



# AsymmeTree User Manual

David Schaller

[sdavid@bioinf.uni-leipzig.de](mailto:sdavid@bioinf.uni-leipzig.de)

# Contents

<b>1</b>	<b>Introduction</b>	<b>2</b>
<b>2</b>	<b>Installation</b>	<b>2</b>
2.1	Easy Installation with pip . . . . .	2
2.2	Installation with the setup file . . . . .	2
2.3	Dependencies . . . . .	2
<b>3</b>	<b>Usage</b>	<b>3</b>
3.1	Tree Data Structures . . . . .	3
3.2	Simulation of Species and Gene Trees . . . . .	4
3.2.1	Species Trees . . . . .	4
3.2.2	Gene Trees . . . . .	5
3.2.3	Assignment of Variable Evolution Rates . . . . .	6
3.2.4	Distance Matrix and Noise . . . . .	7
3.3	Simulation of Sequences . . . . .	8
3.3.1	Substitution Model . . . . .	8
3.3.2	Indel Model . . . . .	9
3.3.3	Heterogeneity Model . . . . .	10
3.3.4	The Class <code>Evolver</code> . . . . .	10
3.4	Simulation of Genomes . . . . .	12
3.5	Best Match Inference . . . . .	13
3.6	Supertree Computation . . . . .	14
3.7	Cograph Editing and ParaPhylo . . . . .	14
<b>A</b>	<b>Subpackages and Modules</b>	<b>19</b>

# 1 Introduction

AsymmeTree is an open-source Python library for the simulation and analysis of phylogenetic scenarios. It includes a simulator for species and gene trees with asymmetric evolution rates, tools for the inference and analysis of phylogenetic best matches [14, 15] (resp. best hits) from known gene trees or evolutionary distances. Moreover, it includes an algorithm to compute supertrees [11] and a method to estimate rooted species trees from an ensemble of orthology/paralogy relations [19].

The library, and especially the simulator, was primarily designed to be able to validate mathematical concepts and test inference methods for various steps on the way to more realistically available data, i.e., dated gene trees, additive distances of gene sets, noisy distances and finally sequences. Both nucleotide and amino acid sequence simulation with or without indels are supported. In both cases, several substitution models are available.

The software is hosted on [GitHub](#) and also available via The Python Package Index (PyPI). Please feel free to report bugs or make suggestions for improvement in the [Issues](#) section of the GitHub repository.

If you use AsymmeTree in your project or code from it, please cite:

Peter F. Stadler, Manuela Geiß, David Schaller, Alitzel López Sánchez, Marcos González Laffitte, Dulce I. Valdivia, Marc Hellmuth, Maribel Hernández Rosales (2020). **From pairs of most similar sequences to phylogenetic best matches.** *Algorithms for Molecular Biology*. doi: 10.1186/s13015-020-00165-2. [35].

# 2 Installation

AsymmeTree requires Python 3.5 or higher. Python 2 is not supported.

## 2.1 Easy Installation with pip

The `asymmetree` package is available on The Python Package Index (PyPI):

```
pip install asymmetree
```

For details about how to install Python packages see [here](#).

## 2.2 Installation with the setup file

Alternatively, you can download or clone the repo, go to the root folder of package and install it using the command:

```
python setup.py install
```

## 2.3 Dependencies

AsymmeTree has several dependencies (which are installed automatically when using `pip` or the `setup.py`):

- [NetworkX](#)

- [SciPy and NumPy](#)
- [Matplotlib](#)

To use the tree reconstruction method for best match inference and the C++ implementation of the quartet method [35], the software [RapidNJ](#) [34], resp., [qinfer](#) must be installed. I recommend that you compile these tools on your machine, place the binaries into a persistent location and add this location to your PATH environment variable. Both tools do *not* need to be installed for most functionalities including the simulator.

## 3 Usage

AsymmeTree is divided into several subpackages and modules, an overview of which is given in Appendix A. The library interface functions are described in the following sections and can be imported directly from the respective subpackage (see examples).

The term ‘color’ regularly appears in the library and refers to the reconciliation of gene trees with species trees. In particular, the terms ‘color’ and ‘species’/‘genome’ in which a gene resides are used interchangeably. The emphasizes the close connection to vertex-colored graphs such as (colored) Best Match Graphs [14].

### 3.1 Tree Data Structures

The two classes **Tree** and **PhyloTree** (inherits from **Tree**) implement tree data structures which are essential for most of the modules in the package. The latter contains converters and parsers for the Newick format and a NetworkX graph format.

The vertices of a **PhyloTree** instance are of type **PhyloTreeNode** and contain the following attributes:

<b>ID</b>	vertex ID ( <b>int</b> )
<b>label</b>	label ( <b>str</b> ), in gene trees: "S" for speciation, "D" for duplication, "H" for horizontal gene transfer, "*" for loss
<b>color</b>	only gene trees; species in which the gene resides, i.e., ID of some vertex in a species tree, <b>int</b> for extant genes, can be of type <b>tuple</b> (of two <b>ints</b> ) for inner and loss vertices
<b>tstamp</b>	time stamp of the event ( <b>double</b> )
<b>dist</b>	evolutionary distance or divergence time from the parent vertex ( <b>double</b> )
<b>tranferred</b>	only gene trees; indicates whether the edge from the parent is the transfer edge from an HGT event; 1 if yes and 0 otherwise

Both species and gene trees can be converted into Newick format using the function `to_newick()` of the **PhyloTree** class. In case of a gene tree, the color is represented in brackets, e.g.

```
>>> '(3<1>:0.534,2<2>:0.762)S<0>:0.273'
```

To suppress this, use `to_newick(color=False)`. Likewise, to suppress the distances, you can use `to_newick(distance=False)`. The function `PhyloTree.parse_newick()` can

handle this customized format as well as the standard Newick format.

