1          Version dated: December 17, 2017

2

# Fast Bayesian Inference of Phylogenetic Models Using Parallel Likelihood Calculation and Adaptive Metropolis Sampling

VENELIN MITOV[1,2], TANJA STADLER[1,2]

[1]*Swiss Federal Institute of Technology in Zurich, Switzerland;*

[2]*Swiss Institute of Bioinformatics, Switzerland*

**Corresponding authors:** Venelin Mitov, Department of Biosystem Sciences and Engineering, Swiss Federal Institute of Technology, Mattenstrasse 26, CH-4058 Basel, Switzerland; E-mail: vmitov@gmail.com.

Tanja Stadler, Department of Biosystem Sciences and Engineering, Swiss Federal Institute of Technology, Mattenstrasse 26, CH-4058 Basel, Switzerland; E-mail: tanja.stadler@bsse.ethz.ch.

## ABSTRACT

Phylogenetic comparative methods have been used to model trait evolution, to test selection versus neutral hypotheses, to estimate optimal trait-values, and to quantify the rate of adaptation towards these optima. Several authors have proposed algorithms calculating the likelihood for trait evolution models, such as the Ornstein-Uhlenbeck (OU)

20  process, in time proportional to the number of tips in the tree. Combined with
21  gradient-based optimization, these algorithms enable maximum likelihood (ML) inference
22  within seconds, even for trees exceeding 10,000 tips. Despite its useful statistical
23  properties, ML has been criticised for being a point estimator prone to getting stuck in
24  local optima. As an elegant alternative, Bayesian inference explores the entire information
25  in the data and compares it to prior knowledge but, usually, needs much longer time, even
26  on small trees. Here, we propose an approach to use the full potential of ML and Bayesian
27  inference, while keeping the runtime within minutes. Our approach combines (i) a new
28  algorithm for parallel traversal of the lineages in the tree, enabling parallel calculation of
29  the likelihood; (ii) a previously published method for adaptive Metropolis sampling. In
30  principle, the strategy of (i) and (ii) can be applied to any likelihood calculation on a tree
31  which proceeds in a pruning-like fashion, leading to enormous speed improvements. We
32  implement several variants of the parallel algorithm in the form of a generic C++ library,
33  "SPLiTTree", capable to choose automatically the optimal algorithm for a given task and
34  computing platform. We give examples of models of discrete and continuous trait evolution
35  that are amenable to parallel likelihood calculation. As a complete showcase, we implement
36  the phylogenetic Ornstein-Uhlenbeck mixed model (POUMM) in the form of an easy-to-use
37  and highly configurable R-package that calls the library as a back-end. In addition to the
38  above-mentioned usage of comparative methods, POUMM allows to estimate non-heritable
39  variance and phylogenetic heritability. Using SPLiTTree, calculating the POUMM
40  likelihood on a 4-core SIMD-enabled processor is up to $10\times$ faster than serial
41  implementations written in C and hundreds of times faster than serial implementations
42  written in R. By combining SPLiTTree likelihood calculation with adaptive Metropolis
43  sampling, the time for Bayesian POUMM inference on a tree of ten thousand tips is
44  reduced from several days to a few minutes.

45  Keywords: Parallel tree traversal, post-order traversal, pre-order traversal, pruning,

46  discrete character, continuous trait, phylogenetic comparative models, Brownian motion,

47  measurement error, stabilizing selection, environmental contribution

<div align="center">

48  ## INTRODUCTION

</div>

49      The past decades have seen active developement of phylogenetic comparative

50  models of trait evolution, progressing from null neutral models, such as single-trait

51  Brownian motion (BM), to complex multi-trait models incorporating selection, interaction

52  between trait values and diversification, and co-evolution of multiple traits (O'Meara 2012;

53  Manceau et al. 2016). Recent works have shown that, for a broad family of phylogenetic

54  comparative models, the likelihood of an observed tree and data conditioned on the model

55  parameters can be computed in time proportional to the size of the tree (FitzJohn 2012;

56  Ho and Ané 2014; Goolsby et al. 2016; Manceau et al. 2016). This family includes

57  Gaussian models like Brownian motion and Ornstein-Uhlenbeck phylogenetic models as

58  well as some non-Gaussian models like phylogenetic logistic regression (Paradis and Claude

59  2002; Ives and Garland 2010; Ho and Ané 2014). All of these likelihood calculation

60  techniques rely on post-order tree traversal also known as *"pruning"* (Felsenstein 1973,

61  1981, 1983). Using pruning algorithms for likelihood calculation in combination with a

62  gradient-based optimization method (Boyd and Vandenberghe 2004), maximum likelihood

63  model inference runs within seconds on contemporary computers, even for phylogenies

64  containing many thousands of tips (Ho and Ané 2014). Other important features of the

65  maximum likelihood estimate (MLE) are its simple interpretation as the point in

66  parameter space maximizing the probability of the observed data under the assumed

67  model, and its theoretical properties making it ideal for hypothesis testing and for model

68  selection via likelihood ratio tests and information criteria. However, major disadvantages

69  of MLE are that, being a point estimate, it does not allow to explore the likelihood surface;

70  gradient based optimization, while fast, is prone to getting stuck in local optima.

71      As an elegant alternative, Bayesian approaches such as Markov Chain Monte Carlo

72  (MCMC) allow to incorporate prior knowledge in the model inference and provide posterior

73  samples and high posterior density (HPD) intervals for the model parameters (Slater et al.

74  2012a; FitzJohn 2012). In contrast with ML inference, though, Bayesian inference methods

3

<sub>75</sub> require many orders of magnitude more likelihood evaluations, in order to obtain a

<sub>76</sub> valuable posterior sample. This presents a bottleneck in Bayesian analysis, in particular,

<sub>77</sub> when faced with large phylogenies of many thousands of tips, such as transmission trees

<sub>78</sub> from large-scale epidemiological studies, e.g. Hodcroft et al. (2014). While big data

<sub>79</sub> provides sufficient statistical power to fit a complex model, the time needed to perform a

<sub>80</sub> full scale Bayesian inference often limits the choice to a faster but less informative

<sub>81</sub> ML-inference, or a Bayesian inference on a simplified model.

<sub>82</sub> In this article, we propose a general approach allowing to use ML and Bayesian

<sub>83</sub> inference to their full potential, even for complex models and for very large trees exceeding

<sub>84</sub> millions of tips. To achieve this goal, our approach combines two ideas: (i) the pruning

<sub>85</sub> algorithm for likelihood calculation can be accelerated by orders of magnitude through

<sub>86</sub> parallel processing (traversal) of the independent lineages in the tree; (ii) the number of

<sub>87</sub> iterations needed for MCMC convergence can be reduced several times by the use of

<sub>88</sub> adaptive Metropolis sampling (Vihola 2012; Scheidegger 2017). A nice property of the

<sub>89</sub> parallel pruning algorithm is that its parallel efficiency increases with the number of tips in

<sub>90</sub> the tree as well as with the complexity of the pruning operation. Thus, for large trees and

<sub>91</sub> complex models, the parallel speed-up is practically limited by the number of available

<sub>92</sub> processing cores and other hardware resources such as the memory bandwidth.

<sub>93</sub> We provide SPLiTTree: a C++ library for Serial and Parallel Lineage Traversal of

<sub>94</sub> Trees. The library is designed to be generic with respect to the tree topology, the type of

<sub>95</sub> data, the model parameters and the specific visit-node operation. It has been tested on

<sub>96</sub> large ($N > 100,000$) balanced and highly unbalanced trees. Noticing that there is no

<sub>97</sub> one-size-fits-all strategy for optimal parallel speed-up (e.g. the parallelization may be

<sub>98</sub> useless or even slowing-down the performance in the case of a ladder tree or in the case of a

<sub>99</sub> memory intensive visit-node operation), SPLiTTree implements several parallelization

<sub>100</sub> strategies (e.g. classical serial traversal versus queue-based and range-based parallel

traversal) and is capable to automatically select the fastest strategy for a given tree and computing hardware. In this way, the library allows the user to neglect complex technical aspects of parallel programming and focus on implementing the visit-node operation specific for the application of interest.

We begin with a formulation of a general framework for parallel post-order tree traversal (pruning). Then, we give examples of models of discrete and continuous trait evolution, amenable to using the library for parallel likelihood calculation. These examples represent mere adaptations of previously published pruning algorithms to the terms of the framework (Felsenstein 1983; Pagel 1994; Ho and Ané 2014).

We go on with a showcase demonstrating the combined use of parallel pruning and adaptive Metropolis sampling on a univariate phylogenetic Ornstein-Uhlenbeck mixed model (POUMM). Previously, we and other authors have used the POUMM in heritability analysis of epidemiological data (Mitov and Stadler 2016; Bachmann et al. 2017; Bertels et al. 2017; Blanquart et al. 2017). Here, we describe the pruning formulation of the POUMM likelihood, validate the technical correctness of the implementation and compare the performance of this implementation to several other implementations of OU-models of evolution including FitzJohn (2012) and Pennell et al. (2014). We report significant speed-ups, both for the time for single likelihood calculation as well as for the overall Bayesian inference of the model.

# Setup

Through the rest of the article we will use the following notation. Given is a rooted phylogenetic tree $\mathcal{T}$ with a total of $M$ nodes, including $N < M$ tips denoted $1, ..., N$, $M - N - 1$ internal nodes denoted $N + 1, ..., M - 1$, and a root node denoted $M$ (Fig. 1). Without restrictions on the tree topology, non-ultrametric trees (i.e. tips have different heights) and polytomies (i.e. nodes with any finite number of descendants) are accepted.

5

126 We denote by $\mathcal{T}_i$ the subtree rooted at node $i$. For any tip or internal node $i$, we denote its

127 parent node by $Parent(i)$. For any node $j$, we denote by $Desc(j)$ the set of its direct

128 descendants ($Desc(j) = \phi$ if $j$ denotes a tip). Furthermore, for any $i \in Desc(j)$, we denote

129 by $t_i$ the length of the branch leading to $i$ and by $h_i$ the height of $i$, i.e. the sum of branch

130 lengths from the root to $i$. The mean height of all tips in the tree is denoted by $\bar{h}$. For two

131 tips, $i$ and $k$, we denote by $h_{(ik)}$ the height to their most recent common ancestor (mrca),

132 and by $d_{ik}$ the sum of branch lengths on the path from $i$ to $k$ (also called phylogenetic or

133 patristic distance between $i$ and $k$). Associated with each node $i$ there is an input data in

134 the form of a single or multivariate categorical or numerical value denoted $z_i$. For tips, $z_i$

135 can be partially unobserved (having NA entries), while for internal nodes or the root it can

136 also be fully absent (NULL). We denote by $\mathbf{z}_i$ the sub-vector of input data for the nodes in

137 $\mathcal{T}_i$. Associated with each node, $i$, there is a vector of model parameters, $\Theta_i$. We use bold

138 style $\mathbf{t}$, $\mathbf{z}$ and $\boldsymbol{\Theta}$ when denoting the vectors of all branch lengths, input data and

139 parameters.

