



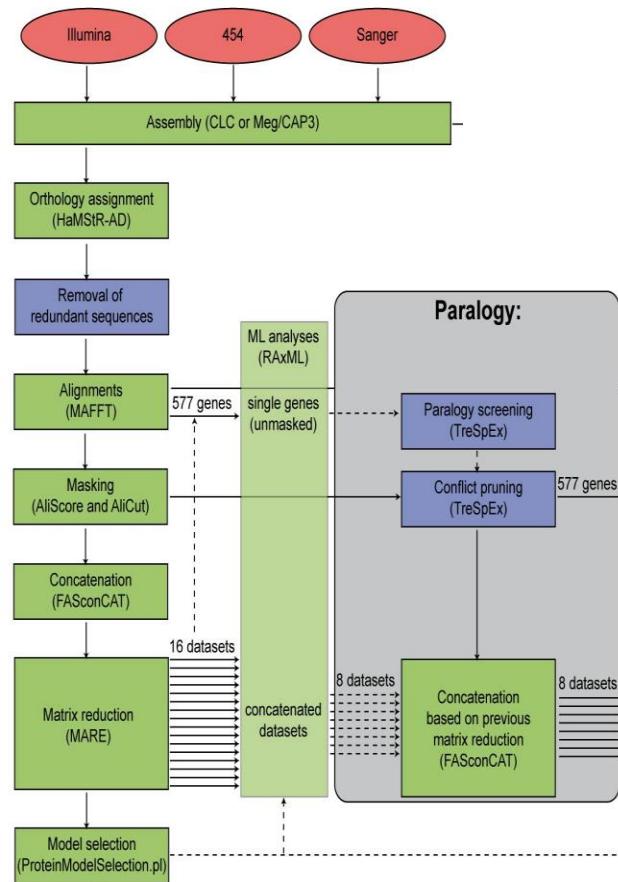
Paralogy



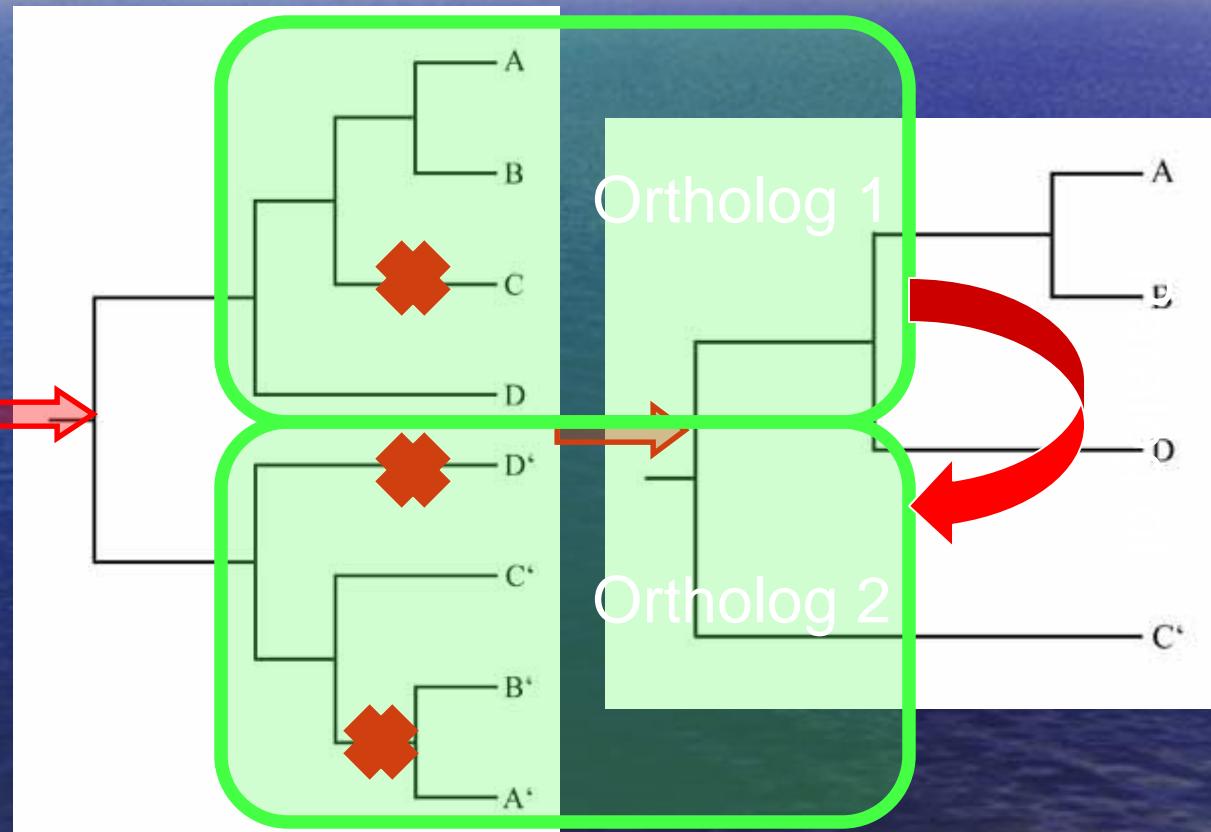
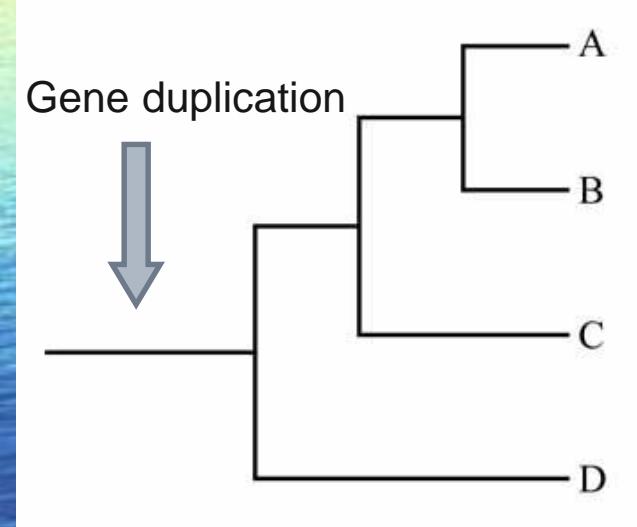
**Torsten Struck – t.h.struck@nhm.uio.no
NHM UiO**

