

The RAxML v8.2.X Manual

by Alexandros Stamatakis
Heidelberg Institute for Theoretical Studies
July 2012

Structure of this manual

- . About RAxML
- . (etting Hel)
- . RAxML * eb+servers and (,
- . o/nloading RAxML
- . Oom)iling RAxML
- . RAxML Likelihood -alues 1 diosyn"rasies
- . Alignment in)ut 2ile 2ormats
- . The RAxML o) tions
- X. 3ut)ut 2iles
- X. Oom)uting TO and O values
- X. Sim)le RAxML Analyses
- X. A Sim)le Heterota"mous Model
- X. 2re4uently Asked 5uestions

I. About RAxML

RAxML (Randomized Accelerated Maximum Likelihood) is a program for sequential and parallel Maximum Likelihood based inference of large phylogenetic trees. It can also be used for post-analysis of sets of phylogenetic trees, analyses of alignments and evolutionary placement of short reads.

It has originally been derived from Fast. RAxML /hi"h in turn /as derived from Joe Felsenstein's dnaml /hi"h is)art of the ; H<L ;)a" kage.

When using RAxML please cite:

A. Stamatakis= >RAxML -ersion 8= A tool for ; hylogeneti" Analysis and ; ost+Analysis of Large ; hylogenies>. n *Bioinformatics* 2012; 28(12):2362-2368. doi:10.1093/bioinformatics/bts344. link=<http://bioinformatics.oxfordjournals.org/content/early/2012/12/12/bioinformatics.bts344.full>

II. Getting help

RAxML support is provided via the RAxML (oogle grou) at= <https://groups.google.com/forum/#!topic/RAxML> Note that (oogle grou)s have a sear" h !un"tion

Thus% before)osting to the RAxML google grou)=

1. Sear" h the grou) to see if your issue has not already been discussed
2. If you don't /ant to get a rude reply% read this manual first
3. Read a standard textbook about phylogeneti"s su" h as Jiheng <angle's excellent *Computational Molecular Evolution*. If you haven't read one you should rather not be using RAxML.

Do never send emails with RAxML questions to A. Stamatakis directly!

A step by step tutorial with some basic commands is available at <http://www.o.hits.org/exelixis/ebsoft/are@raxml@handsKon.html>

Additional helpful links and documentation is available at the RAxML software page <http://www.o.hits.org/exelixis/ebsoft/are@raxml@index.html>

III. RAxML web-servers and Graphical User Interfaces

* While there exist several web-servers that allow you to run RAxML, we are directly involved in running three of them.

1. The Online server <http://www.phylo.org/subsections/online>
2. The web-server at Vital IT in Switzerland <http://embnet.vital-it.ch/raxml+bb>
3. A dedicated server for the Evolutionary Analysis Algorithm <http://www.a.hits.org/raxml>

There is no official graphical user interface supported by me, but a GUI has been developed by researchers at the research museum in Frankfurt which is available here <http://sourceforge.net/projects/raxmlgui>

Note that we will not provide any sort of support for the GUI, as you need to contact the original authors for this.

IV. Downloading RAxML

RAxML is open source under GPL. It is distributed via Alexis' github repository <https://github.com/stamatak/standard-RAxML> where you can always download the most up to date version. Make sure to **watch** the github repository to remain up to date regarding code changes.

* We do not provide any support whatsoever for previous versions of the code.

Version numbers follow the notation **x.y.z** where **x** changes with major code reorganizations, **y** changes when new features are added and **z** changes with bug fixes.

V. Compiling RAxML

RAxML comes in a lot of different flavors.

It can run sequentially in parallel using either MPI (message passing interface) or pthreads for multi-core shared memory systems.

It also has a hybrid combined pthreads / MPI parallelization that uses MPI to distribute bootstrap replicates or independent tree searches to different shared memory nodes in a cluster while it uses pthreads to parallelize the likelihood calculations of single tree searches. * We call this coarse grain MPI and fine grain pthreads parallelism.

Thus before compiling you need to know on what kind of system you intend to execute RAxML. Also note that the MPI version only implements a subset of the RAxML functionality, it can only distribute different tree searches to processors.

Another important thing to consider prior to compiling is what your target processor architecture is. Modern x86 processors are very powerful because they have so-called vector instructions.

. ending on ho/ ne/ your)ro"essor is it /ill su))ort SSEB ve"tor instru"tions% or% i! ne/er also the faster A-X or the even faster A-X2 ve"tor instru"tions. These instru"tions are used by RAxML to substantially a""elerate the likelihood and)arsimony "om)utations it "ondu"ts.

Thus% you should al/ays try to "om)ile the "ode in a /ay that best ex)loits the "a)abilities o! your O; , 6s8. Note that% even most modern la)to)s have more than & O; , @"ore% hen"e you /ill)robably al/ays /ant to "om)ile the ; Threads version o! RAxML.

9o/ let's have a look at the RAxML sour"e "ode% /hen you do/nload it it /ill be in a file "alled=

```
standard-RAxML-8.0.0.tar.gz
```

un"om)ress it by ty)ing=

```
gunzip standard-RAxML-8.0.0.tar.gz
tar xf standard-RAxML-8.0.0.tar
```

and "hange into the dire"tory that "ontains the sour"e files and list the "ontents=

```
cd standard-RAxML-8.0.0/
ls
```

There is a subdire"tory "alled usefulScripts that "ontains a "ou)le o!)erl s"ri)ts !or various little RAxML tasks.

9ext% let's list the Makefiles=

```
ls Makefile.*
```

/hi" h /ill generate a listing looking like this=

Makefiles !or se4quential version% hybrid M; @; threads version% se4quential version !or MAOs using the "lang "om)iler% M; version% ; Threads version% ; Threads version !or MAOs using "lang that all rely on the most re"ent A-X2 ve"tor instru"tions=

```
Makefile.AVX2.gcc
Makefile.AVX2.HYBRID.gcc
Makefile.AVX2.mac
Makefile.AVX2.MPI.gcc
Makefile.AVX2.PTHREADS.gcc
Makefile.AVX2.PTHREADS.mac
```

Oorres)onding Makefiles using A-X instru"tions=

```
Makefile.AVX.gcc
Makefile.AVX.HYBRID.gcc
Makefile.AVX.mac
Makefile.AVX.MPI.gcc
Makefile.AVX.PTHREADS.gcc
Makefile.AVX.PTHREADS.mac
```

Oorres)onding Makefiles not using any ve"tor instru"tions 6only re4quired /hen your O; , is more than P+ years old6. The file "alled= Makefile.QuartetMPI.gcc is a dedi"ate Makefile that im)lements a M;)aralleli7ation o! the Quartet evaluation !un"tionality% !or details see the se"tion des"ribing the "ommand line arguments6

)er!orman"e de"reaseG

The M; version is for executing really large)rodu"tion runs 6i.e. &\$\$ or &\$\$\$\$ bootstra)s8 on a Linux "luster. <ou "an also)er!orm multi)le in!eren"es on larger datasets in)arallel to lnd a best+ kno/n ML tree !or your dataset.

2inally% the ra)id RS algorithm and the asso"iated ML sear"h have also been)aralleli7ed /ith M; .

Warning: Redu"ed !un"tionality o! M; versionG The "urrent M; version only /orks)ro)erly i! you s)e"ily the -# or -N o)tion in the "ommand line% sin"e it has been designed to do multi)le in!eren"es or ra)id@standard RS 6bootstra)8 sear"hes in)arallelG

2or all remaining o)tions% the usage o! this ty)e o! "oarse+grained)arallelism does not make mu"h senseG

Processor Affinity and Thread Pinning with the PThreads Version

An im)ortant as)e"t i! you /ant to use the ; Threads version o! the)rogram is to lnd out ho/ your o)erating system@lat!orm handles)ro"essor aLnity o! threads. * ithin the shared+memory or multi+core "ontext)ro"essor aLnity means that i! you run!or instan"e? threads on a ?+/ay O; , or ? "ores ea"h individual thread should al/ays run on the same O; , % that is% thread0 on CPU0% thread1 on CPU1 et".

This is im)ortant !or eL "ien"y% sin"e "a"he entries "an be "ontinuously re+used i! a thread% /hi" h /orks on the same)art o! the shared memory s)a"e% remains on the same O; , . ! threads are moved around% !or instan"e% thread0 is initially exe"uted on CPU0 but then on CPU4 et". the "a"he memory o! the O; , /ill have to be re+lled every time a thread is moved. * ith)ro"essor aLnity enabled%)er!orman"e im)rovements o! PS have been measured on suL "iently large and thus memory+intensive datasets.

There is a !un"tion that /ill automati"ally)in threads to O; , s% that is% enlor"e thread aLnity% under L9, X@, 9 X. Re"ause this !un"tion might o""asionally "ause some error messages during "om)ilation due to)ortability issues it is disabled by de!ault. To enable it% you /ill need to "omment out this Nag=

```
#define _PORTABLE_PTHREADS
```

in the axml.c sour"e l le.

How many Threads shall I use?

t is im)ortant to kno/ that the)arallel eL "ien"y o! the ; Threads version o! RAxML de)ends on the alignment length. 9ormally% you /ould ex)e"t a)arallel)rogram to be"ome !aster as you in"rease the number o! "ores@)ro"essors you are using. This is ho/ever not generally true% be"ause the more)ro"essors you use% the more a""umulated time they send /aiting !or the in)ut to be)arsed and "ommuni"ating /ith ea"h other. n "om)uter s"ien"e this)henomenon is kno/ as Amdahl's la/ 6see <http://en.iki.edu.org/wiki/AmdahlsKla/>8.

Thus% i! you run RAxML /ith B2 instead o! & thread this does not mean that it /ill automati"ally be"ome B2 times !aster% it may a"tually even be"ome slo/er. As already mentioned% the)arallel eL "ien"y% that is% /ith ho/ many threads@ores you "an still exe"ute it eL "iently in)arallel de)ends on the alignment length% or to be more)re"ise on the number o! distin"t)atterns in your alignment. This number is)rinted by RAxML to the terminal and into the RAxML_info.runID l le and look like this=

Alignment has 70 distinct alignment patterns

As a rule o! thumb lld use one "ore@thread)er P\$\$. 9A site)atterns% i.e.% i! you have less% than it's)robably better to lust use the se4quential version. Single+gene . 9A alignments /ith around &\$\$\$ sites "an be analy7ed /ith 2 or at most ? threads. Thus% the more)atterns your alignment has% the

more threads@ores you can use eL "iently.

Also note that eL "ien"y varies depending on the type of data or more precisely the number of states in your data e.g. ? in . 9A% 2\$ for proteins. The more states you have the fewer site patterns you need per thread@ore for RAxML to execute eL "iently in parallel. This is because there is more combinatorial work for more mathematical operations to be done per site pattern as the number of states increases. * with protein data you thus require less sites per thread for RAxML to run eL "iently. Thus an MSA with &\$\$\$ protein site patterns may still run eL "iently when using &' "ores@threads.

Finally parallel eL "ien"y also depends on the rate heterogeneity model. * with the (AMMA model that entails more combinations you can thus typically use more threads than with the OAT model that only executes approximately 1/4 of the combinations the (AMMA model requires.

Note that these are just very rough rules of thumb you need to test what the optimal setting is for your dataset

VI. RAxML Likelihood Values & Idiosyncrasies

It is very important to note that the likelihood values produced by RAxML can not be directly compared to likelihood values of other ML programs. However the likelihood values of the current version are very similar to those obtained by other programs with respect to previous releases of RAxML (usually between 0.1 & 0.2 log likelihood units of those obtained e.g. by PHML or GARL).

The above of course refers to evaluating the likelihood on identical trees in terms of tree searches the programs will yield different tree topologies and hence different likelihoods in most cases anyway.

Note that the deviations between PHML@RAxML and GARL likelihood values can sometimes be larger because GARL uses a slightly different procedure to compute empirical base frequencies. Erik Jönvall's personal communication many years ago while the method in RAxML is exactly the same as implemented in PHML.

These deviations between RAxML, PHML on the one side and GARL on the other side appear to be larger on long multi-gene alignments. Also note that even likelihood values obtained by different RAxML versions especially should not be directly compared with each other either.

This is due to frequent code and data structure changes in the likelihood function implementation and model parameter optimization procedures.

Thus if you want to compare topologies obtained by different ML programs with respect to their likelihood make sure that you optimize branch lengths and model parameters of identical topologies with one and the same program.

This can be done by either using the respective RAxML option -f e8 or e.g. the corresponding option in PHML (see <http://www.atg-montpellier.fr/hymle>).

Differences in Likelihood scores

In theory all ML programs implement the same mathematical function and should thus yield the same likelihood score for a fixed model and a given tree topology.

However if you try to implement a numerical function on a finite machine you will unavoidably obtain rounding errors. Even if you change the sequence for if it is changed by the compiler/high it usually is of some operations applied to floating point or double precision arithmetic in our computer you will probably get different results.

In my experiments I have observed differences among identical likelihood values between GARL 5.99%; PHML; RAxML (every program showed a different value).

<ou "an also ex)eriment /ith this by removing the g"" o)timi7ation Nag -02 in one o! the RAxML Makel les. This /ill yield mu"h slo/er "ode% that is in theory mathemati"ally e4uivalent to the o)timi7ed "ode% but /ill yield slightly diMerent likelihood s"ores% due to re+ordered Noating)oint o)erations.

My)ersonal o)inion is that the to)ologi"al sear" h6number o! to)ologies analy7ed8 is mu"h more im)ortant than exa"t likelihood s"ores to obtain *good* l nal ML trees.

Es)e"ially on large trees /ith more than &\$\$\$ se4uen"es the diMerent"es in likelihood s"ores indu"ed by the to)ology are usually so large% that a very rough)arameter o)timi7ation /ith an e)silon 6RAxML -e o)tion8 o! & log likelihood unit% i.e.% i! the diMerent"e e)silon bet/een t/o su""essive model)arameter o)timi7ation iterations is smaller than &.\$ /e sto) the o)timi7ation% /ill already "learly sho/ the diMerent"es.

9ote that% i! you)er!orm a bootstra) analysis you don:t need to /orry too mu"h about likelihood values any/ay% sin"e usually you are only interested in the bootstra)ed to)ologies.

The CAT model of rate heterogeneity

The name o! this model has "aused a lot o! "on!usion be"ause there is a OAT model also im)lemented in ; hylorayes 6see [http://megasun.b" h.umontreal."a@.eo\)le@artillot@//@index.htm](http://megasun.b)

, n!ortunately% /as not a/are o! this /hen introdu"ed the OAT model in RAxML% the OAT model in RAxML is something "om)letely diMerent6 Ho/ever% de"ided not to "hange the name !or ba"k/ard "om)atibility su" h that ne/ RAxML version kee) /orking /ith old /ra))er s"ri)ts.

Warning: t is not a good idea to use the OAT a))roximation o! rate heterogeneity on datasets /ith less than P\$ taxa. n general there /ill not be enough data)er alignment "olumn available to reliably estimate the)er-site rate)arameters.

OAT has been designed to a""elerate the "om)utations on large datasets /ith many taxa6 ; lease read the res)e"tive)a)er

[http://ieeex\)lore.ieee.org@x\)l@login.lS\)Ct\)D1arnumberD&'BQBPB1urlDhtt\)SBAS22S22ieeex\)lore.ieee.orgS22x\)lS22absKall.lS\)SB2arnumberSB.&'BQBPB](http://ieeex)lore.ieee.org@x)l@login.lS)Ct)D1arnumberD&'BQBPB1urlDhtt)SBAS22S22ieeex)lore.ieee.orgS22x)lS22absKall.lS)SB2arnumberSB.&'BQBPB)

to understand ho/ OAT /orks% /hat the rate "ategories are 6they are "on"e)tu ally diMerent !rom the dis"rete rate "ategories o! the (AMMA model8 and /hat the limitations o! this method are.

The GTRCAT a))roximation is a "om)utational /orkWaround !or the /idely used (eneral Time Reversible model o! nu"leotide substitution under the (amma model o! rate heterogeneity. OAT servers the analogous)ur)ose% that is% to a""ommodate sear"hes that in"orate rate heterogeneity.

The a!orementioned)a)er des"ribes /hat GTRCAT is and /hy don:t like GTRGAMMA des)ite the !a"t that it is a beauti!ul (reek letter.

The main idea behind GTRCAT is to allo/ !or integration o! rate heterogeneity into)ylogeneti" analyses at a signi!antly lo/er "om)utational "ost 6about ? times !aster8 and memory "onsum)tion 6? times lo/er8.

Essentially% GTRCAT re)resents a rather un+mathemati"al *quick & dirty* a))roa" h to ra)idly navigate into)ortions o! the tree s)a"e% /here the trees s"ore /ell under GTRGAMMA.

Ho/ever% due to the /ay individual rates are o)timi7ed and assigned to rate "ategories in GTRCAT% the likelihood values "om)uted under GTRCAT are "om)letely meaningless.

Warning: never "om)are alternative tree to)ologies using their OAT+based likelihood s"ores6

<ou /ill)robably obtain a biased assessment o! trees. This is the reason /hy GTRCAT is "alled

approximation instead of model. The same applies to the OAT approximation when used with AA data or any other data type that RAxML supports.

In general, the OAT approximation of rate heterogeneity works very well on datasets with more than 10 taxa. For "conducting tree searches under a model that accommodates rate heterogeneity among sites. In other words, if you score the trees obtained under OAT using (AMMA you will usually obtain equally good trees as with full searches under (AMMA at a substantially lower computational cost. Also, (AMMA may not work for numerical reasons on very large trees with more than 10,000 taxa (see <http://www.biomedcentral.com/submit?X=2&P=2&X=8>).

Another detailed study of a OAT-like model is conducted in the FastTree paper, see <http://www.losone.org/article.php?id=doi:10.1186/1471-2148-5-147>

VII. Alignment input File Formats

Alignment & Tree Input Formats

Alignments: The input alignment format of RAxML is relaxed interleaved or sequential ;H<L ; or 2ASTA. *Relaxed* means that sequence names can be of variable length between 1 to 255 characters.

If you need longer taxon names you can add the constant `#define nmlength 256` in the source file `raxml.h` appropriately.

Moreover, RAxML is less sensitive with respect to the ;H<L ; formatting (tabs, insets, etc.) of interleaved ;H<L ; files.

