

SOFTWARE

Open Access

# $\pi$ BUSS: a parallel BEAST/BEAGLE utility for sequence simulation under complex evolutionary scenarios

Filip Bielejec<sup>1\*</sup>, Philippe Lemey<sup>1</sup>, Luiz Max Carvalho<sup>2</sup>, Guy Baele<sup>1</sup>, Andrew Rambaut<sup>3</sup> and Marc A Suchard<sup>4,5</sup>

## Abstract

**Background:** Simulated nucleotide or amino acid sequences are frequently used to assess the performance of phylogenetic reconstruction methods. BEAST, a Bayesian statistical framework that focuses on reconstructing time-calibrated molecular evolutionary processes, supports a wide array of evolutionary models, but lacked matching machinery for simulation of character evolution along phylogenies.

**Results:** We present a flexible Monte Carlo simulation tool, called  $\pi$ BUSS, that employs the BEAGLE high performance library for phylogenetic computations to rapidly generate large sequence alignments under complex evolutionary models.  $\pi$ BUSS sports a user-friendly graphical user interface (GUI) that allows combining a rich array of models across an arbitrary number of partitions. A command-line interface mirrors the options available through the GUI and facilitates scripting in large-scale simulation studies.  $\pi$ BUSS may serve as an easy-to-use, standard sequence simulation tool, but the available models and data types are particularly useful to assess the performance of complex BEAST inferences. The connection with BEAST is further strengthened through the use of a common extensible markup language (XML), allowing to specify also more advanced evolutionary models. To support simulation under the latter, as well as to support simulation and analysis in a single run, we also add the  $\pi$ BUSS core simulation routine to the list of BEAST XML parsers.

**Conclusions:**  $\pi$ BUSS offers a unique combination of flexibility and ease-of-use for sequence simulation under realistic evolutionary scenarios. Through different interfaces,  $\pi$ BUSS supports simulation studies ranging from modest endeavors for illustrative purposes to complex and large-scale assessments of evolutionary inference procedures. Applications are not restricted to the BEAST framework, or even time-measured evolutionary histories, and  $\pi$ BUSS can be connected to various other programs using standard input and output format.

**Keywords:** Simulation, Monte Carlo, Phylogenetics, BEAST, BEAGLE, Evolution

## Background

Recent decades have seen extensive development in phylogenetic inference, resulting in a myriad of techniques, each with specific properties concerning evolutionary model complexity, inference procedures and performance both in terms of speed of execution and estimation accuracy. With the development of such phylogenetic inference methods comes the need to synthesize evolutionary data in order to compare estimator performance

and to characterize strengths and weaknesses of different approaches (e.g. [1,2]). Whereas the true underlying evolutionary relationships between biological sequences are generally unknown, Monte Carlo simulations allow generating test scenarios while controlling for the evolutionary history as well as the tempo and mode of evolution. This has been frequently used to compare the performance of tree topology estimation (e.g. [3]), but it also applies to evolutionary parameter estimation and ancestral reconstruction problems (e.g. [4]). In addition, Monte Carlo sequence simulation has proven useful for assessing model adequacy (e.g. [5]) and for testing competing

\*Correspondence: filip.bielejec@rega.kuleuven.be

<sup>1</sup>Department of Microbiology and Immunology, Rega Institute, KU Leuven, Leuven, Belgium

Full list of author information is available at the end of the article

evolutionary hypotheses (e.g. [6]). It is therefore not surprising that several general sequence simulation programs have been developed (e.g. Seq-Gen [7]), but also inference packages that do not primarily focus on tree reconstruction, such as PAML [8] and HyPhy [9], maintain code to simulate sequence data under the models they implement.

As a major application of phylogenetics, estimating divergence times from molecular sequences requires an assumption of roughly constant substitution rates throughout evolutionary history [10]. Despite the restrictive nature of this molecular clock assumption, its application in a phylogenetic context has profoundly influenced modern views on the timing of many important events in evolutionary history [11]. Following a long history of applying molecular clock models on fixed tree topologies, the Bayesian Evolutionary Analysis by Sampling Trees (BEAST) package [12] fully integrates these models, including more realistic relaxed clock models [13,14], in a phylogenetic inference framework. Despite its popularity this framework has lacked a flexible and efficient simulation tool. Here, we address this pitfall by introducing a parallel BEAST/BEAGLE utility for sequence simulation ( $\pi$ BUSS) that integrates substitution models, molecular clock models, tree-generative (coalescent or birth-death) models and trait evolutionary models in a modular fashion, allowing the user to simulate sequences under different parameterizations for each module.

$\pi$ BUSS readily incorporates the temporal dimension of evolution through the possibility of specifying different molecular clock model. Further, many models and data types available for BEAST inference are matched by their simulation counter-parts in  $\pi$ BUSS, including relatively specific processes, such as for discrete phylogeography with rate matrices that can be sparse or non-reversible [15] that are generally beyond the scope of most sequence simulation tools. The BEAST- $\pi$ BUSS connection is further reinforced by the fact that  $\pi$ BUSS can easily generate simulation specification in XML format for BEAST. Finally, we implement the core simulation routine within the BEAST code-base to ensure a shared XML syntax between the two packages and to allow for joint simulation and inference analysis using a single input file.

