



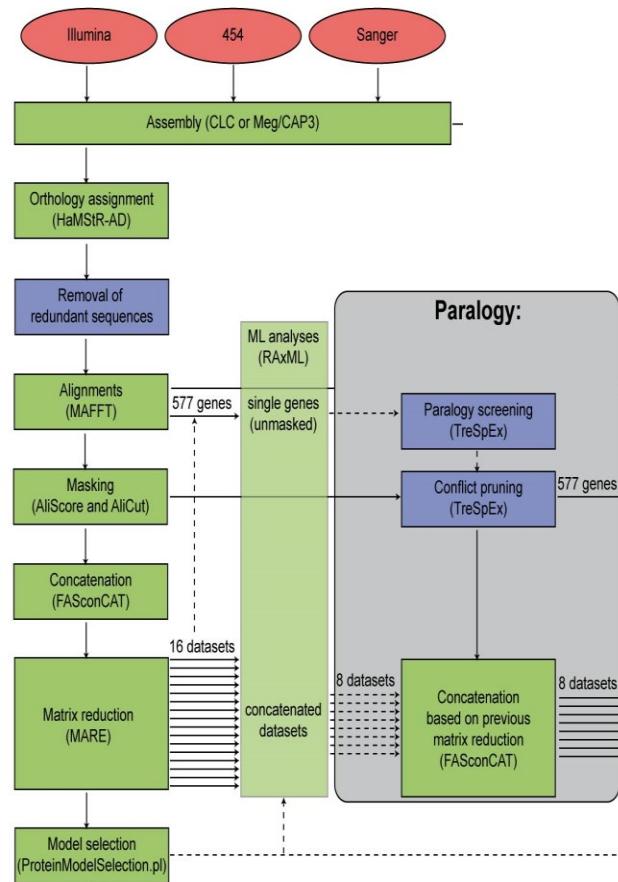
# Paralogy



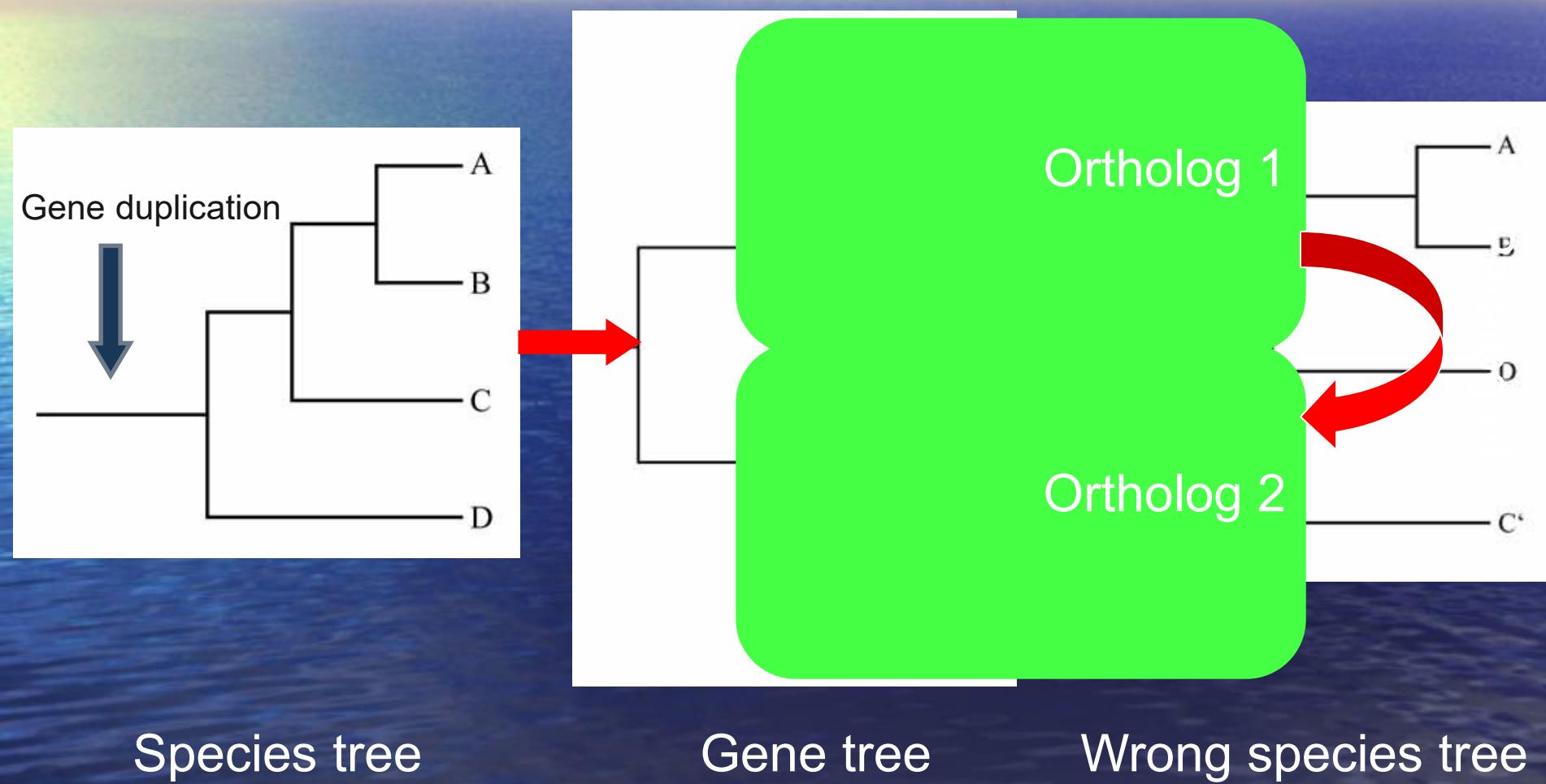
**Torsten Struck – [t.h.struck@nhm.uio.no](mailto:t.h.struck@nhm.uio.no)**  
**NHM UiO**