Trees: The input tree format is Newick. The RAxML input trees must not always be comprehensive, i.e., need not contain all taxa of the alignment.

See <http://www.evolution.genetics.washington.edu/hyli/newicktree.html> for details on the Newick format.

Alignment Error Checking

Many alignments need to be checked for the following errors/insufficiencies before running an analysis with RAxML or any other phylogenetic inference program.

RAxML automatically analyzes the alignment and checks for the following errors=

1. identical sequence names appearing multiple times in an alignment, this can easily happen when you export a standard ;H<L ; file from some tool which truncates the sequence names to 8 or 15 characters.
2. identical sequences that have different names but are exactly identical. This mostly happens when you excluded some hard-to-align alignment regions from your alignment and does not make sense to use.
3. undetermined columns that contain only ambiguous characters that will be treated as missing data, i.e., columns that entirely consist of X, ?, *, - for AA data and N, O, x, ?, - for tRNA data (analogous for other data types)
4. undetermined sequences that contain only ambiguous characters (see above) that will be treated as missing data.

Prohibited characters in taxon names are names that contain any form of whitespaces like blanks, tabulators, and carriage returns as well as one of the following prohibited characters: : or () or [] .

n "ase that RAxML de"ts identi"al se4uen"es and/or undetermined "olumns and /as exe"uted% e.g.% /ith -n alignmentName it /ill automati"ally /rite an alignment file "alled alignmentName.reduced /ith identi"al se4uen"es and/or undetermined "olumns removed.

! this is de"ted for a)artitioned model analysis a res)e"tive model file modelFileName.reduced /ill also be /ritten. n "ase RAxML en"ounters identi"al se4uen"e names or undetermined se4uen"es or illegal "hara"ters in taxon names it /ill exit /ith an error and you /ill have to l x your alignment.

VIII. The RAxML options

The single by !ar most im)ortant "ommand is the RAxML hel) o)tion that dis)lays all o)tions.

also !re4uently use it be"ause "an not remember them all. n the !ollo/ing /ill assume that /e are !ust using the se4uential unwe"tori7ed "ode% that is% the raxmlH; O exe"utable.

/ill dis"uss ea"h o)tion and)rovide a sim)le usage exam)le for it.

Thus% to get on+line hel) ty)e=

```
raxmlHPC -h
```

and you /ill get the !ollo/ing% very long listing% that /ill be dis"ussed at length belo/=

```
raxmlHPC    -s sequenceFileName -n outputFileName -m substitutionModel
             [-a weightFileName] [-A secondaryStructureSubstModel]
             [-b bootstrapRandomNumberSeed] [-B wcCriterionThreshold]
             [-c numberOfCategories] [-C] [-d] [-D]
             [-e likelihoodEpsilon] [-E excludeFileName]
             [-f a|A|b|B|c|C|d|D|e|E|F|G|g|H|h|i|I|j|J|k|m|n|N|o|p|q|r|R|s|S|t|T|u|
             v|V|w|W|x|y]
             [-F]
             [-g groupingFileName] [-G placementThreshold] [-h] [-H]
             [-i initialRearrangementSetting]
             [-I autoFC|autoMR|autoMRE|autoMRE_IGN]
             [-j] [-J MR|MR_DROP|MRE|STRICT|STRICT_DROP|T_<PERCENT>] [-k] [-K]
             [-L MR|MRE|T_<PERCENT>] [-M]
             [-o outGroupName1[,outGroupName2[,...]]] [-O]
             [-p parsimonyRandomSeed] [-P proteinModel]
             [-q multipleModelFileName] [-r binaryConstraintTree]
             [-R binaryModelParamFile] [-S secondaryStructureFile]
             [-t userStartingTree]
             [-T numberOfThreads] [-u] [-U] [-v] [-V] [-w outputDirectory]
             [-W slidingWindowSize]
             [-x rapidBootstrapRandomNumberSeed] [-X] [-y]
             [-Y quartetGroupingFileName|ancestralSequenceCandidatesFileName]
             [-z multipleTreesFile]
             [-#|-N numberOfRuns|autoFC|autoMR|autoMRE|autoMRE_IGN]
             [--mesquite][--silent][--no-seq-check][--no-bfgs]
             [--asc-corr=stamatakis|felsenstein|lewis]
             [--flag-check]
             [--auto-prot=ml|bic|aic|aicc ]
             [--epa-keep-placements=number]
             [--epa-accumulated-threshold=threshold]
             [--epa-prob-threshold=threshold]
```

```

[--JC69][--K80]
[--set-thread-affinity]
[--bootstop-perms=number]
[--quartets-without-replacement]
[--print-identical-sequences]

```

- a Specify a column weight file name to assign individual weights to each column of the alignment. Those weights must be integers separated by any type and number of white-spaces within a separate file.

In addition, there must be as many weights as there are columns in your alignment. The contents of an example file would look like this=

```

5 1 1 2 1 1 1 1 1 1 2 1 1 3 1 1 1 1 1 1 1 1
1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 4 1 1 1 4 1 1

```

Example: `raxmlHPC -a wgtFile -s alg.phy -p 12345 -m GTRCAT -n TEST`

- A Specify one of the secondary structure substitution models implemented in RAxML. The same nomenclature as in the PHASE manual is used, available models: S6A, S6B, S6C, S6D, S6E, S7A, S7B, S7C, S7D, S7E, S7F, S16, S16A, S16B

DEFAULT: 16-state GTR model (S16)

Note that partitioning does not work with secondary structure models. That is a secondary structure can only be superimposed to a single partition. Also note that you need to also specify a file defining the secondary structure via the `-S` option.

Example: `raxmlHPC -S secondaryStructureFile -s alg.phy -A S7D -p 12345 -m GTRGAMMA -n TEST`

- b Specify an integer number (random seed) and turn on bootstrapping

DEFAULT: OFF

This option allows you to turn on non-parametric bootstrapping. To allow for reproducibility of runs in the sequential program you have to specify a random number seed e.g. `-b 123476`.

Note however that parallel bootstraps with the parallel version `raxmlHPC+M` are not reproducible despite the fact that you specify a random number seed.

Example: `raxmlHPC -b 12345 -p 12345 -# 100 -s alg -m GTRCAT -n TEST`

- B specify a floating point number between 0.0 and 1.0 that will be used as cutoff threshold for the MR-based bootstopping criteria. The recommended setting is 0.03.

DEFAULT: 0.03 (recommended empirically determined setting)

This setting allows to specify a threshold for the so-called bootstopping criteria that will automatically determine if you have conducted enough bootstrap replicate searches for obtaining stable support values. Note that this only has an effect if you use the bootstopping criteria that rely on building majority rule consensus trees for determining convergence. The option will not have an effect when the recommended "frequency" bootstopping criterion is being used. Please also read and cite the corresponding paper=

<http://link.springer.com/handle/1007/1234567890>

Example: `raxmlHPC -B 0.02 -b 12345 -p 12345 -# AUTOMR -s alg -m GTRCAT -n TEST`

- c Specify number of distinct rate categories for RAXML when model of rate heterogeneity is set to CAT. Individual per-site rates are categorized into numberOfCategories rate categories to accelerate computations.

DEFAULT: 25

Example: `raxmlHPC -c 40 -p 12345 -s alg -m GTRCAT -n TEST`

Warning: Note that the setting of -c has no effect whatsoever on the number of discrete rate categories that are used to approximate the Gamma distribution of rate heterogeneity! RAXML always uses 4 discrete rate categories to approximate Gamma!

- C Enable verbose output for the "-L" and "-f i" options. This will produce more, as well as more verbose output files

DEFAULT: OFF

The above option will generate verbose output and additional output files when RAXML is used to compute the TO and O measures as introduced by Saliotis and Rokas=<http://www.nature.com/nature/journal/v407/n7457/full/nature04288.html>. The method is described in more detail in the following paper: Saliotis and Rokas=[http://embio.oxfordjournals.org/content/early/2006/04/20/embio.msu.04288.abstract?keyref=1111keyD%20P2u\(9x\\$H7R23/](http://embio.oxfordjournals.org/content/early/2006/04/20/embio.msu.04288.abstract?keyref=1111keyD%20P2u(9x$H7R23/)

Example: `raxmlHPC -L -m GTRCAT -L MRE -z treeSet -n TEST`

- d start ML optimization from random starting tree

DEFAULT: OFF

This option allows you to start the RAXML search with a complete random starting tree instead of the default randomized stepwise addition Maximum parsimony starting tree. On smaller datasets (around 100 taxa) it has been observed that this might sometimes yield topologies of distinct local likelihood maxima which better correspond to empirical expectations.

It has also been observed that this sometimes yield better with more diverse with starting trees for the analysis of broad phylogenomic alignments that have a very strong phylogenetic signal.

Example: `raxmlHPC -d -p 12345 -s alg -m GTRGAMMA -n TEST`

- D ML search convergence criterion. This will break off ML searches if the relative Robinson-Foulds distance between the trees obtained from two consecutive lazy SPR cycles is smaller or equal to 1%. Usage recommended for very large datasets in terms of taxa. On trees with more than 500 taxa this will yield execution time improvements of approximately 50% while yielding only slightly worse trees.

DEFAULT: OFF

* When enabling this option RAXML will store ML and standard deviation bootstrap search results early when the R2 distance between the trees generated by two consecutive "yes" or "no" subtree pruning moves is smaller than 0.5. This leads to substantial speed improvements while the decrease in log likelihood scores is only very small. The option has been tested on several datasets and the results have been included in the following book chapter:

A. Stamatakis: "Hylogenetic Search Algorithms for Maximum Likelihood". In M. Elloumi, A. Jomayya, editors. Algorithms in Computational Biology. Lecture Notes in Computer Science, Springer, 2008.

Example: `raxmlHPC -D -p 12345 -s alg -m GTRCAT -n TEST`

-e set model optimization precision in log likelihood units for final optimization of tree topology

DEFAULT: 0.1 for models not using proportion of invariant sites estimate
0.001 for models using proportion of invariant sites estimate

This allows you to specify a high likelihood difference between the model parameters will be optimized when RAXML uses the GTR model or when you must evaluate a tree with the -f option or a bunch of trees with the -f n option or similar options. This has shown to be useful to quickly evaluate the likelihood of a bunch of large local trees or more than 1000 taxa because it will run much faster.

Typically use e.g. -e 1.0 or -e 2.0 in order to rapidly compare different local tree topologies based on their likelihood values.

Note that the log likelihood differences are typically far larger than 0.5 or 2.0 log likelihood units. The default setting is 0.5 log likelihood units which proves to be sufficient in most practical cases.

Example: `raxmlHPC -f e -e 0.00001 -s alg -t tree -m GTRGAMMA -n TEST`

-E specify an exclude file name, that contains a specification of alignment positions you wish to exclude.

Format is similar to Nexus, the file shall contain entries like "100-200 300-400", to exclude a single column write, e.g., "100-100", if you specify a partition file via "-q", an appropriately adapted partition file will also be written.

This option will make RAXML write a reduced alignment file without the excluded columns that are subsequently used for the analysis you actually want to conduct.

If you use a partitioned model and a proportionately adapted model file will also be written. Note that excluding sites with RAXML is a **two-step** procedure. You will first have to invoke RAXML once using the -E option to generate a H<L formatted file with the desired sites excluded. Then you will have to invoke RAXML again on the alignment file that was written as `alg.excludeFile` in the example below to do an actual tree search for instance.

Example: `raxmlHPC -E excludeFile -s alg -m GTRCAT -q part -n TEST`

In this case the alignment file with columns excluded will be named `alg.excludeFile` and the partition file with the specified columns excluded is named `part.excludeFile`

If you want to exclude sites 1-100 note that the borders that is sites 1-100 and 101 are included and will be removed and sites 101-200 the corresponding exclude file will be

text file would contain the following lines below:

100-199
200-299

-f select algorithm:

The option is a very important one because in many cases it allows you to select what kind of algorithm RAxML shall execute. If you don't specify -f at all RAxML will execute the standard hill climbing algorithm by default. The individual options are listed below.

-f a rapid Bootstrap analysis and search for best-scoring ML tree in one program run

Tell RAxML to conduct a rapid Bootstrap analysis and search for the best-scoring ML tree in one single program run.

Example: `raxmlHPC -f a -p 12345 -s alg -x 12345 -# 100 -m GTRCAT -n TEST`

-f A compute marginal ancestral states on a ROOTED reference tree provided via -t

This option allows you to compute marginal ancestral states on a given fixed and rooted reference tree. If you don't know what marginal ancestral states are please read Jiheng Zhang's book on Computational Molecular Evolution.

Example: `raxmlHPC -f A -t testTree -s testData -m GTRGAMMA -n TEST`

A small properly formatted test alignment would look like this:

```
4 4
t1 ACGT
t2 AATT
t3 ATAT
t4 CCGT
```

and a test tree:

```
((t1,t2),(t3,t4));
```

RAxML will then output the rooted binary tree again but with inner node labels that must be used to associate ancestral sequences with their corresponding positions in the tree. It also writes to file the marginal probabilities for each inner node label as well as one containing guesses by taking the maximum probability for the ancestral sequences.

-f b draw bipartition information on a tree provided with -t (typically the best-known ML tree) based on multiple trees (e.g., from a bootstrap) in a file specified by -z

Example: `raxmlHPC -f b -t ref -z trees -m GTRCAT -n TEST`

-f B optimize branch length scaler and other model parameters (GTR, alpha, etc.) on a tree provided with -t. The input tree passed via -t needs to contain branch

lengths. The branch lengths will not be optimized, just scaled by a single common value.

This is a slightly idiosyncratic option that was mainly developed to work on "partitioning" with the partitionfinder tool (see <http://www.robertlanier.com/partitionfinder>).

The "one" of branch length scaling is similar to the one implemented in MrBayes. That is, we assume that all partitions have the same branch lengths proportionally to each other. However, for each partition p we do a maximum likelihood estimate of a single scaling factor $s[p]$ by dividing the branch lengths. Essentially this represents a parameterization trade-off between estimating a separate set of branch lengths for each partition (heavy increase in the number of parameters) and doing one common branch length estimate across all partitions (possible underparameterization). Note that this works only for evaluating a given fixed tree with branch lengths and not for tree searching.

Example: `raxmlHPC -s alg.phy -t tree -m GTRGAMMA -f B -q part -n TEST`

`-f c` check if the alignment can be properly read by RAXML

Example: `raxmlHPC -f c -m GTRCAT -s alg -n TEST`

`-f C` ancestral sequence test for my PhD student Jiajie Zhang, users will also need to provide a list of taxon names via `-Y` separated by white-spaces

This is an algorithm we tried to implement for post-analyzing viral phylogenies. The idea was to build a test that can determine if certain sequences in the viral tree are truly ancestral that is, if they should essentially be located at or right next to an inner node of the tree.

For this we built the following test algorithm:

(Given a fixed phylogenetic tree (usually the best-known ML tree) for a set of sequences and given a list of candidate ancestral sequences via `-Y` for each of these candidate sequences we compare the likelihood of the fixed tree and the given optimal branch length for that putative ancestral sequence with the likelihood of the fixed tree and an essentially zero branch length for the ancestral candidate.

Then we apply the test to all three branches the candidate ancestral sequence is attached to that is the branch it is directly attached to as well as the two branches to the right and left of the node to which the candidate ancestral sequence is attached to.

The likelihoods of the different branch length durations (the three almost-zero branch lengths) are compared using the Shimodaira-Hasegawa test to the original optimal likelihood. The Shimodaira-Hasegawa test might not be the most modern test available but was sufficient for conducting initial tests. It turned out that the test worked very well on simulated data but not at all for real data. Thus we decided to abandon further testing but nonetheless keep the option in RAXML.

Example: `raxmlHPC -f C -Y candidateTaxonNames -m GTRGAMMA -t tree -s alg -n TEST`

`-f d` new rapid hill-climbing

DEFAULT: ON

This is the default RAXML tree search algorithm and is substantially faster than the original

search algorithm. It takes some shortcuts but yields trees that are almost as good as the ones obtained from the full search algorithm. The algorithm is described and assessed in this paper here: <http://link.springer.com/article/10.1007/s11268-005-9027-2>

Example: `raxmlHPC -f d -m GTRCAT -p 12345 -s alg -n TEST`

-f D execute one or more rapid hill-climbing searches that will also generate RELL bootstraps

This is an implementation of the technique introduced in the following paper: <http://www.ncbi.nlm.nih.gov/pubmed/2678809> Some initial tests have shown that the RAXML implementation yields results that are well correlated with standard bootstrap values. Using this option is recommended when you have very large trees and don't have the resources/time for computing bootstrap estimates.

Example: `raxmlHPC -f D -m GTRCAT -p 12345 -s alg -n TEST`

-f e optimize model parameters+branch lengths for given input tree
You need to provide an input tree via the `-t` option.

Example: `raxmlHPC -f e -t ref -m GTRGAMMA -s alg -n TEST`

-f E execute very fast experimental tree search, at present only for testing

This option will execute a very fast tree search algorithm that will not try as hard to optimize the likelihood. It is intended for very large trees and follows a similar logic as the FastTree program: <http://meta.mi.robersonline.org/fasttree/>

Example: `raxmlHPC -f E -m GTRCAT -p 12345 -s alg -n TEST`

-f F execute fast experimental tree search, at present only for testing

This option will execute a fast tree search algorithm that will not try as hard to optimize the likelihood. It is intended for very large trees and follows a similar logic as the FastTree program: <http://meta.mi.robersonline.org/fasttree/>

Example: `raxmlHPC -f F -m GTRCAT -p 12345 -s alg -n TEST`

-f g compute per site log Likelihoods for one or more trees passed via `-z` and write them to a file in treepuzzle format that can be read by CONSEL.
The model parameters will be estimated on the first tree only!

Example: `raxmlHPC -f g -s alg -m GTRGAMMA -z trees -n TEST`

-f G compute per site log Likelihoods for one or more trees passed via `-z` and write them to a file in treepuzzle format that can be read by CONSEL.
The model parameters will be re-estimated for each tree!

Example: `raxmlHPC -f G -s alg -m GTRGAMMA -z trees -n TEST`

-f h compute log likelihood test (SH-test) between best tree passed via `-t` and a bunch of other trees passed via `-z`
The model parameters will be estimated on the first tree only!

Example: `raxmlHPC -f h -t ref -z trees -s alg -m GTRGAMMA -n TEST`

-f H compute log likelihood test (SH-test) between best tree passed via -t and a bunch of other trees passed via -z
The model parameters will be re-estimated for each tree !

