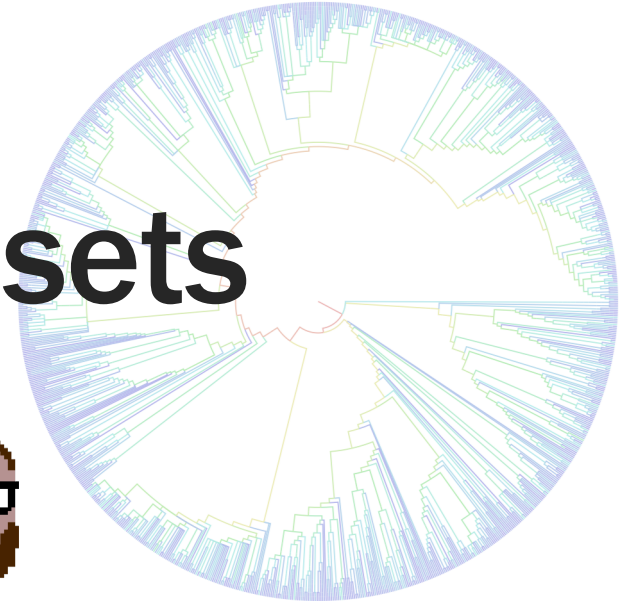
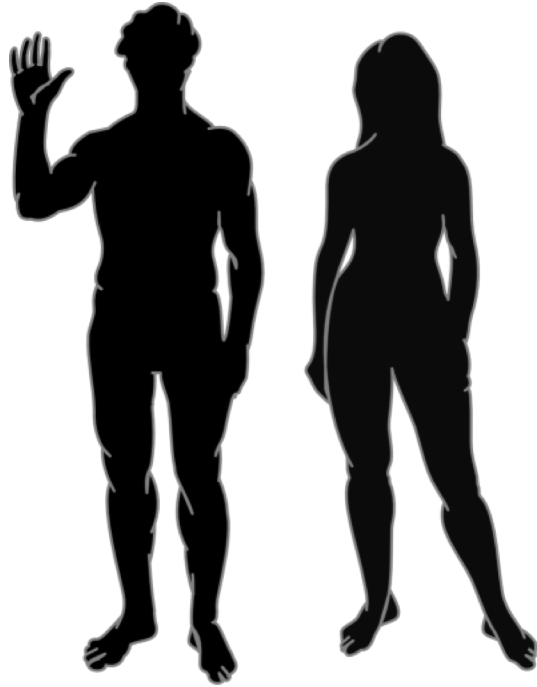
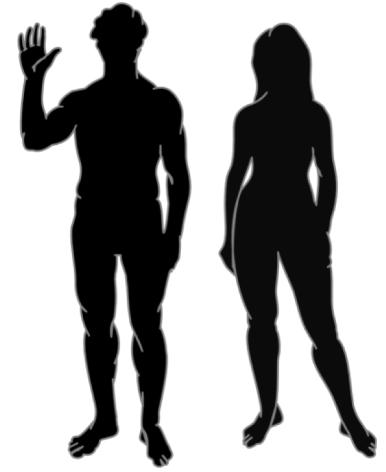


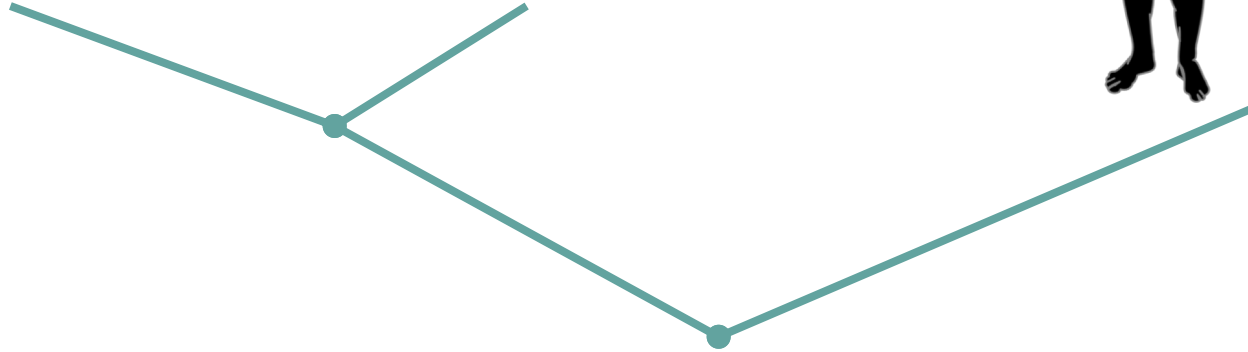
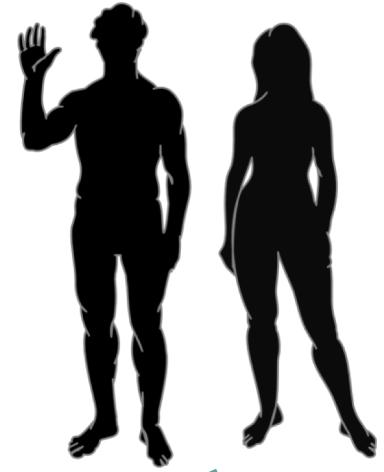
# Conflict and concordance in phylogenomic datasets

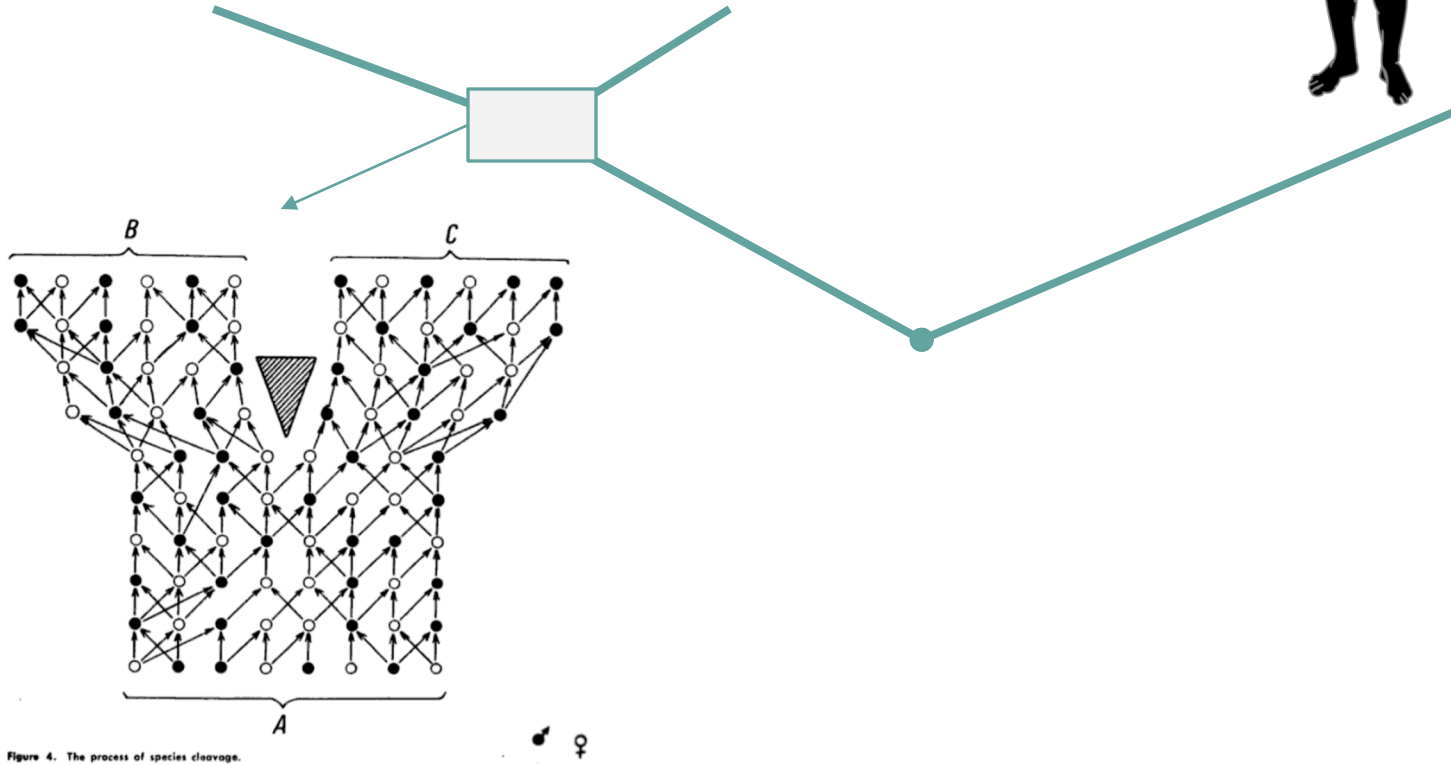
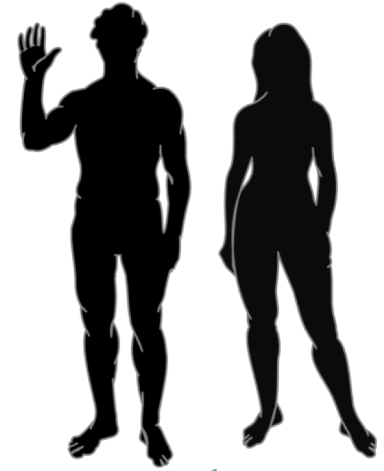
**DR. STEPHEN A. SMITH**  
UNIVERSITY OF MICHIGAN