© nhm.uio.no

# *Going beyond the standard*

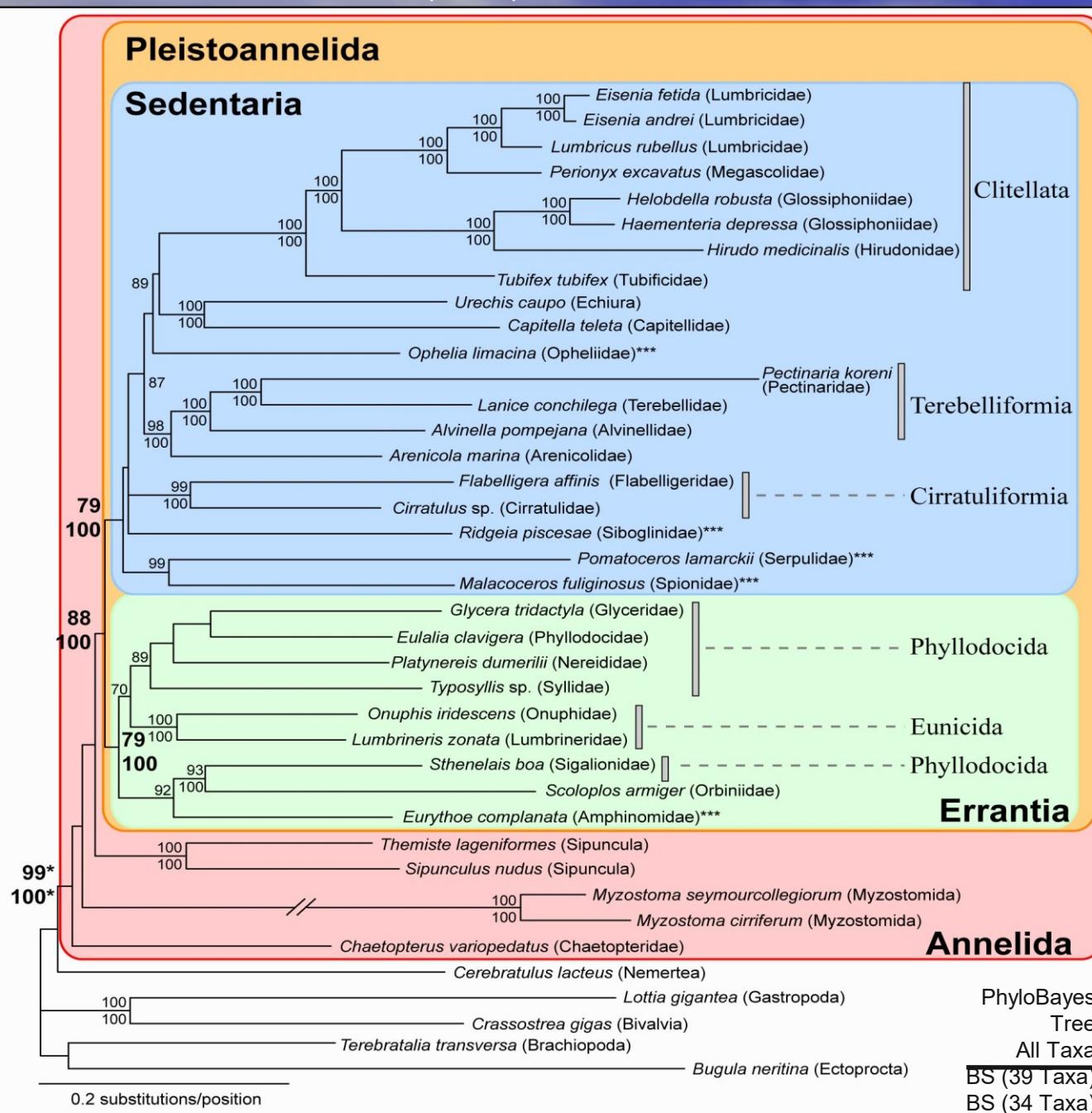


# *Problem of paralogy*

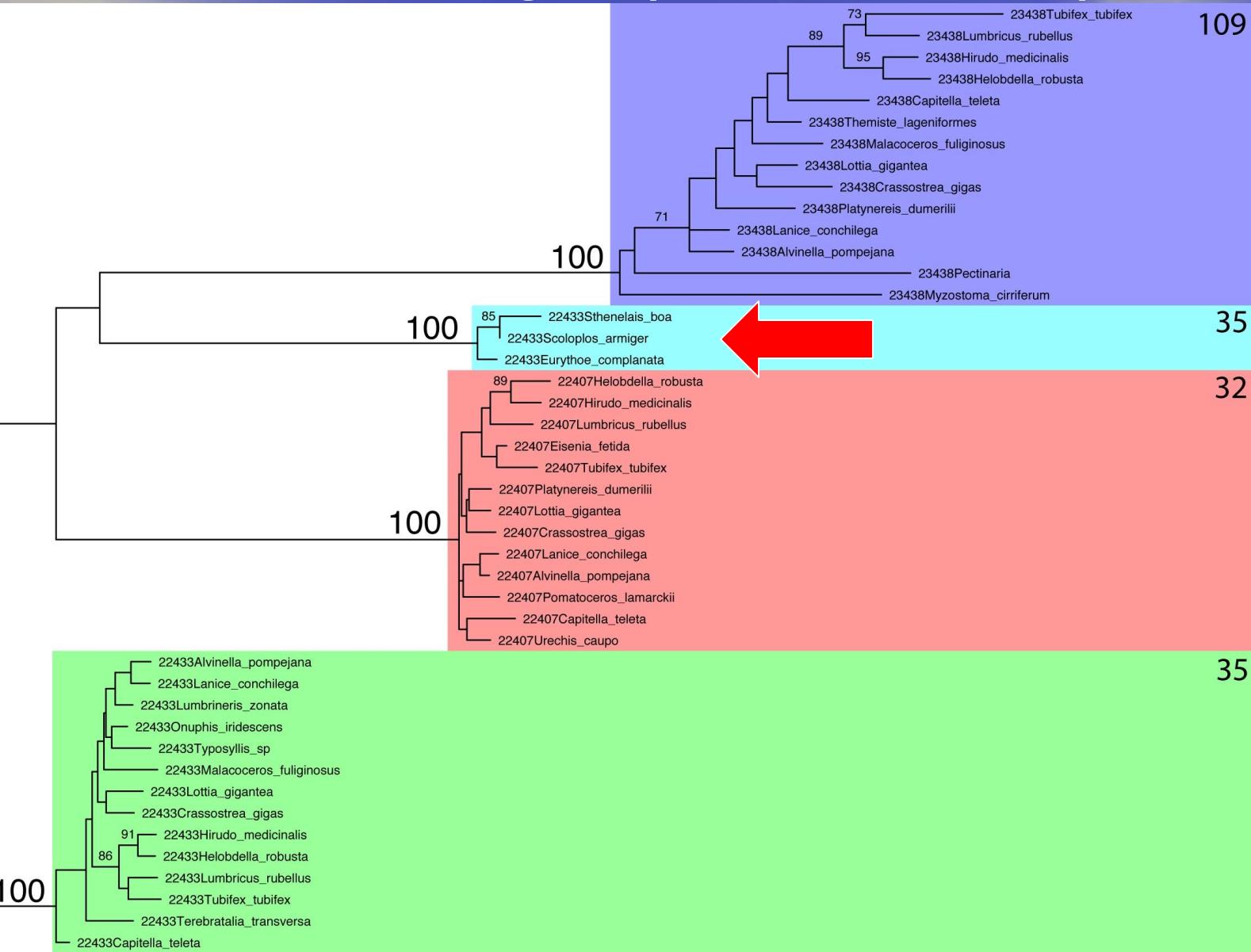


# Phylogenomics

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



# Paralogs (32/109/35)



# Paralogs (32/109/35)

Blast against *Bos taurus* transcriptome  
= Proteasome subunit alpha

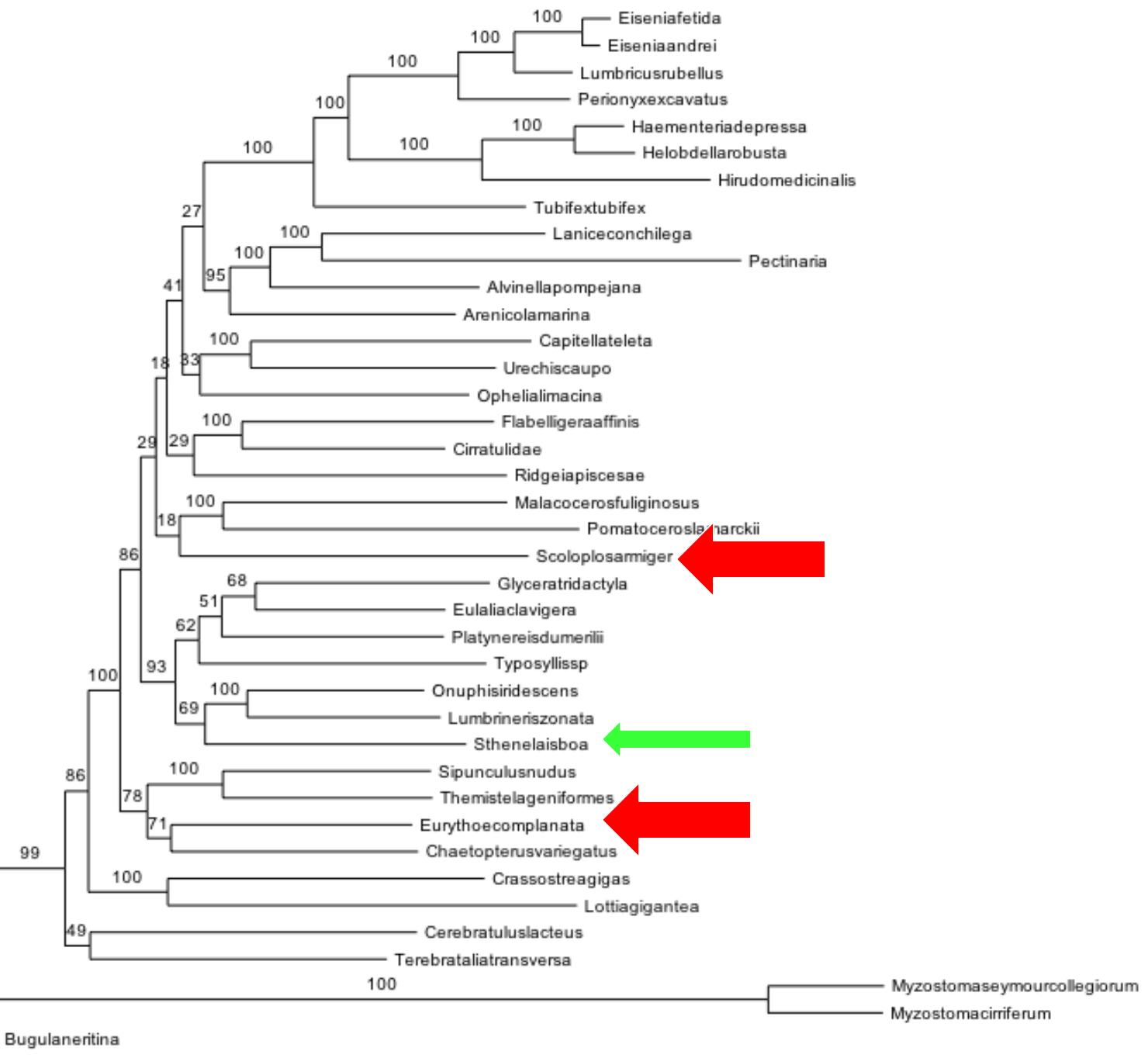