Example: `raxmlHPC -f H -t ref -z trees -s alg -m GTRGAMMA -n TEST`

-f i calculate IC and TC scores (Salichos and Rokas 2013 <http://www.nature.com/nature/journal/v425/n7203/full/nature06899.html>) on a reference tree provided with -t based on multiple trees (e.g., from a bootstrap) in a file specified by -z

The method is described in more detail in the following paper: Salichos and Rokas=

<http://embe.oxfordjournals.org/content/early/2013/05/22/oxfordjournals.molbev.a0052231.abstract>

The method is then extended here <http://dx.doi.org/10.1093/molbev/msu011> for partial gene trees.

Warning: The default topology calculations are not done exactly as described in our MRE paper. To have RAXML do exactly what is described in the paper, please edit the file `bipartitionList.c` by commenting out or removing the line
`#define BIP_FILTER`

Example: `raxmlHPC -f i -m GTRCAT -t referenceTree -z bootstrapTrees -n TEST`

-f I a simple tree rooting algorithm for unrooted trees. It roots the tree by rooting it at the branch that best balances the subtree lengths (sum over branches in the subtrees) of the left and right subtree. A branch with an optimal balance does not always exist! You need to specify the tree you want to root via -t.

Example: `raxmlHPC -f I -m GTRCAT -t unrootedTree -n TEST`

-f j generate a bunch of bootstrapped alignment files from an original alignment file. You need to specify a seed with -b and the number of replicates with -#

Example: `raxmlHPC -f j -b 12345 -# 100 -s alg -m GTRCAT -n TEST`

Warning: Note that if you are generating those RS replicate alignments for a dataset for which you subsequently intend to conduct a partitioned analysis you will also need to specify the partition file here via -q because RAXML re-samples sites on a per-partition basis

-f J Compute SH-like support values on a given tree passed via -t.

This option will "omute sh+like su))ort values as described here <http://www.ncbi.nlm.nih.gov/pubmed/2522088> on a given tree. The input tree is typically the best-known ML tree found by a RAxML analysis. Zee) in mind that for applying the SH-like test the tree needs to be 99 nearest neighbor interchange optimal. Thus RAxML will initially try to apply 99 moves to further improve the tree and then "omute the SH test for each inner branch of the tree.

Example: `raxmlHPC -f J -p 12345 -m GTRGAMMA -s alg -t tree -n TEST`

-f k Fix long branch lengths in partitioned data sets with missing data using the branch length stealing algorithm.
This option only works in conjunction with "-t", "-M", and "-q".
It will print out a tree with shorter branch lengths, but having the same likelihood score.

This option tries to fix the issue of very long branch lengths in partitioned datasets with missing data. For each partition and each branch of a given tree it checks if there is data available for that partition on both sides of the tree as defined by the branch. If this is not the case this means that there is only missing data on one side. In such a case we usually have a very long branch length. The algorithm fixes this issue by stealing branch lengths from those partitions that have data on both sides of the branch under consideration. If there are several other partitions that have data it computes a weighted average for the stolen branch length based on the site counts in each partition from which it stole a branch length. The nice property of this algorithm is that by changing the branch lengths in this way the likelihood of the tree is not changed.

Example: `raxmlHPC -f k -m GTRGAMMA -s alg -q part -M -t tree -n TEST`

A new tree file containing a shorter tree with stolen branch lengths will be written to a file named `RaxML_stolenBranchLengths.TEST`

-f m compare bipartitions between two bunches of trees passed via "-t" and "-z" respectively. This will return the Pearson correlation between all bipartitions found in the two tree files. A file called `RAxML_bipartitionFrequencies.outputFileName` will be printed that contains the pair-wise bipartition frequencies of the two sets

Example: `raxmlHPC -f m -t trees1 -z trees2 -s alg -m GTRCAT -n TEST`

-f n compute the log likelihood score of all trees contained in a tree file provided by -z
The model parameters will be estimated on the first tree only!

Example: `raxmlHPC -f n -z trees -s alg -m GTRGAMMA -n TEST`

-f N compute the log likelihood score of all trees contained in a tree file provided by -z
The model parameters will be re-estimated for each tree!

Example: `raxmlHPC -f N -z trees -s alg -m GTRGAMMA -n TEST`

-f o old and slower rapid hill-climbing without the heuristic cutoff described in <http://link.springer.com/article/10.1007/s00438-006-0122-2>

If you use this option you will typically get slightly better likelihood scores while the run times are expected to increase by a factor of 2 to 3.

Example: `raxmlHPC -f o -p 12345 -m GTRCAT -s alg -n TEST`

-f p perform pure stepwise MP addition of new sequences to an incomplete starting tree and exit

Example: `raxmlHPC -f p -t ref -p 12345 -s alg -m GTRCAT -n TEST`

-f q fast quartet calculator

This option will calculate the likelihood scores of all 4 quartets for a given input alignment. Initially RAxML will construct a randomized starting tree or you can also pass a given tree (e.g. the best-known ML tree) to RAxML via `-t`. This tree is used to estimate model parameters that will then remain fixed during the 4 quartet calculations. For each 4 quartet that is evaluated only the branch lengths will be optimized.

The algorithm itself has three flavors: 1) randomly evaluate a set of 4 quartets, 2) evaluate all 4 quartets, or 3) evaluate 4 quartets from a restricted 4 quartet constraint tree (i.e. that must contain exactly 1 monophyletic and multiple branching groups).

For details on the input format for the 4 quartet constraint file see the `-y` option.

To not evaluate all possible 4 quartets (keep in mind that this number grows quickly as a function of the taxa in the alignment) you have to use either `-#` or `-N` to specify how many 4 quartets you want to randomly evaluate.

Example 1: `raxmlHPC -m GTRGAMMA -t tree -s alg -f q -p 12345 -n TEST`

The above example will evaluate all possible 4 quartets.

Example 2: `raxmlHPC -m GTRGAMMA -t tree -s alg -f q -p 12345 -N 100 -n TEST`

The above example will evaluate 100 randomly chosen 4 quartets.

The output is written to a file named `RaxML_quartets.TEST`

It looks like this-

Taxon names and indexes:

```
Cow 1
Carp 2
Chicken 3
Human 4
Loach 5
Mouse 6
Rat 7
Seal 8
Whale 9
Frog 10
```

```

1 2 | 3 4: -2640.417355
1 3 | 2 4: -2640.810748
1 4 | 2 3: -2638.359608
1 2 | 3 5: -2552.181311
1 3 | 2 5: -2543.674965
...

```

First each taxon name is assigned a number then the 4 quartets are represented as bipartitions. For instance 1 2 | 3 4: -2640.417355 is the quartet ((Cow,Carp), (Chicken, Human)) and has a likelihood of -2640.417355.

The Mj version for this serial option that can be compiled with the corresponding Makefiles will distribute quartet evaluations across processors with Mj.

-f r compute pairwise Robinson-Foulds (RF) distances between all pairs of trees in a tree file passed via -z. If the trees have node labels represented as integer support values the program will also compute two flavors of the weighted Robinson-Foulds (WRF) distance

Example: `raxmlHPC -m GTRCAT -z trees -f r -n TEST`

The two flavors of the weighted R2 distance that are computed are=

1. Just add the support of those bipartitions contained in one tree but not the other
2. add the support of those bipartitions contained in one tree but not the other and also add the difference in support for each shared bipartition

The output looks as follows (not including weighted R2 distances)

```
0 1: 8 0.125000
```

Here 0 and 1 denote that this is the distance between tree 0 and 1 in the tree file. The first and second tree in there. Here 8 is the plain R2 distance and 0.12500 is the normalized R2 distance corresponding to 2.PS.

The normalized relative R2 distance is calculated by dividing the plain R2 by $2(n-3)$ where n is the number of taxa.

The value of the normalized R2 distance ranges between 0.0 and 1.0 and is interpreted as the percentage of splits (bipartitions) that are unique to one of the two trees. In the above example 2.PS of the splits are unique to either tree 0 or tree 1.

-f R compute all pairwise Robinson-Foulds (RF) distances between a large reference tree passed via -t and many smaller trees (that must have a subset of the taxa of the large tree) passed via -z. This option is intended for checking the plausibility of very large phylogenies that can not be inspected visually any more.

Example: `raxmlHPC -f R -m GTRCAT -t hugeTree -z manySmallTrees`

-f s split up a multi-gene partitioned alignment into the respective subalignments

Example: `raxmlHPC -f s -q part -s alg -m GTRCAT -n TEST`

Warning: Note that the subalignments will be named by the names you have given to

the individual partitions. If your partition file contains two partitions called part1 and part2 this command will generate two alignment files called part1.phy and part2.phy

`-f S` compute site-specific placement bias using a leave one out test inspired by the evolutionary placement algorithm

The leave-one-out approach for assessing site-specific incongruence with the underlying tree is based on the evolutionary placement algorithm% see <http://sysbio.oxfordjournals.org/content/30/20/short>

Using a given reference tree (typically the best-known ML tree) a leave-one-out test is conducted by running a single taxon t at a time% removing it from the tree and subsequently re-inserting it again into its original position. For each taxon t we use a sliding window of size w in the default case $w=100$ of consecutive sites in the alignment and compute the maximum likelihood placement in the tree with respect to those w sites only.

This is done repeatedly for each alignment site by moving the window over the alignment on a site by site basis.

Once we have obtained the best placements for each sliding window starting position for a taxon t we compute a score $S[i][t]$ for each site.

The score $S[i][t]$ is the mean distance in terms of number of nodes in the tree between the respective best placements for all sliding windows that comprise site i and the original placement of taxon t .

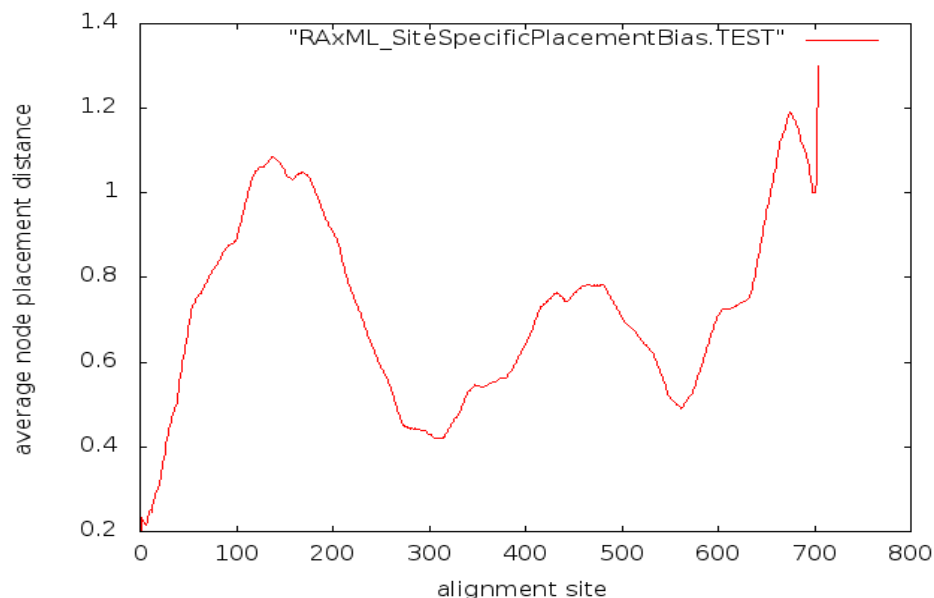
Finally we calculate a global $S[i]$ score for each site by computing the mean of the $S[i][t]$ computed for each taxon t .

Thus RAxML will return an list of n $S[i]$ average node distances% where n is the number of sites in the alignment.

Example: `raxmlHPC -f S -s alg -t tree -m GTRGAMMA -n TEST`

The main idea behind this algorithm is to provide a means for assessing the phylogenetic variability of different areas/ parts of a gene in order to decide which part of the gene (e.g. a & S gene) to amplify and sequence using qPCR.

The option above will generate an output file called `RaxML_SiteSpecificPlacementBias.TEST` which you can then plot using a tool like gnuplot to obtain the following plot=



The)lot indi"ates that the alignment site bet/een)osition 2\$\$ and ' \$\$ /ill yield the most "ongruent signal /ith the tree "al"ulated on the entire alignment. 3verall% the mean node)la"ement distan"es are relatively lo/% meaning that this alignment has a relatively stable)hylogenet"i" signal over its entire length.

-f t do randomized tree searches on one fixed starting tree

This "ommand /ill exe"ute a s)e"il ed number 6by -N or -#8 o! randomi7ed tree sear"hes. The diMer"e to the standard sear"h algorithm is% that% given a starting tree it /ill a))ly S; R 6subtree)runing re+gra!ting8 moves in a randomi7ed order and thus may lend better trees.

Example: raxmlHPC -f t -t -m GTRCAT -s alg -p 12345 -n TEST

-f T do final thorough optimization of ML tree from rapid bootstrap search in stand-alone mode

This o)tion allo/s to do a more thorough tree sear"h that uses less la7y% i.e.% more exhaustive S; R moves% in stand+alone mode. This algorithm is ty)i"ally exe"uted in the very end o! a sear"h done via -f a.

Example: raxmlHPC -p 12345 -m GTRGAMMA -s alg -f T -t tree -n TEST

-f u execute morphological weight calibration using maximum likelihood, this will return a weight vector. you need to provide a morphological alignment and a reference tree via -t

This o)tion /ill determine to /hi"h degree the sites o! a mor)hologi"al alignment are "ongruent /ith a given mole"ular releren"e tree. As out)ut it /ill generate a RAXML /eight l le that reNe"ts the degree o! "ongruen"e and that "an be subse4uently read via the -a o)tion to "ondu"t a more in!ormed evolutionary)la"ement o! !ossils using the evolutionary)la"ement algorithm 6-f v o)tion8. 2or more details% see the "orres)onding)a)ers= Evolutionary)la"ement <http://sysbio.oxfordjournals.org/content/54B/2Q&.short> 2ossil "alibration [http://ieeex.lore.ieee.org/x/login.js?CtD1arnumberDPP8'QBQ1urlDhttSBAS2S22ieeex\)lore.ieee.orgS22x\)lsS22absKall.js\)SB2arnumberSB.PP8'QBQ](http://ieeex.lore.ieee.org/x/login.js?CtD1arnumberDPP8'QBQ1urlDhttSBAS2S22ieeex)lore.ieee.orgS22x)lsS22absKall.js)SB2arnumberSB.PP8'QBQ)

Example: raxmlHPC -f u -m BINGAMMA -t molecularTree -s morphologicalAlignment

-f v classify a bunch of environmental sequences into a reference tree using thorough read insertions you will need to start RAXML with a non-comprehensive reference tree and an alignment containing all sequences (reference + query)

For details on the design and performance of this algorithm see <http://sysbio.oxfordjournals.org/content/320/short>
 RAXML will produce a "ou"le of idiosyncratic output files for the elements but also an output file according to the common file standard defined by Erik Matsen for his pplacer (see <http://matsen.ih.rutgers.edu/la>) program that is similar to the E;A Evolutionary
 ; la"ement Algorithm and myself.
 The common file format is described in this paper here [http://www.losone.org/article.php?id=SBAdoiS22&\\$&BX&S22journal.\)one.\\$B&\\$Q](http://www.losone.org/article.php?id=SBAdoiS22&$&BX&S22journal.)one.$B&$Q)

Example: `raxmlHPC -f v -s alg -t referenceTree -m GTRCAT -n TEST`

If you frequently want to insert different sequences into the same reference tree there is a shortcut to make this more efficient because the model parameters and branch lengths are estimated from scratch for the reference tree each time when you invoke the program. If you want to avoid this you proceed as follows.

(Generate a binary file containing the model parameters for the reference tree

Example 1: `raxmlHPC -f e -m GTRGAMMA -s referenceAlignment -t referenceTree
-n PARAMS`

This will generate a file named `RAXML_binaryModelParameters.PARAMS` which can then be read by the E;A to avoid re-estimating parameters

Example 2: `raxmlHPC -f v -R RAXML_binaryModelParameters.PARAMS
-t RAXML_result.PARAMS -s alg -m GTRGAMMA -n TEST2`

Warning: Note that when using binary model files The models of rate heterogeneity in examples 1 & 2 e.g. GTRGAMMA or GTRCAT must be the same. Secondly when using GTRCAT you must disable pattern compression via the `-H` flag in both runs. Thirdly the reference tree read in via `-t` in example 2 must have branch lengths according to the model specified via `-m`. The branch lengths of the reference tree are not re-estimated by the "all in example 2. It is very easy to make a mistake and read in a tree whose branch lengths were estimated under GTRGAMMA while you are using GTRCAT. All of the above only applies when using binary model files through the `-R` option.

`-f V` classify a bunch of environmental sequences into a reference tree using thorough read insertions you will need to start RAXML with a non-comprehensive reference tree and an alignment containing all sequences (reference + query)

This is an experimental extension of the E;A for including short reads in multi-gene datasets. It uses some computational shortcuts to make the element faster. It has been used for the following paper: <http://www.nature.com/nmeth/journal/v5/n2/full/nmeth.210B.html>

Example: `raxmlHPC -f V -q part -s alg -t referenceTree -m GTRCAT -n TEST`

Warning: this is a test implementation for more efficient handling of multi-gene 1 / whole genome datasets

-f w compute ELW test on a bunch of trees passed via **-z**
The model parameters will be estimated on the first tree only!

Uses the Extended Likelihood * eights test by Strimmer and Rambaut% see <http://www.bi.nlm.nih.gov/pubmed/128728> on a set of trees.

Example: `raxmlHPC -f w -m GTRGAMMA -s alg -z trees -n TEST`

-f W compute ELW test on a bunch of trees passed via **-z**
The model parameters will be re-estimated for each tree

Example: `raxmlHPC -f W -m GTRGAMMA -s alg -z trees -n TEST`

-f x compute pair-wise ML distances, ML model parameters will be estimated on an MP starting tree or a user-defined tree passed via **-t**, only allowed for GAMMA-based models of rate heterogeneity

This option will first compute the pair-wise distances between the sequences of an alignment using a maximum likelihood estimate. To obtain a model parameter estimate (TR% al) has the parameter estimated initially either read in a user supplied tree (e.g. the best-known ML tree) or generate a randomised starting tree. Then compute the parsimony starting tree if no input tree is provided. It then optimises model parameters on this tree and then computes all pair-wise distances.

