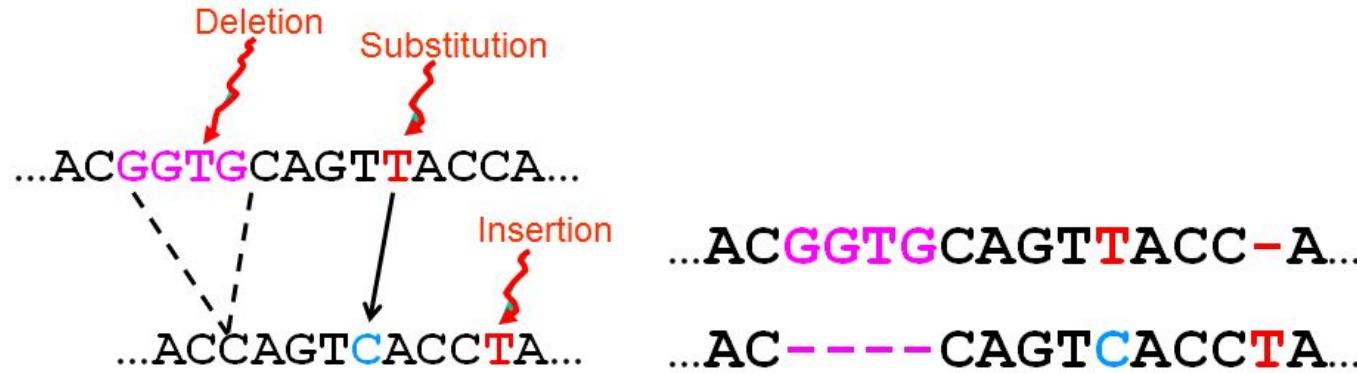


# Multiple Sequence Alignment

Lisa Pokorny & Marina Marcet-Houben

# Multiple Sequence Alignment (MSA)



## The true multiple alignment

- Reflects historical substitution, insertion, and deletion events
- Defined using transitive closure of pairwise alignments computed on edges of the true tree

# Multiple Seq Alignment (MSA)

Standard two-phase approach: 1<sup>st</sup> ALIGNMENT (positional homology)

S1 = AGGCTATCACCTGACCTCCA  
S2 = TAGCTATCACGACCGC  
S3 = TAGCTGACCGC  
S4 = TCACGACCGACA

S1 = -AGGCTATCACCTGACCTCCA  
S2 = TAG-CTATCAC--GACCGC--  
S3 = TAG-CT-----GACCGC--  
S4 = -----TCAC--GACCGACA

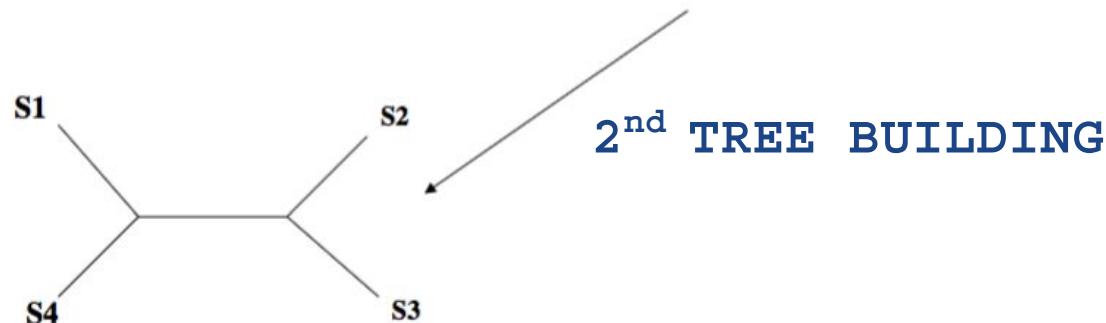
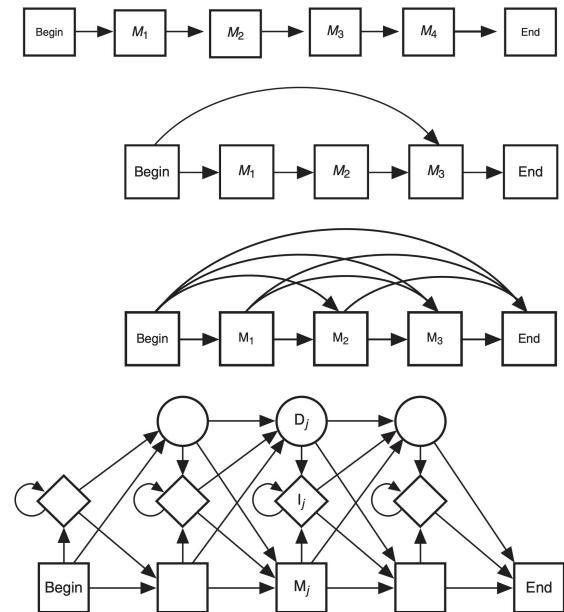


Fig. 9.1. Warnow. 2017. *Computational Phylogenetics. An Introduction to Designing Methods for Phylogeny Estimation*. CUP.

# Optimization Problems & MSA Methods (MSAMs)

- Sum-of-Pairs Alignment (SOP)
- Tree Alignment (TL) and Generalized TL
- Sequence Profiles
- Profile Hidden Markov Models (HMM)
- Reference-based Alignments
- Template-based Methods
- Seed Alignment Methods
- Weighted-Homology Pair Methods
- Progressive Methods
- Divide-and-Conquer Methods
- Co-estimation of Alignments and Trees
- Structure Informed Methods, etc.



# Phylogenetic Tree Estimation w/o Alignment?

*Syst. Biol.* 66(2):218–231, 2017  
© The Author(s) 2016. Published by Oxford University Press, on behalf of the Society of Systematic Biologists. All rights reserved.  
For Permissions, please email: journals.permissions@oup.com  
DOI:10.1093/sysbio/syw074  
Advance Access publication September 14, 2016

## Phylogenetic Tree Estimation With and Without Alignment: New Distance Methods and Benchmarking

MARCIN BOGUSZ AND SIMON WHELAN\*

Department of Evolutionary Biology, Evolutionary Biology Centre, Uppsala University, Norbyvägen 18D, 752 36 Uppsala, Sweden  
\*Correspondence to be sent to: Department of Evolutionary Biology, Evolutionary Biology Centre, Uppsala University, Norbyvägen 18D, 752 36 Uppsala, Sweden; E-mail: simon.whelan@ebc.uu.se.

Received 3 February 2016; reviews returned 12 June 2016; accepted 23 August 2016

Associate Editor: David Bryant

**Abstract.**—Phylogenetic tree inference is a critical component of many systematic and evolutionary studies. The majority of these studies are based on the two-step process of multiple sequence alignment followed by tree inference, despite persistent evidence that the alignment step can lead to biased results. Here we present a two-part study that first presents PaHMM-Tree, a novel neighbor joining-based method that estimates pairwise distances without assuming a single alignment. We then use simulations to benchmark its performance against a wide-range of other phylogenetic tree inference methods, including the first comparison of alignment-free distance-based methods against more conventional tree estimation methods. Our new method for calculating pairwise distances based on statistical alignment provides distance estimates that are as accurate as those obtained using standard methods based on the true alignment. Pairwise distance estimates based on the two-step process tend to be substantially less accurate. This improved performance carries through to tree inference, where PaHMM-Tree provides more accurate tree estimates than all of the pairwise distance methods assessed. For close to moderately divergent sequence data we find that the two-step methods using statistical inference, where information from all sequences is included in the estimation procedure, tend to perform better than PaHMM-Tree, particularly full statistical alignment, which simultaneously estimates both the tree and the alignment. For deep divergences we find the alignment step becomes so prone to error that our distance-based PaHMM-Tree outperforms all other methods of tree inference. Finally, we find that the accuracy of alignment-free methods tends to decline faster than standard two-step methods in the presence of alignment uncertainty, and identify no conditions where alignment-free methods are equal to or more accurate than standard phylogenetic methods even in the presence of substantial alignment error. [Alignment-free; distance-based phylogenetics; pair Hidden Markov Models; phylogenetic inference; statistical alignment.]

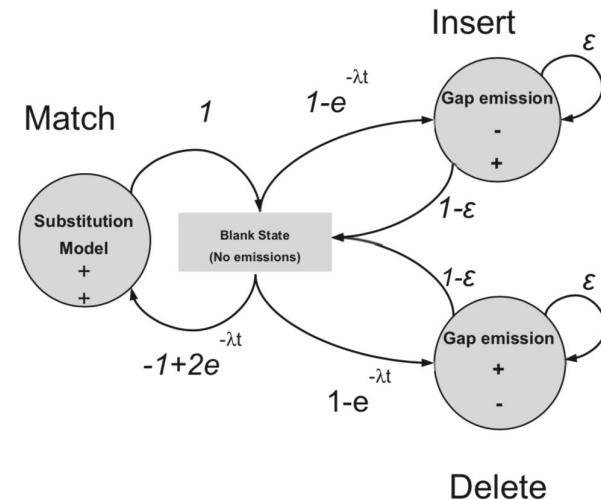
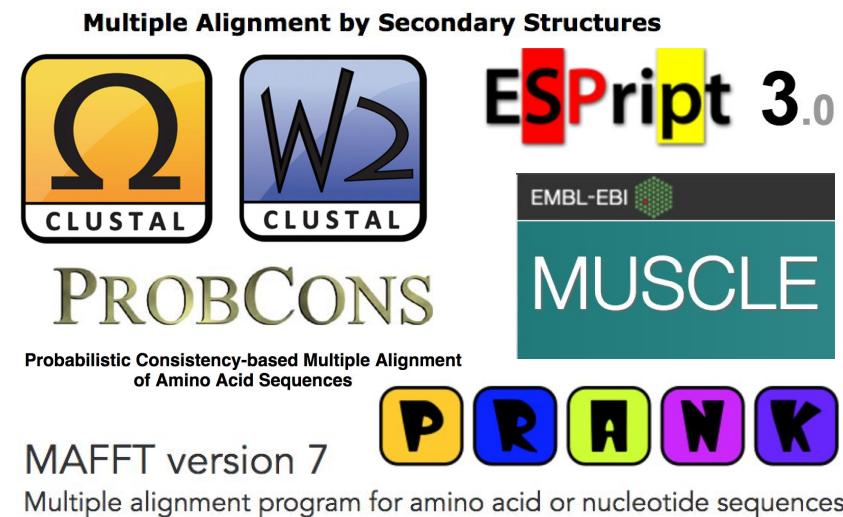
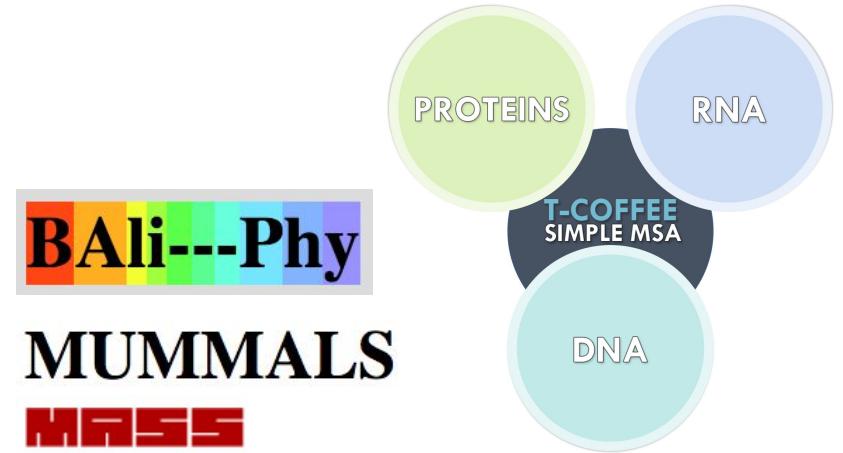