reference taxa	32 (23428)	109 (22407)	35 (22433)
Helobdella	PSMA5	PSMA3	PSMA7/8
Capitella	PSMA5	PSMA3	PSMA7/8
Lottia	PSMA5	PSMA3	PSMA7/8

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

Blast of *Bos taurus* sequence against *Helobdella* transcriptome:  
PSMA2 (NM\_001034662): *Helobdella* sequence 185484 e^-39  
That is the sequence of orthology group 22433

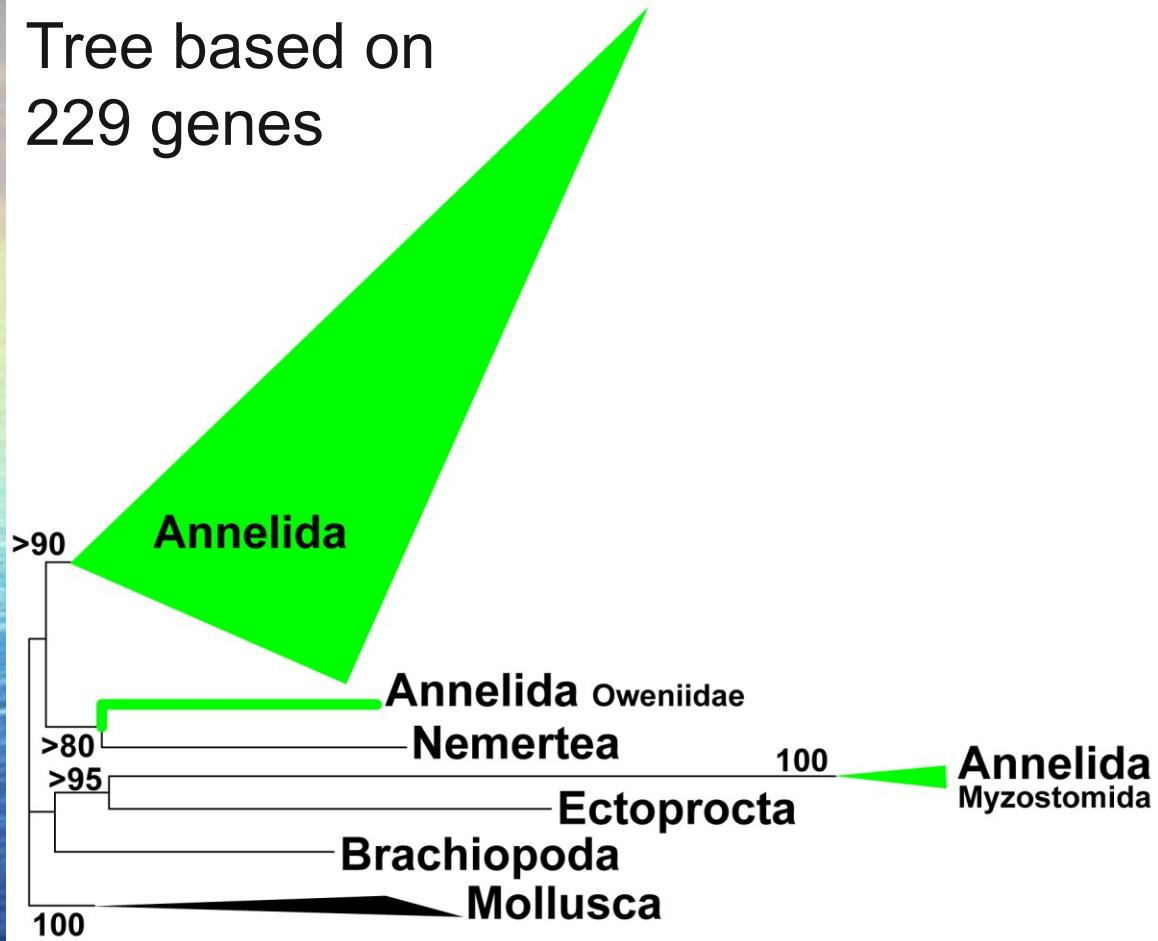
# What has happened?