## Implementation

Through different implementations, we support several sequence simulation procedures that balance between ease-of-use and accessibility, to model complexity. On the one hand, the core simulation routine can be performed following specifications in an XML input file that is understood by BEAST (Figure 1A). This procedure provides the most comprehensive access to the  $\pi$ BUSS arsenal of models, but may require custom XML editing. On the other hand,  $\pi$ BUSS also represents a stand-alone package that conveniently wraps the simulation

routines in a user-friendly graphical user interface (GUI), allowing users to set up and run simulations by loading input, selecting models from drop-down lists, setting their parameter values, and generating output in different formats (Figure 1B). To facilitate scripting,  $\pi$ BUSS is further accessible through a command-line interface (CLI), with options that mirror the GUI. The simulation routines are implemented in Java and interface with the Broad-platform Evolutionary Analysis General Likelihood Evaluator (BEAGLE) high-performance library [16] through its application programming interface (API) for computationally intensive tasks.

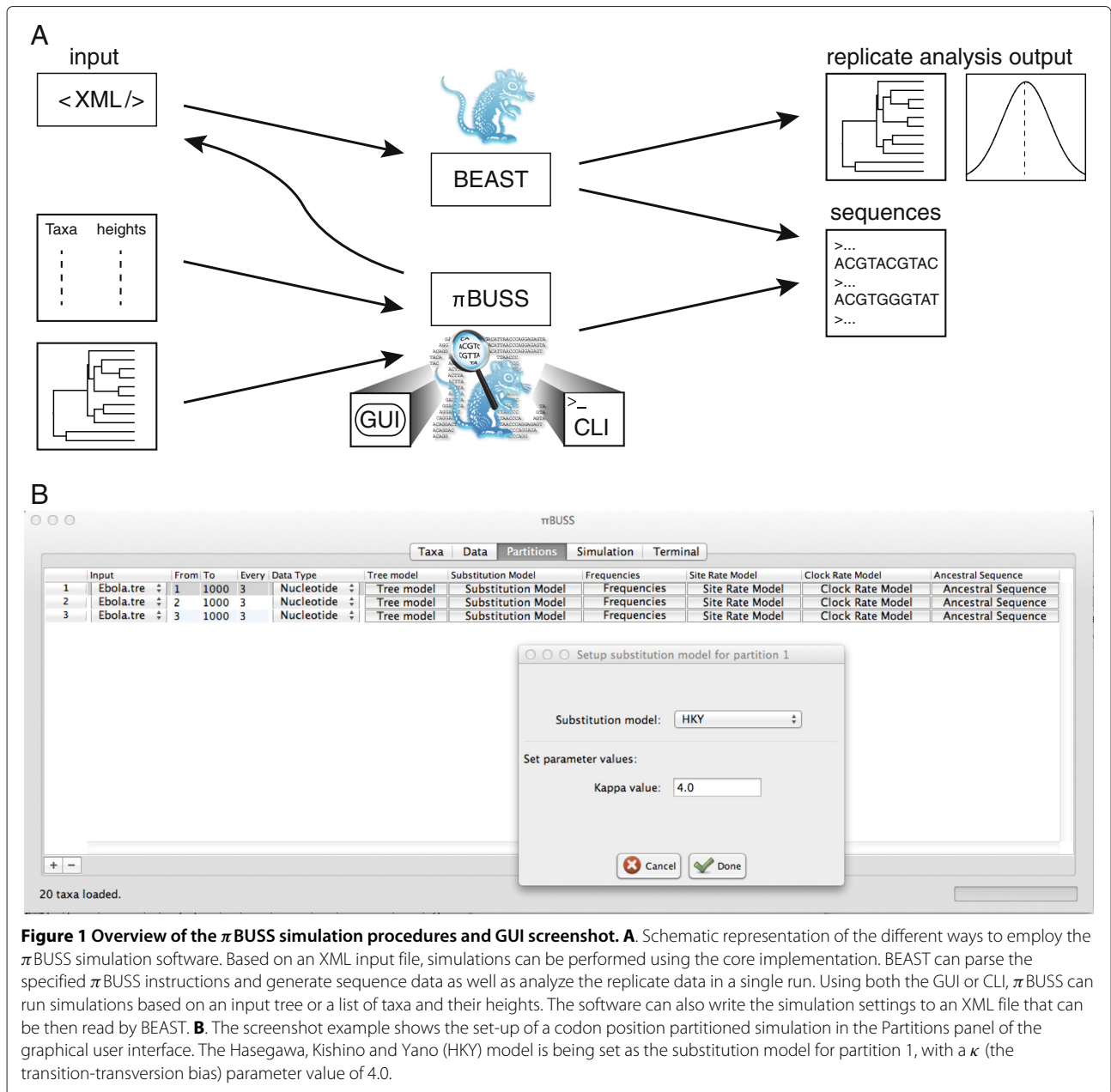
The core of  $\pi$ BUSS consists of a recursive tree-traversal that is independent of the BEAST inference machinery. The algorithm simulates discrete state realizations by visiting the tree nodes in pre-order fashion, i.e., parental nodes are visited before child nodes. When a child node is visited,  $\pi$ BUSS samples its state from the conditional probabilities of changing to state  $j$  given state  $i$  at the parental node. For a branch length  $t$  and clock rate  $r$ , the finite-time transition probability matrix  $\mathbf{P}(r \times t)$  is calculated through the eigen-decomposition of the infinitesimal rate matrix  $\mathbf{Q}$  along that branch. For a review of methods to numerically approximate a matrix exponential, we refer to [17]. By sharing the set of XML parsers with BEAST, we simplify the simultaneous development of both packages and facilitate the ability to perform joint simulation and inference analyses.