Fig. 1. Bogusz & Whelan. 2017. *Syst. Biol.* 66(2):218–231

# Comparing MSAMs

Tool	Options	Algorithm	Alphabet
ClustalW	Defaults	Progressive	Amino Acid
Muscle	Defaults	Progressive (iterative)	Amino Acid
MAFFT	Defaults	Progressive (iterative)	Amino Acid
ProbCons	Defaults	Consistency	Amino Acid
ProbAlign	Defaults	Consistency	Amino Acid
Mummals	Defaults	Consistency/Structure	Amino Acid
Dialign-TX	Defaults	Greedy/Progressive	Amino Acid
Prank (AA)	+F (AA)	"Phylogenetically-aware"	Amino Acid
Prank	+F -codon	"Phylogenetically-aware"	Codon
BAli-Phy	Model M0	Statistical Alignment	Codon
BAli-Phy samples	Model M0	Statistical Alignment	Codon
BAli-Phy integrated	Model M0	Statistical Alignment	Codon



## Mean dist btw MSAMs

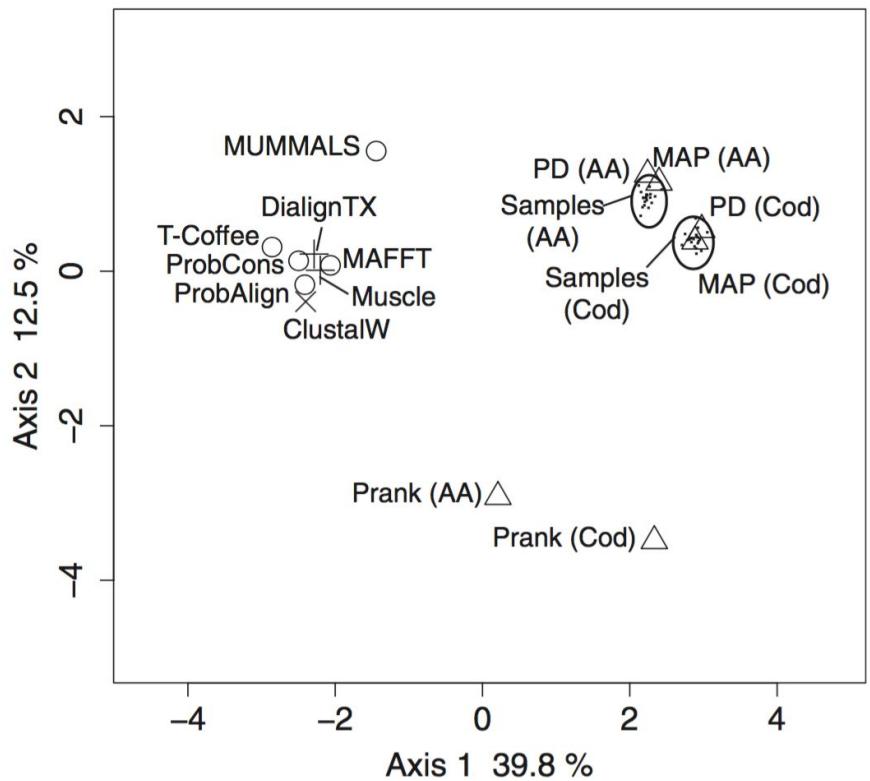


Fig. 1. Blackburne & Whelan. 2013. *Mol. Biol. Evol.* 30(3):642–653.

## Dist btw trees w $\neq$ MSAMs

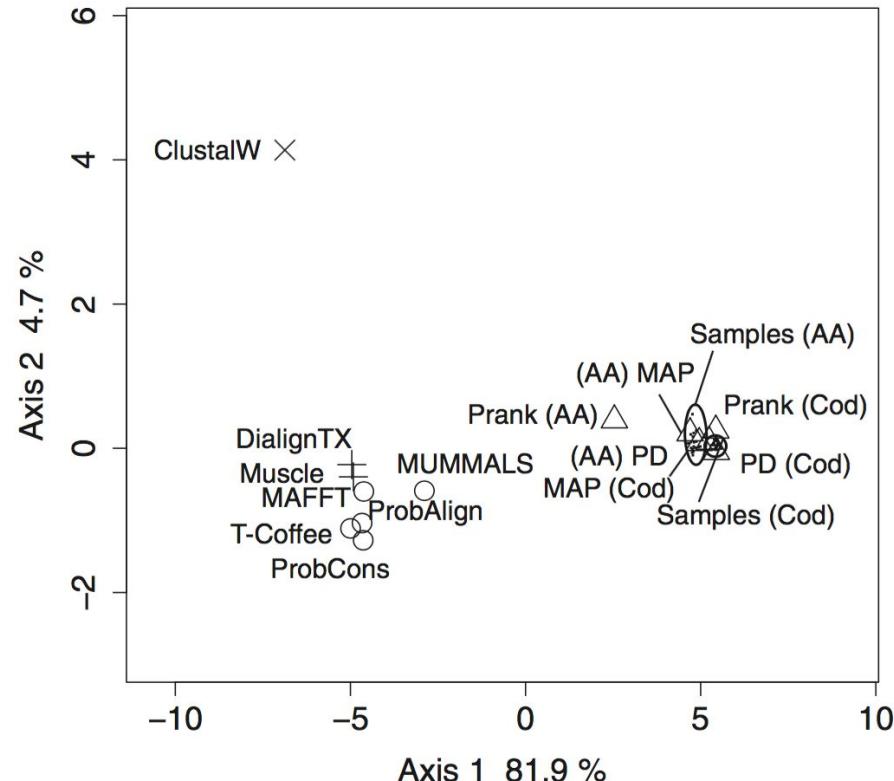
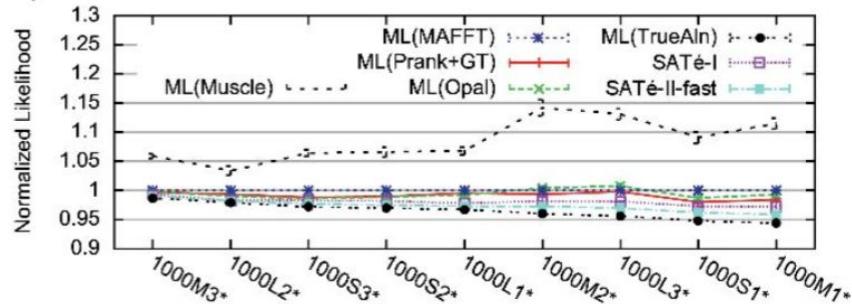


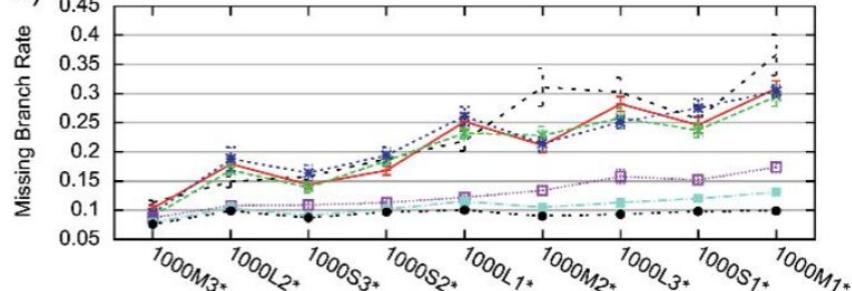
Fig. 2. Blackburne & Whelan. 2013. *Mol. Biol. Evol.* 30(3):642–653.

# Even more MSAMs comparisons

a)

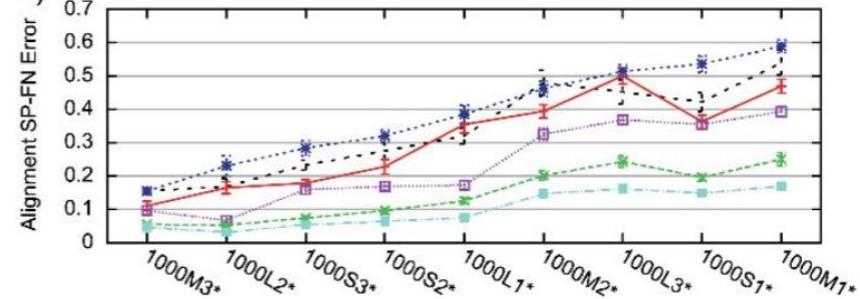


b)



c)

Alignment SP-FN Error



d)

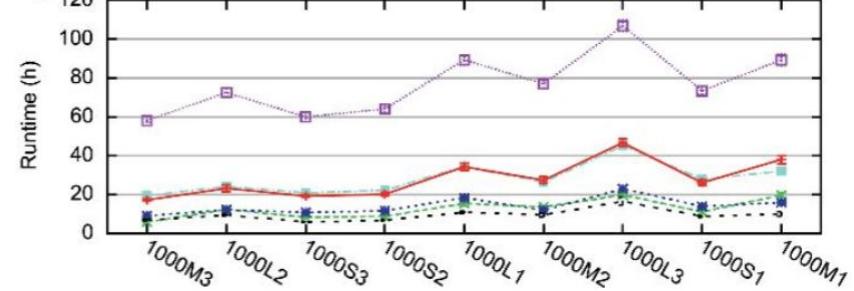
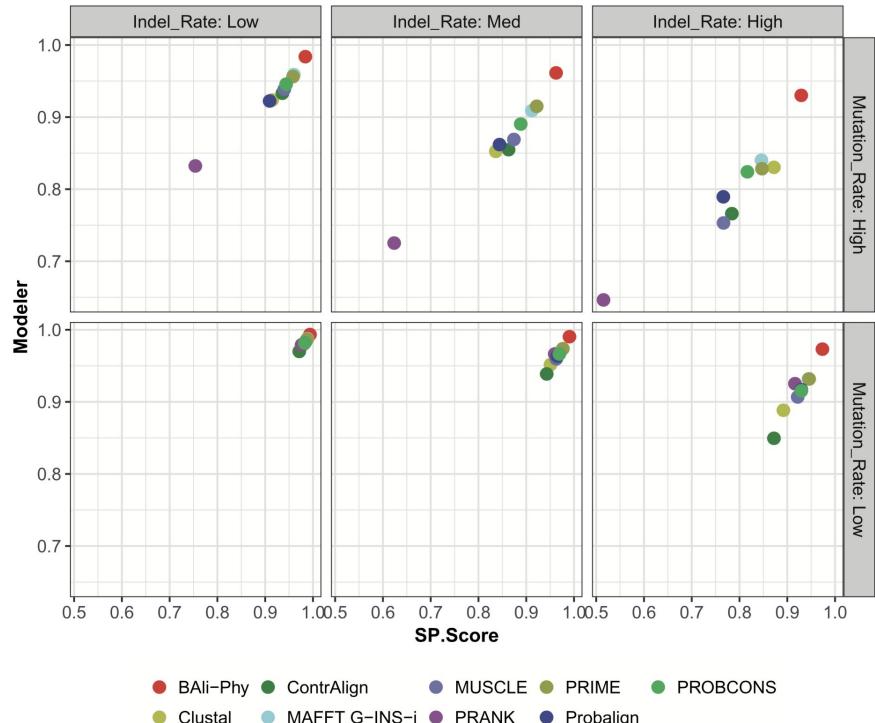


Fig. 4. Liu et al. 2013. *Syst. Biol.* 61(1):90–106.

