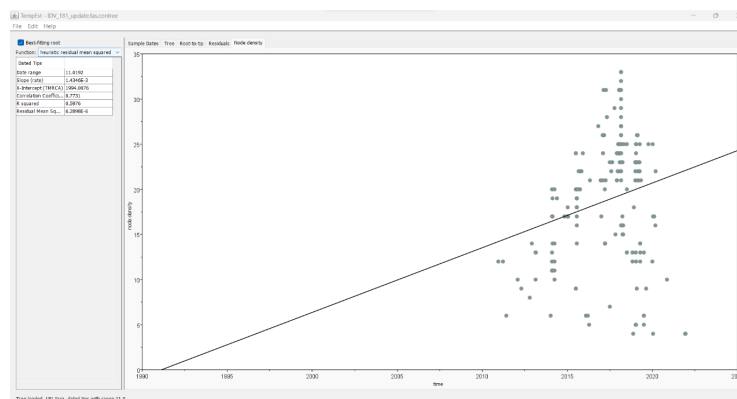
Figure 3.15: TEMPEST 9



Figure 3.16: TEMPEST 10

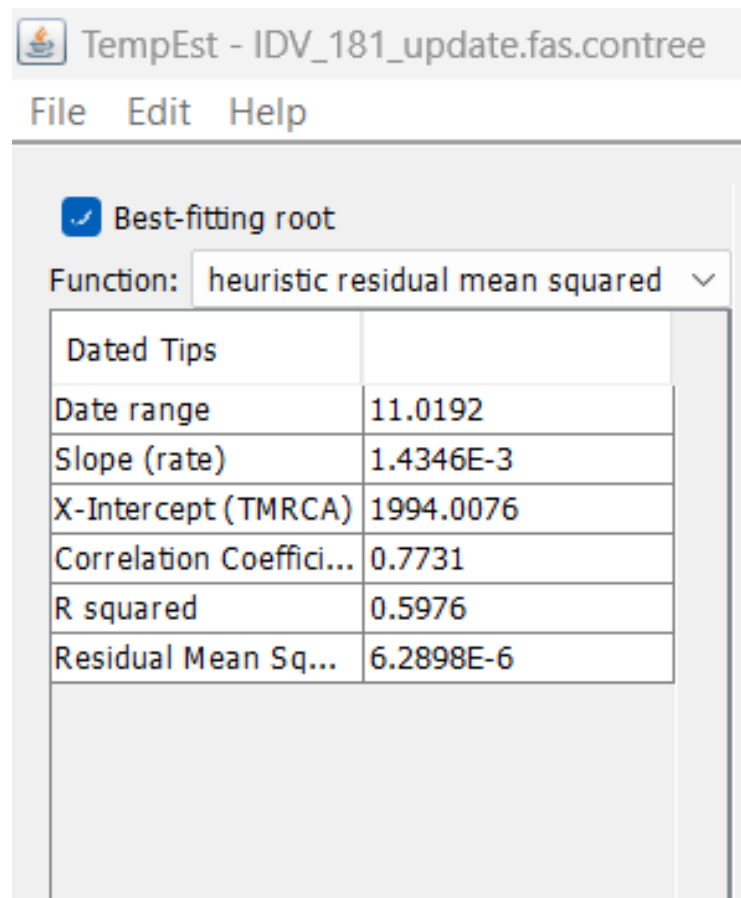- On the left-hand side of the window there is a table of statistics

Figure 3.17: TEMPEST 11

- Slope(rate) is an estimate of the rate of evolution in substitutions per site per year.

- $X$-Intercept is an estimate of the date of the root of the tree.

### 3.3.3   BEAUti

Following are the steps for generating a BEAST XML input file:

1. Loading the NEXUS file: To load a NEXUS format alignment, simply select the Import Alignment option from the File menu.
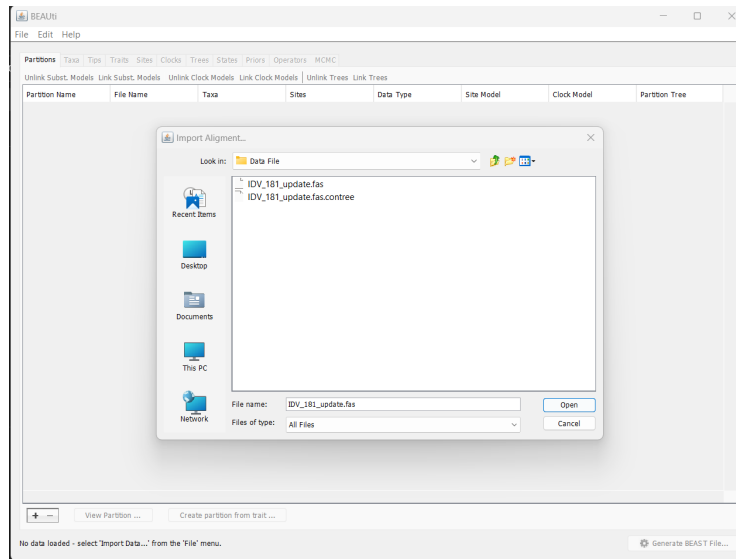
Figure 3.18: BEAUti 1

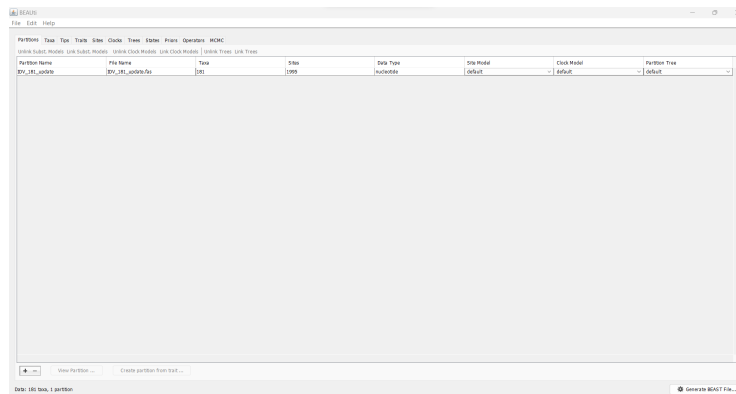Once loaded, the alignment will be displayed in the main window in a table:



Figure 3.19: BEAUti 2

2. Setting the dates of the taxa:

   In the Tips options, you will see a table with all of the taxa that were in alignment. This panel allows you to give the taxa dates.