Example: `raxmlHPC -f x -p 12345 -s alg -m GTRGAMMA -n TEST`

-f y classify a bunch of environmental sequences into a reference tree using parsimony you will need to start RAXML with a non-comprehensive reference tree and an alignment containing all sequences (reference + query)

This command is analogous to the **-f v** command with the only difference that it uses parsimony (or) parsimony reads into the reference tree. Because it uses parsimony it is orders of magnitude faster than the likelihood-based parsimony. You can use it (or) parsimony millions of reads into extremely large reference trees.

Example: `raxmlHPC -f y -m GTRCAT -s alg -t referenceTree -n TEST`

Warning: This algorithm is based on parsimony. Thus a read may be placed into more than one equally parsimonious position on the reference tree. Be careful to take all equally parsimonious placements into account when analysing the results

-F enable ML tree searches under CAT model for very large trees without switching to GAMMA in the end (saves memory). This option can also be used with the GAMMA models in order to avoid the thorough optimization of the best-scoring ML tree in the end.

DEFAULT: OFF

* When conducting analyses under the CAT model of rate heterogeneity in RAXML the program will in most cases try to evaluate the (AMMA-based score and do some additional optimisation of the tree under (AMMA in the very end of the search. If you want to avoid this because of time and/or memory constraints (AMMA needs four times more memory than CAT) you can use this option.

Example: `raxmlHPC -F -m GTRCAT -p 12345 -s alg -N 10 -n TEST`

- g specify the file name of a multifurcating constraint tree this tree does not need to be comprehensive, i.e. must not contain all taxa

This option frequently causes confusion among users for two main reasons: first, the constraint tree you provide is not comprehensive, that is, it does not contain all taxa in the alignment; the taxa not included in the constraint are free, that is, they are allowed to fall into any part of the tree. Second, users often state that the RAxML tree does not correspond to their constraint tree. This is mostly due to the tree visualization tools that are used, depending on whether they are rooted or unrooted. Remember that maximum likelihood trees are always unrooted; if the tree is rooted it may look as if it does not comply with the constraint. Please double-check the tree before posting to the RAxML google group. Finally, also note that any multi-lurcations in the input tree will be resolved via a maximum likelihood search.

Example: `raxmlHPC -g constraintTree -m GTRGAMMA -p 12345 -s alg -n TEST`

- G enable the ML-based evolutionary placement algorithm heuristics by specifying a threshold value (fraction of insertion branches to be evaluated using thorough insertions under ML).

This option can be used in conjunction with the +v and +l options to speed up evolutionary placements of short reads. A setting of 0.5 of branches considered for thorough insertions yields a good accuracy-speed trade-off. For more details you may actually consider to read the paper: <http://sysbio.oxfordjournals.org/content/58/2/20>

Example: `raxmlHPC -f v -G 0.1 -s alg -t referenceTree -m GTRCAT -n TEST`

- h Display this help message.

Example: `raxmlHPC -h`

- H Disable pattern compression.

DEFAULT: ON

This option allows to disable alignment site compression in the input alignment. If alignments contain identical sites they can be compressed to become shorter and thus speed up the likelihood calculations. This option has been developed mostly for internal use in my lab.

Example: `raxmlHPC -H -s alg -p 12345 -m GTRCAT -n TEST`

Warning: using this option may considerably slow down RAxML

- i Initial rearrangement setting for the subsequent application of topological changes phase

This option allows to specify the radius (number of nodes away from the original rooting position) to which rooted subtrees will be re-inserted in the course of S/R subtree

)running re-grafting moves during the tree search. By default this setting will be determined automatically by RAXML. The example below shows it is a radius of 1 that is subtrees will be inserted into all branches that are between 0.5 and 1 nodes away from their original position.

Example: `raxmlHPC -i 10 -s alg -m GTRCAT -p 12345 -n TEST`

-I a posteriori bootstrapping analysis. Use:

-I autoFC for the frequency-based criterion
 -I autoMR for the majority-rule consensus tree criterion
 -I autoMRE for the extended majority-rule consensus tree criterion
 -I autoMRE_IGN for metrics similar to MRE, but include bipartitions under the threshold whether they are compatible or not. This emulates MRE but is faster to compute.

You also need to pass a tree file containing several bootstrap replicates via `-z`

This option allows to carry out the bootstrap "convergence" test that is the test that determines if you have "obtained sufficient support for getting stable support values" a posteriori that is after you have already "obtained" or "instantiated" bootstrap replicates. This option is particularly useful when you are doing large-scale tree searches on supercomputers. Here you would initially run `P$ replicates` and "check if they are sufficient". Then you'd run another `P$` and assess if the \$\$\$ trees you have now are sufficient. My favorite Navors in terms of "convergence" criteria are: autoMRE and if the tree is very large in terms of number of taxa (more than 10000) autoMRE_IGN to save time. Before using this option you should read the "converging" paper (http://link.springer.com/handle/10.1007/978-1-4020-8843-1_2)

Example: `raxmlHPC -m GTRCAT -z RAXML_bootstrap.BS -I autoMRE -n TEST`

Note that if you have several files "containing bootstrap trees" they can easily be "concatenated" by using the Linux `cat` command.

-j **Example:** `cat RAXML_bootstrap.BS_1 RAXML_bootstrap.BS_2 > allBootstraps`
 Specifies that intermediate tree files shall be written to file during the standard ML and BS tree searches.

DEFAULT: OFF

This will similarly print out a "couple of intermediate trees during the tree search and not the final tree only. The intermediate trees are written to files "named" `RAXML_checkpoint.TEST.0`, `RAXML_checkpoint.TEST.1`, etc.

Example: `raxmlHPC -j -m GTRGAMMA -s alg -p 12345 -n TEST`

-J Compute majority rule consensus tree with "-J MR" or extended majority rule consensus tree with "-J MRE" or strict consensus tree with "-J STRICT". For a custom consensus threshold $\geq 50\%$, specify `T<NUM>`, where $100 \geq \text{NUM} \geq 50$.

Options "-J STRICT_DROP" and "-J MR_DROP" will execute an algorithm that identifies dropsets which contain rogue taxa as proposed by Pattengale et al. in the paper "Uncovering hidden phylogenetic consensus".

You will also need to provide a tree file containing several UNROOTED trees

via "-z"

This option actually triggers two different sets of algorithms: one set for building consensus trees and one set for finding rogue taxa. Also note that these algorithms have been parallelized so that if you use the ;Threads version they will run faster.

You should read and cite the following papers when using these options:

For consensus trees=

<http://www.sciencedirect.com/science/article/pii/S0881815505000000>

For identifying rogue taxa=

<http://www.elsevier.com/locate/jmb.2005.10.001>
<http://www.elsevier.com/locate/jmb.2005.10.001>

The two examples below demonstrate an extended majority rule consensus tree and a XPS majority rule consensus tree.

Example: `raxmlHPC -J MRE -z trees -m GTRCAT -n TEST`

Example: `raxmlHPC -J T_75 -z trees -m GTRCAT -n TEST`

The example below identifies rogues using the aforementioned dropset method and using a majority rule threshold. For details please read the paper.

Example: `raxmlHPC -J MR_DROP -z trees -m GTRCAT -n TEST`

-k Specifies that bootstrapped trees should be printed with branch lengths. The bootstraps will run a bit longer, because model parameters will be optimized at the end of each replicate under GAMMA or GAMMA+P-Invar respectively.

DEFAULT: OFF

Example: `raxmlHPC -b 12345 -p 12345 -# 100 -k -s alg -m GTRGAMMA -n TEST`

-K Specify one of the multi-state substitution models (max 32 states) implemented in RAXML. Available models are: ORDERED, MK, GTR

DEFAULT: GTR model

Evidently for this option to have an effect you need to have an alignment containing multi-state characters. If you have several partitions that consist of multi-state characters the model specified via +Z will be applied to all models. Thus it is not possible to assign different models to different multi-state partitions.

Example: `raxmlHPC -p 12345 -m MULTIGAMMA -s multiStateAlignment -n TEST`

-L Compute consensus trees labeled by IC supports and the overall TC value as proposed in Salichos and Rokas 2013. Compute a majority rule consensus tree with -L MR or an extended majority rule consensus tree with -L MRE.

For a custom consensus threshold $\geq 50\%$, specify -L T_<NUM>, where 100 \geq

NUM >= 50.

You will of course also need to provide a tree file containing several UNROOTED trees via -z!

This option allows to compute the TO and O metrics on a consensus tree as described by Sali and Rokas in this paper here=

[http://www.nature.com/nature/journal/v40X/nX/??Q&A=fullnature&2&B\\$.html](http://www.nature.com/nature/journal/v40X/nX/??Q&A=fullnature&2&B$.html)

The method is described in more detail in the following paper /note /ith Sali and Rokas=

[http://embio.oxfordjournals.org/content/early/2006/02/28/embio.oxfordjournals.org.a1111111keyD.P2u\(9x\\$H7R23/](http://embio.oxfordjournals.org/content/early/2006/02/28/embio.oxfordjournals.org.a1111111keyD.P2u(9x$H7R23/)

Warning: The default TO/O calculations are not done exactly as described in our MRE paper. To have RAXML do exactly what is described in the paper) leased edit file bipartitionList.c by commenting out or removing the line
#define BIP_FILTER

There is a dedicated section further down in this manual with detailed instructions.

Example: `raxmlHPC -m GTRCAT -L MRE -z trees -n TEST`

-m Model of Binary (Morphological), Nucleotide, Multi-State, or Amino Acid Substitution:

This is probably the most complex and confusing command line option because it can be used to specify a lot of things. A very important thing to be aware of is that in case you use a partitioned alignment the data type in the partitions and the specified models of substitution to be used are actually specified in the partition file. Thus! you use the -q option the only information that will be extracted from this string here that is passed via -m is which model of rate heterogeneity is going to be used. Note that the rate heterogeneity model is always applied to all data partitions when you can not have one partition be analyzed under CAT and another one be analyzed under GAMMA.

In general the term CAT in the strings always indicates that you want to use the CAT model of rate heterogeneity. The string CATI indicates that after a tree search under CAT you want to evaluate the final trees under a GAMMA plus proportion of invariable sites estimate instead of the default pure GAMMA.

By appending x to the model strings you indicate that you want to use a maximum likelihood estimate for the base frequencies.

By appending the string ASC_ you indicate that you want to apply an ascertainment bias correction to the likelihood calculations. This is useful for binary/morphological datasets that only contain variable sites (the identical morphological features are usually not included in the alignments) when you need to correct for this see e.g. <http://sysbio.oxfordjournals.org/content/55/4/688.full>

For a data set this option might be useful when you analyze alignments of sequences that also don't contain constant sites. Note that for mathematical and numerical reasons you can not apply an ascertainment bias correction to datasets or partitions that contain constant sites. In this case RAXML will exit with an error.

Also note how ambiguous sites are defined in RAXML. If you have a site looking like this=

AAAA---MR

this site is still considered to be an invariable site since – are considered as in most likelihood-based models as completely undetermined characters and thus – may as well as R could be. In other words, RAXML will not allow you to conduct an analysis with an ascertainment bias correction if it encounters such a site in your MSA.

* We also observed that some S9 datasets do not require a model of rate heterogeneity according to the standard model tests. * Hence this is the case you can use –m ASC_GTRCAT in combination with the –v option. This will make RAXML execute an inference under a Jukes (TR model) without any correction for rate heterogeneity.

Also note that while doing the Jukes site rate model of rate heterogeneity (the OAT model) with an ascertainment bias correction, we recommend that you either use OAT in combination with –v (no rate heterogeneity) or the Jukes (gamma model of rate heterogeneity) with ascertainment bias correction –m ASC_GTRGAMMA etc.

Finally, note that as of version 8.2.X of RAXML you will need to specify the type of ascertainment bias correction you want to use, we know of three distinct corrections. For details, please refer to the --asc-corr option.

BINARY:

- | | |
|--------------------|---|
| -m BINCAT[X] | Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under BINGAMMA, depending on the tree search option.
With the optional "X" appendix you can specify a ML estimate of base frequencies. |
| -m BINCATI[X] | Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under BINGAMMAI, depending on the tree search option.
With the optional "X" appendix you can specify a ML estimate of base frequencies. |
| -m ASC_BINCAT[X] | Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under BINGAMMA, depending on the tree search option.
With the optional "X" appendix you can specify a ML estimate of base frequencies.
The ASC prefix will correct the likelihood for ascertainment bias. You will also need to specify the correction type via --asc-corr! |
| -m BINGAMMA[X] | GAMMA model of rate heterogeneity (alpha parameter will be estimated). With the optional "X" appendix you can specify a ML estimate of base frequencies. |
| -m ASC_BINGAMMA[X] | GAMMA model of rate heterogeneity (alpha parameter will be |

estimated). The ASC prefix will correct the likelihood for ascertainment bias.
 With the optional "X" appendix you can specify a ML estimate of base frequencies. You will also need to specify the correction type via --asc-corr!
 -m BINGAMMAI[X] Same as BINGAMMA, but with estimate of proportion of invariable sites. With the optional "X" appendix you can specify a ML estimate of base frequencies.

NUCLEOTIDES:

- m GTRCAT[X] GTR + Optimization of substitution rates + Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated under GTRGAMMA, depending on the tree search option. With the optional "X" appendix you can specify a ML estimate of base frequencies.
- m GTRCATI[X] GTR + Optimization of substitution rates + Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated under GTRGAMMAI, depending on the tree search option. With the optional "X" appendix you can specify a ML estimate of base frequencies.
- m ASC_GTRCAT[X] GTR + Optimization of substitution rates + Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated under GTRGAMMA, depending on the tree search option. With the optional "X" appendix you can specify a ML estimate of base frequencies.
 The ASC prefix will correct the likelihood for ascertainment bias. You will also need to specify the correction type via --asc-corr!
- m GTRGAMMA[X] GTR + Optimization of substitution rates + GAMMA model of rate heterogeneity (alpha parameter will be estimated). With the optional "X" appendix you can specify a ML estimate of base frequencies.
- m ASC_GTRGAMMA[X] GTR + Optimization of substitution rates + GAMMA model of rate heterogeneity (alpha parameter will be estimated). The ASC prefix will correct the likelihood for ascertainment bias.
 With the optional "X" appendix you can specify a ML estimate of base frequencies. You will also need to specify the correction type via --asc-corr!
- m GTRGAMMAI[X] Same as GTRGAMMA, but with estimate of proportion of invariable sites. With the optional "X" appendix you can specify a ML estimate of base frequencies.

MULTI-STATE:

- m MULTICAT[X]** Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under MULTIGAMMA, depending on the tree search option.
- With the optional "X" appendix you can specify a ML estimate of base frequencies.
- m MULTICATI[X]** Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under MULTIGAMMAI, depending on the tree search option.
- With the optional "X" appendix you can specify a ML estimate of base frequencies.
- m ASC_MULTICAT[X]** Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under MULTIGAMMA, depending on the tree search option.
- With the optional "X" appendix you can specify a ML estimate of base frequencies. The ASC prefix will correct the likelihood for ascertainment bias. You will also need to specify the correction type via --asc-corr!
- m MULTIGAMMA[X]** GAMMA model of rate heterogeneity (alpha parameter will be estimated). With the optional "X" appendix you can specify a ML estimate of base frequencies.
- m ASC_MULTIGAMMA[X]** GAMMA model of rate heterogeneity (alpha parameter will be estimated). The ASC prefix will correct the likelihood for ascertainment bias. With the optional "X" appendix you can specify a ML estimate of base frequencies. You will also need to specify the correction type via --asc-corr!
- m MULTIGAMMAI[X]** Same as MULTIGAMMA, but with estimate of proportion of invariable sites. With the optional "X" appendix you can specify a ML estimate of base frequencies.
- You can use up to 32 distinct character states to encode multi-state regions, they must be used in the following order:
- 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, A, B, C, D, E, F, G, H, I, J, K, L, M, N, O, P, Q, R, S, T, U, V
- i.e., if you have 6 distinct character states you would use 0, 1, 2, 3, 4, 5 to encode these.
- The substitution model for the multi-state regions can be selected via the "-K" option

AMINO ACIDS:

- m PROTCATmatrixName[F|X]** specified AA matrix + Optimization of substitution rates + Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate

categories for greater computational efficiency. Final tree might be evaluated automatically under PROTGAMMAmatrixName[F|X], depending on the tree search option.
With the optional "X" appendix you can specify a ML estimate of base frequencies.

-m PROTCATmatrixName[F|X] specified AA matrix + Optimization of substitution rates + Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under PROTGAMMAmatrixName[F|X], depending on the tree search option. With the optional "X" appendix you can specify a ML estimate of base frequencies.

-m ASC_PROTCATmatrixName[F|X] specified AA matrix + Optimization of substitution rates + Optimization of site-specific evolutionary rates which are categorized into numberOfCategories distinct rate categories for greater computational efficiency. Final tree might be evaluated automatically under PROTGAMMAmatrixName[F|X], depending on the tree search option.
With the optional "X" appendix you can specify a ML estimate of base frequencies. The ASC prefix will correct the likelihood for ascertainment bias. You will also need to specify the correction type via --asc-corr!

-m PROTGAMMAmatrixName[F|X] specified AA matrix + Optimization of substitution rates + GAMMA model of rate heterogeneity (alpha parameter will be estimated).
With the optional "X" appendix you can specify a ML estimate of base frequencies.

-m ASC_PROTGAMMAmatrixName[F|X] specified AA matrix + Optimization of substitution rates + GAMMA model of rate heterogeneity (alpha parameter will be estimated). The ASC prefix will correct the likelihood for ascertainment bias.
With the optional "X" appendix you can specify a ML estimate of base frequencies. You will also need to specify the correction type via --asc-corr!

-m PROTGAMMAmatrixName[F|X] Same as PROTGAMMAmatrixName[F|X], but with estimate of proportion of invariable sites.
With the optional "X" appendix you can specify a ML estimate of base frequencies.

Available AA substitution models:

DAYHOFF, DCMUT, JTT, MTREV, WAG, RTREV, CPREV, VT, BLOSUM62, MTMAM, LG, MTART, MTZOA, PMB, HIVB, HIVW, JTTDCMUT, FLU, STMTREV, DUMMY, DUMMY2, AUTO, LG4M, LG4X, PROT_FILE, GTR_UNLINKED, GTR

With the optional "F" appendix you can specify if you want to use empirical base frequencies. AUTOF and AUTOX are not supported any more, if you specify AUTO it will test prot subst. models with and without empirical

base frequencies now!