# Modeler Precision vs. Recall



# Expansion Ratios

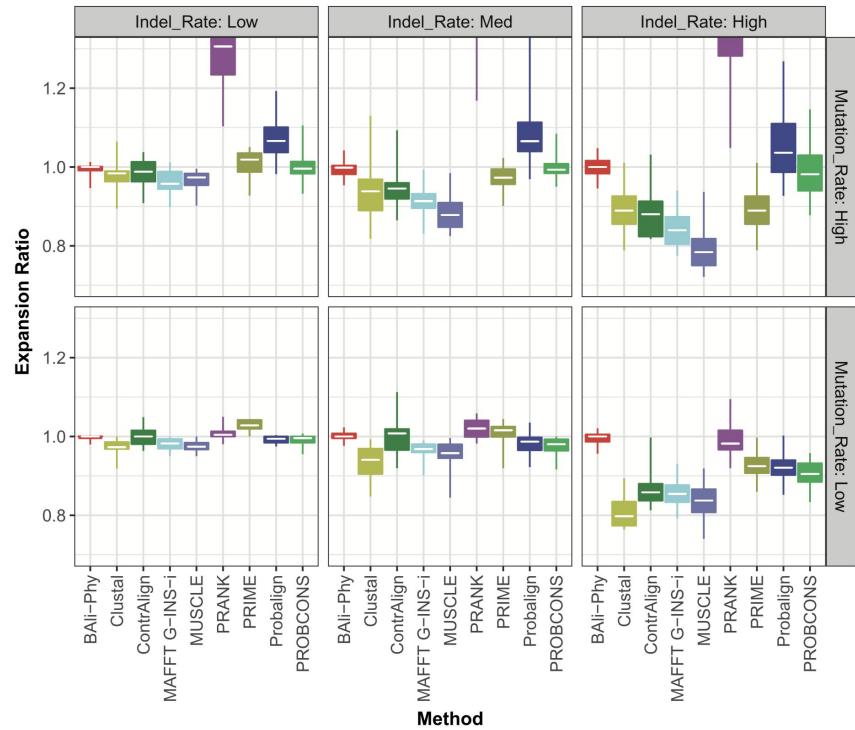
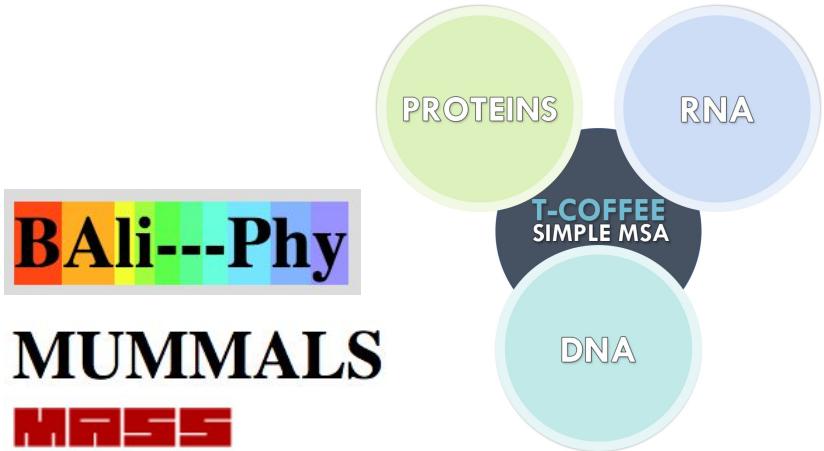


Fig. 6. Nute et al. 2018. *Syst. Biol.* <https://doi.org/10.1093/sysbio/syy068>

Fig. 7. Nute et al. 2018. *Syst. Biol.* <https://doi.org/10.1093/sysbio/syy068>



Multiple Alignment by Secondary Structures



CLUSTAL



CLUSTAL



EMBL-EBI



PROBCONS

Probabilistic Consistency-based Multiple Alignment  
of Amino Acid Sequences

MAFFT version 7

Multiple alignment program for amino acid or nucleotide sequences



# Know Your Limits

**Data type** → DNA vs. RNA, coding vs. non-coding nucleotides (wobble bp), AAs, proteins, etc.

**Data properties** → substitution ( $\neq$  mutation) rate strength ( $\uparrow$  vs.  $\downarrow$ ), indel size and rate (% gap & gap length), pairwise sequence identity (PID), etc.

**Data matrix properties** → # of tips, # of sequences, (alignment length  $\propto$ ) data matrix weight, e.g., light (K, M) vs. heavy (G, T), etc.

**CPU time and RAM memory** → computing resources available

# Divide and Conquer Method: PASTA

**PASTA** estimates **alignments and ML trees** from unaligned sequences using an **iterative approach**. In each iteration, it first estimates a multiple sequence alignment using the current tree as a guide and then estimates an ML tree on (a masked version of) the alignment. By default, PASTA performs 3 iterations, but a host of options enable changing that behavior. In **each iteration**, a **divide-and-conquer** strategy is used for estimating the alignment. The **set of sequences** is divided into smaller subsets, **each** of which is aligned using an external alignment tool (default is **MAFFT**). These **subset alignments** are then **pairwise merged** (by default using **Opal**) and finally the **pairwise merged alignments** are **merged** into a **final alignment using a transitivity merge** technique. The division of the dataset into smaller subsets and selecting which alignments should be pairwise merged is guided by the tree from the previous iteration. The first step therefore needs an initial tree.

Acknowledgment: The current **PASTA** code is heavily based on the **SATé** code developed by Mark Holder's group at KU.

# PASTA

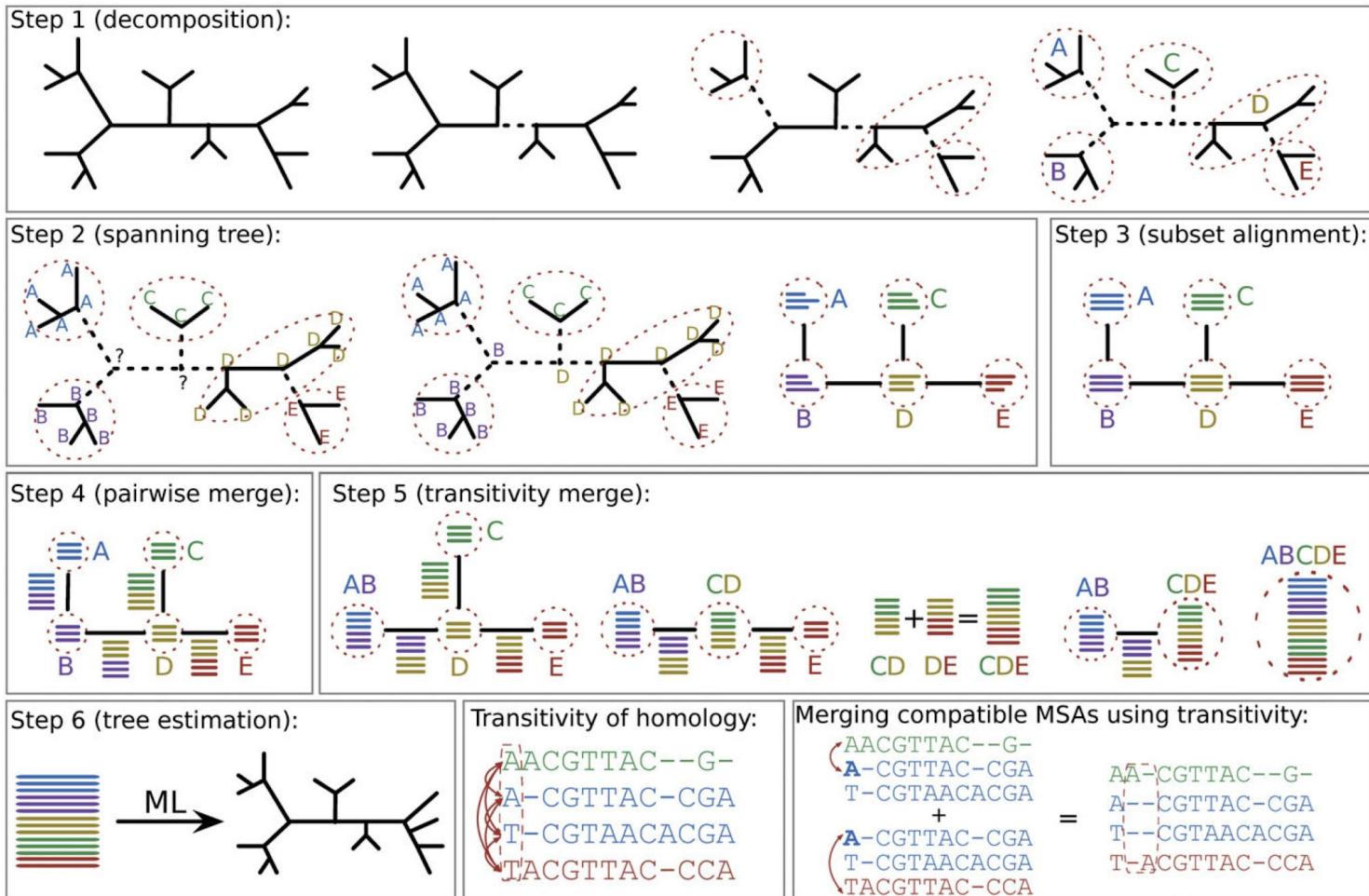


Fig. 1. Mirarab et al. 2015. *J. Comp. Biol.* 22(5):377–386.