## Program input

The core implementation of the software can be invoked by loading an XML file with simulation settings in the BEAST software. The simulation procedure requires a user-specified tree topology or a set of taxa with their heights (inversely proportional to their sampling time) for which a tree topology can be simulated using a coalescent model. Setting all heights to 0 would be equivalent to contemporaneously-sampled taxa. In  $\pi$ BUSS, such a tree can be loaded in NEXUS or NEWICK format, or a taxa list can be set-up in the Data panel for subsequent coalescent simulation of the genealogy. Creating the latter is further facilitated by the ability to load a tab-delimited file with a set of taxa and their corresponding heights. The input tree or taxon list can also be specified through the command-line interface of  $\pi$ BUSS.

## Program output

$\pi$ BUSS generates sequence output in FASTA or NEXUS format but it also supports XML output of the simulation settings. The XML provides a notation for the models used, it can also be used to store a record of the settings. Similar to BEAute for BEAST,  $\pi$ BUSS can generate an xml template for editing more complex simulations, or this can be amended with BEAST analysis settings



**Figure 1 Overview of the  $\pi$ BUSS simulation procedures and GUI screenshot. A.** Schematic representation of the different ways to employ the  $\pi$ BUSS simulation software. Based on an XML input file, simulations can be performed using the core implementation. BEAST can parse the specified  $\pi$ BUSS instructions and generate sequence data as well as analyze the replicate data in a single run. Using both the GUI or CLI,  $\pi$ BUSS can run simulations based on an input tree or a list of taxa and their heights. The software can also write the simulation settings to an XML file that can be then read by BEAST. **B.** The screenshot example shows the set-up of a codon position partitioned simulation in the Partitions panel of the graphical user interface. The Hasegawa, Kishino and Yano (HKY) model is being set as the substitution model for partition 1, with a  $\kappa$  (the transition-transversion bias) parameter value of 4.0.

order to directly analyze the generated sequence data, which avoids writing to an intermediate file. The tutorial hosted on  $\pi$ BUSS webpage provides examples of these possibilities.

### Models of evolution

$\pi$ BUSS is capable of generating trees from a list of taxa using simple coalescent models, including a constant population size or exponential growth model. The software supports simulation of nucleotide, amino acid and codon data along the simulated or user-specified phylogeny using standard substitution models. For nucleotide

data, the Hasegawa, Kishino and Yano model (HKY; [18]), the Tamura Nei model (TN93; [19]) and the general time-reversible model (GTR; [20]) can be selected from a drop-down list, and more restrictive continuous-time Markov chain (CTMC) models can be specified by tailoring parameters values. Coding sequences can be simulated following the Goldman and Yang model of codon evolution (GY94; [21]), which is parameterized in terms of a non-synonymous and synonymous substitution rate ratio ( $dN/dS$  or  $\omega$ ) and a transition/transversion rate ratio ( $\kappa$ ) or following the Muse and Gaut model (MG94; [22]). Several empirical amino acid substitution models

are implemented, including the Dayhoff [23], JTT [24], BLOSUM [25], WAG [26] and LG [27] model. Equilibrium frequencies can be specified for all substitution models as well as among-site rate heterogeneity through the widely-used discrete-gamma distribution [28] and proportion of invariant sites [29].

An important feature of  $\pi$ BUSS is the ability to set up an arbitrary number of partitions for the sequence data and associate independent substitution models to them. Such settings may reflect codon position-specific evolutionary patterns or approximate genome architecture with separate substitution patterns for coding and non-coding regions. Partitions may also be set to evolve along different phylogenies, which could be used, for example, to investigate the impact of recombination or to assess the performance of recombination detection programs in specific cases. Finally, partitions do not need to share the exact same taxa (e.g. reflecting differential taxon sampling), and in partitions where a particular taxon is not represented the relevant sequence will be padded with gaps.

$\pi$ BUSS is equipped with the ability to simulate evolutionary processes on trees calibrated in time units. Under the strict clock assumption, this is achieved by specifying an evolutionary rate parameter that scales each branch from time units into substitution units.  $\pi$ BUSS also supports branch-specific scalars drawn independently and identically from an underlying distribution (e.g. log normal or inverse Gaussian distributions), modeling an uncorrelated relaxed clock process [13]. Simulations do not need to accommodate an explicit temporal dimension and input trees with branch lengths in substitution units will maintain these units with the default clock rate of 1 (substitution/per site/per time unit).

The data types and models described above are available through the  $\pi$ BUSS GUI or CLI, but additional data types and more complex models can be specified directly in an XML file. This allows, for example, simulating any discrete trait, e.g. representing phylogeographic locations, under reversible and nonreversible models [15,30], with potentially sparse CTMC matrices [15], as well as simulating a combination of sequence data and such traits. As an example of available model extensions is the ability to specify different CTMC matrices over different time intervals of the evolutionary history, allowing for example to model changing selective constraints through different codon model parameterizations or seasonal migration processes for viral phylogeographic traits [31].