Moreover, phylogenetic trees can easily be serialized in `json` or `pickle` (Python's serialization library) format:

```
# tree is of type 'PhyloTree'
tree.serialize('tree.pickle')
# or
tree.serialize('tree.json')
```

By default, the serialization format is inferred from the file extension. Alternatively, it can be specified as keyword argument, e.g. `mode="json"`. To load a tree that was serialized this way, use:

```
from asymmetree import PhyloTree

tree = PhyloTree.load('tree.json')
```

## 3.2 Simulation of Species and Gene Trees

The subpackage `asymmetree.treeevolve` contains modules for the simulation of dated species and gene trees. In terms of divergence time, these trees define an ultrametric on the set of their (extant) leaves. Gene trees, furthermore, can be manipulated with a realistic rate heterogeneity among their branches resulting in general additive distances.

### 3.2.1 Species Trees

The function `simulate_species_tree(N)` simulates a dated species tree with `N` leaves (i.e. recent species) using the specified model. The following models are available:

"innovation"	Innovation model Keller-Schmidt and Klemm [25], if not specified the divergence time between the (planted) root and the leaves will be normalized to unity
"yule"	standard Yule model [43], default birth rate is 1.0
"BDP"	constant-rate birth-death process [see e.g. 26, 18], default birth rate is 1.0 and death rate is 0.0
"EBDP"	episodic birth-death process, algorithm of [36]

The following keyword parameters (with their default value) are available:

<code>model="innovation"</code>	model for the species tree simulation, currently only the 'innovation model' is available
<code>non_binary_prop=0.0</code>	probability that an inner edge is contracted, results in a non-binary tree
<code>planted=True</code>	add a planted root that has the first true speciation node as its single neighbor, this way duplication (and loss) events can occur before the first speciation event in a subsequent gene tree simulation
<code>remove_extinct=False</code>	remove all branches leading to losses, only relevant for models with death events

<code>rescale_to_height=</code>	specify the divergence time between the (planted) root and the
<code>None</code>	leaves i.e. the final height of the dated tree

For any model, the root of the resulting tree has the maximal time stamp and all (extant) species have time stamp 0.0. The episodes of the "EBDP" model must be supplied as a list of tuples/lists where each episode has the structure

```
(birthrate, deathrate, proportion_of_survivors, time_stamp).
```

Note that the first elements in this list correspond to the most recent ones, and that the first episode should have a time stamp of 0.0. Example usage:

```
import asymmetree.treeevolve as te

S1 = te.simulate_species_tree(10, planted=True, non_binary_prob=0.2)
print(S1.to_newick())

S1 = te.simulate_species_tree(10, model='EBDP',
                              episodes=[(1.0, 0.3, 0.8, 0.0),
                                         (0.9, 0.4, 0.6, 0.3)])
print(S2.to_newick())
```

### 3.2.2 Gene Trees

Dated gene trees are simulated along a given species tree *S* with a variant of the birth-death, resp., Gillespie algorithm [26, 17]. To this end, use the class `GeneTreeSimulator` or the function `simulate_dated_gene_tree(S, **kwargs)` with a species tree of type `PhyloTree`. The following parameters are available:

<code>DLH_rates=</code> <code>(0.0, 0.0, 0.0)</code>	a tuple of three floats, rates for duplication, loss and horizontal gene transfer (HGT) events
<code>dupl_polytomy=0.0</code>	allows non-binary duplication events by specifying the parameter $\lambda$ of a Poisson distribution (copy number = 2 + drawn number)
<code>prohibit_extinction=</code> <code>"per_species"</code>	avoid loss events for genes that are the last survivor in their species branch (" <code>per_species</code> "), the last survivor of the whole family (" <code>per_family</code> "); or no constraints ( <code>False</code> )

For the constraints to avoid extinction, the loss rate in the respective branches are temporarily set to zero. Example usage:

```
import asymmetree.treeevolve as te

# S is a species tree of type PhyloTree
tree_simulator = te.GeneTreeSimulator(S)
tree = tree_simulator.simulate(DLH_rates=(1.0, 0.5, 0.1))

# or
tree = te.simulate_dated_gene_tree(S, DLH_rates=(1.0, 0.5, 0.1))
```

The function `observable_tree(tree)` returns the observable part of a gene tree, i.e., it copies the tree, removes all branches that lead to loss events only and suppresses all

inner nodes with only one child. It also removes the planted root. Example usage:

```
# observable gene tree
ogt = ts.observable_tree(tree)
```

### 3.2.3 Assignment of Variable Evolution Rates

The module `EvolutionRates` contains functions to model realistic (asymmetric) evolution rates for a given gene tree. Moreover, correlation of the evolution rates between genes of the same (and closely related) species is introduced (autocorrelation, [28]). The function `assign_rates(T, S)` takes a gene tree `T` and the **corresponding** species tree `S` as input, and manipulates the branch length of the gene tree. The following keyword parameters (with their default values) are available:

<code>base_rate=1.0</code>	starting value for the substitution rate (per time unit) and expected value for conserved genes
<code>autocorr_factors=None</code>	a dictionary containing rate factors for the edges of the species tree
<code>autocorr_variance=0.0</code>	variance factor for a lognormal distribution that controls autocorrelation between genes of the same (and closely related) species, the higher the lower the autocorrelation; only relevant if <code>autocorr_factors</code> is not directly supplied
<code>gamma_param=(0.5, 1.0, 2.2)</code>	parameter the for Gamma distribution ( <code>a</code> , <code>loc</code> , <code>scale</code> ) from which rate factors for divergent are drawn, the default values are chosen to fit observed asymmetries between paralogs in yeast data [3]
<code>CSN_weights=(1, 1, 1)</code>	weights for choice between conservation, subfunctionalization and neofunctionalization after a duplication event
<code>inplace=True</code>	manipulate edge lengths ( <code>dist</code> ) of the gene tree in-place, otherwise copy the tree

It is recommended to apply the rate assignment to the true gene tree that still contains loss events. Note that the rates are used to manipulate the `dist` attributes in the gene tree and not returned explicitly. Example usage:

```
import asymmetree.treeevolve as te

S = te.simulate_species_tree(10)

# true gene tree (with losses)
tree = te.simulate_dated_gene_tree(S, DLH_rates=(0.5, 0.5, 0.5))
te.assign_rates(tree, S, base_rate=1,
                autocorr_variance=0.2,
                gamma_param=(0.5, 1.0, 2.2))

# or in a single step
tree = simulate_gene_trees(S, N=1,
                           DLH_rates=(0.5, 0.5, 0.5)
                           base_rate_distr=('constant', 1.0),
                           autocorr_variance=0.2,
                           gamma_param=(0.5, 1.0, 2.2))
```

```
# observable gene tree
ogt = ts.observable_tree(tree)
```