# PASTA for nucleotides and AA

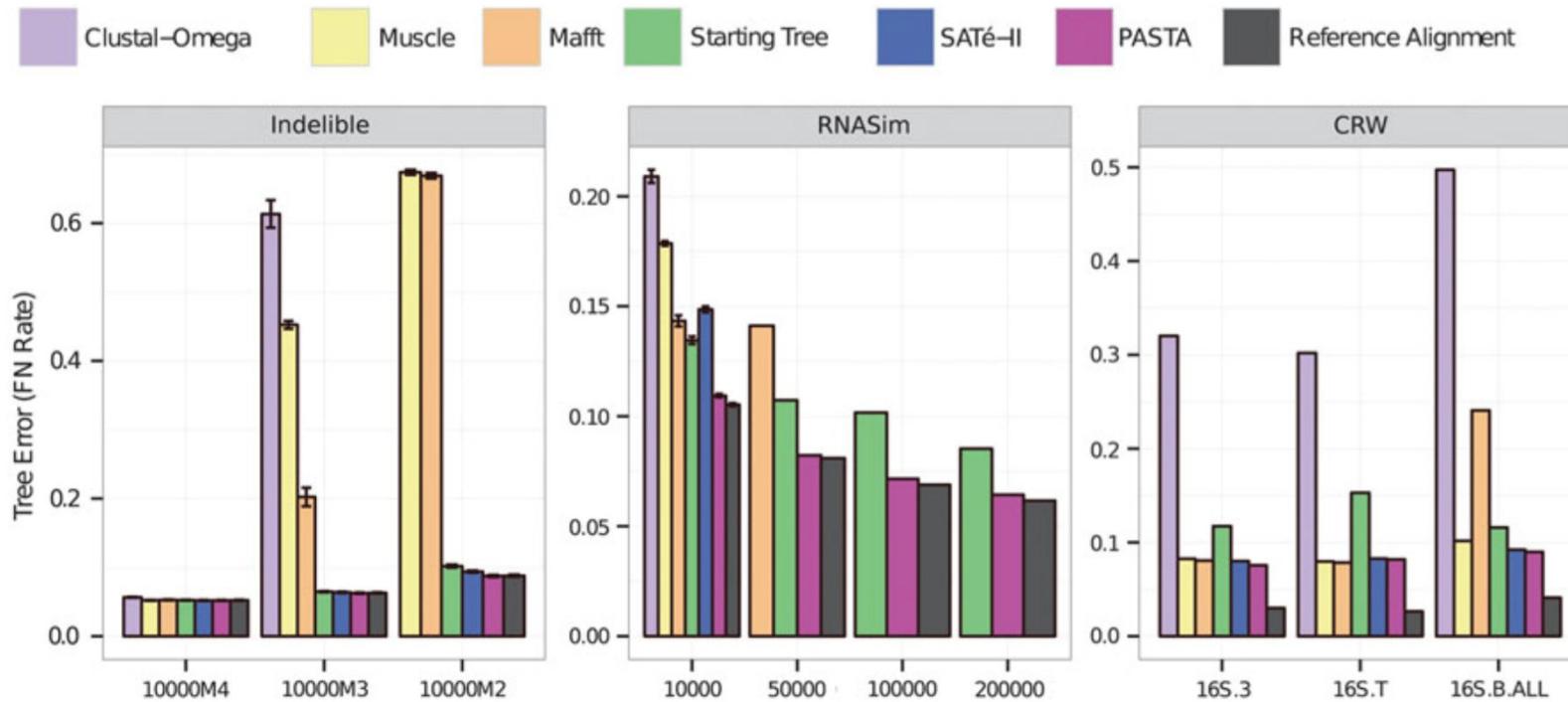
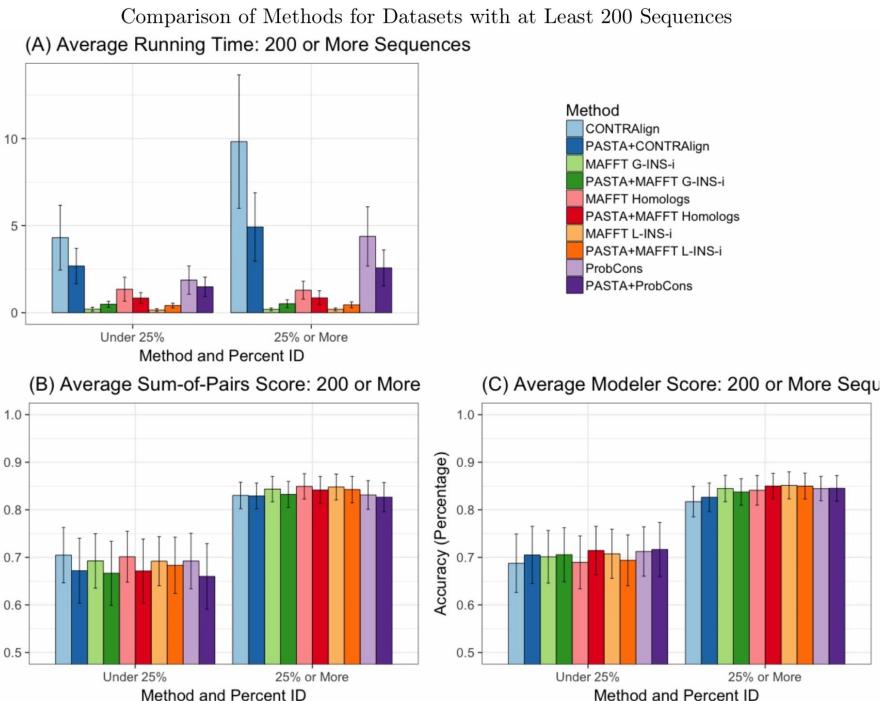
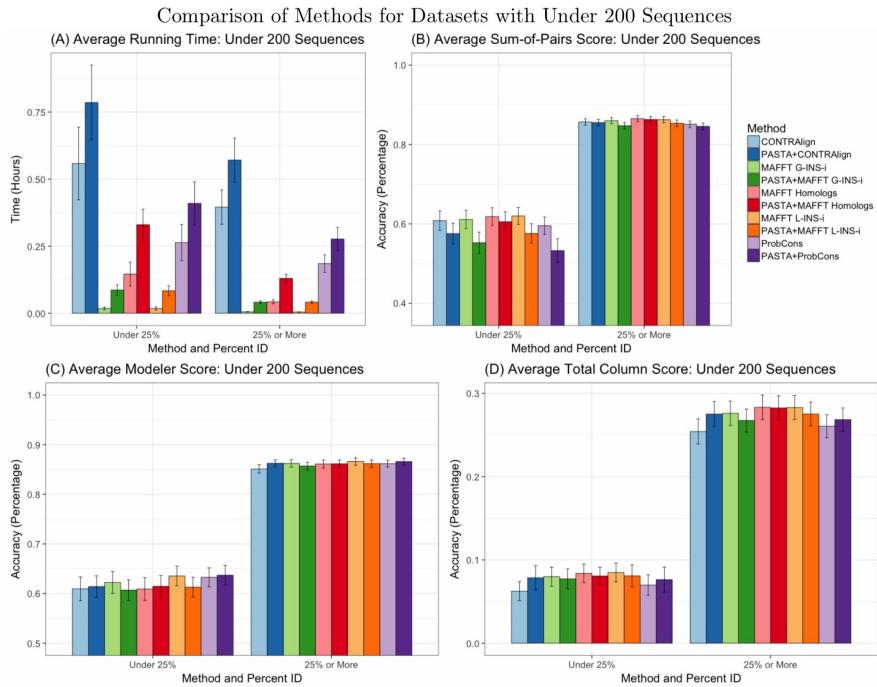


Fig. 2. Mirarab et al. 2015. *J. Comp. Biol.* 22(5):377–386.

# PASTA for proteins



# Running PASTA (from Command-line)

If your installation is successful, you should be able to run **PASTA** by running the following command from any location. Open up a terminal window and type:

```
run_pasta.py --help
```

Running **PASTA** with the `--help` option produces the list of options available in **PASTA**. **PASTA** automatically picks its algorithmic settings based on your input, so you can ignore most of these options (but `-d` is essential if you have anything other than DNA sequences). The basic command-line usage you need to know is:

```
run_pasta.py -i input_fasta_file
```

# Running PASTA (from Command-line)

The `-i` option is used to specify the input sequence file. The input file needs to be in the relaxed FASTA format. This command will start **PASTA** and will run it on your input file.

For a test run, use the `cd` command to go to the `data` directory under your **PASTA** installation directory. From there, run

```
run_pasta.py -i small.fasta
```

This will start **PASTA** and will finish quickly (30 seconds to 5 minutes based on your machine). Read **PASTA** output and make sure it finishes without producing any errors. If **PASTA** runs successfully, it produces a multiple sequence alignment and a tree, which we will explore in the next step.

# Inspecting the Output of PASTA

The two main outputs of **PASTA** are an alignment and a tree. The tree is saved in a file called `[jobname].tre` and the alignment file is named `[jobname].marker001.small.aln`. The `[jobname]` is a prefix which is by default set to `pastajob`, but can be changed by the user (see option `-j` below). When you start **PASTA**, if your output directory (which is by default where your input sequences are) already contains some files with the `pastajob` prefix, then the `pastajob1` prefix is used, and if that exists, `pastajob2` is used, and so forth. Thus the existing files are never overwritten. The name of your job and therefore the prefix used for output files can be controlled using the `-j` argument for command-line or the "Job Name" field on the GUI.

Tree Viewing Software (TVS) → [https://en.wikipedia.org/wiki/List\\_of\\_phylogenetic\\_tree\\_visualization\\_software](https://en.wikipedia.org/wiki/List_of_phylogenetic_tree_visualization_software)

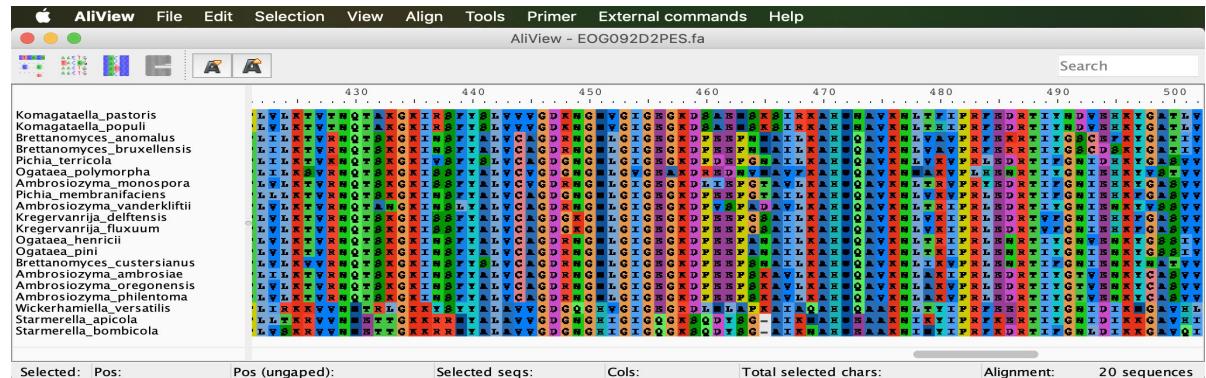
Alignment Viewing Software (AVS), e.g., <http://doua.prabi.fr/software/seaview> or <http://www.ormbunkar.se/aliview/>

# Light-weight AVS — AliView & SeaView

**AliView** → Larsson, A. (2014). AliView: a fast and lightweight alignment viewer and editor for large data sets.

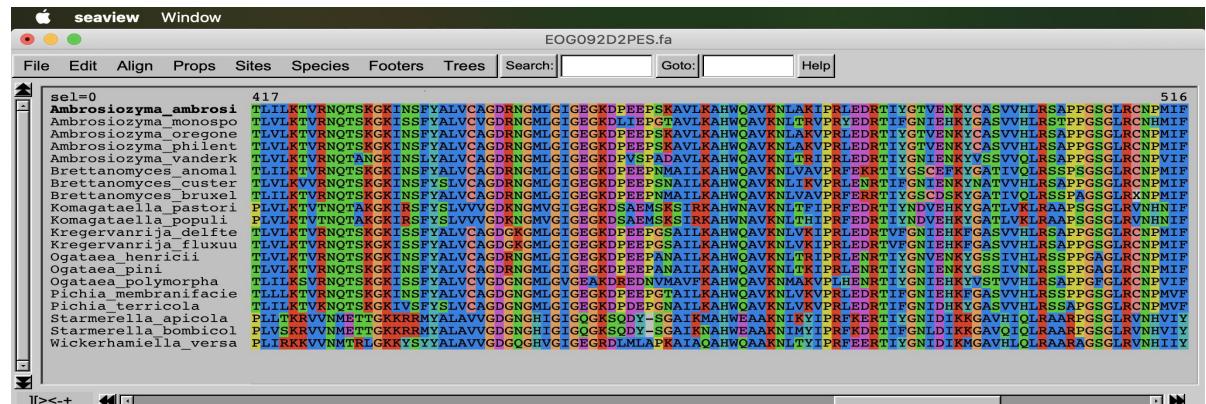
*Bioinformatics* 30(22):3276–3278.

<http://dx.doi.org/10.1093/bioinformatics/btu531>



**SeaView** → Gouy M., Guindon S. & Gascuel O. (2010) SeaView version 4 : a multiplatform graphical user interface for sequence alignment and phylogenetic tree building. *Molecular Biology and Evolution* 27(2):221–224.

<https://academic.oup.com/mbe/article/27/2/221/970247>



# Inspecting the Output of PASTA

## What about bootstrapping?

**PASTA** does not perform bootstrapping. The tree outputted by **PASTA**, depending on the options used, might include support values on the branches. These are *not bootstrap* support values. Instead, they are **SH-like local** support values computed by **FastTree**, and are generally believed to be not as reliable as bootstrap support values. In our experience they tend to overestimate support. Thus, if you want to have support values that can be trusted, we suggest that you use the **PASTA** alignment and an external tool (e.g., **RAXML**) for bootstrapping. If your alignment is too big for bootstrapping using **RAXML**, you can always use **FastTree** or **IQ-tree** for bootstrapping.

## Comparing alignments

When two alignments are generated on the same set of sequences, one can ask how similar they are. We have a tool called **FastSP** <http://www.cs.utexas.edu/~phylo/software/fastsp/> that compares two alignments and tells you how similar or different they are. **FastSP** does not require installation. You can just download it and run it (Java is required). Assuming **FastSP** is located at `~/bin/`, you can, e.g., compare the reference alignment and your estimated alignment:

```
java -jar ~/bin/FastSP_1.6.0.jar -r 16S.E.ALL.reference.fasta -e pastajob.marker001.16S.E.ALL.unaligned.aln
```

# Running PASTA (from Command-line)

You can script a while loop in bash to run **PASTA** on multiple fasta files. First open a text editor

```
nano pasta_loop.sh
```

Write your bash script

```
#!/bin/bash
while read targetname;
do
    python ABSOLUTE_PATH_HERE/run_pasta.py -i "$targetname".fasta -j $targetname
done < targetlist.txt
```

Close **CTRL+X** and save your script. This script assumes all target files are in the same folder in fasta format. It also assumes that folder contains a text file listing all targets. From there, run

```
bash pasta_loop.sh
```

# Understanding and Using PASTA Options

The command line allows you to alter the behavior of the algorithm using a variety of configuration options. Running **PASTA** with the `-h` option lists all the options that can be provided to the command-line (see below for the most important ones). In addition to the command-line itself, **PASTA** can read the options from one or more configuration files. The configuration files have the following format:

```
[commandline]
option-name = value
```

```
[state]
option-name = value
```

Note that as mentioned before, with every run, **PASTA** saves the configuration file for that run as a temporary file called `[jobname]_temp_pasta_config.txt` in your output directory. You can view one of these files in a Text editor for better understanding the format of the configuration file.