## Results and discussion

We have developed a new simulation tool, called  $\pi$ BUSS, that we consider to be a rejuvenation of Seq-Gen [7], with several extensions to better integrate with the BEAST inference framework. Compared to Seq-Gen and other

simulation software (Table 1),  $\pi$ BUSS covers a relatively wide range of models while, similar to Mesquite, offering a cross-platform, user-friendly GUI.  $\pi$ BUSS is implemented in the Java programming language, and therefore requires a Java runtime environment, and depends on the BEAGLE library. Although speed is unlikely to be an impeding factor in most simulation efforts, the core implementation using the BEAGLE library provides substantial increases in speed for large-scale simulations, in particular when invoking multi-core architecture to produce highly partitioned synthetic sequence data.

### Program validation

We validate  $\pi$ BUSS in several ways. First, we compare the expected site probabilities, as calculated using tree pruning recursion [56], with the observed counts resulting from  $\pi$ BUSS simulations. To this purpose, we calculate the probabilities for all  $4^3$  possible nucleotide site patterns observed at the tips of a particular 3-taxon topology using an HKY model with a discrete gamma distribution to model rate variation among sites. We then compare these probabilities to long-run ( $n = 100,000$ ) site pattern frequencies simulated under this model and observe good correspondence in distribution (Pearson's  $\chi^2$  test,  $p = 0.42$ ).

We also perform simulations over larger trees and estimate substitution parameters (e.g.  $\kappa$  in the HKY model) using BEAST for a large number of replicates. Not only do the posterior mean estimates agree very well with the simulated values, but we also find close to nominal coverage, and relatively small bias and variance (mean squared error). These good performance measures have also recently been demonstrated for more complex substitution processes [31].

### Example application

We illustrate the use of simulating sequence data along time-calibrated phylogenies to explore the limitations of estimating old divergence times for rapidly-evolving viruses. Wertheim and Kosakovsky Pond [57] examine the evolutionary history of Ebola virus from sequences sampled over the span of three decades. Although maintaining remarkable amino acid conservation, the authors estimate nucleotide substitution rates on the order of  $10^{-3}$  substitutions/per site/per year and a time to most recent common ancestor (tMRCA) of about 1,000 years ago. These estimates suggest a strong action of purifying selection to preserve amino acid residues over longer evolutionary time scales, which may not be accommodated by standard nucleotide substitution models. The authors demonstrate that accounting for variable selective pressure using codon models can result in substantially older origins in such cases.

**Table 1 Comparison between a selection of sequence simulation packages**

Program	Evolutionary modelling						Implementation			
	Codons <sup>1</sup>	Amino acids <sup>2</sup>	Indels	Partitions	Molecular clocks	Ancestral sequences	Coalescent models <sup>3</sup>	GUI	Multi-core	Cross-platform <sup>4</sup>
$\pi$ BUSS	X	X		X	X	X	X	X	X	X
Seq-Gen [7]		X		X						X
indel-Seq-Gen2 [32]		X	X	X		X				X
PhyloSim [33]	X	X	X	X		X				X
Recodon [34]	X					X	X		X	X
NetRecodon [35]	X					X	X		X	X
Indelible [36]	X	X	X	X				X		
DAWG [37]			X			X	X			X
Mesquite [38]						X	X	X		X
Rose [39]			X			X				
Evolver [8]	X	X		X		X				X
ProteinEvolver [40]		X				X	X		X	X
ALF [41]	X	X	X	X		X	X	X		X
GenomePop [42]	X					X	X <sup>5</sup>			X
SIMCOAL [43]						X	X			X
SIMPROT [44]		X	X			X		X		X

<sup>1</sup> $\pi$ BUSS: GY94, MG94; PhyloSim: GY94 x M0 - M4; Recodon: GY94 x M0, M1, M7, M8; NetRecodon: GY94 x M0, M1, M7, M8; Indelible: GY94 x M0 - M10; Evolver: GY94 x M0, M1, M2, M3, M7, M8; ALF: GY94 x M0, M1, M7 and M8; GenomePop: MG94.

<sup>2</sup> $\pi$ BUSS: BLOSUM [25], CPREV [45], Dayhoff [23], FLU [46], JTT [24], LG [27], MTREV [47], WAG [26]; Seq-Gen: JTT, WAG, PAM [48], BLOSUM, MTREV; indel-Seq-Gen2: PAM, JTT, MTREV, CPREV; PhyloSim: CPREV, JTT, LG, MTART [49], MTMAM [50], MTREV24 [51], MTZOA [52], PAM, WAG; Indelible: Dayhoff, JTT, WAG, VT [53], LG, BLOSUM, MTMAM, MTREV, MTART, CPREV, RTREV [54], HIVb [55], HIVw [55]; Evolver: Dayhoff, JTT, WAG, MTMAM, MTREV; ProteinEvolver: BLOSUM, CPREV, Dayhoff, HIVb, HIVw, JTT, Jones [24], LG, MTART, MTMAM, MTREV24, RTREV, VT, WAG; ALF: PAM, GTT, LG, WAG; SIMPROT: PAM, JTT, PMB.