Please note that for partitioned models you can in addition specify the per-gene AA model in the partition file. Also note that if you estimate AA GTR parameters on a partitioned dataset, they will be linked (estimated jointly) across all partitions to avoid over-parametrization

Warning: * When using the **LG4X** (protein substitution model) please keep in mind that this model actually has more free parameters than a normal fixed protein substitution model. Thus for each partition that evolves under LG4X there are 4 additional free parameters for the 4 sites and 4 for the rates. Note that it is only 4 and 4 rate parameters because the 4 sites need to sum to 1. The 4th parameter is thus determined by the other 3 and the sum of the 4 parameters is 1. The 4th parameter also needs to be 1. (see page 202 in the reference) <http://mbe.oxfordjournals.org/content/20/2/101.full.pdf>. Also note that in RAXML LG4X can be used with empirical base frequencies drawn from the data (LG4X2) and with a ML estimate of base frequencies (LG4X4). Unlike in the published LG4X model the base frequencies are the same **shared/linked** for all 4 substitution matrices of the LG4X models and there are 4 additional parameters in these models. * When using this modified form of LG4X please make sure to state it explicitly when you publish results and keep in mind that they have more parameters when you do AIC test etc. Also these options have been implemented for "convenience" in RAXML but do not correspond to the spirit of the above paper where each of the 4 matrices has its own set of base frequencies. *Olivier Gascuel asked me to explicitly state that he does not like the following models: LG4XF, LG4MF, LG4XX, LG4MX to be used and I think he has a point!*

The example below will execute a simulation search under the OAT approximation of rate heterogeneity on a .9A dataset and evaluate the final tree under (AMMA).

Example: `raxmlHPC -m GTRCAT -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the OAT approximation of rate heterogeneity on a .9A dataset but evaluate the final tree under (AMMAU; + nvar).

Example: `raxmlHPC -m GTRCATI -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the OAT approximation of rate heterogeneity on a .9A dataset and do a maximum likelihood estimate of the base frequencies instead of using empirical base frequencies.

Example: `raxmlHPC -m GTRCATX -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the (AMMA model of rate heterogeneity on a .9A dataset.

Example: `raxmlHPC -m GTRGAMMA -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the (AMMA model of rate heterogeneity and "correct" for ascertainment bias on a .9A dataset.

Example: `raxmlHPC -m ASC_GTRGAMMA -s alg -p 12345 -n TEST`

The example below will execute a simulation search under a Jukes (TR model) no rate heterogeneity and "correct" for ascertainment bias on a .9A dataset.

Example: `raxmlHPC -m ASC_GTRCAT -V -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the (AMMA model) of rate heterogeneity on a binary dataset.

~~PROTGAMMAILGX~~ **Example:** `raxmlHPC -m BINGAMMA -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the (AMMA model) of rate heterogeneity on a protein dataset using the substitution matrix and the base frequencies that come with the *A(model.

Example: `raxmlHPC -m PROTGAMMAWAG -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the (AMMA model) of rate heterogeneity on a protein dataset using empirical base frequencies and the L(substitution model.

Example: `raxmlHPC -m PROTGAMMALGF -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the (AMMAU; + nvar model) of rate heterogeneity on a protein dataset using a maximum likelihood estimate of the base frequencies and the L(substitution model.

Example: `raxmlHPC -m PROTGAMMAILGX -s alg -p 12345 -n TEST`

The example below will execute a simulation search under the (AMMA model) of rate heterogeneity on a protein dataset using empirical base frequencies and estimating a (TR model) of amino acid substitution.

Warning: This may be very slow and you might overparameterize the model since you need to do a maximum likelihood estimate of 80 rate parameters in the (TR matrix.

Example: `raxmlHPC -m PROTGTRGAMMA -s alg -p 12345 -n TEST`

The example below will automatically determine which is the best of the one with the highest likelihood score on the parsimony starting tree) protein substitution model for your dataset using the base frequencies that come with the models. It will choose among the following models= DAYHOFF, DCMUT, JTT, MTREV, WAG, RTREV, CPREV, VT, BLOSUM62, MTMAM, LG, MTART, MTZOA, PMB, HIVB, HIVW, JTTDCMUT, FLU, DUMMY, DUMMY2. These models will not be considered= LG4M, LG4X, PROT_FILE, GTR_UNLINKED, GTR!

RAXML will not also automatically test models with and without empirical base frequencies. The criterion for making this choice can not be selected via the

files

A weighted average of the branch lengths is computed by using the respective partition lengths

DEFAULT: OFF

This will enable an independent partition estimate of branch lengths. Note that this increases dramatically the number of free parameters in your model. An unrooted binary tree has $2n-3$ branch lengths where n is the number of taxa. If you use this option instead of $2n-3$ parameters you will get $p(2n-3)$ parameters where p is the number of partitions. RAxML will also run notably slower if you use this option.

Apart from the normal output tree% RAxML will also generate output files for each partition individually that allow to recover the partition branch lengths if needed.

Example: `raxmlHPC -p 12345 -M -m GTRGAMMA -s alg -q part -n TEST`

-n Specifies the name of the output file.

This option has to be always specified. The arbitrary name passed via `-n` will be appended to all RAxML output files such that you know which files have been generated by which invocation.

If you intend to do multiple runs that write files into the same directory with the same name specified by `-n` the program will exit with an error to prevent you overwriting output files from a previous run.

Example: `raxmlHPC -p 12345 -m GTRGAMMA -s alg
-n MySuperDuperRAxMLOutputFileName`

-o Specify the name of a single outgroup or a comma-separated list of outgroups, e.g., `-o Rat` or `-o Rat,Mouse`, in case that multiple outgroups are not monophyletic the first name in the list will be selected as outgroup, don't leave spaces between taxon names!

If there is more than one outgroup a heuristic for monophyly of the outgroups in the respective output tree will be performed. If the outgroup is not monophyletic the tree will be rooted at the first outgroup in the list and a respective warning will be printed.

Remember in mind that outgroups are just a tree drawing option% the trees remain% essentially unrooted trees.

Example: `raxmlHPC -s alg -p 12345 -m GTRGAMMA -o Rat,Mouse -n TEST`

A comment on using outgroups: How I would do it.

In general I'd avoid using outgroups in the initial ML and bootstrap analyses. I'd proceed as follows and provide the rationale for this below:

1. Build a tree (ML/URS search) only on the ingroup
2. Then use E; A (see= [http://sysbio.oxfordjournals.org/content/\\$B20&short](http://sysbio.oxfordjournals.org/content/$B20&short)) to place the outgroups onto the tree a posteriori.

This has several advantages=

1. You can use different outgroups
2. You don't have to re-run the entire analysis if the outgroup somehow perturbed the

analysis of the ingroup). There have been debates on this issue.

B. To avoid the outgroup ambiguity of the ingroup because of the way the evolutionary likelihood algorithm has been built.

?. You get likelihood probabilities for each outgroup e.g. how likely the outgroup is to fall into one branch or another. This also tells you if your outgroup is good (high probability for being placed into one side of the branch or if it is bad suggests across the tree).

- O Disable check for completely undetermined sequence in alignment. The program will not exit with an error message when "-O" is specified.

DEFAULT: check enabled

In general it does not make sense to analyze alignments where one or more sequences consist of completely undetermined characters e.g. -, N, ?, X etc. because you don't have any information about these sequences and they will simply randomly scatter throughout the tree. Normally RAXML will exit with an error message when it detects such sequences. This option here can disable this behavior but be warned that you really need to know what you are doing when using this option.

Example: `raxmlHPC -s alg -p 12345 -m GTRGAMMA -O -n TEST`

- p Specify a random number seed for the parsimony inferences. This allows you to reproduce your results and will help me debug the program.

For all options/algorithms in RAXML that require some sort of randomization this option must be specified. Make sure to pass different random number seeds to RAXML and not only &2B?P as have done in the examples. * then not specifying it when it is required by RAXML the program will exit with a respective error message.

In the example below for the ML tree the random number seed is required for randomized stepwise addition order parsimony starting tree that is computed prior to the actual ML optimization.

Example: `raxmlHPC -s alg -p 2352890 -m GTRGAMMA -n TEST`

- P Specify the file name of a user-defined AA (Protein) substitution model. This file must contain 420 entries, the first 400 being the AA substitution rates (this must be a symmetric matrix) and the last 20 are the empirical base frequencies

Note that the substitution model you specify via `-m PROTGAMMAWAG` in the example below will be ignored here. RAXML will only extract the information that this is a protein data alignment and that you want to use the GAMMA model of rate heterogeneity from the string and will ignore WAG.

If you want to use your own protein substitution models for partitioned datasets the substitution model files that have the same format as described here need to be specified in the partition file.

Example: `raxmlHPC -s alg -p 12345 -P myProteinModel -m PROTGAMMAWAG -n TEST`

- q Specify the file name which contains the assignment of models to alignment partitions for multiple models of substitution. For the syntax of this file please consult the manual.

This option allows you to specify the regions of your alignment for which an individual model of nucleotide substitution should be estimated. They can also contain different types of data.

This will typically be useful to infer trees for long bin terms of basepairs multi-gene alignments. For example, TRAM is used to estimate individual parameters (TR rates and empirical base frequencies) which will be estimated and optimized for each partition.

Since RAXML can handle alignments consisting of different data types (binary data, A data, protein data etc.) you must specify the type of data for each partition in the partition file before the partition name.

For A data this just means that you have to add DNA to each line in the partition file. For AA data this is done by specifying the respective AA substitution matrix (e.g. WAG or LG) you want to use for a partition. For binary data you'd specify BIN and for multi-state data MULTI.

For if you want to analyze a specific partition using an ascertainment bias correction you need to re-encode ASC_ to the data type (e.g. ASC_WAG or ASC_DNA).

For a specific partition you want to use a maximum likelihood estimate for the base frequencies you'd add a) end x to the data type name (e.g. DNAX or LGX).

You want to use empirical base frequencies instead of the default redefined ones that ship with the models you'd write (e.g. WAG or JTT).

You want to use a protein substitution model of your own for a specific partition it must be in the same format as specified for the -P option. Instead of specifying the data type with WAG for instance you will need to specify the protein substitution model file name in square brackets (e.g. [myProteinSubstitutionModelFileName]).

You want to do a partitioned analysis of concatenated AA and A partitions you can either specify -m GTRGAMMA or e.g. -m PROTGAMMAWAG. The only thing that will be extra added from the string passed via +m is the model of rate heterogeneity you want to use.

You have a file A alignment with &\$\$\$b) from t/o genes gene& positions &WP\$\$\$ and gene2 positions P&W&\$\$\$ the information in the partition file should look as follows=

```
DNA, gene1 = 1-500
DNA, gene2 = 501-1000
```

You want an ML estimate of frequencies for the first partition it will look like this=

```
DNAX, gene1 = 1-500
DNA, gene2 = 501-1000
```

To analyze the first partition with ascertainment bias correction you need to write=

```
ASC_DNA, gene1 = 1-500
DNA, gene2 = 501-1000
```

gene& is started through the alignment (e.g. positions &W2\$\$\$ and 8\$W&\$\$\$ you specify this with=

```
DNA, gene1 = 1-200, 800-1,000
DNA, gene2 = 201-799
```

You can also assign distinct models to the codon positions i.e. if you want a distinct model to be estimated for each codon position in gene& you can specify=

```
DNA, gene1codon1 = 1-500\3
DNA, gene1codon2 = 2-500\3
DNA, gene1codon3 = 3-500\3
DNA, gene2 = 501-1000
```

! you only need a distinct model for the Brd "odon)osition you "an /rite=

```
DNA, gene1codon1andcodon2 = 1-500\3, 2-500\3
DNA, gene1codon3 = 3-500\3
DNA, gene2 = 501-1000
```

As already mentioned for AA data you must specify the transition matrices for each partition=

```
JTT, gene1 = 1-500
WAGF, gene2 = 501-800
WAG, gene3 = 801-1000
```

The AA substitution model must be the first entry in each line and must be separated by a comma from the gene partition name just like the . 9A token above. <ou "an not assign different models of rate heterogeneity to different partitions i.e. it will be either OAT % (AMMA % (AMMA et". for all)artitions as specified with -m .

! you want to use a)rotein model of your own for)artition one you'd /rite=

```
[prot~myProtenSubstitutionModelFileName], gene1 = 1-500
WAGF, gene2 = 501-800
WAG, gene3 = 801-1000
```

To use a maximum likelihood estimate of the base frequencies for all)artitions you'd /rite=

```
JTTX, gene1 = 1-500
WAGX, gene2 = 501-800
WAGX, gene3 = 801-1000
```

Finally if you have a "concatenated . 9A and AA alignment with . 9A data at)ositions &W P\$\$ and AA data at P\$\$&\$\$ /ith the * A(model the)artition file should look as follows=

```
DNA, gene1 = 1-500
WAG, gene2 = 501-1000
```

A)artition file for binary data would look like this=

```
BIN, gene1 = 1-500
BIN, gene2 = 501-1000
```

and for multi-state data

```
MULTI, gene1 = 1-500
MULTI, gene2 = 501-1000
```

Example: raxmlHPC -s alg -m GTRGAMMA -q part -p 12345 -n TEST

-r Specify the file name of a binary constraint tree.
 This tree does not need to be comprehensive, i.e. must not contain all taxa

This option allows you to pass a binary constraint backbone tree in NEXUS format to RAxML.
 Note that using this option only makes sense if this tree contains less taxa than the input alignment. The remaining taxa will initially be added by using parsimony criterion. Once a comprehensive tree with all taxa has been obtained it will be optimized under ML respecting the restrictions of the binary constraint tree.

Thus option will not change the structure of the binary backbone tree you provide as input at all. What it will do is to just add the taxa not contained in the binary backbone tree to the best positions based on their likelihood. The result is a fully resolved binary tree.

Example: `raxmlHPC -s alg -m GTRGAMMA -r constr -p 12345 -n TEST`

-R Specify the file name of a binary model parameter file that has previously been generated with RAxML using the `-f e` tree evaluation option. The file name should be: `RAxML_binaryModelParameters.runID`

This option can be used to save time for some other RAxML options by not re-estimating the model parameters of a fixed tree again but only doing this once. The binary model parameter input file is always automatically generated by the `-f e` tree evaluation function.

It is useful with the `-f v` option when you want to place different reads into the same reference tree and with the quartet evaluation function `-f q` in particular the parallel version of it.

Example: `raxmlHPC -f v -R RAxML_binaryModelParameters.PARAMS
 -t RAxML_result.PARAMS -s alg -m GTRCAT -n TEST2`

-s Specify the name of the alignment data file in PHYLIP or FASTA format

Specify the name of the alignment data file which can be in relaxed NEXUS format. Relaxed means that you don't have to worry if the sequence file is interleaved or sequential and that the taxon names are too long. What you do need to worry about though is that there always needs to be a space between the taxon name and the sequence.

Sequence names can be of variable length between 1 and 255 characters. If you need longer taxon names you can adapt the constant `#define nmlength 256` in file `axml.h` appropriately. Moreover RAxML is not sensitive with respect to the formatting of tabs and spaces of interleaved NEXUS files.

RAxML can also parse FASTA format. If RAxML notices that it cannot parse a phylogenetic format it will try to parse the alignment file as FASTA format. So far it has been able to parse all FASTA files.

A plethora of small examples input files for different data types can be found in the step-by-step online tutorial: <http://o.h+its.org/exelixis/ebolit/are@raxml@handsKon.html>

Example: `raxmlHPC -s alignment.fasta -m GTRGAMMA -p 12345 -n TEST`

-S Specify the name of a secondary structure file. The file can contain "." for alignment columns that do not form part of a stem and characters "()<>[]{}" to define stem regions and pseudoknots

Secondary structure models for an RNA alignment works slightly differently than passing other data types to RAXML because we read in a plain RNA alignment and then need to tell RAXML by an additional text file that is passed via `-S /hi" h RNA alignment sites need to be grouped@evolve together.`

* We do this in a standard bracket notation written into a plain text file e.g. our . RNA test alignment has ' \$ sites thus our secondary structure file needs to contain a string of ' \$ "characters like this one=

```
.....((.....)).....(.....).....
```

The `||` symbol indicates that this is just a normal RNA site while the brackets indicate stems. Evidently the number of opening and closing brackets must match. In addition it is also possible to specify pseudo knots with additional symbols `<>[]{}|` for instance=

```
.....((.....)).....{....(....}....).....
```

In terms of models there are '+state% X+state and '&'+state models for accommodating secondary structure that are specified via `-A`.

Available models are S6A, S6B, S6C, S6D, S6E, S7A, S7B, S7C, S7D, S7E, S7F, S16, S16A, S16B. The default is the (TR &'+state model 6-A S16B. In RAXML the same nomenclature as in ;HASE is used so please consult the ;hase manual at <http://intranet.cs.man.ac.uk/ai//Software/phase/phase-2.0-manual.pdf> for a nice and detailed description of these models.

Example: `raxmlHPC -m GTRGAMMA -p 12345 -S secondaryStructure.txt
-s rna.phy -n TEST`

A common question is whether secondary structure models can also be partitioned. This is presently not possible. However you can partition the underlying RNA data e.g. use t/o partitions for our . RNA dataset as before. * hat RAXML will do internally though is to generate a third partition for secondary structure that does not take into account that distinct secondary structure sites may be part of different partitions of the alignment.

-t Specify a user starting tree file name in Newick format

Specifies a user starting tree file name which must be in Newick format. The branch lengths of that tree will generally be ignored for most options and re-estimated instead by RAXML. Note that you can also specify a non-comprehensive (not containing all taxa in the alignment) starting tree. This might be useful if newly aligned/sequenced taxa have been added to your alignment and you want to extend the current tree.

Initially taxa will be added to the tree using parsimony. The comprehensive tree will then be optimized under ML.

Example: `raxmlHPC -s alg -m GTRGAMMA -t tree -p 12345 -n TEST`

-T PTHREADS VERSION ONLY! Specify the number of threads you want to run. Make sure to set `"-T"` to at most the number of CPUs you have on your machine, otherwise, there will be a huge performance decrease!