The function `simulate_gene_trees(S)` combines the simulation of dated gene trees and the rate assignment into one step. If `N=1`, a single gene tree is returned. Otherwise, a list of gene trees is returned that shared the same rate factors for the branches in the species tree (autocorrelation factors) in the rate assignment procedure. Moreover, a distribution for the base rate (assigned the planted edge of the gene tree) can be specified with the parameter `base_rate_distr`. Available are e.g. ("`constant`", `rate`) and ("`gamma`", `shape`, `scale`).

### 3.2.4 Distance Matrix and Noise

Distances derived from (real-life) gene or protein sequences are always burdened with noise. Such data can either be modeled by simulating sequences, or by disturbing the distances specified by a given tree directly. The latter alternatively is described briefly in this section.

The additive (i.e. noiseless) distance from an **observable** gene tree can be computed using the function `distance_matrix()` of a `PhyloTree` instance. It returns a tuple containing a list of leaves in the tree (corresponding to the row/column order) and the distance matrix as a 2-dimensional `numpy` array.

```
# tree is an observable gene tree
leaves, D = tree.distance_matrix()
leaf_index_dict = {leaf: i for i, leaf in leaves}
```

In the next step, noise can be introduced into a distance matrix using the `NoisyMatrix` module. Random noise can be simulated with the function `noisy_matrix(orig_matrix, sd)`. The following parameters are available (keyword arguments are indicated by their default value):

<code>orig_matrix</code>	original matrix to be disturbed
<code>sd</code>	standard deviation of a normal distribution with mean 1 from which noise factors are drawn
<code>metric_repair="reject"</code>	method to ensure that the resulting distance matrix is still a metric, available are the rejection of noise steps that violate the metric property (" <code>reject</code> "), the decrease-only metric repair (" <code>DOMR</code> ") and the general metric repair (" <code>general</code> ") algorithm

Alternatively, the function `convex_linear_comb(D1, D2)` can be used to simulate systematically biased noise by computing a linear convex combination with a disturbance matrix. The function thus takes two distance matrices (`numpy` arrays) not necessarily of the same size as input and disturbs them with one another. The contribution of the respective disturbance matrix is controlled by the keyword parameter `alpha` (default is 0.5). If the keyword parameter `first_only` is `True`, only the first disturbed matrix is returned. Otherwise, both are returned in a tuple.



### 3.3 Simulation of Sequences

AsymmeTree supports the simulation of nucleic and amino acid sequences using time-continuous Markov models, as usually applied for this purpose [for textbooks, see e.g. 12, 41, 42]. The subpackage `asymmetree.seqevolve` contains the modules and functions for this task.

The class `Evolver` takes several model as parameters for its initialization:

- a **substitution** model (`SubstModel`, required) for the substitution of single bases or amino acids,
- an **indel** model (`IndelModel`, optional) for the simulation of insertions and/or deletion, and
- a model for rate **heterogeneity** (`HetModel`, optional) among the sites of the sequence under evolution

#### 3.3.1 Substitution Model

A substitution model usually comprises an exchangeability matrix  $S$  and a vector  $\pi$  containing the equilibrium frequencies of the alphabet  $A$  of nucleobases, amino acids et cetera. From this, the rate matrix  $Q$  can be computed as  $S\Pi$  where  $\Pi = \text{diag}\{\pi_1, \dots, \pi_{|A|}\}$  [41]. The substitution probability matrix, in turn, is given by

$$P = e^{Qt}$$

which is computed efficiently by AsymmeTree using matrix diagonalization.

The following models for nucleotide (**n**) and amino acid (**a**) substitution are currently available (codon models are not supported at the moment):

model	type	reference	required parameters ( <code>kwargs</code> )
JC96	n/a	Jukes & Cantor 1969 [24]	-
K80	n	Kimura 1980 [27]	<code>kappa</code> (transition/transversion rate ratio)
GTR	n	general time-reversible model (GTR) 1986 [37]	<code>abcdef</code> (list of rates (a) $C \leftrightarrow T$ , (b) $A \leftrightarrow T$ , (c) $G \leftrightarrow T$ , (d) $A \leftrightarrow C$ , (e) $C \leftrightarrow G$ , (f) $A \leftrightarrow G$ ); <code>f</code> (list of equilibrium frequencies $A/C/G/T$ )
DAYHOFF	a	Dayhoff 1978 [10]	-
BLOSUM62	a	BLOSUM62 1992 [20]	-
JTT	a	Jones, Taylor & Thornton 1992 [23]	-
WAG	a	Whelan & Goldman 2001 [38]	-
LG	a	Le & Gascuel 2008 [29]	-
CUSTOM	n/a	-	<code>filename</code> (path to a file with a model in PAML [40] format)

Note that a custom substitution model can be specified via `model_name="CUSTOM"`. In this case, the path to the model in PAML [40] format must be supplied. Moreover, the model type (n/a) must fit this model. Example usage:

```
from asymmetree.seqevolve import SubstModel

# non-empirical model
subst_model_jc69 = SubstModel('n', 'K80', kappa=2.0)

# empirical model
subst_model_jtt = SubstModel('a', 'JTT')

# custom model
subst_model_custom = SubstModel('a', 'CUSTOM', filename='path/to/custom.paml')
```

### 3.3.2 Indel Model

Insertions and deletions are modeled based on Dawg [4]. An indel model requires sitewise rates for insertion `insertion_rate` and deletion `deletion_rate`. Moreover, the following parameters are available (with default values):