<sup>3</sup> $\pi$ BUSS: demography; Recodon: recombination, migration, demography; NetRecodon: recombination, migration, demography; Mesquite: speciation; ProteinEvolver: recombination, migration, demography; SIMCOAL: demography and migration.

<sup>4</sup>PhyloSim: R package; Indelible: Executables for Windows and MacOS; ALF: Web interface; GenomePop: Executables for Windows and Linux; SIMCOAL: Executables for Windows; SIMPROT: GUI, Web interface; SIMPROT: Executables for Windows and Linux.

<sup>5</sup>Forward simulation including recombination, demography and migration.

Here, we explore the effect of temporally varying selection pressure throughout evolutionary history on estimates of the tMRCAs using nucleotide substitution models. In particular, we model a process that is characterized by increasingly stronger purifying selection as we go further back in time. To this purpose, we set up an 'epoch model' that specifies different GY94 codon substitution processes along the evolutionary history [31], and parameterize them according to a log-linear relationship between time and  $\omega$ . Specifically, we let the process transition from  $\omega = 1.0, 0.2, 0.1, 0.02, 0.01, 0.002,$  and  $0.001$  at time = 10, 50, 100, 500, 1000 and 5000 years in the past, respectively. We simulate a constant population size genealogy of 50 taxa, sampled evenly during a time interval of 25 years, and simulate sequences according to the time-heterogeneous codon substitution process with a constant clock rate of  $3 \times 10^{-3}$  codon substitutions/codon site/year. We simulate 100 replicates over genealogies with varying tMRCAs, by generating topologies under different population sizes parameterized by the product of effective population size ( $N_e$ ) and generation time scaled in years ( $\tau$ ):

$$N_e \times \tau = 1, 5, 10, 50, 100, 500, 1000$$

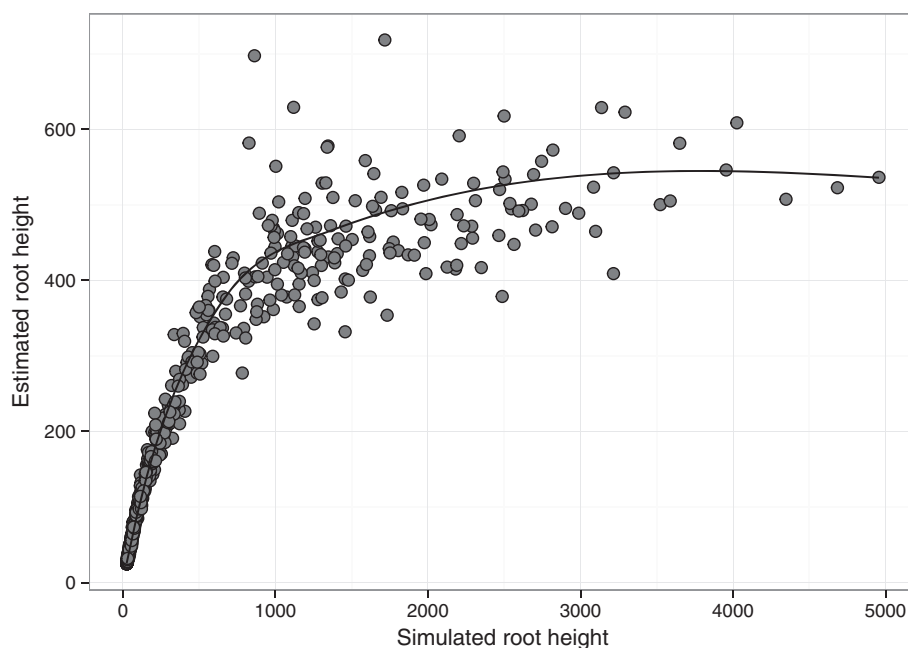
We note that under this model, trees with tMRCAs of about 10,000 years still result in sequences with a noticeable degree of homology (resulting in a mean amino acid distance of about 0.5, which is in the same range of the mean amino acid distance for sequences representative of the primate immunodeficiency virus diversity).

Using a constant  $\omega$  of 0.5 on the other hand results in fairly randomized sequences. We subsequently analyze the replicate data using a codon position partitioned nucleotide substitution model in BEAST and plot the correspondence between simulated and estimated tMRCAs in Figure 2.

Our simulation exercise shows that a linear relationship between simulated and estimated tMRCAs only holds for 100 to 200 years in the past, and estimates quickly level off after about 1000 years in the past. This can be explained by the unaccounted decline in amino acid substitutions and saturation of the synonymous substitutions as we go further back in time. Although we are not claiming that evolution occurs quantitatively or even qualitatively according to the particular process we simulate under, and we ignore other confounding factors (such as potential selective constraints on non-neutral synonymous sites), this simulation does conceptualize some of the limitations to estimating ancient origins for rapidly evolving viruses that experience strong purifying selection over longer evolutionary time scales.