   - Here, the year that each virus was isolated is given at the end of the labels and hence the order is last and is separated by.
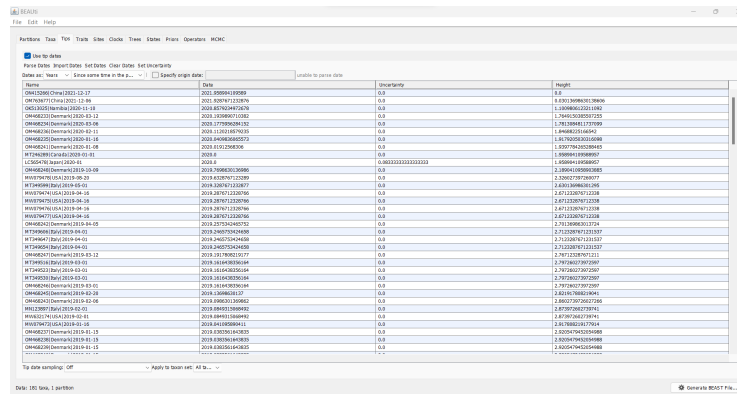
- Parse calendar date with variable precision: This parses a very specific way of specifying dates which are by far the most reliable. This format is the equivalent of using yyyy-MM-dd in the box above. This format has an additional advantage in that the day can be omitted if it is not known. Both the day and the month can be omitted if only the year is known.

  When using this option cases where the day or day and month are omitted are noted and the uncertainty is recorded in the precision column.



Figure 3.20: BEAUti 3

- The parsed dates are as follows:

  The Height column lists the ages of the tips relative to time 0.

Figure 3.21: BEAUti 4

3. Setting the evolutionary model:

In the Sites tab, you can set the model of molecular evolution for the sequence data you have loaded.

- Substitution models specify the evolutionary process that describes how the data was generated on the tree.
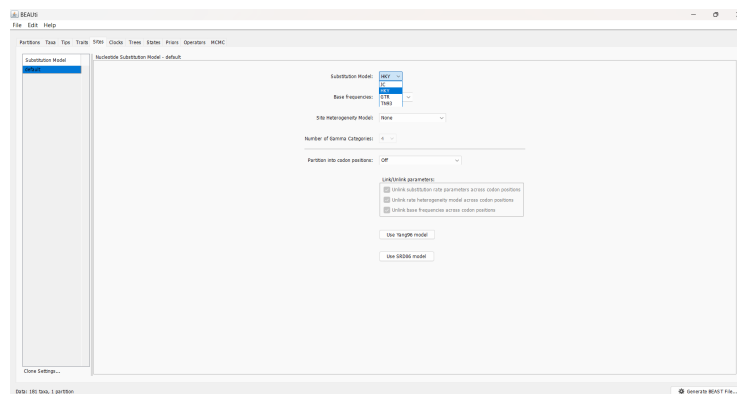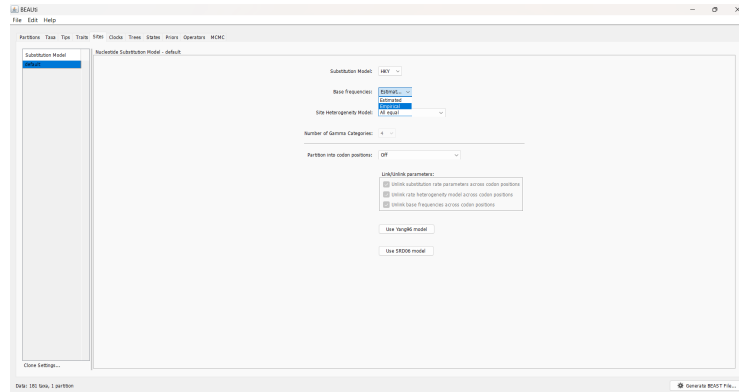


Figure 3.22:   BEAUti 5

- Base frequencies:



Figure 3.23: BEAUti 6

- Site Heterogeneity Model:

  Gamma distribution ($G$): gamma-distributed rate variation among sites

  The proportion of invariable sites ($I$): invariant sites across all taxa
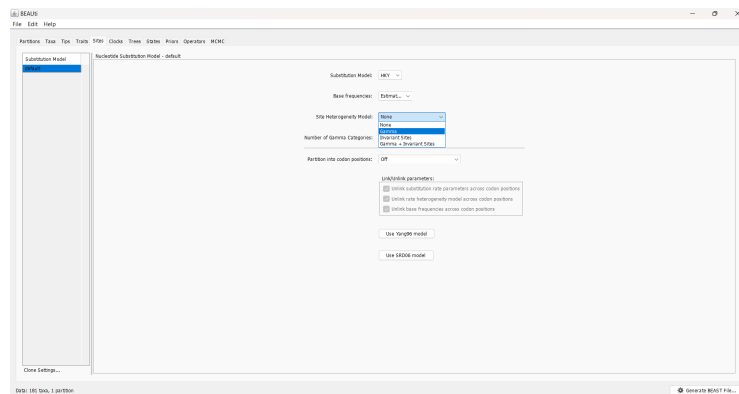


Figure 3.24: BEAUti 7

4. Setting the molecular clock model:

   In the next tab Clock we set the model of a molecular clock we will
   use. BEAST exclusively uses molecular clock models so that trees have
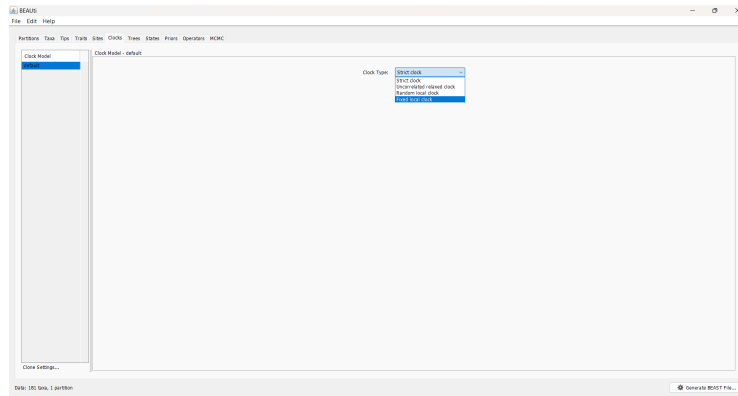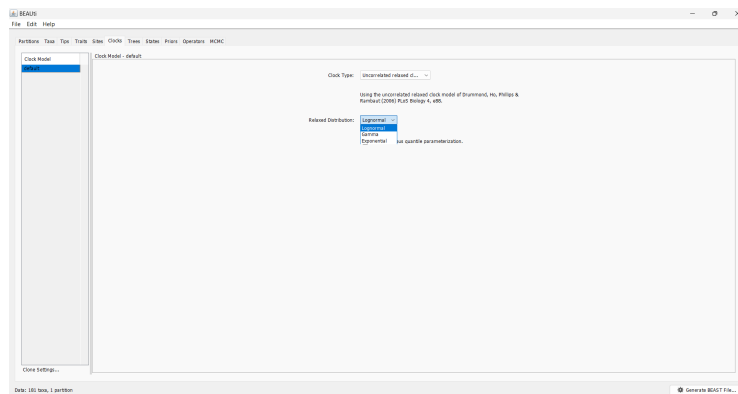   a timescale.