# Understanding and Using PASTA Options

**PASTA** can read multiple configuration files. Configuration files are read in the order they occur as arguments (with values in later files replacing previously read values). Options specified in the command line are read last. Thus these values "overwrite" any settings from the configuration files.

The following is a list of important options used by **PASTA**. Note that by default **PASTA** picks these parameters for you, and thus you might not need to ever change these (with the important exception of the `-d` option):

- **Initial tree:** As mentioned before, **PASTA** needs an initial tree for doing the first round of the alignment. Here is how the initial tree is picked.
  - If a starting tree is provided using the `-t` option, then that tree is used.

```
run_pasta.py -i small.fasta -t small.tree
```

# Understanding and Using PASTA Options

**PASTA** can read multiple configuration files. Configuration files are read in the order they occur as arguments (with values in later files replacing previously read values). Options specified in the command line are read last. Thus these values "overwrite" any settings from the configuration files.

The following is a list of important options used by **PASTA**. Note that by default **PASTA** picks these parameters for you, and thus you might not need to ever change these (with the important exception of the `-d` option):

- **Initial tree:** As mentioned before, **PASTA** needs an initial tree for doing the first round of the alignment. Here is how the initial tree is picked.
  - If a starting tree is provided using the `-t` option, then that tree is used.
  - If the input sequence file is already aligned and `--aligned` option is provided, then **PASTA** computes a ML tree on the input alignment and uses that as the starting tree.
    - If the input sequences are not aligned (or if they are aligned and `--aligned` is not given), **PASTA** uses the following procedure for estimating the starting alignment and tree. It 1) randomly selects a subset of 100 sequences, 2) estimates an alignment on the subset using the subset alignment tool (default **MAFFT-L-insi**), 3) builds a **HMMER** model on this "backbone" alignment, 4) uses **hmmpg** to align the remaining sequences into the backbone alignment, 5) runs **FastTree** on the alignment obtained in the previous step.

# Understanding and Using PASTA Options

- **Data type:** **PASTA** does not automatically detect your data type. Unless your data is DNA, you need to set the data type using `-d` command. Your options are DNA, RNA, and PROTEIN.

```
run_pasta.py -i BBA0067-half.input.fasta -t BBA0067-half.startingtree.tre -d PROTEIN
```

# Understanding and Using PASTA Options

- **Data type:** **PASTA** does not automatically detect your data type. Unless your data is DNA, you need to set the data type using `-d` command. Your options are DNA, RNA, and PROTEIN.
- **Tree estimation tool:** the default tool used for estimating the phylogenetic tree in **PASTA** is **FastTree**. The only other option currently available is **RAxML**. You can set the tree estimator to **RAxML** using the `--tree-estimator` option. However, Be aware that **RAxML** takes much longer than FastTree. If you really want to have a **RAxML** tree, we suggest obtaining one by running it on the final **PASTA** alignment. You can change the model used by FastTree (default: `-nt -gtr -gamma` for nt and `-wag -gamma` for aa) or **RAxML** (default `GTRGAMMA` for nt and `PROTWAGCAT` for AA) by updating the `[model]` parameter under `[FastTree]` or `[RAxML]` header in the input configuration file. The model cannot be currently updated in the command line directly as an option.
- **Subset alignment tool:** the default tool used for aligning subsets is **MAFFT**, but you can change it using the `--aligner` option. We strongly suggest alignment subset size should always be no more than 200 sequences, because for subsets that are larger than 200, the most accurate version of **MAFFT** (`-linsi`) is not used.
- **Pairwise merge tool:** the default merger too is **Opal**. You can change it using `--merger` option. If you have trouble with **Opal** (java version, memory, etc.) using **Muscle** should solve your problem and in our experience, it doesn't really affect the accuracy by a large margin.

# Understanding and Using PASTA Options

- **CPUs:** **PASTA** tries to use all the available cpus by default. You can use `--num_cpus` to adjust the number of threads used.

```
run_pasta.py -i small.fasta --num_cpus 1
```

# Understanding and Using PASTA Options

- **CPUs:** **PASTA** tries to use all the available cpus by default. You can use `--num_cpus` to adjust the number of threads used.
- **Number of iterations:** the simplest option that can be used to set the number of iterations is `--iter-limit`, which sets the number of iterations **PASTA** should run for. You can also set a time limit using `--time-limit`, in which case, **PASTA** runs until the time limit is reached, and then continues to run until the current iteration is finished, and then stops. If both options are set, **PASTA** stops after the first limit is reached. The remaining options for setting iteration limits are legacies of **SATé** and are not recommended.
- **Masking:** Since **PASTA** can produce very gappy alignments, it is a good idea to remove sites that are almost exclusively gaps before running the ML tree estimation. By default, **PASTA** removes sites that are more than 99.9% gaps. You can change that using the `--mask-gappy-sites` option. For example, using `--mask-gappy-sites 10` would remove sites that are gaps for all sequences except for (at most) 10 sequences. Increasing the masking can make **PASTA** a bit faster and can potentially reduce the memory usage. But it could also have a small effect on the final tree. If unsure, leave the option unchanged. Note that the final alignment outputted by **PASTA** is NOT masked, but masked versions of the output are also saved as temporary files (see below).

# Understanding and Using PASTA Options

- **Maximum subset size:** two options are provided to set the maximum subset size: `--max-subproblem-frac` and `--max-subproblem-size`. The `--max-subproblem-frac` option is a number between 0 and 1 and sets the maximum subset size as a fraction of the entire dataset. The `--max-subproblem-size` option sets the maximum size as an absolute number. When both numbers are provided (in either a configuration file or the command line), the *LARGER* number is used. This is an unfortunate design (legacy of SATé) and can be quite confusing. Please always double check the actual subset size reported by PASTA and make sure it is the value intended. The default subset sizes should work just fine. In our limited experiments, we have noticed that reducing the maximum subset size from 200 to 100 for very large datasets increases speed with little or no effect on the final alignments.
- **Temporary files:** PASTA creates many temporary files, and deletes most at the end. You can control the behavior of temporary files using few options: `--temporaries` sets the directory where temp files are created, `-k` instructs PASTA to keep temporary files, and `--keepalignmenttemps` will keep even more temporary files. Note that these are different from the temporary files created in the output directory (which are always kept).
- **Dry run:** The `--exportconfig` option can be used to just create a config file and exit without actually running PASTA. This is useful for making sure the configurations are correct before actually running the job.

# Running PASTA Using Configuration Files

The configurations used for running **PASTA** are all saved to a configuration file, and also, **PASTA** can be run using a configuration file. These configuration files are useful for multiple purposes. For example, if you want to reproduce a **PASTA** run, or if you want to report the exact configurations used. Always make sure to keep the produced configuration files for future reference. Note however, that configuration files can be used as input only using command-line.

Let's open `myjob_temp_pasta_config.txt` under the data directory and take a look at it. Notice that the options we referred to are all mentioned here.

Now imagine that we wanted to instruct **PASTA** to use the JTT model instead of WAG for a protein run. Here is how we can accomplish that. Copy the `myjob_temp_pasta_config.txt` file as a new file (e.g. `cp myjob_temp_pasta_config.txt jtt_config.txt`). Then open `jtt_config.txt` using a text editor of your choice. Find `model = -wag -gamma -fastest` under the `[FastTree]` header. Remove the `-wag` option and save the config file. Note that the default model in FastTree is JTT, and therefore, when the `-wag` is removed, it automatically switches to using JTT. To run **PASTA** using this new configuration file, run:

```
run_pasta.py jtt_config.txt
```

# Running PASTA Using Configuration Files

**Adding custom parameters to aligners:** It is also possible to add custom parameters to alignment and merge tools. To do so, you need to use the config file. Under each alignment tool in the config file, you can add an `args` attribute and list all the attributes you want to pass to that tool. For example, to run **MAFFT** with your choice of gap penalty value, edit the config file under the `[mafft]` heading to something like:

```
[mafft]
path = [there will be a path here to your pasta directory]/bin/mafft
args = --op 0.2 --ep 0.2
```

and use this config file to run **PASTA**.

Note that **PASTA** does not try to understand these extra parameters you pass to external tools. It simply appends these parameters to the end of the command it executes.

At this stage, if you have input files that you like to have analyzed, you know enough to start doing that.

Email: [pasta-users@googlegroups.com](mailto:pasta-users@googlegroups.com) for all issues.

# MAFFT

**MAFFT** is a multiple sequence alignment program for unix-like operating systems. It offers a range of multiple alignment methods, **L-INS-i** (accurate; recommended for <200 sequences), **FFT-NS-2** (fast; recommended for >2,000 sequences), etc. Accuracy-oriented methods:

\***L-INS-i** (probably most accurate; recommended for <200 sequences; iterative refinement method incorporating local pairwise alignment information):

```
mafft --localpair --maxiterate 1000 input [> output]  
linsi input [> output]
```

\***G-INS-i** (suitable for sequences of similar lengths; recommended for <200 sequences; iterative refinement method incorporating global pairwise alignment information):

```
mafft --globalpair --maxiterate 1000 input [> output]  
ginsi input [> output]
```

\***E-INS-i** (suitable for sequences containing large unalignable regions; recommended for <200 sequences):

```
mafft --ep 0 --genafpair --maxiterate 1000 input [> output]  
einsi input [> output]
```

# Bali-Phy

## BAli---Phy

[Welcome!]

- Introduction
- References

Download

Examples

Documentation

- Users Guide
- Tutorial
- Manual pages
- Developers Guide

Contact

- Bug reports
- Mailing lists

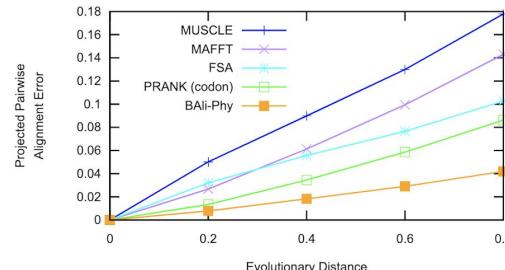
## Introduction

01/18/19: **BAli-Phy 3.4.1 released** - [Download](#)

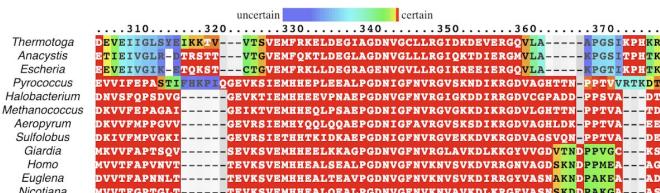
16-state RNA stem models + character sets + bug fixes ([release notes](#))

*BAli-Phy* is software by Ben Redelings that estimates multiple sequence alignments and evolutionary trees from DNA, amino acid, or codon sequences. It uses likelihood-based evolutionary models of substitutions and insertions and deletions to place gaps. It has been used in published analyses on data sets up to 117 taxa.