## A GENERAL FRAMEWORK FOR PARALLEL TREE TRAVERSAL

141 Let $F_{\mathcal{T}}(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$ be a function of the branch lengths, the input data and the

142 parameters. A post-order tree traversal algorithm can be used to calculate $F_{\mathcal{T}}$ if, for all

143 subtrees $\mathcal{T}_j$ of $\mathcal{T}$, there exist functions $S_j(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$, hereby called "states", satisfying the

144 following rules:

145 (1) $F_{\mathcal{T}}(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$ can be calculated from $S_M(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$;

146 (2) For each node $j \in \{1, ..., M\}$, there exists a (recursive) relationship $R_j$ between $S_j$

147 and the set of states at $j$'s descendants, such that:

$$S_j(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}) = R_j\bigg( \big\{ S_i(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}) : i \in Desc(j) \big\}, \mathbf{t}, \mathbf{z}, \boldsymbol{\Theta} \bigg). \tag{1}$$
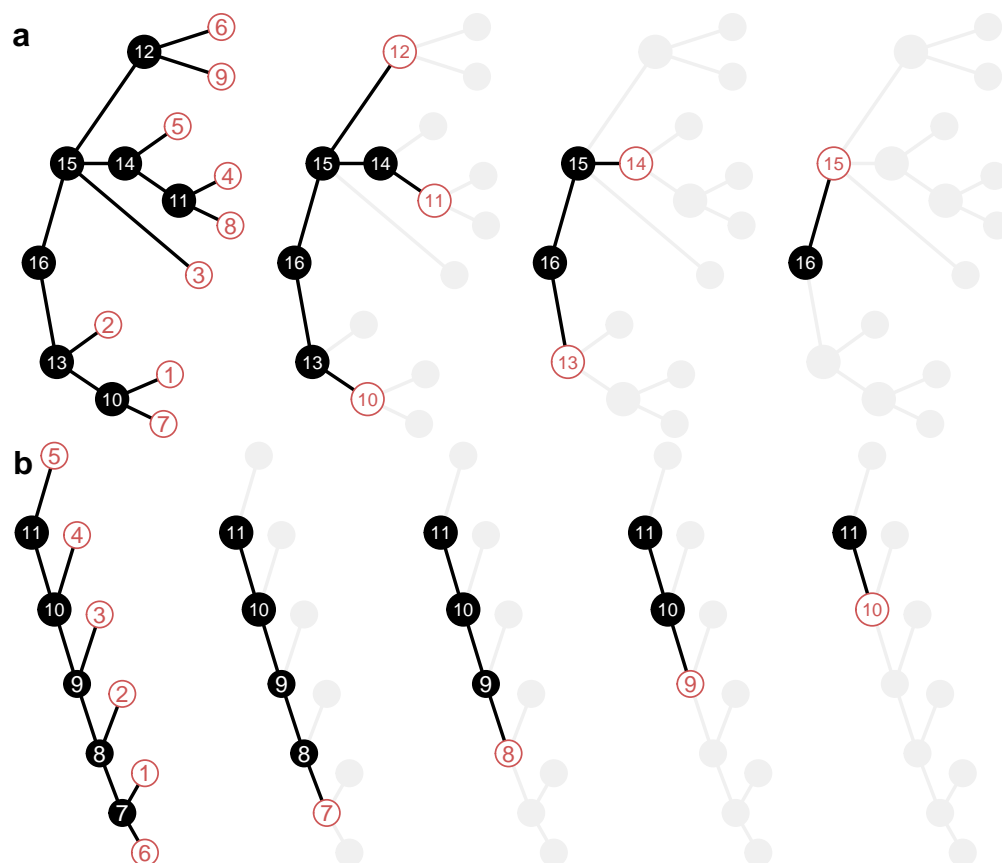
Figure 1: **Parallel pruning.** The trees from left to right depict generations of nodes that can be processed in parallel. The processing of a node consists in calculating its state based on the input data, the branch lengths and the states of the node's direct descendants (eq. 1). black: nodes having one or more non-processed descendants; red: nodes ready to be processed; grey: nodes processed in a previous generation. a) a balanced tree; b) a ladder.

148    We note that analogical terms can be defined for pre-order tree traversal. In this

149   case the target functions are values $Z_{\mathcal{T},j}(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$ corresponding to the nodes $j \in \{1, ..., M\}$,

150   and rule (2) is updated to:

151    (2')  $Z_M$ can be calculated from the input data. For each node $j \in \{1, ..., M-1\}$, there

152         exists a (recursive) relationship $R_j'$ between $Z_j$ and $Z_{Parent(j)}$, such that:

$$Z_j(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}) = R_j'\big(Z_{Parent(j)}(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}), \mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}\big). \tag{2}$$

153   The states, i.e. the values of the functions $S_j$ and $Z_j$, may be deterministic or stochastic

154   functions of the input tree and data. They can be real numbers, vectors, matrices or higher

155   order combinations thereof.

156       For the rest of the article, we focus on parallel post-order tree traversal or

157   "pruning", noting that the algorithms for parallel pre-order traversal are simple analogies.

158   The SPLiTTree library implements both traversal types.

159       Rule (2) ensures that calculating the state of a node $j$ can be done independently

160   from the calculation of any other node $k$, provided that neither $j$ is an ancestor of $k$, nor $k$

161   is an ancestor of $j$. Based on this observation, we develop two alternative parallel

162   algorithms for calculating the root state $S_M$ described in the following subsections.

## Queue-based parallel pruning

163

164       It is possible to parallelize the computation of the states $S_j$ across multiple

165   computing threads using a first-in-first-out list (queue) of the nodes in the tree (algorithm

166   1). Initially, the queue is filled with all tips in the tree and a counter with the number of

167   direct descendants is set for each internal or root node. Then, each thread takes a node $i$

168   from the front of the queue, calculates its state and decrements the counter of $Parent(i)$. If

169   the counter of $Parent(i)$ has become zero, $Parent(i)$ is added to the queue, so that it will

8

170   be processed as soon as a free thread picks it from the queue. Assuming an unlimited

171   number of threads and a negligible cost of the queue- and the counter- operations,

172   algorithm 1 guarantees that a node will be processed immediately after all of its direct

173   descendants have been processed. Thus, in theory, algorithm 1 maximizes the parallel

174   execution. However, an implementation of the atomic operations on the queue and the

175   counters would have to rely on a thread synchonization mechanism such as a mutex, which

176   can be slow on some systems. Thus, a decent parallelization speed-up would only be

177   possible if the overall cost of synchronization is insignificant compared to the functions $R_j$.

## Range-based parallel pruning

179      We consider an alternative of algorithm 1 minimizing the synchronization overhead.

180   This approach consists in splitting the tree into "generations" of nodes, such that nodes

181   within a generation can be processed in random order and in parallel, but only if all

182   generations containing descendants of these nodes have already been processed (fig. 1). A

183   "master" thread is responsible for launching a team of "worker" threads on each

184   generation, starting from a generation of all tips, then taking their parents, and so on until

185   reaching the root of the tree. To be efficient, this procedure requires that the data

186   associated with the nodes in a generation occupy a consecutive region in the address-space.

187   This eliminates the need for synchronization between the worker threads, because each

188   worker thread can deduce its own portion based on its thread-id and the address-range of

189   the generation. To orchestrate the worker teams, the master thread only needs to keep

190   account of the address-ranges. Technically, this is accomplished by iterating over a vector

191   of offsets (algorithm 2).

192      In algorithm 2, the number of synchronization points is reduced to the number of

193   generations, $K$. In balanced trees, $K$ would increase logarithmically with $N$ and, for big

194   $N$, the tree would be split into a few generations of many nodes (fig. 1a). If there are no

9

---

**Algorithm 1:** Queue-based parallel pruning

---

**Input:** $\mathcal{T}$, $\mathbf{t}$, $\mathbf{z}$, $\boldsymbol{\Theta}$
**Output:** $S_M(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$
    /* a vector of $M$ states                                                      */

**1**   $State \longleftarrow [...]_M$;
    /* a vector of the numbers of remaining descendants for each node    */

**2**   $NumDesc \longleftarrow \big[|Desc(i)| : i \in \{1, ..., M\}\big]$;
    /* initiate $Queue$ with all tips:                                     */

**3**   $Queue \longleftarrow [1, ..., N]$;

**4**   **begin** Parallel block

**5**      **while** *(TRUE)* **do**
          /* if $Queue$ is empty, thread waits.                               */

**6**          $j \longleftarrow$ PopFirst($Queue$);

**7**          $State[j] \longleftarrow R_j\bigg(\big\{State[i] : i \in Desc(j)\big\}, \mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}\bigg)$;

**8**          **if** *(j < M)* **then**
              /* the root has not been processed yet.                       */

**9**             $NumDesc[Parent(j)] \longleftarrow NumDesc[Parent(j)] - 1$;

**10**           **if** *(NumDesc[Parent(j)] == 0)* **then**
               /* If $Queue$ is currently empty a waiting thread will be notified.           */

**11**                AddLast($Queue$, $Parent(j)$);

**12**          **else**
              /* the root has been processed.                              */
              /* Notify waiting threads by adding a stopping node-id to $Queue$.           */

**13**             AddLast($Queue$, $M + 1$);
              /* All work done, exit the loop.                              */

**14**             break;

**15**   **return** $State[M]$;

---

---

**Algorithm 2:** Range-based parallel pruning

---

**Input:** $\mathcal{T}$, **t**, **z**, $\boldsymbol{\Theta}$

**Output:** $S_M(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$

**Data:**

```
/* A pre-calculated vector with starting offsets for each generation:
   */
```

1   $Range = \big[0, N, N + |G_1|, N + |G_1| + |G_2|, ..., M - 1, M\big]_{K+1}$;

```
/* a vector of M elements                                              */
```

2   $State \longleftarrow [0, ..., 0]_M$;

```
/* The master thread iterates over the generations:                    */
```

3   **foreach** $k \in \{1, ..., K\}$ **do**

    
```
/* The master thread starts a team of worker threads running equal
   portions of the following loop:                                      */
```

4      **foreach** $j \in \{Range[k] + 1, ..., Range[k+1]\}$ **do**

5         $State[j] \longleftarrow R_j\Big(\big\{State[i] : i \in Desc(j)\big\}, \mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}\Big)$;

6   **return** $State[M]$;

---

other blocking resources, such as memory bandwidth, the parallel speed-up would be nearly as high as the number of processing cores and the synchronization overhead would be negligible. Conversely, in strongly unbalanced trees, $K$ would tend to increase linearly with $N$ and the tree would be split into many generations of a few nodes (fig. 1b). This would result in low parallel speed-up and excessive synchronization cost. Also noteworthy is the fact that algorithm 2 reduces the number of synchronization points at the cost of some parallelization. If each worker thread gets assigned to an approximately equal number of nodes in a generation and if a few of the nodes take much longer time to process than the rest, then most of the worker threads would have to wait until the last node in the generation has been processed. These subtleties of the tree and input data indicate that there is no "one size fits all" strategy when it comes to optimizing the parallelization speed-up. The parallel pruning framework provides two ways to deal with these: (a) allowing the user to choose a parallelization mode before executing a pruning operation on a given tree and data; (b) providing a mode "auto", in which the framework compares the

11

<sup>209</sup> execution time of different pruning algorithms during the first several calls on a given tree

<sup>210</sup> and data, choosing the fastest one for all subsequent calls.

<sup>211</sup> <div align="center">*The SPLiTTree library*</div>

<sup>212</sup> We provide SPLiTTree in the form of an open source C++ library licensed under

<sup>213</sup> version 3.0 of the GNU Lesser General Public License (LGPL v3.0) and available on

<sup>214</sup> https://github.com/venelin/SPLiTTree.git. In its current implementation, the library uses

<sup>215</sup> the C++11 language standard, the standard template library (STL) and the OpenMP

<sup>216</sup> standard for parallel processing. The library is designed as a set of C++ template classes,

<sup>217</sup> generic with respect to the application specific details, such as the types of input data,

<sup>218</sup> model parameters and definitions of the node states, $S_i$, and visit-node functions, $R_i$. The

<sup>219</sup> library defines two layers (fig. 2):

<sup>220</sup> • a framework layer defining the main logical and data structures. These include a

<sup>221</sup> linear algorithm for initial reordering and splitting of the input tree into generations

<sup>222</sup> of nodes, which can be visited in parallel, both during a post-order as well as

<sup>223</sup> pre-order traversal, and a growing collection of pre-order and post-order traversal

<sup>224</sup> algorithms, targeting different parallelization modes (e.g. queue-based versus

<sup>225</sup> range-based parallelization) on different computing devices (currently implemented

<sup>226</sup> for CPUs only).

<sup>227</sup> • a user layer at which the user of the library must write a

<sup>228</sup> `CustomTraversalSpecification` class defining all typedefs and methods of the

<sup>229</sup> interface `TraversalSpecification`. The most important of entity to define is the

<sup>230</sup> method `VisitNode` specifying the function $R_i$ for each tip or internal node, $i$.