<code>length_distr=</code> ("zipf", 1.821)	distribution of indel length, default value for zipf distribution cf. [5]
<code>min_length=1</code>	integer min. value at which the specified distribution is truncated, must be less than the expected value of the distribution, <code>None</code> means no limit
<code>max_length=None</code>	integer max. value at which the specified distribution is truncated, must be greater than the expected value of the distribution, <code>None</code> means no limit

The following models for indel length are available:

model	distribution	required parameters (kwargs, with default values)
zipf	Zipf	<code>a=1.821</code> (parameter <i>a</i> of the Zipf distribution)
negative_binomial	negative binomial	<code>r=1.0</code> and <code>q=0.5</code> (parameters <i>r</i> and <i>q</i> of the negative binomial distribution)

Example usage:

```
from asymmetree.seqevolve import IndelModel

# sitewise insertion rate 0.01 and deletion rate 0.008
indel_model = IndelModel(0.01, 0.008, length_distr=('zipf', 2.0))
```

### 3.3.3 Heterogeneity Model

Selective pressure usually varies among the sites of a sequence under evolution. To model this, rate factors  $r$  for single sites or groups of sites are commonly drawn from a Gamma distribution ('+ $\Gamma$ ') with mean 1 and parameter  $\alpha$  [4, 13, 9]. The rate matrix  $rQ$  is then used instead of  $Q$ . Note that smaller values for  $\alpha$  correspond to higher heterogeneity.

AsymmeTree supports two modes of the '+ $\Gamma$ '-model. You can specify a number of classes to which the sites are assigned randomly and uniformly distributed. Sites of the same class share a common factor  $r$ . The other possibility is a sitewise heterogeneity, i.e., every site has its own rate. In both cases, the rate or class membership is inherited from the parent sites during the evolution along a tree. Note that the sitewise mode is expected to have a longer running time.

An other aspect of among site heterogeneity is the modeling of invariant sites ('+ $I$ '), i.e., sites that never mutate at all as a result of very strong selective pressure. The (expected) proportion  $p$  of invariant sites can be specified by the user, and sites are assigned as "invariant" with probability  $p$ . Note that  $p > 0$  affects the overall substitution rate. In other words, the rates of the non-invariant sites are **not** adjusted to compensate the decreased number of expected substitution over all sites.

To summarize, the following parameters are available for the class `HetModel` (keyword arguments are indicated by their default value):

<code>alpha</code>	parameter $\alpha$ of the Gamma distribution
<code>classes=5</code>	number of classes; sites in the same class share the same rate factor
<code>sitewise=False</code>	if <code>True</code> , factors are drawn sitewise; the number of classes is ignored in this case
<code>invariant=0.0</code>	(expected) proportion $p$ of invariant sites

Note that the '+ $I$ '-model can be used without the '+ $\Gamma$ '-model by setting `classes=1` (the single class will have a factor of 1) and `invariant` to some proportion greater than 0. Example usage:

```
from asymmetree.seqevolve import HetModel

# +G
het_model_g = HetModel(1.5, classes=10)

# +G +I
het_model_gi = HetModel(2.0, sitewise=True, invariant=0.2)

# +I (value of alpha is irrelevant)
het_model_i = HetModel(1.0, classes=1, invariant=0.25)
```

### 3.3.4 The Class Evolver

The class `Evolver` evolves a sequence according to the specified models (see previous sections) along a phylogenetic tree. In `AsymmeTree`, the `dist` attribute of the vertex  $v$  in an edge  $uv$  of the tree ( $u$  is closer to the root) is always used as the **expected number of substitutions** along this edge. Thus, PAM distances as e.g. used optionally in [9] are

not supported. The `dist` attribute is also used as the duration of the Markov process in which insertions and deletions are drawn.

The following parameters are available for the initialization of an `Evolver` instance (keyword arguments are indicated by their default value):

<code>subst_model</code>	substitution model; instance of <code>SubstModel</code>
<code>indel_model=None</code>	model for insertions and deletions; instance of <code>IndelModel</code>
<code>het_model=None</code>	model for among site heterogeneity and invariant sites; instance of <code>HetModel</code>
<code>jump_chain=False</code>	if <code>True</code> , an alternative Gillespie-like [17] algorithm is applied for the substitution process instead of the computation of $P = e^{Qt}$

Once the `Evolver` is initialized, its function `evolve_along_tree()` can be called to evolve a sequence along a tree. The following parameters are available for this function (keyword arguments are indicated by their default value):

<code>tree</code>	phylogenetic tree; instance of <code>PhyloTree</code>
<code>start_length=200</code>	length of the root sequence which is randomly drawn from the equilibrium frequencies in the specified substitution model
<code>start_seq=None</code>	root sequence ( <code>str</code> ); must be compatible with the specified substitution model ( <code>model_type="n"/"a"</code> ); if supplied, the <code>start_length</code> attribute is ignored

The sequences of a simulation run are returned by this function (and also available via the attribute `sequences` as long as the function has not been called again), which is a `dict` containing the nodes (inner and leaf nodes) as keys and instances of type `EvoSeq` as values. The latter can be converted into `str` using `subst_model.to_sequence(evo_seq)`.

The function `true_alignment()` can be used to compute (and optionally write into a file) the ‘true’ alignment of a simulation run. The following keyword parameters are available:

<code>include_inner=True</code>	if <code>True</code> , include also inner nodes in the alignment; otherwise only sequences of leaf nodes are aligned
<code>write_to=None</code>	path and filename for the output
<code>alignment_format="phylip"</code>	format of the alignment file; available are "phylip", "clustal" and "pretty"

Example usage of the class `Evolver`:

```
from asymmetree.seqevolve import Evolver, SubstModel, IndelModel, HetModel

# specify models (only 'SubstModel' is mandatory)
```

```

subst_model = SubstModel('a', 'WAG')
indel_model = IndelModel(0.01, 0.01, length_distr=('zipf', 1.821))
het_model = HetModel(2.0, classes=7, invariant=0.1)

# evolve sequences along a phylogenetic tree 'tree'
evolver = Evolver(subst_model, indel_model=indel_model, het_model=het_model)
evolver.evolve_along_tree(tree, start_length=150)

# print resulting sequences
for node, evo_seq in evolver.sequences.items():
    print(node.label, subst_model.to_sequence(evo_seq))

# output true alignment
evolver.true_alignment(write_to='path/to/alignment.phylip',
                       alignment_format='phylip')

```

### 3.4 Simulation of Genomes

The class `GenomeSimulator` combines multiple steps described in the previous section in order to conveniently simulate whole genomes/proteomes. An instance of this class is initialized with a species tree (of type `PhyloTree`), and optionally the path to an output directory (`outdir=...`) if the user wants to save the results. The gene trees and the sequences are simulated in subsequent steps using the classes' function

- (i) `simulate_gene_trees(N, **kwargs)`, and
- (ii) `simulate_sequences(subst_model, **kwargs)`.

