



# Chapter 7

## Multiple Sequence Alignment for Large Heterogeneous Datasets Using SATé, PASTA, and UPP

Tandy Warnow and Siavash Mirarab

### Abstract

The estimation of very large multiple sequence alignments is a challenging problem that requires special techniques in order to achieve high accuracy. Here we describe two software packages—PASTA and UPP—for constructing alignments on large and ultra-large datasets. Both methods have been able to produce highly accurate alignments on 1,000,000 sequences, and trees computed on these alignments are also highly accurate. PASTA provides the best tree accuracy when the input sequences are all full-length, but UPP provides improved accuracy compared to PASTA and other methods when the input contains a large number of fragmentary sequences. Both methods are available in open source form on GitHub.

**Key words** Multiple sequence alignment, PASTA, SATé, UPP, Ensembles of Hidden Markov Models

---

### 1 Introduction

Multiple sequence alignment (MSA) is one of the more complex bioinformatics tasks, and a precursor to many downstream analyses, including protein structure and function prediction, phylogeny estimation, orthology prediction, and even genome assembly. This chapter focuses mainly on the use of multiple sequence alignment for phylogeny estimation, and in particular on the challenge of computing alignments on large, heterogeneous datasets, where standard off-the-shelf methods have low accuracy. Of particular relevance is the challenge of computing multiple sequence alignments of datasets that exhibit sequence length heterogeneity, where all standard methods have particularly poor accuracy.

In 2009, SATé (Simultaneous Alignment and Tree Estimation) was developed to enable the co-estimation of alignments and trees on large challenging datasets [1]. SATé used a combination of divide-and-conquer (where alignments are computed on subsets using standard MSA methods and then merged into an alignment on the full dataset; *see* Fig. 1a) and iteration (where each iteration computes a new alignment based on the tree from the prior

iteration, and then a new tree is computed on the new alignment) in order to obtain highly accurate alignments on large datasets. SATé-II was developed in 2012 [2] to improve on the accuracy and scalability of SATé; it used a modified decomposition strategy but otherwise had the same structure as SATé. SATé-II was able to run on much larger datasets than SATé, but was still limited to approximately 50,000 sequences. Finally, in 2014, the algorithmic design was changed again to produce PASTA [3, 4]. The objective in the design modification was to enable analyses of even larger datasets, but these changes also improved accuracy. Thus, PASTA, which mainly differs from SATé-II in its merging step (Fig. 1b), has the best accuracy and scalability of these three methods.

In 2015, we discovered that PASTA was unable to produce highly accurate alignments when the input dataset has many fragmentary sequences. To address this challenge, we developed UPP (Ultra-large alignments using Phylogeny-aware Profiles [5]), a new technique for alignment estimation that is based on a machine learning technique we developed, called an Ensemble of Profile Hidden Markov Models [6–8]. UPP uses PASTA to compute a “backbone alignment” of a subset of the input sequences (restricted to just the full-length sequences) and then adds the remaining sequences to the backbone alignment using a computed Ensemble of Profile Hidden Markov Models. UPP provides advantages over PASTA for datasets with fragmentary sequences, but PASTA has advantages over UPP when all the sequences are full-length. Like PASTA, UPP is able to compute highly accurate alignments on ultra-large datasets, including those with 1,000,000 sequences.

This chapter describes, at a very high level, how the PASTA and UPP algorithms operate, and provides some guidance on how to use these methods to obtain the best accuracy.<sup>1</sup> More information on how to run these methods can be obtained from the tutorials for PASTA and UPP available at the GitHub sites for these methods [10, 11].

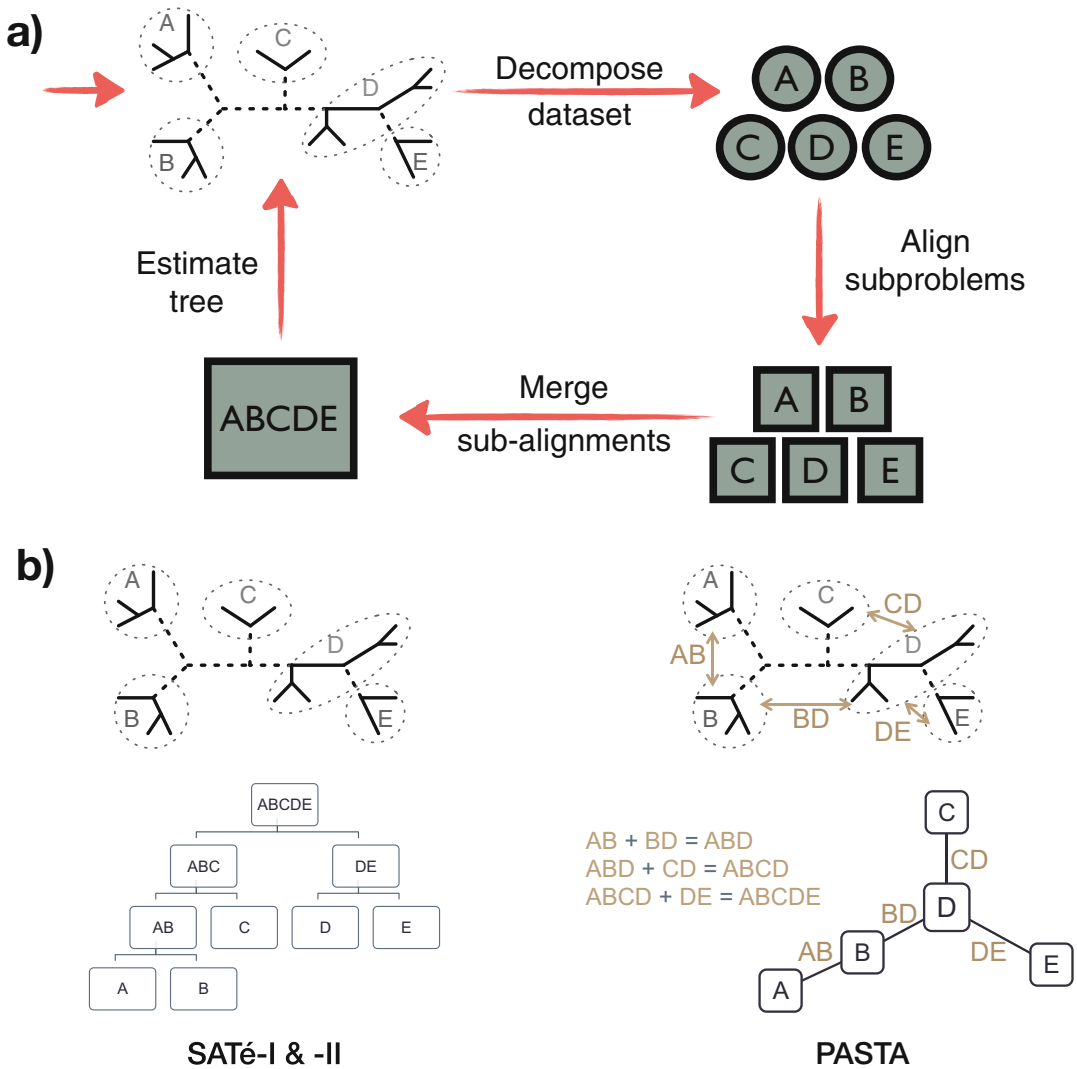
---

## 2 SATé and PASTA

SATé [1], SATé-II [2], and PASTA [4] are methods for computing multiple sequence alignments and trees from unaligned sequences (*see Note 20*). They all have the same basic algorithmic strategy (Fig. 1), and so can be considered to be members of the same basic paradigm; however, SATé-II was designed to improve on SATé (now called SATé-I) and PASTA was designed to improve on