## Conclusion

$\pi$ BUSS provides simulation procedures under many evolutionary models or combinations of models available in the BEAST framework. This feature facilitates the evaluation of estimator performance during the development of novel inference techniques and the generation of predictive distributions under a wide range of evolutionary scenarios that remain critical for testing



**Figure 2** Correspondence between simulated and estimated tMRCAs when purifying selection increases back in time in simulated data sets.

competing evolutionary hypotheses. Combinations of different evolutionary models can be accessed through a GUI or CLI, and further extensions can be specified in XML format with a syntax familiar to the BEAST user community. Analogous to the continuing effort to support model set-up for BEAST in BEAUti, future releases of  $\pi$ BUSS aim to provide simulation counterparts to the BEAST inference tools, both in terms of data types and models, while also maintaining general purposes simulation capabilities. Interesting targets include discrete traits, which can already be simulated through XML specification, continuously-valued phenotype data [58] and indel models. Finally,  $\pi$ BUSS provides opportunities to pursue further computational efficiency through parallelization on advancing computing technology. We therefore hope that  $\pi$ BUSS will further stimulate the development of sequence and trait evolutionary models and contribute to advancement of our knowledge about evolutionary processes.

## Availability and requirements

**Project name:**  $\pi$ BUSS;

**Project home page:** [www.rega.kuleuven.be/cev/ecv/software/pibuss](http://www.rega.kuleuven.be/cev/ecv/software/pibuss);

**Operating system(s):** Platform independent;

**Programming language:** Java;

**Other requirements:** Java 1.5 or higher, BEAGLE library;

**License:** GNU Lesser GPL;

**Any restrictions to use by non-academics:** None.

Source code of the parallel BEAST/BEAGLE utility for sequence simulation is freely available as part of the BEAST Google Code repository: [www.code.google.com/p/beast-mcmc/](http://www.code.google.com/p/beast-mcmc/).

The Broad-platform Evolutionary Analysis General Likelihood Evaluator (BEAGLE) library has its both source code and binary installers available from [www.code.google.com/p/beagle-lib](http://www.code.google.com/p/beagle-lib).

Scripts and input files required for repeating the simulation study presented in **Example application** are hosted at [www.github.com/phylogeography/DeepRoot](http://www.github.com/phylogeography/DeepRoot).

## Competing interests

The authors declare that they have no competing interests.

## Authors' contributions

FB designed and implemented the software. Large portions of code extend or are based on interfaces designed and developed by AR and MAS. PL and GB conceived the original idea and helped with the design of the software. PL designed the simulation study employing the software. LMC wrote the software tutorial, tested and commented on the software. All authors contributed to the writing of this manuscript. All authors read and approved the final manuscript.

## Acknowledgments

The research leading to these results has received funding from the European Research Council under the European Community's Seventh Framework Programme (FP7/2007-2013) under Grant Agreement no. 278433-PREDEMICS and ERC Grant agreement no. 260864, the Wellcome Trust (grant no. 092807)

and the National Institutes of Health (R01 GM086887 and R01 HG006139) and the National Science Foundation (IIS 1251151 and DMS 1264153). The National Evolutionary Synthesis Center (NESCent) catalyzed this collaboration through a working group (NSF EF-0423641). LMC would like to thank the Program for Scientific Computing staff for operational support.

## Author details

<sup>1</sup>Department of Microbiology and Immunology, Rega Institute, KU Leuven, Leuven, Belgium. <sup>2</sup>Program for Scientific Computing (PROCC), Fundação Oswaldo Cruz, Rio de Janeiro, Brazil. <sup>3</sup>Institute of Evolutionary Biology, University of Edinburgh, Edinburgh, UK. <sup>4</sup>Departments of Biomathematics and Human Genetics, David Geffen School of Medicine at UCLA, University of California, Los Angeles, CA, 90095, USA. <sup>5</sup>Department of Biostatistics, UCLA Fielding School of Public Health, University of California, Los Angeles, CA, 90095, USA.