The first step (i) takes the same keyword parameters as input as the function `simulate_gene_trees()` in Section 3.2.3, where `N` is the number of gene families to be simulated. Thus, rates for the three event types (`DLH_rates`), autocorrelation (`autocorr_variance`), the distribution of base rates (`base_rate_distr`) etc. can be specified.

The second step (ii) simulates the sequences along the observable part (without loss branches) of the simulated gene trees. The function takes the following parameters as input (keyword arguments are indicated by their default value):

<code>subst_model</code>	substitution model; instance of <code>SubstModel</code>
<code>indel_model=None</code>	model for insertions and deletions; instance of <code>IndelModel</code>
<code>het_model=None</code>	model for among site heterogeneity and invariant sites; instance of <code>HetModel</code>
<code>root_genome=None</code>	list of sequences for the roots of the gene trees; must contain the same number of <code>str</code> sequences as trees were simulated in step (i) i.e. <code>N</code> ; sequences must be compatible with the specified substitution model ( <code>model_type="n"/"a"</code> )
<code>length_distr=</code> <code>("constant", 200)</code>	distribution of the length of the root sequences if <code>root_genome</code> is not supplied; available are <code>("constant", n)</code> , <code>("gamma", shape, scale)</code> and <code>("gamma_mean", mean)</code> , for the latter the shape parameter $k = 1$
<code>min_length=10</code>	minimal length at which the distribution of lengths is truncated; must be less than the mean of this distribution
<code>max_length=None</code>	maximal length at which the distribution of lengths is truncated; must be greater than the mean of this distribution

<code>write_fastas=True</code>	if <code>True</code> and an output directory was specified, write the <b>sequences</b> (one file <b>per species</b> ) into the directory <code>fasta_files</code> in the output directory
<code>write_alignments=True</code>	if <code>True</code> and an output directory was specified, write the <b>true alignments</b> (one file <b>per gene tree</b> ) into the directory <code>alignments</code> in the output directory

After step (i), the lists of full and observable gene trees are accessible via the attributes `true_gene_trees` and `observable_gene_trees`, respectively. Moreover, the full gene trees are serialized into the directory `true_gene_trees` if an output directory was specified. After step (ii), the lists of sequence dictionaries are accessible via the attribute `sequence_dicts`.

Example usage:

```
from asymmetree.genome import GenomeSimulator
from asymmetree.treeevolve import simulate_species_tree
from asymmetree.seqevolve import SubstModel, IndelModel

# simulate a species tree
species_tree = simulate_species_tree(10, model="yule")

# specify models (heterogeneity omitted)
subst_model = SubstModel('n', 'JC69')
indel_model = IndelModel(0.015, 0.012)

# simulate 50 gene families incl. sequences
gs = GenomeSimulator(species_tree, outdir='path/to/genome_directory')
gs.simulate_gene_trees(50, DLH_rates=(1.0, 0.5, 0.0),
                      base_rate_distr=('gamma', 1.0, 1.0),
                      prohibit_extinction='per_species')
gs.simulate_sequences(subst_model,
                     indel_model=indel_model,
                     length_distr=('constant', 200))

# results have been written to directories 'true_gene_trees',
# 'fasta_files' and 'alignments' in 'path/to/genome_directory'
```

### 3.5 Best Match Inference

Phylogenetic best matches of a gene  $x$  of species  $X$  are defined as those genes  $y$  of another species  $Y \neq X$  that share the lowest common ancestor with  $x$  in the gene tree among all genes in that species [14, 16, 15]. In contrast, two genes are orthologs if their last common ancestor was a speciation event. Orthology and reciprocal best matches are closely related [15].

The subpackage `asymmetree.best_matches` contains functions to compute both relations from a given gene tree or to estimate them from distance data on a set of genes [35]. If the true (observable) gene tree is known (as e.g. the case in simulations), best matches and orthologs can be computed using the module `TrueBMG`. The functions `bmg_from_tree()` and `orthology_from_tree()` return the respective graph representation as `NetworkX (Di)Graphs`:

```
from asymmetree.best_matches import orthology_from_tree, bmg_from_tree

# tree is an observable gene tree
```

```

orthology_graph = orthology_from_tree(tree)

# Best Match Graph and Reciprocal Best Match Graph as a tuple
bmg, rbmg = bmg_from_tree(tree, supply_rbmg=True)

```

If only distance data is available, best matches have to be estimated. AsymmeTree currently implements three different methods that are described by Stadler et al. [35]:

- Extended Best Hits (module `ExtBestHits`)
- Neighborjoining [32] and midpoint rooting (module `TreeReconstruction`, requires the installation and accessibility of `RapidNJ` [34])
- Quartet method (module `Quartets`, Python implementation and wrapper for the C++ tool `qinfer`)

Please see the file `examples/best_match_infer.py` in the GitHub repo for an example usage of these modules following the simulation of gene tree scenarios.