---

<sup>1</sup> This chapter is an update of [9], a previous article for *Methods in Molecular Biology*, which focused on using SATé [1, 2] for co-estimation of alignments and trees.



**Fig. 1 (a)** The general divide-and-conquer strategy used in SATé-I, SATé-II, and PASTA. In each iteration, using the current tree, sequences are divided into smaller subsets, each subset is aligned, alignments of subsets are merged, a new tree is inferred, and a new iteration starts. **(b)** SATé and PASTA differ mainly in how they merge sub-alignments. SATé uses a hierarchical approach, where the hierarchy reflects the tree, and uses external methods like Opal or Muscle to merge alignments. PASTA, on the other hand, uses a spanning tree, computed from the phylogeny, to compute a set of pairwise alignment mergers, which are then combined using transitivity

SATé-II, with the result that PASTA dominates the other methods with respect to accuracy, running time, memory usage, and scalability to large datasets. For example, PASTA has been able to compute alignments and trees on up to 1,000,000 sequences, but SATé-I and SATé-II have not been able to analyze datasets of this

size. Furthermore, PASTA, which is the method of choice in this family of methods, has an active user community (e.g., Google group [pasta-users@googlegroups.com](mailto:pasta-users@googlegroups.com)) (*see* **Notes 1–3**).

## 2.1 *Iterative Divide-and-Conquer Strategy*

Each of these methods has the same basic structure. They begin by computing a quick alignment and tree, for example, using the fast maximum likelihood (ML) heuristic FastTree-2 [12] on a fast alignment, such as Clustal-Omega [13], and then they iterate between computing a new alignment using the current tree and computing an ML tree on the new alignment. The number of iterations (*see* **Note 10**) can be selected by the user or the user can simply run the method until some stopping criterion is met (e.g., the ML score stops improving). The final alignment/tree pair is then returned.

As noted, each iteration uses the tree from the previous iteration to compute a new alignment, and then a new ML tree is computed on the new alignment. The key to computing the new alignment is divide-and-conquer: the current tree is used to decompose the sequences into disjoint subsets, new alignments are computed on the subsets using a selected “subset aligner,” and then the subset alignments are merged together into an alignment on the full dataset (Fig. 1a).

The only difference between SATé-I and SATé-II is that SATé-II enables the user to specify how large the subsets can be, and it modified the decomposition strategy so that the subsets do not exceed the specified maximum size; this change enables SATé-II to analyze larger datasets and results in improved accuracy compared to SATé-I. The major difference between SATé-II and PASTA is how the subset alignments are merged into a single alignment on the full dataset (Fig. 1b). The change in the sub-alignment merging strategy (in addition to other smaller changes, such as using a new method to obtain initial alignments) enables PASTA to analyze larger datasets than SATé-II and also improves its accuracy. In fact, PASTA can compute alignments on up to 1,000,000 sequences, and neither SATé-I nor SATé-II can analyze datasets of this size. Thus, PASTA strictly dominates SATé-I and SATé-II in terms of accuracy, scalability, and speed.

By design, PASTA is fundamentally a method for enabling a selected MSA method to be run only on subsets of bounded size (where the bound is selected by the user). Furthermore, PASTA provides many choices for the subset aligner, including MAFFT [14], Clustal-Omega, Opal [15], Prank [16], and Muscle [17], and additional subset aligner methods for protein sequences (*see* **Note 8**).

The rest of this section is described in terms of how to use the PASTA GUI (which is similar to the SATé GUI). However, the command line version of PASTA enables other options than the GUI, and so the advanced users should not restrict themselves to the GUI (*see* **Notes 1–3**).

## 2.2 PASTA Parameters

The PASTA GUI, shown in Fig. 2, shows the choices that the user has in running PASTA.

- **Aligner:** This option specifies the method used to compute alignments on subsets; the default is MAFFT, but other alignment methods are available. *See Notes 6–8.*
- **Merger:** This option allows the user to choose the method to merge pairs of alignments during its approach for combining

The screenshot displays the PASTA GUI with the following sections and settings:

- External Tools:**
  - Aligner: MAFFT
  - Merger: MUSCLE
  - Tree Estimator: FASTTREE
  - Model: GTR+G20
- Sequences and Tree:**
  - Sequence file ...: [Empty text box]
  - Multi-Locus Data: ☐
  - Data Type: DNA
  - Initial Alignment: ☐ Use for initial tree
  - Tree file (optional) ...: [Empty text box]
- Workflow Settings:**
  - Algorithm: ☐ Two-Phase (not PASTA)
  - Post-Processing: ☐ Extra RAXML Search
- Job Settings:**
  - Job Name: pastajob
  - Output Dir.: [Empty text box]
  - CPU(s) Available: 1
  - Max. Memory (MB): 1024
- PASTA Settings:**
  - Max. Subproblem: ☒ Percentage (50) ☐ Size (200)
  - Decomposition: MinCluster
  - Time Limit (hr): 24
  - ☒ Iteration Limit: 100
  - Return: Final

At the bottom, there is a 'Start' button and a status bar showing 'PASTA 1.8.5, 2013-2017', 'Running Log (2019-07-18 09:27:11 PST)', and 'PASTA Ready!'.