<sup>231</sup> The bridge between the two layers is provided by an object of the `TraversalTask`

<sup>232</sup> template class. Once the `TraversalSpecification` implementation has been written, the
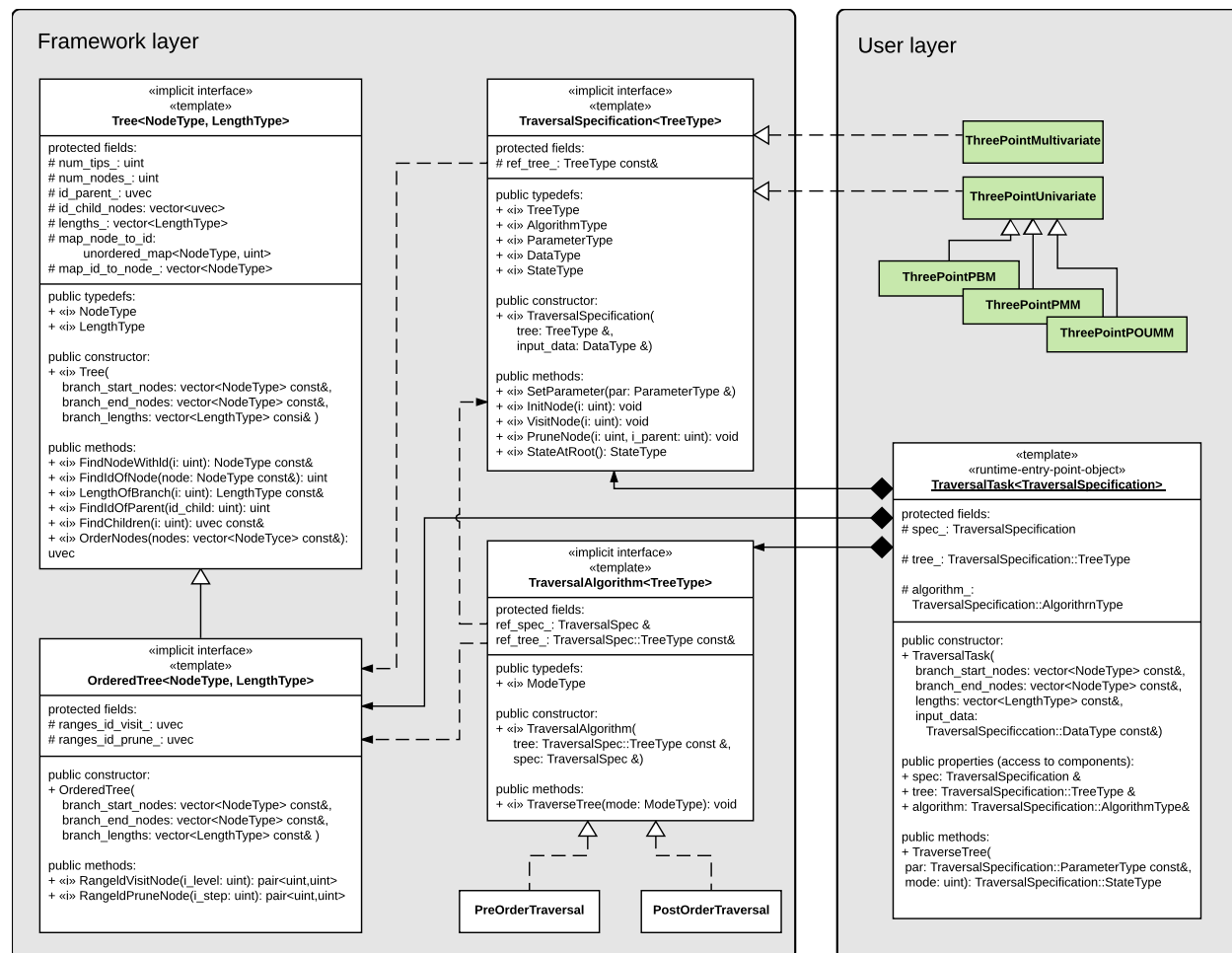
<div align="center">12</div>

Figure 2: **A UML class diagram of the SPLiTTree library** In the framework layer, the class `TraversalSpecification` defines the application-specific data types and logic; the class `Tree` serves as a base-class implementing common tree operations, such as constructing a tree from a list of branches, checking the validity of the input (e.g. lack of cycles or isolated branches), finding the parent and the descendants of a node, etc; the class `OrderedTree` maintains the order of the nodes in a tree so that they can be split in contiguous generations for parallel post-order or pre-order traversal; the class `TraversalAlgorithm` serves as a base class and an implicit interface for its two subclasses implementing the two supported types of tree traversal: `PreOrderTraversal` and `PostOrderTraversal`. At the user layer are the user defined implementations of the `TraversalSpecification` interface (shown in green) and an instance of the `TraversalTask` class, which constructs all necessary internal objects and serves as a runtime entry point to the framework.

<sub>233</sub> user instantiates a `TraversalTask` object passing the tree and the input data as arguments.

<sub>234</sub> This triggers the creation of the internal objects of the framework, i.e. an `OrderedTree`

<sub>235</sub> object maintaining the order in which the nodes are processed and a `PreOrderTraversal`

<sub>236</sub> or a `PostOrderTraversal` object implementing different parallelization modes of the two

<sub>237</sub> traversal types. In the ideal use-case, the `TraversalTask`'s `TraverseTree()` method will

<sub>238</sub> be called repeatedly, varying the model parameters, the input data and branch lengths on a

<sub>239</sub> fixed tree topology. This encompasses all scenarios where a model is fitted to a fixed tree

<sub>240</sub> and data, e.g. ML or Bayesian PCM inference. In the current version of the framework any

<sub>241</sub> change in the tree topology invalidates the associated `TraversalTask` object, so a new one

<sub>242</sub> has to be created. Thus, in the case of inferring a tree topology, the overall performance

<sub>243</sub> would depend on the frequency of changes in the topology per traversal and the time for

<sub>244</sub> re-creating the `TraversalTask` object relative to the total time of one traversal.


<sub>245</sub> ## *Examples*


<sub>246</sub> In this section, we show example usages of the parallel traversal framework. In each

<sub>247</sub> of these examples, we solve a particular problem, such as calculating the likelihood of a

<sub>248</sub> continuous Markov model for a categorical or a continuous trait. In terms of the

<sub>249</sub> framework, the task boils down to formulating the node states $S_j(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta})$ and the recursive

<sub>250</sub> functions $R_j$ satisfying rules (1) and (2).

<sub>251</sub> *Example 1: Models of categorical trait evolution.—*

<sub>252</sub> Consider a trait taking values in $\{0, 1\}$ evolving independently along the lineages of

<sub>253</sub> a phylogenetic tree, $\mathcal{T}$ with branch lengths $\mathbf{t}$. A continuous-time Markov model can be

<sub>254</sub> used to characterize the transitions of the trait value along each branch (Felsenstein 1983;

<sub>255</sub> Pagel 1994). This model assumes constant rates of change from 0 to 1, $q_{01}$ and from 1 to 0,

<sub>256</sub> $q_{10}$, representing the probability that the change has occurred during an infinitesimal

14

interval of time. These rates are used to define a rate matrix:

$$\mathbf{Q} = \begin{bmatrix} -q_{01} & q_{01} \\ q_{10} & -q_{10} \end{bmatrix}. \tag{3}$$

Given $\mathbf{Q}$, the transition probability matrix $\mathbf{P}(t)$ for an arbitrary long period $t$ is given by

$$\mathbf{P}(t) = \begin{bmatrix} P_{00}(t) & P_{01}(t) \\ P_{10}(t) & P_{11}(t) \end{bmatrix} = \mathbf{C} \begin{bmatrix} e^{\lambda_1 t} & 0 \\ 0 & e^{\lambda_2 t} \end{bmatrix} \mathbf{C}^{-1} \tag{4}$$

where $\lambda_i$ are the eigenvalues of $\mathbf{Q}$ and $\mathbf{C}$ is a matrix, which's $i^{th}$ column represents the $i^{th}$ eigenvector of $\mathbf{Q}$ (Pagel 1994). Assuming that the value at the root is known to be $z_M$, we want to find the probability with which the model specified by the parameters $\mathbf{\Theta} = (q_{01}, q_{10})$ generates an $N$-vector of values, $\mathbf{z}$ observed at the tips. This represents the conditional likelihood $L_{\mathcal{T}}(\mathbf{t}, \mathbf{z}, \mathbf{\Theta}, z_M)$. The pruning algorithm for calclulating $L$ relies on calculating the "fragmentary" likelihood $L_i(b) = P(\mathbf{z}_i|z_i = b; \mathbf{\Theta})$ for each node $i$ and each $b \in \{0, 1\}$ (Felsenstein 1983). In terms of the framework, we define the state $S_j(\mathbf{t}, \mathbf{z}, \mathbf{\Theta})$ of a node $j$ as the pair $< L_j(0), L_j(1) >$. Following eq. 4 in (Felsenstein 1983), the recursive $R_j$ are given by:

$$S_j(\mathbf{t}, \mathbf{z}, \mathbf{\Theta}) = \begin{cases} \langle \delta(z_j = 0), \delta(z_j = 1) \rangle & \text{if } j \text{ is a tip} \\ \left\langle \prod_{i \in Desc(j)} \left[ \sum_{z_i} P_{0z_i}(t_i) L_i(z_i) \right], \prod_{i \in Desc(j)} \left[ \sum_{z_i} P_{1z_i}(t_i) L_i(z_i) \right] \right\rangle & \text{if } j \text{ is internal,} \end{cases} \tag{5}$$

where we use the Kronecker delta function $\delta(x = y)$ equalling to 1 if $x = y$ and 0, otherwise. In the above eqation 5, the values $L_i(z_i)$ are available from the descendants'

states $S_i$. Finally, the conditional likelihood $L_{\mathcal{T}}(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}, z_M)$ is given by $L_M(z_M)$, which is one of the two members in $S_M$.

The above model can be extended to a multivariate case, such as calculating the probability of a nucleotide or aminoacid sequence alignment as is the case in (Felsenstein 1983). Suppose that there are $p$ nucleotide sites, which are evolving independently. Then, the state for a node $j$ would represent a $p \times 4$ matrix

$$S_j(\mathbf{t}, \mathbf{z}, \boldsymbol{\Theta}) = \begin{bmatrix} L_j^{(1)}(A) & L_j^{(1)}(C) & L_j^{(1)}(T) & L_j^{(1)}(G) \\ \vdots & \vdots & \vdots & \vdots \\ L_j^{(p)}(A) & L_j^{(p)}(C) & L_j^{(p)}(T) & L_j^{(p)}(G) \end{bmatrix}, \tag{6}$$

where the letters $A$, $C$, $T$ and $G$ denote the nucleotides and the superscript in parentheses denotes a site in the alignment. To define the recursive functions $R_j$, equation 5 can be extended to accomodate one row of $S_j$ (four possible values instead of two) and evaluated $p$ times to obtain the full state.

The model can also be extended to support correlated evolution between the sites. As shown in (Pagel 1994), this involves extending the rate matrix $\mathbf{Q}$ to embed transition rates between pairs, triplets or higher order combinations of sites in the sequence. Accounting for correlated evolution between combinations of sites dramatically increases the computational complexity, but does not present a conceptual change from the point of view of the pruning operation. Thus, accommodating such models in the framework, although involved technically, should not present a conceptual challenge.

We should not omit mentioning other software libraries implementing parallel likelihood computation of different Markov models of sequence evolution. For example, several high level tools for ML and Bayesian tree inference, e.g. Drummond et al. (2012); Bouckaert et al. (2014); Ronquist and Huelsenbeck (2003), use the library BEAGLE which distributes the computation for the independent sites of the sequence alignment among

16

293 multiple CPU or GPU cores (Ayres et al. 2012). SPLiTTree operates on a different level,

294 namely, it parallelizes the computation for independent lineages in the tree. Both

295 approaches are interesting because they fit well to different sizes of the input data - while

296 BEAGLE achieves significant parallel speed-ups in long alignments comprising many

297 thousands nucleotide or codon columns (Ayres et al. 2012), SPLiTTree is better suited to

298 shorter alignments of potentially many thousands of species.