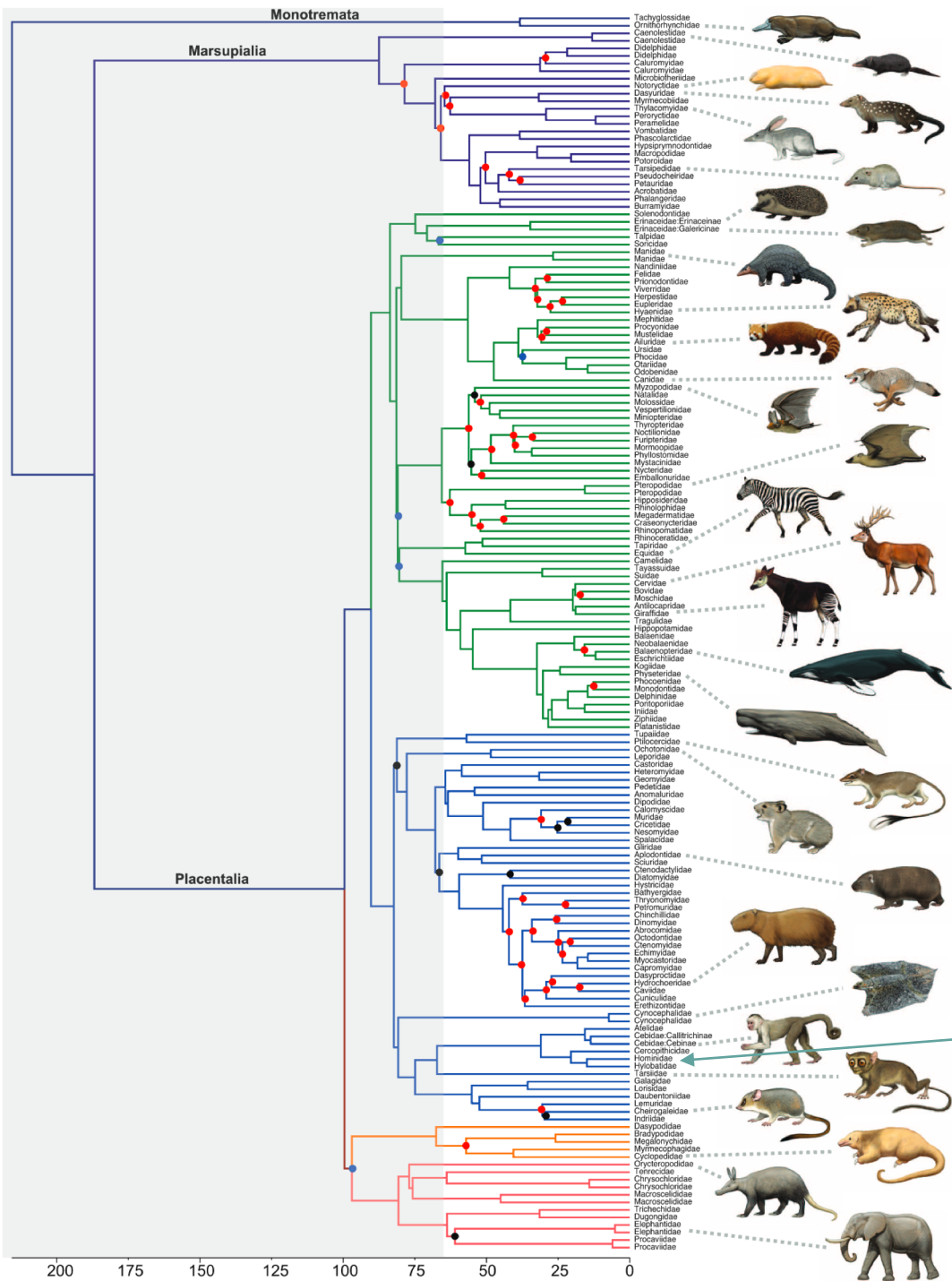


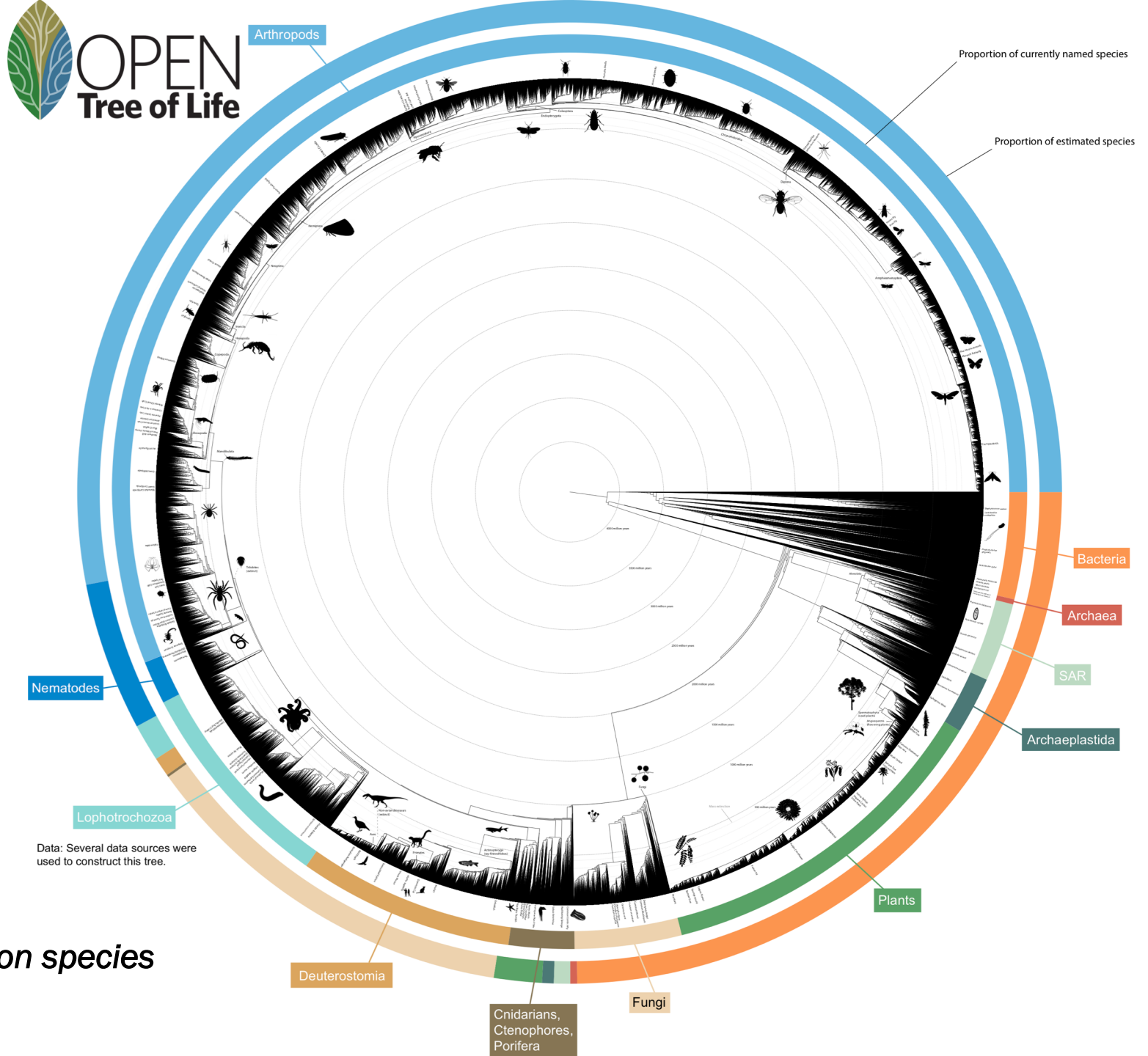












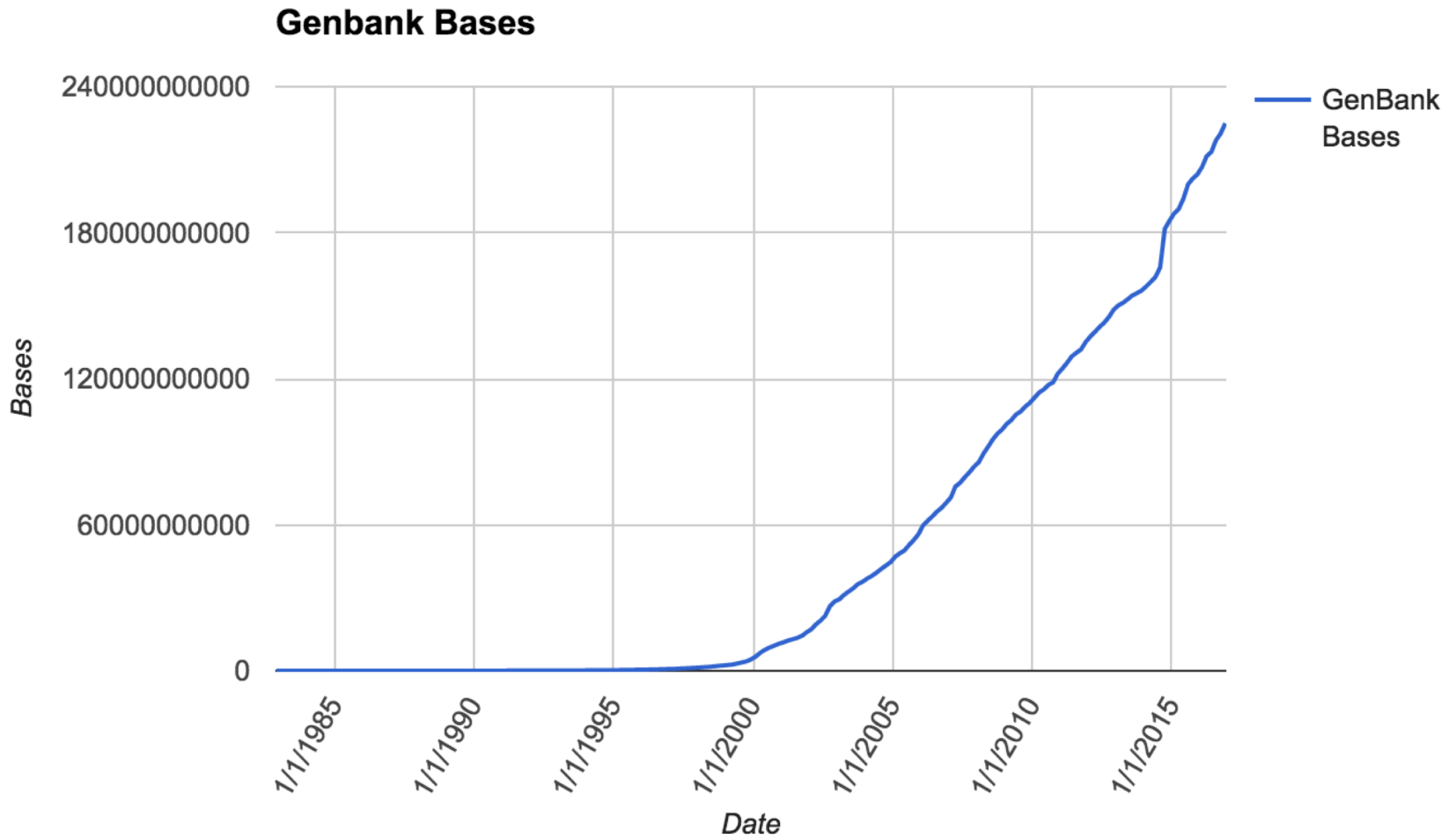
Data: Several data sources were used to construct this tree.

**~2.3 million species**

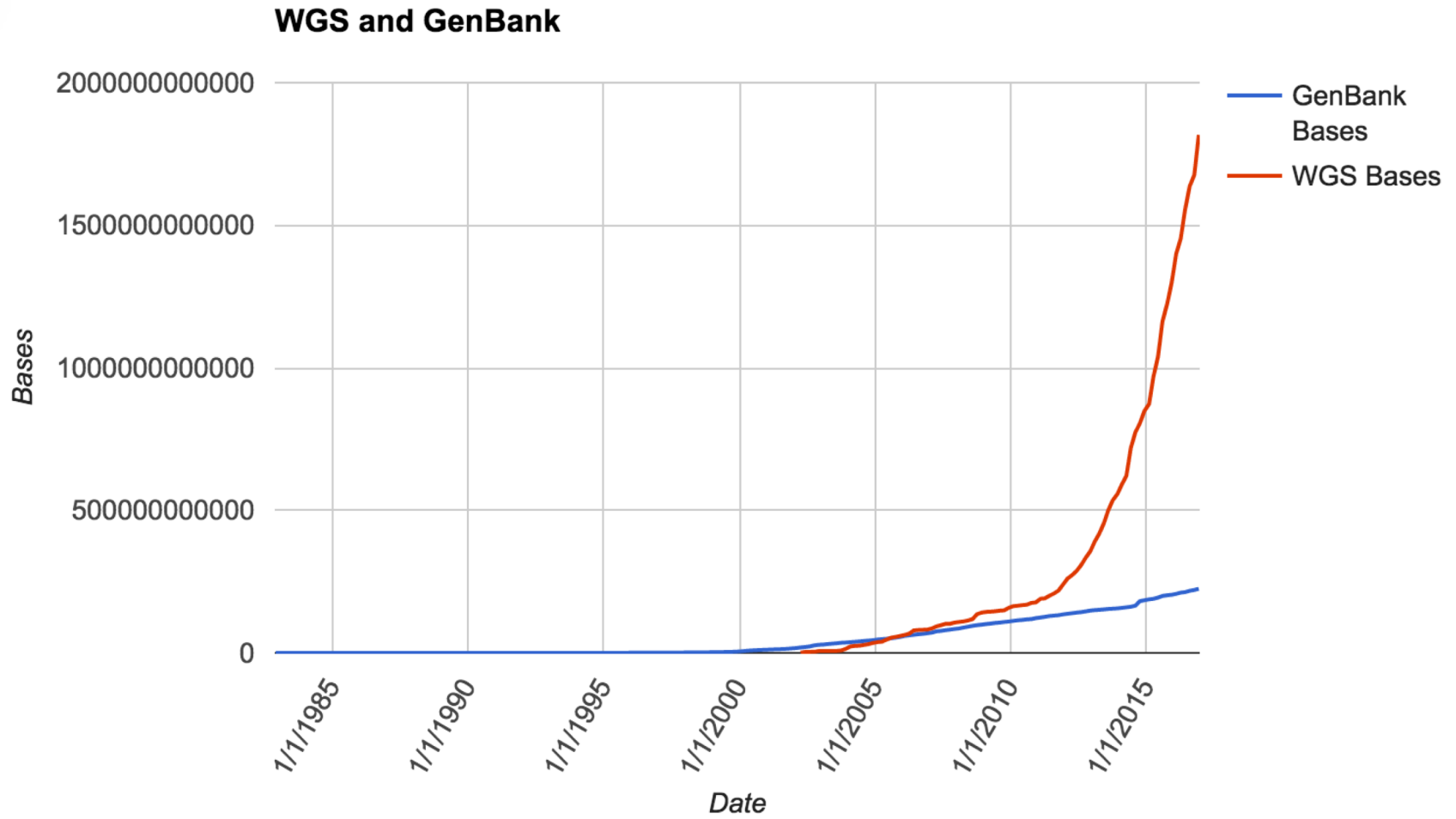




# Molecular data availability



# Molecular data availability (with whole genomes)



## Article Contents

Abstract

Supplementary data

ACCEPTED MANUSCRIPT

## 10KP: A Phylodiverse Genome Sequencing Plan

Shifeng Cheng , Michael Melkonian, Stephen A Smith, Samuel Brockington, John M Archibald, Pierre-Marc Delaux, Fay-wei Li, Barbara Melkonian, Evgeny V Mavrodiev, Wenjing Sun ... [Show more](#)

[Author Notes](#)

*GigaScience*, giy013, <https://doi.org/10.1093/gigascience/giy013>

**Published:** 20 February 2018



Freshwater algae in the genus *Zygnema* would be one target of sequencing project. NORBERT HÜLSMANN/FICKR (CC BY-NC-SA 2.0)

## Plant scientists plan massive effort to sequence 10,000 genomes

By **Dennis Normile** | Jul. 27, 2017, 8:00 AM

# Large datasets with many genes

## Typical phylogenetic analyses

- 1-10 genes
- 17 genes. Plants (Soltis et al. 2011)
- 19 genes. Birds (Hackett et al. 2007)

## Transcriptomic and genomic phylogenetic analyses

- 140 genes. Metazoa (Dunn et al. 2008)
- 242 genes. Metazoa (Ryan et al. 2013)
- 248 genes. Turtles (Chiari et al., 2012)
- 1185 genes. Molluscs (Smith et al. 2011)
- 1720 genes. Rice (Cranston et al. 2007)
- 2970 genes. Seed plants (Lee et al. 2011)
- >8000 genes. Birds (Jarvis et al. 2014)
- 259 genes. Birds (Prum et al. 2015)
- 859 genes. Seed plants (Wickett et al. 2014)

# Concatenate genes to get more information

Gene 1

Gene 2

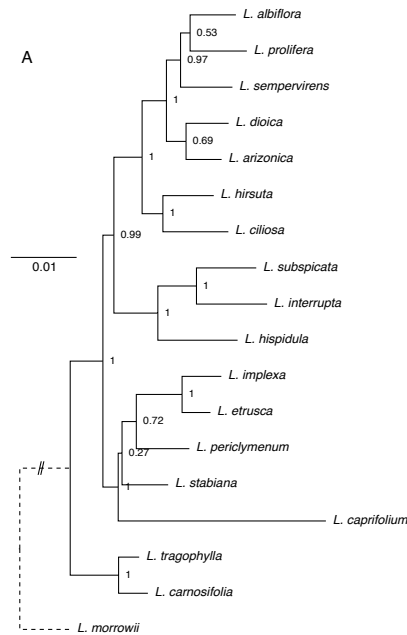
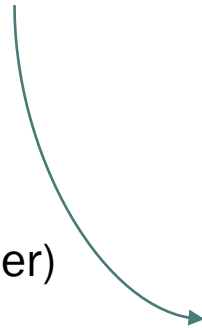
Gene 3

Gene 4

Gene 5

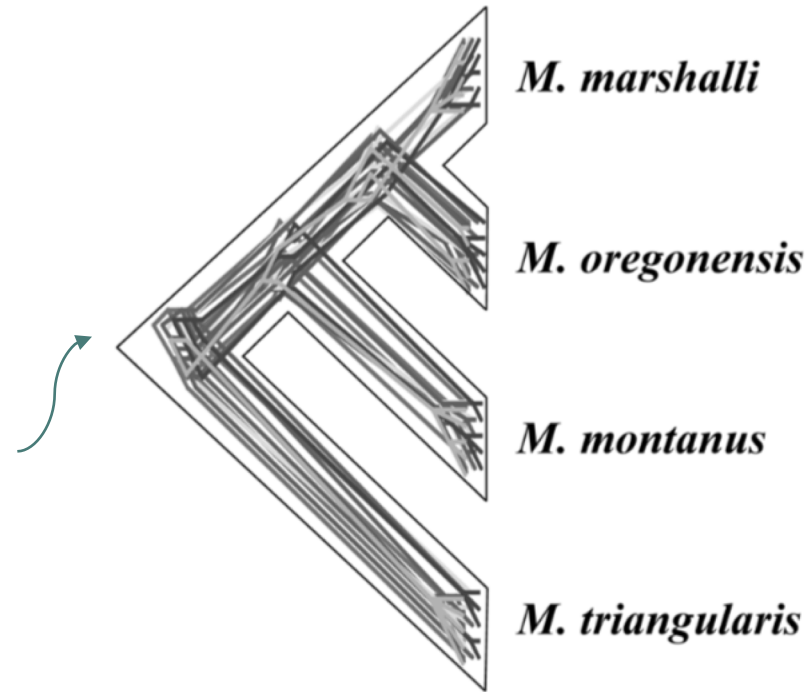
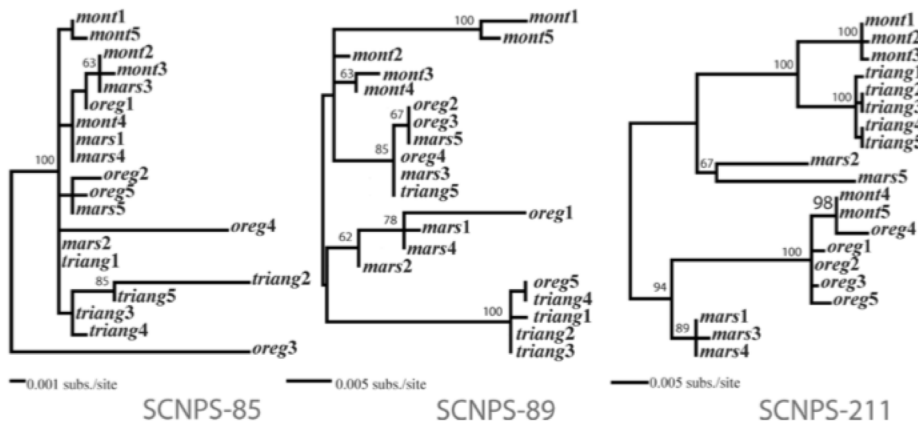
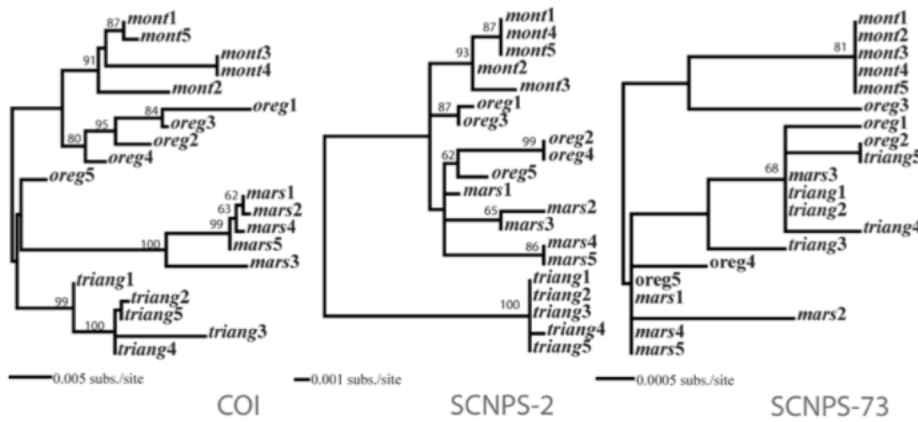


