

1 babette: BEAUti 2, BEAST2 and Tracer for R

2 Richèl J.C. Bilderbeek¹ and Rampal S. Etienne¹



3 ¹Groningen Institute for Evolutionary Life Sciences, University of
4 Groningen, Groningen, The Netherlands

5 February 26, 2018

6 Summary

7 **1.** In the field of phylogenetics, BEAST2 is one of the most widely
8 used software tools. It comes with the graphical user interfaces BEAUti
9 2, DensiTree and Tracer, to create BEAST2 configuration files and to in-
10 terpret BEAST2's output files. However, when many different alignments
11 or model setups are required, a workflow of graphical user interfaces is
12 cumbersome.

13 **2.** Here, we present a free, libre and open-source package, **babette**:
14 'BEAUti 2, BEAST2 and Tracer for R', for the R programming language.
15 **babette** creates BEAST2 input files, runs BEAST2 and parses its results,
16 all from an R function call.



17 **3.** We describe **babette**'s usage and  the novel functionality it provides
18 compared to the original tools and  some examples.

19 **4.** As **babette** is designed to be of high quality and extendable, we
20 conclude by describing the further development of the package.
21

22 **Keywords:** computational biology, evolution, phylogenetics, BEAST2, R

23 1 Introduction

24 Phylogenies are commonly used to explore evolutionary hypotheses. Not only
25 can phylogenies show us how species (or other evolutionary units) relate to
26 each other, but we also estimate relevant parameters such as extinction and
27 speciation rates.

28  There are many phylogenetics tools available to obtain an estimate of the
29 phylogenetic tree of a given set of species. BEAST2 [8] is one of the most widely
30 used ones. It creates a posterior of jointly-estimated phylogenies and model
31 parameters, from one or more DNA, RNA or  amino acid alignments (see figure
32 1 for an overview of the workflow). It has been a graphical and a command-line
33 interface, that both need a configuration file containing alignments and model
34 parameters.

BEAST2 is bundled with BEAUti 2 [11] ('BEAUti' from now on), a desktop application to create a BEAST2 configuration file. BEAUti has a user-friendly graphical user interface, with helpful and readable default settings. As such, BEAUti is an attractive alternative to manual and error-prone editing of BEAST2 configuration files.

However, BEAUti cannot be called from a command-line script. This implies that when the user wants to explore the consequences of various settings, this must be done manually. This is the common workflow when using a few alignments and doing a superficial analysis of sensitivity of the reconstructed tree to model settings. For exploring many trees (for instance from simulations) and for more thorough sensitivity analysis, one would like to loop through multiple (simulated) alignments, nucleotide substitution models, clock models and tree priors. One such tool to replace BEAUti is **BEASTmaster** [17], which focuses on morphological traits and tip-dating, but also supports DNA data. **BEASTmaster**, however, takes hundreds of lines to R code and a Microsoft Excel file.

BEAST2 is also associated with Tracer [21] and DensiTree [7]. Both are desktop applications to analyze the output of BEAST2, each with a user-friendly graphical user interface. Tracer's purpose is to analyze the parameter estimates generated from a BEAST 2 run. It shows, among others, the effective sample size (ESS) and time series ('the trace', hence the name) of each variable in the MCMC run. Both ESS and trace are needed to assess the strength of the inference. DensiTree visualizes the phylogenies of a BEAST2 posterior, with many options to improve the display of many phylogenies.

However, for exploring the output of many BEAST2 runs, one would like a script to collect all parameters' ESSes, parameter traces and posterior phylogenies. There is no single package that offers a complete solution, but examples of R packages that offer a partial solution are **BEASTmaster**, **rBEAST** [22] and **RBeast** [13].

Here, we present **babette**: BEAUti 2, BEAST2 and Tracer for R, which creates BEAST2 configuration files, runs BEAST2, and analyzes its results, all from an R function call. This will save time, tedious mouse clicking and reduces the chances of errors in such repetitive actions. The interface of **babette** mimics the tools it is based on. This familiarity helps both beginner and experienced BEAST2 users to make the step from those tools to **babette**. **babette** enables the creation of a single-script pipeline from sequence alignments to posterior analysis in R.

babette is the first R package that unifies the full workflow of working with BEAST2. Whilst **BEASTmaster** needs hundreds of lines to create a BEAST2 configuration file, with **babette** the same can be created by a simple one-liner. Unlike **BEASTmaster**, **rBEAST** and **RBeast**, **babette** implements all of Tracer's primary functionality. There exists no package that calls BEAST2 from R. **babette** sets up BEAST2, runs it, and analyzes its results, making it a complete solution for using BEAST2 from R.