### 3.6 Supertree Computation

The module `BuildST` contains an implementation of the BuildST algorithm described by Deng and Fernández-Baca [11] to compute a supertree from a given list of tree based on the leaf labels. The algorithm uses the dynamic graph data structure described by Henzinger and King [21] and Holm et al. [22]. The latter can also be used separately:

```

from asymmetree.tools import HDTGraph

```

The class `BuildST` is initialized with a list of trees that are of type `Tree` (thus also `PhyloTree` is allowed). The method `run()` then returns a supertree if the trees in the list are compatible *and* they overlap in their sets of leaf labels. More precisely, the graph on the set of input trees, in which two trees are connected by an edge if and only if they have at least one leaf label in common, must be connected. Example usage:

```

from asymmetree.tools import BuildST

# tree_list is a list of 'Tree' instances
st_builder = BuildST(tree_list)
supertree = st_builder.run()

if supertree:
    print(supertree.to_newick())
else:
    print('could not build a supertree')

```

### 3.7 Cograph Editing and ParaPhylo

The subpackages `asymmetree.cograph` and `asymmetree.proteinortho` contain heuristics for cograph editing and a method to compute rooted species tree from orthology/paralogy relations. The latter is a reimplementaion of [ParaPhylo](#) [19] which uses heuristics for the NP-hard steps instead of exact ILP solutions. For cograph editing, the  $\mathcal{O}(n^2)$  algorithm (where  $n$  is the number of vertices in a connected graph) by Crespelle [7] is applied. For the Maximum Consistent Triple Set problem, tree different heuristics are available:

BPMF	Best-Pair-Merge-First [39] (modified for weighted triples)
MINCUT	Aho's BUILD with weighted MinCut [1, 2]
GREEDY	a greedy approach based on Aho's BUILD

The class `TreeReconstructor` in the module `SpeciesTreeFromParalogs` computes a species tree after it is provided with one or more NetworkX graphs that represent (estimated) orthology relations. To this end, the nodes in these graphs must have a "color" attribute, since these will be the leaf labels in the reconstructed species tree. Example usage:

```
from asymmetree.paraphylo import TreeReconstructor

tree_reconstr = TreeReconstructor()

# ortho_relations is a list of orthology relations
for graph in ortho_relations:
    tree_reconstr.add_ortho_graph(graph)

# finally an estimate the species tree can be computed
tree = tr.build_species_tree(mode='BPMF')
print(tree.to_newick())
```

The module `SpeciesTreeFromProteinOrtho` contains functions to estimate a species tree from a `ProteinOrtho` [30] output file. For example, the function `reconstruct_from_proteinortho(filename, triple_mode="BPMF")` takes the filename to the output file and optionally the triple heuristic as input, and returns a tuple consisting of the estimated species tree (`PhyloTree`) and a Newick representation (`str`) containing support values for the inner nodes.

## References

- [1] Aho, A. V., Sagiv, Y., Szymanski, T. G., and Ullman, J. D. Inferring a Tree from Lowest Common Ancestors with an Application to the Optimization of Relational Expressions. *SIAM Journal on Computing*, 10(3):405–421, August 1981. ISSN 0097-5397, 1095-7111. doi: 10.1137/0210030.
- [2] Byrka, J., Guillemot, S., and Jansson, J. New results on optimizing rooted triplets consistency. *Discrete Applied Mathematics*, 158(11):1136–1147, June 2010. ISSN 0166218X. doi: 10.1016/j.dam.2010.03.004.
- [3] Byrne, K. P. and Wolfe, K. H. Consistent Patterns of Rate Asymmetry and Gene Loss Indicate Widespread Neofunctionalization of Yeast Genes After Whole-Genome Duplication. *Genetics*, 175(3):1341–1350, March 2007. ISSN 0016-6731, 1943-2631. doi: 10.1534/genetics.106.066951.
- [4] Cartwright, R. A. DNA assembly with gaps (Dawg): Simulating sequence evolution. *Bioinformatics*, 21(Suppl 3):iii31–iii38, November 2005. ISSN 1367-4803, 1460-2059. doi: 10.1093/bioinformatics/bti1200.
- [5] Chang, M. S. and Benner, S. A. Empirical Analysis of Protein Insertions and Deletions Determining Parameters for the Correct Placement of Gaps in Protein Sequence Alignments. *Journal of Molecular Biology*, 341(2):617–631, August 2004. ISSN 00222836. doi: 10.1016/j.jmb.2004.05.045.
- [6] Corneil, D. G., Perl, Y., and Stewart, L. K. A Linear Recognition Algorithm for Cographs. *SIAM Journal on Computing*, 14(4):926–934, November 1985. ISSN 0097-5397, 1095-7111. doi: 10.1137/0214065.
- [7] Crespelle, C. Linear-Time Minimal Cograph Editing. 2019.