299 *Example 2: Models of continuous trait evolution.—*

300 Ho and Ané (2014) noticed that the computational complexity in multivariate

301 Gaussian and some non-Gaussian models concentrates in the calculation of determinants

302 $|\mathbf{V_\Theta}|$ and quadratic quantities of the form $\mathbf{Q_\Theta} = \mathbf{X'_\Theta} \mathbf{V_\Theta^{-1}} \mathbf{Y_\Theta}$, where $\mathbf{V_\Theta}$ represents the

303 variance covariance matrix expected under the model specified by $\mathbf{\Theta}$ and the matrices $\mathbf{X_\Theta}$

304 and $\mathbf{Y_\Theta}$ represent centered observed data at the tips in the tree. For example, in the case

305 of Brownian motion and Ornestein-Uhlenbeck models, the log-likelihood function is equal

306 to the log-density of a multivariate Gaussian distribution:

$$\ln f(\mathbf{z}|\mathbf{\Theta}) = -\frac{1}{2} \left( N \ln(2\pi) + \ln |\mathbf{V_\Theta}| + (\mathbf{z} - \mu_\mathbf{\Theta})' \mathbf{V_\Theta^{-1}} (\mathbf{z} - \mu_\mathbf{\Theta}) \right), \qquad (7)$$

307 where $\mathbf{V_\Theta} = \mathbf{\Sigma}$ and $\mathbf{V_\Theta} = \mu$ (table 1).

308 Ho and Ané (2014) developed a pruning algorithm which allows to calculate $|\mathbf{V_\Theta}|$

309 and $\mathbf{Q_\Theta}$ simultaneously and without constructing or allocating the matrix $\mathbf{V_\Theta}$ in memory,

310 provided $\mathbf{V_\Theta}$ has a "3-point structure". Then, they showed several examples of Gaussian

311 models such as Brownian motion and Ornestein-Uhlenbeck, as well as non-Gaussian models,

312 such as phylogenetic logistic and Poisson regression, where $\mathbf{V_\Theta}$ is or can be "converted" to

313 a 3-point structured matrix (discussed later). Adapting the notation from (Ho and Ané

314 2014, p. 399), we define the node states as $S_j(\mathbf{t}, \mathbf{z}, \mathbf{\Theta}) = \langle p_{A,j}, p_j, \hat{\mu}_{Y,j}, \tilde{\mu}'_{X,j}, \ln |\mathbf{V}|_j, \mathbf{Q}_j \rangle$.

315 The recursive functions $R_j$ follow immediately from points 1 and 2 of the algorithm (Ho

17

Table 1: Population properties at the tips of the phylogeny under BM and OU models and their mixed counterparts The acronyms are: PBM - Phylogenetic Brownian motion (without non-heritable component); PMM - Phylogenetic Mixed Model (adding a non-heritable component to PBM); POU - Phylogenetic Ornstein-Uhlenbeck (without non-heritable component), also known as "Hansen's model" or Single Stationary Peak (SSP); POUMM - Phylogenetic Ornstein-Uhlenbeck Mixed Model (adding a non-heritable component to the POU model. Expressions for the OU-models were adapted from (Hansen 1997). $\mu_{\mathbf{\Theta},i}$: expected value at tip $i$; $\mathbf{V}_{\mathbf{\Theta},ii}$: expected variance for tip $i$; $\mathbf{V}_{\mathbf{\Theta},ij}$: expected covariance of the values of tips $i$ and $j$.

| | PBM | PMM | POU | POUMM |
|---|---|---|---|---|
| $\mathbf{\Theta}$: | $< g_M, \sigma >$ | $< g_M, \sigma, \sigma_e >$ | $< g_M, \alpha, \theta, \sigma >$ | $< g_M, \alpha, \theta, \sigma, \sigma_e >$ |
| $\mu_{\mathbf{\Theta},i}$: | $g_M$ | $g_M$ | $e^{-\alpha h_i} g_M + \left(1 - e^{-\alpha h_i}\right)\theta$ | $e^{-\alpha h_i} g_M + \left(1 - e^{-\alpha h_i}\right)\theta$ |
| $\mathbf{V}_{\mathbf{\Theta},ii}$: | $\sigma^2 h_i$ | $\sigma^2 h_i + \sigma_e^2$ | $\frac{\sigma^2}{2\alpha}\left(1 - e^{-2\alpha h_i}\right)$ | $\frac{\sigma^2}{2\alpha}\left(1 - e^{-2\alpha h_i}\right) + \sigma_e^2$ |
| $\mathbf{V}_{\mathbf{\Theta},ij}$: | $\sigma^2 h_{(ij)}$ | $\sigma^2 h_{(ij)}$ | $\frac{\sigma^2}{2\alpha} e^{-\alpha d_{ij}}\left(1 - e^{-2\alpha h_{(ij)}}\right)$ | $\frac{\sigma^2}{2\alpha} e^{-\alpha d_{ij}}\left(1 - e^{-2\alpha h_{(ij)}}\right)$ |

316  and Ané 2014):

$$
\begin{cases}
S_j(\mathbf{t}, \mathbf{z}, \mathbf{\Theta}) = \langle\ p_{A,j} = 0, \\
\qquad p_j = \frac{1}{t_j}, \\
\qquad \hat{\mu}_{Y,j} = \mathbf{y}_{\mathbf{\Theta},j}, \\
\qquad \tilde{\mu}'_{X,j} = \mathbf{x}'_{\mathbf{\Theta},j}, \qquad\qquad\qquad \text{if } j \leq N \\
\qquad \ln|\mathbf{V}|_j = \ln t_j, \\
\qquad \mathbf{Q}_j = \mathbf{x}'_{\mathbf{\Theta},j}\mathbf{y}_{\mathbf{\Theta},j}\ \rangle \\
S_j(\mathbf{t}, \mathbf{z}, \mathbf{\Theta}) = \langle\ p_{A,j} = \sum_{i \in Desc(j)} p_i, \\
\qquad p_j = \frac{p_{A,j}}{1 + t_j p_{A,j}}, \\
\qquad \hat{\mu}_{Y,j} = \sum_{i \in Desc(j)} \frac{p_i}{p_A}\hat{\mu}_{Y,i}, \\
\qquad \tilde{\mu}'_{X,j} = \sum_{i \in Desc(j)} \frac{p_i}{p_A}\tilde{\mu}'_{X,i}, \qquad\qquad \text{otherwise.} \\
\qquad \ln|\mathbf{V}|_j = \sum_{i \in Desc(j)} \ln|\mathbf{V}|_i + \ln(1 + t_j p_{A,j}), \\
\qquad \mathbf{Q}_j = \sum_{i \in Desc(j)} \mathbf{Q}_i \frac{1}{2}\ln(1 + t_j p_{A,j}) \qquad\qquad \rangle
\end{cases}
\tag{8}
$$

The caveat in applying the 3-point algorithm is that except for BM models and OU models (ultrametric trees only), the matrix $\mathbf{V_\Theta}$ does not necessarily satisfy the 3-point condition (Ho and Ané 2014). As the authors show, it is still possible to use the algorithm in that case, provided that $\mathbf{V_\Theta}$ satisfies a "generalized 3-point condition" (Ho and Ané 2014). More precisely, in most of their examples, the authors showed that there exist a transformation of the branch lengths, $\tilde{\mathbf{t}}$, diagonal matrices $\mathbf{D}_1$ and $\mathbf{D}_2$ with non-zero diagonal elements and a 3-point structured matrix $\tilde{\mathbf{V}}_\Theta$, such that $\tilde{\mathbf{V}}_\Theta$ is equal to the variance-covariance on the tree $\tilde{\mathcal{T}}$ with the transformed branch lengths and $\mathbf{V_\Theta} = \mathbf{D}_1 \tilde{\mathbf{V}}_\Theta \mathbf{D}_2$. If so, the algorithm is applied to $\tilde{\mathbf{V}}_\Theta$ using $\tilde{\mathbf{t}}$ and transformed data $\tilde{\mathbf{X}} = \mathbf{D}_2^{-1} \mathbf{X}$, $\tilde{\mathbf{Y}} = \mathbf{D}_1^{-1} \mathbf{Y}$. Then the quadratic form of interest, $\mathbf{Q_\Theta}$, would be equal to the resulting quadratic form at the root, $\mathbf{Q}_M$ and the determinant $|\mathbf{V_\Theta}|$ is obtained by the formula:

$$|\mathbf{V_\Theta}| = |\mathbf{D}_1||\tilde{\mathbf{V}}_\Theta||\mathbf{D}_2| \tag{9}$$

We note that finding a suitable transformation of the branch lengths remains a model specific task. We give an example of such branch transformation in the next showcase, referring the reader to (Ho and Ané 2014) for further examples.

# A showcase: The phylogenetic Ornstein-Uhlenbeck mixed model

Here, we describe a phylogenetic Ornstein-Uhlenbeck mixed model of continuous trait evolution, which we and other authors have used previously to analyze the evolution of set-point viral load in HIV patients (Mitov and Stadler 2016; Bachmann et al. 2017; Bertels et al. 2017; Blanquart et al. 2017). To calculate the likelihood, we will use the fact that the variance covariance matrix of the POUMM has a generalized 3-point structure.

19

339    We will use a simulation based method to validate the technical correctness of the model.

340    Then, we will report a benchmark comparing the times for likelihood calculation between

341    different serial and parallel pruning implementations in the SPLiTTree library, as well as

342    two third-party serial pruning implementations of the same likelihood calculation. Finally,

343    we will report the performance gain for the Bayesian inference of the model upon

344    combining parallel pruning with adaptive Metropolis sampling on a real dataset.

## *The model*

345

346    Consider a real valued trait evolving independently along the lineages of a phylogenetic

347    tree, $\mathcal{T}$ with branch lengths $\mathbf{t}$. The phylogenetic Ornstein-Uhlenbeck mixed model

348    (POUMM) decomposes the trait value as a sum of a non-heritable component, $e$, and a

349    genetic component, $g$, which (i) evolves continuously according to an Ornstein-Uhlenbeck

350    (OU) process along branches; (ii) gets inherited by the branches descending from each

351    internal node. In biological terms, $g$ is a genotypic value (Lynch and Walsh 1998) that

352    evolves according to random drift with stabilizing selection towards a global optimum; $e$ is

353    a non-heritable component, which can be interpreted in different ways, depending on the

354    application, i.e. a measurement error, an environmental contribution, a residual with

355    respect to a model prediction, or the sum of all these. The OU-process acting on $g$ is

356    parameterized by an initial genotypic value at the root, $g_M$, a global optimum, $\theta$, a

357    selection strength, $\alpha > 0$, and a random drift unit-time standard deviation, $\sigma$. Denoting by

358    $W_t$ the standard Wiener process (Grimmett and Stirzaker 2001), the evolution of the

359    trait-value, $z(t)$, along a given lineage of the tree is described by the equations:

$$z(t) = g(t) + e \tag{10}$$

$$dg(t) = \alpha[\theta - g(t)]dt + \sigma dW_t \tag{11}$$

$$g(0) = g_M, \tag{12}$$

360 The stochastic differential equation 11 defines the OU-process, which represents a random

361 walk tending towards the global optimum $\theta$ with stronger attraction for bigger difference

362 between $g(t)$ and $\theta$ (Ornstein and Zernike 1919; Uhlenbeck and Ornstein 1930). The model

363 assumptions for $e$ are that they are independent and identically distributed (i.i.d.) normal

364 with mean 0 and standard deviation $\sigma_e$ at the tips. Any process along the tree that gives

365 rise to this distribution at the tips may be assumed for $e$. For example, in the case of

366 epidemics, a newly infected individual is assigned a new $e$-value which represents the

367 contribution from its immune system and this value can change or remain constant

368 throughout the course of infection. In particular, the non-heritable component $e$ does not

369 influence the behavior of the OU-process $g(t)$. Thus, if we were to simulate trait values $z$

370 on the tips of a phylogenetic tree $\mathcal{T}$, we could first simulate the OU-process from the root

371 to the tips to obtain $g$, and then add the white noise $e$ (i.e. an i.i.d. draw from a normal

372 distribution) to each simulated $g$ value at the tips. The POUMM represents an extension

373 of the phylogenetic mixed model (PMM) (Lynch 1991; Housworth et al. 2004), since, in the

374 limit $\alpha \to 0$, the OU-process converges to a Brownian motion (BM) with unit-time

375 standard deviation $\sigma$. Both, the POUMM and the PMM, define an expected multivariate

376 normal distribution for the trait values at the tips. The mean vectors and the

377 variance-covariance matrices of these distributions are written in table 1. Note that the

378 trait expectation and variance for a tip $i$ depends on its height $(h_i)$, and the trait

379 covariance for a pair of tips $(ij)$ depends on the height of their mrca $(h_{(ij)})$, and, in the

21

380 case of POUMM, on their patristic distance ($d_{ij}$) (table 1).