Figure 3.25: BEAUti 8



Figure 3.26: BEAUti 9

5. Setting the tree prior: In this panel, one can set the model that provides a prior on the tree and some choices about the starting tree in the MCMC run.
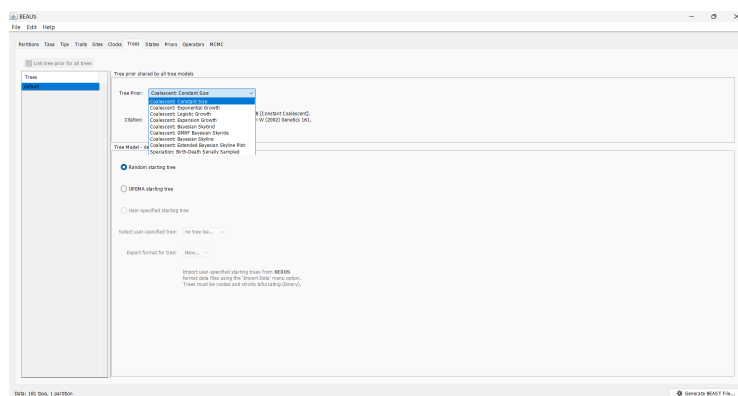


Figure 3.27: BEAUti 10

6. Setting MCMC options:

   - The next tab, MCMC, provides more general settings to control the length of the MCMC run and the file names. The Chain Length is the number of steps the MCMC will make in the chain before finishing and it depends on the size of the data set, the complexity of the model, and the quality of the answer required. The next options specify how often the parameter values in the Markov chain should be displayed on the screen and recorded in the log file
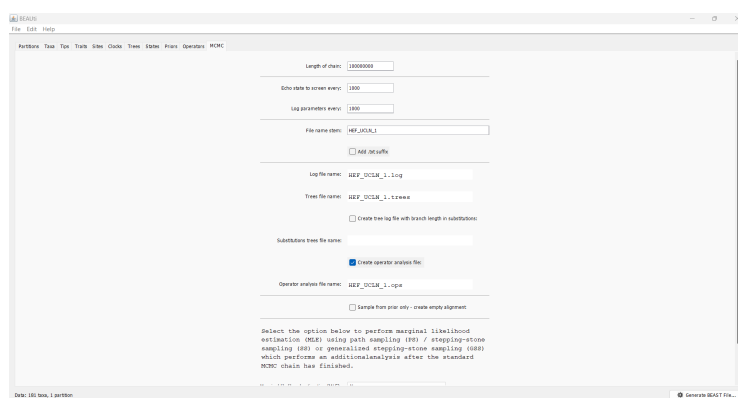


Figure 3.28: BEAUti 11

- This option allows to performance of marginal likelihood estimation using path sampling (PS) / stepping-stone (SS) or generalized stepping-stone (GSS) which performs an additional analysis after the standard MCMC has finished.
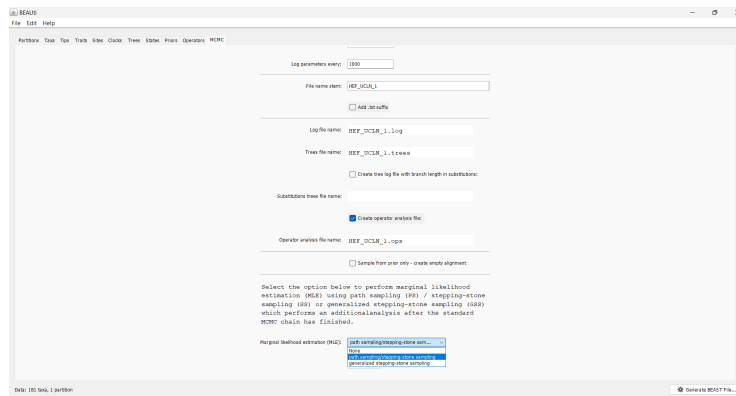


Figure 3.29: BEAUti 12

7. Generating the BEAST XML file: We are now ready to create the BEAST XML file. Select Generate XML from the File menu. It ends with '.xml'.
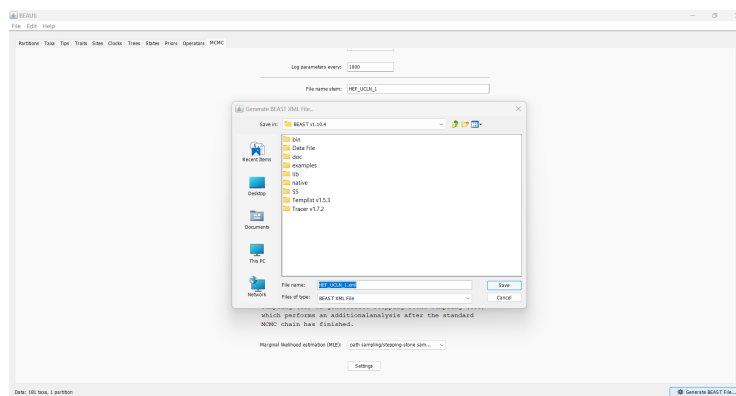


Figure 3.30: BEAUti 13

### 3.3.4 BEAST

- Select the XML file you created in BEAUti and press Run.