More information  
(all the genes together)



# Combining gene trees (using conflict as information)

Each gene contributes individually



Assume one underlying tree

# Triumphant phylogenomics

---



EXPLORING UNDERLYING MOLECULAR  
PATTERNS AND PROCESSES

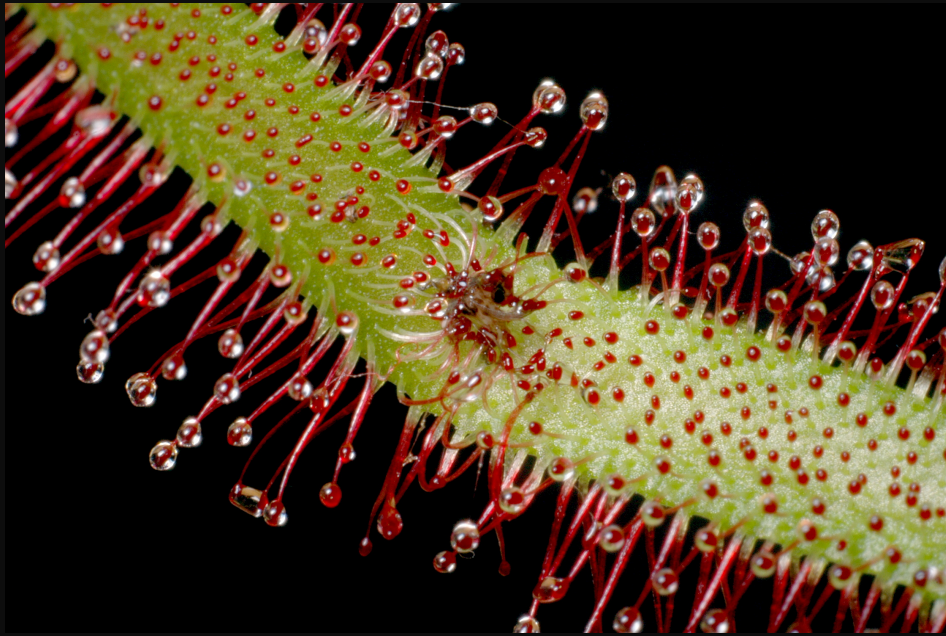
# Caryophyllales

- >12,500 species in 39 families
- extreme disparity in life history and ecology

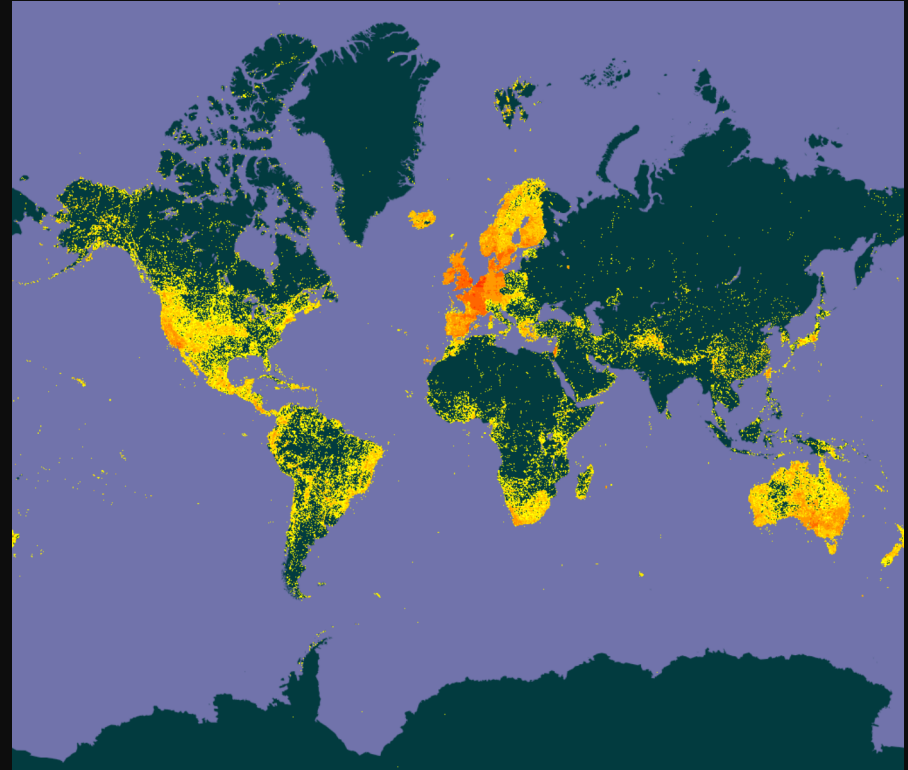




# Carnivory



# Cold environments



# Photosynthetic modifications

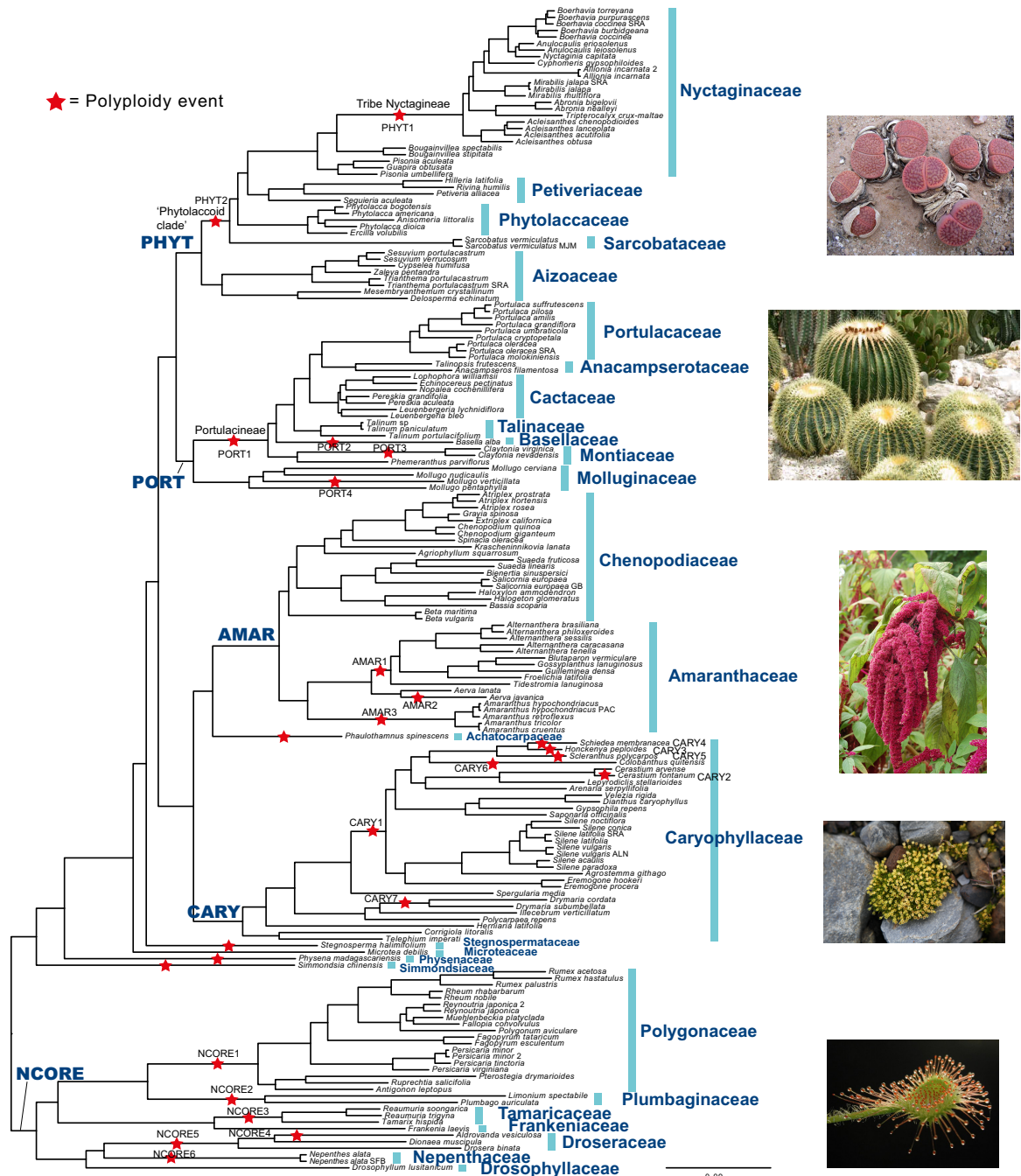


# Morphological modifications



# Genome duplications in Caryophyllales

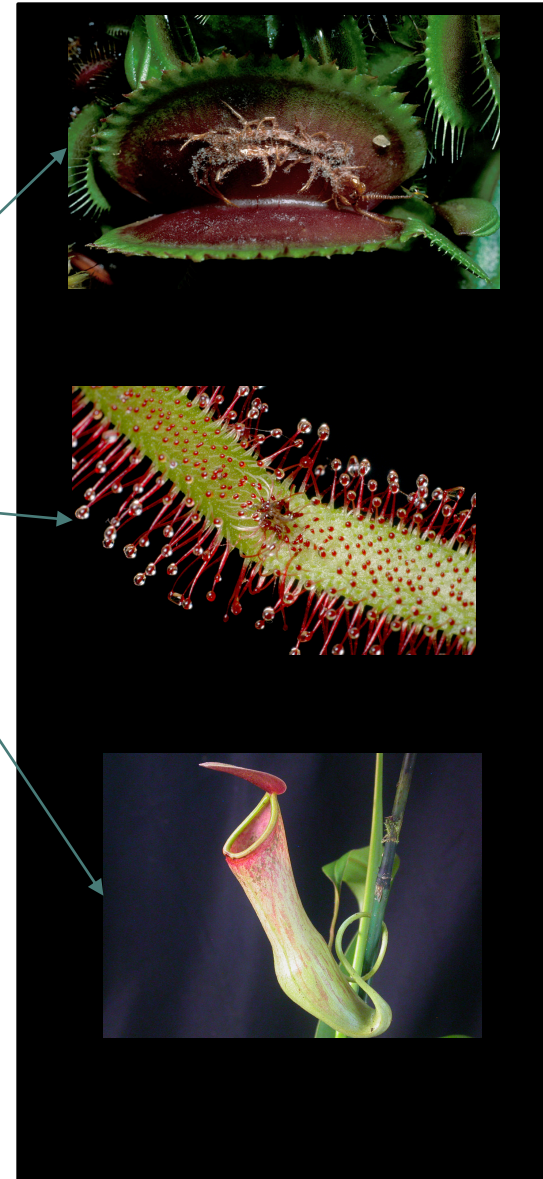
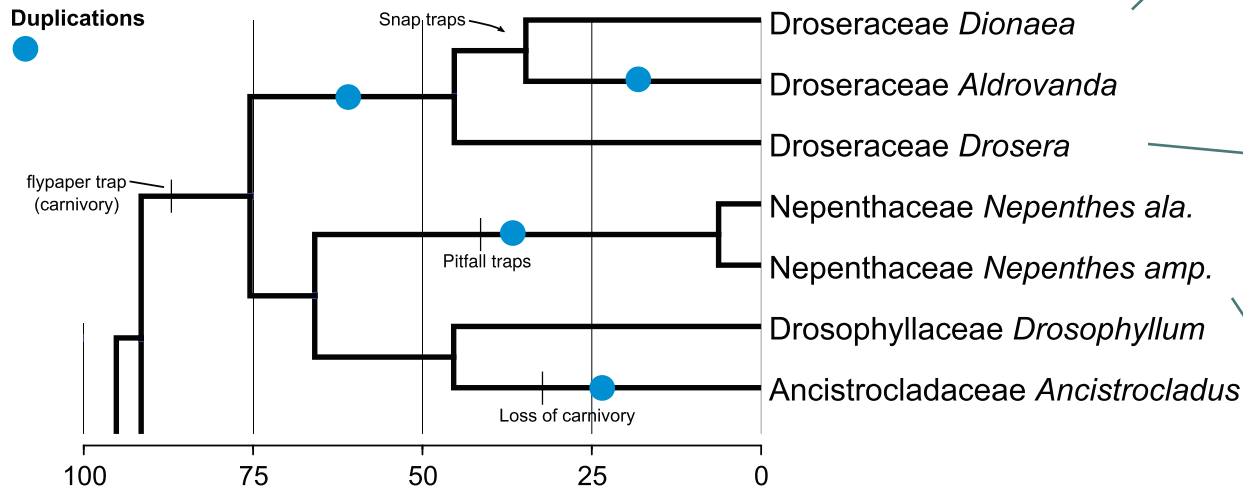
- 168 species
- 736 gene regions used for inference (from thousands identified)
- 26 duplications
- Yang et al. 2017



# Duplications in carnivores

Walker et al. 2017

4 duplications associated with the evolution of carnivory

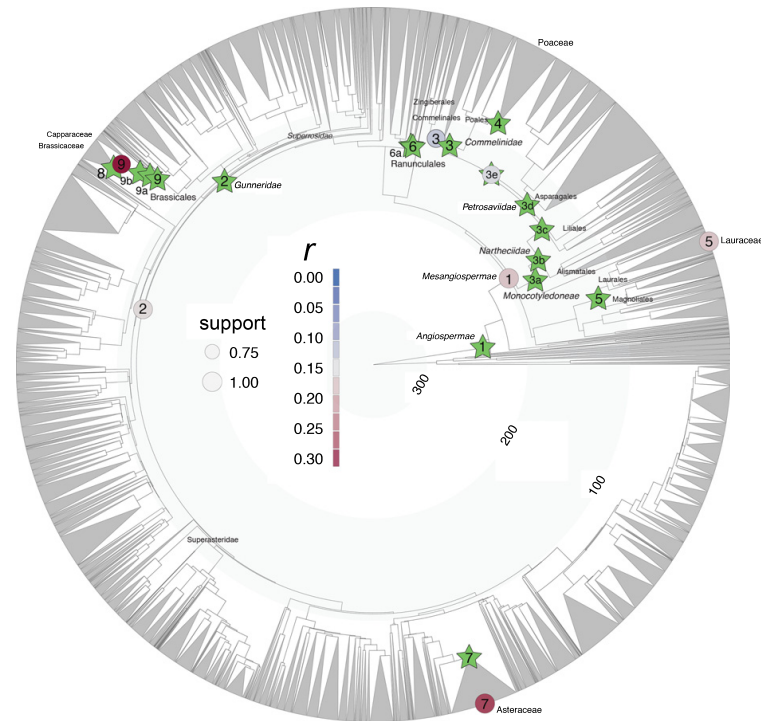


# Duplications and diversification

- It has been suggested that **diversification** is associated with **genome duplications**

*Are there diversification shifts associated with genome duplication?*

*Are there climatic shifts associated with genome duplication?*



# Caryophyllales: annual mean temperature



5036 taxa

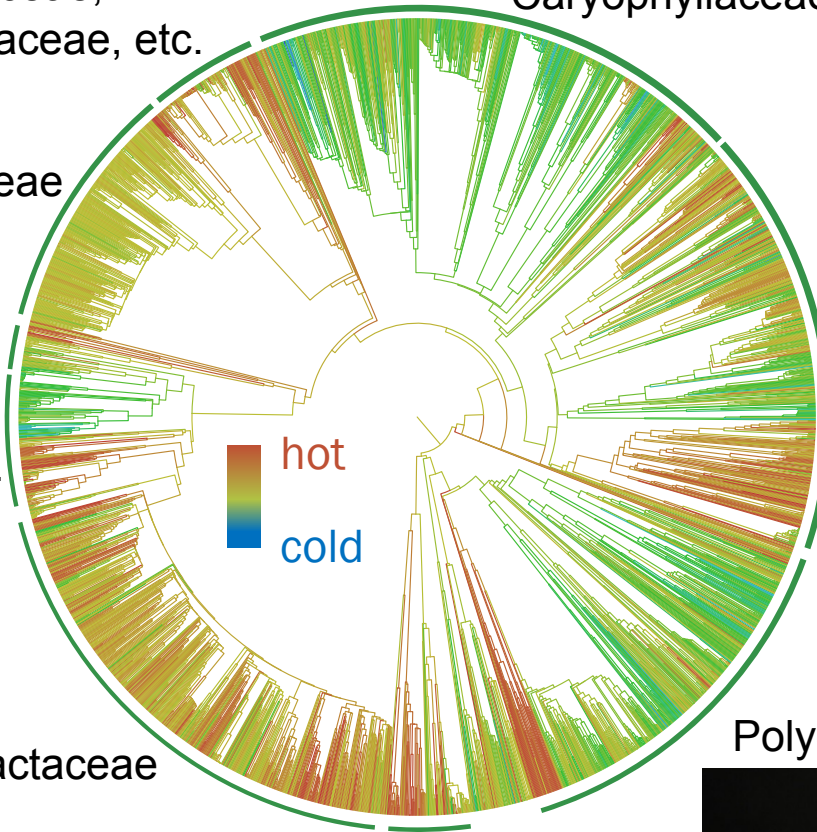
Amaranthaceae



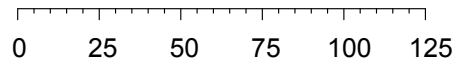
Polygonaceae



Carnivorous



hot  
cold



Nyctaginaceae,  
Phytolaccaceae, etc.