	PSMA7/8	PSMA2
Helobdella	present	absent
Eurythoe	absent	present

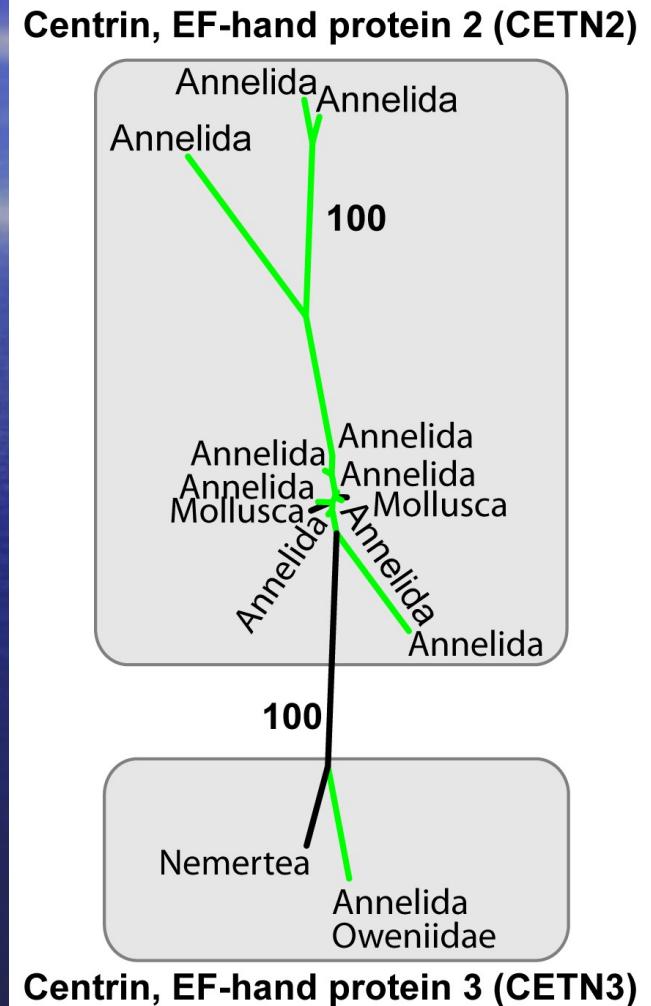


# Problem of paralogy

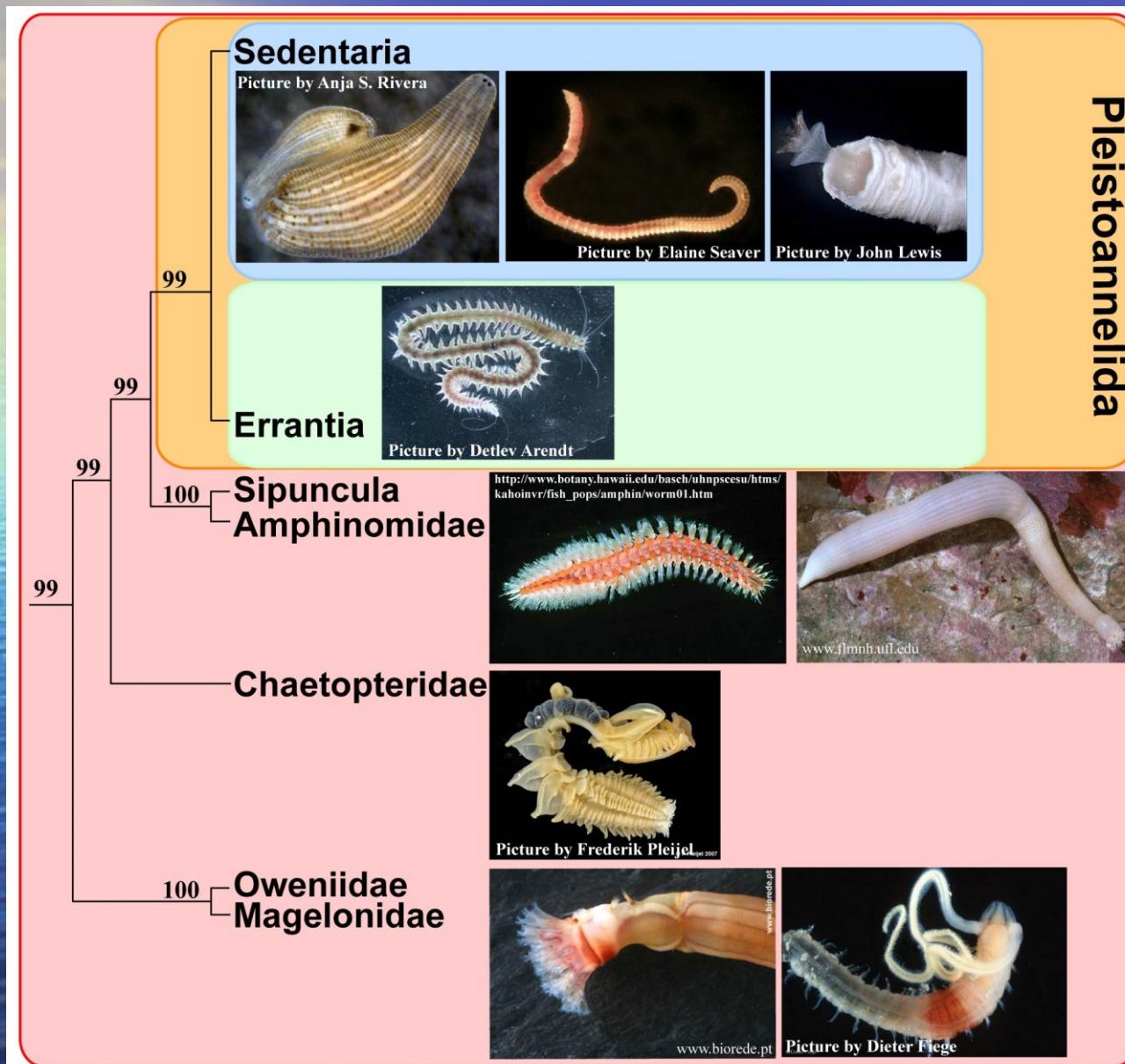
Tree based on  
229 genes



Struck (2013) PLoS one



# More data and taxa of Annelida



28 new NGS datasets  
620 orthologous genes  
155422 aa

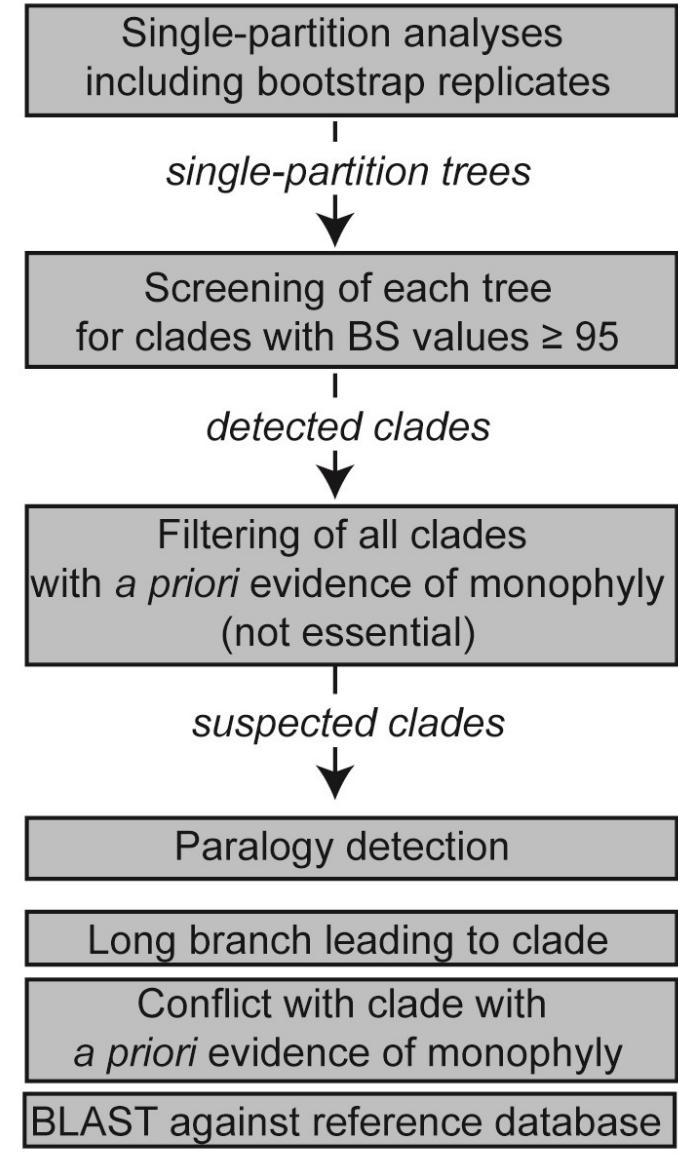
RAxML