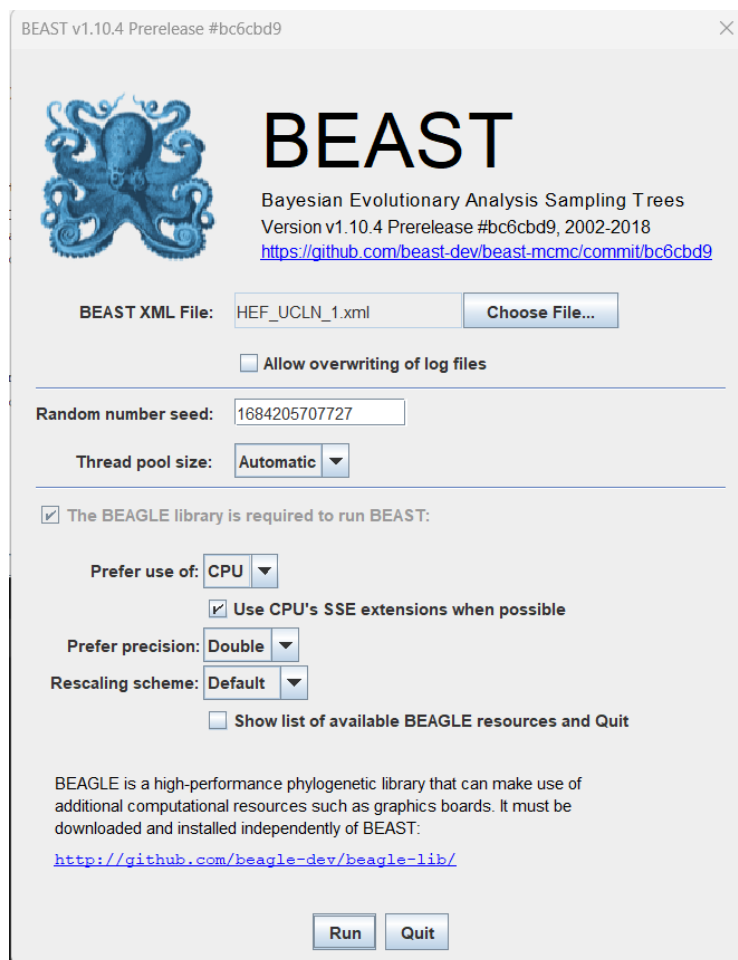Figure 3.31: BEAST 1

- All citations relevant for the analysis are mentioned at the start of the run.

```
                    BEAST v1.X, 2002-2102
            Bayesian Evolutionary Analysis Sampling Trees
                     Designed and developed by
         Alexei J. Drummond, Andrew Rambaut and Marc A. Suchard

                     Department of Computer Science
                        University of Auckland
                       alexei@cs.auckland.ac.nz

                    Institute of Evolutionary Biology
                        University of Edinburgh
                          a.rambaut@ed.ac.uk

                    David Geffen School of Medicine
                  University of California, Los Angeles
                          msuchard@ucla.edu

                      Downloads, Help & Resources:
                          http://beast.community

         Source code distributed under the GNU Lesser General Public License:
                     http://github.com/beast-dev/beast-mcmc
```

Figure 3.32: BEAST 2

- The first column is the 'state' number — in this case it is incrementing by 1000 so between each of these lines it has made 1000 operations. The screen log shows only a few of the metrics and parameters but it is also recording a log file to disk with all of the results in it (along with a '.trees' file containing the sampled trees for these states). After a few thousand states it will start to report the number of hours per million states.
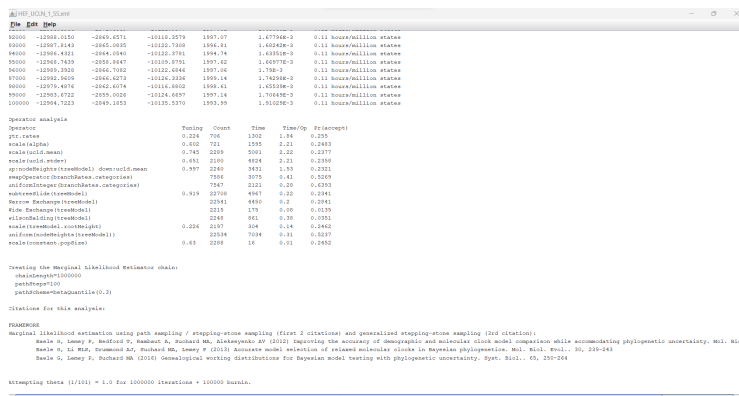


Figure 3.33: BEAST 3

Figure 3.34: BEAST 2

## 3.3.5 LOGCOMBINER

- Selecting log files from "file type" menu. The input files are those log files generated from the above step. Choosing an output file will save the combined log files into it. Then select "Run".
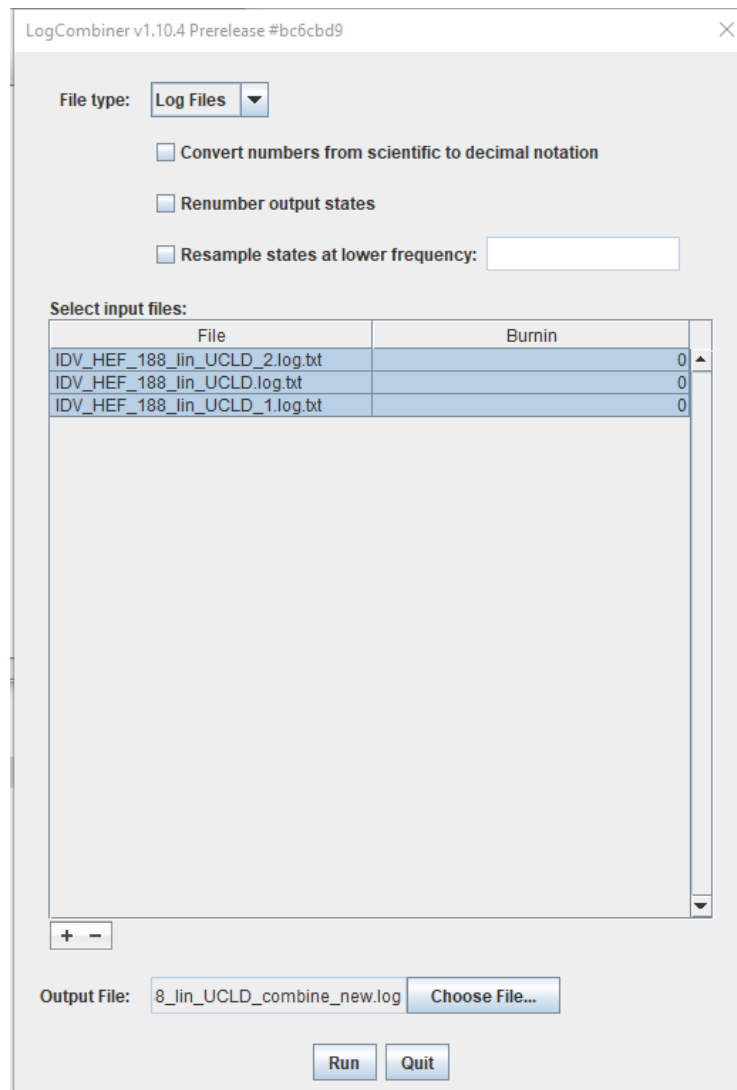
Figure 3.35: LOGCOMBINER 1