Aizoaceae

Molluginaceae

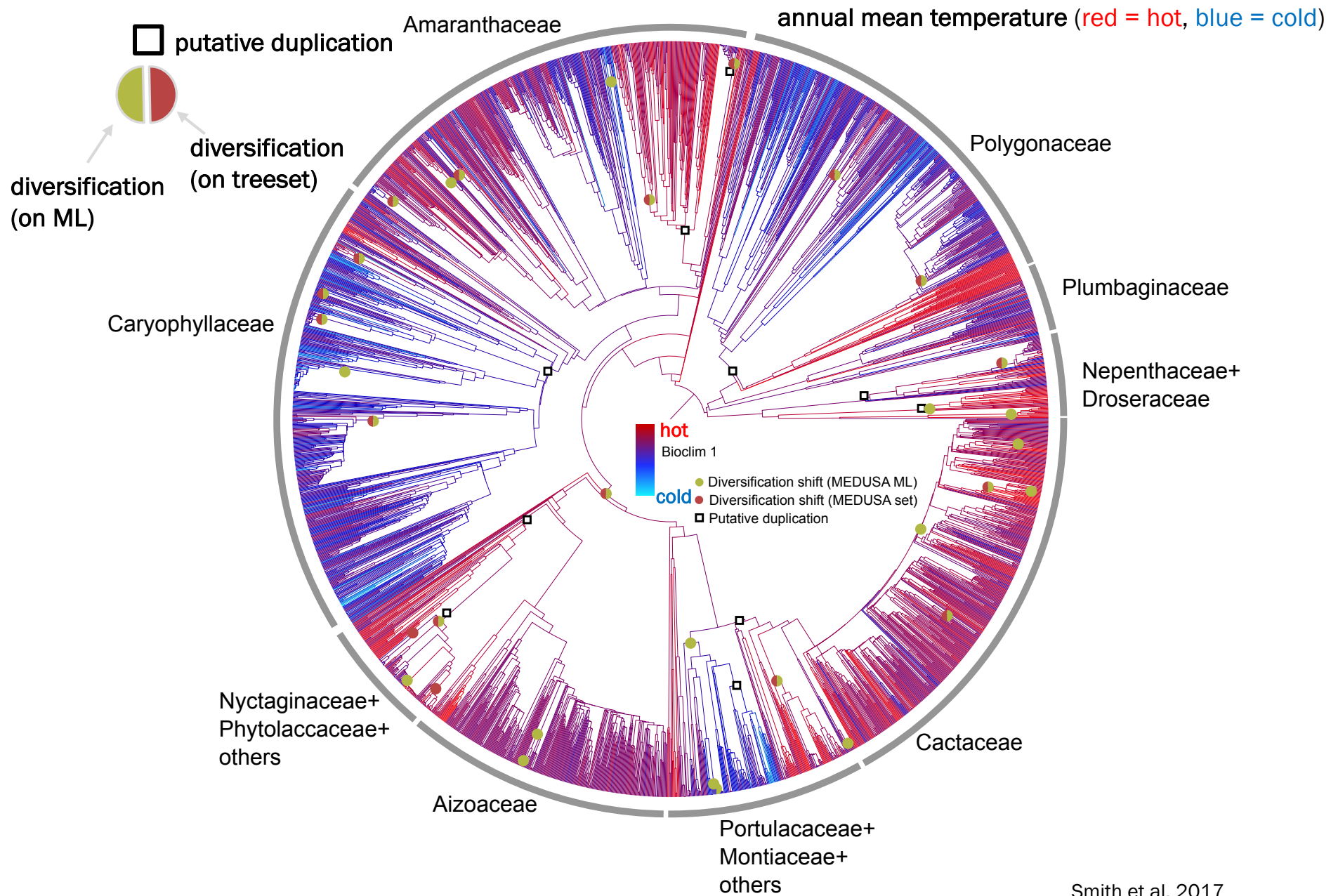
Portulacaceae,  
Montiaceae, etc.



Cactaceae

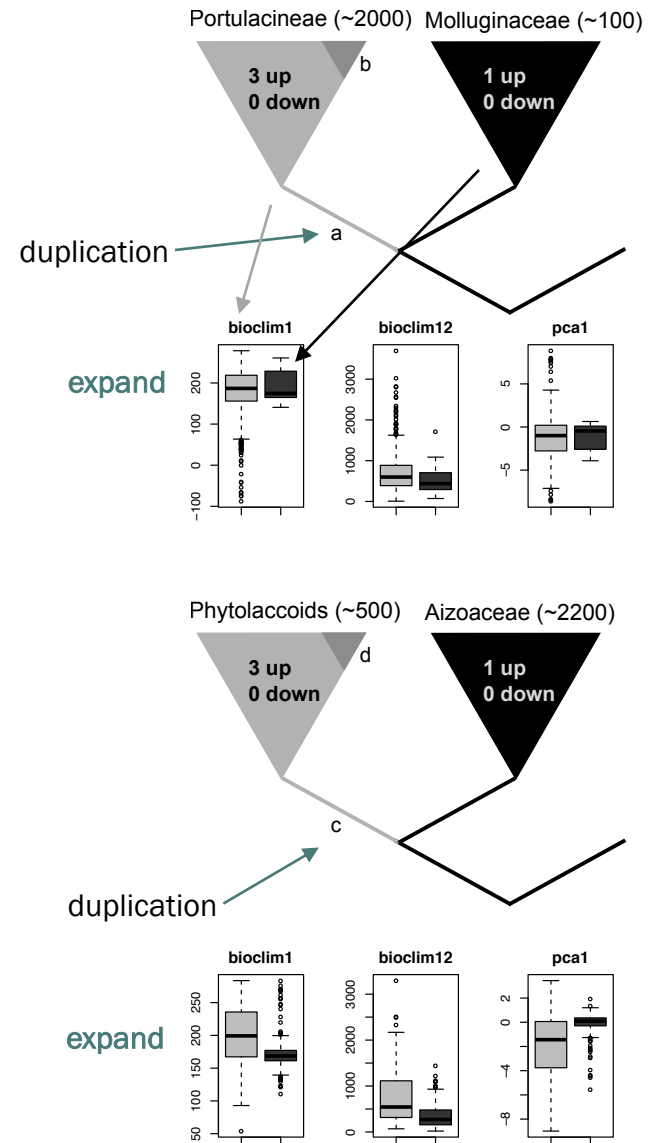






# Summary of biological results

- Duplications **are not** associated with faster speciation
- Duplications occur **before** increases in speciation
- *Many duplications associated with expanded climate ranges*

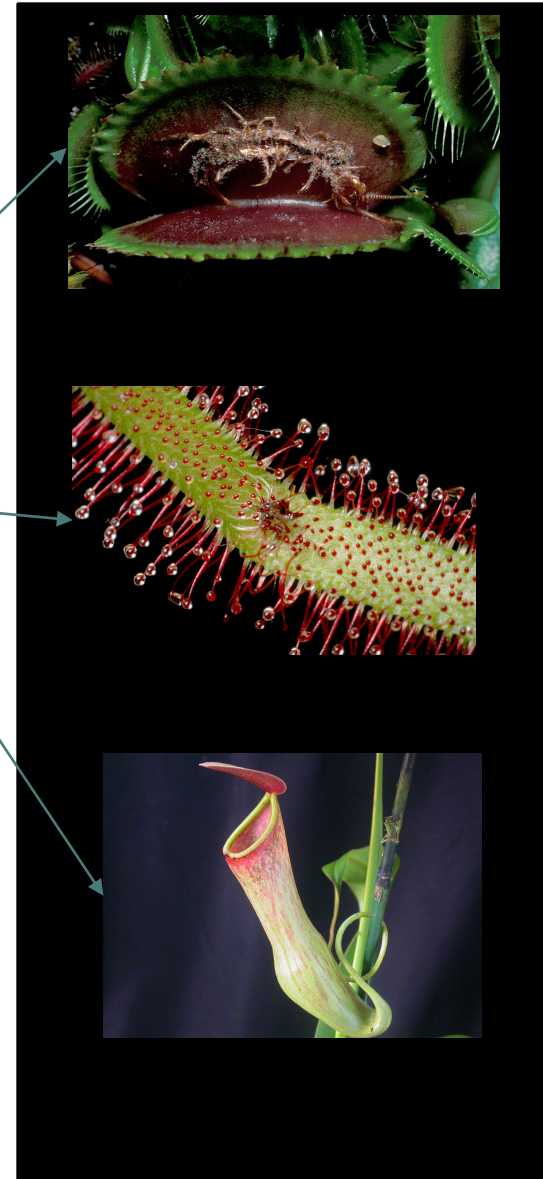
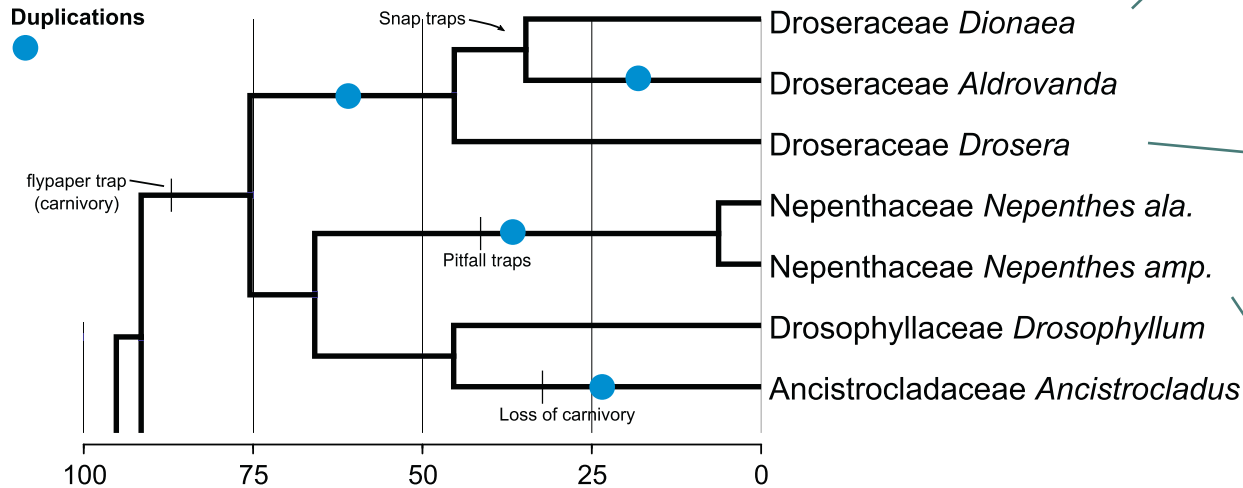


**There are cracks on the  
horizon...**

# Duplications in carnivores

Walker et al. 2017

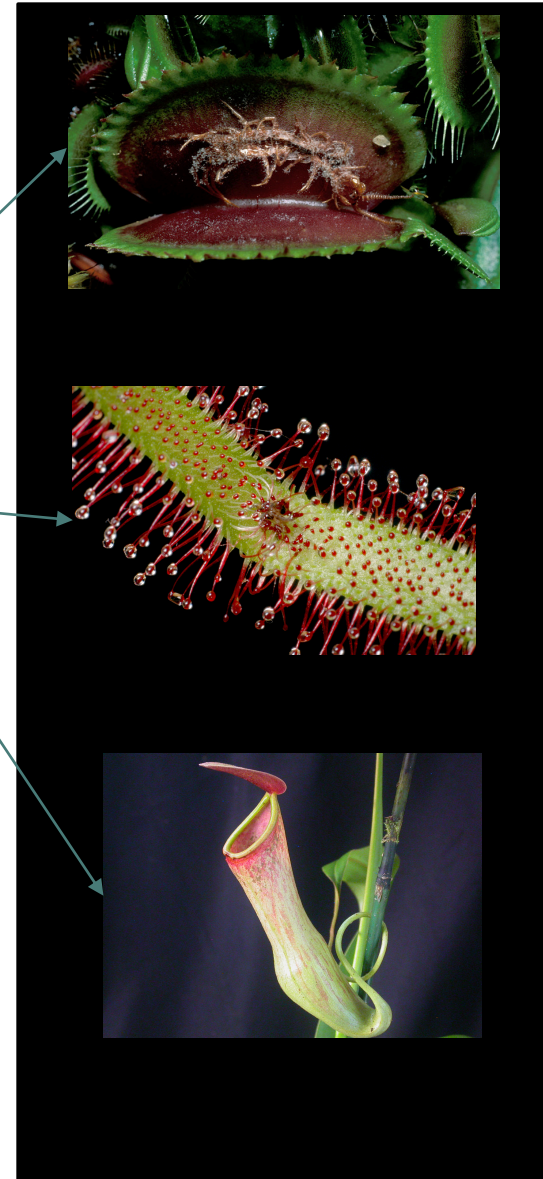
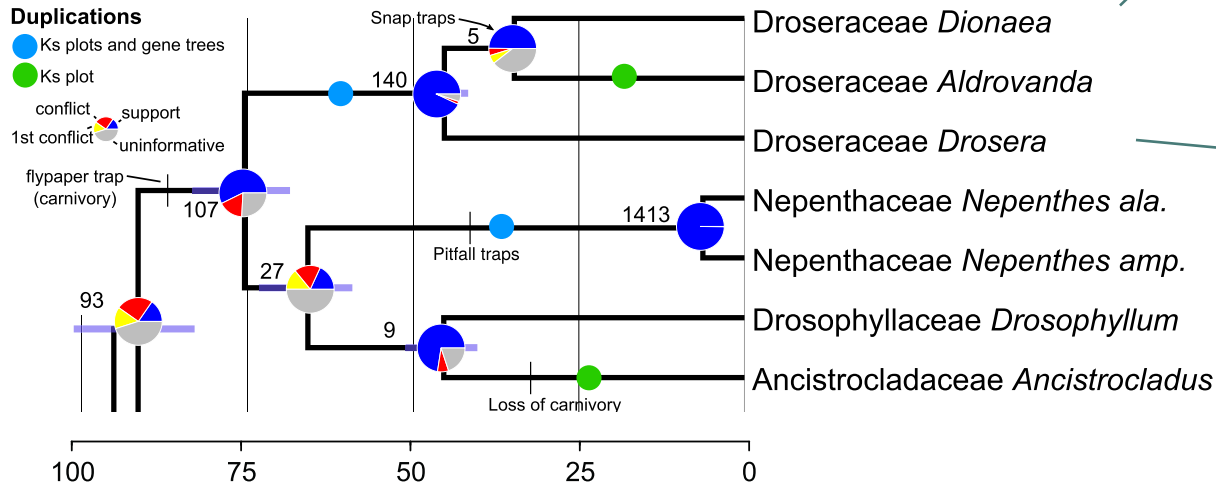
7 duplications associated with the evolution of carnivory