© nhm.uio.no

Going beyond the standard



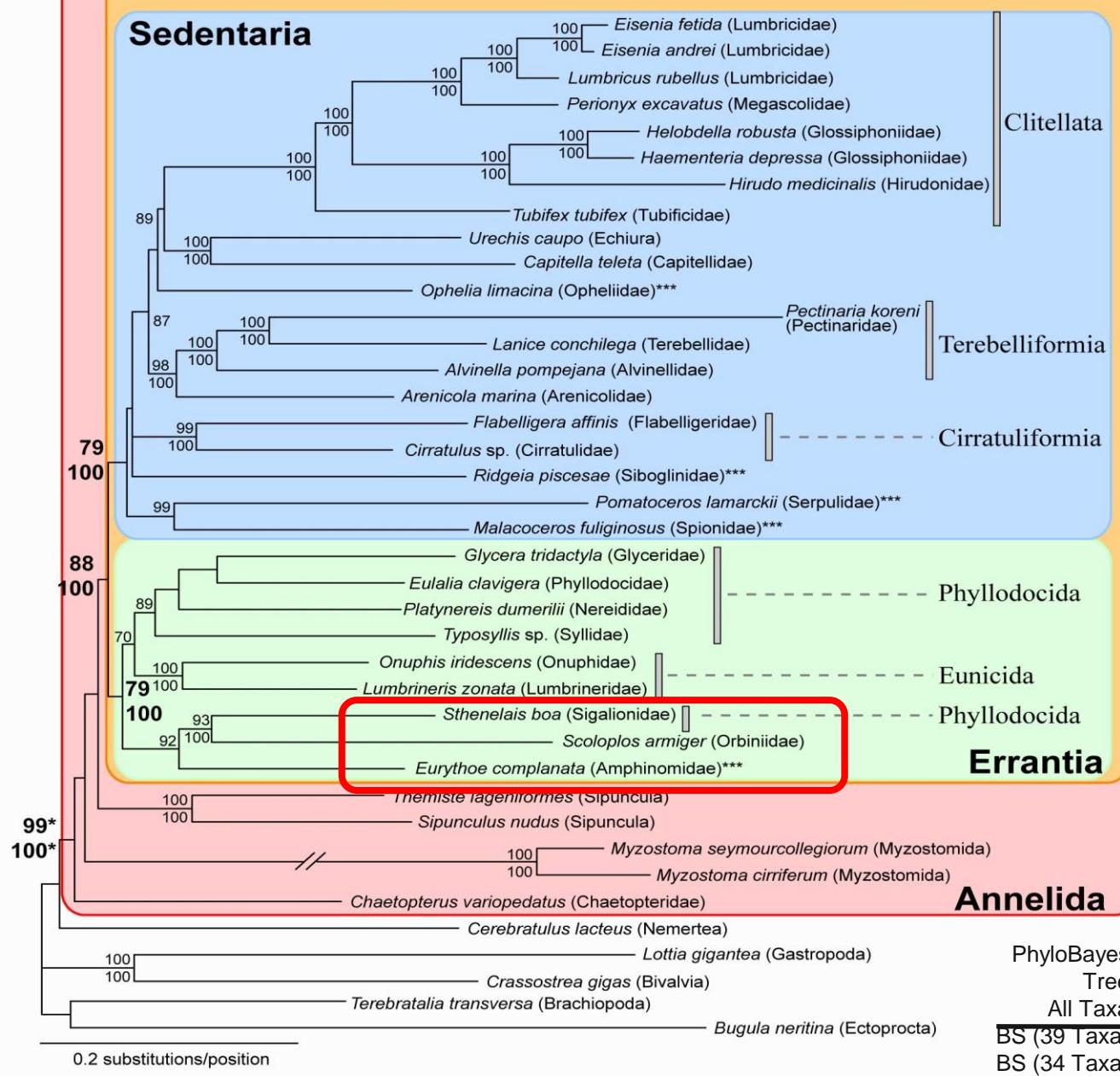
Problem of paralogy



Species tree

Gene tree

Wrong species tree

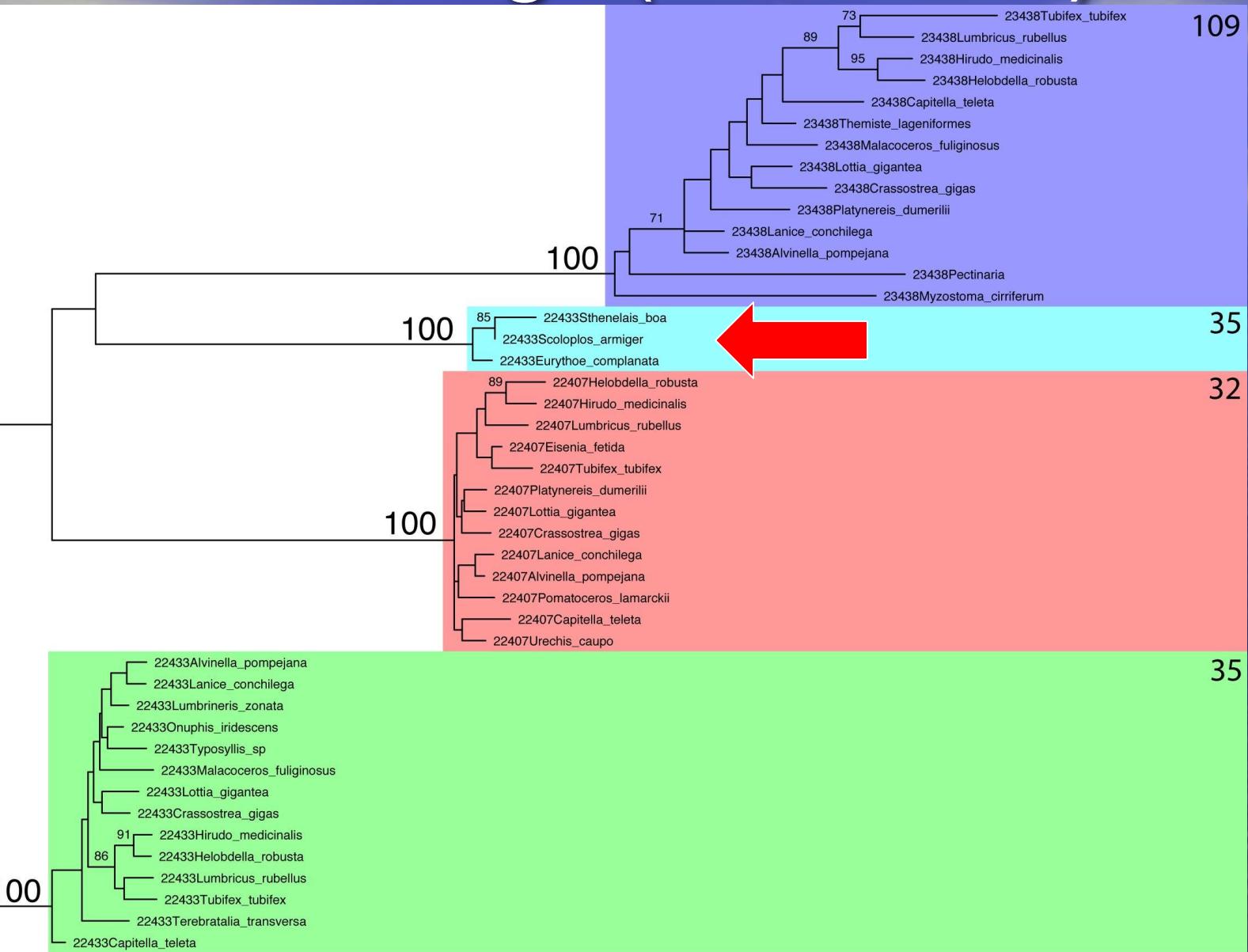
Pleistoannelida**Sedentaria****Phylogenomics**

17 new EST-Datasets
Approx. 1,300 clones
231 ortholog genes
47,952 aa-positions
41.7% coverage per taxon
20 polychaete families
[inclusive of Siboglinidae (Pogonophora)]
Clitellata, Sipuncula, Echiura, Myzostomida

PhyloBayes/RAXML
Only bootstrap values $\geq 70\%$
Sensitivity analyses
*** unstable taxa

*BS-values without Myzostomida due to long-branch issues in RAXML

Paralogs (32/109/35)



Paralogs (32/109/35)

Blast against *Bos taurus* transcriptome
= Proteasome subunit alpha

reference taxa	32 (23428)	109 (22407)	35 (22433)
Helobdella	PSMA5	PSMA3	PSMA7/8
Capitella	PSMA5	PSMA3	PSMA7/8
Lottia	PSMA5	PSMA3	PSMA7/8

Critical taxa	35 (22433)
Eurythoe	PSMA2
Scoloplos	PSMA2
Sthenelais	PSMA2

Blast of *Bos taurus* PSMA2 sequence against the *Helobdella* transcriptome:
PSMA2 (NM_001034662): *Helobdella* sequence 185484 e^-39
That is the sequence of orthology group 22433 (gene number 35)

What has happened?

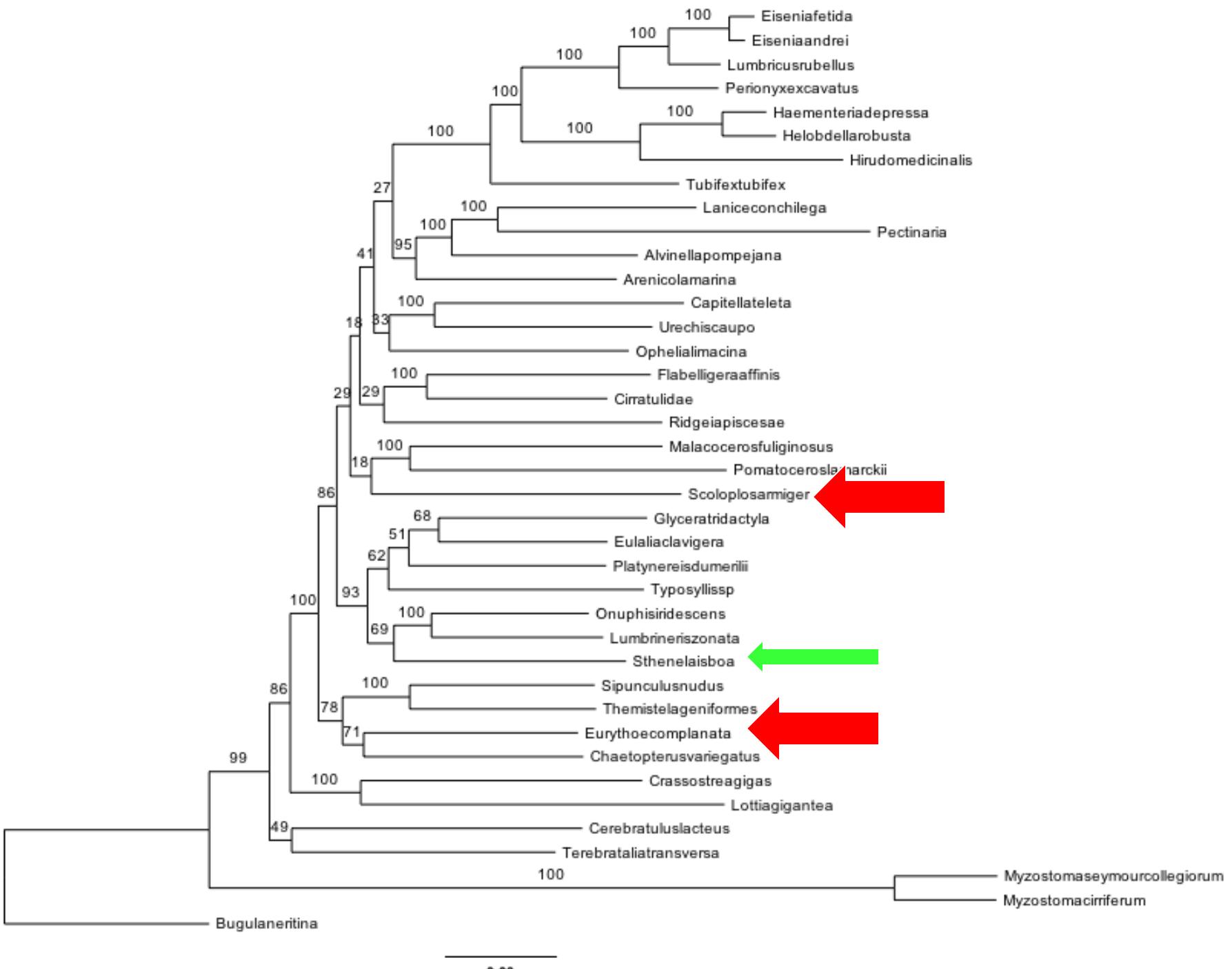
	PSMA7/8	PSMA2
Helobdella	present	absent
Eurythoe	absent	present

The table illustrates the presence (green) or absence (red) of two genes, PSMA7/8 and PSMA2, in two species, Helobdella and Eurythoe. The 'present' status is indicated by a green box, while the 'absent' status is indicated by a red box. Red circular 'no' symbols are placed over the 'absent' cells, and green arrows point from the 'present' cells to their respective 'absent' counterparts.