78 2 Description

79 **babette** is written in the R programming language [20] and enables the full
80 BEAST2 workflow from an R function call, in a similar way to what BEAUti,
81 DensiTree and Tracer do. **babette**'s main function is `run_beast2`, which con-
82 figures BEAST2, runs it and parses its output. `run_beast2` needs at least the
83 name of a FASTA file containing a DNA alignment. The default settings for
84 the other arguments of `run_beast2` are identical to BEAUti's and BEAST2's
85 default settings. Per alignment, a site model, clock model and tree prior can
86 be chosen. Multiple alignments can be used, each with its own (unlinked) site
87 model, clock model and tree prior.

88 **babette** currently has 61 exported functions to set up a BEAST2 configura-
89 tion file. **babette** is an alternative for a majority of BEAUti use cases, **but**
90 **does not yet support the full functionality of BEAUti**. Because of BEAUti's
91 high number of plugins, **babette** uses a software architecture that ex-
92 tended. Furthermore, **babette** has 7 exported function to run and help run
93 BEAST2. One function is used to run BEAST2, others allow the user to check
94 if a BEAST2 configuration file is indeed valid. Finally, **babette** has 20 exported
95 function to parse the BEAST2 output files and analyze the created posterior.
96 **babette** gives the same ESSes and summary statistics as Tracer. The data is
97 formatted as such, that it can easily be visualized using `ggplot2` (for a trace,
98 similar to Tracer) or `phangorn` [23] (for the phylogenies in a posterior, similar
99 to DensiTree).

100 Currently, **babette** does not replace all functionality in BEAUti, as it does
101 not provide 3 out of 7 tree priors, nor does it support RNA alignments or
102 use of morphological data. The many plug-ins of BEAUti are not yet sup-
103 ported by **babette**. **babette** does not support all command-line arguments of
104 BEAST2, does not provide the more specialized Tracer analysis options, nor is it
105 as feature-rich in plotting options as DensiTree. Up until now, the **babette** fea-
106 tures implemented are those requested by users. Further extension of **babette**
107 will be based on future user requests.

108 3 Usage

109 In R, the functions of a package need to be loaded in the global namespace first:

```
110 library(babette)
```

111 BEAUti, and likewise **babette**, needs at least a FASTA filename to produce a
112 BEAST2 configuration file. In BEAUti, this is achieved by loading a FASTA file,
113 then saving an output file using a common save file dialog. After this, BEAST2
114 needs to be applied to the created configuration file. It creates multiple files
115 storing the posterior. These output files must be parsed by either Tracer or
116 DensiTree. In **babette**, all this is achieved by:

```
117 out <- run_beast2("anthus_aco.fas")
```

118 This code will create a (temporary) BEAST2 configuration file, from the
 119 FASTA file with name `anthus_aco.fas`, using the same default settings as
 120 BEAUti, which are, among others, a Jukes-Cantor site model, a strict clock,
 121 and a Yule birth tree prior. `babette` will then execute BEAST2 using that
 122 file, and parses the output. The returned data structure, named `out`, is a
 123 list of parameter estimates (called `estimates`), posterior phylogenies (called
 124 `anthus_aco_trees`), and MCMC operator performance (`operators`). An ex-
 125 ample of using a different site model, clock model and tree prior is:

```
126 out <- run_beast2(  
127   "anthus_aco.fas",  
128   site_models = create_hky_site_model(),  
129   clock_models = create_rln_clock_model(),  
130   tree_priors = create_bd_tree_prior()  
131 )
```