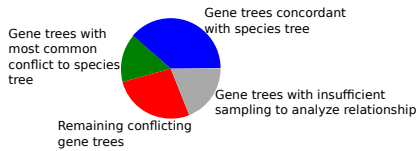


# Duplications in carnivores

Walker et al. 2017

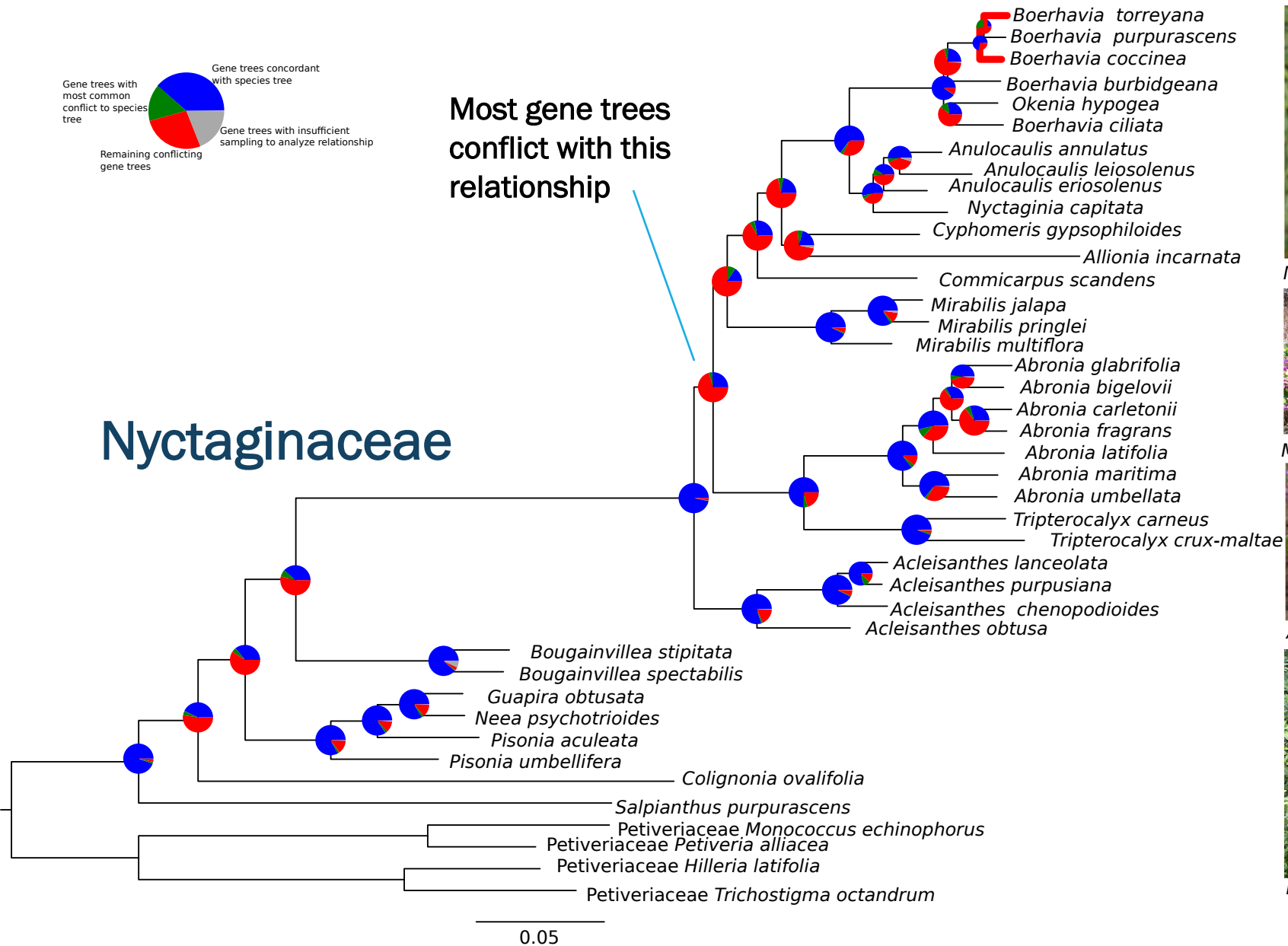
7 duplications associated with the evolution of carnivory





Most gene trees conflict with this relationship

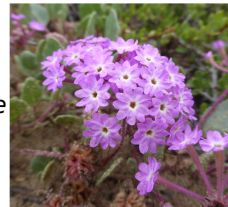
# Nyctaginaceae



Nyctaginia capitata



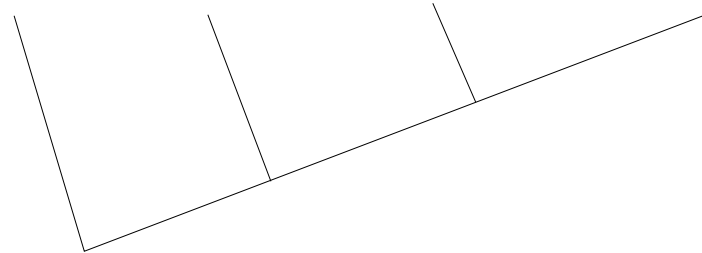
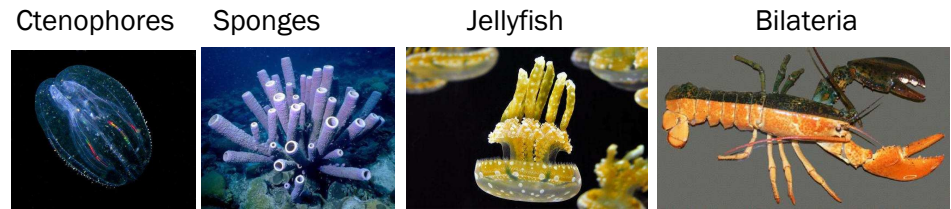
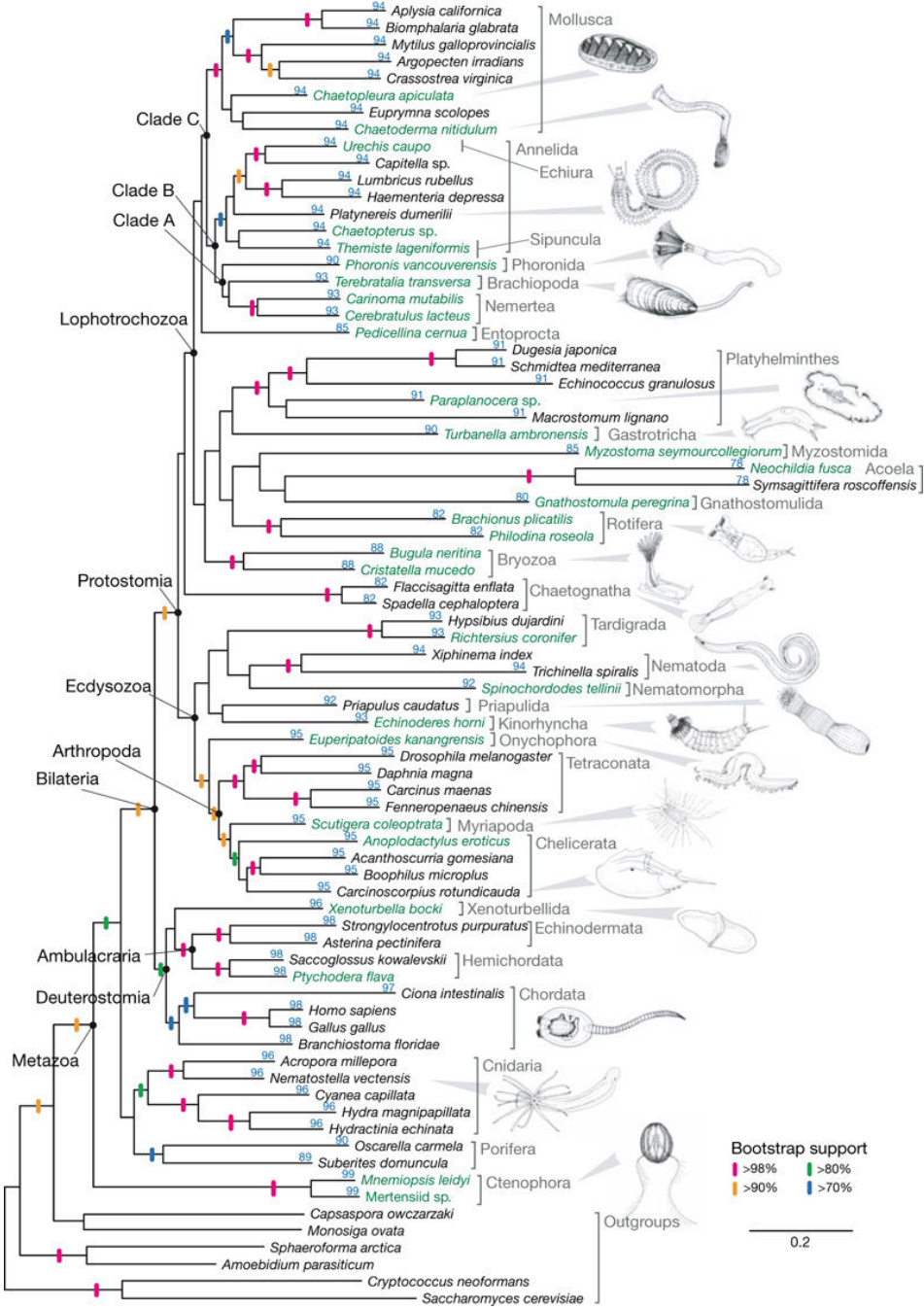
Mirabilis multiflora



Abronia umbellata



Pisonia umbellifera



Dunn et al. 2008

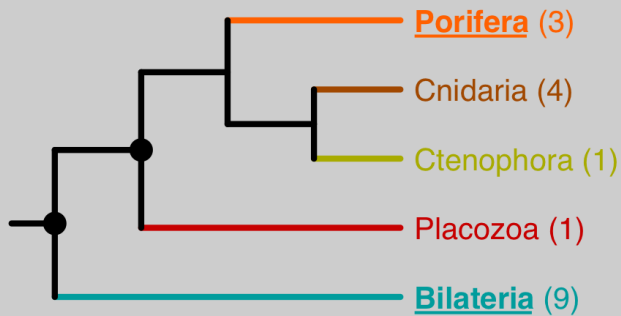
# Alternative views

Just showed you this one



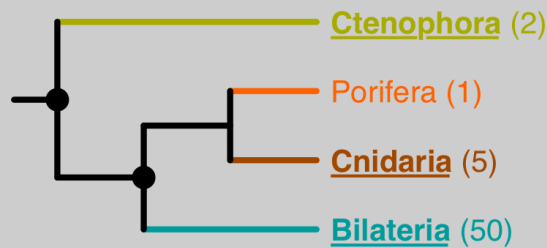
## Schierwater et al. (2009)

15 mitochondrial and 34 nuclear genes  
GTR and WAG models



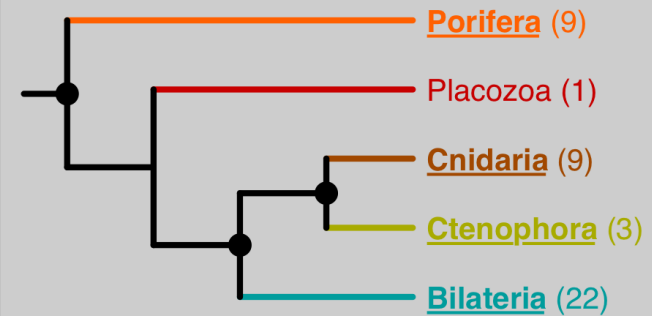
## Dunn et al. (2008)

