

# Ngesh: a Python library for synthetic phylogenetic data

Tiago Tresoldi<sup>1, 2</sup>

<sup>1</sup> Department of Linguistics and Philology, Uppsala University <sup>2</sup> Department of Linguistic and Cultural Evolution, Max Planck Institute for the Science of Human History

DOI: [10.21105/joss.03173](https://doi.org/10.21105/joss.03173)

## Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Mark A. Jensen](#) ↗

## Reviewers:

- [@DavidNickle](#)
- [@rvosa](#)

Submitted: 24 March 2021

Published: 26 May 2021

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

This work presents [ngesh](#), a Python library for simulating phylogenetic trees and data, designed for usage in development, debugging, and benchmarking of analysis pipelines and methods for phylogenetic inference, particularly in historical linguistics and stemmatics. The package generates reproducible stochastic simulations of evolution according to various criteria, including character mutation rates and probability of horizontal transfer, and its results can include the simulation of inadequate data compilation and sampling. Different output formats are supported, both for visualization (such as plain text and with integrated graphical viewers) and for software interoperability (such as Newick and NEXUS).

## Background

Computational phylogenetics is being increasingly accepted in fields beyond biology, such as historical linguistics ([Bouckaert et al., 2012](#)) and stemmatics ([Robinson, 2016](#)). Stochastic simulations, long advocated for natural sciences in general ([Bailey, 1964](#)) and genetics in specific ([Foote et al., 1999](#); [Harmon, 2019](#)), are not used enough in these fields. However, they are very desirable, allowing to evaluate evolutionary analogies, models, and performance through vast amounts of simulated histories, without limits imposed by data availability and collection time, with quantifiable precision of results. Simulations can also be used to perform fuzzy testing of software and to support studies on which evolutionary models, processes, and evolutionary parameters better match the observed phenomena.

The [ngesh](#) library is a tool that allows to perform such simulations, designed for easy integration into phylogenetic pipelines. It can generate reproducible trees and correlated data following both user-established parameters, such as ratios of birth and death, and constraints, such as branch lengths and minimum number of taxa. The library can label taxa progressive enumeration or with random names that are easy to pronounce (e.g. “Sume” and “Fekobir”) or which imitate the binominal nomenclature (e.g. “Sburas wioris” and “Zurbata pusso”). Character evolution related to the tree topology can likewise be simulated, including *ex novo* mutations and horizontal gene transfers. Results can be manipulated in diverse manners, for example by pruning extinct leaves or simulating uneven sampling. The simulated trees are standard ETE3 objects ([Huerta-Cepas et al., 2016](#)) and may be exported into different formats such as Newick trees, ASCII-art representation, and tabular lists.

## Statement of need

The library addresses the need of more tools to investigate and teach phylogenetics in historical linguistics and stemmatics. As a building block for evaluating pipelines of analysis, it is an

alternative to the basic technique of randomizing taxa placement in existing cladograms, and to simpler tools such as the one by Noutahi (2017) or the populate() method of ETE3's Tree class (Huerta-Cepas et al., 2016). It compares to popular alternatives, such as the TreeSim (Stadler, 2011) and geiger packages (Pennell et al., 2014) and the rtree() function of ape (Paradis & Schliep, 2018), by using default parameters that produce trees closer to those found in the fields of its intended scope, also in terms of file formats that better fit existing pipelines.

## Installation, Usage, & Examples

Users can install the library with the standard pip tool for managing Python packages. Trees can be generated from the command-line, defaulting to small phylogenies in Newick format:

```
$ ngesh
(Ukis:1.11985,(Koge:0.880823,(Rozkob:0.789548,(Meu:0.706601,
(((Felbuh:0.189693,Kefa:0.189693)1:0.117347,((Epib:0.153782,
Vugog:0.153782)1:0.0884745,Puluk:0.242256)1:0.0647836)1:0.0469885,
Efam:0.354028)1:0.352573)1:0.0829465)1:0.0912757)1:0.23903);
```

The tool supports both configuration files and command-line flags that take precedence over the former. Here we specify a model to generate Nexus data for a reproducible Yule tree, with a birth rate of 0.75, at least 5 leaves, "human" labels, and 20 presence/absence features:

```
$ cat my_tree.conf
[Config]
labels=human
birth=0.75
death=0.0
output=nexus
min_leaves=5
num_chars=20
$ ngesh -c my_tree.conf --seed 12345
begin data;
  dimensions ntax=6 nchar=33;
  format datatype=standard missing=? gap=-;
  matrix
Buza      111110110111011011010101000100110
Lenlar    111111010110111101100010010011001
Mukom     111110111011011011101001000100110
Pagil     111110110111011011100100100100110
Suglu     111110110111011011100011001001010
Wite      111110110111011011100101000100110
;
end;
```