T is set to \$ by default the ; Threads version will produce an error if you do not set `-T` to at least 2.

Example: `raxmlHPC-PTHREADS -T 4 -s alg -m GTRGAMMA -p 12345 -n TEST`

-u use the median for the discrete approximation of the GAMMA model of rate heterogeneity

DEFAULT: OFF

! you s)e"ily this o)tion% RAxML /ill use the median instead o! the mean !or the dis"rete (amma model o! rate heterogeneity. Ty)i"ally% using the median /ill yield slightly better likelihood s"ores.

Example: raxmlHPC -p 12345 -u -s alg -m GTRGAMMA -n TEST

-U Try to save memory by using SEV-based implementation for gap columns on large gappy alignments. The technique is described here: <http://www.biomedcentral.com/1471-2105/12/470>

This will only work for DNA and/or PROTEIN data and only with the SSE3 or AVX-vextorized version of the code.

This o)tion "an hel) you to save memory and)otentially also time on large ga))y multi+ gene alignments /ith missing data% in)arti"ular in "ases /here you have a ty)i"al ga))y)artitioned alignment 6data !or "ertain taxa missing !or "ertain genes@)artitions8. The amount o! saved memory is roughly)ro)ortional to the)ro)ortion o! missing data.

Example: raxmlHPC -p 12345 -U -s alg -m GTRGAMMA -n TEST

-v Display version information

. is)lays the RAxML version you are using.

Example: raxmlHPC -v

-V Disable rate heterogeneity among sites model and use one without rate heterogeneity instead. Only works if you specify the CAT model of rate heterogeneity.

DEFAULT: use rate heterogeneity

This o)tion "an be used i! your dataset better !ts a model /ithout rate heterogeneity. This is rare% but su"h datasets do exist.

Example: raxmlHPC -m GTRCAT -s alg -p 12345 -V -n TEST

-w FULL (!) path to the directory into which RAxML shall write its output files

DEFAULT: current directory

9ame o! the /orking dire"tory /here RAxML shall /rite its out)ut !les to. 9ote that you need to s)e"ily the !ull)ath and not the relative)ath6

Example: raxmlHPC -m GTRCAT -p 12345 -s alg
-w /home/stamatak/Desktop/myAnalysis -n TEST

-W Sliding window size for leave-one-out site-specific placement bias algorithm only effective when used in combination with -f S

DEFAULT: 100 sites

. el nes the sliding /indo/ si7e !or the +! S o)tion% &\$ \$ usually yields relatively smooth
)lots.

Example: `raxmlHPC -f S -s alg -t tree -m GTRGAMMA -N 100 -W 200 -n TEST`

-x Specify an integer number (random seed) and turn on rapid bootstrapping
CAUTION: unlike in previous versions of RAXML will conduct rapid BS
replicates under the model of rate heterogeneity you specified via -m and
not by default under CAT

This /ill invoke the ra)id bootstra))ing algorithm des"ribed in
<http://sysbio.oxfordjournals.org/content/PX/P/XP8.short>

Example: `raxmlHPC -x 12345 -p 12345 -# 100 -m GTRCAT -s alg -n TEST`

<ou "an also "ombine this /ith the bootsto))ing o)tion 6bootstra) "onvergen"e "riterion8%
e.g.=

Example: `raxmlHPC -x 12345 -p 12345 -# AUTO_MRE -m GTRCAT -s alg -n TEST`

or also have RAXML sear"h !or the best+s"oring ML tree alter the bootstra) sear"hes using
the -f a o)tion=

Example: `raxmlHPC -x 12345 -p 12345 -# AUTO_MRE -m GTRCAT -s alg -f a
-n TEST`

-X Same as the "-y" option below, however the parsimony search is more
superficial.
RAXML will only do a randomized stepwise addition order parsimony tree
reconstruction without performing any additional SPRs.
This may be helpful for very broad whole-genome datasets, since this can
generate topologically more different starting trees.

DEFAULT: OFF

Example: `raxmlHPC -X -m GTRCAT -s alg -p 12345 -n TEST`

-y If you want to only compute a parsimony starting tree with RAXML specify -y,
the program will exit after computation of the starting tree

DEFAULT: OFF

! you /ant to only "om)ute a randomi7ed)arsimony starting tree /ith RAXML and not
exe"ute an ML analysis o! the tree s)e"i!y -y. The)rogram /ill exit alter "om)utation o!
the starting tree.

Example: `raxmlHPC -y -m GTRCAT -s alg -p 12345 -n TEST`

-Y Pass a quartet grouping file name defining four groups from which to draw
quartets
The file input format must contain 4 groups in the following form:
(Chicken, Human, Loach), (Cow, Carp), (Mouse, Rat, Seal), (Whale, Frog);

Only works in combination with -f q !

The example below will randomly draw 4 quartets from the reduced quartet group above and evaluate their likelihood.

Example: `raxmlHPC -f q -m GTRGAMMA -t tree -Y quartetFile -n 100 -n TEST`

-z Specify the file name of a file containing multiple trees e.g. from a bootstrap that shall be used to draw bipartition values onto a tree provided with **-t**.
It can also be used, for instance to compute per site log likelihoods in combination with **-f g** and to read a bunch of trees for a couple of other options (**-f h**, **-f m**, **-f n**).

This option is required in combination with a lot of other RAXML options in essence every time you need to read in a bunch of trees. Below is an example where you use the option to draw RS support values on a best-known ML tree.

Example: `raxmlHPC -f b -m GTRCAT -t mlTree -z bootstrapTrees -n TEST`

-#|-N Specify the number of alternative runs on distinct starting trees
In combination with the **"-b"** option, this will invoke a multiple bootstrap analysis
Note that **"-N"** has been added as an alternative since **-#** sometimes caused problems with certain MPI job submission systems, since **-#** is often used to start comments.
If you want to use the bootstopping criteria specify **-# autoMR** or **-# autoMRE** or **-# autoMRE_IGN** for the majority-rule tree based criteria (see **-I** option) or **-# autoFC** for the frequency-based criterion.
Bootstopping will only work in combination with **-x** or **-b**

DEFAULT: 1 single analysis

Specify the number of alternative runs on distinct starting trees e.g. **-# 10** or **-N 10** is specified RAXML will compute & distinct ML trees starting from & distinct randomized maximum parsimony starting trees.

In combination with the **-b** option this will invoke a multiple bootstrap analysis.

In combination with **-x** this will invoke a rapid RS analysis and combined with **-f a -x a** rapid RS search and thereafter a thorough ML search on the original alignment.

* We introduced **-N** as an alternative to **-#** since the special character **#** seems to sometimes cause problems with certain batch job submission systems.

In combination with **-f j** this will generate number of runs bootstrapped alignment files.

Note that several other program options such as the quartet evaluation algorithm **-f q** option or the site-specific element bias algorithm **-f s** option require this option to be specified as well. If you forgot to specify this option RAXML will tell you that you need to do so.

The example below will do 2 independent ML tree searches on 2 randomized sets in addition order parsimony tree.

Example: `raxmlHPC -p 12345 -s alg -n TEST -m GTRGAMMA -# 20`

--mesquite Print output files that can be parsed by Mesquite.

DEFAULT: Off

Example: `raxmlHPC -p 12345 -s alg --mesquite -n TEST -m GTRGAMMA`

`--silent` Disables printout of warnings related to identical sequences and entirely undetermined sites in the alignment. The program might run faster when this is enabled.

DEFAULT: Off

Example: `raxmlHPC -p 12345 --silent -s alg -n TEST -m GTRGAMMA`

`--no-seq-check` Disables checking the input MSA for identical sequences and entirely undetermined sites. Enabling this option may save time, in particular for large phylogenomic alignments. Before using this, make sure to check the alignment using the `-f c` option!

DEFAULT: Off

Example: `raxmlHPC -p 12345 -s alg --no-seq-check -n TEST -m GTRGAMMA`

`--no-bfgs` Disables automatic usage of BFGS method to optimize GTR rates on unpartitioned DNA datasets. Using BFGS can improve speeds for model optimization by up to 30%. It's enabled by default when analyzing single-partition DNA datasets.

DEFAULT: BFGS on

Example: `raxmlHPC -p 12345 -s alg --no-bfgs -n TEST -m GTRGAMMA`

`--asc-corr` Allows to specify the type of ascertainment bias correction you wish to use. There are three types available:

`--asc-corr=lewis`: the standard correction by Paul Lewis

`--asc-corr=felsenstein`: a correction introduced by Joe Felsenstein that allows to explicitly specify the number of invariable sites (if known) one wants to correct for.

`--asc-corr=stamatakis`: a correction introduced by myself that allows to explicitly specify the number of invariable sites for each character (if known) one wants to correct for.

Let's start with the Lewis correction that is easy to handle=

Example: `raxmlHPC -p 12345 -s alg -n TEST -m ASC_GTRGAMMA --asc-corr=lewis`

The above will run a standard tree search with the ascertainment bias correction by Paul Lewis.

The two other options are a bit more tricky to use. They have been designed for cases where the exact number of invariable sites for the dataset is actually known (Felsenstein option) or the exact frequencies of each invariable character state is known (Stamatakis option). Note that in such cases the number of invariable sites can be larger than the

number of variable sites in the alignment. This is not the case for the Le/ is "correction" because the maths break down.

For those two "corrections" we thus need to tell RAXML that the number of invariable sites for each partition of felsenstein is and that the numbers of invariable sites per state for each partition is of stamatakis. This is best handled via the partition file when even if you only have one partition you nonetheless need to specify a partition file via -q to pass this information to RAXML.

Example 2: `raxmlHPC -p 12345 -s alg -n TEST -m ASC_GTRGAMMA --asc-corr=stamatakis -q part`

Example 3: `raxmlHPC -p 12345 -s alg -n TEST -m ASC_GTRGAMMA --asc-corr=felsenstein -q part`

Below we give an example of such a partition file=

```
[asc-p1.txt], ASC_DNA, p4=1-1000
[asc-p2.txt], ASC_DNA, p5=1001-1965
```

Here the first entries in each line provide the name of the file containing the invariable site counts or frequencies for each partition that are stored in plain text files called p1.txt and p2.txt. Note that the respective file names must be preceded by the asc key/word and the ~ separator symbol.

Thus each partition is associated with a file containing these counts.

For the Felsenstein "correction" these files would look like this=

```
>1.txt=
1000
```

```
>2.txt=
2000
```

* With that we tell RAXML that the likelihood first partition shall be "corrected" for 1000 invariable sites while the likelihood of the second partition shall be "corrected" for 2000 invariable sites. Note that each file p1.txt and p2.txt must only contain a single line of rising integer values.

For my "correction" the files would look as follows for 9A data=

```
>1.txt=
250 300 400 100
```

```
>2.txt=
500 300 200 200
```

* With that we tell RAXML that the likelihood of the first partition shall be "corrected" for 250 sites consisting of 300 sites consisting of 400 sites of 100 sites of 100. For the second partition the likelihood is "corrected" for 500 invariable sites of 300 of 200 of 200.

The difference among the two "corrections" is that with the Felsenstein "correction" we know how many invariable sites there are but not their composition so we "correct" for the absence of \$\$\$\$ invariable sites in the first partition that would consist of As or Os or (s or

Ts. My "orre"tion "an be used /hen /e do kno/ the exa"t !re4uen"ies o! invariable site)atterns.

WARNING: * hen generating these as"ertainment bias "orre"tion l les 6e.g.%p1.txt and p2.txt8 !or my "orre"tion)ay)arti"ular attention to the order o! states that "orres)onds to the order o! !re4uen"ies o! invariable sites)er state=

BINARY: 0,1
DNA: A, C, G, T
PROTEIN: A, R, N, D, C, Q, E, G, H, I, L, K, M, F, P, S, T, W, Y, V

--flag-check When using this option, RAxML will only check if all command line flags specified are available and then exit with a message listing all invalid command line flags or with a message stating that all flags are valid.

This o)tion is intended !or B^d)arty /ra))er so!t/are that invokes RAxML. t allo/s !or "he"king be!orehand i! "ertain o)tions are available in the RAxML version being used or not. The o)tion /ill !ust analy7e the "ommand line !or valid@invalid Nags and then exit /ith an a))ro)riate message.

Example 1: raxmlHPC --flag-check -f a -s alg -p 12345 -N 20 -x 12345
-m GTRGAMMA -n T1

Output: All options supported

Example 2: raxmlHPC --flag-check -f a -s alg -p 12345 -N 20 -x 12345
-m GTRGAMMA --nonsense-option -n T1

Output: Option --nonsense-option not supported

--auto-prot=ml|bic|aic|aicc When using automatic protein model selection you can chose the criterion for selecting these models. RAxML will test all available prot subst. models except for LG4M, LG4X, and GTR-based models with and without empirical base frequencies. You can chose between ML score based selection and the BIC, AIC, and AICc criteria.

DEFAULT: ml

Example: raxmlHPC -s alg -p 12345 -m PROTGAMMAAUTO --auto-prot=bic -n T1

The above invo"ation /ill sele"t bet/een)rotein models using the Bayesien n!ormation Oriterion.

2or understanding the !ollo/ing three o)tions that all refer to the standard)la"ement out)ut !ormat 6.jplace out)ut l le o! the E; A8 /e have del ned%it /ill be very hel)!ul i! you read the "orres)onding)a)er=

[http://www.losone.org/arti"le@n!oSBAdoiS22&\\$.&BX&S22!ournal.\)one.\\$B&\\$\\$Q](http://www.losone.org/arti)

--epa-keep-placements=number specify the number of potential placements you want to keep for each read in the EPA algorithm .

Note that, the actual values printed will also depend on the settings for

`--epa-prob-threshold=threshold !`

DEFAULT: 7

Example: `raxmlHPC -f v --epa-keep-placements=100 -t tree -m GTRCAT -s alg
-n T1`

Here% RAXML /ill)rint at most &\$\$)la"ements to the .jplace)la"ement lle%
de)ending on the setting o! --epa-prob-threshold.

`--epa-prob-threshold=threshold` specify a percent threshold for including potential placements of a read depending on the maximum placement weight for this read. If you set this value to 0.01 placements that have a placement weight of 1 per cent of the maximum placement will still be printed to file if the setting of `--epa-keep-placements` allows for it

DEFAULT: 0.01

Example: `raxmlHPC -f v --epa-prob-threshold=0.05 -t tree -m GTRCAT -s alg
-n T1`

Here% RAXML /ill)rint those)la"ements to the .jplace)la"ement lle that have a likelihood /eight o! at least PS o! that o! the maximum)la"ement /eight !or the s)e"il " read and i! the number o!)la"ements to)rint as s)e"il ed by Wepa-keep-placements has not been ex"eeded.

`--epa-accumulated-threshold=threshold` specify an accumulated likelihood weight threshold for which different placements of read are printed to file. Placements for a read will be printed until the sum of their placement weights has reached the threshold value. Note that, this option can neither be used in combination with `--epa-prob-threshold` nor with `--epa-keep-placements`!

Example: `raxmlHPC -f v --epa-accumulated-threshold=0.95 -t tree -m GTRCAT -s
alg -n T1`

Here% RAXML /ill)rint those)la"ements o! a read to the .jplace)la"ement lle until their a""umulated likelihood /eights has rea"hed a value o! \$.QP 6out o! a total o! &.\$8.

`--JC69` specify that all DNA partitions will evolve under the Jukes-Cantor model, this overrides all other model specifications for DNA partitions.

DEFAULT: Off

Example: `raxmlHPC -p 12345 -m GTRGAMMA -s alg --JC69 -n T1`

`--K80` specify that all DNA partitions will evolve under the K80 model, this overrides all other model specifications for DNA partitions.

DEFAULT: Off

Note that% the out)ut o! the)rogram might look a bit /eird% sin"e unlike in the del nition o! the model% RAXML a"tually estimates the rates !rom A O (and O O T /hile all other rates are set to &.\$. Note that% this does not matter% sin"e the rates in the rate matrix are relative

rates[the results likelihoods will be the same.

Example: `raxmlHPC -p 12345 -m GTRGAMMA -s alg --K80 -n T1`

`--HKY85` specify that all DNA partitions will evolve under the HKY85 model, this overrides all other model specifications for DNA partitions.

DEFAULT: Off

Note that the output of the program might look a bit weird since unlike in the definition of the model RAxML actually estimates the rates from A to C and C to T while all other rates are set to 1. Note that this does not matter since the rates in the rate matrix are relative rates[the results likelihoods will be the same.

Example: `raxmlHPC -p 12345 -m GTRGAMMA -s alg --HKY85 -n T1`

`--set-thread-affinity` specify that thread-to-core affinity shall be set by RAxML for the hybrid MPI-PThreads version

DEFAULT: Off

You might have to use this option for the hybrid MPI+Threads version of RAxML on some cluster installations such that the node interacts properly with the scheduling system. On other systems however enabling this option may yield problems. This unfortunately a matter of trial and error that is due to how different schedulers handle hybrid MPI+Threads nodes

Example: `raxmlHPC-HYBRID -m GAMMA -s alg -p 12345 -N 20 --set-thread-affinity -n T1`

`--bootstop-perms=number` specify the number of permutations to be conducted for the bootstopping/bootstrap convergence test. The allowed minimum number is 100!

DEFAULT: 100

You might want to set this to 1000 or higher to reduce the variation in the number of bootstrap trees you calculate on noisy or large datasets. The variance of the number of bootstrap trees that bootstrap runs with different random number seeds will generate will decrease by increasing the number of permutations.

Example: `raxmlHPC -B 0.02 --bootstop-perms=1000 -b 12345 -p 12345 -# AUTOMR -s alg -m GTRCAT -n TEST`

`--quartets-without-replacement` specify that quartets are randomly subsampled, but **without** replacement.

DEFAULT: random sampling **with** replacements

Example: `raxmlHPC -s alg -p 12345 -m GTRGAMMA -f q -N 50 --quartets-without-replacement -n TEST`

`--print-identical-sequences` specify that RAxML shall automatically generate a .reduced alignment with all undetermined columns removed, **but without** removing exactly identical sequences

DEFAULT: identical sequences will also be removed in the .reduced file

Example: `raxmlHPC -s alg -p 12345 -m GTRGAMMA --print-identical-sequences -n TEST`

IX. Output Files

. ending on the sear" h)arameter settings RAxML /ill /rite a number o! out)ut l les. The most im)ortant l les% a run named -n exampleRun /ill /rite% are listed belo/. Sin"e the number and names o! out)ut l les vary de)ending on the RAxML o)tions you used% the easiest /ay to list all out)ut l les o! the run is to ty)e the !ollo/ing Linux "ommand=

```
ls *exampleRun*
```

RaxML_info.exampleRun: "ontains in!ormation about the model and algorithm used and ho/ RAxML /as "alled. The l nal (AMMA+based likelihoods8 as /ell as the al)ha sha)e)arameter6s8 are)rinted to this l le. n addition% ! the rearrangement setting /as determined automati"ally 6-i has not been used8 the rearrangement setting !ound by the)rogram /ill be indi"ated.