# *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*

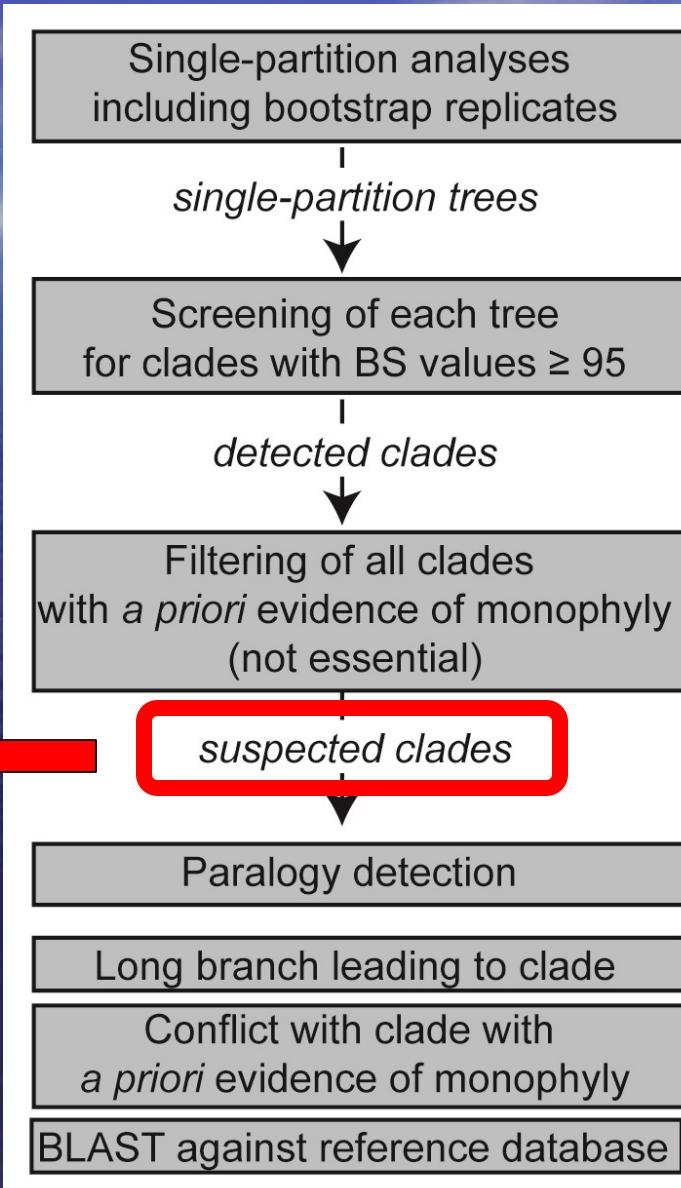
Paralogous sequences  
can be problematic in  
large-scale analyses  
of hundreds of genes



# *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

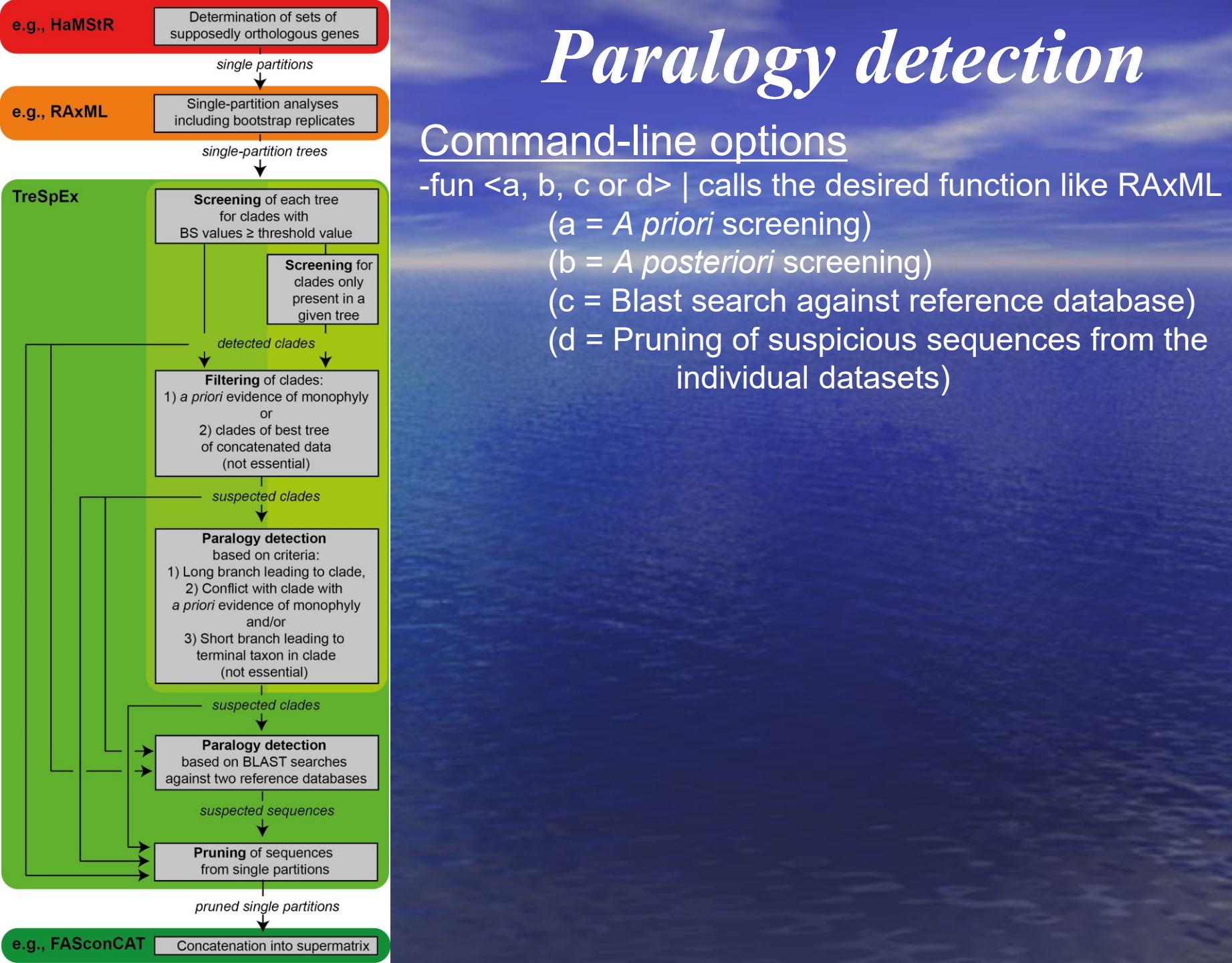
### TreSpEx—Detection of Misleading Signal in Phylogenetic Reconstructions Based on Tree Information