**High alignment accuracy:** Redelings (2014) showed that *BAli-Phy* had 3.5 times fewer alignment errors than MUSCLE and MAFFT on simulated data:



**Eliminate bias:** Fletcher and Yang (2010) showed that relying on a ClustalW alignment estimate could lead to a 99% false-positive rate in detecting positive selection. In general, inferring evolutionary trees, branch lengths, or positive selection from a single alignment can lead to bias if the alignment is ambiguous. *BAli-Phy* solves this problem by using MCMC and Bayesian methods to estimate evolutionary trees, positive selection, and branch lengths while averaging over alternative alignments.



# Profile HMM Methods: SEPP, TIPP, UPP, & HIPPI

**SEPP** stands for "SATé-enabled Phylogenetic Placement", and addresses the problem of phylogenetic placement of short reads into reference alignments and trees.

**TIPP** stands for "Taxonomic Identification and Phylogenetic Profiling", and addresses the problem of taxonomic identification and abundance profiling of metagenomic data.

**UPP** stands for "Ultra-large alignments using Phylogeny-aware Profiles", and addresses the problem of **alignment of very large datasets, potentially containing fragmentary data**. UPP can align datasets with up to 1,000,000 sequences.

**HIPPI** stands for "Highly Accurate Protein Family Classification with Ensembles of HMMs", and addresses the problem of classifying query sequences to protein families.

# UPP

UPP is a modification of SEPP for performing **alignments of ultra-large** and **fragmentary datasets**. UPP operates in four steps:

- In the first step, UPP partitions set S into a **backbone set** and a **query set** and computes an alignment and tree on the backbone set using PASTA, which is a direct improvement to SATé.
- In the next step, UPP decomposes the **backbone alignment** into an ensemble of profile Hidden Markov Models (HMMs).
- The third step in UPP searches for the **best alignment** of the **query sequence** to each HMM.
- The final step inserts the **query sequence** into the **backbone alignment** using the best scoring HMM.

Our study shows that UPP results in accurate alignments, and that ML trees estimated on the alignments are also highly accurate. UPP has **good accuracy** on datasets that contain **fragmentary sequences**.

# UPP

To run **UPP**, invoke the `run_upp.py` script from the `bin` sub-directory of the location in which you installed the Python packages.

To see options for running the script, use the command: `python <bin>/run_upp.py -h`

The general command for running **UPP** is: `python <bin>/run_upp.py -s <unaligned_sequences>`

This will run **UPP(Default)**. This will automatically select up to 1,000 sequences to be in the backbone set, generate a **PASTA alignment** and **tree**, and then align the remaining sequences to the backbone alignment.

# UPP

To run **UPP**, invoke the `run_upp.py` script from the `bin` sub-directory of the location in which you installed the Python packages.

To see options for running the script, use the command: `python <bin>/run_upp.py -h`

The general command for running **UPP** is: `python <bin>/run_upp.py -s <unaligned_sequences>`

This will run **UPP(Default)**. This will automatically select up to 1,000 sequences to be in the backbone set, generate a **PASTA alignment** and **tree**, and then align the remaining sequences to the backbone alignment.

**UPP** can also be run using a configuration file. To run using a configuration file, run: `python <bin>/run_upp.py -c sample.config`

To run **UPP(Fast)**, run: `python <bin>/run_upp.py -s input.fas -B 100`

**UPP** currently assumes that the input sequences are nucleotide sequences. To select the input data type, run: `python <bin>/run_upp.py -s input.fas -m [dna|rna|amino]`

# Trimming

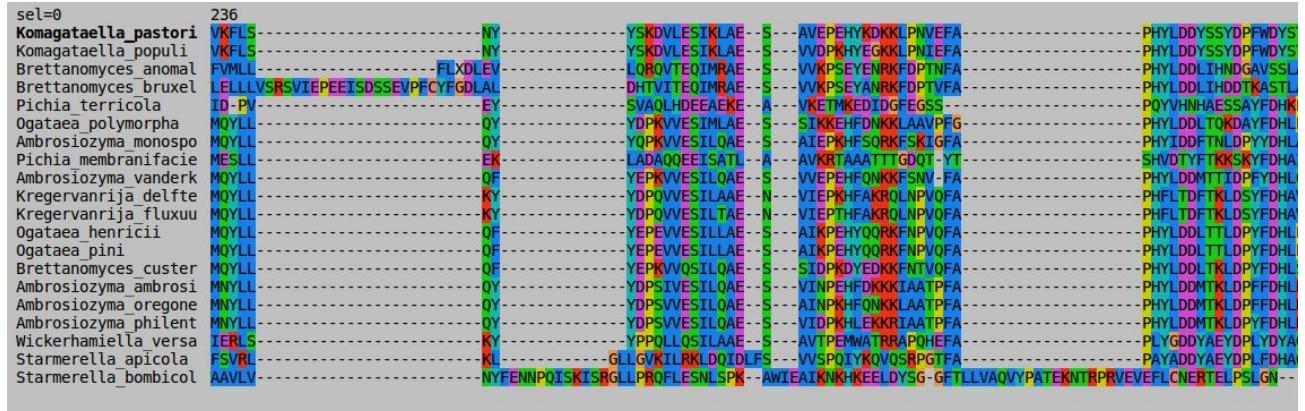
Lisa Pokorny & Marina Marcet-Houben

## Multiple sequence alignments can have many different forms

Conserved

Komagataella_pastori	EDAKKEEAIVRHDMVAHVHTFGKTCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Komagataella_populi	EDAKKEEAIVRHDMVAHVHTFGKTCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ogataea_polymorpha	EAAKKEEAIVRHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ogataea_henricii	EAAKKEEAIVRHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ogataea_pini	EAAKKEEAIVRHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ambrosiozyma_monospo	EAAKKEEAIVRHDMVAHVHTGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ambrosiozyma_vanderk	EAAKKEEAIVRHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ambrosiozyma_ambrosi	EAAKKEEAIVRHDMVAHVHTGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ambrosiozyma_oregone	EAAKKEEAIVRHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Ambrosiozyma_philent	EAAKKEEAIVRHDMVAHVHTGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Kregervanrija_delfte	EIAKIEESKVRDHDMVAHVHTGQTCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Kregervanrija_fluxuu	EIAKIEESKVRDHDMVAHVHTGOTCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Pichia_membranifacie	EAAKIEESKVRDHDMVAHVHFGETCPEAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Pichia_terricola	EDAKKEESAVRHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Brettanomyces_custer	EAAKKEEARVRDHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Brettanomyces_anomal	EAAKKEEARVRDHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Brettanomyces_bruwel	EAAKKEEARVRDHDMVAHVHFGETCPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Wickerhamiella_versa	WAASKOEAEIVRHDMVAHVHEFGVECPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Starmerella_apicola	EATKOEAEIVRHDMVAHVHFGECPAAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL
Starmerella_bombicola	EAGKKOEAEIVRHDMSHVHQYGLLEPAAGIHLGATSCVITDADLIFLIPKLVNVIDRLSKFALEYKDLPLVGWTHFOPAQLTTVGKRSTLWLQELLWLDRNMRPARNDIGLRGAKGTGTOASFSL

Noisy



# Why can alignments be noisy?

Biological reasons: If we compare sequences from proteins from distantly related species, there is high chance that only the functional part of the protein is well conserved in terms of sequences. Other parts, such as loops, are more likely to have altered their amino acid sequence, both in terms of amino acid content and with the presence of indels. Even when the prediction of the multiple sequence alignment is correct, it may negatively affect the inference of the phylogenetic tree.

## Errors:

- Errors derived from genome assembly problems
- Errors derived from gene prediction
- Errors derived from alignment of multiple sequences

Different trimming programs aim to solve some of the problems caused by the presence of gaps or badly aligned regions.

Biological reasons

Errors in multiple sequence alignments

Traditional block-based methods  
(Gblocks, trimAI, BMGE, Zorro, ...)

Errors in sequencing

Errors in gene predictions

Segment based methods  
(Prequal, hmmcleaner)

## To trim or not to trim

The trimming of multiple sequence alignments is still a controversial topic. Some authors claim that it's necessary to improve phylogenetic reconstruction and the detection of evolutionary events.

### **Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments.**

Talavera G<sup>1</sup>, Castresana J.

Whereas some others claim it either does not have any effect or it make the problem worse.

### **Current Methods for Automated Filtering of Multiple Sequence Alignments Frequently Worsen Single-Gene Phylogenetic Inference**

Ge Tan,<sup>1,2,3</sup> Matthieu Muffato,<sup>4</sup> Christian Ledergerber,<sup>1</sup> Javier Herrero,<sup>4,5</sup> Nick Goldman,<sup>4</sup> Manuel Gil,<sup>6,7</sup> and Christophe Dessimoz<sup>5,4,\*</sup>

Table 4  
Topological accuracy of single-gene phylogenies

VERTEBRATA		
version	mean	loss (%)
RAW	65.64%	NA
RAW (long)	68.85%	NA
HMM	65.23%	5.61
HMM-L	65.41%	4.22
HMM-LS	65.58%	2.68
PREQUAL	65.73%	3.06
BMGE	64.83%	4.83
TrimAl	65.28%	1.8
HMM Random	64.56%	5.61
HMM + BMGE	63.94%	13.38
HMM + TrimAl	62.90%	13.76
MIN	68.71%	0.71
MIN + HMM	68.67%	6.43

There's little difference.

BMC Evol Biol. 2019 Jan 11;19(1):21. doi: 10.1186/s12862-019-1350-2.

## Evaluating the usefulness of alignment filtering methods to reduce the impact of errors on evolutionary inferences.

Di Franco A<sup>1</sup>, Poujol R<sup>2</sup>, Baurain D<sup>3</sup>, Philippe H<sup>4,5</sup>.

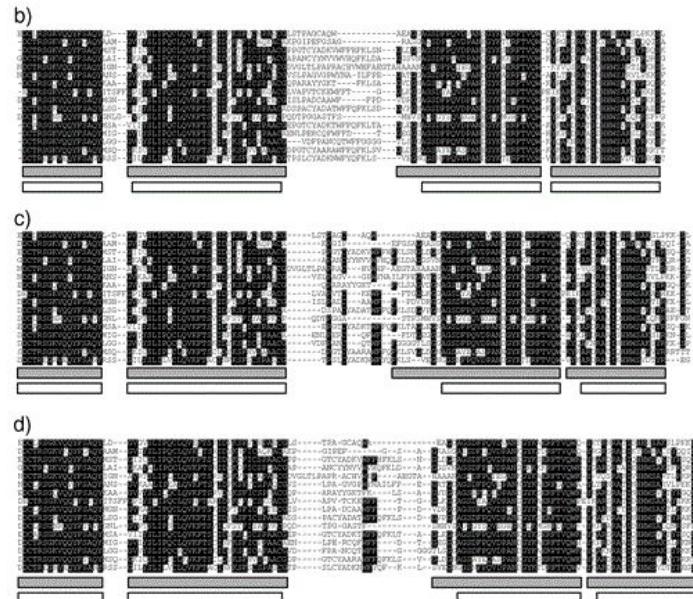
The probable truth: it depends on the dataset and the methodology used.

<b>Program</b>	<b>Number of citations (Google scholar)</b>
BMGE	381
trimAl	1737
Gblocks	5736

Trimming alignments tends to be part of a normal phylogenetic reconstruction pipeline.

The first one: GBLOCKS (<http://molevol.cmima.csic.es/castresana/Gblocks.html>)

Gblocks selects blocks in a similar way as it is usually done by hand but following a reproducible set of conditions. The selected blocks must fulfill certain requirements with respect to the lack of large segments of contiguous nonconserved positions, lack of gap positions and high conservation of flanking positions, making the final alignment more suitable for phylogenetic analysis.



The white and grey blocks under the alignments represent the parts of the alignment that Gblocks would keep using a more relaxed and a more stringent approach.

## How to run Gblocks: website

Gblocks Server

Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis

**About the Gblocks Server**

Version 0.91b, January 2002

Copyright © Jose Castresana

Gblocks eliminates poorly aligned positions and divergent regions of a DNA or protein alignment so that it becomes more suitable for phylogenetic analysis. This server implements the most important features of the Gblocks program to make its use as simple as possible without loosing the functionality that it is necessary in most of the cases. Other options can be changed in the stand-alone program. You can see here an example output file showing the blocks selected from a protein alignment. Further information can be found in the online documentation. Please see the Gblocks page for citations.