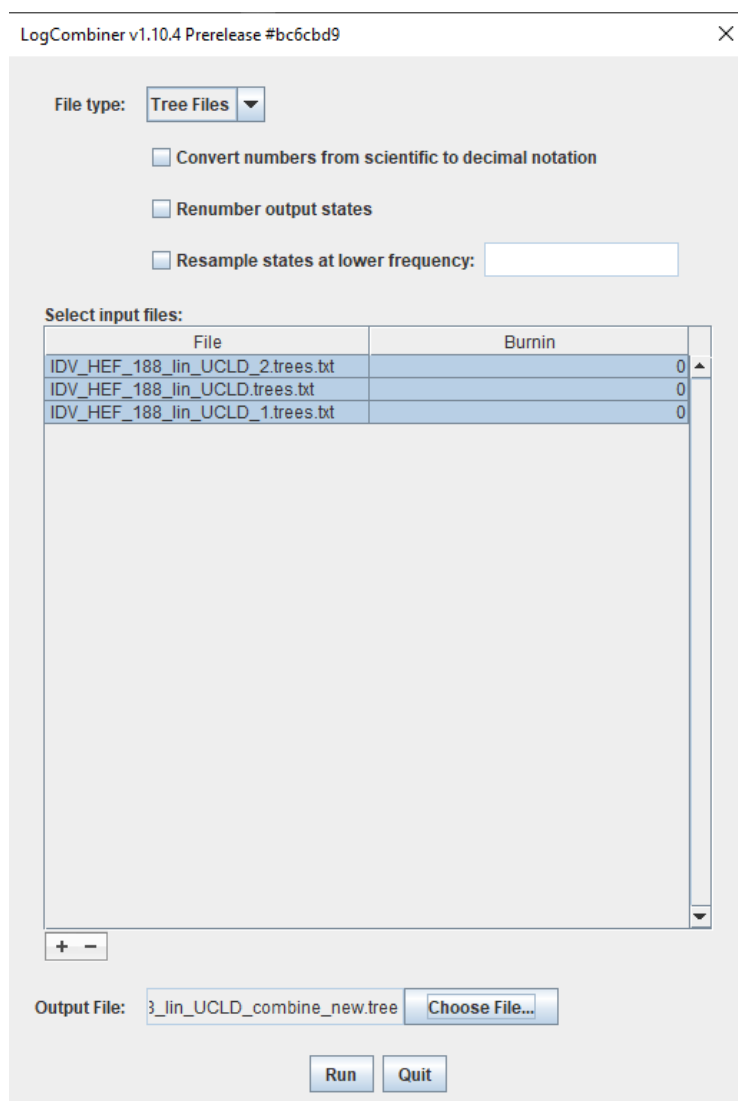- Similarly, for tree files.

Figure 3.36: LOGCOMBINER 2

- When the program says 'Finished', the combined log file you selected, above, will be ready. Close the program by selecting the Quit option. Once LogCombiner has finished you can analyze the combined log file in Tracer.
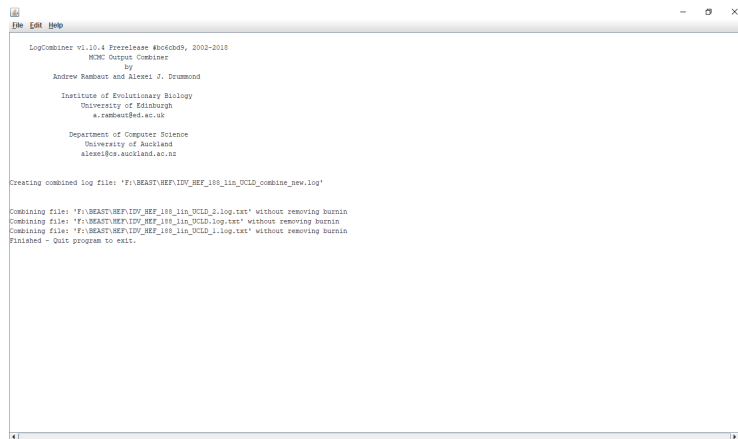
- Similarly, for tree files.
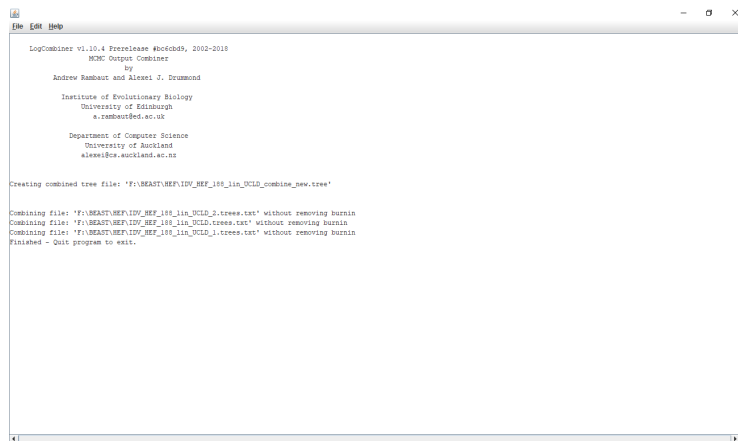
Figure 3.37: LOGCOMBINER 3



Figure 3.38: LOGCOMBINER 4

### 3.3.6 TRACER

- We start by loading the output .log files of both replicates of the same BEAST XML into Tracer 1.7. To load the log file(s), select the Open option from the File menu or drag and drop the log file into the Tracer window. The files will load and you will be presented with a window similar to the one below.

- As with loading a single log file, the name of the log file loaded and the traces that it contains can be seen on the left hand side. When the different files loaded contain the same set of logged parameters, then a
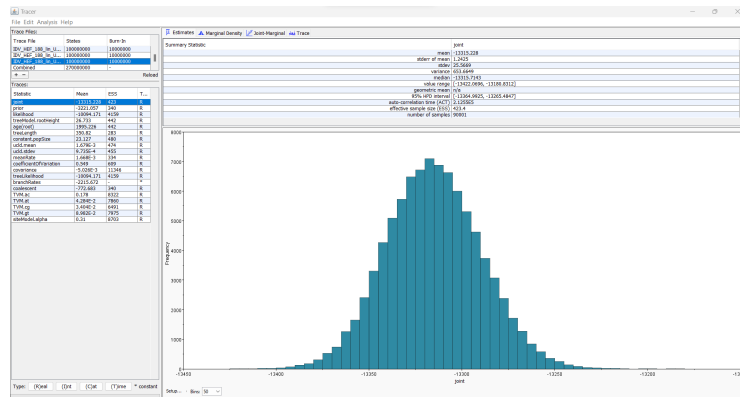
Figure 3.39: TRACER 1

Combined trace will automatically appear (with a number of iterations equal to the sum of the two traces minus their burn-in length).