Torsten H. Struck

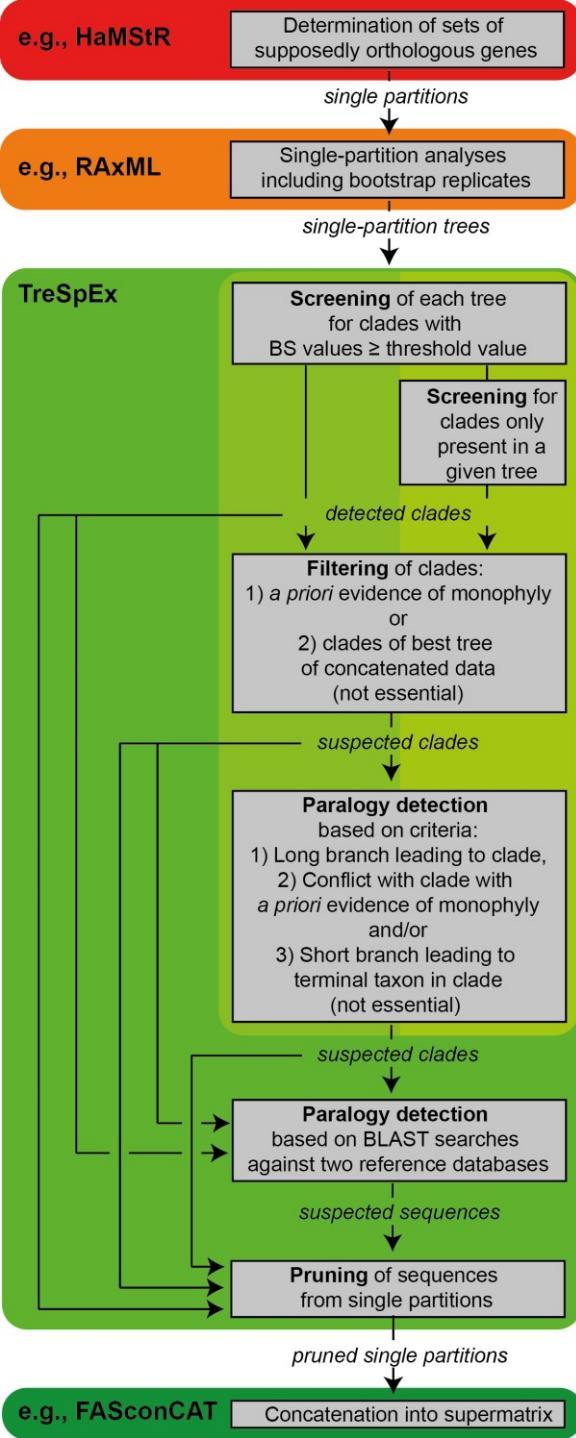
Zoological Research Museum Alexander Koenig, Bonn, Germany.

#### 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

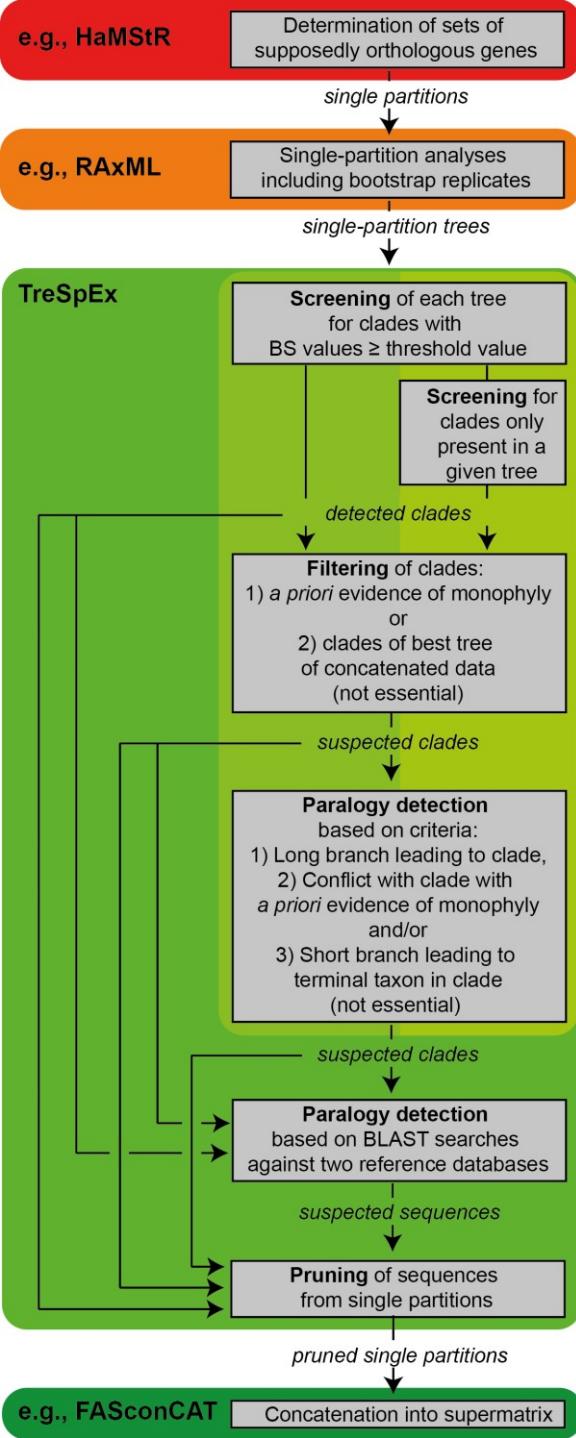


## Command-line options for fun a

-ipt <file name> | file containing the list of names of tree files in Newick format with branch length and bootstrap values  
 -lowbs <value> and -upbs <value> | provide the lower and upper threshold of bootstrap values for the screening (e.g. from 95 to 100 or from 95 to 95)  
 -path <path> | specifies the path to the files (optional)

### *For filtering clades:*

-gts <Y or N> | to mask clades and additionally in the same folder a file or files beginning with *GroupedTaxa* and ending with .txt (e.g., GroupedTaxaFamily.txt) has to be provided. A nested approach is NOT possible.