381 ## *Calculating the likelihood*

382 The POUMM likelihood is defined as the multivariate probability density of an

383 observed vector $\mathbf{z}$ at the tips of $\mathcal{T}$ for given model parameters $\boldsymbol{\Theta} = < g_M, \alpha, \theta, \sigma, \sigma_e >$:

$$\ell\ell(\boldsymbol{\Theta}) = \ln(f(\mathbf{z}|\mathcal{T}, \mathbf{t}, \boldsymbol{\Theta})). \tag{13}$$

384 The probability density function, $f$ is multivariate Gaussian with mean vector $\mu_{\boldsymbol{\Theta}}$

385 and variance-covariance matrix $\mathbf{V}_{\boldsymbol{\Theta}}$ written in table 1. Since $\mathbf{V}_{\boldsymbol{\Theta}}$ has a generalized 3-point

386 structure (Ho and Ané 2014), we can apply the recursion in eq. 8, upon a transformation

387 of the branch lengths and the data. This is obtained through adapting the transformation

388 for an non-mixed OU-model in a ultrametric tree (Ho and Ané 2014) to accommodate the

389 non-heritable variance:

$$\tilde{t}_i = \frac{\sigma^2}{2\alpha} \left[ e^{2\alpha T} \left( e^{2\alpha h_i} - e^{2\alpha h_{Parent(i)}} \right) \right] + \frac{\sigma_e^2}{e^{2\alpha u_i}} \delta(i \leq N) \quad \text{for } i \in \{1, ..., M-1\} \tag{14}$$

390

$$\tilde{\mathbf{X}}_i = \tilde{\mathbf{Y}}_i = \frac{z_i - \mu_i}{e^{\alpha u_i}} \quad \text{for } i \in \{1, ..., N\}, \tag{15}$$

391 where $\tilde{\mathbf{X}}$ and $\tilde{\mathbf{Y}}$ are identical $N$-vectors, $T$ is the maximum tip-height in the tree and

392 $u_i = T - h_i$ for $i \in \{1, ..., N\}$. After running the post-order traversal, using eq. 8 as a

393 visit-node operation, we apply eq. 9, to obtain $|\mathbf{V}_{\boldsymbol{\Theta}}|$ and eq. 7 to obtain the log-likelihood.

394 We note that the branch transformation (eq. 14) can be done "locally" on every

395 branch, using pre-calculated heights of the parent and daughter nodes connected by the

396 branch. Thus, it is safe to include the transformation in the visit-node operation and the

22

<sup>397</sup> parallelization of pruning would not suffer. Otherwise, the transformation would have had

<sup>398</sup> to be done in a preprocessing step. Again, this is a model specific consideration.

### *Combined maximum likelihood and Bayesian inference of the model*

<sup>400</sup>     We implement maximum likelihood and Bayesian inference of the POUMM

<sup>401</sup> parameters, $\Theta$, using the L-BFGS-R convex optimization algorithm (R-function `optim`)

<sup>402</sup> and a variant of the Random Walk Metropolis (RWM) Markov Chain Monte Carlo

<sup>403</sup> (MCMC) sampling (Metropolis et al. 1953). This combined inference capitalizes on two

<sup>404</sup> practical ideas:

<sup>405</sup>   • A MCMC has higher chance to converge to the target posterior distribution faster if

<sup>406</sup>      it has been started from a previously estimated MLE;

<sup>407</sup>   • If an MCMC encounters a point in the parameter space that has higher likelihood

<sup>408</sup>      than a previously inferred MLE, running maximum likelihood optimization from that

<sup>409</sup>      point is more likely to find the global likelihood optimum.

<sup>410</sup>     An important step in RWM is the choice of a proposal (jump) distribution shape

<sup>411</sup> matrix used as a scaling factor on each next proposal in the Metropolis algorithm.

<sup>412</sup> Choosing the shape matrix with respect to the scale and the correlation structure of the

<sup>413</sup> parameter space minimizes the number of iterations needed for MCMC convergence and

<sup>414</sup> mixing. Thus, numerous variants of the RWM have been proposed, performing "on-the-fly"

<sup>415</sup> adaptation of the shape matrix based on what has been "learned" about the parameter

<sup>416</sup> space from the past RWM iterations (Haario et al. 2001; Vihola 2012). Of these variants,

<sup>417</sup> we chose the adaptive Metropolis sampling with coerced acceptance rate, because it is

<sup>418</sup> shown to be robust with respect to the posterior distribution, it performs a relatively cheap

<sup>419</sup> adaptation of the shape (Vihola 2012) and it has an implementation in the R within the

<sup>420</sup> package `adaptMCMC` (Scheidegger 2017).

23

<sub>421</sub> The fitting of the POUMM model was implemented as a pipeline including the

<sub>422</sub> following steps:

<sub>423</sub> 1. Perform three MLE searches using the R-function `optim` and the L-BFGS-B method

<sub>424</sub> (Byrd et al. 1995), starting from three randomly chosen points in parameter space;

<sub>425</sub> 2. Run three MCMC chains as follows: (i) a chain sampling from the prior distribution;

<sub>426</sub> (ii) a chain sampling from the posterior distribution and started from the MLE found

<sub>427</sub> in step 1; (iii) a chain sampling from the posterior distribution and started from a

<sub>428</sub> random point in parameters space.

<sub>429</sub> 3. If the parameter tuple of highest likelihood sampled in the MCMC has a likelihood

<sub>430</sub> higher than the MLE found in step 1, repeat the MLE search starting from that

<sub>431</sub> parameter tuple;

<sub>432</sub> By running MLE first and starting an MCMC chain from the MLE candidate, we

<sub>433</sub> increase the chance that at least one of the MCMCs would converge faster to the posterior

<sub>434</sub> distribution. By comparing the posterior samples from two MCMCs initiated from

<sub>435</sub> different starting points, it can be assessed whether the MCMCs have converged to the true

<sub>436</sub> posterior. We do this quantitatively by the use of the Gelman-Rubin convergence

<sub>437</sub> diagnostic (Brooks and Gelman 1998) implemented in the R-package coda (Plummer et al.

<sub>438</sub> 2006). Values of the Gelman-Rubin (G.R.) statistic significantly different from 1 indicate

<sub>439</sub> that at least one of the two posterior samples deviates significantly from the true posterior

<sub>440</sub> distribution. By visual comparison of posterior density with prior desnity plots, it is

<sub>441</sub> possible to assess whether the data contains information differring from the prior

<sub>442</sub> knowledge for a given parameter. In step 3, we capitalize on the chance that the MCMCs

<sub>443</sub> have explored a wider region of the parameter space than the likelihood optimization.

<sub>444</sub> *The POUMM R-package*

24

445 We implement the model in the form of an R-package called POUMM, which

446 embeds the SPLiTTree library as an Rcpp module. Before model fitting, the user can

447 choose from different POUMM parametrizations and prior settings (function

448 `specifyPOUMM`). A set of standard generic functions, such as `plot`, `summary`, `logLik`, `coef`,

449 etc., provide means to assess the quality of a fit (i.e. MCMC convergence, consistence

450 between ML and MCMC fits) as well as various inferred properties, such as high posterior

451 density (HPD) intervals (more details in the package user guide).

## *Technical correctness*

453 To validate the correctness of the Bayesian POUMM implmentation, we used the

454 method of posterior quantiles (Cook et al. 2006). In this method, the idea is to generate

455 samples from the posterior quantile distributions of selected model parameters (or

456 functions thereof) by means of numerous "replications" of simulation followed by Bayesian

457 parameter inference. In each replication, "true" values of the model parameters are drawn

458 from a fixed prior distribution and trait-data is simulated under the model specified by

459 these parameter values. We perform these simulations on a fixed tree of size $N = 4000$.

460 Then, the to-be-tested software is used to produce a posterior distribution of parameters

461 based on the simulated trait-data. Next, the posterior quantiles of the "true" parameter

462 values (or functions thereof) are calculated from the corresponding posterior samples

463 generated by the to-be-tested software. By running multiple independent replications on a

464 fixed prior, it is possible to generate large samples from the posterior quantile distributions

465 of the individual model parameters, as well as any derived quantities. Assuming correctness

466 of the simulations, any statistically significant deviation from uniformity of these posterior

467 quantile samples indicates an error in the to-be-tested software (Cook et al. 2006).

468 Two phylogenetic trees were used for the simulations:

469    • Ultrametric (BD, $N = 4000$) - an ultrametric birth-death tree of 4000 tips generated

470      using the TreeSim R-package (Stadler et al. 2013; Boskova et al. 2014) (function call:

471      `sim.bd.taxa(4000, lambda = 2, mu = 1, frac = 1, complete = FALSE)`);

472    • Non-ultrametric (BD, $N = 4000$) - a non-ultrametric birth-death tree of 4000 tips

473      generated using the TreeSim R-package (Stadler et al. 2013; Boskova et al. 2014)

474      (function call: `sim.bdsky.stt(4000, lambdasky = 2, deathsky = 1,`

475      `timesky=0)`).

476      Simulation scenarios of 2000 replications were run using the prior distribution

477 $< g_M, \alpha, \theta, \sigma, \sigma_e > \sim \mathcal{N}(5, 25) \times \text{Exp}(0.1) \times \mathcal{U}(2, 8) \times \text{Exp}(0.4) \times \text{Exp}(1)$. The goal of using

478 this prior was to explore a large enough subspace of the POUMM parameter space, while

479 keeping MCMC convergence and mixing within reasonable time (runtime up to 30 minutes

480 for two MCMCs of $10^6$ adaptive Metropolis iterations at target acceptance rate of 1%).

481 From the above prior, we drew a sample of $n = 2000$ parameter tuples, $\{\mathbf{\Theta}^{(1)}, ..., \mathbf{\Theta}^{(n)}\}$,

482 which were used as replication seeds. For a given $\mathbf{\Theta}^{(i)}$, we simulate genotypic values

483 $\mathbf{g}^{(i)}(\mathcal{T}, \mathbf{\Theta}^{(i)})$ according to an OU-branching process with initial value $g_M^{(i)}$ and parameters

484 $\alpha^{(i)}$, $\theta^{(i)}$, $\sigma^{(i)}$. Then, we add random white noise ($\sim \mathcal{N}(0, \sigma_e^{2(i)})$) to the genotypic values at

485 the tips, to obtain the final trait values $\mathbf{z}^{(i)}$.

486      For the two simulated trees, we executed a total of $2 \times 2000 = 4000$ replications.

487 The resulting posterior quantile distributions for the each tree are shown on Fig. 3. We

488 notice that the posterior quantiles for all relevant parameters are uniformly distributed.

489 This is confirmed visually by the corresponding histograms (fig. 3), as well as statistically,

490 by a non-significant p-value from a Kolmogorov-Smirnov uniformity test at the 0.01 level.