Received: 10 October 2013 Accepted: 24 April 2014

Published: 7 May 2014

## References

1. Arenas M: **Simulation of molecular data under diverse evolutionary scenarios.** *PLoS Comput Biol* 2012, **8**(5):e1002495.
2. Hoban S, Bertorelle G, Gaggiotti OE: **Computer simulations: tools for population and evolutionary genetics.** *Nat Rev Genet* 2011, **13**(2):110–122.
3. Stamatakis A: **An efficient program for phylogenetic inference using simulated annealing.** In *Parallel and Distributed Processing Symposium, 2005. Proceedings. 19th IEEE International.* New York, USA: IEEE; 2005.
4. Blanchette M, Diallo AB, Green ED, Miller W, Haussler D: **Computational reconstruction of ancestral DNA sequences.** *Methods Mol Biol* 2008, **422**:171–184.
5. Brown JM, ElDabaje R: **PuMA: Bayesian analysis of partitioned (and unpartitioned) model adequacy.** *Bioinformatics* 2009, **25**(4):537–538.
6. Goldman N: **Statistical tests of models of DNA substitution.** *J Mol Evol* 1993, **36**(2):182–198.
7. Rambaut A, Grass NC: **Seq-Gen: an application for the Monte Carlo simulation of DNA sequence evolution along phylogenetic trees.** *Comput Appl Biosci* 1997, **13**(3):235–238.
8. Yang Z: **PAML 4: Phylogenetic Analysis by Maximum Likelihood.** *Mol Biol Evol* 2007, **24**(8):1586–1591.
9. Kosakovsky Pond SL, Frost SDW, Muse SV: **HyPhy: hypothesis testing using phylogenies.** *Bioinformatics* 2005, **21**(5):676–679.
10. Zuckerkandl E, Pauling LB: *Molecular Disease, Evolution, and Genetic Heterogeneity.* New York: Academic Press; 1962.
11. Arbogast BS, Edwards SV, Wakeley J, Beerli P, Slowinski JB: **Estimating divergence times from molecular data on phylogenetic and population genetic timescales.** *Annu Rev Ecol Evol Systemat* 2002, **33**:707–740.
12. Drummond AJ, Suchard MA, Xie D, Rambaut A: **Bayesian phylogenetics with BEAUti and the BEAST 1.7.** *Mol Biol Evol* 2012, **29**(8):1969–1973.
13. Drummond AJ, Ho SYW, Phillips MJ, Rambaut A: **Relaxed phylogenetics and dating with confidence.** *PLoS Biol* 2006, **4**(5):e88.
14. Drummond A, Suchard M: **Bayesian random local clocks, or one rate to rule them all.** *BMC Biol* 2010, **8**:114.
15. Lemey P, Rambaut A, Drummond AJ, Suchard MA: **Bayesian Phylogeography Finds Its Roots.** *PLoS Comput Biol* 2009, **5**(9):e1000520.
16. Ayres DL, Darling A, Zwickl DJ, Beerli P, Holder MT, Lewis PO, Huelsenbeck JP, Ronquist F, Swofford DL, Cummings MP, Rambaut A, Suchard MA: **BEAGLE: an application programming interface and high-performance computing library for statistical phylogenetics.** *Syst Biol* 2012, **61**:170–173.
17. Moler C, Loan CV: **Nineteen dubious ways to compute the exponential of a matrix.** *SIAM Rev* 1978, **20**:801–836.
18. Hasegawa M, Kishino H, Yano Ta: **Dating of the human-ape splitting by a molecular clock of mitochondrial DNA.** *J Mol Evol* 1985, **22**:160–174.
19. Tamura K, Nei M: **Estimation of the number of nucleotide substitutions in the control region of mitochondrial DNA in humans and chimpanzees.** *Mol Biol Evol* 1993, **10**(3):512–526.
20. Tavaré S: **Some probabilistic and statistical problems in the analysis of DNA sequences.** *Lect Math Life Sci (American Mathematical Society)* 1986, **17**:57–86.