Figure 3.40: TRACER 2

- Selecting the Combined trace allows to explore a concatenation of the log files. Tracer will plot a (marginal posterior) histogram for the selected statistic and also give you summary statistics such as the mean and median. The 95% HPD interval stands for highest posterior density interval and represents the most compact interval on the selected parameter that contains 95% of the posterior probability. In the top right of the window is a table of calculated statistics for the selected trace



Figure 3.41: TRACER 3

### 3.3.7 TREE ANNOTATOR

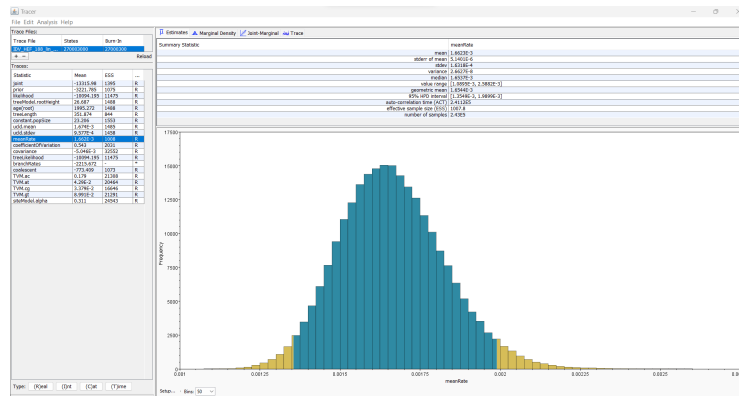- This is a post-analysis program that will produce a summary tree from the output of BEAST. It's simply not feasible to inspect every tree that was visited during the BEAST analysis, hence we will create a consensus tree summarizing the posterior tree distribution. Within the BEAST package, this is done by constructing a maximum clade credibilty (MCC) tree using the program TreeAnnotator.

  Upon running TreeAnnotator, you will be presented with the following window: Typically, 10% of the total number of iterations is used as the burn-in for analysis, provided that this is sufficient to have made it past the actual burn-in phase (which can be inspected/checked in Tracer). For the Input Tree File, select the file which was generated during the BEAST run needs to be selected. For the Output File, no such file can be selected but rather its file name needs to be entered manually in the following window:
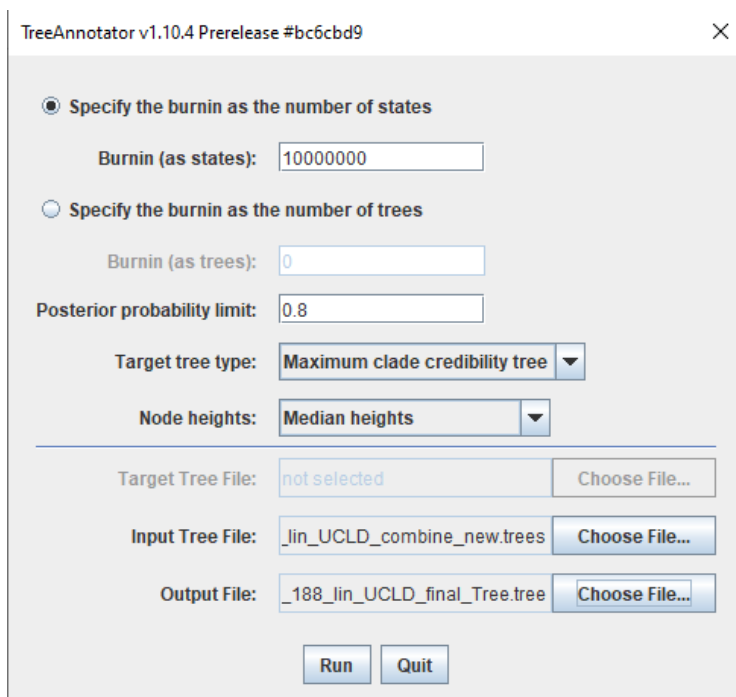


Figure 3.42: TREE ANNOTATOR 1

48

- After all the required settings have been entered, you can start constructing the MCC tree by clicking 'Run. Progress can monitored in the following window, at the end of which you will be asked to quit the program and the MCC tree will have been written to the selected output file:



Figure 3.43: TREE ANNOTATOR 2

## 3.3.8 FIG TREE

The final step is to load the constructed MCC tree into FigTree, which allows visualizing of the tree and accompanying summary information produced by TreeAnnotator. After starting FigTree, simply go to File and Open

- The following tree is the tree with node ages.



Figure 3.44: FIGTREE 1

- The following tree mentions the mean substitution rate.



Figure 3.45: FIGTREE 2

• The following tree mentions the 95% HPD intervals



Figure 3.46: FIGTREE 3

# Chapter 4

# Results and Conclusions

## 4.1 Results

- The following table summarizes our analysis which includes the best models for each gene along with the substitution rate and the time for a most recent ancestor.

| Dataset | Gene name | Substitution Model | Clock | $R^2$ | Correlation coefficient |
|---------|-----------|--------------------|-------|-------|-------------------------|
| 1 | PB2 | TIM+F+G4 | UCLN | 0.7447 | 0.863 |
| 2 | PB1 | TIM+F+G4 | STRICT | 0.0019 | 0.045 |
| 3 | P3 | TIM+F+I+G4 | STRICT | 0.4904 | 0.7003 |
| 4 | HEF | TVM+F+G4 | UCLN | 0.5927 | 0.7699 |
| 5 | NP | TN+F+G4 | STRICT | 0.6533 | 0.8082 |
| 6 | P42 | GTR+F+G4 | STRICT | 0.0640 | 0.2531 |
| 7 | NS | HKY+F+G4 | UCED | 0.2622 | 0.512 |
| 7.1 | NS1 | HKY+F+G4 | UCLN | 0.2939 | 0.5421 |
| 7.2 | NS2 | TN93+F+G4* | UCLN | 0.1826 | 0.4273 |