491 This observation validates the technical correctness of the software.
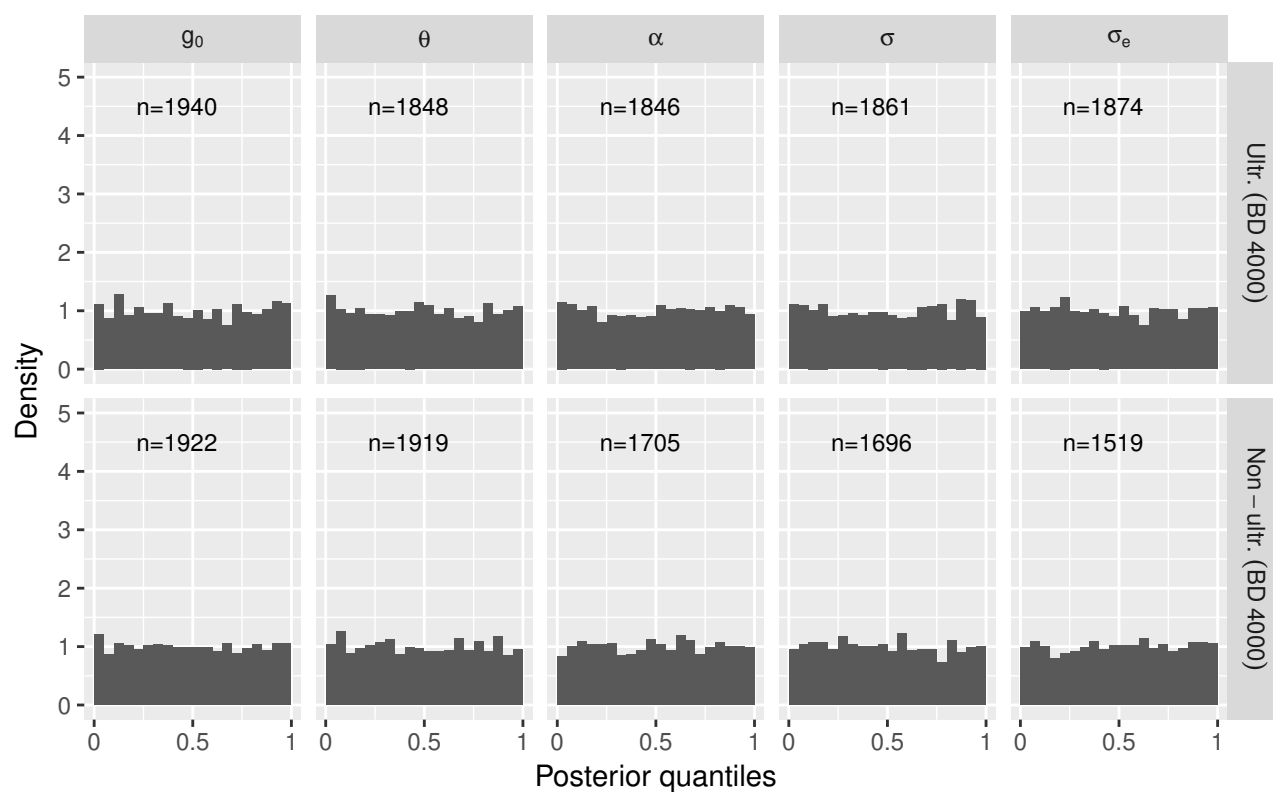
492                   *Performance on simulated data*

Figure 3: Posterior quantiles from simulations on a ultrametric and a non-ultrametric tree ($N = 4000$). The number $n$ at the top of each histogram denotes the number of replications out of 2000 which reached acceptable MCMC convergence and mixing after one million iterations. Uniformity was confirmed using a Kolmogorov-Smirnov test which was insignificant for all parameters (P-value above 0.1).

493  We ran a performance benchmark on a MacBook Pro laptop (Retina, 15-inch, Late

494  2013) running a 2.3GHz Intel(R) Core i7 processor with 4 physical cores. We used the

495  R-package `apTreeshape` (Bortolussi et al. 2012) to generate tree topologies of sizes

496  $N \in \{100, 1000, 10,000, 100,000\}$. To generate the trees, we used the function

497  `rtreeshape()` with a `biased` model. A parameter $p$ in this model controls the

498  disproportion of branching rates for the left and right lineages starting from a given parent

499  node. For each $N$, we used four settings for $p$ as follows:

500  1. $p = 0.5$ corresponding to equal left and right branching rates and resulting in

501  balanced trees;

502  2. $p = 0.1$ corresponding to unbalanced trees in which one of any two sibling branches

503  (sharing the same parent node) splits at rate $p = 0.1$, while the other splits at rate

504  $p' = 1 - p = 0.9$ (time units are arbitrary, so we can assume that the rates correspond

505  to splitting probabilities per unit time).

506  3. $p = 0.01$ corresponding to very unbalanced trees (splitting rates of $p = 0.01$ and

507  $p' = 0.99$ for any couple of sibling branches;

508  4. $p = 0.01/N$ corresponding to a ladder-like tree (see fig. 1b).

509  This resulted in a total of 16 tree topologies. For each topology, random branch

510  lengths were assigned overwriting the default branch lengths of 1 assigned by

511  `rtreeshape()`. For each tree, we generated random trait-values using random parameters

512  of the POUMM model.

513  We compared serial and parallel likelihood calculation within the POUMM package

514  to serial pruning implementations provided in the R-packages geiger (Pennell et al. 2014),

515  and diversitree (FitzJohn 2012). All packages including C or C++ code were compiled

516  from source-code using the R-command `install.packages('package-directory',`

517  `repos=NULL, type='source')`, and the same C++ compiler and compiler arguments

518  (version 16.0.0 of the Intel compiler, command `icpc` with options `-O2 -march=native`).

519  *Time for preprocessing the tree.—*

520    Each of the tested packages implements a preprocessing step initializing cached

521  data-structures that are re-used during likelihood calculation. In the case of `POUMM`, this is

522  the constructor of the class `TraversalTask` (fig. 2); in the case of `diversitree`, this is the

523  function `make.ou`; in the case of `geiger`, this is the internal function `bm.lik`. We note that

524  the time for creating the cache structure is not important in scenarios of fitting

525  comparative models to a fixed tree and data (created once, at the beginning of the

526  inference process). These times become important in the case when the tree topology is

527  inferred together with the model parameters from trait and sequence alignment data.

528    We measured the preprocessing time on the 16 trees (table 2). The times scaled

529  linearly with the size of the tree for the packages `POUMM` and `diversitree`. For these two

530  packages the time was not affected by the unbalancedness of the tree. For `geiger`, we

531  observed higher time complexity, both in $N$ as well as in the unbalancedness (longer times

532  for unbalanced trees). For $N = 100,000$ and $p = 0.01/N$, both, `diversitree` and `geiger`

533  failed with a `stack-overflow` error. The relatively short `POUMM` times indicate that

534  `SPLiTTree` could potentially be used for phylogenetic inference models.

535  *Time for one POUMM likelihood calculation.—*

536    We distinguish the different implementations according to three criteria:

537  • Mode: denotes whether the implementation is single threaded using one physical core

538    of the CPU (serial) or multi-threaded, running 8 virtual threads on 4 physical cores

539    (parallel);

540  • Order: denotes the order in which the prune-able nodes are processed. We tested

29

Table 2: Times for tree-preprocessing in milliseconds.

| N | Implementation | p=0.5 | p=0.1 | p=0.01 | p=0.01/N |
|---|---|---|---|---|---|
| 100 | geiger | 5 | 6 | 9 | 9 |
| 100 | diversitree | 4 | 4 | 4 | 4 |
| 100 | POUMM | 2 | 2 | 2 | 1 |
| 1,000 | geiger | 18 | 26 | 78 | 414 |
| 1,000 | diversitree | 20 | 20 | 22 | 30 |
| 1,000 | POUMM | 3 | 2 | 3 | 3 |
| 10,000 | geiger | 358 | 449 | 1,345 | 355,396 |
| 10,000 | diversitree | 207 | 211 | 227 | 1,338 |
| 10,000 | POUMM | 14 | 13 | 13 | 15 |
| 100,000 | geiger | 20,215 | 21,629 | 36,349 | - |
| 100,000 | diversitree | 2,421 | 2,619 | 2,883 | - |
| 100,000 | POUMM | 130 | 131 | 131 | 140 |

three possible orders: postorder (Mode=serial only) - the nodes are processed sequentially without paying attention to their allocation in the memory - no SIMD operations are possible; queue-based (Mode=parallel only) - the nodes are processed according to their entering order in the queue (see algorithm 1) - no SIMD operations are possible, synchronized access to the queue; range-based (Mode=parallel only) - the nodes in each pruning generation are processed in order of their allocation in memory - SIMD operations are possible, no need for a synchronized access to a queue (see algorithm 2).

- Implementation: the R-package and the back-end used (R or C++).

To measure the likelihood calculation times we used the R-package microbenchmark (Mersmann 2015) with argument `times` set to 100. The resulting average times in milliseconds are shown on fig. 4.

On small trees of 100 tips, the fastest implementations were the serial C++ implementations from the packages POUMM and diversitree (about 0.03 ms); the POUMM range-based parallel implementation was nearly as fast on balanced trees ($p = 0.5$) but was progressively slower on unbalanced trees. The geiger implementation was
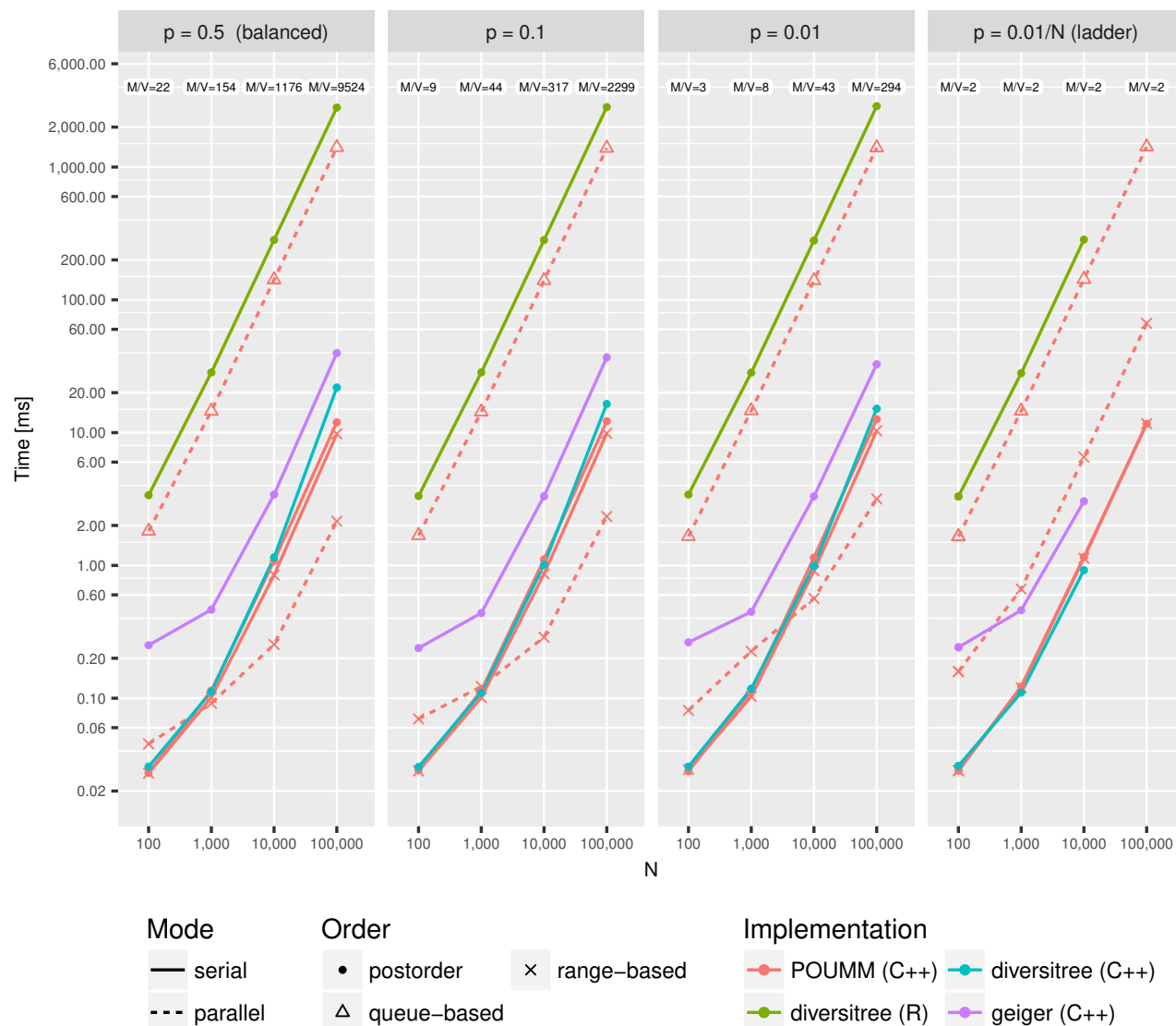
30

Figure 4: Likelihood calculation times for R and C++ implementations of the pruning algorithm. The labels $M/V$ on top denote the average number of nodes per generation (Visit-range) (see algorithm 2).