# *Paralogy detection*

# Command-line options for fun a

## *Criteria-based sorting:*

-possc <0,1,2 or 3> | strong conflict to a priori clade;  
0 = turned off;  
1, 2 or 3 = turned on & position of priority provided  
-poslb <0,1,2 or 3> | long branch leading to clade  
-possb <0,1,2 or 3> | short branch leading to terminal taxon

## *Additional options for –possc:*

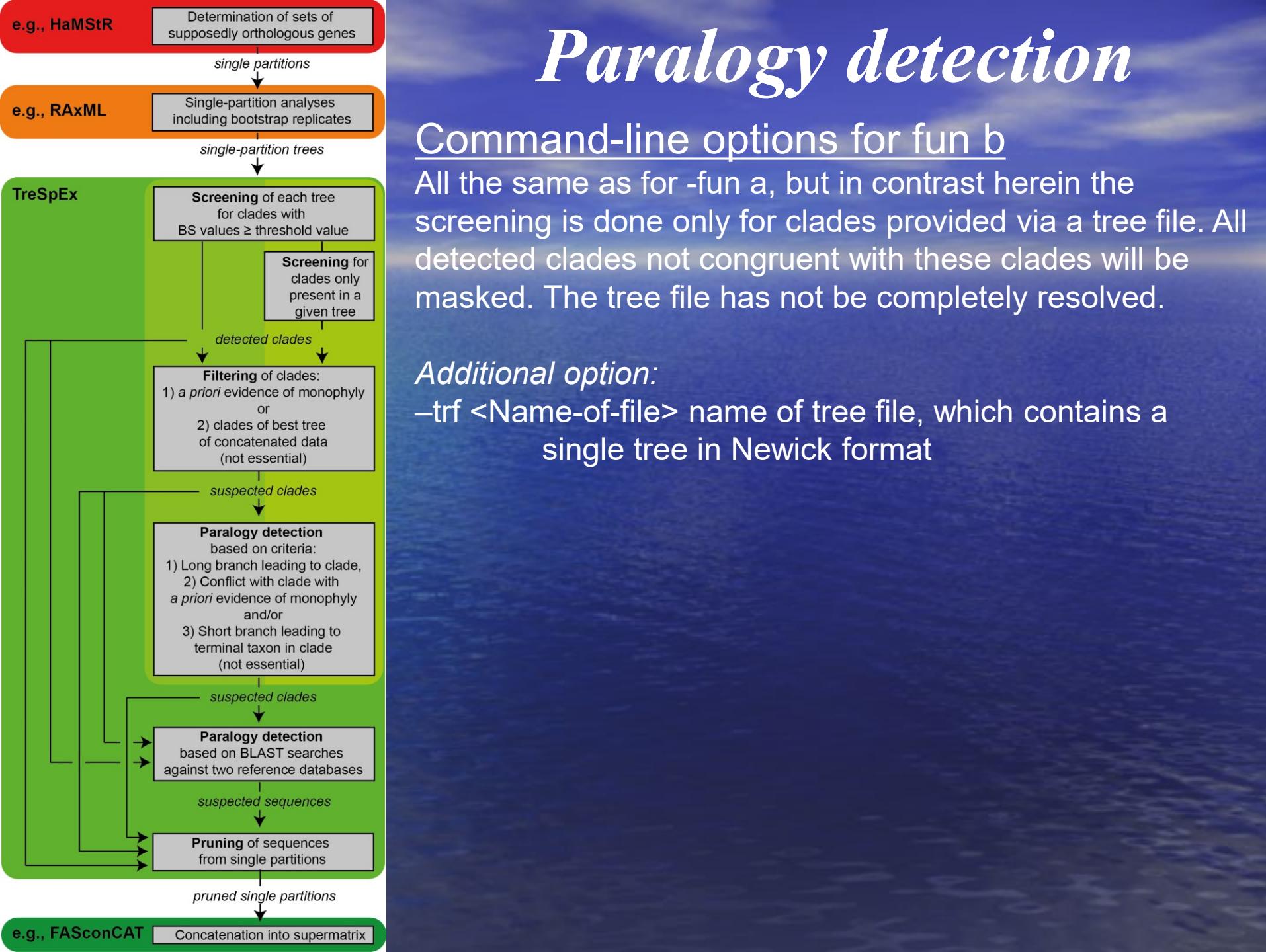
-gts Y is required plus the accompanying file or files