| Dataset | Gene name | MEAN NSR (sub/site/year) | NSR (95% HPD Interval) | tMRCA |
|---|---|---|---|---|
| 1 | PB2 | $1.35 \times 10^{-3}$ | $(1.122 \times 10^{-3}, 1.5935 \times 10^{-3})$ | 1997.47 |
| 2 | PB1 | $1.16 \times 10^{-3}$ | $(9.9688 \times 10^{-4}, 1.3341 \times 10^{-3})$ | 1997.04 |
| 3 | P3 | $1.31 \times 10^{-3}$ | $(1.1203 \times 10^{-3}, 1.4832 \times 10^{-3})$ | 1997.53 |
| 4 | HEF | $1.66 \times 10^{-3}$ | $(1.4047 \times 10^{-3}, 2.0118 \times 10^{-3})$ | 1996.09 |
| 5 | NP | $1.47 \times 10^{-3}$ | $(1.2389 \times 10^{-3}, 1.6881 \times 10^{-3})$ | 1998.97 |
| 6 | P42 | $1.35 \times 10^{-3}$ | $(1.106 \times 10^{-3}, 1.6039 \times 10^{-3})$ | 1997.30 |
| 7 | NS | $1.31 \times 10^{-3}$ | $(9.0492 \times 10^{-4}, 1.7197 \times 10^{-3})$ | 1997.36 |
| 7.1 | NS1 | $1.32 \times 10^{-3}$ | $(9.572 \times 10^{-4}, 1.7028 \times 10^{-3})$ | 1998.62 |
| 7.2 | NS2 | $1.63 \times 10^{-3}$ | $(1.0584 \times 10^{-3}, 2.2291 \times 10^{-3})$ | 2003.16 |

indicate that In the NS2 gene, TPM2+F+G4 is the actual best model but this model is not available in BUAUTi software hence we use the second-best model which is TN93+F+G4.

- **Model Explanation :**

- $F$: Empirical Base Frequency.

- $G4$: Site heterogeneity model in Gamma and number of Gamma category is 4

- UCLN: Uncorrelated Clock Lognormal Distribution

- UCED: Uncorrelated Clock Exponential Distribution

**Substitution model:**

- TIM: Transition model

- GTR: General time reversible

- TN: Tamura and Nei plot

- HKY: Hasegawa-Kishino-Yano

- TN93: Tamura and Nei plot 93

- TVM: Transversion-Varying Model

## 4.2 Conclusion on the basis of Molecular Phylogenetic Tree

The following tree best summarizes our Bayesian phylogenetic analysis.

- The following tree is the tree with node ages.



Figure 4.1: Best tree with node ages
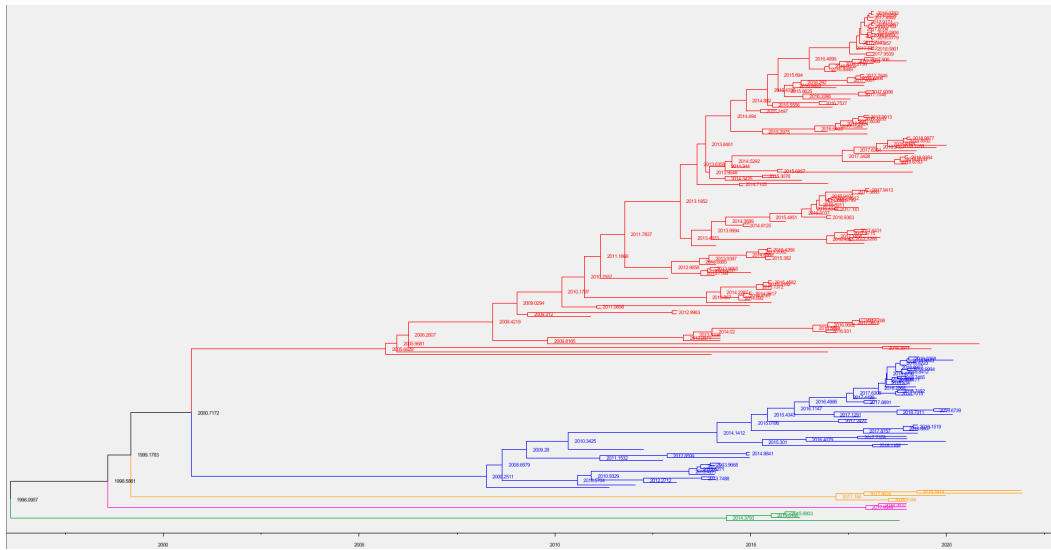
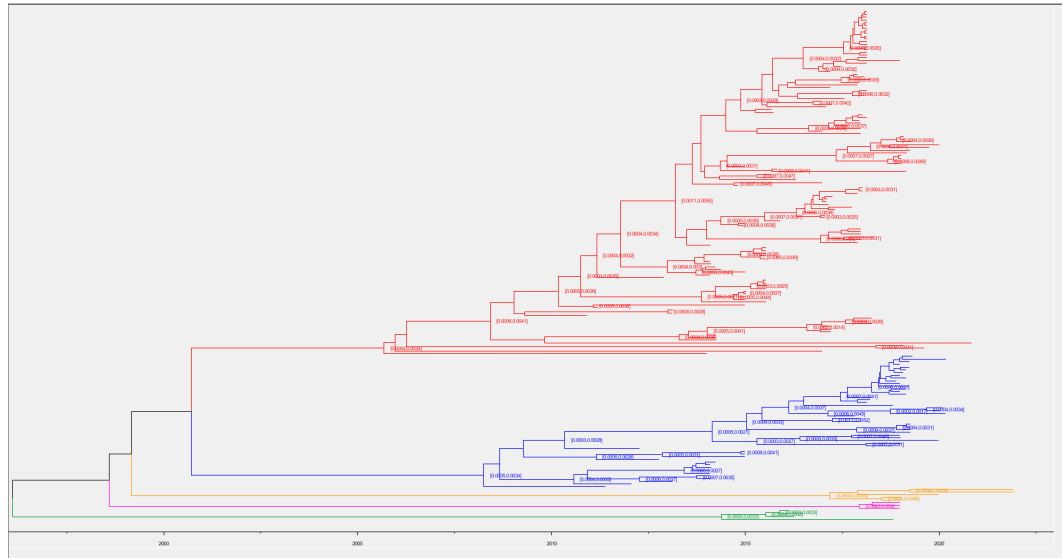- The following tree mentions the mean substitution rate.