This is the most important output file be"ause it tells you /hat RAxML did and is al/ays /ritten irres)e"tive o! the "ommand line o)tion. n addition it "ontains in!ormation about all other out)ut l les that /ere /ritten by your run.

RAxML_log.exampleRun: A l le that)rints out the time% likelihood value o! the "urrent tree and number o! the "he"k)oint l le 6! the use o! "he"k)oints has been s)e"il ed8 alter ea" h iteration o! the sear" h algorithm. n the last line it also "ontains the l nal likelihood value o! the l nal tree to)ology. This l le is not /ritten i! multi)le bootstra)s are exe"uted% i.e. +# and -b have been s)e"il ed. n "ase o! a multi)le in!eren"e on the original alignment 6-# o)tion8 the Log+2iles are numbered a""ordingly.

RAxML_result.exampleRun: Oontains the l nal tree to)ology o! the "urrent run. This l le is also /ritten alter ea" h iteration o! the sear" h algorithm% su" h that you "an restart your run /ith -t in "ase your "om)uter "rashed. This l le is not /ritten i! multi)le bootstra)s are exe"uted% i.e. -# and -b have been s)e"il ed.

RAxML_parsimonyTree.exampleRun: "ontains the randomi7ed)arsimony starting tree i! the)rogram has not been)rovided a starting tree by -t. Ho/ever% this l le /ill not be /ritten i! a multi)le bootstra) is exe"uted using the -# and -b o)tions.

RAxML_randomTree.exampleRun: "ontains the "om)etely random starting tree i! the)rogram /as exe"uted /ith -d.

RAxML_checkpoint.exampleRun.checkpointNumber: ;rinted i! you s)e"il ed by -j that "he"k)oints shall be /ritten. Ohe"k)oints are numbered !rom \$ to n /here n is the number o! iterations o! the sear" h algorithm. Moreover% the "he"k)oint l les are additionally numbered i! a multi)le in!eren"e on the original alignment has been s)e"il ed using -#. *riting o! "he"k)oint l les is disabled /hen a multi)le bootstra) is exe"uted.

RAxML_bootstrap.exampleRun: ! a multi)le bootstra) is exe"uted by -# and -b or -x all l nal

bootstrapped trees will be written to this one single file.

`RAXML_bipartitions.exampleRun`: If you used the `-f b` option this file will contain the input tree with "original" values from \$ to &\$\$ drawn on its nodes. It is also printed when `-f a -x` have been specified at the end of the analysis the program will draw the RS support values on the best tree found during the ML search.

`RAXML_bipartitionsBranchLabels.exampleRun`: Contains the same information as the file above but support values are "correctly" displayed as Greek branch labels and not node labels. **Support values always refer to branches/splits of trees and never to nodes of the tree.** Note that some tree viewers have problems displaying branch labels, they are however part of the standard Greek format.

`RAXML_bipartitionFrequencies.exampleRun`: Contains the pairwise bipartition frequencies of all trees contained in files passed via `-t` and `-z` when the `-f m` option has been used.

`RAXML_perSiteLLs.exampleRun`: Contains the per-site log likelihood scores in TreeU77le format for usage with O39SEL (<http://www.is.titech.ac.jp/~shimo/rog/onsel>). This file is only printed when `-f g` is specified.

`RAXML_bestTree.exampleRun`: Contains the best-scoring ML tree of a thorough ML analysis.

`RAXML_distances.exampleRun`: Contains the pairwise ML-based distances between all taxa-pairs in the alignment. This file is only printed when the `-f x` option is used.

X. Computing TC and IC values

In the following I provide some more detailed examples for computing the TC and TO metrics from SaliHos and Rokas 2008 (<http://www.ncbi.nlm.nih.gov/pubmed/2815288>).

The method is described in more detail in the following paper: <http://mbe.oxfordjournals.org/content/early/2010/02/25/mbe.msb.2010.001111> (Dreier et al. 2010).

Also note that as of version 8.2.1, RAXML now performs TC calculations on "colletions" of partial gene trees. The "corrections" for partial gene trees are described in this paper here: <http://dx.doi.org/10.1093/molbev/msu011>.

Warning: The default TO calculations are not done exactly as described in our MRE paper. To have RAXML do exactly what is described in the paper, please edit file `bipartitionList.c` by commenting out or removing the line `#define BIP_FILTER`.

The way it is implemented in RAXML by default and where it differs from the version described in the paper regards the "calculation" of the OA and TOA scores. When there are several "contending" bipartitions, the standard method takes all of them into account for computing the OA score as long as their frequency exceeds a frequency threshold of 0.5.

It does however not account for the fact that some of these "contending" bipartitions, with respect to the reference bipartition, do not "contend" with each other but are in fact "compatible" with each other. Thus, by default RAXML only uses those bipartitions that "contend" with the reference bipartition that are also mutually compatible with each other to calculate the OA score. Thus, less bipartitions will typically be used to calculate the OA and consequently the TOA scores. My personal opinion is that this is more reasonable since we are not counting "contests" several times in case these "contests" emerge from compatible bipartitions.

(given a set of gene trees) RAxML "can directly" calculate a majority rule consensus MR in RAxML terminology as well as an extended majority rule consensus tree MRE in RAxML terminology on this set that has every internode that is internal branch annotated by their respective O and OA scores.

For instance to compute the O and OA scores for a given set of gene trees on a MRO tree you could type=

```
raxmlHPC -L MR -z 1070_yeast_genetrees.tre -m GTRCAT -n T1
```

/here

```
-L MR specifies that the scores will be displayed on a MR tree that is computed by RAxML
-z 1070_yeast_genetrees.tre specifies the filename that contains the set of gene trees
6/hi" are the maximum likelihood trees from the &%X$ yeast genes analyzed by Sali"hos%
and Rokas 2$&B% and /hi" are provided as supplementary data to this manuscript
-m GTRCAT is an arbitrary substitution model (this will have no effect whatsoever% but is
required as input to RAxML
-n T1 is the run . that is appended to output files.
```

RAxML will automatically build the MR tree% annotate it with the O and OA scores% and report both in an output file named RAxML_MajorityRuleConsensusTree_IC.T1% /hi" will look like this=

```
(Scer,Spar,(Smik,(Skud,(Sbay,(Scas,(Cgla,(Kpol,(Zrou,((Clus,((Psti,((Ctro,
(Calb,Cdub):1.0[0.95,0.95]):1.0[0.77,0.77],
(Cpar,Lelo):1.0[0.76,0.76]):1.0[0.75,0.75]):1.0[0.11,0.11],
(Cgui,Dhan):1.0[0.02,0.07]):1.0[0.02,0.08]):1.0[0.97,0.97],((Sklu,
(Kwal,Kthe):1.0[0.97,0.97]):1.0[0.32,0.23],
(Agos,Klac):1.0[0.08,0.08]):1.0[0.04,0.10]):1.0[0.59,0.47]):1.0[0.02,0.02]):1.0[0.1
1,0.11]):1.0[0.02,0.02]):1.0[0.97,0.97]):1.0[0.05,0.14]):1.0[0.30,0.27]):1.0[0.54,0
.54]);
```

For each internode or internal branch of the "constructed MR tree" RAxML will assign an length[x,y] branch label% /here length corresponds to the branch length (because this is a MRE tree% all internal branch lengths have been arbitrarily set to 1.0 by default% x corresponds to the O score and y to the OA score.

RAxML will also "calculate" the TO and TOA scores for the MR tree% as well as the relative TO and TOA scores that are normalized by the maximum possible TO and TOA scores for a fully bifurcating tree from the same number of taxa.

The scores are displayed in the terminal output and in the RAxML_info.runID standard output file associated with the run (in this case RAxML_info.T1) and will look like this=

```
Tree certainty for this tree: 7.642240
Relative tree certainty for this tree: 0.382112
```

```
Tree certainty including all conflicting bipartitions (TCA) for this tree: 7.580023
Relative tree certainty including all conflicting bipartitions (TCA) for this tree: 0.379001
```

(given a set of gene trees) RAxML "can also directly" calculate an extended MR tree on this set that has every internode that is internal branch annotated by their respective O and OA scores. The "artificially" computed+intensive inference of extended MRO trees (finding the optimal extended MRO tree is in fact% 9; +hard[; hillis% and * arno/ &QQ' 8 relies on RAxML's fast parallel implementation. Thus if you use the ; Threads version of RAxML% this part will run in parallel.

To "om)ute O% OA% TO% and TOA s"ores on an extended MR tree you /ould ty)e=

```
raxmlHPC -L MRE -z 1070_yeast_genetrees.tre -m GTRCAT -n T2
```

RAxML "an "om)ute MR and extended MR trees% using both fully bilur"ating and)artially resolved@multiur"ating trees as an in)ut. RAxML "an also "om)ute stri"ter MR trees /ith arbitrary threshold settings that range bet/een P& and &\$\$S. 2or instan"e% by ty)ing

```
raxmlHPC -L T_75 -z 1070_yeast_genetrees.tre -m GTRCAT -n T3
```

RAxML /ill dis)lay O% OA% TO and TOA s"ores on a MR tree that only in"ludes those bi)artitions that have] XPS su))ort.

* e have also im)lemented an o)tion 6-f i8 that allo/s the user to "al"ulate and dis)lay O% OA% TO% and TOA s"ores onto a given% stri"tly bilur"ating releren"e tree 6!or exam)le% the best+kno/n ML tree8.

This is analogous to the standard -f b o)tion in RAxML that dra/s bootstra) su))ort values !rom a set o! bootstra) trees onto a releren"e)hylogeny. The o)tion "an be invoked by ty)ing

```
raxmlHPC -f i -t yeast_concatenationtree.tre -z 1070_yeast_genetrees.tre -m GTRCAT -n T4
```

Note that% the tree "ontained in l le yeast_concatenationtree.tre needs to be stri"tly bilur"ating and "ontain bran"h lengths. n this exam)le% the yeast_concatenationtree.tre l le is the best+kno/n maximum likelihood tree re"overed by "on"atenation analysis o! the &\$\$X\$ yeast genes !rom the)a)er.

, sing this "ommand% RAxML /ill annotate the tree in yeast_concatenationtree.tre /ith the O and OA s"ores% and re)ort both in an out)ut l le named RAxML_IC_Score_BranchLabels.T4% /hi"n /ill look like this=

```
(((((Clus:0.47168135428609103688((((Lelo:0.30356174702769450624,Cpar:0.2549087423
9480920682):0.13023178275857649755[0.76,0.76],
(Ctro:0.18383414558272206940(Calb:0.04124660275465741321,Cdub:0.0429080158839683228
9):0.14526604486383792869[0.95,0.95]):0.12355825028654655873[0.77,0.77]):0.17335821
030783615804[0.75,0.75],Psti:0.42255112174261910685):0.07862882822310976461[0.11,0.
11],
(Cgui:0.45961028886034632768,Dhan:0.28259245937168109286):0.05586015476156453580[0.
02,0.07]):0.08116340505230199009[0.02,0.08]):1.03598510402913923656[0.97,0.97],
((Agos:0.53332956655591512440,Klac:0.47072785596320687596):0.08132006357704427146[0.
08,0.08],
((Kthe:0.17123899487739652203,Kwal:0.17320923240031221857):0.25620117495110567019[0.
97,0.97],Sklu:0.24833228915799765435):0.05646992617871094550[0.32,0.23]):0.0523630
6187235122145[0.04,0.10]):0.10686517691208799463[0.59,0.47],Zrou:0.4130783368556378
2877):0.03792570537296727218[0.02,0.02],Kpol:0.43287284049576529865):0.045603416931
36910068[0.11,0.11],Cgla:0.49584136365135367264):0.04363310339731014259[0.02,0.02],
Scas:0.37212829744050218705):0.29362133996280515014[0.97,0.97],
(Skud:0.06926467973344750673,(Smik:0.06535810850036427588,
(Scer:0.04285848856634000975,Spar:0.03030513540244994877):0.02506719066056842596[0.
54,0.54]):0.02459323291555862850[0.30,0.27]):0.02524223867026276907[0.05,0.14],Sbay
:0.06506923220637816918);
```

2or ea"n internode or internal bran"n o! this out)ut tree RAxML /ill assign a length[x,y] bran"n label% /here length "orres)onds to the original bran"n length in the in)ut tree)assed via -t% x "orres)onds to the O s"ore and y to the OA s"ore. RAxML /ill also dis)lay the TO and TOA s"ores o! this tree both in the terminal out)ut and in the RAxML_info.T4 out)ut l le asso"iated /ith the run.

It should further be noted that the O and OA scores are represented as branch labels since as is the case for bootstrap support values information associated to internodes splits/bifurcations of a tree always refers to branches and not nodes.

Each tree viewer (e.g. Endros) or <http://ab.inl.uni-tuebingen.de/soft/areadendros/> enables that "an" (possibly) parse the Newick tree format is able to display these branch labels.

The rationale for not providing O and OA scores as node labels is that some viewers may not (possibly) rotate the node labels when the tree is re-rooted by the user which will lead to an erroneous branch to O and branch to OA score association.

* When calculating O and OA scores on extended MR trees or when drawing O and OA scores onto a given reference tree it may occur that the bifurcation that has been included in the tree has more support than one or more "non-tipping" bifurcations. In this case RAxML will re-root O and OA scores on the inferred tree with negative signs.

Finally we have implemented a verbose output option that allows users to further scrutinize particularly interesting "non-tipping" bifurcations.

Verbose mode is activated by adding the -C command line switch to any of the above examples.

In verbose mode RAxML will generate two additional types of output files: One set of files containing one included bifurcation and the corresponding "non-tipping" bifurcations in Newick format called RAxML_verboseIC.runID.0 ... RAxML_verboseIC.runID.N-1 where N is the number of bifurcations in the tree and an output file that lists all bifurcations included and "non-tipping" in a ;H<L;+like format called RAxML_verboseSplits.runID.

For example by adding -C to the previous command

```
raxmlHPC -f i -t yeast_concatenationtree.tre -z 1070_yeast_genetrees.tre -m GTRCAT -n T5 -C
```

will produce 2\$ files: One for each of the 2\$ bifurcations present in the yeast_concatenationtree.tre named RAxML_verboseIC.T5.0, RAxML_verboseIC.T5.1, ..., RAxML_verboseIC.T5.19

For example the RAxML_verboseIC.T5.0 file will look like this=

```
((Cpar, Lelo),(Scer, Smik, Skud, Cgla, Kpol, Zrou, Kwal, Kthe, Agos, Klac, Clus, Cgui, Psti, Ctro, Calb, Cdub, Dhan, Sklu, Scas, Sbay, Spar));
((Cpar, Ctro, Calb, Cdub),(Scer, Smik, Skud, Cgla, Kpol, Zrou, Kwal, Kthe, Agos, Klac, Clus, Cgui, Psti, Lelo, Dhan, Sklu, Scas, Sbay, Spar));
```

Here the first Newick string represents the bifurcation that was included in the yeast_concatenationtree.tre and all following Newick strings represent the corresponding "non-tipping" bifurcations in descending order of their frequency of occurrence in the case of the RAxML_verboseIC.T5.0 file the first bifurcation which is included in the yeast_concatenationtree.tre "non-tips" with only one other bifurcation which is listed as the second bifurcation.

Analogously the output file that lists all bifurcations included and "non-tipping" in a ;H<L;+like format RAxML_verboseSplits.T5 looks like this=

1. Scer
2. Smik
3. Skud
4. Cgla
5. Kpol
6. Zrou

7. Kwal
 8. Kthe
 9. Agos
 10. Klac
 11. Clus
 12. Cgui
 13. Psti
 14. Cpar
 15. Lelo
 16. Ctro
 17. Calb
 18. Cdub
 19. Dhan
 20. Sklu
 21. Scas
 22. Sbay
 23. Spar

partition:

-----	-----	----**	-----	---	956/89.345794/0.761406
-----	-----	----*-	***--	---	39/3.644860/0.761406

partition:

-----	-----	-----	-**--	---	1051/98.224299/0.949483
-----	-----	-----	**---	---	6/0.560748/0.949483

.
 .
 .

partition:

--***	*****	*****	*****	**--	641/59.906542/0.303620
-**--	-----	-----	-----	-*-	148/13.831776/0.303620
-*--*	*****	*****	*****	*--	114/10.654206/0.303620

partition:

--***	*****	*****	*****	**--	825/77.102804/0.545775
-**--	-----	-----	-----	-**	87/8.130841/0.545775

bestW scoring ML tree on the original alignment.

To just do a RS search you could type=

```
raxmlHPC -x 12345 -p 12345 -# 100 -m GTRGAMMA -s ex_al -n TEST
or using the bootstrap "convergence" criterion=
```

```
raxmlHPC -x 12345 -p 12345 -# autoMRE -m GTRGAMMA -s ex_al -n TEST
```

This will "ondu"t rapid bootstraping and an ML search under the (AMMA model of rate heterogeneity.

So if you want to run a full analysis i.e. RS and ML search type=

```
raxmlHPC -f a -x 12345 -p 12345 -# 100 -m GTRGAMMA -s ex_al -n TEST
```

This will first "ondu"t a RS search and once that is done a search for the bestW scoring ML tree. Such a program run will return the bootstrapped trees 6RAXML_bootstrap.TEST8 the best scoring ML tree 6RAXML_bestTree.TEST8 and the RS support values drawn on the best scoring tree as node labels 6RAXML_bipartitions.TEST8 as well as more "correctly" sorted values refer to branches as branch labels 6RAXML_bipartitionsBranchLabels.TEST8.

Finally note that by increasing the number of RS replicates via -# you will also make the ML search more thorough since for ML optimization every Pth RS tree is used as a starting point to search for ML trees.