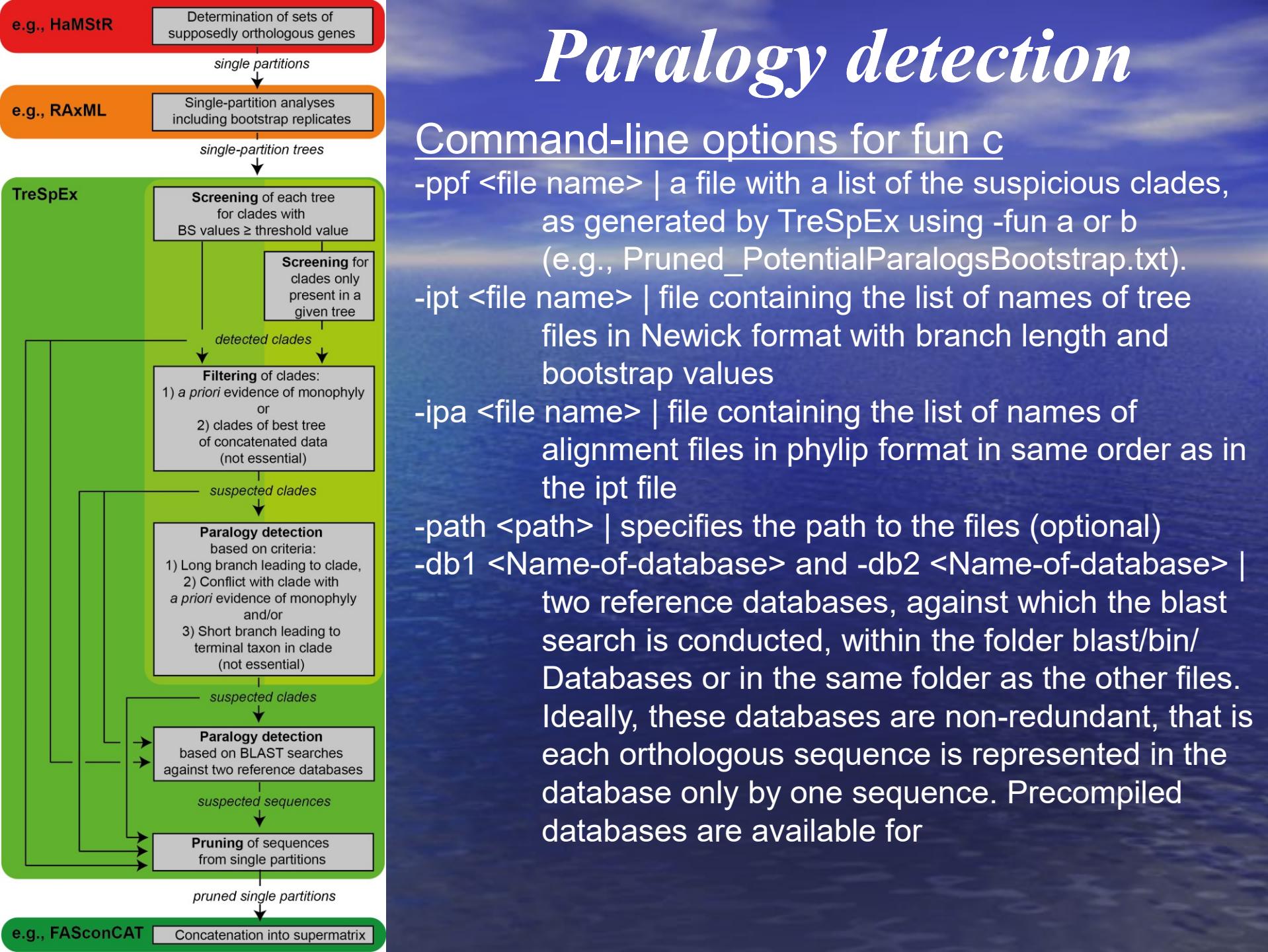
## *Additional options for –poslb.*

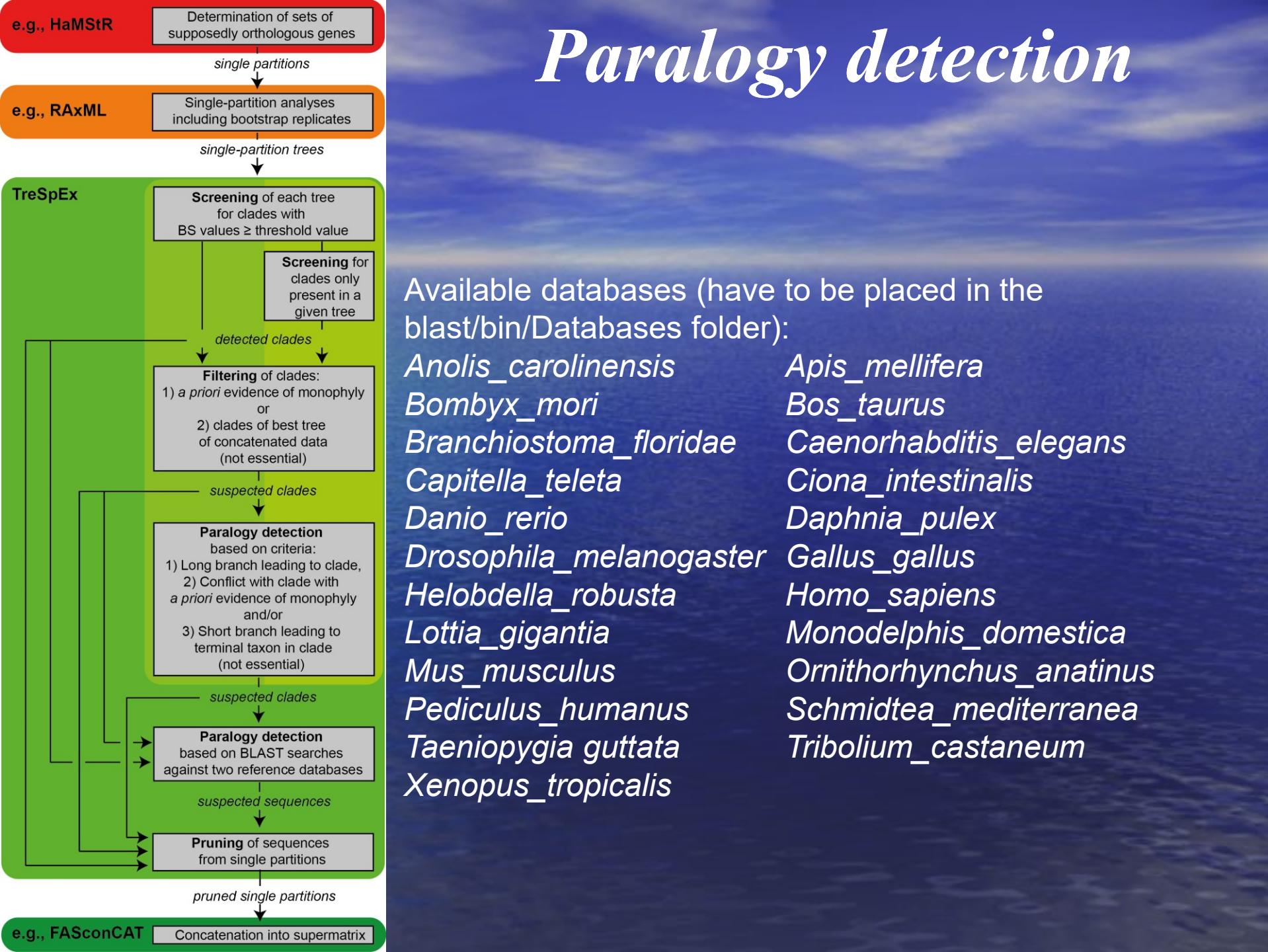
**-lowbl <value>** and **-upbl <value>** | lower and upper limit of the ratio the leading branch can be longer than the average of all internal branches

### *Additional options for –possb:*

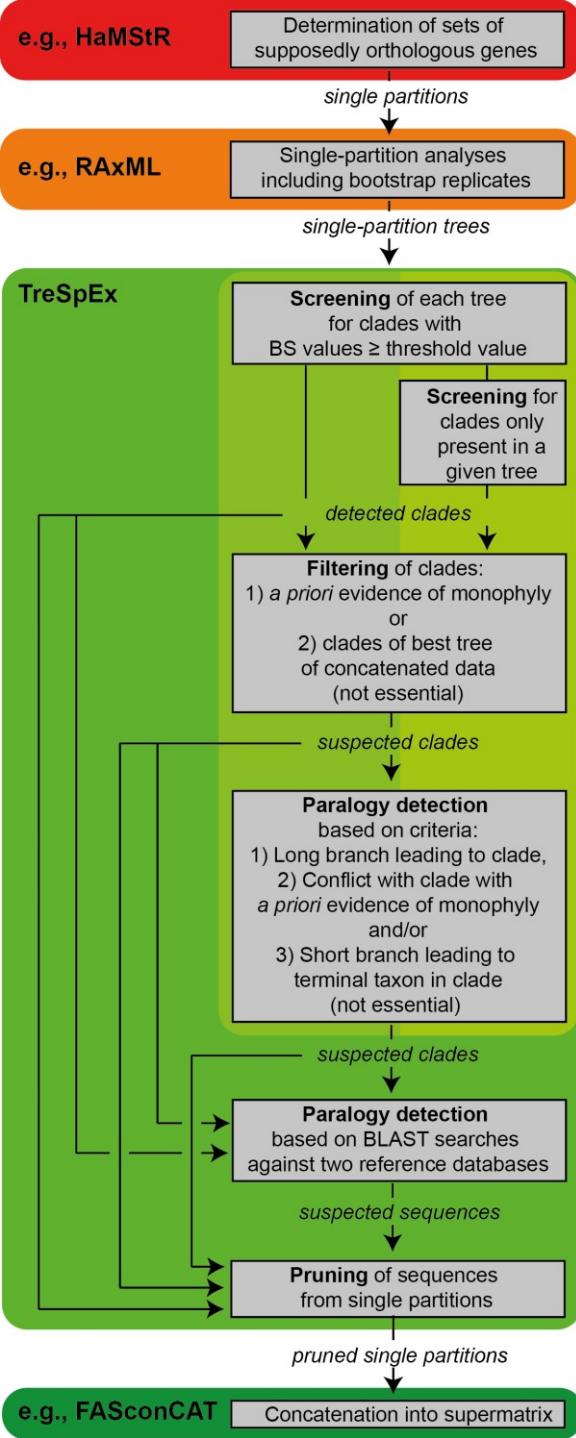
**-blt <value>** | upper limit of the terminal branch length  
**-maxtaxa <value>** | maximum number of taxa present in a detected clade with short branches







# Paralogy detection



## Command-line options for fun c

**-evaluate <value>** | The threshold of the e value for the blast searches (e.g., -evaluate 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)