**Gblocks Server**

Paste an alignment in NBRF/PIR or FASTA format:

Or upload an alignment file:  
 No file selected.

Type of sequence:

DNA  Protein  Codons

Options for a less stringent selection:

Allow smaller final blocks  
 Allow gap positions within the final blocks  
 Allow less strict flanking positions

Options for a more stringent selection:

Do not allow many contiguous nonconserved positions

The is an on-line server that you can use if you only want to trim one alignment.



At the end of the alignment representation there's a link to obtain the trimmed alignment.

## How to run Gblocks: command line

```
mmarcet@saturn:~/Desktop/evomics$ Gblocks
*****
GBLOCKS 0.91b
SELECTION OF CONSERVED BLOCKS FROM MULTIPLE ALIGNMENTS
FOR THEIR USE IN PHYLOGENETIC ANALYSIS
*****
```

**o. Open File**      Used to upload your alignment file (Fasta or NBRF/PIR format)

**b. Block Parameters**

**s. Saving Options**

**g. (Get Blocks)**

**q. Quit**

Your Choice: █

Your Choice: o

Please enter a DNA or protein alignment in NBRF/PIR or FASTA format or a paths file.  
File Name: cytochrome.alg.forward.fasta

96 sequences and 1543 positions in the first alignment file:  
cytochrome.alg.forward.fasta

Once the alignment has been introduced it will go back to the main menu and  
you can choose option g which will execute the program with the default  
parameters.

Your Choice: g

cytochrome.alg.forward.fasta

Original alignment: 1543 positions

Gblocks alignment: 148 positions (9 %) in 8 selected block(s)

## Example Gblocks -

A- Start Gblocks and open the file called EOG092D2PES.alg (use the o option). Then use g to obtain the blocks with the default parameters. Pay attention to how long your resulting alignment is and which percentage of the alignment has been kept with the default parameters.

If you go to the folder where the sequences are, you will see that there are now two new files: EOG092D2PES.alg-gb contains the trimmed alignment while EOG092D2PES.alg-gb.htm has a html representation of the alignment and the blocks that have been kept. Open it and have a look. Afterwards, rename the file to EOG092D2PES.default.htm so that it is not lost.

Option b in the Gblocks menu allows you to change the parameters by which the blocks are defined.

B.- Use this option to adjust the minimum number of sequences needed to define a conserved position and the minimum number of sequences to define a flanking region. Assign a random number within the accepted scope. How does this affect your alignment?

Again have a look at the results. Do you see a change between the blocks that have been selected? Again rename the file so that it is not lost

C.- Now decrease the length of your conserved blocks. Did that have an effect on your alignment? What happens if you increase it?

Check out the results again and compare them with the previous runs of Gblocks.

D.- These alignments are in general pretty conserved. What do you think would happen if we were comparing more distantly related species?



**A tool for automated alignment trimming**

trimAI was initially born because Gblocks could be too restrictive when automatically building thousands of alignments. Unlike Gblocks, trimAI implements different trimming strategies based on gap content, similarity or consistency across different alignment methods. It also implements a conservation score which always ensures that a percentage of the alignment is conserved.

(<https://github.com/scapella/trimai>)

trimAI has a wide array of options, including user defined trimming parameters:

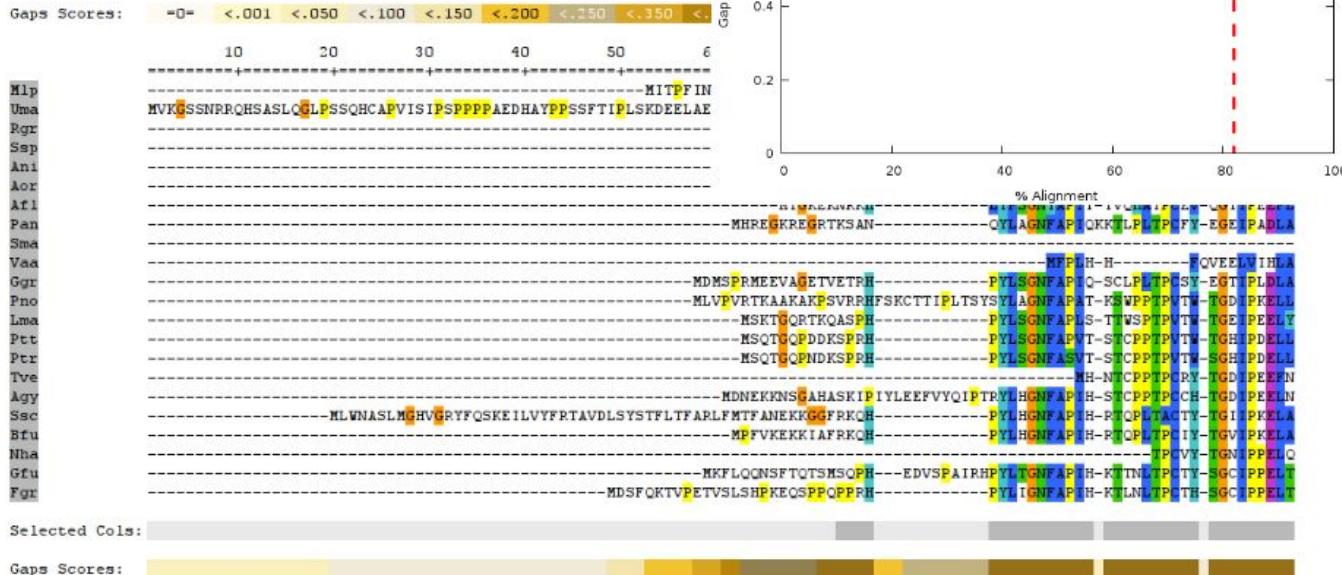
```
-gt -gapthreshold <n>      1 - (fraction of sequences with a gap allowed). Range: [0 - 1]
-st -simthreshold <n>      Minimum average similarity allowed. Range: [0 - 1]
-ct -conthreshold <n>      Minimum consistency value allowed. Range: [0 - 1]
-cons <n>                  Minimum percentage of the positions in the original alignment to conserve. Range: [0 - 100]
```

And automated methods that predict the best parameters for a given alignment:

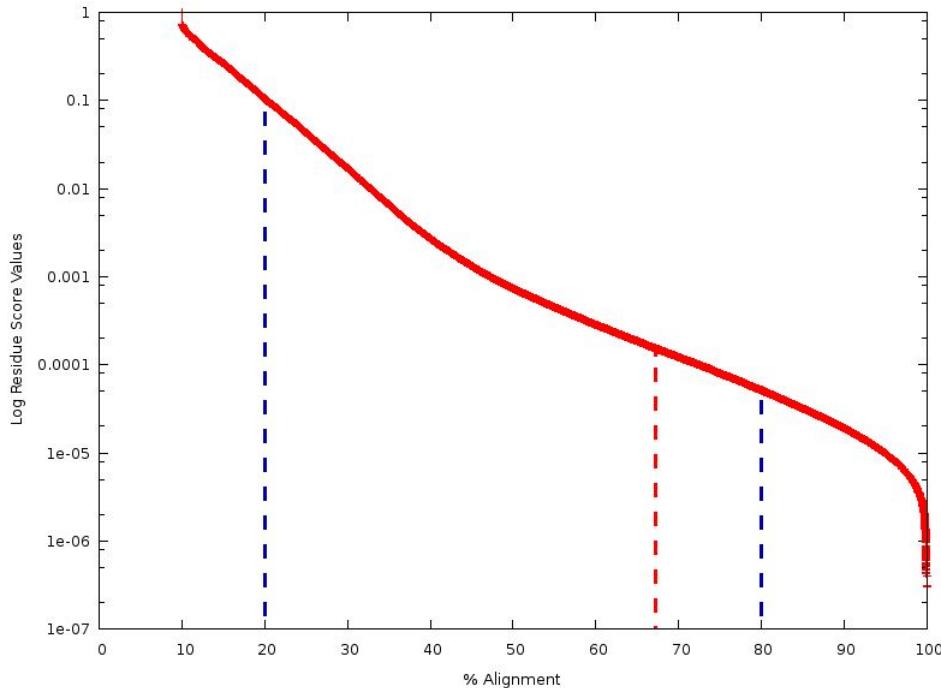
```
-gappyout                 Use automated selection on "gappyout" mode. This method only uses information based on gaps' distribution. (see User Guide).
-strict                   Use automated selection on "strict" mode. (see User Guide).
-strictplus                Use automated selection on "strictplus" mode. (see User Guide).
                           (Optimized for Neighbour Joining phylogenetic tree reconstruction).
```

Automated methods:

## -gappyout



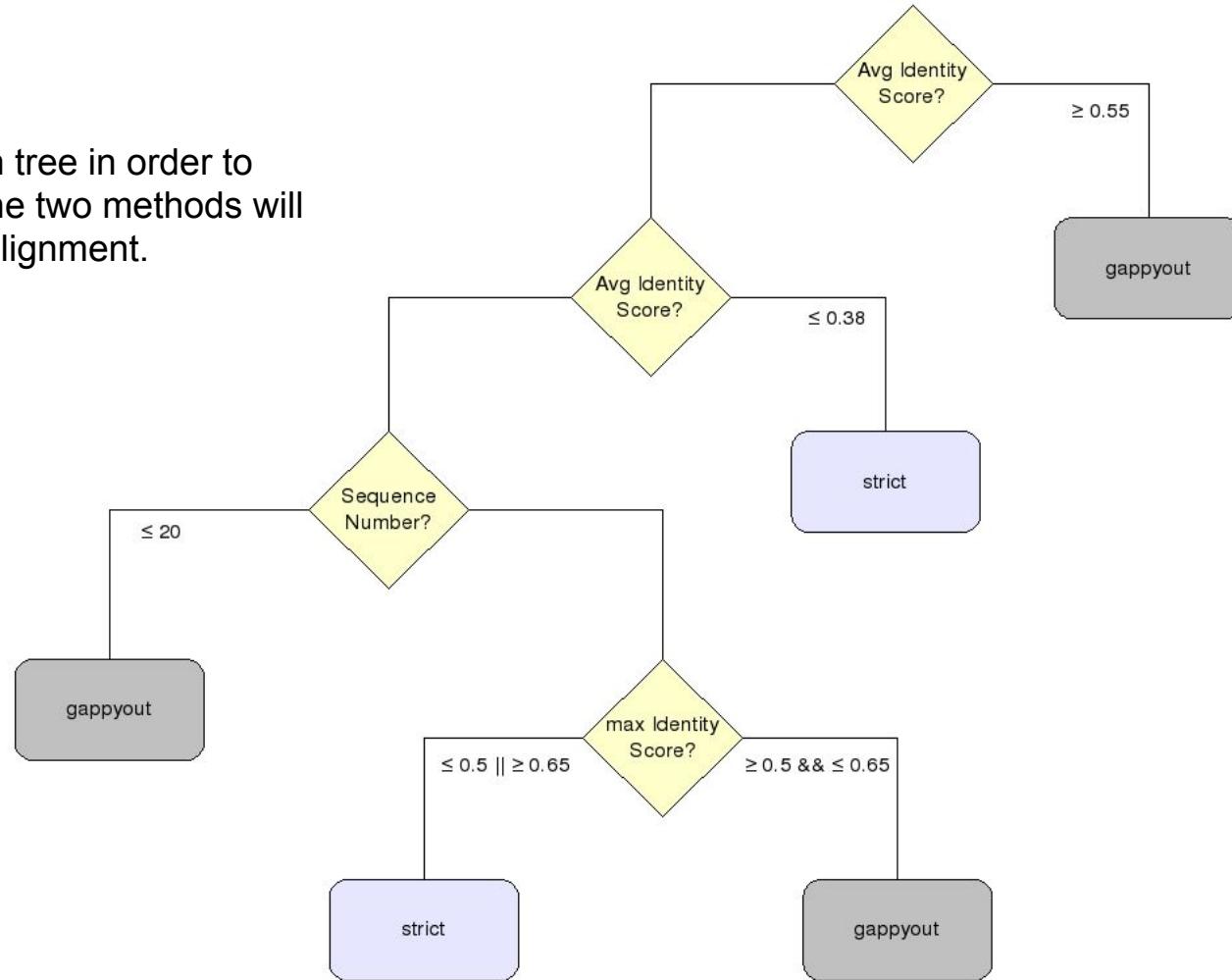
## **-strict -strictplus**



Gappy-out + trimming by similarity scores → they will only delete blocks of data so if one column has been marked to be deleted but it is surrounded by non-marked columns it will be kept in the alignment. The two methods differ on how they define the block size.