The problem of reciprocal lack

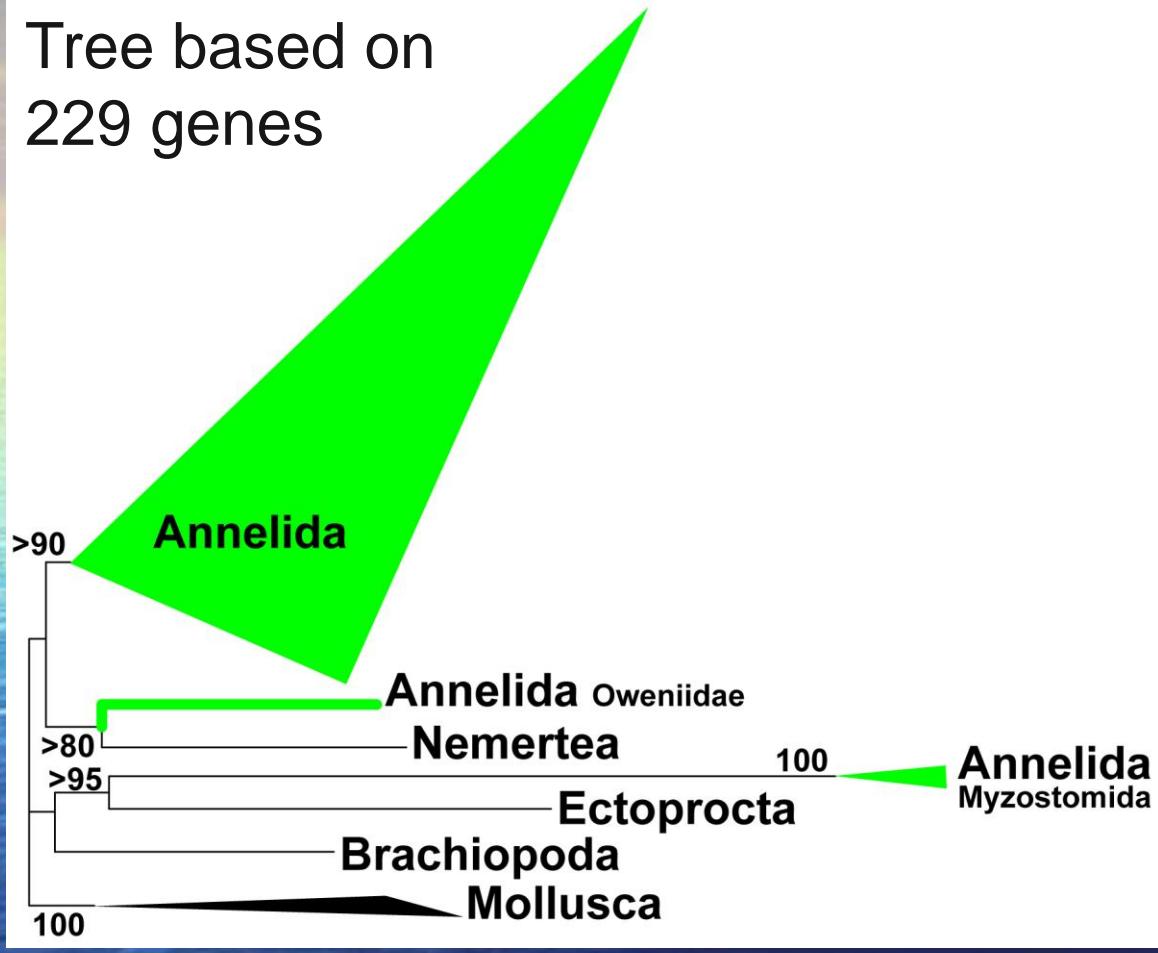
A wide-angle photograph of a calm ocean under a vast, cloudy sky. The horizon is visible in the distance, where the ocean meets a sky filled with various shades of blue, white, and yellow clouds. The lighting suggests either a sunrise or sunset, with warm tones on the left side of the frame.

Does it matter?