6 mitochondrial and 144 nuclear genes  
WAG and CAT models

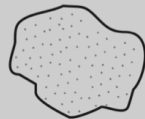


## Philippe et al. (2009)

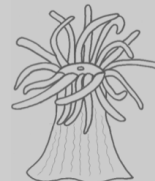
128 nuclear genes  
CAT model



Porifera



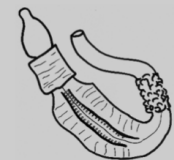
Placozoa



Cnidaria



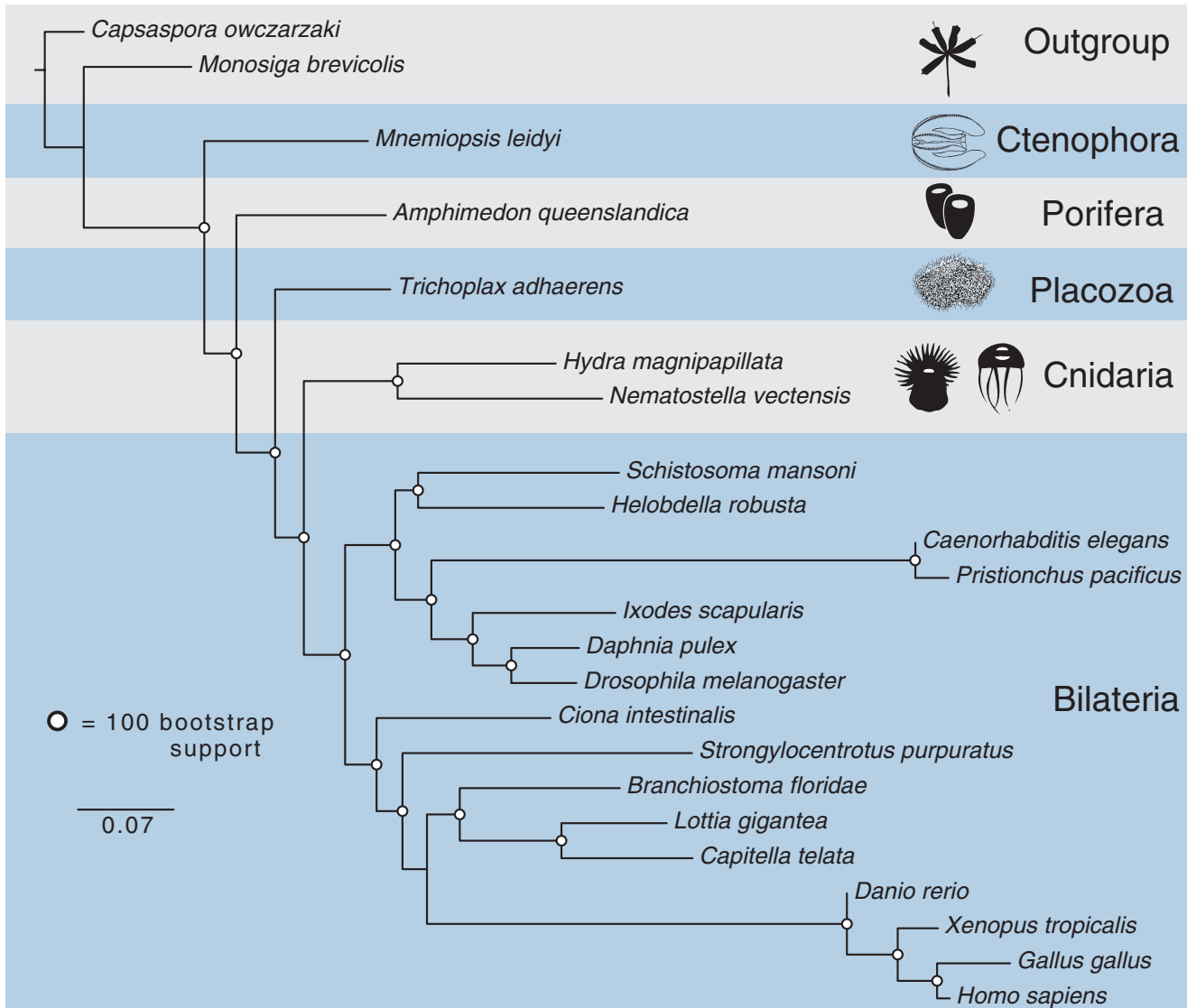
Ctenophora



Bilateria



# Metazoan phylogeny using 242 genes



← Sequenced the whole genome

*Is data the problem?*

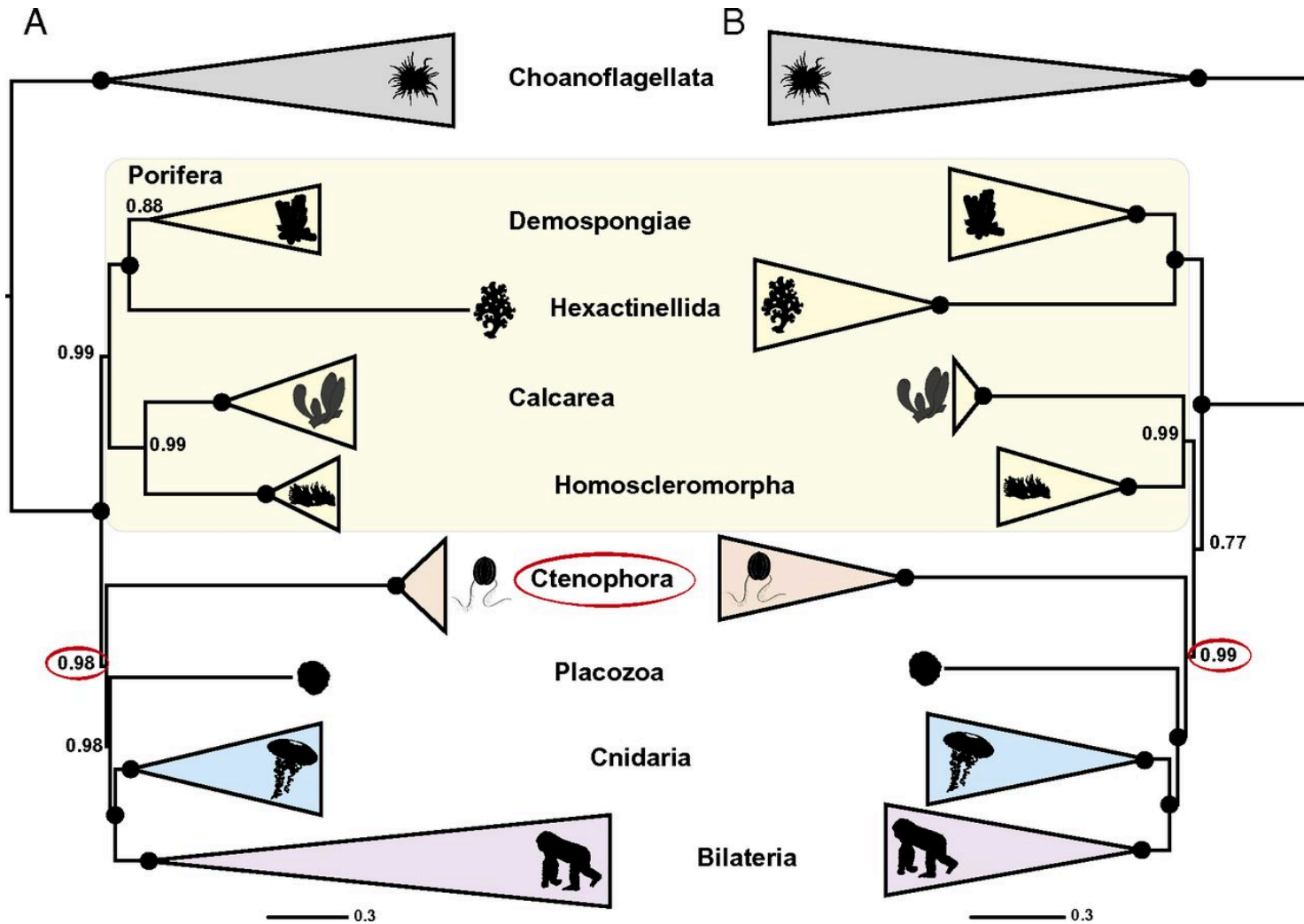
Ryan et al. 2013

# Genomic data do not support comb jellies as the sister group to all other animals

Davide Pisani, Walker Pett, Martin Dohrmann, Roberto Feuda, Omar Rota-Stabelli, Hervé Philippe, Nicolas Lartillot, and Gert Wörheide

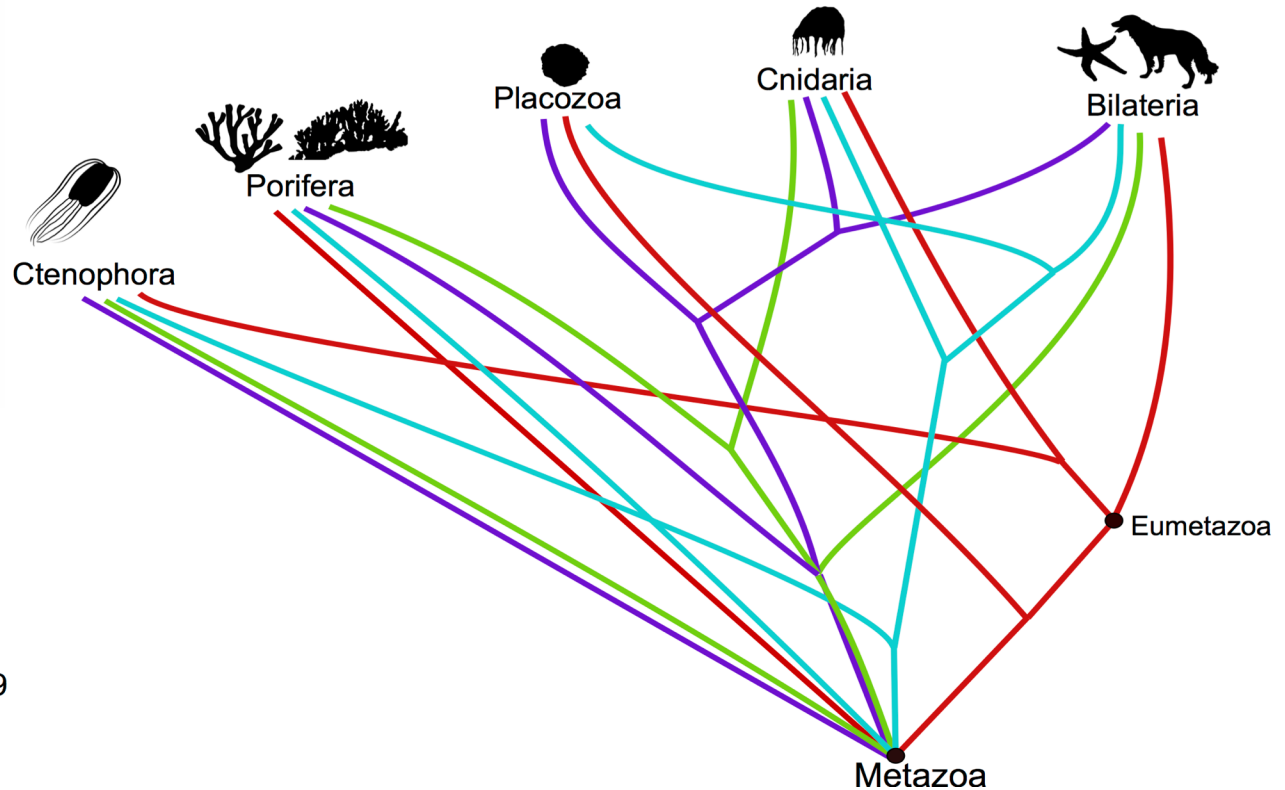
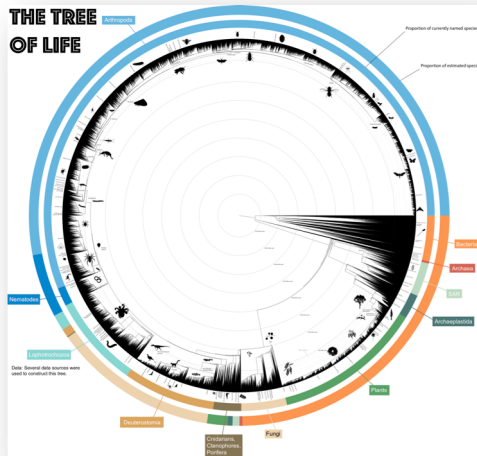
PNAS 2015 December, 112 (50) 15402-15407. <https://doi.org/10.1073/pnas.1518127112>

*Are molecular models the problem?*



# Building a more complete view of the Tree of Life

Open Tree of Life tried to accommodate this but is there a way to resolve any of this *confidently*?



Dunn et al., 2008

Pick et al. 2010

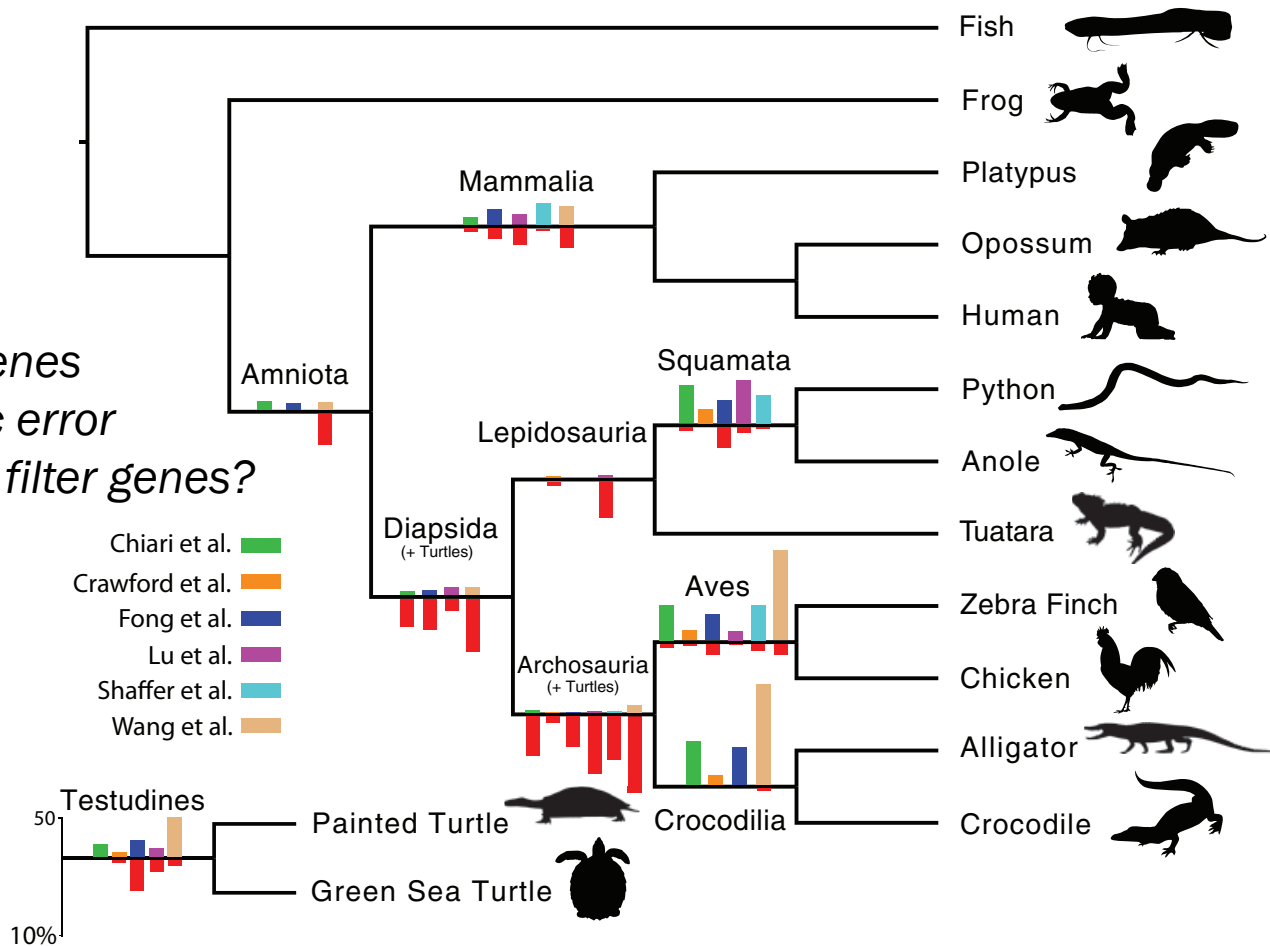
Ryan et al. 2013

Phillipe et al. 2009

## Bayes Factors Unmask Highly Variable Information Content, Bias, and Extreme Influence in Phylogenomic Analyses

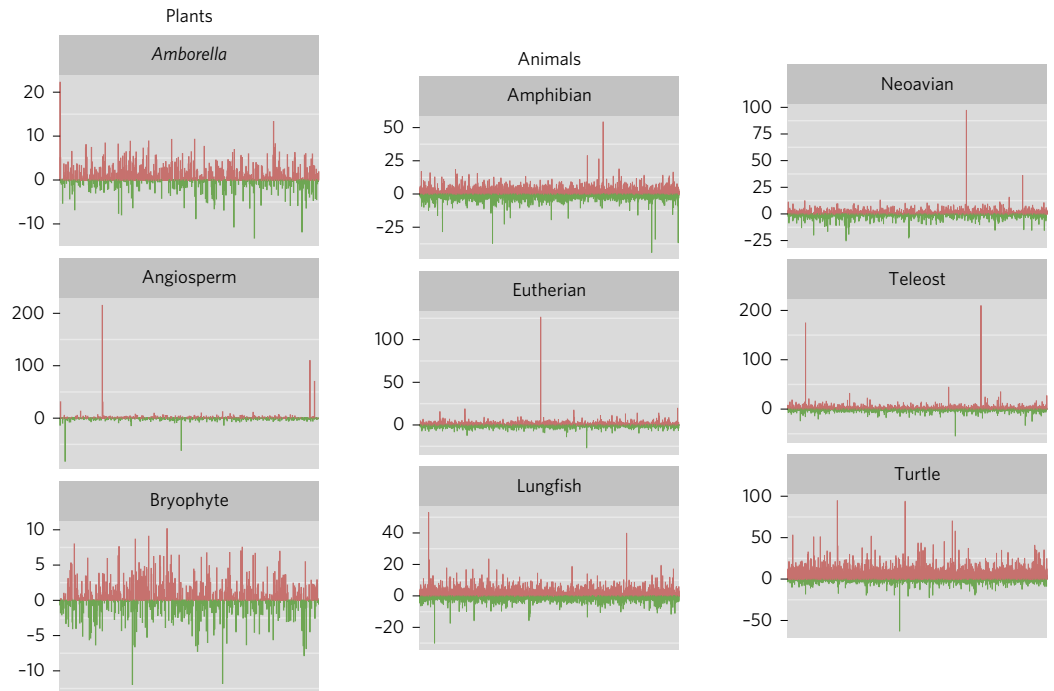
JEREMY M. BROWN<sup>1,\*</sup> AND ROBERT C. THOMSON<sup>2</sup>