# -automated1

Will use a decision tree in order to choose which of the two methods will work best on the alignment.



# readAI: Reformatting MSAs

One of the main problems of alignments is the fact that different formats exist, and there may not be a match between the output format of an alignment program and the input format the next program needs.

```
#NEXUS
[!Imported PHYLIP file "030103phylip.phy" (Fri Jan 03 12:43:38 2003)]
Begin data;
  Dimensions ntax=29 nchar=949;
  Format datatype=nucleotide gap=- missing=? matchchar=. interleave;
  options gapmode=missing;
  Matrix
  487GJS AACGTTACCAAAACTCTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  484GJS AACGTTACCAAAACTCTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  476GJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  481GJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  497GJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  501GJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  477GJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  493GJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  486GJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  TmnuGJS AACGTTACCAAAACTGTGGCCTGGGGGGAAAAA--TTCCATCGCCCCGGG
  Tharsia GCCCATTCTACGGAGACATCATTGAGACACCCGCTGGTATGGCAGACT
  Tvirnes GCCCATTCTACGGAGACATCATTGAGACACCCGCTGGTATGGCAGACT
  Thematium GCCCATTCTACGGAGACATCTTCGAGACGCCGGTGAATTGGCAGCT
  473GJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  Hplululi GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  153DGJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  130DGJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  135DGJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  139DGJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  147GJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  138DGJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  491GJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  460GJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  467GJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  124DGJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  150DGJS GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  croceum GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
  polyp poment GCTCATTTAACGAGACGATCTTCGAGACGCCGGTGAATTGGCAGCT
;
```

End;  
begin sets;  
charpartition genes = ITS:1-383, EF:384-495, ECH:496-; ;  
end;

NEXUS format

FASTA format

```
>TRY2_RAT/24-239
-----IVGGYTCQENSVPYQVSLNSGY-----HFC
GGS LI-----NDQ-WV-VSAAHCYKS-----RIQVRLGE-HNINVLEG N-----HFC
-----EQFVNAAKIIKHPNFRKT-L-----NNDIMI LKLS
SP-----VKLNARVATVALPS---SCA---PAGTQLCLISGWGN-----TLSSGV-----
-----NEPDLLQ-CLDAP-LLPQADCEAS---YPGK-----ITDNMVCVGFL-----
-EGG-KDSCQGDGGPPVVCNGE-----LQGIVSWG-YCALPDN-----PGVYTKVCNVY
VDWI-----
```

>Q16LB2\_AEDAE/136-374
-----ILNGIEADLEDFPYL GALALLDNYT-----STVSYRC
GANLI-----SDR-FM-LTAHCLFG-----KQAIHV RMGTL SLDNPDED
-----APVIIGVERVFHR NYTRRPIT-----RN DIALIKLN
RT-----VVDFLIPVCLYT-----EQNDP-LPTVPLTIAGWGG-----NDSAS
-----LMSSSL-KASVT-TYERDECNSL---LAKKI-----VRLSN DQLCALGRSEF
NDGLRN DTCVGDSGGPLEL SIGR-----RKYIVGLTST-IVCGNE-F---PSIYTRISQF
IDWI-----

```
| 150 1075
Phy0007P00_GIBZE
Phy0007P00_GIBZE
Phy0008000D_SCHPO
Phy0009QV7 ASPFU
Phy0009QV7 ASPFU
Phy0009FRU5 ASPFU
Phy0009FRU5 ASPFU
Phy0009F06_ ASPFU
Phy0009F06_ ASPFU
Phy0009UZH_ ASPFU
Phy0009UZH_ ASPFU
Phy0009V06_ ASPFU
Phy0009V06_ ASPFU
Phy0009W15_ ASPFU
Phy0009W15_ ASPFU
Phy0009X01_ EMENI
Phy0009X01_ EMENI
Phy0009YXQ_ EMENI
Phy0009YXQ_ EMENI
Phy0009ZYA_ EMENI
Phy0009ZYA_ EMENI
Phy0009G15_ EMENI
Phy0009G15_ EMENI
Phy0009G06G_ EMENI
Phy0009G06G_ EMENI
Phy0009G07_ EMENI
Phy0009G07_ EMENI
Phy0009G12_ EMENI
Phy0009G12_ EMENI
Phy0009G28_ EMENI
Phy0009G28_ EMENI
Phy0009G57Q_ ASPNG
Phy0009G57Q_ ASPNG
Phy0009C9H_ ASPNG
Phy0009C9H_ ASPNG
Phy0009CB2J_ ASPNG
Phy0009CB2J_ ASPNG
Phy0009CBXK_ ASPNG
Phy0009CBXK_ ASPNG
Phy0009GEPP_ ASPNG
Phy0009GEPP_ ASPNG
```

-----PG-A-FIP-WDKYGLL RVC G-----N-----VGTSV R
-----LR-----R-----HHSKH-R-TCL-OFASLKR-SIGS VD SM-----GESCKAP-Q
MFHT-----AL-N-PRA-A-----TGEER-G
-----MDP-----N-----NSSSST-QLSPK VLS PRITKA-----P-R-----T
-----P-----S-----PTAER-SVNSKV A/PBSANPS-----S-W-----T-S
-----RTS-LQ-PG-SSK-YINL-NERVYK A M-----Q-HQL QPS-----S-PQR P-K
MFHT-----EG-P-GAS-APAK-----GRE-----RQSVG-R
-----MSGPHEEL SAD-----YQLQS GGYFH YARS P-L
-----MA-----DY-QN-VRPL R-E
-----MAS-T
-----G-----SMSEQ-R-----PSEPS-TPGCK1P1P RVSOLR-----A-V-----PAGPAEP-S
-----H-----DASVOPQDHLVLA-----EA-A-NAS-KOLESGPRNPKSV-----T-----GR TSGP-B
MFIA-----EA-A-NAS-KOLESGPRNPKSV-----T-----LTNS-AV-
-----PNDIR-A-LIT-GASCGIGAACAHOLHLA-----M-----SSRGP T-R
-----QOYEALQVSP1LRQ RSTL-A-VIG-GDL LQNGHSA-----SEAY-----GSVFRS R-H
-----M-----SKNMG S-T
MFHT-----EG-P-AAS-APARA-----GRE-----RQSSIG-R
-----TDPH-----GAQPGIHLASRLQGR-K
-----DIQS R-R
-----MYRS-----S-----RHV R-K
MASA-----AD-----EDDDSSFF-E-DHD-ASPGHDLMMKDAL-----GDGP-----DPLPQ-K
-----MNTHREGEPLA-K
-----MAS-T
-----SEV-

PHYLIP format

# readAI: Reformatting MSAs

readAI is a sister program to trimAI that allows us to convert alignment between each other.

**readal -in [input file] -format -out [output file]**

Input file → Alignment file

Output file → Resulting file

Format → Can be any of the formats that readAI has and that you wish to use as output:  
Fasta, phylip, mega, nexus, clustal,...

Use one of the previous alignments to try out the following things

1.- Open the alignment file (EOG092D2PES.alg) and check in which format it has been generated.

Now use readal to (make sure each result is in a different file):

- Change the format of the current alignment to fasta format
- Change the format of the current alignment to nexus format
- Change the format of the current alignment to clustal format
- Use the -onlyseqs option

Open the different files and notice the differences between the alignment formats.

*Tip: readal is run like this:*

```
readal -in alignment_file -out trimmed_alignment_file -format  
FORMAT_NAME
```

You can check out all the formats supported by readal by typing: readal -h

2.- Use trimAI to trim the alignment (EOG092D2PES.alg) according to a gap threshold using the following parameters:

- A gap threshold of 0.1 (-gt 0.1)
- A gap threshold of 0.5 (-gt 0.5)
- A gap threshold of 0.9 (-gt 0.9)

Make sure that the output of your alignment is in phylip format. Now you can visualize each alignment either using a text editor or using seaview. Which of the previous commands deletes the largest amount of columns?

3.- Now use the -gt 0.5 command but add a conservation score of different values: 30, 50 and 80 (-cons option). Again make sure that your output alignment is in phylip format. Which effect does it have on the trimmed alignment?

4.- Now instead of using the gap threshold, we'll be using the similarity threshold (-st). Repeat the trimming of the original alignment using different similarity thresholds (0.1, 0.5 and 0.9). Again, how does the alignment trimming vary? Which approach is more aggressive? How can you make sure you don't lose all the alignment?

6.- Now we are going to use the automated trimming methods. Trim your alignment using:

- Use the different automated trimming methods: -gappyout, -strict, -strictplus, -automated1
- Use the more radical methods to delete all the columns with gaps in your alignment: -nogaps

Of all the trimming strategies you've tried, which is the best one? Can you know?

## BMGE (Block mapping and gathering through entropy)

BMGE is a trimming method that bases its trimming on the calculation of the entropy generated from moving between the different states found in each column in an alignment. It compares these entropy value to standard substitution matrices to see whether the entropy values have biological meaning. For each column of the alignment a score is calculated. BMGE then removes the blocks with high entropy values (poorly conserved regions).

BMGE is a command line program that runs on Java, so you'll have to have Java installed for it to run.

The easy way to run BMGE is simply calling:

```
BMGE -i EOG092D2PES.alg -t AA -o EOG092D2PES.BMGE.alg
```

This will result in a trimmed alignment

You can obtain a visual output by using the option -oh

```
BMGE -i EOG092D2PES.alg -t AA -oh EOG092D2PES.BMGE.html
```

BMGE also implements other trimming methods such as gap based:

```
BMGE -i EOG092D2PES.alg -t AA -h 1 -w 1 -g 0.1 -o  
EOG092D2PES.gapBased.phy
```

Where -h makes sure there is no entropy trimming and -w says that the sliding window should be of one so that all columns are considered for trimming.

## Prequal

Another kind of trimming tools are those that try to address errors caused by miss-assemblies or gene prediction errors.

Different to the previous programs they work on the multi-sequence fasta and they try to identify non-homologous regions within the sequences included in the multi-fasta. These non-homologous regions are then masked or altogether removed from the sequence prior to the alignment.

How to run Prequal:

A.- Note that this time the input is an un-aligned multi fasta file.

```
prequal inputFile
```

The result will be a filtered fasta file where some parts of the alignments have been deleted and some others are masked (see stretches of X)

# Alignment and trimming challenge

Lisa Pokorny & Marina Marcet-Houben

In phylogenomics we will not work with one single set of sequences that have to be aligned and trimmed. But rather with a large set. So, how can we work with them? Which programs can we use? How can we adjust the parameters to create good alignments?

## Untangling the early diversification of eukaryotes: a phylogenomic study of the evolutionary origins of Centrohelida, Haptophyta and Cryptista

Article (PDF Available) in [Proceedings of the Royal Society B: Biological Sciences](#) 283(1823):20152802 · January 2016 with 730 Reads  
DOI: [10.1098/rspb.2015.2802](https://doi.org/10.1098/rspb.2015.2802)

<https://datadryad.org/resource/doi:10.5061/dryad.rj87v>