**Fig. 2** PASTA graphical user interface (GUI). The GUI shows the major algorithmic choices in running PASTA; *see* Subheading 2.2. The **EXTERNAL TOOLS** determine how subsets are aligned, how these subset alignments are merged, and how trees are computed on the merged alignment in each iteration (which is determined both by the tree estimator method and the sequence evolution model). **SEQUENCES AND TREE** indicate the type of data, and also allow the user to provide an initial alignment and tree. The **PASTA SETTINGS** specify the maximum size for the subsets, the type of decomposition used in decomposing the dataset into subsets, how many iterations to perform, and whether to return the tree from the last iteration or the tree (from among all the iterations) with the best maximum likelihood score. Finally, **WORKFLOW SETTINGS** allow the user to just perform a two-phase analysis (first compute an alignment and then a tree) instead of using PASTA, or to run RAXML [18] on the final alignment returned by PASTA

subset alignments into an alignment on the full dataset; the choice is between Opal and Muscle. The merger technique in PASTA is only used on *pairs* of alignments, and the final alignment is then constructed from these merged pairs using transitivity. *See Note 9.*

- **Tree Estimator:** This option allows the user to choose between RAxML and FastTree-2, two heuristics for maximum likelihood, for computing trees in each iteration. The default is FastTree-2; *see Note 11.*
- **Model:** This option specifies the sequence evolution model, but this depends on the data type (RNA, DNA, or protein) as well as the tree estimator (RAxML or FastTree-2); *see Notes 12–14.*
- **Data Type:** This section allows the user to specify the data type (DNA, RNA, or protein). The default for the data type is DNA, and unless the user specifies otherwise, the analysis will be performed as though the data are DNA. *See Note 6.*
- **Initial alignment:** This is an optional command that allows the user to provide a pre-computed alignment to PASTA, for use in computing the first tree. *See Note 4.*
- **Tree file:** This is an optional command that allows the user to provide a pre-computed tree to PASTA, for use in computing the first decomposition and subsequent alignment. *See Note 4.*
- **Max. Subproblem:** The options here let the user specify the maximum subset size, either as a percentage of the full set of sequences or as a fixed number (i.e., “size”). *See Notes 5, 7, and 15.*
- **Decomposition:** This option specifies both which edges to remove (MinCluster, centroid edge or longest edge) in computing the decomposition of the sequences into subsets (the default is MinCluster) and how many iterations to perform (either a specific number of iterations or a maximum amount of time). The MinCluster decomposition, which minimizes the number of subsets with a bounded size [19], began with version 1.8.0, and it further improves alignment accuracy.
- **Return:** This option allows the user to decide whether to return the tree from the last iteration or the tree with the best maximum likelihood score of all the trees from all the iterations. *See Note 16.*

The most important considerations in running PASTA are (1) what alignment method to use to compute alignments on subsets, (2) how small to make the subsets, and (3) how many iterations to run. A good default for the subset aligner is MAFFT [14]. When MAFFT is used to align subsets, then limiting the subsets to 200 sequences makes it feasible to run the more computationally intensive variants of MAFFT (such as MAFFT L-INS-i

and MAFFT G-INS-i), which improves accuracy; reducing the maximum subset size will tend to reduce the running time while increasing the maximum subset size will tend to increase the running time (*see* **Note 7** for a discussion about the impact on accuracy). Other methods, such as BALi-Phy [20, 21], can also be used to align subsets, as shown in [22].

The number of iterations is also important, and previous studies have shown that accuracy improves substantially in the first few iterations, and then alignment and tree accuracy seem to stabilize. While three iterations seem to be sufficient for high accuracy in most conditions, it seems possible that more iterations could improve accuracy for challenging datasets. However, increasing the number of iterations also increases the running time. Hence, this is an issue that involves a potential tradeoff between time and accuracy; *see* **Note 4**.

The tree estimation method used in PASTA within each iteration also impacts accuracy, and the default is FastTree-2. Most users will prefer to use other methods than FastTree-2 for the final tree, and PASTA enables the user to perform a final RAXML analysis on the final tree. This is advisable whenever the dataset is not so large that RAXML is infeasible. *See* **Notes 11** and **14**.

### 2.3 PASTA Output

PASTA produces both an alignment and a tree. In addition to the final alignment and tree, PASTA outputs alignments and trees generated in each iteration as temporary files. It also outputs a config file recording all the settings used.

For large datasets (many thousands of sequences), PASTA tends to produce very long and gappy alignments because it is conservative in inferring homologies. The gappy alignments (partially a natural consequence of large datasets and partially a consequence of the algorithmic design of PASTA) seem to not hurt PASTA's ability to produce very accurate trees, but the alignments will certainly look strange to some users (*see* discussion in [23] about the preference among some users for less gappy alignments). Furthermore, whether accurate or not, these gappy alignments can also cause difficulties in some subsequent analyses. For example, phylogenetic inference can become slow given long alignments, and the inclusion of gappy sites may not result in improved phylogenetic accuracy. To speed up the tree estimation stage within each iteration, PASTA alignments are first masked to remove all sites that are at least 99.9% gapped, before trees are computed. This default setting for masking sites inside PASTA can be modified using the `--mask-gappy-sites` option. Removing gappy sites from the final alignment generated by PASTA can be done using the `run_seqtools.py` script, which is packaged with and installed together with PASTA. More aggressive filtering would further reduce the running time for phylogenetic inference, but could also have negative consequences for accuracy; *see* discussion in [24].

## 2.4 Websites for PASTA

- PASTA is available in open source form on GitHub at [10] and a protein version is available at [25]. The software is developed under the GNU Public License (GPL).
- A version of PASTA for use with BALi-Phy is available at [26].
- All questions and inquiries should be addressed to our user email group: <https://pasta-users@googlegroups.com>, with posts available at <https://groups.google.com/forum/#!forum/pasta-users>.
- A PASTA tutorial is available at <https://github.com/smirarab/pasta/blob/master/pasta-doc/pasta-tutorial.md>.

---

## 3 UPP

As we have found, PASTA produces highly accurate alignments and trees, and improves on the accuracy of other methods for ultra-large datasets with high rates of heterogeneity. However, when the input dataset has many fragmentary sequences, then PASTA does not have good accuracy. Furthermore, no standard alignment method has good accuracy when fragments are included. However, UPP [5] is an alternative approach that has good accuracy, and is the focus of this section.

### 3.1 Ensembles of Profile Hidden Markov Models (HMMs)

UPP builds on PASTA to improve its ability to align datasets with fragmentary sequences using the “ensembles of HMMs” technique, which we now describe in the context of working with multiple sequence alignments.

A profile Hidden Markov Model (HMM) [27] is a probabilistic graphical model that has vertices and directed edges, with a single vertex for the start state, a single vertex for the end state, and additional vertices corresponding to match states, insertion states, and deletion states. With the exception of the insertion states (which can have self-loops), there are no directed cycles in a profile HMM. Each directed edge  $e = v \rightarrow w$  in the profile HMM is annotated with a real number  $p_e$  where  $p_e$  is the probability of moving from  $v$  to  $w$ . Finally, the insertion states and match states emit letters (e.g., nucleotides or amino acids) from a probability distribution. Thus, when tracing a path through a profile HMM, and selecting the letters to be emitted by the visited match and insertion states, a sequence is produced.