Despite the benefit of a stand-alone tool, the package is designed to be run as a library. The two primary functions are gen\_tree(), which returns a random tree, and add\_characters(), which adds character evolution data to a tree. Users can generate random trees without character information or simulate character evolution within existing trees, including non-simulated ones.

```
>>> import ngesh
>>> tree = ngesh.gen_tree(1.0, 0.5, max_time=0.3, labels="bio",
```

```

seed="135")

>>> print(tree)

    /-Lubedsas larpes
  --|
    | /-Rasso wimapudda
    \-|
      \-Sbaes rapis
>>> print(tree.write())
(Lubedsas larpes:0.201311,(Rasso wimapudda:0.0894405,Sbaes rapis:0.0894405)
1:0.11187);
>>> tree = ngesh.add_characters(tree, 15, 2.0, 0.5)

```

56 Besides the `write()` method above, which outputs Newick trees, results can be exported  
 57 in NEXUS format with `tree2nexus()` and in a tabular output expected by tools such as  
 58 BEASTling (Maurits et al., 2017), with `tree2wordlist()`.

## 59 Code and Documentation Availability

60 The ngesh source code is available on GitHub at <https://github.com/tresoldi/ngesh>.  
 61 User documentation is available at <https://ngesh.readthedocs.io/>.

## 62 Acknowledgements

63 The author has received funding from the Riksbankens Jubileumsfond (ID: MXM19-1087:1,  
 64 “Cultural Evolution of Texts”) and from the European Research Council (ERC) under the  
 65 European Union’s Horizon 2020 research and innovation programme (No. ERC #715618,  
 66 “Computer-Assisted Language Comparison”).

## 67 References

- 68 Bailey, N. T. J. (1964). *The elements of stochastic processes with applications to the natural*  
 69 *sciences*. John Wiley & Soins.
- 70 Bouckaert, R., Lemey, P., Dunn, M., Greenhill, S. J., Alekseyenko, A. V., Drummond, A.  
 71 J., Gray, R. D., Suchard, M. A., & Atkinson, Q. D. (2012). Mapping the origins and  
 72 expansion of the indo-european language family. *Science*, 337(6097), 957–960.
- 73 Foote, M., Hunter, J. P., Janis, C. M., & Sepkoski, J. J. (1999). Evolutionary and preserva-  
 74 tional constraints on origins of biologic groups: Divergence times of eutherian mammals.  
 75 *Science*, 283(5406), 1310–1314. <https://doi.org/10.1126/science.283.5406.1310>
- 76 Harmon, L. J. (2019). *Phylogenetic comparative methods*. University of Idaho.
- 77 Huerta-Cepas, J., Serra, F., & Bork, P. (2016). ETE 3: Reconstruction, Analysis, and  
 78 Visualization of Phylogenomic Data. *Molecular Biology and Evolution*, 33(6), 1635–1638.  
 79 <https://doi.org/10.1093/molbev/msw046>
- 80 Maurits, L., Forkel, R., Kaiping, G. A., & Atkinson, Q. D. (2017). BEASTling: A software  
 81 tool for linguistic phylogenetics using BEAST 2. *PloS One*, 12(8).

- 82 Noutahi, M.-R. (2017). *How to simulate a phylogenetic tree?* [https://mrnoutahi.com/2017/](https://mrnoutahi.com/2017/12/05/How-to-simulate-a-tree/)  
83 [12/05/How-to-simulate-a-tree/](https://mrnoutahi.com/2017/12/05/How-to-simulate-a-tree/)
- 84 Paradis, E., & Schliep, K. (2018). Ape 5.0: An environment for modern phylogenetics and  
85 evolutionary analyses in R. *Bioinformatics*, 35, 526–528.
- 86 Pennell, M., Eastman, J., Slater, G., Brown, J., Uyeda, J., FitzJohn, R., Alfaro, M., & Harmon,  
87 L. (2014). Geiger v2.0: An expanded suite of methods for fitting macroevolutionary models  
88 to phylogenetic trees. *Bioinformatics*, 30, 2216–2218.
- 89 Robinson, P. (2016). The digital revolution in scholarly editing. *Ars Edendi Lecture Series*, 4,  
90 181–207.
- 91 Stadler, T. (2011). Simulating trees with a fixed number of extant species. *Systematic*  
92 *Biology*, 60(5), 676–684. <https://doi.org/10.1093/sysbio/syr029>

DRAFT