132 This code uses an HKY site model, a relaxed log-normal clock model and
 133 a birth-death tree prior, each with their default settings in BEAUti. Table
 134 1 shows an overview of all functions to create site models, clock models and
 135 tree priors. Note that the arguments' names `site_models`, `clock_models` and
 136 `tree_priors` are plural, as each of these can be (a list of) one or more elements.
 137 Each of these arguments must have the same number of elements, so that each
 138 alignment has its own site model, clock model and tree prior. An example of
 139 two alignments, each with its own site model, is:


```
140 out <- run_beast2(  
141   c("anthus_aco.fas", "anthus_nd2.fas"),  
142   site_models = list(  
143     create_tn93_site_model(),  
144     create_gtr_site_model()  
145   )  
146 )
```


147 `babette` also uses the same default prior distributions as BEAUti for each
 148 of the site models, clock models and tree priors. For example, by default,
 149 a Yule tree prior assumes that the birth rate follows a uniform distribution,
 150 from minus infinity to plus infinity. This assumption implies that negative and
 151 positive birth rates are just as likely, where a negative birth rate is biologically
 152 impossible (note that in practice, this usually works out just fine). One may
 153 prefer an exponential distribution instead, as this would assume only positive
 154 birth rates, and makes high birth rates unlikely. To do this in `babette`:

```
155 out <- run_beast2(  
156   "anthus_aco.fas",  
157   tree_priors = create_yule_tree_prior(  
158     birth_rate_distr = create_exp_distr()  
159   )  
160 )
```


161 Within this same example, one may specify the initial shape parameters
 162 of the exponential distribution. In BEAST2's implementation, an exponential

163 distribution has one shape parameter: its mean, which can be set to any value
 164 with BEAUti. Within **babette**, to set the initial mean value of the exponential
 165 distribution, do:

```
166 out <- run_beast2(  
167   "anthus_aco.fas",  
168   tree_priors = create_yule_tree_prior(  
169     birth_rate_distr = create_exp_distr(  
170       mean =  rate_mean_param(value = 1.0)  
171     )  
172   )  
173 )
```

174 Our initial motivation to create **babette** was that we wanted to fix the crown
 175 age of a phylogeny. BEAUti assumes that a phylogeny has a crown age that
 176 needs to be jointly-estimated with the phylogeny and other parameters. It does
 177 not allow for fixing the crown age. Without **babette**, one needs to manually edit
 178 the BEAST2 configuration file, which is tedious and prone to errors.  Making the
 179 crown ages is especially useful for theoretical experiments, as this allows for one
 180 less source of variation. This is how to specify a fixed crown age with **babette**:

```
181 out <- run_beast2(  
182   "anthus_aco.fas",  
183   posterior_crown_age = 15  
184 )
```

185 **babette** allows for the same functionality as Tracer. Tracer works on the
 186 values of the parameter estimates sampled in the BEAST2 run. This is called
 187 the "trace" (hence the name). The start of the trace is usually discarded, as
 188 an MCMC algorithm (such as used by BEAST2) first has to converge to its
 189 equilibrium. The start of the trace, called the "burn-in", will be removed, as its
 190 parameter estimates are not representative. By default, Tracer discards the first
 191 10% of all the parameter estimates. To remove the burn-in from all parameter
 192 estimates in **babette**, the following code can be used: 

```
193 traces <- remove_burn_ins(out$estimates)
```

194 Tracer shows the ESSes of each posterior's variables. These ESSes are im-
 195 portant to determine the strength of the inference. As a rule of thumb, an
 196 ESS of 200 is acceptable for any parameter estimate. To calculate the effective
 197 sample sizes (of all estimated variables) in **babette**:

```
198 esses <- calc_esses(  
199   traces,  
200   sample_interval = 1000  
201 )
```

202 Tracer displays multiple summary statistics for each estimated variable: the
 203 mean and its standard error, standard deviation, variance, median, mode, geo-
 204 metric mean, 95% highest posterior density interval, auto-correlation time and
 205 effective sample size. It displays these statistics per variable. In **babette**, these
 206 summary statistics are collected for all estimated parameters at once:

Name	Description
<code>run_beast2</code>	Run BEAST2
<code>create_gtr_site_model</code> <code>create_hky_site_model</code> <code>create_jc69_site_model</code> <code>create_tn93_site_model</code>	Create a GTR site model Create an HKY site model Create a Jukes-Cantor site model Create a TN93 site model
<code>create_rln_clock_model</code> <code>create_strict_clock_model</code>	Create a relaxed log-normal clock model Create a strict clock model
<code>create_bd_tree_prior</code> <code>create_cbs_tree_prior</code> <code>create_ccp_tree_prior</code> <code>create_cep_tree_prior</code> <code>create_yule_tree_prior</code>	Create a birth-death tree prior Create a coalescent Bayesian skyline tree prior Create a coalescent constant-population tree prior Create a coalescent exponential-population tree prior Create a Yule tree prior
<code>create_beta_distr</code> <code>create_exp_distr</code> <code>create_gamma_distr</code> <code>create_inv_gamma_distr</code> <code>create_laplace_distr</code> <code>create_log_normal_distr</code> <code>create_normal_distr</code> <code>create_one_div_x_distr</code> <code>create_poisson_distr</code> <code>create_uniform_distr</code>	Create a beta distribution Create an exponential distribution Create a gamma distribution Create an inverse gamma distribution Create a Laplace distribution Create a log-normal distribution Create a normal distribution Create a 1/X distribution Create a Poisson distribution Create a uniform distribution

Table 1: babette’s main functions

```

207 sum_stats <- calc_summary_stats(
208   traces,
209   sample_interval = 1000
210 )

```


211 `babette` allows for the same functionality as `DensiTree`. `DensiTree` displays
212 the phylogenies in a posterior at the same time scale, drawn one over one an-
213 other, allowing to see the uncertainty in topology and branch lengths. Within
214 the object `out`, the posterior phylogenies are stored as `anthus_aco_trees`, and
215 can be plotted as such:

```

216 plot_densitree(out$anthus_aco_trees)

```

217 4 babette resources

218 `babette` is free, libre and open source software available from the official R pack-
219 age archive at <http://cran.r-project.org/src/contrib/PACKAGES.html#babette>
220 and is licensed under the GNU General Public License v3.0.
221  `babette` uses the Travis CI [2] continuous integration service, which is known
222 to significantly increase the number of bugs exposed [24] and increases the speed

at which new features are added [24]. **babette** has a 100% code coverage, which correlates with code quality [16, 10]. **babette** follows Hadley Wickham’s style guide [25], which improves software quality [12]. **babette** depends on multiple packages, which are **AR** [8], **beautier** [5], **beastier** [4], **devtools** [28], **geiger** [15], **ggplot2** [26], **knitr** [30], **phangorn** [23], **rmarkdown** [3], **seqinr** [9], **stringr** [27], **testit** [29] and **tracerer** [6].

babette’s development takes place on GitHub [1], <https://github.com/richelbilderbeek/babette>, which accommodates collaboration [19] and improves transparency [14]. **babette**’s GitHub facilitates feature requests and has guidelines how to do so.

babette’s documentation is extensive. All functions are documented in the package’s internal documentation. For quick use, each exported function shows a minimal example. For easy exploration, each exported function’s documentation links to related functions. Additionally, **babette** has a vignette that demonstrates extensively how to use it. The GitHub documentation helps to get started, with a dozen examples of BEAUti screenshots with equivalent **babette** code.

5 Citation of babette

Scientists using **babette** in a published paper can cite this article, and/or cite the **babette** package directly. To obtain this citation from within an R script, use:

```
> citation("babette")
```

6 Acknowledgements

Thanks to Yacine Ben Chehida and Paul van Els for supplying their BEAST2 use cases. Thanks to Tonel Herrera-Alsina, Raphael Scherrer and Giovanni Laudanno for their contributions to this package and article.

References

- [1] Github. <https://github.com/>.
- [2] Travis CI. <https://travis-ci.org/>.
- [3] JJ Allaire, Yihui Xie, Jonathan McPherson, Javier Luraschi, Kevin Ushey, Aron Atkins, Hadley Wickham, Joe Cheng, and Winston Chang. *rmarkdown: Dynamic Documents for R*, 2017. R package version 1.8.
- [4] Richel J.C. Bilderbeek. *beastier: BEAST2 from R*. R package version 1.0, <https://github.com/richelbilderbeek/beastier>.