Profile HMMs are a major part of many bioinformatics analyses, and one of the interesting uses is to add sequences into multiple sequence alignments. In what follows, we describe how profile Hidden Markov Models can be used specifically for multiple sequence alignment; *see* [28] for additional details and discussion.



To add a sequence  $s$  into a multiple sequence alignment  $A$ , a profile HMM is built for  $A$ , and then an optimal path (e.g., a maximum likelihood path) through the model is found for  $s$ . Once the path is found, it defines a way of adding  $s$  into the alignment  $A$ . Note that this addition does *not* define an alignment between  $A$  and those letters in  $s$  that are mapped to insertion states. Thus, when using HMMs to extend  $A$  to include  $s$ , some parts of  $s$  may remain *unaligned*. Besides finding the best alignment, given the sequence  $s$  and a profile HMM, the fit between the profile HMM and  $s$  can be calculated in various ways, including finding the overall probability that the profile HMM would generate  $s$ . The HMMER3 [29] suite of tools provides a particular implementation of the general profile HMM concept and includes many further optimizations, both for accuracy and speed. HMMER3 includes tools for all these analyses (i.e., building profile HMMs from alignments, scoring the fit between a profile HMM and a sequence, and finding the best path through the model for the sequence) [29, 30].

An *ensemble of profile HMMs* is a collection of profile HMMs that are built using a multiple sequence alignment  $A$ , with each profile HMM in the set based on just a subset of the sequences in the set. Thus, the match states in each of the profile HMMs in the set correspond to sites in  $A$ . Now, given a sequence  $s$ , the profile HMM in the collection that has the best fit to  $s$  can be found, the best path through the model can be computed, and thus the sequence  $s$  can be added to the alignment  $A$ . Thus, an ensemble of profile HMMs can also be used to represent the alignment  $A$  and then used to add new sequences to  $A$ . Here we will show how UPP uses an ensemble of profile HMMs to compute multiple sequence alignments, noting also that ensembles of HMMs have been used for phylogenetic placement [6], taxonomic identification of metagenomic data [7], and classification of protein sequences into families and superfamilies [8].

### 3.2 UPP's Algorithmic Protocol

In essence, UPP is a combination of PASTA (which it uses to construct a multiple sequence alignment on a subset of the input sequences) with a way of computing an ensemble of profile HMMs, which it then uses to add the remaining sequences into the PASTA alignment. Here we describe this process as operating in four steps. The first two steps can be omitted if the user wishes to provide UPP with a pre-computed backbone alignment and tree; *see Note 17*.

1. Given a set  $S$  of unaligned sequences, UPP begins by identifying those sequences to be part of “backbone alignment.” This is performed first by restricting  $S$  to just those sequences with length within 25% of the median sequence length, and then randomly selecting a set of sequences from that set. The number of sequences in the backbone and restrictions on what

sequences can be included in the backbone can be modified by using a configuration file or input options (`-B`, `-M`, `-T`, and `-L`).

2. UPP uses PASTA to compute a multiple sequence alignment  $A'$  and tree  $T$  on  $S'$ , which are then referred to as the backbone alignment and backbone tree.
3. UPP builds an ensemble of profile HMMs to represent the multiple sequence alignment  $A'$  on  $S'$ : it uses the tree  $T$  to break the set of sequences into disjoint subsets of bounded size (using the same centroid edge decomposition as in SATé-II), and then computes a profile HMM for each of the subsets (i.e., for the rows of the alignment  $A'$  defined by the sequences in the subset). The set of profile HMMs it creates is the ensemble of profile HMMs used in the next step.
4. The remaining sequences (i.e., the ones that are not in the backbone alignment) are added to  $A'$  using the ensemble of profile Hidden Markov Models computed in **Step 3**, thus producing a multiple sequence alignment  $A$  on  $S$ .

As shown in [5], UPP produces more accurate alignments than PASTA and other multiple sequence alignment methods when the input set  $S$  has many fragmentary sequences, and trees computed on the alignment are more accurate than trees computed on the other alignments in the presence of fragmentation.

### 3.3 UPP's Parameters

The most important algorithmic options in using UPP are (a) which sequences to put in the backbone subset  $S'$ , (b) which method to use to compute an alignment on  $S'$ , and (c) which algorithmic parameters to use for building the ensemble of profile Hidden Markov Models.

For which sequences to put in  $S'$ , there are two decisions that need to be made: first, which sequences are close enough to full-length to be considered, and second, how many of these sequences to use for the backbone alignment. The default UPP operates as follows: it computes the median sequence length of the input sequences and considers any sequence within 25% of this length to be “full-length.” Then, UPP selects a random subset of the “full-length” sequences to include in the backbone alignment, with the default setting for the size of this set being the minimum of  $\{1000, N\}$ , where  $N$  is the number of “full-length” sequences. Changing the number of sequences to put in the backbone set can affect accuracy and running time, and is discussed in *see Note 18*.

For how to compute the backbone alignment, the default is to use PASTA, and this is certainly appropriate when  $S'$  is large. However, when  $S'$  is small enough, then other methods can potentially provide improved accuracy compared to PASTA. For example, BALi-Phy [20, 21] and other statistical methods could be used to compute an alignment  $A'$  on  $S'$ . Once the backbone alignment is

built, a backbone tree is also needed, which can be estimated using fast ML heuristics, such as FastTree-2 (e.g., as outputted by PASTA).

There are several algorithmic options for building the ensemble of profile HMMs on  $A'$ , which we briefly discuss here. Recall that an ensemble of profile HMMs is a collection of profile HMMs, where each of the profile HMMs is constructed on a subset of the sequences in the backbone alignment. To add a sequence  $s$  into the backbone alignment,  $s$  is scored with respect to each profile HMM in the collection, and the profile HMM with the best score is selected. Thus, every sequence  $s$  that is not in the backbone alignment must be scored against every profile HMM in the collection. Although we observe that typically accuracy is increased by having a large number of profile HMMs, this also increases the running time. Thus, there is a potential tradeoff between accuracy and running time. The default in UPP produces 10 profile HMMs, which provides an improvement over a single profile HMM and (obviously) also increases the running time. See **Note 19** for additional considerations for this algorithmic setting.

Although modifications to the default settings can result in improved accuracy or speed, the default settings for UPP are sufficient to improve on PASTA if the proportion of fragmentary sequences is large enough. Detailed information on how to adjust the settings of UPP are given in its README file at <https://github.com/smirarab/sepp/blob/master/README.UPP.md>.

### 3.4 Websites for UPP

- The UPP software is available in open source form on GitHub at [11], and is part of the SEPP [6] distribution (which has code for various methods that use ensembles of profile HMMs). UPP is available as Python code.
- A tutorial on UPP is available at <https://github.com/smirarab/sepp/blob/master/tutorial/upp-tutorial.md>
- The UPP users group forum is available at <https://groups.google.com/forum/#!forum/ensemble-of-hmms>.