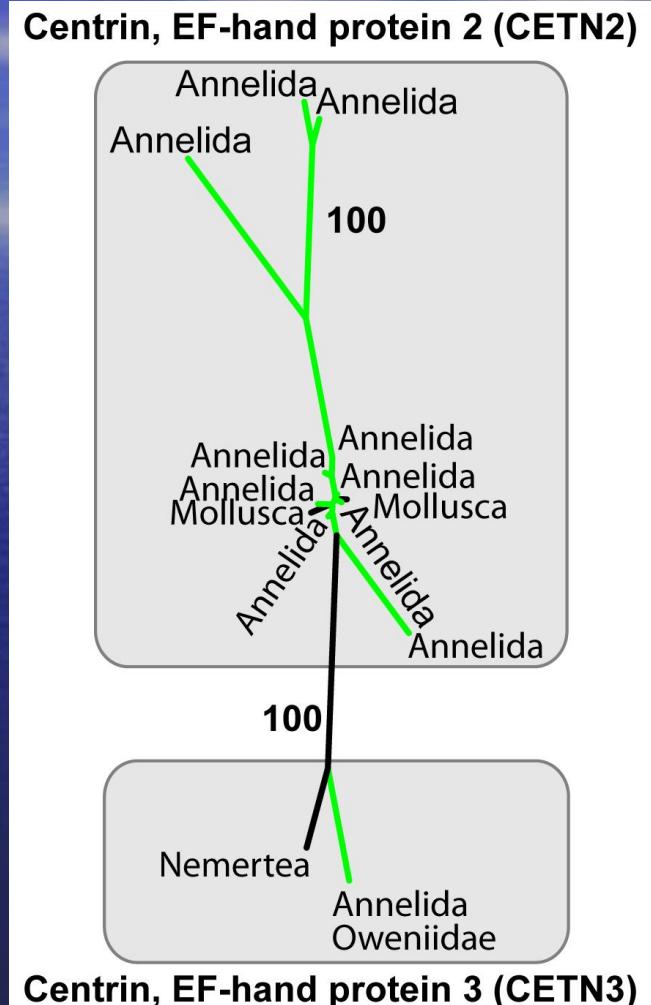


Problem of paralogy

Tree based on
229 genes

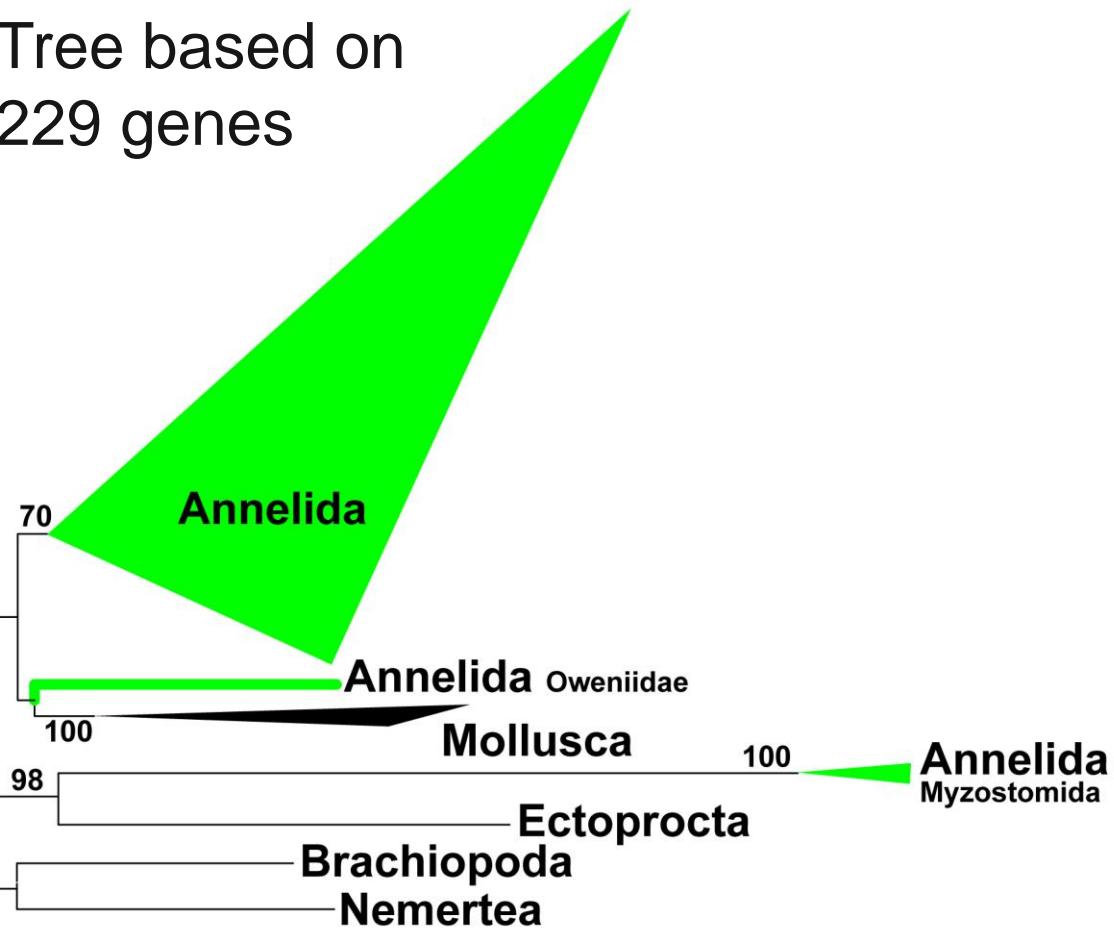


Struck (2013) PLoS one

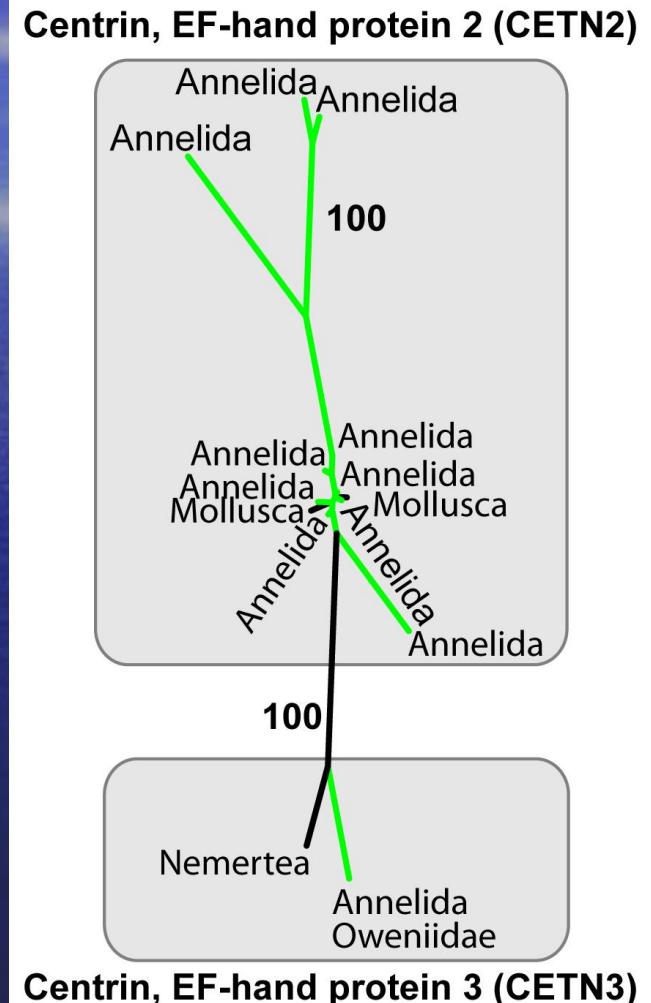


Problem of paralogy

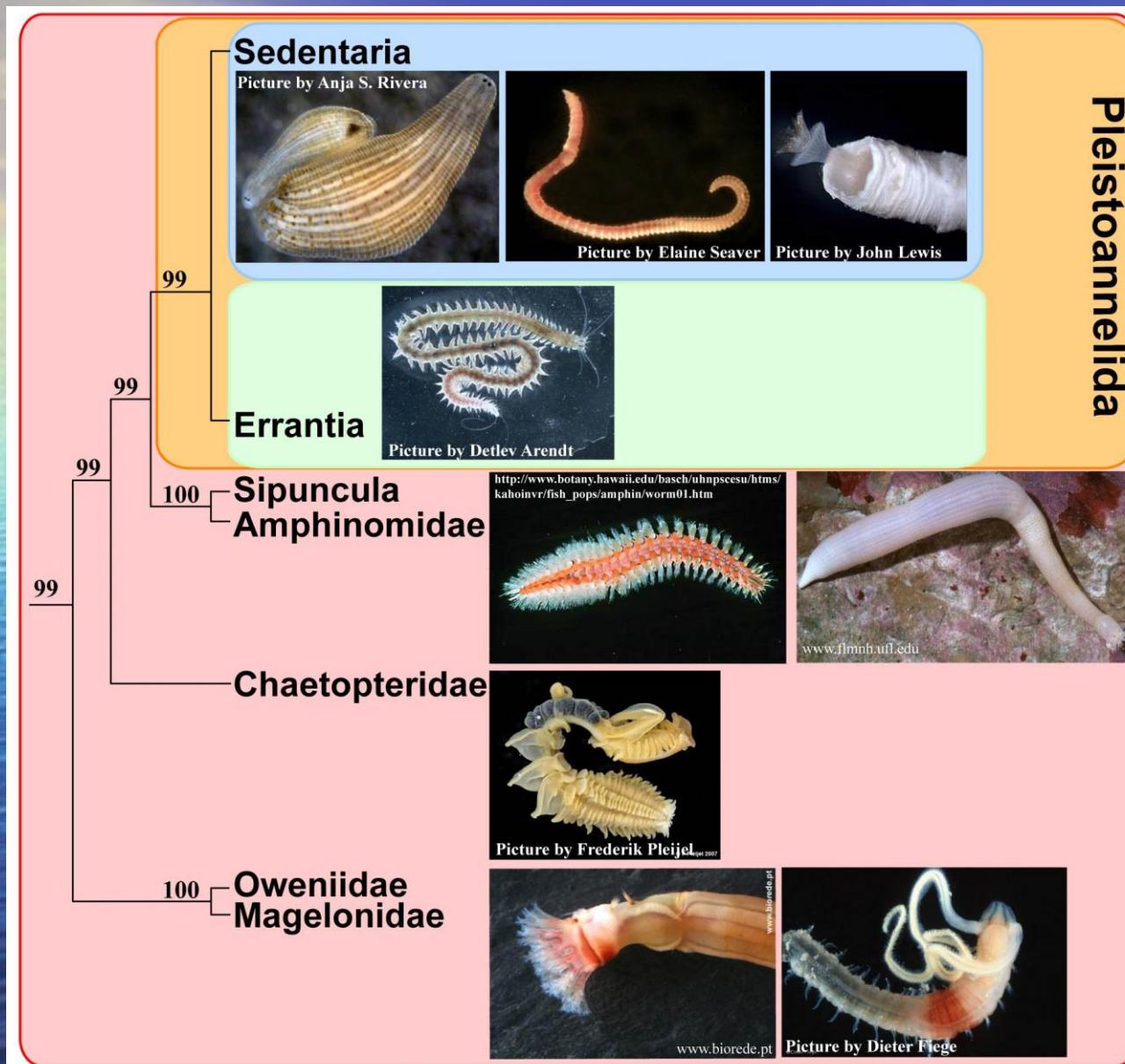
Tree based on
229 genes



Struck (2013) PLoS one



More data and taxa of Annelida



28 new NGS datasets
620 orthologous genes
155422 aa

RAxML

Hughes et al. (2018)

- paralogy originating from inferred WGDs in ancestral vertebrates or teleosts
- two sets of topological constraints (monophyly of all teleosts; monophyly of Ostariophysi)
- topology tests (AU tests) of the constrained tree against the unconstrained ML topology
- rejection of constraint topology ($p < 0.05$) was expected in the presence of paralogy

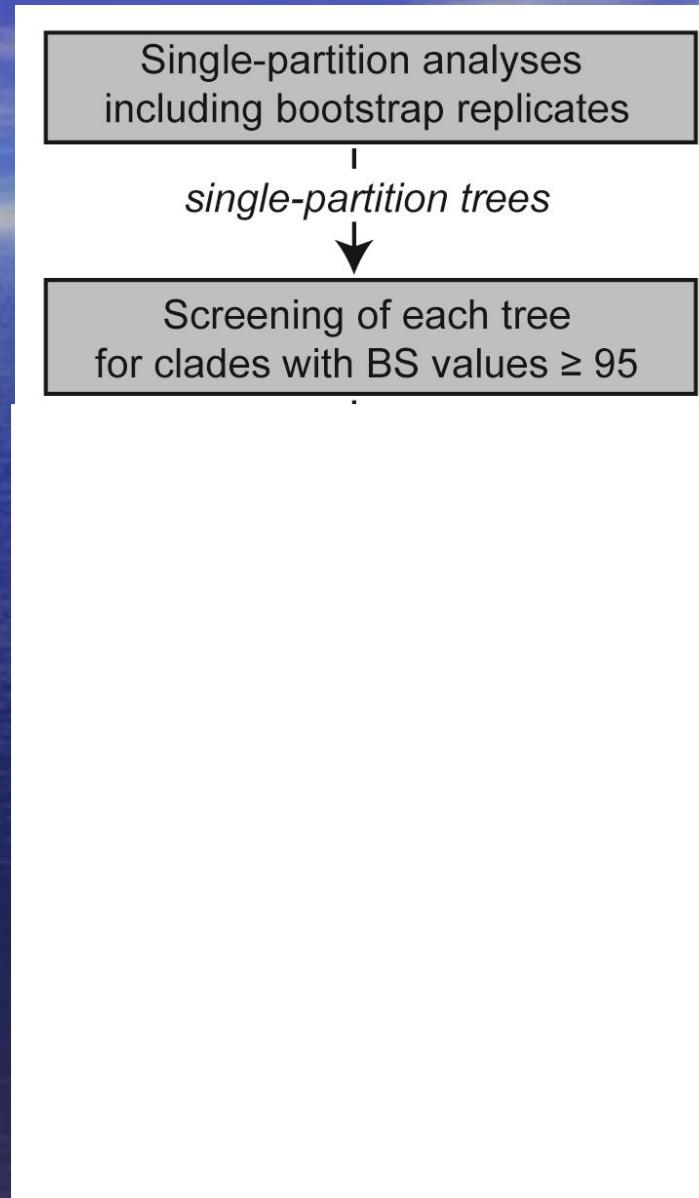
Detection of paralogy

Paralogous sequences
can be problematic in
large-scale analyses
of hundreds of genes

Single-partition analyses
including bootstrap replicates

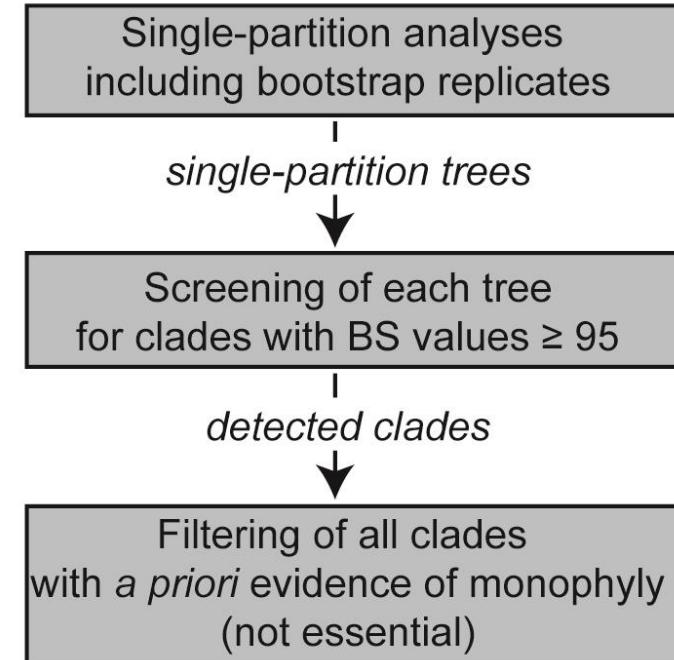
Detection of paralogy

Paralogous sequences
can be problematic in
large-scale analyses
of hundreds of genes



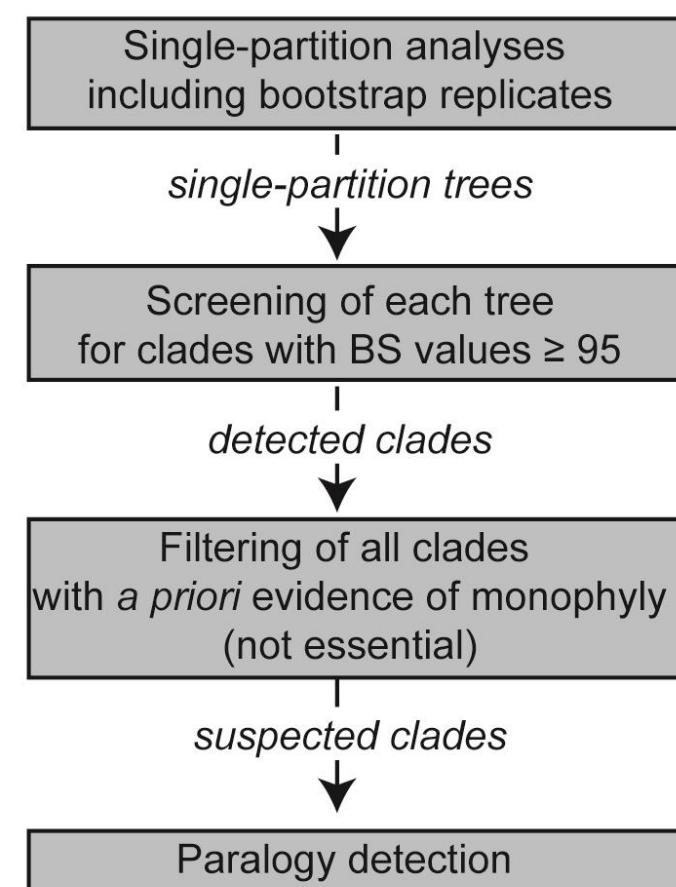
Detection of paralogy

Paralogous sequences
can be problematic in
large-scale analyses
of hundreds of genes