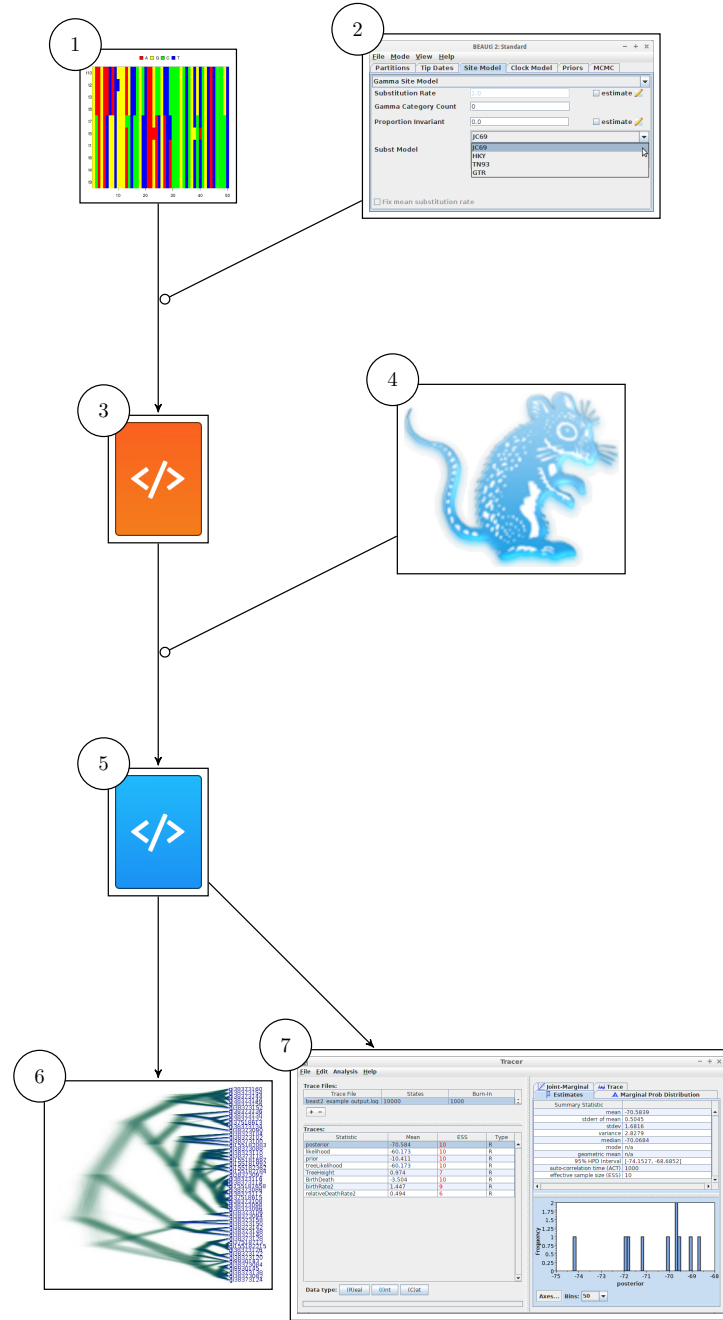


Figure 1: Workflow using GUI tools. From an alignment (1) and BEAUti (2), a BEAST2 configuration file (3) is created. BEAST2 (4) uses that file to infer a posterior, storing it in multiple files (5). These results are visualized using DensiTree (6) and Tracer (7). babette allows for the same workflow, all from an R function call.