---

## 4 Discussion and Summary

UPP and PASTA are two methods for large-scale multiple sequence alignment that provide improved accuracy over standard methods when datasets are large and heterogeneous. UPP provides a specific advantage over PASTA when the dataset has fragmentary sequences and PASTA provides advantages when all the sequences are full-length. UPP and PASTA are available in open source form in order to encourage further development by the research community. Furthermore, each method is designed to improve scalability of

MSA methods, which are run only on subsets of the input sequence set. Therefore, as new MSA methods are developed, PASTA and UPP can be extended to use these new methods.

PASTA is described here as a method for co-estimating alignments and trees, but it is not a statistical co-estimation method in the sense that BALi-Phy [20, 21] and StatAlign [31] are. However, PASTA can run on very large datasets while truly statistical co-estimation methods are limited to fairly small datasets (perhaps 100 sequences). Furthermore, PASTA and UPP have been used with BALi-Phy to compute subset alignments [22], thus enabling BALi-Phy to scale (in some sense) to very large datasets (e.g., up to 10,000 sequences!).

Although the discussion here was largely based on using PASTA within the GUI, the command line version enables additional settings that can provide improved accuracy. This was intentional, as the GUI is the easiest way to become familiar with PASTA, and the GUI version provides the same advantages as the command line version over other methods on large datasets. However, advanced users should use the command line version, which allows the algorithmic settings to be modified in additional ways.

We set out to discuss multiple sequence alignment for the purpose of tree estimation. PASTA is specifically designed to co-estimate alignments and trees (in an iterative fashion), so that the final tree is produced by running a maximum likelihood heuristic (either RAxML or FastTree-2) on the final alignment. Some consideration, therefore, should be made for how to compute trees from these improved alignments. While we focused on maximum likelihood under standard sequence evolution models, other approaches could be used, including Bayesian estimation (e.g., MrBayes [32] and BEAST [33, 34]), distance-based estimation (e.g., FastME [35]), and parsimony analyses (e.g., TNT [36] and PAUP\* [37]). Bayesian or maximum likelihood analyses under non-standard sequence evolution models may also be necessary, especially for datasets that span large evolutionary distances where violations of the usual model assumptions (stationarity, time reversibility, and homogeneity) are likely to occur [38–40]. Divide-and-conquer phylogeny estimation, where the set of species is divided into smaller, more homogeneous subsets, and then trees on the subsets are computed and combined into a tree on the full dataset (e.g., DACTAL [41], constrained-INC [42, 43], NJMerge [44], and TreeMerge [45]), may provide an improvement in tree accuracy for those datasets that violate the standard model assumptions but are too large for methods that are based on more complex models.

Finally, although UPP was able to produce better alignments than PASTA for datasets with a high number of fragmentary sequences, the construction of trees from such datasets presents additional challenges, even given error-free alignments [46]. One

possible direction is to use phylogenetic placement, where an initial tree is built using the full-length sequences and then the fragmentary sequences are added to the tree [6], but other approaches may provide better accuracy. Thus, tree estimation on large heterogeneous datasets will need to be revisited, in order to achieve the goal of accurate inference of large phylogenies.

---

## 5 Notes

We now give some high-level advice on using PASTA and UPP. The first 16 notes are for PASTA, the next three notes are for UPP, and the final note is common to both methods. The reader will benefit from consulting the GitHub sites for these methods (and in particular the tutorials and READMEs at those sites). PASTA users should also read the Notes section in [9] for advice about using PASTA (which is built on the SATé codebase, so that much of the advice for SATé is relevant to PASTA).

1. If you have a MAC, then installing PASTA by downloading the MAC application .dmg from the GitHub site is easy, but it only allows you to use the GUI (which is not always the most up-to-date version of PASTA). If you prefer to use the command line or do not have a MAC, you will need to install PASTA using some other process. The PASTA GitHub site provides details on how to do these installations, and the PASTA users group can help with installation issues.
2. PASTA has been mainly developed and tested for Linux and MAC; as a result, Windows users will generally have more difficulty and will need to rely on virtualization (through virtual images or docker images provided on the website).
3. PASTA utilizes FASTA-formatted sequence files and Newick-formatted tree files. See the PASTA README for details about allowed characters in the input data.
4. PASTA uses iteration as well as divide-and-conquer to improve alignment accuracy compared to standard MSA methods. The main algorithmic parameters (i.e., how small to make the subsets, how to compute subset alignments, how to compute trees on the alignments and which sequence evolution models to use) impact the accuracy that can be obtained in each iteration, but also impact running time. In general, our recommendation is to use the best method you can afford to run that still allows PASTA to perform at least three iterations (and more iterations, when time permits). This will allow the alignment produced by PASTA to have very good accuracy, and a final tree can then be computed on the PASTA alignment using more computationally intensive tree estimation methods. Much of

the discussion below about how to set the algorithmic parameters reflects this point. Similarly, if desired, the final alignment/tree pair produced by PASTA can be given as input to PASTA (see Sequences and Tree in the PASTA GUI, in Fig. 2), if additional PASTA iterations using more computationally intensive approaches are desired.

5. PASTA and SATé were designed to enable improved accuracy on large datasets, but they have also been used to compute alignments on small datasets (e.g., the avian datasets in [47] with fewer than 50 taxa). The PASTA default setting automatically adjusts the subset size appropriately for small datasets. For example, on sufficiently small datasets, PASTA may set the maximum subset size to as much as 50% of the number of sequences in the input.
6. PASTA (in command line mode) does not automatically detect the data type (DNA, RNA, or proteins), and the default setting is DNA. Therefore, if your data are not DNA sequences, you should make sure to specify the type explicitly, as otherwise the behavior of PASTA can be unpredictable (and the resultant alignment and tree may have poor accuracy).
7. As mentioned above, MAFFT is the default technique for aligning subsets, and works well for both proteins and nucleotides. However, when aligning proteins, other subset alignment methods can also have good accuracy, and are enabled in [25]. As mentioned earlier, when MAFFT is used to align subsets, limiting the subsets to 200 sequences makes it feasible to run the most accurate (but also most computationally intensive) variants of MAFFT, such as MAFFT L-INS-i and MAFFT G-INS-i. Changing the subset size will change the final alignment. Our studies (published and unpublished) have revealed inconsistent trends regarding the impact of the subset size parameter. However, at this time, based on the preponderance of the evidence, we suggest using the default settings, which puts the alignment subset size at 200, when using MAFFT as the subset aligner. The interested user may wish to explore the impact of changing alignment subset size, for those datasets that are small enough to allow such exploratory data analysis.
8. For protein alignment, PASTA enables the use of additional subset aligners MAFFT-G-INS-i, MAFFT-homologs, CONTRAlign (version 1) [48], and PROBCONS [49, 50]. To use MAFFT-homologs and CONTRAlign (available only in command line), the user must take additional steps during installation, as detailed in the most up to date README file. If you wish to use MAFFT-Homologs as the subset aligner, you should use the version of PASTA available at [25].