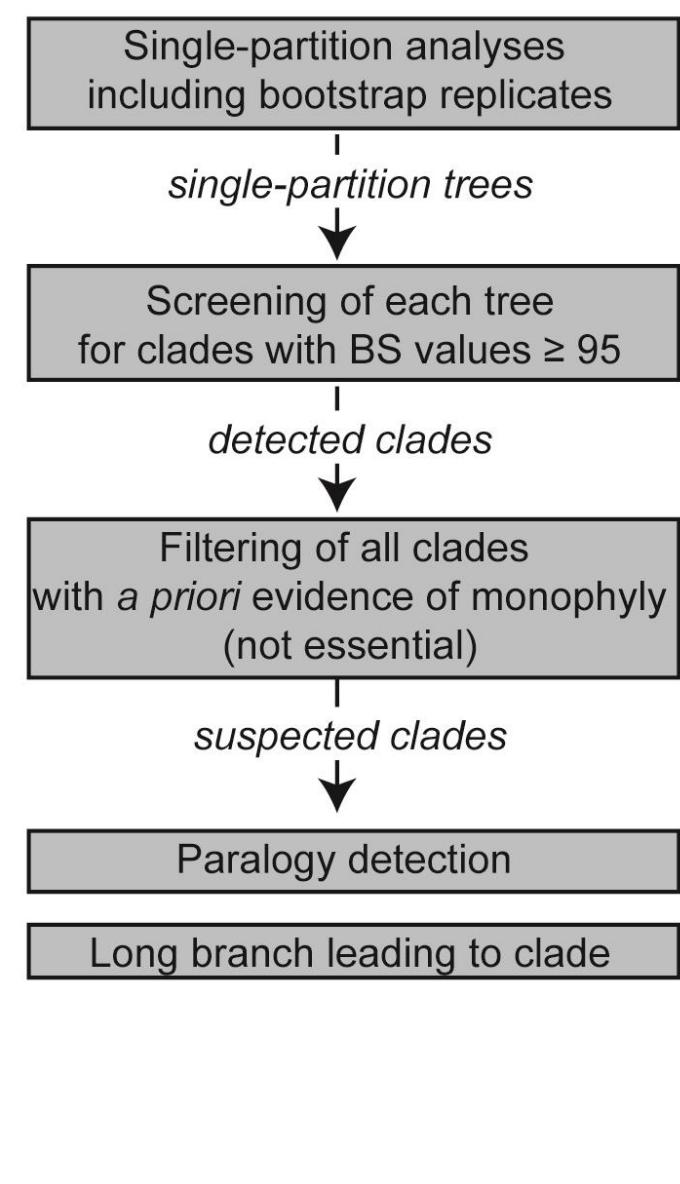
Detection of paralogy

Paralogous sequences
can be problematic in
large-scale analyses
of hundreds of genes



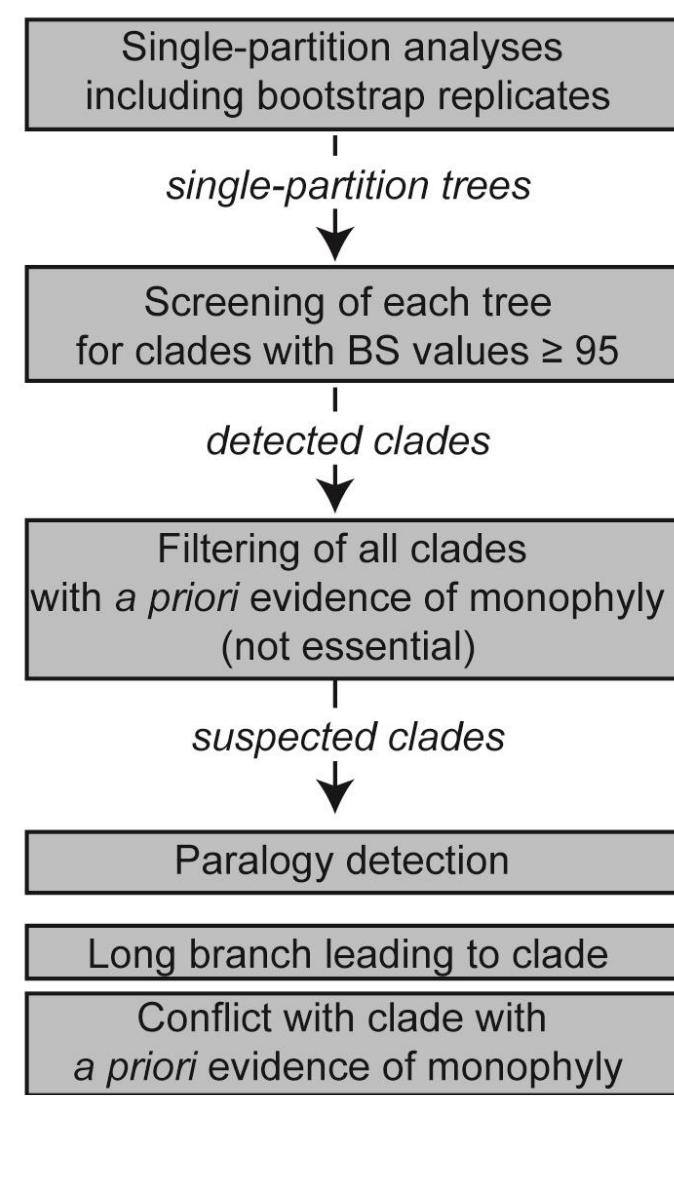
Detection of paralogy

Paralogous sequences
can be problematic in
large-scale analyses
of hundreds of genes



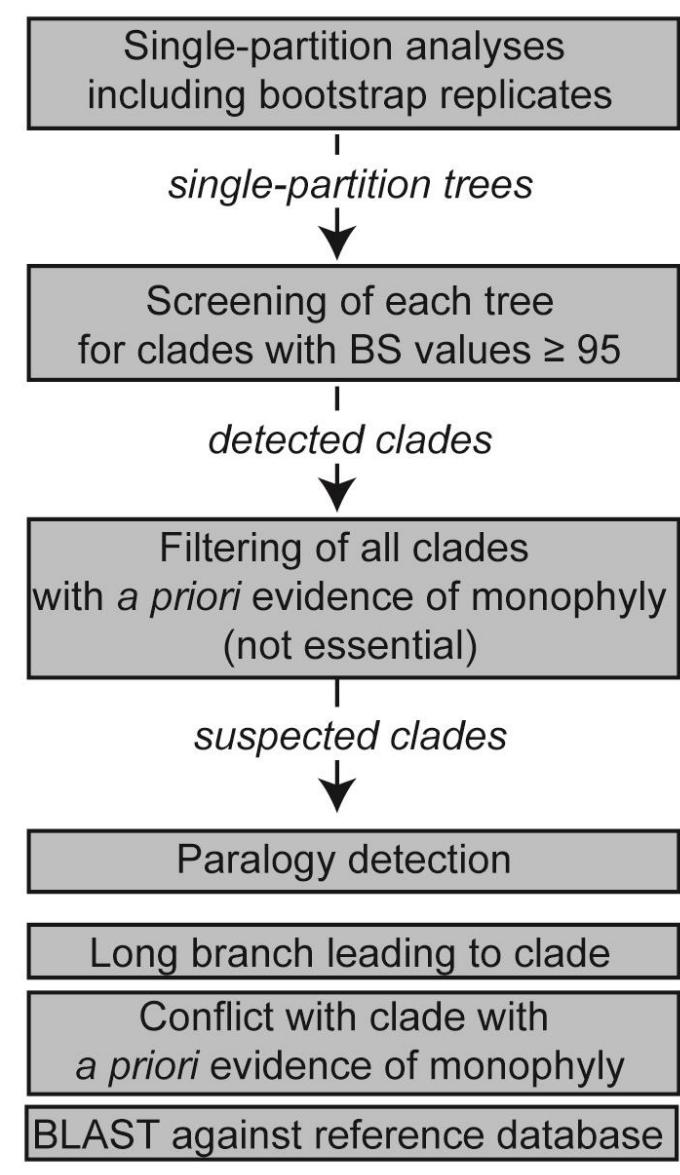
Detection of paralogy

Paralogous sequences
can be problematic in
large-scale analyses
of hundreds of genes



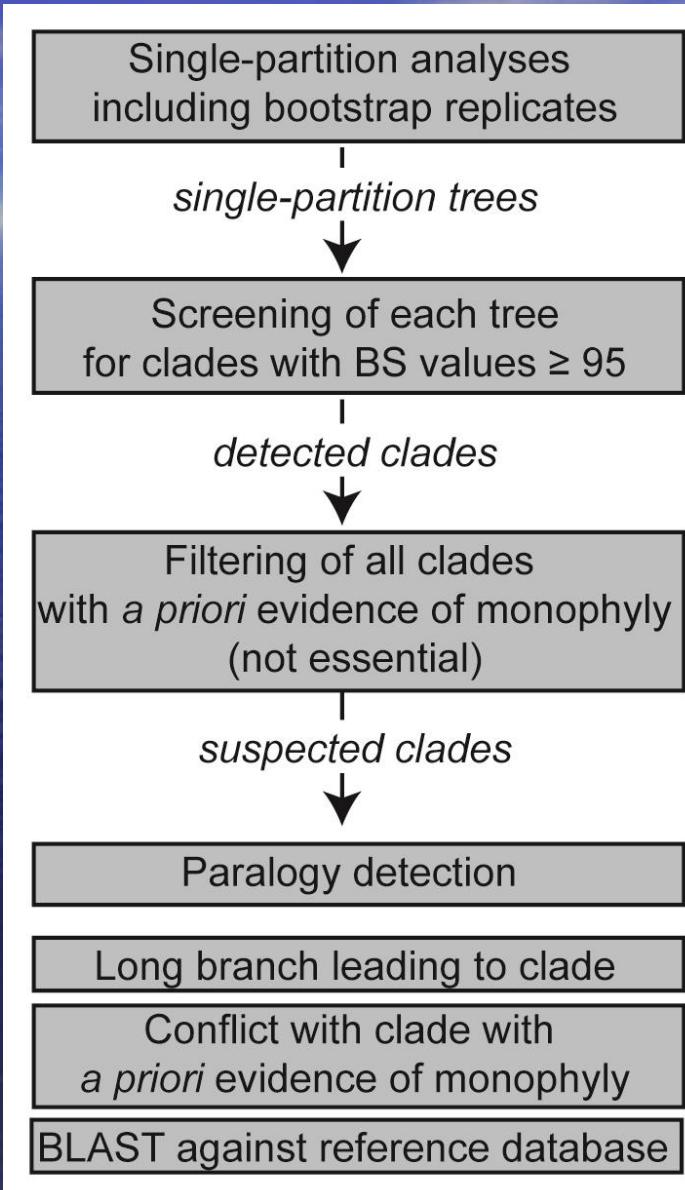
Detection of paralogy

Paralogous sequences
can be problematic in
large-scale analyses
of hundreds of genes