- [8] Criscuolo, A. Empirical Models of Amino Acid Substitution. <http://giphy.pasteur.fr/empirical-models-of-amino-acid-substitution/>.
- [9] Dalquen, D. A., Anisimova, M., Gonnet, G. H., and Dessimoz, C. ALF—A Simulation Framework for Genome Evolution. *Molecular Biology and Evolution*, 29(4):1115–1123, April 2012. ISSN 1537-1719, 0737-4038. doi: 10.1093/molbev/msr268.
- [10] Dayhoff, M. and Schwartz, R. A model for evolutionary change in proteins. In *Atlas of Protein Sequence and Structure*, pages 345–352. National Biomedical Research Foundation, Washington D.C., 1978.
- [11] Deng, Y. and Fernández-Baca, D. Fast Compatibility Testing for Rooted Phylogenetic Trees. page 12 pages, 2016. doi: 10.4230/LIPICS.CPM.2016.12.
- [12] Felsenstein, J. *Inferring Phylogenies*. Sinauer Associates, Sunderland, Mass, 2004. ISBN 978-0-87893-177-4.
- [13] Fletcher, W. and Yang, Z. INDELible: A Flexible Simulator of Biological Sequence Evolution. *Molecular Biology and Evolution*, 26(8):1879–1888, August 2009. ISSN 0737-4038, 1537-1719. doi: 10.1093/molbev/msp098.
- [14] Geiß, M., Chávez, E., González Laffitte, M., López Sánchez, A., Stadler, B. M. R., Valdivia, D. I., Hellmuth, M., Hernández Rosales, M., and Stadler, P. F. Best match graphs. *Journal of Mathematical Biology*, 78(7):2015–2057, June 2019. ISSN 0303-6812, 1432-1416. doi: 10.1007/s00285-019-01332-9.
- [15] Geiß, M., Laffitte, M. E. G., Sánchez, A. L., Valdivia, D. I., Hellmuth, M., Rosales, M. H., and Stadler, P. F. Best match graphs and reconciliation of gene trees with species trees. *Journal of Mathematical Biology*, January 2020. ISSN 0303-6812, 1432-1416. doi: 10.1007/s00285-020-01469-y.
- [16] Geiß, M., Stadler, P. F., and Hellmuth, M. Reciprocal best match graphs. *Journal of Mathematical Biology*, 80(3):865–953, February 2020. ISSN 0303-6812, 1432-1416. doi: 10.1007/s00285-019-01444-2.
- [17] Gillespie, D. T. A general method for numerically simulating the stochastic time evolution of coupled chemical reactions. *Journal of Computational Physics*, 22(4):403–434, December 1976. ISSN 00219991. doi: 10.1016/0021-9991(76)90041-3.
- [18] Hagen, O. and Stadler, T. *TreeSim* GM : Simulating phylogenetic trees under general Bellman–Harris models with lineage-specific shifts of speciation and extinction in R. *Methods in Ecology and Evolution*, 9(3):754–760, March 2018. ISSN 2041-210X, 2041-210X. doi: 10.1111/2041-210X.12917.
- [19] Hellmuth, M., Wieseke, N., Lechner, M., Lenhof, H.-P., Middendorf, M., and Stadler, P. F. Phylogenomics with paralogs. *Proceedings of the National Academy of Sciences*, 112(7):2058–2063, February 2015. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.1412770112.
- [20] Henikoff, S. and Henikoff, J. G. Amino acid substitution matrices from protein blocks. *Proceedings of the National Academy of Sciences*, 89(22):10915–10919, November 1992. ISSN 0027-8424, 1091-6490. doi: 10.1073/pnas.89.22.10915.
- [21] Henzinger, M. R. and King, V. Randomized dynamic graph algorithms with polylogarithmic time per operation. In *Proceedings of the Twenty-Seventh Annual ACM Symposium on Theory of Computing - STOC '95*, pages 519–527, Las Vegas, Nevada, United States, 1995. ACM Press. ISBN 978-0-89791-718-6. doi: 10.1145/225058.225269.
- [22] Holm, J., de Lichtenberg, K., and Thorup, M. Poly-logarithmic deterministic fully-dynamic algorithms for connectivity, minimum spanning tree, 2-edge, and biconnectivity. *Journal of the ACM*, 48(4):723–760, July 2001. ISSN 00045411. doi: 10.1145/502090.502095.
- [23] Jones, D. T., Taylor, W. R., and Thornton, J. M. The rapid generation of mutation data matrices from protein sequences. *Bioinformatics*, 8(3):275–282, 1992. ISSN 1367-4803, 1460-2059. doi: 10.1093/bioinformatics/8.3.275.
- [24] Jukes, T. H. and Cantor, C. R. Evolution of Protein Molecules. In *Mammalian Protein Metabolism*, pages 21–132. Elsevier, 1969. ISBN 978-1-4832-3211-9. doi: 10.1016/B978-1-4832-3211-9.50009-7.

- [25] Keller-Schmidt, S. and Klemm, K. A model of macroevolution as a branching process based on innovations. *Advances in Complex Systems*, 15(07):1250043, October 2012. ISSN 0219-5259, 1793-6802. doi: 10.1142/S0219525912500439.
- [26] Kendall, D. G. On the Generalized "Birth-and-Death" Process. *The Annals of Mathematical Statistics*, 19(1):1–15, March 1948. ISSN 0003-4851. doi: 10.1214/aoms/1177730285.
- [27] Kimura, M. A simple method for estimating evolutionary rates of base substitutions through comparative studies of nucleotide sequences. *Journal of Molecular Evolution*, 16(2):111–120, December 1980. ISSN 0022-2844.
- [28] Kishino, H., Thorne, J. L., and Bruno, W. J. Performance of a Divergence Time Estimation Method under a Probabilistic Model of Rate Evolution. *Molecular Biology and Evolution*, 18(3):352–361, March 2001. ISSN 1537-1719, 0737-4038. doi: 10.1093/oxfordjournals.molbev.a003811.
- [29] Le, S. Q. and Gascuel, O. An Improved General Amino Acid Replacement Matrix. *Molecular Biology and Evolution*, 25(7):1307–1320, April 2008. ISSN 0737-4038, 1537-1719. doi: 10.1093/molbev/msn067.
- [30] Lechner, M., Findeiß, S., Steiner, L., Marz, M., Stadler, P. F., and Prohaska, S. J. Proteinortho: Detection of (Co-)orthologs in large-scale analysis. *BMC Bioinformatics*, 12(1), December 2011. ISSN 1471-2105. doi: 10.1186/1471-2105-12-124.
- [31] Lechner, M., Hernandez-Rosales, M., Doerr, D., Wieseke, N., Thévenin, A., Stoye, J., Hartmann, R. K., Prohaska, S. J., and Stadler, P. F. Orthology Detection Combining Clustering and Synteny for Very Large Datasets. *PLoS ONE*, 9(8):e105015, August 2014. ISSN 1932-6203. doi: 10.1371/journal.pone.0105015.
- [32] Saitou, N. and Nei, M. The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Molecular Biology and Evolution*, 4(4):406–425, July 1987. ISSN 0737-4038. doi: 10.1093/oxfordjournals.molbev.a040454.
- [33] Schaller, D., Geiß, M., Stadler, P. F., and Hellmuth, M. Complete Characterization of Incorrect Orthology Assignments in Best Match Graphs. *arXiv:2006.02249 [cs, math, q-bio]*, June 2020.
- [34] Simonsen, M., Mailund, T., and Pedersen, C. N. S. Rapid Neighbour-Joining. In Crandall, K. A. and Lagergren, J., editors, *Algorithms in Bioinformatics*, volume 5251, pages 113–122. Springer Berlin Heidelberg, Berlin, Heidelberg, 2008. ISBN 978-3-540-87360-0 978-3-540-87361-7. doi: 10.1007/978-3-540-87361-7\_10.
- [35] Stadler, P. F., Geiß, M., Schaller, D., López Sánchez, A., González Laffitte, M., Valdivia, D. I., Hellmuth, M., and Hernández Rosales, M. From pairs of most similar sequences to phylogenetic best matches. *Algorithms for Molecular Biology*, 15(1):5, December 2020. ISSN 1748-7188. doi: 10.1186/s13015-020-00165-2.
- [36] Stadler, T. Simulating trees with a fixed number of extant species. *Systematic Biology*, 60(5):676–684, October 2011. ISSN 1076-836X. doi: 10.1093/sysbio/syr029.
- [37] Tavaré, S. Some Probabilistic and Statistical Problems in the Analysis of DNA Sequences. *Lectures on Mathematics in the Life Sciences*, 17:57–86, 1986.
- [38] Whelan, S. and Goldman, N. A General Empirical Model of Protein Evolution Derived from Multiple Protein Families Using a Maximum-Likelihood Approach. *Molecular Biology and Evolution*, 18(5): 691–699, May 2001. ISSN 1537-1719, 0737-4038. doi: 10.1093/oxfordjournals.molbev.a003851.
- [39] Wu, B. Y. Constructing the Maximum Consensus Tree from Rooted Triples. *Journal of Combinatorial Optimization*, 8(1):29–39, March 2004. ISSN 1382-6905. doi: 10.1023/B:JOCO.0000021936.04215.68.
- [40] Yang, Z. PAML: A program package for phylogenetic analysis by maximum likelihood. *Bioinformatics*, 13(5):555–556, 1997. ISSN 1367-4803, 1460-2059. doi: 10.1093/bioinformatics/13.5.555.
- [41] Yang, Z. *Computational Molecular Evolution*. Oxford Series in Ecology and Evolution. Oxford University Press, Oxford, 2006. ISBN 978-0-19-856699-1 978-0-19-856702-8.