9. PASTA allows two methods for merging pairs of alignments—Opal and Muscle. The choice between the two methods does not have a large impact on accuracy, provided that the subsets are not too small (because when the subsets are very small, then the returned alignment is largely based on the technique used to merge pairs of alignments).
10. Although the default setting for PASTA sets the number of iterations to three, additional iterations could lead to improved accuracy under some conditions. In general, using additional iterations has shown some improvement in tree and alignment accuracy, but the optimal number of iterations is an under-explored topic. We therefore recommend that the user consider enabling additional iterations, when time permits, and explore the set of alignment/tree pairs that are returned in these iterations.
11. PASTA has two methods for tree estimation that are used in each iteration. The default is FastTree-2, but RAxML is also allowed. In our experience, RAxML is much more computationally intensive than FastTree-2, making FastTree-2 a better choice on large datasets, since many iterations can be run if FastTree-2 is used instead of RAxML (see previous Note). When the number of sequences is small enough, then adding a final RAxML run (with the post-processing command in the GUI) is recommended, since RAxML generally produces better ML scores and can, in some conditions, improve the tree accuracy (although there are many conditions where the improvement in ML score does not correspond to an improvement in tree topology accuracy [51]). However, when the number of sequences is large then we do not recommend having PASTA automatically perform a RAxML analysis on the final alignment, as this can be too computationally intensive. Instead, for very large datasets, we recommend the following approach: let PASTA perform its iterations using FastTree-2, save the final PASTA alignment, and then separately compute a tree on the final PASTA alignment using the preferred software (e.g., RAxML, or potentially some other method) and selected sequence evolution model. In this way, PASTA can be used to produce a highly accurate alignment, and then the best tree accuracy (and associated numeric parameters) can be obtained using a separate tree estimation phase.
12. The set of possible sequence evolution models depends on the tree estimation method (RAxML or FastTree-2) and the type of data (nucleotides or proteins). When PASTA is used with FastTree-2, only a very limited number of models are available (described in the notes below). If the user wishes to select a model for PASTA, then they should obtain a preliminary alignment and then use external software (e.g., ProtTest [52] for



protein datasets and ModelTest [53] or PLTB [54] for nucleotide datasets). However, an alternative approach can also be used: the user can run PASTA using the default model, then use the resultant alignment with the external software to select a substitution model for a final round of tree inference or potentially another iteration of PASTA (using the new model). See discussion in [28] about selecting models for phylogeny estimation.

13. To select a nucleotide sequence evolution model within PASTA, FastTree-2 and RAxML both enable the Generalized Time Reversible (GTR, [55]) model, and each can be used with a selected model for rate variation across sites (with different models depending on what tree estimation method is selected). In addition, FastTree-2 enables the use of the Jukes–Cantor (JC, [56]) model (with two models for rate variation across sites); however, we do not recommend using the JC model unless the data seem to fit the JC model well. FastTree-2 only enables two types of rate variation across sites (CAT and G20, which is an approximation of gamma distributed rates with 20 categories), but RAxML enables rate variation models that include invariable sites. The choice of rate variation model can also impact accuracy, but the more complex models are also more computationally intensive. However, our studies suggest that using simple sequence evolution models within the iterative process may not reduce the alignment accuracy substantially, and a new tree can be estimated on the final alignment using more complex models.
14. For protein alignment, the two tree estimation methods, RAxML and FastTree-2, offer very different sequence evolution models. Specifically, FastTree-2 only offers two protein substitution models (JTT and WAG), each with two site variation models, and RAxML offers 11 protein substitution models, each with four site variation models. Thus, RAxML allows a larger set of protein sequence evolution models than FastTree-2, making RAxML a better method for computing trees than FastTree-2 for proteins. However, here too the benefit from using RAxML within the iterative process may be offset by the extra time used to compute trees with RAxML. Hence, we would suggest instead that FastTree-2 be used as the tree estimation method within the iterative procedure, even for protein sequences. Then, after the PASTA alignment is computed, the user can compute a new tree on the alignment using RAxML or some other software, under the best fitting model.
15. If you wish to use BALi-Phy to align subsets within PASTA, the maximum subset size and the running time for each subset alignment need to be set so that BALi-Phy is able to converge on each subset. This is discussed in [57], and the



software for PASTA using BALi-Phy is available at [26]. However, *see* [58] for a study comparing BALi-Phy and other alignment methods on protein benchmark datasets, which showed differences between performance on simulated and biological datasets.

16. The choice of which alignment/tree pair to return (i.e., whether to return the pair produced in the final iteration or the pair that has the best maximum likelihood score) is an interesting one. In general, we expect little difference in accuracy between the two options, and so the choice may not make much difference in practice. In addition, there is no theoretical basis on which to select the pair that has the best maximum likelihood score [2], since the alignment is allowed to change. For these reasons, and also because the ML score is calculated within PASTA on masked versions of the computed alignments, the default in PASTA is the final alignment/tree pair.
17. The user can provide UPP with a pre-computed backbone alignment and tree (referred to in the UPP tutorial as a “custom seed alignment and tree”); this is a natural approach when using alignments and trees obtained from external sources (such as PFAM [59]) or when alignments and trees have been estimated using additional information (such as secondary or tertiary structure) or by specialized methods not available within UPP.
18. UPP uses PASTA to compute its backbone alignment and tree, but the selection of which sequences are put into the backbone set can be controlled by the user. In the default mode, UPP operates as follows: it computes the median sequence length of the input sequences and considers any sequence within 25% of this length to be “full-length”; the user can modify this approach as needed using options `-M` and `-T`. Once that set of full-length sequences is determined, the user can specify how many of the sequences to include in the backbone alignment using the `-B` option. The default is to take the minimum of  $\{1000, N\}$ , where  $N$  is the number of “full-length” sequences. However, another option is to include all of the full-length sequences (even when this is more than 1000); in our experience, this improves accuracy but may also increase the running time. Furthermore, reducing the number of sequences, even to as low as just 100 (the UPP-fast version), produces a reduction in accuracy (but sometimes only a small reduction, which depends on the heterogeneity in the input dataset) and a dramatic reduction in running time. Hence, there is a potential tradeoff between accuracy and running time that needs to be considered in building the ensemble.