557 nearly an order of magnitude slower (0.2 ms). The POUMM queue-based parallel

558 implementation was nearly 100 times slower (nearly 2 ms), presumably due to the excessive

559 synchronization overhead. The serial R implementation from the diversitree package was

560 the slowest (above 2 ms), which was expected, since the R interpreter is notorious for its

561 slow speed compared to compiled languages like C++.

562 On bigger balanced trees ($N > 100$, $p = 0.5$), the POUMM range-based parallel

563 implementation took over, reaching up to $4\times$ speed-up with respect to the POUMM

564 range-based serial implementation, up to $6\times$ speed-up with respect to the POUMM

565 postorder serial implementation and up to $10\times$ speed-up with respect to the diversitree

566 serial C++ implementation. This reveals a consistent speed-up from SIMD operations for

567 all trees except the ladder tree, where parallelization of the internal nodes is not possible

568 (see fig. 1b). The time for the other serial implementations and the POUMM queue-based

569 parallel implementation scaled up by a factor of 10 for each higher value of $N$.

570 There was no significant difference in the times for the serial implementations and

571 the queue-based parallel implementation when comparing their performance on balanced

572 versus unbalanced trees. For the POUMM range-based parallel implementation, though,

573 the parallel speed-up was progressively less pronounced, in particular, for $N < 1000$ and for

574 ladder trees. Still, the speed-up was very good for $N \geq 10,000$ and $p \geq 0.01$.

575 *Improved MCMC convergence and MLE inference through adaptive Metropolis sampling.—*

576 To measure the MCMC convergence speed-up from the adaptive Metropolis

577 sampling, we reran one simulation scenario (2000 replications on a non-ultrametric tree of

578 4000 tips) with disabled adaptation. As a criterion for convergence, we used the absolute

579 difference from 1 of the Gelman-Rubin convergence diagnostic (Brooks and Gelman 1998)

580 (the closer $|G.R. - 1|$ is to 0, the better the convergence). When enabling adaptive

581 Metropolis sampling, more than 1600 (80%) of the 2000 replications had reached

582 $|G.R. - 1| < 0.01$ after a million iterations. Conversely, when disabling adaptive Metropolis

583 sampling, less than 300 (15%) of the replications had reached $|G.R. - 1| < 0.01$ after a

584 million iterations (the 80% quantile of $|G.R. - 1|$ was equal to 0.57, indicating very power

585 convergence). We also noticed that 1455 out of 2000 replications (73%) of the POUMM

586 inferences with enabled adaptative Metropolis sampling resulted in an improved MLE after

587 running the MCMC chains, compared to 1045 (50%) when disabling adaptation. These

588 observations show that adaptive Metropolis sampling considerably accelerates the MCMC

589 convergence towards the posterior distribution and can be used to improve the MLE

590 inference when using a weak prior or a prior that does not strongly contradict with the

591 evidence (likelihood).

## *Performance on real data*

592

593 This showcase would be incomplete if we don't provide an assessment of the performance of

594 the combined parallel likelihood and adaptive MCMC approach on a real dataset. We have

595 used the POUMM package to estimate the heritability of set-point viral load in a data-set

596 of 8,483 HIV patients. While the results of this analysis have been reported elsewhere

597 (Mitov and Stadler 2016), here, we briefly report the times and the quality statistics for the

598 MCMC inference of the model with and without adaptation.

599 First, we ran the classical RWM Metropolis sampler with a default identity shape

600 matrix for two MCMCs of ten million iterations on the above-mentioned hardware (2.3GHz

601 Intel(R) Core i7 processor with 4 cores), using the fastest (range-based) parallel likelihood

602 calculation. The total time for the two MCMCs was 3:18 hours. The run resulted in poor

603 mixing and very low effective posterior sample size for most of the inferred parameters of

604 the model (fig. 5a,b). The Gelman-Rubin statistic was greater than 1.1 for all parameters

605 and the effective sample size was below 400 for all parameters, falling below 50 for $\alpha$ and $\sigma$.
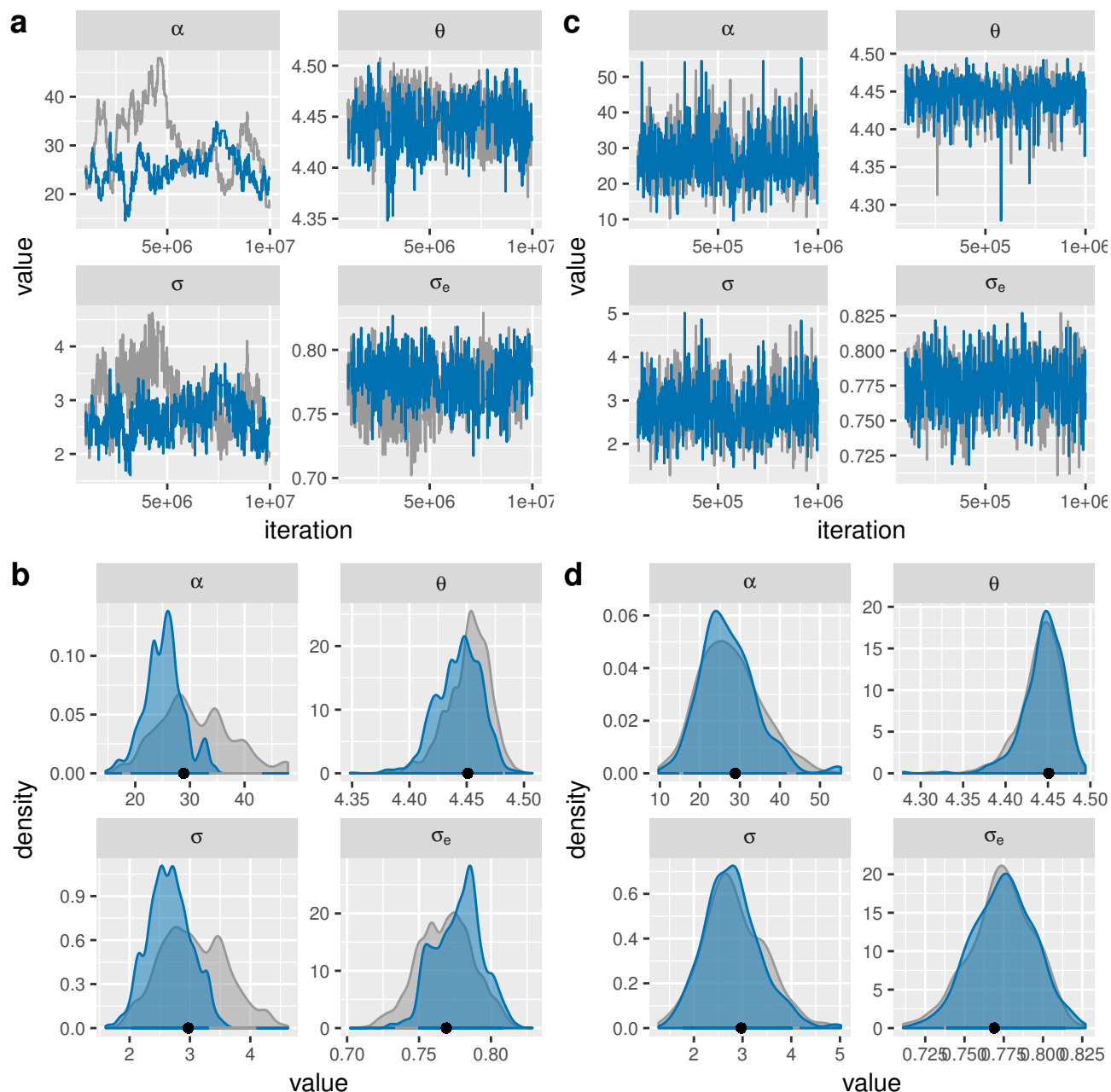
33

Figure 5: Sample trace- and density plots from a POUMM fit to a tree and virulence data 8483 HIV patients (Mitov and Stadler 2016) a,b: no adaptation of the proposal shape matrix (ten million iterations); c,d: on-the-fly adaptation of the proposal shape matrix from the first 100,000 out of one million iterations. The colors correspond to the different chains.

34

606      Next, we ran the adaptive Metropolis sampler for two MCMCs of one million

607  iterations. Adaptations has been enabled only for the first 100,000 iterations in each

608  MCMC. The total runtime was 25 minutes. The two chains mixed very well and the

609  effective sample size for all parameters exceeded 1200 (fig. 5c,d). The difference $|G.R. - 1|$

610  was below 0.01 for all parameters, proving that the MCMCs have converged to the same

611  distribution, which is very likely the true posterior distribution for the model parameters.

## Discussion

613      SPLiTTree can in principle be used for any algorithm that runs a pre-order or

614  post-order tree traversal. However, our benchmarks have shown that a performance gain

615  from parallelization will in most cases depend strongly on the application and the topology

616  of the tree. To cope with this issue, SPLiTTree implements both, serial as well as different

617  parallel modes of the tree traversal and can chose between these implementations during

618  the initial steps of the inference procedure.

619      The examples in this article focused on models of discrete and continuous trait

620  evolution. Another family of models where SPLiTTree could be used are population

621  models, e.g. models with structured populations. For example, when calculating the

622  likelihood for a phylogenetic tree under a structured birth-death model, the calculations

623  proceed in a pruning fashion (Kühnert et al. 2016) and may be improved with respect to

624  speed using our approach. However, the structured coalescent likelihood for a tree is a

625  function of all co-existing lineages even for approximate methods (Müller et al. 2017), and

626  thus a pruning formulation is not available.

627      We did not develop examples of pre-order traversal. One such example is the

628  simulation of traits evolving along the tree, which can be used for validation and

629  approximate inference of phylogenetic models. In complex phylogenetic comparative

630  models, where an exact calculation of the likelihood is elusive or computationally

35

631 intractable, it is possible to use simulations of trait evolution along the tree for

632 approximate likelihood calculation (Kutsukake and Innan 2013) or approximate Bayesian

633 computation (ABC) (Slater et al. 2012b). Both approaches are computationally intensive

634 and could benefit from parallel execution using our framework.

635 The POUMM R-package joins a growing collection of tools implementing phylogenetic

636 OU inference. Among others, these include the R-packages ape v4.0 (Paradis et al. 2004),

637 ouch v2.9-2 (Butler and King 2004), GLSME v1.0.3 (Hansen and Bartoszek 2012),

638 diversitree v0.9-9 (FitzJohn 2012), geiger v2.0.6 (Pennell et al. 2014), surface

639 v0.4-1 (Ingram and Mahler 2013), mvMORPH v1.0.8 (Clavel et al. 2015), bayou v1.0.0

640 (Uyeda et al. 2015), OUwie v1.50 (Beaulieu and OMeara 2016), phylolm v2.5 (Ho and

641 Ané 2014), RPANDA v1.1 (Manceau et al. 2016). Compared to these packages POUMM

642 provides fast Bayesian inference using the combined parallel likelihood and adaptive

643 sampling approach. Another feature of POUMM is that it implements different

644 re-parametrizations of $\Theta = <g_M, \alpha, \theta, \sigma, \sigma_e>$. This is helpful in Bayesian inference,

645 because it allows to express the prior distribution in application-specific terms, such as

646 phylogenetic heritability (Lynch 1991; Housworth et al. 2004).