- [42] Yang, Z. *Molecular Evolution: A Statistical Approach*. Oxford University Press, Oxford, United Kingdom ; New York, NY, United States of America, first edition edition, 2014. ISBN 978-0-19-960261-2 978-0-19-960260-5.
- [43] Yule, G. U. A mathematical theory of evolution, based on the conclusions of Dr. J. C. Willis, F. R. S. *Philosophical Transactions of the Royal Society of London. Series B, Containing Papers of a Biological Character*, 213(402-410):21–87, 1924. ISSN 0264-3960, 2053-9266. doi: 10.1098/rstb.1925.0002.

## A Subpackages and Modules

Packages and Modules	Description
<b>treeevolve</b>	
SpeciesTree	Simulator for dated species trees with different models.
GeneTree	Simulator for dated gene trees, construction of the observable gene tree.
EvolutionRates	Simulation of evolution rate asymmetries, autocorrelation between ancestors and descendants as well as correlation between genes in the same species.
Scenario	Wrapper class for species and gene tree scenarios, computation of the (R)BMG as well as event counts and some statistics.
NoisyMatrix	Generation of a noisy matrix (random perturbation or wrong topology noise).
<b>sequevolve</b>	
Evolver	Includes the class <b>Evolver</b> for the simulation of sequences along a phylogenetic tree.
EvolvingSequence	Data structure for sequences that are under evolution based on a doubly-linked list.
SubstModel	Substitution model for the sequence simulation.
IndelModel	Insertion/deletion (indel) model.
HetModel	Model for rate heterogeneity among the sites of an evolving sequence and for invariant sites.
Alignment	Construction of the true multiple sequence alignment after the simulation of sequences along a tree.
EmpiricalModels	Contains rate matrices and equilibrium frequencies for the empirical substitution models, taken from [8].
<b>tools</b>	
Tree	Includes the basic class <b>Tree</b> , provides functions for tree traversals, Newick parser, etc.
PhyloTree	Includes the class <b>PhyloTree</b> for phylogenetic trees (inherits from <b>Tree</b> ), provides a Newick parser, etc.
BuildST	Includes the class <b>BuildST</b> that computes a supertree from a given list of tress (with overlapping labels) [11].
hdtgraph.DynamicGraph	Dynamic graph data structure described by Holm et al. [22].
DoublyLinkedList	Implementation of a doubly-linked list.
GraphTools	Miscellaneous functions for graphs, e.g. check for graph equality.
DistanceCalculation	Calculation of maximum likelihood distances of pairs of aligned sequences.
<b>best_matches</b>	
TrueBMG	Computation of the true (R)BMG from a gene tree as well as the true orthology relation.
ExtBestHits	Implementation of the <i>Extended Best Hits</i> method, optionally uses <b>qinfer</b> .
TreeReconstruction	Reconstruction of the gene tree with <b>RapidNJ</b> [34] and midpoint rooting.
Quartets	Implementation of <i>Quartet</i> approach with two different methods for outgroup selection, optionally uses <b>qinfer</b> .

*Table continued on next page*

Table continued from previous page

<b>LRTConstructor</b>	Construction of a least resolved tree (LRT) from a BMG via <i>informative triples</i> (optionally uses minimal edge cuts) or from a leaf-colored tree.
<b>Augmentation</b>	Augmentation of the least resolved tree (w.r.t. some BMG) in order to identify all unambiguously false orthology assignments [33].
<b>cograph</b>	
<b>Cograph</b>	Includes the classes <b>Cotree</b> and <b>CotreeNode</b> as well as a generator for random cotrees/cographs.
<b>CographEditor</b>	Implements a heuristic for cograph editing [7].
<b>LinearCographDetector</b>	Implements an $\mathcal{O}(n + m)$ algorithm for cograph detection [6].
<b>paraphylo</b>	
<b>SpeciesTreeFromParalogs</b>	Species tree reconstruction from orthology/paralogy relations. Heuristic version of <b>ParaPhylo</b> [19].
<b>SpeciesTreeFromProteinOrtho</b>	Species tree reconstruction from <b>ProteinOrtho</b> [30, 31] output.
<b>visualize</b>	
<b>GeneTreeVis</b>	Visualization of simulated gene trees (of type <b>PhyloTree</b> ), experimental.