- [5] Richel J.C. Bilderbeek. *beautier: BEAUti 2 from R*. R package version 1.0, <https://github.com/richelbilderbeek/beautier>.
- [6] Richel J.C. Bilderbeek. *tracerer: Tracer from R*. R package version 1.0, <https://github.com/richelbilderbeek/tracerer>.
- [7] Remco Bouckaert and Joseph Heled. Densitree 2: Seeing trees through the forest. *bioRxiv*, page 012401, 2014.
- [8] Remco Bouckaert, Joseph Heled, Denise Kühnert, Tim Vaughan, Chieh-Hsi Wu, Dong Xie, Marc A Suchard, Andrew Rambaut, and Alexei J Drummond. Beast 2: a software platform for bayesian evolutionary analysis. *PLoS Comput Biol*, 10(4):e1003537, 2014.
- [9] D. Charif and J.R. Lobry. SeqinR 1.0-2: a contributed package to the R project for statistical computing devoted to biological sequences retrieval and analysis. In U. Bastolla, M. Porto, H.E. Roman, and M. Vendruscolo, editors, *Structural approaches to sequence evolution: Molecules, networks, populations*, Biological and Medical Physics, Biomedical Engineering, pages 207–232. Springer Verlag, New York, 2007. ISBN : 978-3-540-35305-8.
- [10] Fabio Del Frate, Praerit Garg, Aditya P Mathur, and Alberto Pasquini. On the correlation between code coverage and software reliability. In *Software Reliability Engineering, 1995. Proceedings., Sixth International Symposium on*, pages 124–132. IEEE, 1995.
- [11] Alexei J Drummond, Marc A Suchard, Dong Xie, and Andrew Rambaut. Bayesian phylogenetics with beauti and the beast 1.7. *Molecular biology and evolution*, 29(8):1969–1973, 2012.
- [12] Xuefen Fang. Using a coding standard to improve program quality. In *Quality Software, 2001. Proceedings. Second Asia-Pacific Conference on*, pages 73–78. IEEE, 2001.
- [13] Nuno Faria and Marc A. Suchard. *RBeast*.
- [14] Krzysztof J Gorgolewski and Russell Poldrack. A practical guide for improving transparency and reproducibility in neuroimaging research. *bioRxiv*, page 039354, 2016.
- [15] LJ Harmon, JT Weir, CD Brock, RE Glor, and W Challenger. Geiger: investigating evolutionary radiations. *Bioinformatics*, 24:129–131, 2008.
- [16] Joseph R. Horgan, Saul London, and Michael R Lyu. Achieving software quality with testing coverage measures. *Computer*, 27(9):60–69, 1994.
- [17] Nicholas J. Matzke. *BEASTmasterR: R tools for automated conversion of NEXUS data to BEAST2 XML format, for fossil tip-dating and other uses*. R package version 0.2.

- 294 [18] E. Paradis, J. Claude, and K. Strimmer. APE: analyses of phylogenetics
295 and evolution in R language. *Bioinformatics*, 20:289–290, 2004.
- 296 [19] Yasset Perez-Riverol, Laurent Gatto, Rui Wang, Timo Sachsenberg, Julian
297 Uszkoreit, Felipe Leprevost, Christian Fufezan, Tobias Ternent, Stephen J
298 Eglén, Daniel SS Katz, et al. Ten simple rules for taking advantage of git
299 and github. *bioRxiv*, page 048744, 2016.
- 300 [20] R Core Team. *R: A Language and Environment for Statistical Computing*.
301 R Foundation for Statistical Computing, Vienna, Austria, 2013.
- 302 [21] Andrew Rambaut and Alexei J Drummond. *Tracer v1.4*, 2007. Available
303 from <http://beast.bio.ed.ac.uk/Tracer>.
- 304 [22] Oliver Ratmann. *rBEAST*.
- 305 [23] K.P. Schliep. phangorn: phylogenetic analysis in r. *Bioinformatics*,
306 27(4):592–593, 2011.
- 307 [24] Bogdan Vasilescu, Yue Yu, Huaimin Wang, Premkumar Devanbu, and
308 Vladimir Filkov. Quality and productivity outcomes relating to contin-
309 uous integration in github. In *Proceedings of the 2015 10th Joint Meeting*
310 *on Foundations of Software Engineering*, pages 805–816. ACM, 2015.
- 311 [25] Hadley Wickham. Style guide. <http://r-pkgs.had.co.nz/style.html>.
312 Accessed: 2017-12-20.
- 313 [26] Hadley Wickham. *ggplot2: elegant graphics for data analysis*. Springer
314 New York, 2009.
- 315 [27] Hadley Wickham. *stringr: Simple, Consistent Wrappers for Common*
316 *String Operations*, 2017. R package version 1.2.0.
- 317 [28] Hadley Wickham and Winston Chang. *devtools: Tools to Make Developing*
318 *R Packages Easier*, 2016. R package version 1.12.0.9000.
- 319 [29] Yihui Xie. *testit: A Simple Package for Testing R Packages*, 2014. R
320 package version 0.4, <http://CRAN.R-project.org/package=testit>.
- 321 [30] Yihui Xie. *knitr: A General-Purpose Package for Dynamic Report Gener-*
322 *ation in R*, 2017. R package version 1.17.