Figure 4.2: Best tree with substitution rates per site per year

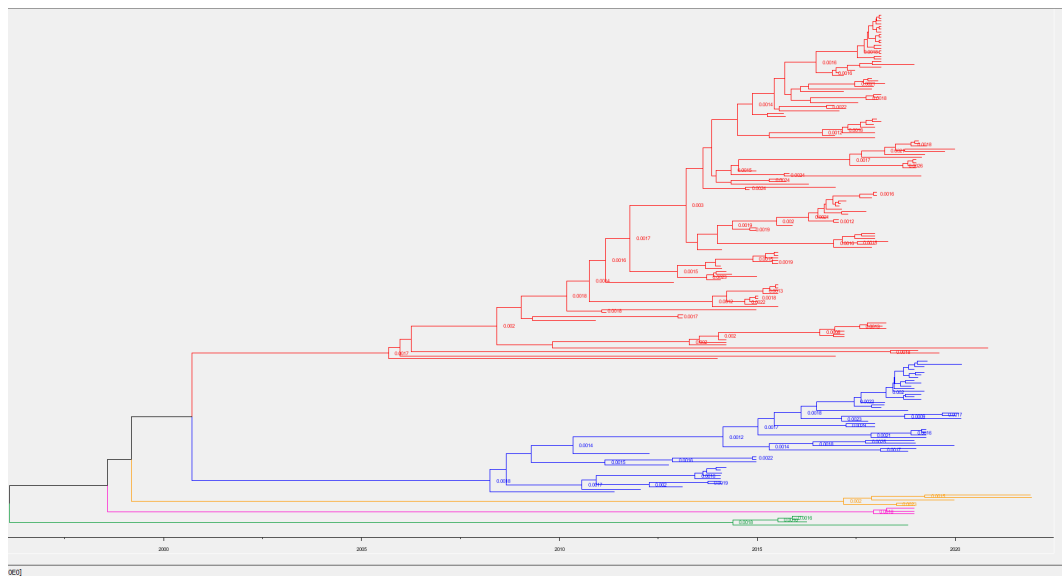- The following tree mentions the 95% HPD intervals



Figure 4.3: Best tree with 95% HPD intervals.

## 4.3   Statistical Conclusions:

- Particularly, for the $4^{th}$ gene (HEF) of IDV, we study its molecular clock behavior by analyzing the rate at which mutations accumulate in the HEF gene sequence and also by comparing the differences in the HEF gene sequences from different strains or isolates of IDV.

- In the case of IDV-HEF, the TVM+F+G4 model was identified as the best fit. The Transversion-Varying Model model (TVM) is a widely used substitution model that allows for different rates of nucleotide substitutions. The F parameter accounts for rate heterogeneity among sites, and the G4 parameter models the rate variation among different evolutionary categories.

- By applying uncorrelated lognormal distribution to the clock model for IDV-HEF, researchers can estimate the rates of molecular evolution and divergence times more accurately. This model accounts for rate variation and provides a more realistic representation of the evolutionary dynamics of the gene sequences.

- For the analysis of HEF, using the BEAST software, the convergence occurred at $10^8$ (MCMC chain length) for all the parameters ($\geq 200$) of Effective Sample Size (ESS).

- The mean rate for HEF, is estimated to be $1.66 \times 10^{-3}$ substitutions per site per year. This rate represents the average rate of evolutionary change in the gene sequence over time.

- Additionally the 95% HPD interval of mean substitution rates, which is a credible interval, is given as ($1.4047 \times 10^{-3}$, $2.0118 \times 10^{-3}$). This interval provides a range of plausible values for the true rate of evolution. It suggests that with 95% credible confidence the actual rate of evolutionary change in gene 4 IDV-HEF falls within this credible interval.

- The analysis of HEF resulted in the identification of five distinct clusters based on the marginal likelihood. This suggests that the data supports the presence of five separate groups or lineages within the gene sequence.

58

# Chapter 5

# Limitations and Scope

## 5.1 Limitations

- Bayesian statistical analysis relies on the availability of high-quality and comprehensive data and hence, problems like underreporting, surveillance gaps or difficulties in identifying and characterizing IDV strains can restrict the accuracy and reliability of the Bayesian analysis.

- Bayesian analysis incorporates prior information or assumptions about the data and model parameters. However, for emerging or newly discovered viruses like IDV, there may be a lack of prior knowledge making it challenging to establish informative prior distributions.

- Bayesian analysis often involves complex mathematical models which require numerous assumptions and they introduce uncertainties and potential biases. The accuracy of the analysis depends on the appropriateness of the chosen model and the quality of parameter estimates.

- Bayesian analysis is usually computationally intensive, particularly when dealing with large datasets or complex models. The analysis may require significant computational resources, time and expertise.

- The accuracy of uncertainty estimates depends on several factors, including the quality of data, model assumptions and parameter estimates and in situations where data or prior knowledge is limited, the uncertainty estimates may be less precise or less reliable.

- The results obtained from Bayesian analysis of IDV evolution may not be readily generalizable to other viruses or populations. Factors such as geographic location, host species and ecological factors can influence the evolutionary dynamics of IDV.

## 5.2   Scope

- Bayesian computational methods can be used to estimate the rate of molecular evolution in IDV which helps in understanding mode of its evolution and helps predict future evolutionary trends.

- Understanding the population dynamics of IDV is crucial for developing effective control and prevention strategies.

- Bayesian methods can be used to study the evolution of IDV and help in the development of improved vaccines and vaccination strategies against IDV.

- Bayesian statistical analysis is essential for understanding the virus's behaviour, improving surveillance and control measures and guiding public health interventions.

## 5.3 Bibliography

- Drummond AJ, Ho SY, Phillips MJ, Rambaut A. Relaxed phylogenetics and dating with confidence. PLoS biology. 2006 May;4(5):e88

- Drummond AJ, Rambaut A. BEAST: Bayesian evolutionary analysis by sampling trees. BMC evolutionary biology. 2007 Dec;7(1):1-8.

- Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior summarization in Bayesian phylogenetics using Tracer 1.7. Systematic biology. 2018 Sep 1;67(5):901-4.

- Ho SY, Duchêne S. Molecular-clock methods for estimating evolutionary rates and timescales. Molecular ecology. 2014 Dec;23(24):5947-65.

- Suchard MA, Lemey P, Baele G, Ayres DL, Drummond AJ, Rambaut A. Bayesian phylogenetic and phylodynamic data integration using BEAST 1.10. Virus evolution. 2018 Jan;4(1):vey016.

- Trifinopoulos J, Nguyen LT, von Haeseler A, Minh BQ. W-IQ-TREE: a fast online phylogenetic tool for maximum likelihood analysis. Nucleic acids research. 2016 Apr 15;44(W1):W232-5.