- *Outlying genes*
- *Systematic error*
- *Should we filter genes?*



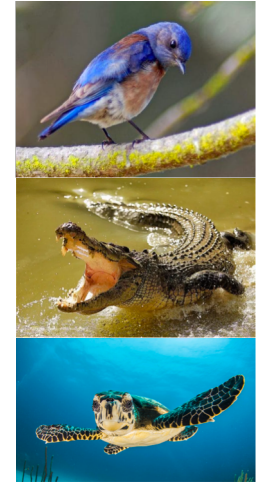
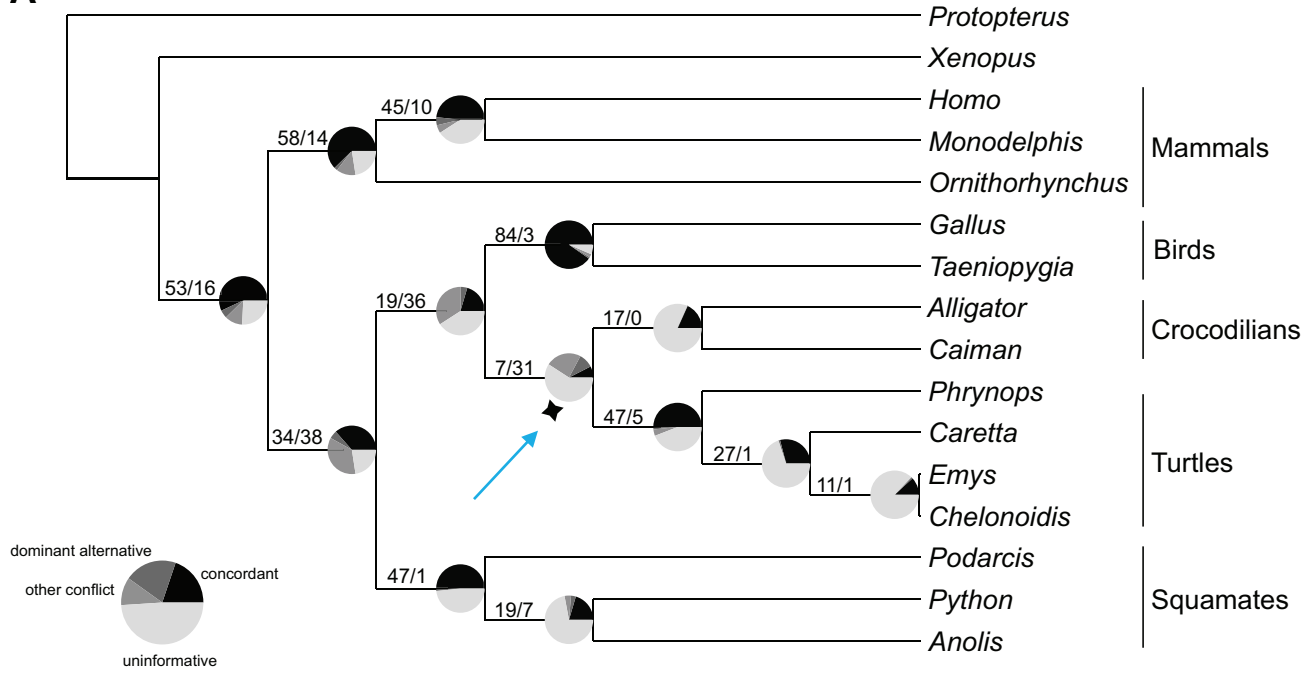
# Contentious relationships in phylogenomic studies can be driven by a handful of genes

Xing-Xing Shen<sup>1</sup>, Chris Todd Hittinger<sup>2</sup> and Antonis Rokas<sup>1\*</sup>

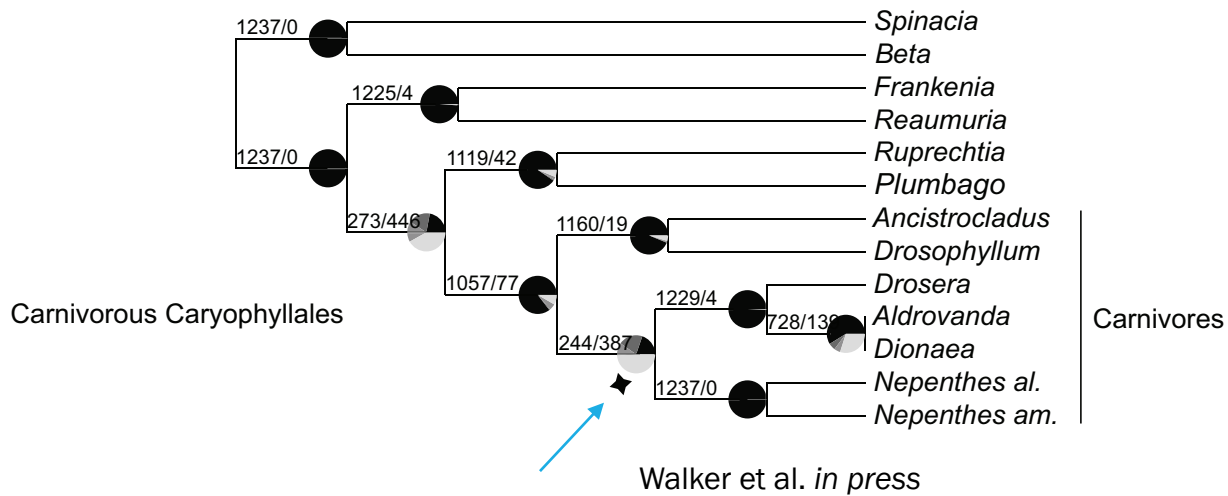


- Outlying genes
- Should we do two (or few) topology tests?

**A**



**B**

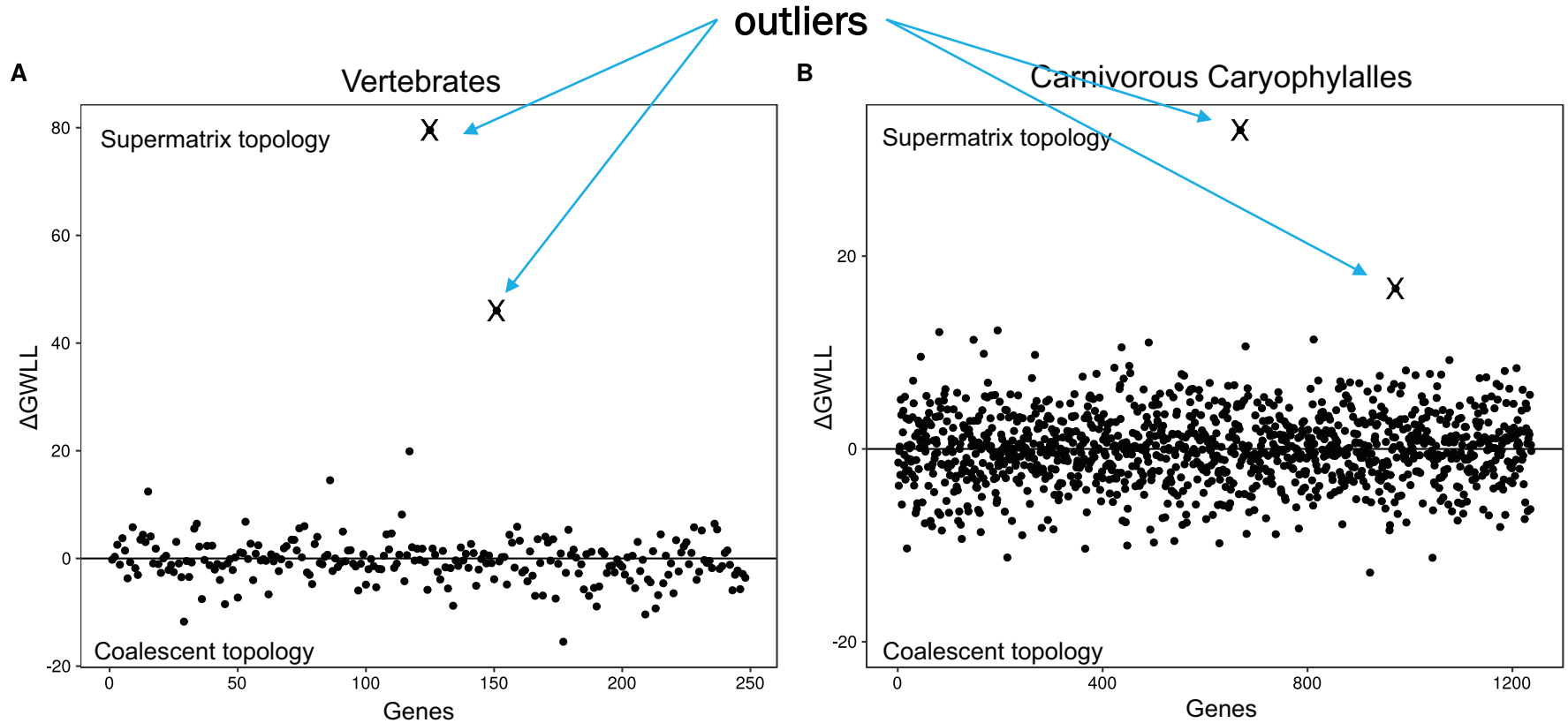


Coalescent topologies were favored for biological reasons

Supermatrix topologies were assumed to reflect error

Two topology comparisons exposed “outlier” genes

When removed, supermatrix topologies match coalescent topologies



Walker et al. *in press*

# Important lessons

Just concatenation and gene tree / species tree (ASTRAL) are probably not going to cut it

## Systematic error

- Researchers need to be significantly more careful about the underlying data
- Check for errors
  - in homology/orthology
  - Alignment
  - Heterogeneity in molecular evolution

## Biological sources of error

- Gene duplication and loss
- ILS
- Hybridization

## Potential “outlier” genes

- One or a few genes can drive phylogenetic inference
- Especially problematic with concatenation

Can we take a different approach to phylogenomic data analysis?



# Edge-based look at plants

Wickett et al. 2014 initial 1KP paper

Conducted ASTRAL and supermatrix analyses on many different datasets

ASTRAL and supermatrix trees are largely congruent but disagree

This didn't settle all the arguments

## Phylotranscriptomic analysis of the origin and early diversification of land plants

Norman J. Wickett<sup>a,b,1,2</sup>, Siavash Mirarab<sup>c,1</sup>, Nam Nguyen<sup>c</sup>, Tandy Warnow<sup>c</sup>, Eric Carpenter<sup>d</sup>, Naim Matasci<sup>e,f</sup>, Saravanaraj Ayyampalayam<sup>g</sup>, Michael S. Barker<sup>f</sup>, J. Gordon Burleigh<sup>h</sup>, Matthew A. Gitzendanner<sup>h,i</sup>, Brad R. Ruhfel<sup>h,j,k</sup>, Eric Wafula<sup>l</sup>, Joshua P. Der<sup>l</sup>, Sean W. Graham<sup>m</sup>, Sarah Mathews<sup>n</sup>, Michael Melkonian<sup>o</sup>, Douglas E. Soltis<sup>h,i,k</sup>, Pamela S. Soltis<sup>h,i,k</sup>, Nicholas W. Miles<sup>k</sup>, Carl J. Rothfels<sup>p,q</sup>, Lisa Pokorny<sup>p,r</sup>, A. Jonathan Shaw<sup>p</sup>, Lisa DeGironimo<sup>s</sup>, Dennis W. Stevenson<sup>s</sup>, Barbara Surek<sup>o</sup>, Juan Carlos Villarreal<sup>t</sup>, Béatrice Roure<sup>u</sup>, Hervé Philippe<sup>u,v</sup>, Claude W. dePamphilis<sup>l</sup>, Tao Chen<sup>w</sup>, Michael K. Deyholos<sup>d</sup>, Regina S. Baucom<sup>x</sup>, Toni M. Kutchan<sup>y</sup>, Megan M. Augustin<sup>y</sup>, Jun Wang<sup>z</sup>, Yong Zhang<sup>y</sup>, Zhijian Tian<sup>z</sup>, Zhixiang Yan<sup>z</sup>, Xiaolei Wu<sup>z</sup>, Xiao Sun<sup>z</sup>, Gane Ka-Shu Wong<sup>d,z,aa,2</sup>, and James Leebens-Mack<sup>g,2</sup>

### Sister to land plants



Zygn-sister  
(Timme et al 2012)



Char-sister  
(Karol et al 2001)

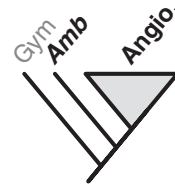


Col-sister  
(Finet et al 2010)

### ANA-grade angiosperms



*Amborella* + *Nuphar*  
(Qiu et al 2010)



Amb-sister  
(Soltis et al 2011)

### Bryophytes



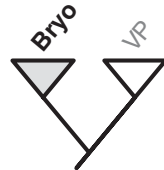
Lv-basal  
(Qiu et al 1998)



Hw-sister  
(Chang & Graham 2011)



Moss + Lv  
(Renzaglia et al 2000)



Bryo monophy.  
(Nishiyama et al 2004)



Hw-basal  
(Nickrent et al 2000)

### Gymnosperms



Gnepine  
(Bowe et al 2000)



Conifers monophy.  
(Lee et al 2011)



Gnetifer  
(Chaw et al 1997)



Gnetales-sister  
(Lee et al 2011)

### Angiosperms



Eudi + Mag  
(Lee et al 2011)



Eudi + Mag/Chl  
(Burleigh et al 2009)



Mono + Eudi  
(Soltis et al 2011)



Mag + Chl, Mono  
(Soltis et al 2000)

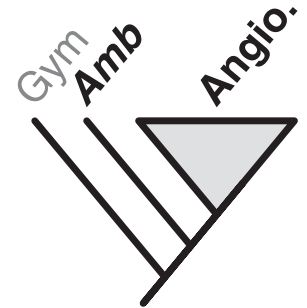


Mag + Chl  
(Soltis et al 2011)