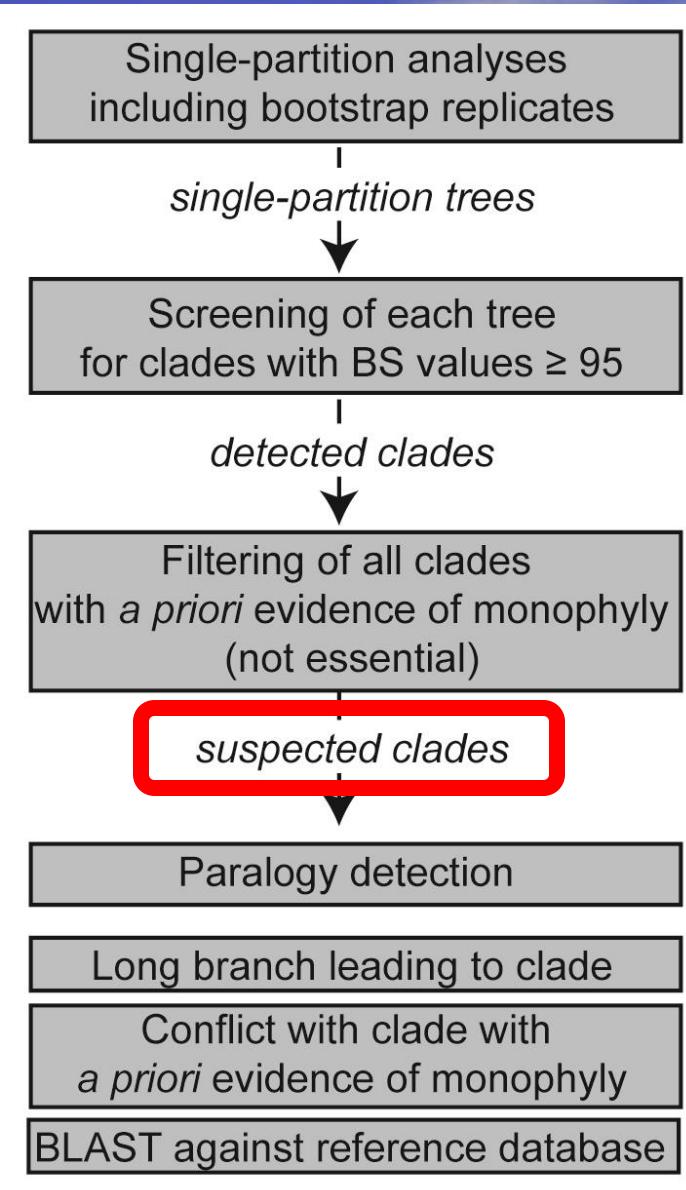
Detection of deviating single genes

Can be used to detect clades in single genes deviating from a given relationship. Based on bootstrap support



Detection of deviating single genes

Can be used to detect clades in single genes deviating from a given relationship. Based on bootstrap support



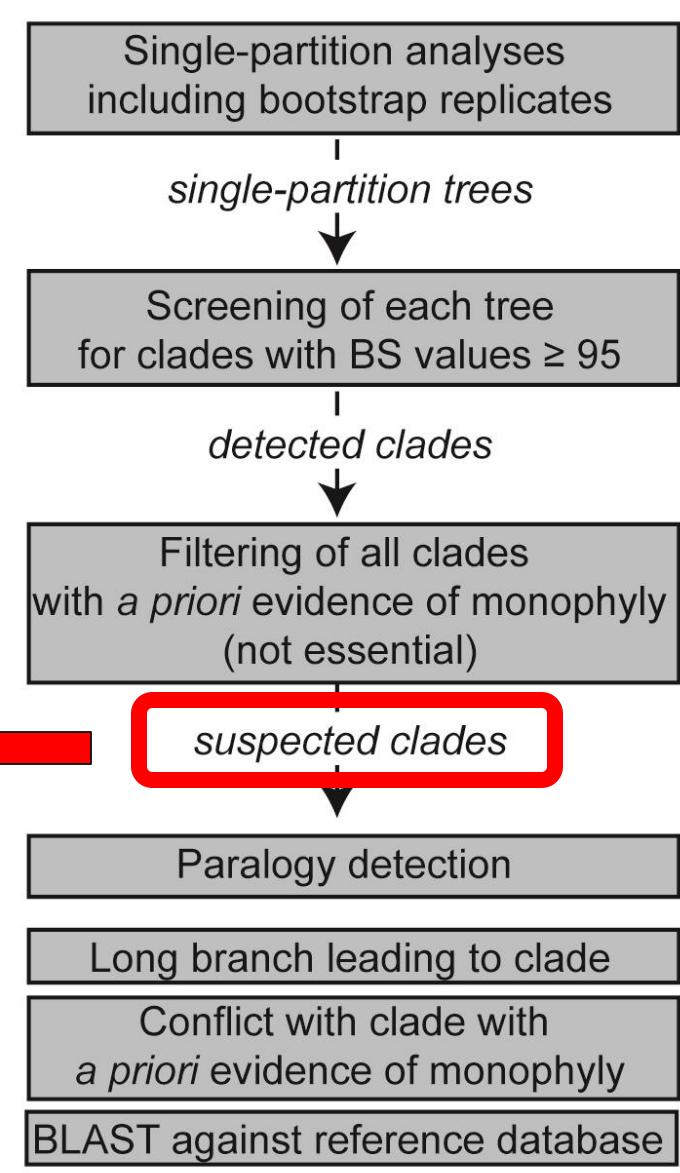
Detection of deviating single genes

Can be used to detect clades in single genes deviating from a given relationship. Based on bootstrap support

Paralogy is only one possible cause of deviation.

Others are:

incomplete lineage sorting,
horizontal gene transfer or
ancestral hybridization.



TreSpEx



Libertas Academica
FREEDOM TO RESEARCH

Open Access: Full open access to
this and thousands of other papers at
<http://www.la-press.com>.

Evolutionary Bioinformatics

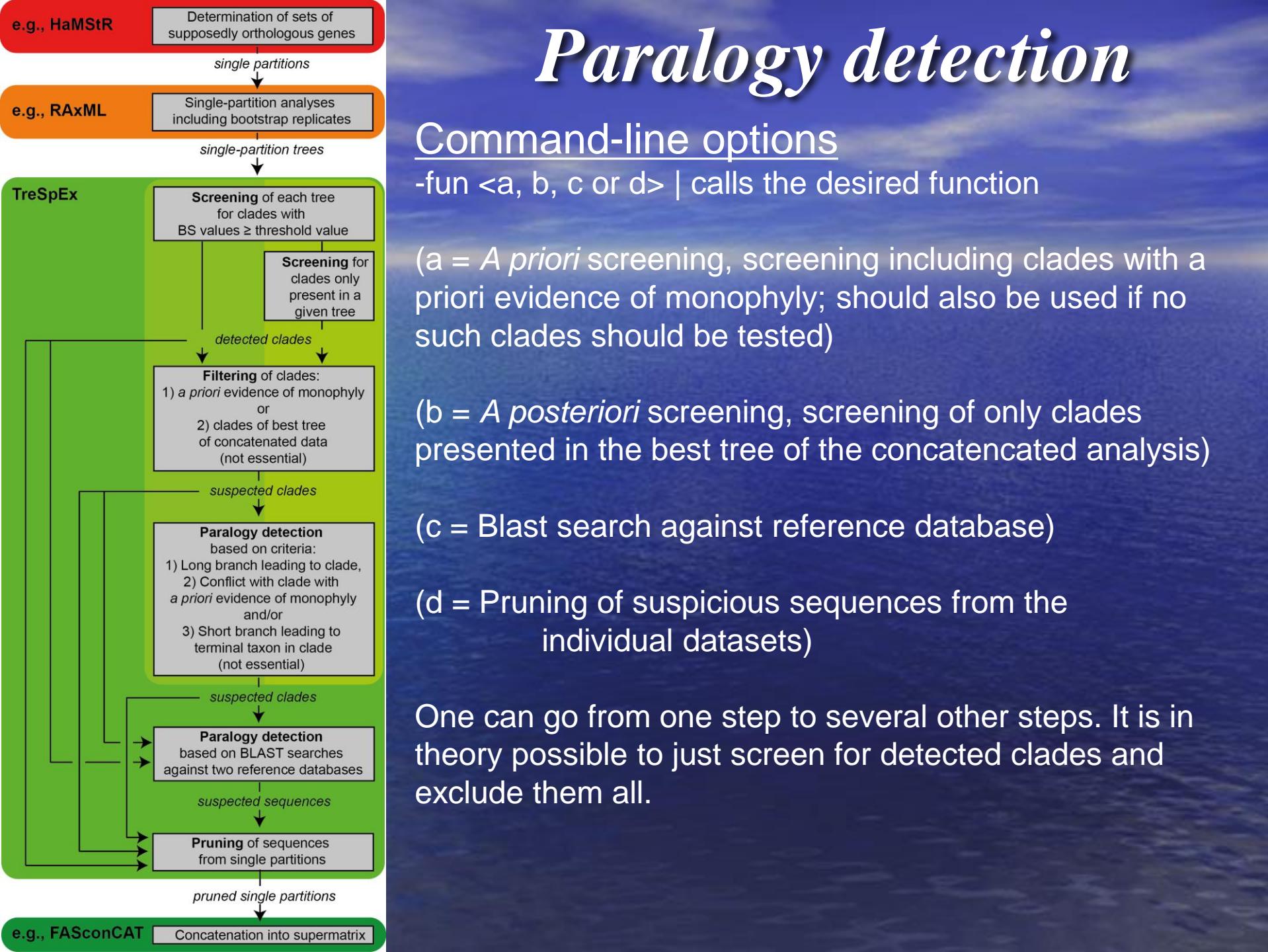
TreSpEx—Detection of Misleading Signal in Phylogenetic Reconstructions Based on Tree Information