* When -# autoMRE is specified RAXML will execute a maximum of &\$\$\$ RS replicate searches but it may of course "converge" earlier.

From what we have observed so far this new ML search algorithm yielded better trees than what is obtained via 2\$ standard ML searches on distinct starting trees for all datasets with ^ &\$\$\$ sequences.

For larger datasets it might be worthwhile to "ondu"t an additional ML search as described below. Just to be sure.

Warning: note that the rapid RS search will "urrently" ignore "ommands" associated to user tree files passed via -t or -z.

However the "onstraint and backbone tree options 6-g and -r8 do work with rapid RS.

The Hard & Slow Way

Describe the observation that the default parameters and the rapid RS and ML algorithm described above work well in most "rational" cases a good thing to do is to adapt the program parameters to your alignment.

This refers to a *good* setting for the rate categories of -m GTRCAT and the initial rearrangement setting.

If you use partitioned models you should add -q partitionFileName to all of the following "ommands.

Getting the Initial Rearrangement Setting right

! you don't see an initial rearrangement setting with the `-i` option the program will automatically determine a good setting based upon the randomized M₁ starting tree. It will take the starting tree and apply many subtree rearrangements with a rearrangement setting of `P%&P%2%2P`. The minimum setting that yields the best likelihood improvement on the starting trees will be used as initial rearrangement setting.

This procedure can have two disadvantages=

- Firstly the initial setting might be very high e.g. `2%` or `2P%` and the program will slow down considerably.

- Secondly a rearrangement setting that yields a high improvement of likelihood scores on the starting tree might let the program get stuck earlier in some local maximum of this behavior could already be observed on a real dataset with about 100 taxa.

Therefore you should run RAxML a couple of times the more the better with the automatic determination of the rearrangement setting and with a predefined value of `&P%` which proved to be sufficiently large and efficient in many practical cases.

In the example below we will do this based on predefined starting trees. So let's first generate a couple of randomized M₁ starting trees.

Note that in RAxML you also always have to specify a substitution model regardless of whether you only want to compute an M₁ starting tree with the `-y` option. Note that we have to pass different random number seeds via `-p` to obtain distinct starting trees here

```
raxmlHPC -y -p 12345 -s ex_al -m GTRCAT -n ST0
...
raxmlHPC -y -p 34556 -s ex_al -m GTRCAT -n ST4
```

Then infer the ML trees for those starting trees using a fixed setting `-i 10`

```
raxmlHPC -f d -i 10 -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST0 -n FI0
...
raxmlHPC -f d -i 10 -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST4 -n FI4
```

and then using the automatically determined setting on the same starting trees=

```
raxmlHPC -f d -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST0 -n AI0
...
raxmlHPC -f d -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST4 -n AI4
```

The setting that yields the best final likelihood scores as automatically computed under the (AMMA model of rate heterogeneity should then be used for subsequent analyses.

Getting the Number of Categories right

Another issue is to get the number of rate categories right. Due to the reduced memory footprint and significantly reduced inference times the recommended model to use with RAxML on large dataset is (TROT if you are doing runs to find the best known ML tree on the original alignment and for bootstrapping).

Thus% you should experiment with a "ou)le o! -c settings and then look /hi" h gives you the best (amma+based likelihood value.

Su))ose that in the)revious se"tion you found that automati"ally determining the rearrangement setting /orks best for your alignment.

<ou should then re+run the analyses /ith distin"t + settings by in"rements o! e.g. &P rate "ategories

```
raxmlHPC -f d -c 10 -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST0 -n C10_0
```

```
...
```

```
raxmlHPC -f d -c 10 -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST4 -n C10_4
```

<ou don:t need to run it /ith the de!ault setting o! -c 25 sin"e you already have that data% su" h that you

"an "ontinue /ith ...

```
raxmlHPC -f d -c 40 -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST0 -n C40_0
```

```
...
```

```
raxmlHPC -f d -c 40 -m GTRCAT -s ex_al -t RAxML_parsimonyTree.ST4 -n C40_4
```

and so on and so forth.

Sin"e the (TROAT a))roximation is still a ne/ "on"e)t little is kno/n about the a))ro)riate setting for -c 25.

Ho/ever% em)iri"ally -c 25 /orked best on &Q real+/orl'd alignments. So testing u) to -c 55 should usually be suL "ient% ex"e)t i! you noti"e a tenden"y for l nal (TR(AMMA likelihood values to further im)rove /ith in"reasing rate "ategory number.

Thus% the assessment o! the *good* -c setting should on"e again be based on the l nal (TR(AMMA likelihood values.

! you don:t have the time or "om)utational)o/er to determine both *good* -c and -i settings you should rather sti"k to determining -i sin"e it has sho/n to have a greater im)a"t on the l nal results.

Also note% that in"reasing the number o! distin"t rate "ategories has a negative im)a"t on exe"ution times. 2inally% i! the runs /ith the automati" determination o! the rearrangement settings from the)revious Se"tion have yielded the best results you should then use exa"tly the same rearrangement settings for ea" h series o! ex)eriments to determine a *good* -c setting.

The automati"ally determined rearrangement settings "an be retrieved from files RAxML_info.AI_0 ... RAxML_info.AI_4

Finding the Best-Known Likelihood tree (BKL)

As already mentioned RAxML uses randomi7ed ste)/ise addition order)arsimony starting trees on /hi" h it then initiates an ML+based o)timi7ation. Those trees are obtained by using a randomi7ed ste)/ise addition se4uen"e to insert one taxon alter the other into the tree. * hen all se4uen"es have been inserted a "ou)le o! subtree rearrangements 6also "alled subtree)uning re+grafting8 /ith a l xed rearrangement distan"e o! 2\$ are exe"uted to further im)rove the M; s"ore.

The "on"e)t to use randomi7ed M; starting trees in "ontrast to the 9# 69eighbor #oining8 starting trees many other ML)ograms use% is regarded as an advantage o! RAxML. This allo/s the)rogram to start ML o)timi7ations o! the to)ology from a distin"t starting)oint in the immense to)ologi"al sear" h s)a"e ea" h time.

There!ore% RAxML is more likely to l nd good ML trees i! exe"uted several times.

This also allows you to build a consensus tree out of the individual tree topologies obtained from each individual run on the original alignment. By this and by combining the individual likelihoods you can get a feeling on how stable the tree is to get caught in local maxima as the search algorithm is on the original alignment.

You can also use the `-f r` option to combine pairwise topological Robinson+Foulds distances between the ML trees you have found.

Thus if you have several computing resources available in addition to bootstrapping you should do multiple inferences by executing 255 inferences in some recent real-world analyses with RaxML on the original alignment.

On smaller datasets or large phylogenomic datasets it will also be worthwhile to use the `-d` option for a couple of runs to see how the program behaves on completely random starting trees. This is where the option as well as the parallel MPI version `raxmlH; MPI; v` come into play.

So to execute a multiple inference on the original alignment on a single processor type=

```
raxmlHPC -f d -p 12345 -m GTRCAT -s ex_al -# 10 -n MultipleOriginal
```

If you have a cluster available you could speedily=

```
raxmlHPC-MPI -f d -m GTRCAT -p 12345 -s ex_al -# 100 -n MultipleOriginal
```

replaced by the respective MPI run-time commands e.g. `mpiexec` or `mpirun` depending on your local installation please check with your local computer scientist.

Bootstrapping with RAXML

To carry out a multiple non-parametric bootstrap with the sequential version of RAXML type=

```
raxmlHPC -f d -m GTRCAT -s ex_al -p 12345 -# 100 -b 12345 -n MultipleBootstrap
```

You have to speedily a random number seed alter `-b` for the random number generator of the bootstraps and `-p` for the random number generator of the parsimony starting trees. This will allow you to generate reproducible results.

To do a parallel bootstrap type=

```
raxmlHPC-MPI -f d -m GTRCAT -s ex_al -# 100 -p 12345 -b 12345 -n MultipleBootstrap
```

once again replaced by the appropriate MPI execution command.

Obtaining Confidence Values

Suppose that you have executed 255 inferences on the original alignment and 1000 bootstrapped runs. You can now use the RAXML `-f b` option to draw the information from the 1000 bootstrapped topologies onto some tree and obtain a topology with support values. From my point of view the most reasonable thing to do is to draw them on the best-scoring ML tree from those 255 runs. Suppose that the best-scoring tree was found in run number 99 and the respective tree file is `RAxML_result.MultipleOriginal.RUN.99`.

If you have executed more than one bootstrap runs with the sequential version of RAXML on distinct computers i.e. 10 runs with 1000 bootstraps on 10 machines you will first have to concatenate the bootstrap files. If your bootstrap result files are named e.g. `RAxML_bootstrap.MultipleBootstrap.0...RAxML_bootstrap.MultipleBootstrap.9` you can

easily concatenate them by using the `cat` command e.g.

```
cat RAxML_bootstrap.MultipleBootstrap.* > RAxML_bootstrap.All
```

In order to get a tree with bootstrap values on it must execute RAxML as indicated below=

```
raxmlHPC -f b -m GTRCAT -z RaxML_bootstrap.All -t RAxML_result.MultipleOriginal.RUN.99 -n BS_TREE
```

This will return the tree passed via `-t` annotated by support values either as branch labels or as node labels. Some tree viewers have problems with displaying trees with node labels in particular if they are used to re-root the tree hence you should use the branch-labeled tree if possible.

You can also compute a majority rule consensus tree out of the bootstrap replicates=

```
raxmlHPC -J MR -m GTRCAT -z RaxML_bootstrap.All -n MR_CONS
```

an extended majority rule consensus tree=

```
raxmlHPC -J MRE -m GTRCAT -z RaxML_bootstrap.All -n MRE_CONS
```

or a strict one=

```
raxmlHPC -J STRICT -m GTRCAT -z RaxML_bootstrap.All -n STRICT_CONS
```

For finding rogue taxa you can use the `-J STRICT_DROP` or `-J MR_DROP` options. Alternatively you can use the web-server implemented by my lab= <http://cbrnr.hiits.org/cbrnr>

XII. A Simple Heterotachous Model

We implemented a simple heterotachous model in RAxML by assigning one (TR) model to terminal branches of the tree and another distinct (TR) model to the inner branches of the tree. The idea was to better capture early rapid radiations with this model but it didn't work well.

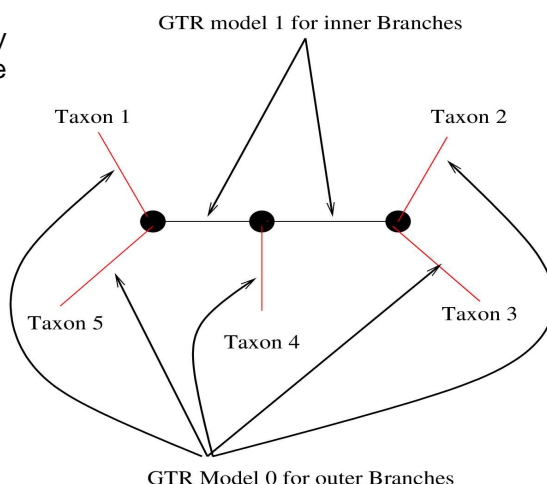
All other parameters such as the alpha shape parameter of the gamma model of rate heterogeneity and the base frequencies are shared across the entire tree. The model parameters of the two (TR) models on the tree are optimized with respect to the overall likelihood of the tree using standard numerical optimization procedures.

This heterotachous model is currently only available in the sequential version of RAxML and only for analyzing 9A sequence data under the gamma model of rate heterogeneity.

There is no command line switch for enabling this model. Instead you will need to recompile RAxML by adding `-D_HET` to the line starting by `CFLAGS =` in the respective Makefile. Then when you specify `-m GTRGAMMA` the search and all other likelihood calculations will be conducted under this heterotachous model.

Keep in mind that when doing model tests etc. having 2 (TR) matrices per partition introduces a total of 24 free parameters for the substitution matrix instead of 12 for a single (TR) matrix.

Our simple heterotachous model is outlined in the figure below.



XIII. Frequently Asked Questions

- Q:** When performing a bootstrap search using a partitioned model does RAxML perform a "conserved bootstrap resampling" i.e. does it resample within genes so that partitions are sustained?
- A:** That is the case. When performing bootstraps on partitioned data sets, bootstrapped alignments will be sampled from within partitions i.e. bootstrapped partitions are sustained and contain exactly the same number of alignment columns as the original partition.
- Q:** Can we use GEX, S-site in our analyses with RAxML?
- A:** Not directly, but my colleague Frank Zou from the University of Zuerich has written a Perl script called `<RAxML2.pl>`. This is a script that reads nexus data files and prepares the necessary input files and command-line options for RAxML. You can download it at <http://www.lut7onilab.net/download>
- Q:** Why don't you like the proportion of nvariable 6; + nvar 8 Sites estimate despite the fact that you implemented it?
- A:** Only implemented 6; + nvar in RAxML to make some users happy but still strongly disagree with its usage.

Personal opinion: It is unquestionable that one needs to incorporate rate heterogeneity in order to obtain *publishable* results. Aside from the *publish-or-perish* argument, there is also strong biological evidence for rate heterogeneity among sites. The rationale for being skeptical about 6; + nvar in RAxML is that all three alternatives (TR, GAMMA, and TROAT) and 6; + nvar represent distinct approaches to incorporate rate heterogeneity. Thus, in principle they account for the same phenomenon by different mathematical means. Also some unpublished concerns have been raised that the usage of 6; + nvar in combination with gamma can lead to a *ping-pong* effect since a change of 6; + nvar leads to a change in gamma and vice versa. This essentially means that those two parameters (i.e. alpha and 6; + nvar) can not be optimized independently from each other and might cause significant trouble and problems during the model parameter everything except tree topology optimization process. I have already observed this when using 6; + nvar in RAxML on a very small AA dataset.

Although this has never been properly documented, several well-known researchers in phylogenetics share this opinion. I quote Jiheng from an email in 2008 regarding this part of the RAxML manual:

I entirely agree with your criticism of the Pinv+Gamma model, even though as you said, it is very commonly used.

Jiheng also addresses the issue in his book on Molecular Evolution (3rd edition, University of Chicago Press, 2004) from pages 110-111.

The model is known as I+G and has been widely used. This model is somewhat pathological as the gamma distribution with alpha already allows for sites with very low rates; as a result, adding a proportion of invariable sites creates a strong correlation between p_0 and alpha, making it impossible to estimate both parameters reliably.

Many analyses have so far not encountered any difficulties with reviews for the real phylogenetic analyses published with colleagues from Biology when we used the TR (AMMA model instead of the more widely spread TRU model.

- Q:** Why does RAXML only implement (TR)-based models of nucleotide substitution?
- A:** For each distinct model of nucleotide substitution RAXML uses a separate highly optimized set of likelihood functions. The idea behind this is that (TR is the most common and general model for real-world DNA analysis. Thus it is better to explicitly implement and optimize this model instead of offering a plethora of distinct models which are only special cases of (TR but are programmed in a generic and thus inefficient way.

Personal opinion: My personal view is that using a simpler model than (TR only makes sense with respect to the computational cost, i.e. it is less expensive to compute. Programs such as Modeltest promote the usage of a simpler model for a sequence alignment if the likelihood of a fixed topology under that simpler model is not significantly worse than that obtained by (TR based on a likelihood ratio test. My experience is that (TR always yields a slightly better likelihood than alternative simpler models. In addition, since RAXML has been designed for the inference of large datasets the danger of over-parameterizing such an analysis is comparatively low. Provided these arguments the design decision was taken to rather implement the most general model explicitly than to provide many inefficient generic implementations of models that are just special cases of (TR. Finally, the design philosophy of RAXML is based upon the observation that a more thorough topological search has a greater impact on final tree quality than modeling details. Thus, the explicit implementation of a rapid search mechanism is considered to be more important than model details.

- Q:** Why has the performance of RAXML mainly been assessed using real-world data?
- A:** *Personal opinion:* Despite the unquestionable need for simulated data and trees to verify and test the performance of current ML algorithms the current methods available for generation of simulated alignments are not very realistic. For example, only few methods exist that incorporate the generation of gaps in simulated alignments. Since the model according to which the sequences are generated on the true tree is reproduced, we are actually assuming that ML exactly models the true evolutionary process which in reality we simply don't know how sequences evolved. The above simplifications lead to perfect alignment data without gaps that evolved exactly according to a reproduced model and thus exhibits a very strong phylogenetic signal in contrast to real data. In addition, the given true tree must not necessarily be the Maximum Likelihood tree. This difference manifests itself in substantially different behaviors of search algorithms on real and simulated data. Typically, search algorithms execute significantly less topological moves on simulated data until convergence as opposed to real data, i.e. the number of successive nearest neighbor exchanges or subtree rearrangements is lower. Moreover, in several cases the likelihood of trees found by RAXML on simulated data was better than that of the true tree. Another important observation is that program performance can be inverted by simulated data. Thus, a program that yields good topological Robinson-Woulds distances on simulated data can in fact perform much worse on real data than a program that does not perform

only on simulated data.

! one is willing to really accept ML as inference criterion on real data one must also be willing to assume that the tree with the best likelihood score is the tree that is closest to the true tree.

My personal conclusion is that there is a strong need to improve simulated data generation and methodology. In addition, the perhaps best way to assess the validity of our tree inference methods consists in an empirical evaluation of new results and insights obtained by real phylogenetic analysis.

This should be based on the prior knowledge of Biologists about the data and the method and scientific benefits attained by the simulation of phylogenies.

Q: * why am I getting weird error messages from the MrBayes version 3.2.1

A: You probably forgot to specify the -# or -N option in the command line which must be used for the MrBayes version 3.2.1 or earlier.

Q: * when using partitioned models, can I link the model parameters of different partitions to be estimated jointly, in a similar way MrBayes does it?

A: Currently not, but the implementation of such an option is planned. However, we are still lacking good criteria and methods that tell us how and why to link/unlink certain parameters across different partitions.

Q: How does RAxML handle gaps?

A: Like most other likelihood-based phylogeny programs it handles them as undetermined characters, that is gaps are treated as Ns. ! this is new to you you may want to read Joe Felsenstein's textbook *Inferring Phylogenies* where he explains why this is a reasonable way to model gaps.

Q: * why am I getting long branch lengths on my phylogenomic datasets with missing data?

A: This is normal and due to the way missing data is modeled, the missing data affects the sequence about which we don't know become very distant for exactly this reason.