## Outlook

647

648 The past decade has seen a rapid advance in the production of multi-core and SIMD

649 processors. At the same time, it appears that the maximum clock frequency of a single

650 processing unit is approaching the maximum achievable for semi-conductor based

651 architectures. This brings the need for development of novel parallel algorithms

652 capitalizing on the multi-core technology. The parallel tree traversal library should enable

653 parallel computation for a vast set of applications facing the challenges of increasing model

654 complexity and volumes of data in phylogenetic analysis.

# Supplementary Material

Data from the performance benchmarks and simulations for technical correctness is available on the dryad database. The POUMM package and user guide is available at `https://CRAN.R-project.org/package=POUMM`. The package also includes the source code of the SPLiTTree framework.

# Funding

# Acknowledgements

$*$

# References

Ayres, D. L., A. Darling, D. J. Zwickl, P. Beerli, M. T. Holder, P. O. Lewis, J. P. Huelsenbeck, F. Ronquist, D. L. Swofford, M. P. Cummings, A. Rambaut, and M. A. Suchard. 2012. BEAGLE: An Application Programming Interface and High-Performance Computing Library for Statistical Phylogenetics. Systematic Biology 61:170–173.

673  Bachmann, N., T. Turk, C. Kadelka, A. Marzel, M. Shilaih, J. Böni, V. Aubert,
674  T. Klimkait, G. E. Leventhal, H. F. Günthard, R. Kouyos, and Swiss HIV Cohort Study.
675  2017. Parent-offspring regression to estimate the heritability of an HIV-1 trait in a
676  realistic setup. Retrovirology 14:33.

677  Beaulieu, J. M. and B. OMeara. 2016. OUwie: Analysis of Evolutionary Rates in an OU
678  Framework .

679  Bertels, F., A. Marzel, G. Leventhal, V. Mitov, J. Fellay, H. F. Günthard, J. Böni, S. Yerly,
680  T. Klimkait, V. Aubert, M. Battegay, A. Rauch, M. Cavassini, A. Calmy, E. Bernasconi,
681  P. Schmid, A. U. Scherrer, V. Müller, S. Bonhoeffer, R. Kouyos, R. R. Regoes, and Swiss
682  HIV Cohort Study. 2017. Dissecting HIV Virulence: Heritability of Setpoint Viral Load,
683  CD4+ T Cell Decline and Per-Parasite Pathogenicity. Molecular biology and evolution .

684  Blanquart, F. o., C. Wymant, M. Cornelissen, A. Gall, M. Bakker, D. Bezemer, M. Hall,
685  M. Hillebregt, S. H. Ong, J. Albert, N. Bannert, J. Fellay, K. Fransen, A. J. Gourlay,
686  M. K. Grabowski, B. Gunsenheimer-Bartmeyer, H. F. G nthard, P. Kivel, R. Kouyos,
687  O. Laeyendecker, K. Liitsola, L. Meyer, K. Porter, M. Ristola, A. van Sighem,
688  G. Vanham, B. Berkhout, P. Kellam, P. Reiss, C. Fraser, and BEEHIVE collaboration.
689  2017. Viral genetic variation accounts for a third of variability in HIV-1 set-point viral
690  load in Europe. Plos Biology 15:e2001855.

691  Bortolussi, N., E. Durand, M. Blum, and O. Francois. 2012. apTreeshape: Analyses of
692  Phylogenetic Treeshape. R package .

693  Boskova, V., S. Bonhoeffer, and T. Stadler. 2014. Inference of Epidemiological Dynamics
694  Based on Simulated Phylogenies Using Birth-Death and Coalescent Models. PLoS
695  Computational Biology (PLOSCB) 10(4) 10:e1003913.

696  Bouckaert, R. R., J. Heled, D. Kühnert, T. G. Vaughan, C.-H. Wu, D. Xie, M. A. Suchard,

697  A. Rambaut, and A. J. Drummond. 2014. BEAST 2 - A Software Platform for Bayesian

698  Evolutionary Analysis. PLoS Computational Biology (PLOSCB) 10(4) 10:e1003537–.

699  Boyd, S. P. and L. Vandenberghe. 2004. Convex Optimization. Cambridge University

700  Press.

701  Brooks, S. P. and A. Gelman. 1998. General methods for monitoring convergence of

702  iterative simulations. Journal of Computational and Graphical Statistics 7:434–455.

703  Butler, M. A. and A. A. King. 2004. Phylogenetic comparative analysis: A modeling

704  approach for adaptive evolution. American Naturalist 164:683–695.

705  Byrd, R. H., P. Lu, J. Nocedal, and C. Y. Zhu. 1995. A limited memory algorithm for

706  bound constrained optimization. SIAM Journal on Scientific Computing 16:1190–1208.

707  Clavel, J., G. Escarguel, and G. Merceron. 2015. mvmorph: an r package for fitting

708  multivariate evolutionary models to morphometric data. Methods in Ecology and

709  Evolution 6:1311–1319.

710  Cook, S. R., A. Gelman, and D. B. Rubin. 2006. Validation of Software for Bayesian

711  Models Using Posterior Quantiles. Journal of Computational and Graphical Statistics

712  15:675–692.

713  Drummond, A. J., M. A. Suchard, D. Xie, and A. Rambaut. 2012. Bayesian phylogenetics

714  with BEAUti and the BEAST 1.7. Molecular biology and evolution 29:1969–1973.

715  Felsenstein, J. 1973. Maximum-likelihood estimation of evolutionary trees from continuous

716  characters. American Journal of Human Genetics 25:471–492.

717  Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood

718  approach. Journal of molecular evolution 17:368–376.

Felsenstein, J. 1983. Statistical Inference of Phylogenies. Journal of the Royal Statistical Society. Series A (General) 146:246.

FitzJohn, R. G. 2012. Diversitree: comparative phylogenetic analyses of diversification in R. Methods in Ecology and Evolution 3:1084–1092.

Goolsby, E. W., J. Bruggeman, and C. Ané. 2016. Rphylopars: fast multivariate phylogenetic comparative methods for missing data and within-species variation. Methods in Ecology and Evolution 8:22–27.

Grimmett, G. and D. Stirzaker. 2001. Probability and Random Processes. Oxford University Press.

Haario, H., E. Saksman, and J. Tamminen. 2001. An adaptive metropolis algorithm. Bernoulli. Official Journal of the Bernoulli Society for Mathematical Statistics and Probability 7:223–242.

Hansen, T. F. 1997. Stabilizing Selection and the Comparative Analysis of Adaptation. Evolution; international journal of organic evolution 51:1341–1351.

Hansen, T. F. and K. Bartoszek. 2012. Interpreting the evolutionary regression: the interplay between observational and biological errors in phylogenetic comparative studies. Systematic Biology 61:413–425.

Ho, L. s. T. and C. Ané. 2014. A linear-time algorithm for Gaussian and non-Gaussian trait evolution models. Systematic Biology 63:397–408.

Hodcroft, E., J. D. Hadfield, E. Fearnhill, A. Phillips, D. Dunn, S. O'Shea, D. Pillay, A. J. L. Brown, o. b. o. t. U. H. D. R. Database, and t. U. C. Study. 2014. The Contribution of Viral Genotype to Plasma Viral Set-Point in HIV Infection. PLoS pathogens 10:e1004112.

Housworth, E. A., E. P. Martins, and M. Lynch. 2004. The phylogenetic mixed model. The American Naturalist 163:84–96.

Ingram, T. and D. L. Mahler. 2013. SURFACE: detecting convergent evolution from comparative data by fitting Ornstein-Uhlenbeck models with stepwise Akaike Information Criterion. Methods in Ecology and Evolution 4:416–425.

Ives, A. R. and T. J. Garland. 2010. Phylogenetic Logistic Regression for Binary Dependent Variables. Systematic Biology 59:9–26.

Kühnert, D., T. Stadler, T. G. Vaughan, and A. J. Drummond. 2016. Phylodynamics with Migration: A Computational Framework to Quantify Population Structure from Genomic Data. Molecular biology and evolution 33:msw064–2116.

Kutsukake, N. and H. Innan. 2013. Simulation-Based Likelihood Approach for Evolutionary Models of Phenotypic Traits on Phylogeny. Evolution; international journal of organic evolution 67:355–367.

Lynch, M. 1991. Methods for the Analysis of Comparative Data in Evolutionary Biology. Evolution; international journal of organic evolution 45:1065–1080.

Lynch, M. and B. Walsh. 1998. Genetics and Analysis of Quantitative Traits. Sinauer Associates Incorporated.

Manceau, M., A. Lambert, and H. Morlon. 2016. A unifying comparative phylogenetic framework including traits coevolving across interacting lineages. Systematic Biology 66:syw115–568.

Mersmann, O. 2015. Accurate Timing Functions [R package microbenchmark version 1.4-2.1] .

Metropolis, N., A. W. Rosenbluth, M. N. Rosenbluth, A. H. Teller, and E. Teller. 1953. Equation of State Calculations by Fast Computing Machines. The Journal of Chemical Physics 21:1087–1092.

Mitov, V. and T. Stadler. 2016. The heritability of pathogen traits - definitions and estimators. bioRxiv Page 058503.

Müller, N. F., D. A. Rasmussen, and T. Stadler. 2017. The Structured Coalescent and Its Approximations. Molecular biology and evolution 34:2970–2981.

O'Meara, B. C. 2012. Evolutionary Inferences from Phylogenies: A Review of Methods. Annual Review of Ecology, Evolution, and Systematics 43:267–285.

Ornstein, L. S. and F. Zernike. 1919. The theory of the Brownian Motion and statistical mechanics. Proceedings of the Koninklijke Akademie Van Wetenschappen Te Amsterdam 21:109–114.

Pagel, M. 1994. Detecting Correlated Evolution on Phylogenies - a General-Method for the Comparative-Analysis of Discrete Characters. Proceedings of the Royal Society B-Biological Sciences 255:37–45.

Paradis, E. and J. Claude. 2002. Analysis of comparative data using generalized estimating equations. Journal of theoretical biology 218:175–185.

Paradis, E., J. Claude, and K. Strimmer. 2004. APE: Analyses of Phylogenetics and Evolution in R language. Bioinformatics 20:289–290.

Pennell, M. W., J. M. Eastman, G. J. Slater, J. W. Brown, J. C. Uyeda, R. G. FitzJohn, M. E. Alfaro, and L. J. Harmon. 2014. geiger v2.0: an expanded suite of methods for fitting macroevolutionary models to phylogenetic trees. Bioinformatics 30:2216–2218.

786  Plummer, M., N. Best, K. Cowles, and K. Vines. 2006. CODA: Convergence Diagnosis and

787     Output Analysis for MCMC. R News 6:7–11.

788  Ronquist, F. and J. P. Huelsenbeck. 2003. MrBayes 3: Bayesian phylogenetic inference

789     under mixed models. Bioinformatics 19:1572–1574.

790  Scheidegger, A. 2017. adaptMCMC. R package .

791  Slater, G. J., L. J. Harmon, and M. E. Alfaro. 2012a. Integrating Fossils With Molecular

792     Phylogenies Improves Inference Of Trait Evolution. Evolution; international journal of

793     organic evolution 66:3931–3944.

794  Slater, G. J., L. J. Harmon, D. Wegmann, P. Joyce, L. J. Revell, and M. E. Alfaro. 2012b.

795     Fitting Models of Continuous Trait Evolution to Incompletely Sampled Comparative

796     Data Using Approximate Bayesian Computation. Evolution; international journal of

797     organic evolution 66:752–762.

798  Stadler, T., D. Kühnert, S. Bonhoeffer, and A. J. Drummond. 2013. Birth-death skyline

799     plot reveals temporal changes of epidemic spread in HIV and hepatitis C virus (HCV).

800     PNAS 110:228–233.

801  Uhlenbeck, G. E. and L. S. Ornstein. 1930. On the Theory of the Brownian Motion.

802     Physical Review 36:823–841.

803  Uyeda, J. C., J. Eastman, and L. Harmon. 2015. bayou: Bayesian Fitting of

804     Ornstein-Uhlenbeck Models to Phylogenies .

805  Vihola, M. 2012. Robust adaptive Metropolis algorithm with coerced acceptance rate.

806     Statistics and Computing 22:997–1008.