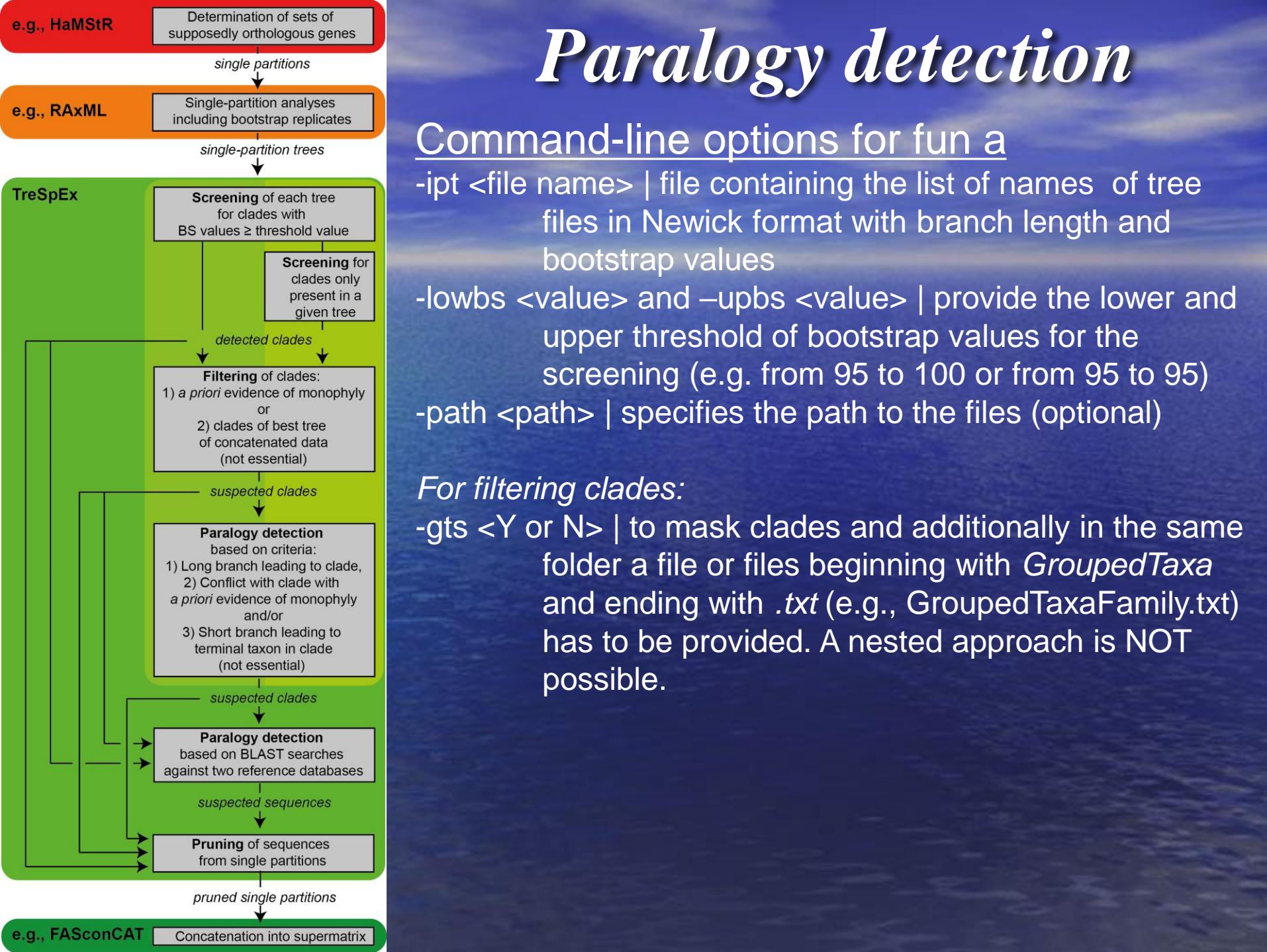
Torsten H. Struck

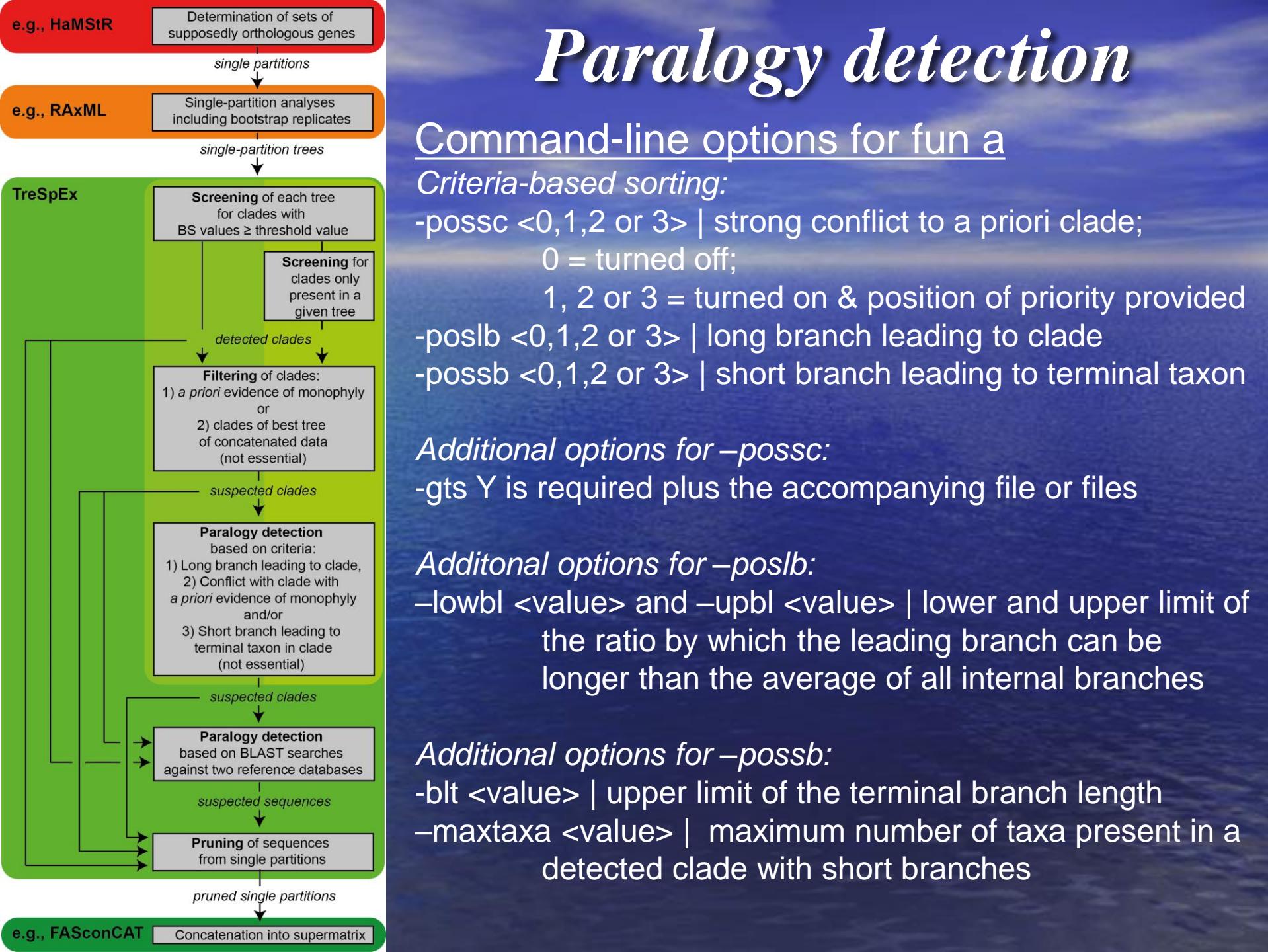
Zoological Research Museum Alexander Koenig, Bonn, Germany.

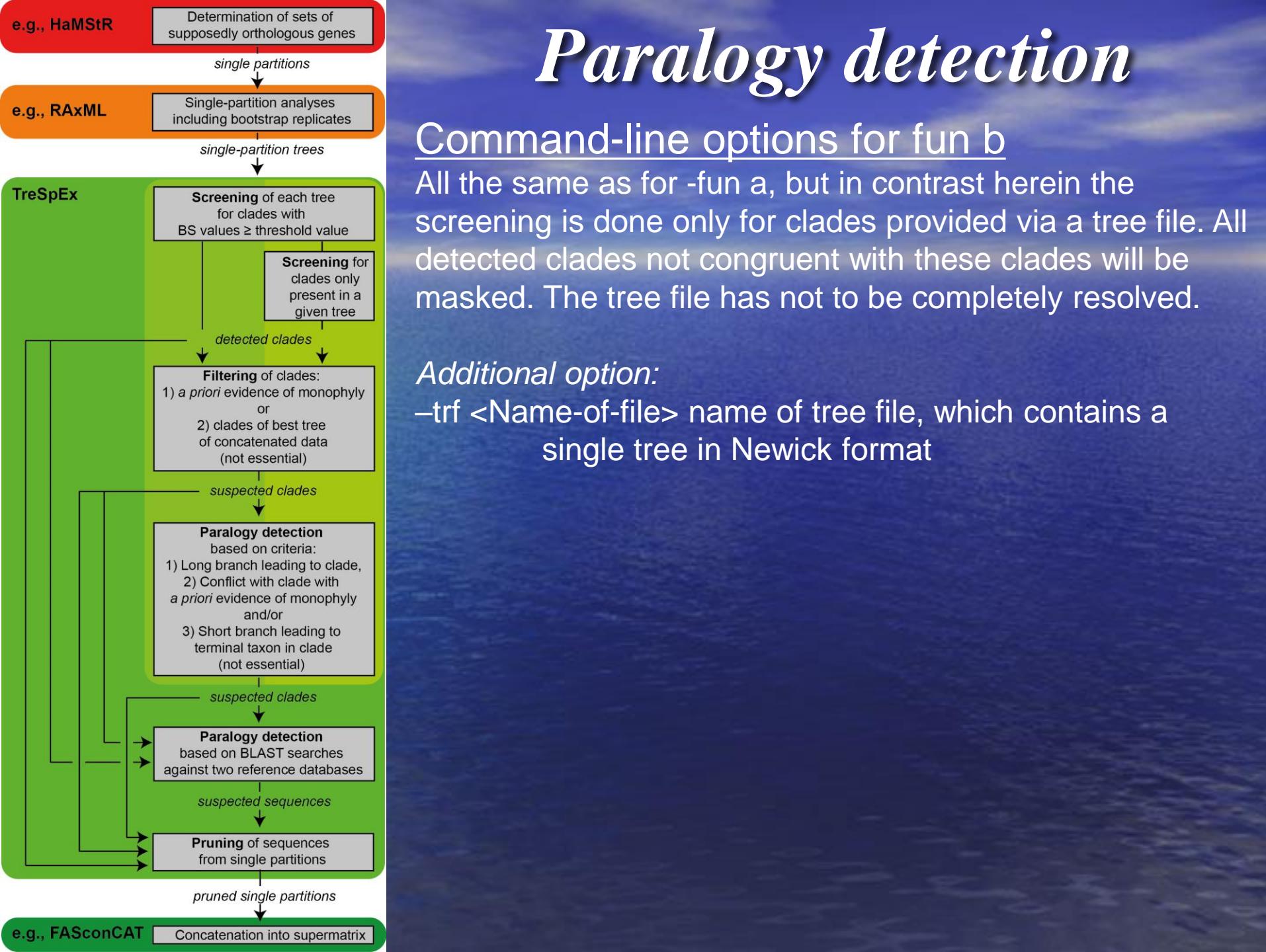
Different methods related to tree-based measurements:

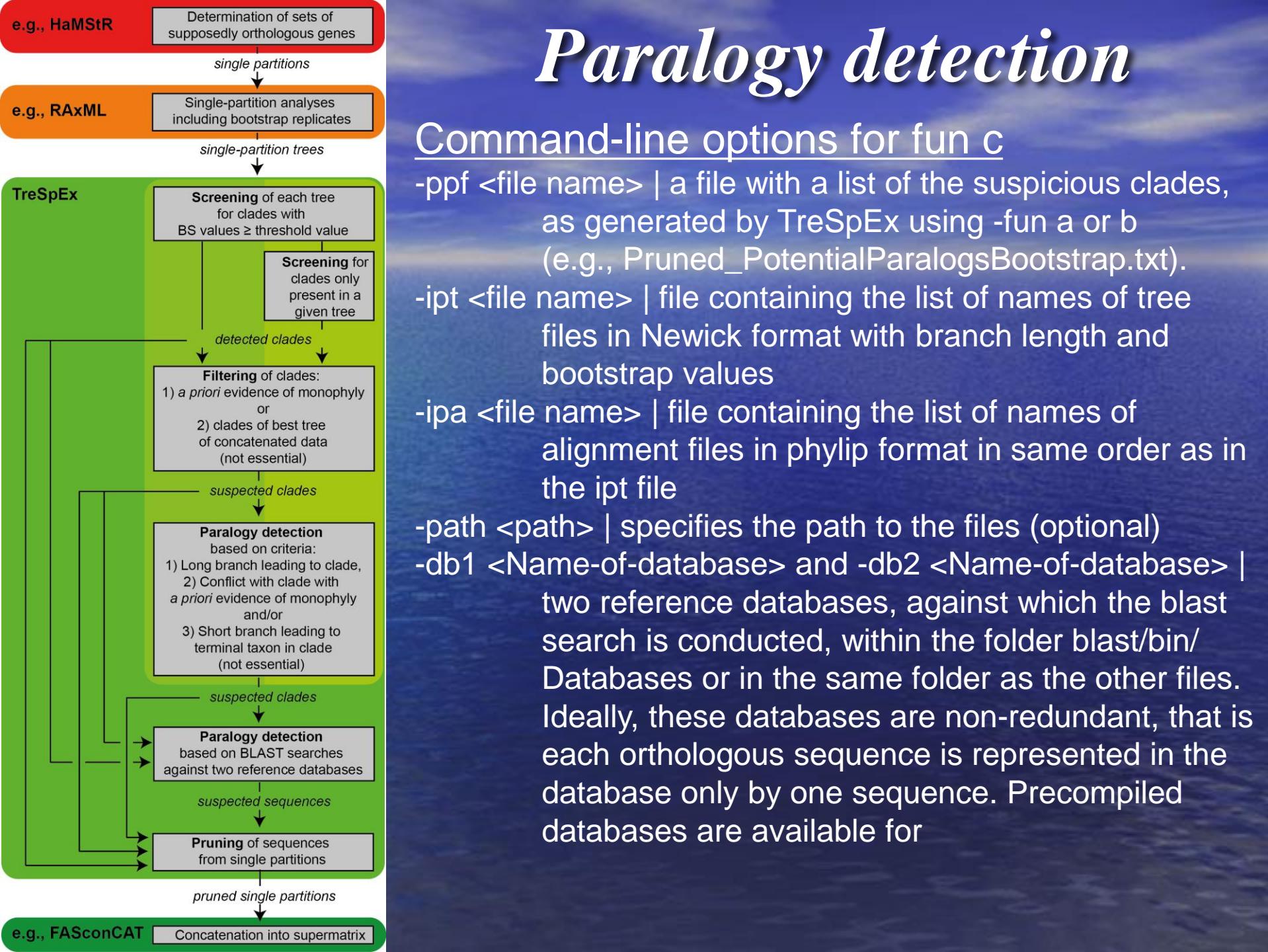
- **Detection and pruning of paralogous sequences**
- Detection of conflict
- Detection of long-branched taxa and partitions
- Detection of saturation and phylogenetic signal



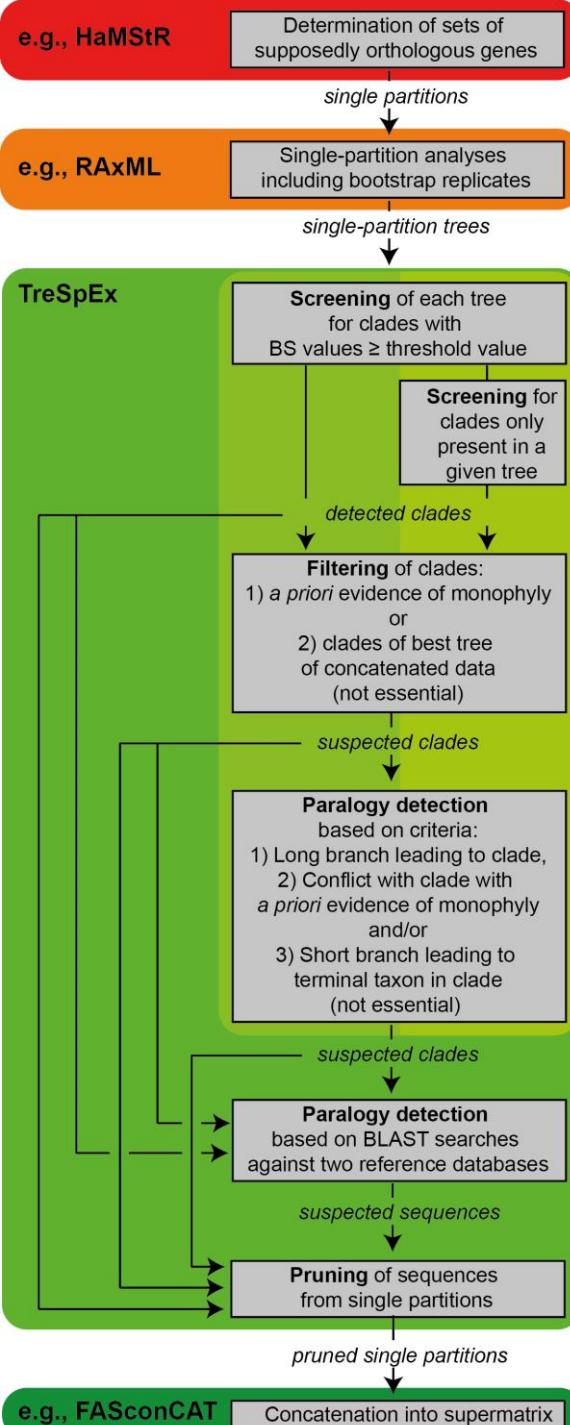








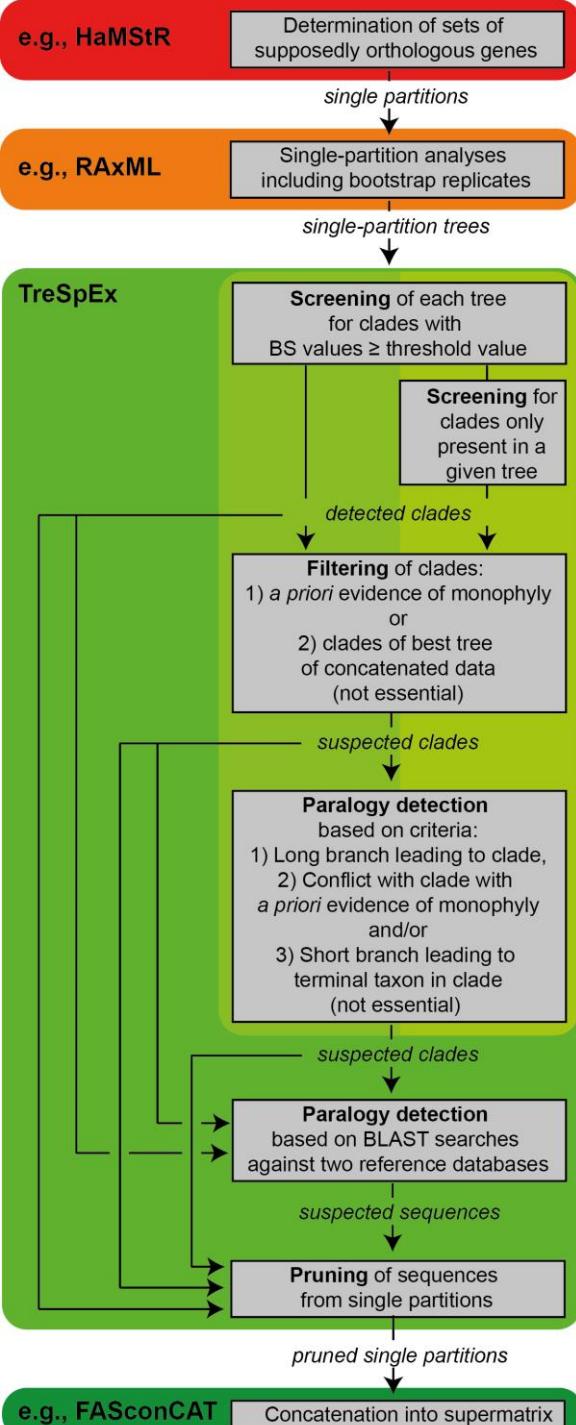
Paralogy detection



Available databases (have to be placed in the blast/bin/Databases folder):

Anolis_carolinensis *Apis_mellifera*
Bombyx_mori *Bos_taurus*
Branchiostoma_floridae *Caenorhabditis_elegans*
Capitella_teleta *Ciona_intestinalis*
Danio_rerio *Daphnia_pulex*
Drosophila_melanogaster *Gallus_gallus*
Helobdella_robusta *Homo_sapiens*
Lottia_gigantia *Monodelphis_domestica*
Mus_musculus *Ornithorhynchus_anatinus*
Pediculus_humanus *Schmidtea_mediterranea*
Taeniopygia_guttata *Tribolium_castaneum*
Xenopus_tropicalis

Paralogy detection



Command-line options for fun c

-evalue <value> | The threshold of the e value for the blast searches (e.g., -evalue 1e-20)

TreSpEx will automatically determine, which blast search has to be conducted given the database and the alignment file and also if blast-searchable database has to be generated.

Sorting options (no hits, certain, no paralogy and uncertain):

-ltp <value> and -utp <value> | lower and upper sorting threshold values in proportion of identical blast results (e.g., -ltp 0.0 -utp 1.0)