19. The default mode for UPP is to create an ensemble of HMMs that has ten (10) profile HMMs. However, changes to this number can be considered with a potential for improved accuracy. In particular, when the input set is highly heterogeneous (as represented by low average sequence similarity), then using a larger number of profile HMMs can improve accuracy; however, the benefit in increasing the number of profile HMMs is reduced when the dataset has high average sequence similarity. Furthermore, increasing the number of profile HMMs automatically increases the running time (as it scales linearly with this number).
20. Errors in the input unaligned sequence data have the potential to reduce the accuracy of the alignment. One way to detect such errors is to use automated methods such as TreeShrink [60]. TreeShrink looks for extremely long branches in the phylogeny to detect potential errors in the data. Thus, TreeShrink can be combined with PASTA in a natural way: remove sequences on long branches from the PASTA alignment and tree (implemented in the `--treeshrink-filter` option), recompute the PASTA alignment, and add back those potentially problematic sequences using UPP. In addition, UPP allows for sequences that are on very long branches to be removed from the backbone set (see `-1`).

---

## Acknowledgements

This paper was supported by NSF grant ABI-1458652 to TW and NSF grant IIS-1845967 to SM.

## References

1. Liu K, Raghavan S, Nelesen S, Linder CR, Warnow T (2009) Rapid and accurate large-scale coestimation of sequence alignments and phylogenetic trees. *Science* 324 (5934):1561–1564
2. Liu K, Warnow T, Holder MT, Nelesen SM, Yu J, Stamatakis AP, Linder CR (2012) SATé-II: very fast and accurate simultaneous estimation of multiple sequence alignments and phylogenetic trees. *Syst Biol* 61(1):90–106
3. Mirarab S, Nguyen N, Warnow T (2014) PASTA: ultra-large multiple sequence alignment. In: International conference on research in computational molecular biology. Springer, Berlin, pp 177–191
4. Mirarab S, Nguyen N, Wang L-S, Guo S, Kim J, Warnow T (2015) PASTA: ultra-large multiple sequence alignment of nucleotide and amino acid sequences. *J Comput Biol* 22:377–386
5. Nguyen N, Mirarab S, Kumar K, Warnow T (2015) Ultra-large alignments using phylogeny aware profiles. *Genome Biol* 16:124. A preliminary version appeared in the Proceedings RECOMB 2015
6. Mirarab S, Nguyen N, Warnow T (2012) SEPP: SATé-enabled phylogenetic placement. In: Pacific symposium on biocomputing, pp 247–58
7. Nguyen N, Mirarab S, Liu B, Pop M, Warnow T (2014) TIPP: taxonomic identification and phylogenetic profiling *Bioinformatics* 30 (24):3548–3555
8. Nguyen N, Nute M, Mirarab S, Warnow T (2016) HIPPI: highly accurate protein family classification with ensembles of hidden Markov

- models. *BMC Bioinformatics* 17(Suppl 10):765
9. Liu K, Warnow T (2014) Large-scale multiple sequence alignment and tree estimation using SATé. In: *Multiple sequence alignment methods*. Springer, Berlin, pp 219–244
10. Mirarab S (2019) Github site for PASTA software. <https://github.com/smirarab/pasta>. Accessed 13 July 2019
11. Mirarab S (2019) Github site for Ensemble of HMM methods (SEPP, TIPP, UPP) software. <https://github.com/smirarab/sepp>. Accessed 13 July 2019
12. Price MN, Dehal PS, Arkin AP (2010) FastTree 2 – approximately maximum-likelihood trees for large alignments. *PLoS ONE* 5(3), e9490. <https://doi.org/10.1371/journal.pone.0009490>
13. Sievers F, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert M, Soding J, Thompson JD, Higgins DG (2011) Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. *Mol Syst Biol* 7:539
14. Katoh K, Toh H (2008) Recent developments in the MAFFT multiple sequence alignment program. *Brief Bioinf* 9(4):286–298
15. Wheeler T, Kecicioglu J (2007) Multiple alignment by aligning alignments. In: *Proceedings of the 15th ISCB conference on intelligent systems for molecular biology*, pp 559–568
16. Löytynoja A, Goldman N (2005) An algorithm for progressive multiple alignment of sequences with insertions. *Proc Nat Acad Sci* 102:10557–10562
17. Edgar RC (2004) MUSCLE: a multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics* 5(113):113
18. Stamatakis A (2006) RAXML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands of taxa and mixed models *Bioinformatics* 22:2688–2690.
19. Balaban M, Moshiri N, Mai U, Mirarab S (2019) TreeCluster: clustering biological sequences using phylogenetic trees. *bioRxiv*, <https://doi.org/10.1101/591388>
20. Suchard MA, Redelings BD (2006) BALi-Phy: simultaneous Bayesian inference of alignment and phylogeny. *Bioinformatics* 22:2047–2048
21. Redelings BD, Suchard MA (2007) Incorporating indel information into phylogeny estimation for rapidly emerging pathogens. *BMC Evol Biol* 7:40
22. Nute M, Warnow T (2016) Scaling statistical multiple sequence alignment to large datasets. *BMC Genomics* 17(10):764
23. Löytynoja A, Goldman N (2008) Phylogeny-aware gap placement prevents errors in sequence alignment and evolutionary analysis. *Science* 320(5883):1632–1635
24. Tan G, Muffato M, Ledergerber C, Herrero J, Goldman N, Gil M, Dessimoz C (2015) Current methods for automated filtering of multiple sequence alignments frequently worsen single-gene phylogenetic inference. *Syst Biol* 64(5):778–791
25. Collins K PASTA for proteins github site. <https://github.com/kodicollins/pasta-databases>
26. Nute M (2019) Github site for PASTA+Bali-Phy. <https://github.com/mgnute/pasta>. Accessed 18 July 2019
27. Durbin R, Eddy S, Krogh A, Mitchison G (1998) *Biological sequence analysis*. Cambridge University Press, Cambridge
28. Warnow T (2018) *Computational phylogenetics: an introduction to designing methods for phylogeny estimation*. Cambridge University Press, Cambridge
29. Eddy SR (2009) A new generation of homology search tools based on probabilistic inference. *Genome Inform* 23:205–211
30. Finn RD, Clements J, Eddy SR (2011) HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res* 39: W29–W37
31. Novák Á, Miklós I, Lyngsoe R, Hein J (2008) StatAlign: an extendable software package for joint Bayesian estimation of alignments and evolutionary trees. *Bioinformatics* 24:2403–2404
32. Huelsenbeck J, Ronquist R (2001) MrBayes: Bayesian inference of phylogeny. *Bioinformatics* 17:754–755
33. Drummond A, Rambaut A (2007) BEAST: Bayesian evolutionary analysis by sampling trees. *BMC Evol Biol* 7:214
34. Bouckaert R, Heled J, Kühnert D, Vaughan T, Wu C-H, Xie D, Suchard MA, Rambaut A, Drummond AJ (2014) BEAST 2: a software platform for Bayesian evolutionary analysis. *PLoS Comput Biol* 10(4):e1003537
35. Lefort V, Desper R, Gascuel O (2015) FastME 2.0: a comprehensive, accurate, and fast distance-based phylogeny inference program. *Mol Biol Evol* 32(10):2798–2800
36. Goloboff P, Farris J, Nixon K (2008) TNT, a free program for phylogenetic analysis. *Cladistics* 24:1–13
37. Swofford DL (1996) *PAUP\*: Phylogenetic analysis using parsimony (and other methods)*, Version 4.0. Sinauer Associates, Sunderland