21. Goldman N, Yang Z: **A codon-based model of nucleotide substitution for protein-coding DNA sequences.** *Mol Biol Evol* 1994, **11**(5):725–736.
22. Muse SV, Gaut BS: **A likelihood approach for comparing synonymous and nonsynonymous nucleotide substitution rates, with application to the chloroplast genome.** *Mol Biol Evol* 1994, **11**(5):715–724.
23. Dayhoff MO, Schwartz RM: **A model of evolutionary change in proteins.** In *Atlas of Protein Sequence and Structure*. Washington, D.C., USA: Citeseer, National Biomedical Research Foundation; 1978.
24. Jones DT, Taylor WR, Thornton JM: **The rapid generation of mutation data matrices from protein sequences.** *Comput Appl Biosci* 1992, **8**(3):275–282.
25. Henikoff S, Henikoff JG: **Amino acid substitution matrices from protein blocks.** *Proc Natl Acad Sci* 1992, **89**(22):10915–10919.
26. Whelan S, Goldman N: **A general empirical model of protein evolution derived from multiple protein families using a maximum-likelihood approach.** *Mol Biol Evol* 2001, **18**(5):691–699.
27. Le SQ, Gascuel O: **An improved general amino acid replacement matrix.** *Mol Biol Evol* 2008, **25**(7):1307–1320.
28. Yang Z: **Among-site rate variation and its impact on phylogenetic analyses.** *Trends Ecol Evol* 1996, **11**(9):367–372.
29. Gu X, Fu YX, Li WH: **Maximum likelihood estimation of the heterogeneity of substitution rate among nucleotide sites.** *Mol Biol Evol* 1995, **12**(4):546–557.
30. Edwards CJ, Suchard MA, Lemey P, Welch JJ, Barnes I, Fulton TL, Barnett R, O'Connell TC, Coxon P, Monaghan N, Valdiosera CE, Lorenzen ED, Willerslev E, Baryshnikov GF, Rambaut A, Thomas MG, Bradley DG, Shapiro B: **Ancient hybridization and an Irish origin for the modern polar bear matriline.** *Curr Biol* 2011, **21**:1251–1258.
31. Bielejec F, Lemey P, Baele G, Rambaut A, Suchard MA: **Inferring heterogeneous evolutionary processes through time: from sequence substitution to phylogeography.** *Syst Biol* 2014. [http://sysbio.oxfordjournals.org/content/early/2014/04/21/sysbio.syu015]
32. Strobe CL, Abel K, Scott SD, Moriyama EN: **Biological sequence simulation for testing complex evolutionary hypotheses: indel-Seq-Gen version 2.0.** *Mol Biol Evol* 2009, **26**(11):2581–2593.
33. Sipos B, Massingham T, Jordan G, Goldman N: **PhyloSim - Monte Carlo simulation of sequence evolution in the R statistical computing environment.** *BMC Bioinformatics* 2011, **12**:104. [http://www.biomedcentral.com/1471-2105/12/104]
34. Arenas M, Posada D: **Recodon: coalescent simulation of coding DNA sequences with recombination, migration and demography.** *BMC Bioinformatics* 2007, **8**:458.
35. Arenas M, Posada D: **Coalescent simulation of intracodon recombination.** *Genetics* 2010, **184**(2):429–437.
36. Fletcher W, Yang Z: **INDELible: a flexible simulator of biological sequence evolution.** *Mol Biol Evol* 2009, **26**(8):1879–1888.
37. Cartwright RA: **DNA assembly with gaps (Dawg): simulating sequence evolution.** *Bioinformatics* 2005, **21**(Suppl 3):i31–i38.
38. Maddison WP, Maddison D: **Mesquite: a modular system for evolutionary analysis.** 2011. [http://mesquiteproject.org]
39. Stoye J, Evers D, Meyer F: **Rose: generating sequence families.** *Bioinformatics* 1998, **14**(2):157–163.
40. Arenas M, Dos Santos HG, Posada D, Bastolla U: **Protein evolution along phylogenetic histories under structurally constrained substitution models.** *Bioinformatics* 2013, **29**(23):3020–3028.
41. Dalquen DA, Anisimova M, Gonnert GH, Dessimoz C: **ALF—a simulation framework for genome evolution.** *Mol Biol Evol* 2012, **29**(4):1115–1123.
42. Carvajal-Rodriguez A: **GENOMEPOP: a program to simulate genomes in populations.** *BMC Bioinformatics* 2008, **9**:223.
43. Excoffier L, Novembre J, Schneider S: **SIMCOAL: a general coalescent program for the simulation of molecular data in interconnected populations with arbitrary demography.** *J Hered* 2000, **91**(6):506–509.
44. Pang A, Smith AD, Nuin PA, Tillier ER: **SIMPROT: using an empirically determined indel distribution in simulations of protein evolution.** *BMC Bioinformatics* 2005, **6**:236.
45. Adachi J, Waddell PJ, Martin W, Hasegawa M: **Plastid genome phylogeny and a model of amino acid substitution for proteins encoded by chloroplast DNA.** *J Mol Evol* 2000, **50**(4):348–358.
46. Dang C, Le Q, Gascuel O, Le V: **FLU, an amino acid substitution model for influenza proteins.** *BMC Evol Biol* 2010, **10**:99. [http://www.biomedcentral.com/1471-2148/10/99]
47. Adachi J, Hasegawa M: **Model of amino acid substitution in proteins encoded by mitochondrial DNA.** *J Mol Evol* 1996, **42**(4):459–468.
48. Dayhoff M, Eck R, (US) NBRF: *Atlas of Protein Sequence and Structure* 1965. t. 1, National Biomedical Research Foundation 1965. [http://books.google.be/books?id=9Hp5nAEACAAJ]
49. Abascal F, Posada D, Zardoya R: **MtArt: a new model of amino acid replacement for Arthropoda.** *Mol Biol Evol* 2007, **24**:1–5.
50. Cao Y, Janke A, Waddell PJ, Westerman M, Takenaka O, Murata S, Okada N, Paabo S, Hasegawa M: **Conflict among individual mitochondrial proteins in resolving the phylogeny of eutherian orders.** *J Mol Evol* 1998, **47**(3):307–322.
51. Adachi J, Hasegawa M: *MOLPHY Version 2.3: Programs for Molecular Phylogenetics Based on Maximum Likelihood*. Tokyo, Japan: Computer science monographs 28, Institute of Statistical mathematics Tokyo; 1996.
52. Rota-Stabelli O, Yang Z, Telford MJ: **MtZoa: a general mitochondrial amino acid substitutions model for animal evolutionary studies.** *Mol Phylogenet Evol* 2009, **52**:268–272.
53. Muller T, Vingron M: **Modeling amino acid replacement.** *J Comput Biol* 2000, **7**(6):761–776.
54. Dimmic MW, Rest JS, Mindell DP, Goldstein RA: **rtREV: an amino acid substitution matrix for inference of retrovirus and reverse transcriptase phylogeny.** *J Mol Evol* 2002, **55**:65–73.
55. Nickle DC, Heath L, Jensen MA, Gilbert PB, Mullins JJ, Kosakovsky Pond SL: **HIV-specific probabilistic models of protein evolution.** *PLoS ONE* 2007, **2**(6):e503.
56. Felsenstein J: **Evolutionary trees from DNA sequences: A maximum likelihood approach.** *J Mol Evol* 1981, **17**:368–376.
57. Wertheim JO, Kosakovsky Pond SL: **Purifying selection can obscure the ancient age of viral lineages.** *Mol Biol Evol* 2011, **28**(12):3355–3365.
58. Lemey P, Rambaut A, Welch JJ, Suchard MA: **Phylogeography takes a relaxed random walk in continuous space and time.** *Mol Biol Evol* 2010, **27**(8):1877–85.

doi:10.1186/1471-2105-15-133

Cite this article as: Bielejec et al.:  $\pi$ BUSS: a parallel BEAST/BEAGLE utility for sequence simulation under complex evolutionary scenarios. *BMC Bioinformatics* 2014 **15**:133.

Submit your next manuscript to BioMed Central and take full advantage of:

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
www.biomedcentral.com/submit