38. Naser-Khdour S, Minh BQ, Zhang W, Stone E, Lanfear R (2019) The prevalence and impact of model violations in phylogenetics. *BioRxiv*. <https://doi.org/10.1101/460121>
39. Crotty SM, Minh BQ, Bean NG, Holland BR, Tuke J, Jermin LS, Haeseler Av (2019) GHOST: recovering historical signal from heterotachously-evolved sequence alignments. *bioRxiv*, <https://doi.org/10.1101/174789>
40. Jermin LS, Catullo RA, Holland BR (2018) A new phylogenetic protocol: dealing with model misspecification and confirmation bias in molecular phylogenetics. *bioRxiv*, <https://doi.org/10.1101/400648>
41. Nelesen S, Liu K, Wang L-S, Linder CR, Warnow T (2012) DACTAL: divide-and-conquer trees (almost) without alignments. *Bioinformatics* 28:i274–i282
42. Zhang Q, Rao S, Warnow T (2019) Constrained incremental tree building: new absolute fast converging phylogeny estimation methods with improved scalability and accuracy. *Algorithms Mol Biol* 14(1):2
43. Le T, Sy A, Molloy EK, Zhang QR, Rao S, Warnow T (2019) Using INC within divide-and-conquer phylogeny estimation. In: *International conference on algorithms for computational biology*. Springer, Berlin, pp 167–178
44. Molloy EK, Warnow T (2018) NJMerge: a generic technique for scaling phylogeny estimation methods and its application to species trees. In: *RECOMB International conference on comparative genomics*. Springer, Berlin, pp 260–276
45. Molloy EK, Warnow T (2019) TreeMerge: a new method for improving the scalability of species tree estimation methods. *Bioinformatics* 35(14):i417–i426
46. Sayyari E, Whitfield JB, Mirarab S (2017) Fragmentary gene sequences negatively impact gene tree and species tree Reconstruction. *Mol. Biol. Evol.* 34(12):3279–3291
47. Jarvis E, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho S, Faircloth BC, Nabholz B, Howard JT, Suh A, Weber CC, daFonseca RR, Li J, Zhang F, Li H, Zhou L, Narula N, Liu L, Ganapathy G, Boussau B, Bayzid MS, Zavidovych V, Subramanian S, Gabaldón T, Capella-Gutiérrez S, Huerta-Cepas J, Rekepalli B, Munch K, Schierup M, Lindow B, Warren WC, Ray D, Green RE, Bruford MW, Zhan X, Dixon A, Li S, Li N, Huang Y, Derryberry EP, Bertelsen ME, Sheldon FH, Brumfield RT, Mello CV, Lovell PV, Wirthlin M, Schneider MPC, Prosdocimi F, Samaniego JA, Velazquez AMV, Alfaro-Núñez A, Campos PF, Petersen B, Sicheritz-Ponten T, Pas A, Bailey T, Scofield P, Bunce M, Lambert DM, Zhou Q, Perelman P, Driskell AC, Shapiro B, Xiong Z, Zeng Y, Liu S, Li Z, Liu B, Wu K, Xiao J, Yinqi X, Zheng Q, Zhang Y, Yang H, Wang J, Smeds L, Rheindt FE, Braun M, Fjeldsa J, Orlando L, Barker FK, Jonsson KA, Johnson W, Koepfli K-P, O'Brien S, Haussler D, Ryder OA, Rahbek C, Willerslev E, Graves GR, Glenn TC, McCormack J, Burt D, Ellegren H, Alstrom P, Edwards SV, Stamatakis A, Mindell DP, Cracraft J, Braun EL, Warnow T, Jun W, Gilbert MTP, Zhang G (2014) Whole-genome analyses resolve early branches in the tree of life of modern birds. *Science* 346(6215):1320–1331
48. Do CB, Gross SS, Batzoglou S (2006) CONTRAlign: discriminative training for protein sequence alignment. In: *Proceedings of the tenth annual international conference on computational molecular biology (RECOMB 2006)*. Springer, Berlin, pp 160–174
49. Do CB, Mahabhashyam MS, Brudno M, Batzoglou S (2005) ProbCons: probabilistic consistency-based multiple sequence alignment *Genome Res* 15(2):330–340
50. Do CB, Mahabhashyam MS, Brudno M, Batzoglou S (2006) ProbCons: probabilistic consistency-based multiple sequence alignment of amino acid sequences. Software available at <http://probcons.stanford.edu/download.html>
51. Liu K, Linder C, Warnow T (2012) RAXML and FastTree: comparing two methods for large-scale maximum likelihood phylogeny estimation. *PLoS ONE* 6(11):e27731
52. Abascal F, Zardoya R, Posada D (2005) ProtTest: selection of best-fit models of protein evolution. *Bioinformatics* 21(9):2104–2105
53. Posada D, Crandall K (1998) Modeltest: testing the model of DNA substitution. *Bioinformatics* 14(9):817–818
54. Hoff M, Orf S, Riehm B, Darriba D, Stamatakis A (2016) Does the choice of nucleotide substitution models matter topologically? *BMC Bioinformatics* 17:143
55. Tavaré S (1986) Some probabilistic and statistical problems in the analysis of DNA sequences. In: *Lectures on mathematics in the life sciences*, vol 17. American Mathematical Society, Providence, pp 57–86
56. Jukes TH, Cantor CR (1969) Evolution of protein molecules. In: Munro HN (ed) *Mammalian protein metabolism*. Academic, New York, pp 21–132
57. Nute M, Warnow T (2016) Scaling statistical multiple sequence alignment to large datasets. *BMC Genomics* 17:764(2016) Special issue for RECOMB-CG 2016. <https://doi.org/10.1186/s12864-016-3101-8>

58. Nute M, Saleh E, Warnow T (2018) Evaluating statistical multiple sequence alignment in comparison to other alignment methods on protein data sets. *Syst Biol* 68(3):396–411
59. Bateman A, Birney E, Cerruti L, Durbin R, Etwiller L, Eddy SR, Griffiths-Jones S, Howe KL, Marshall M, Sonnhammer EL (2002) The Pfam protein families database. *Nucleic Acids Res.* 30:276–280
60. Mai U, Mirarab S (2018) TreeShrink: fast and accurate detection of outlier long branches in collections of phylogenetic trees. *BMC Genomics* 19(S5):272