### ANA-grade angiosperms

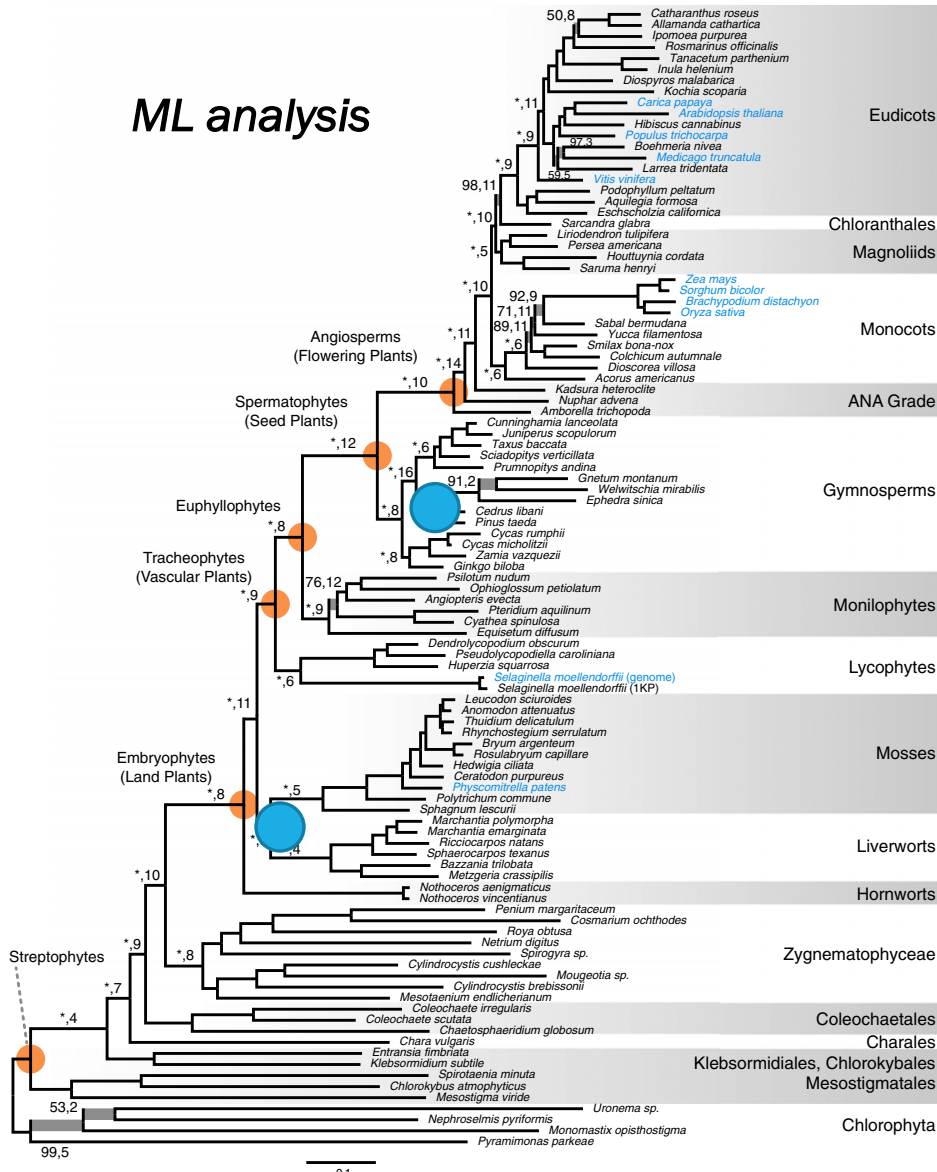


*Amborella* + *Nuphar*  
(Qiu et al 2010)

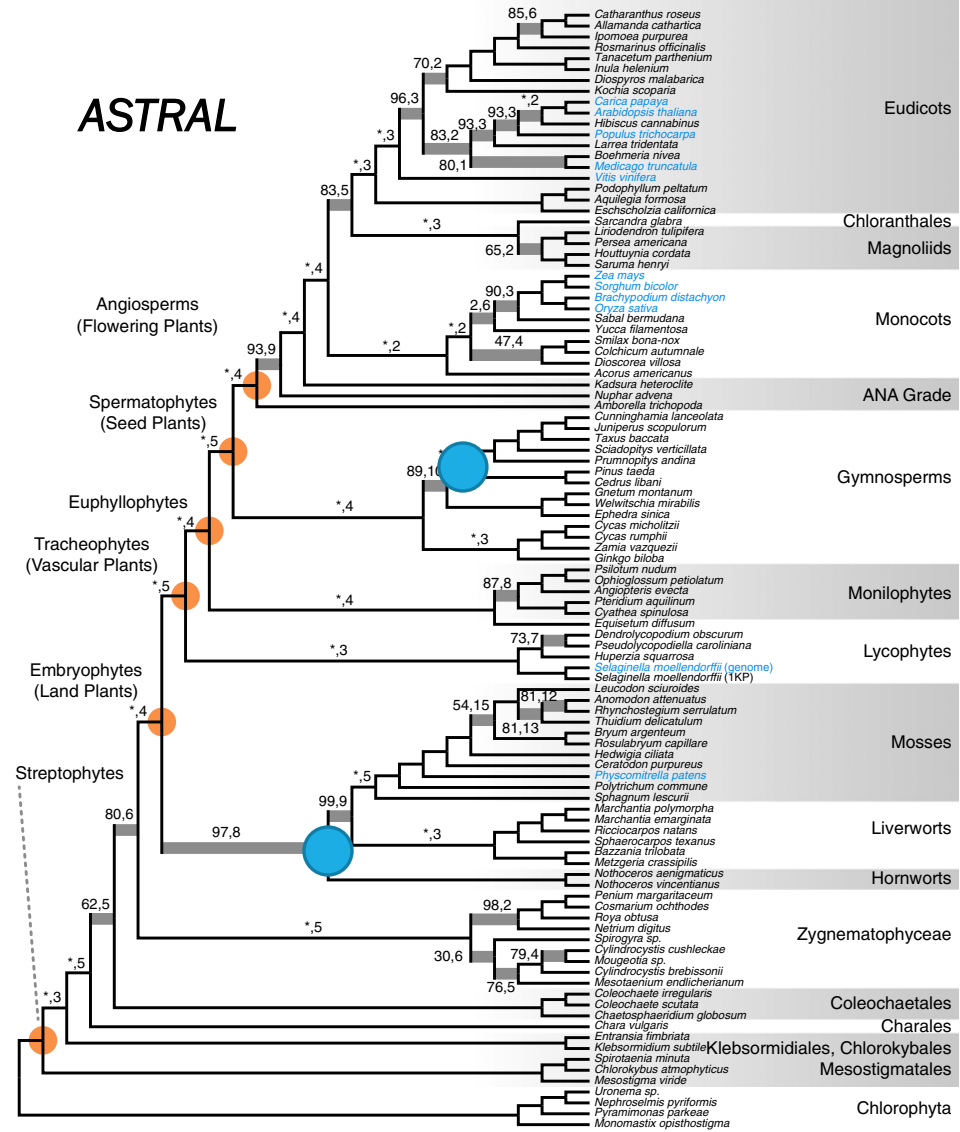


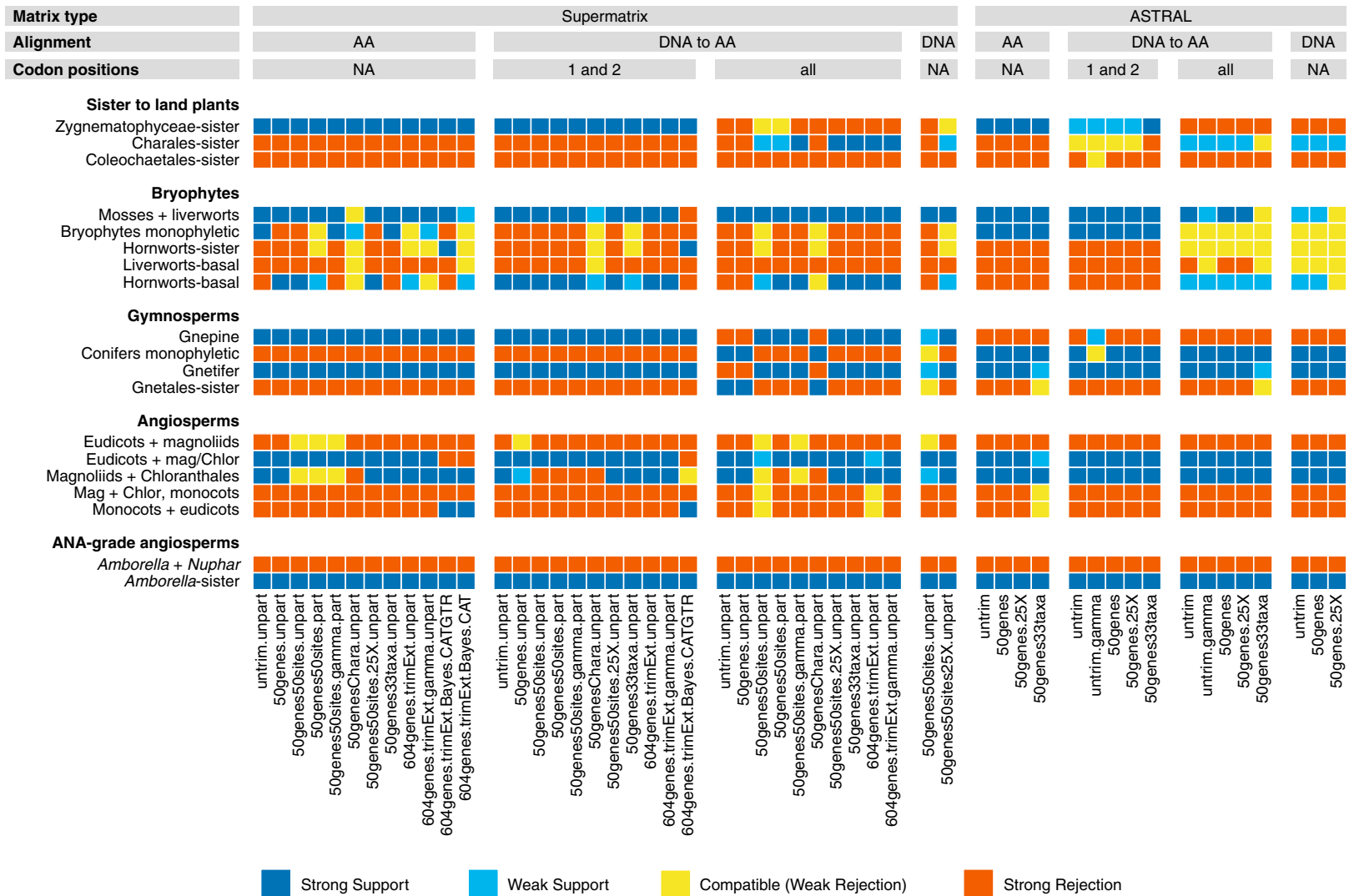
Amb-sister  
(Soltis et al 2011)

# ML analysis



# ASTRAL





# Analyses on these data

Look at the nucleotide data

852 gene regions

3<sup>rd</sup> positions removed because of molecular evolution issues with *Selaginella*

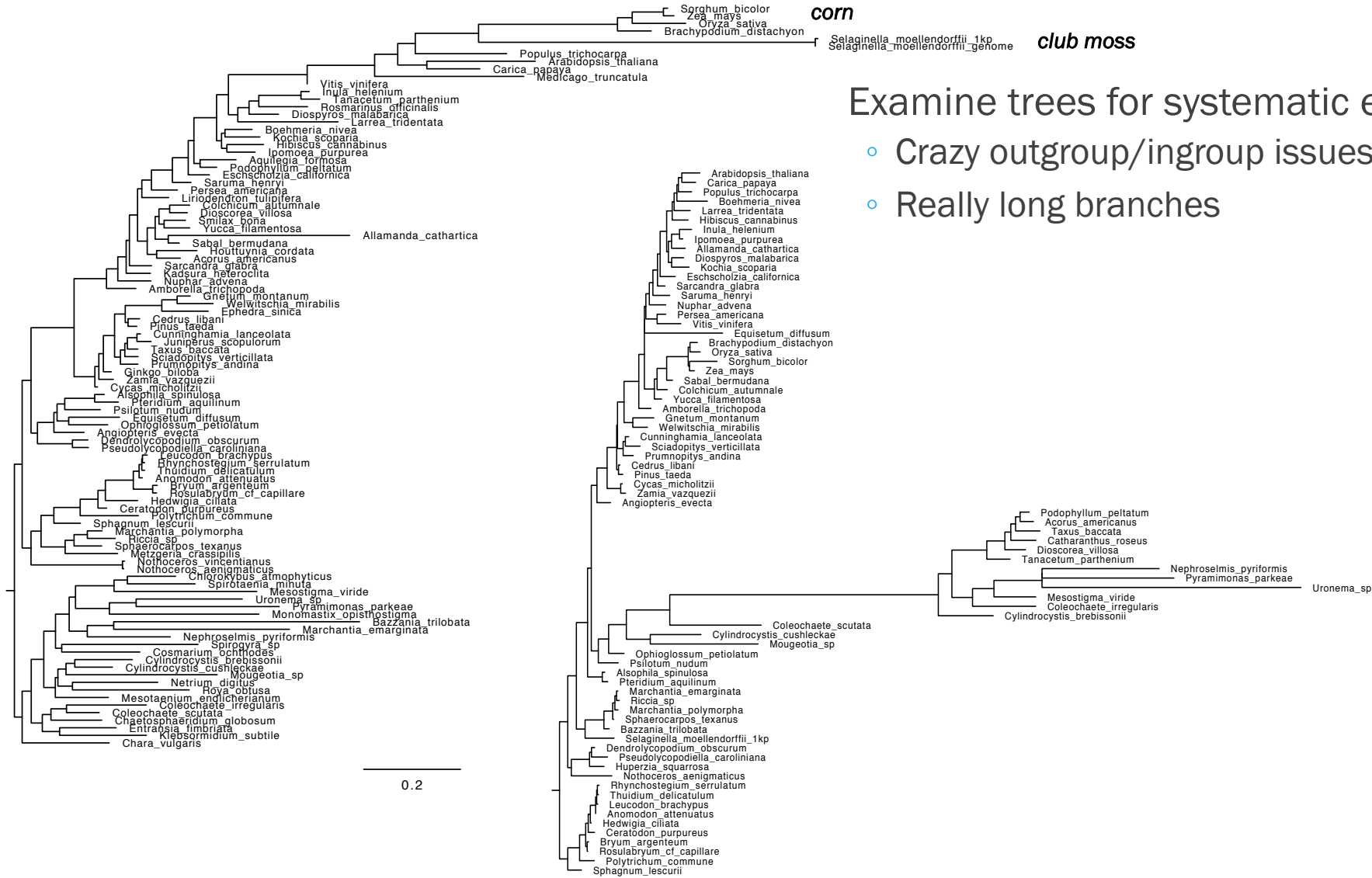
Take an edge-based view of the reconstruction of the relationships without requiring the genes to share topology beyond the edge of interest

Each gene is partially overlapping in taxonomic coverage

Answering the question: *What does the information in **this dataset** suggest for the resolution of several clades?*

Work with: Joseph Walker, Joseph Brown, Nat Hale (and me)

# Gene tree examination exposes errors



Examine trees for systematic error

- Crazy outgroup/ingroup issues
- Really long branches

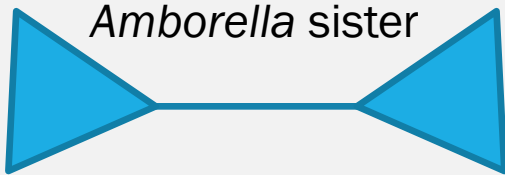
# Analyses

Instead of examining one or a few topologies

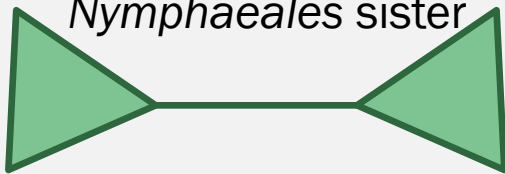
- Check concordance and conflict with ML gene trees (can filter on support)
- Constrain topologies based on the edges we are interested in and their alternatives

## Constraints

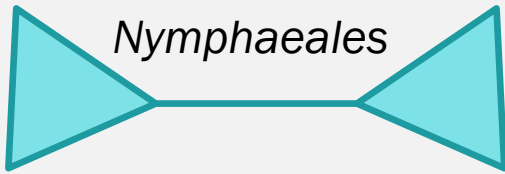
*Amborella* sister



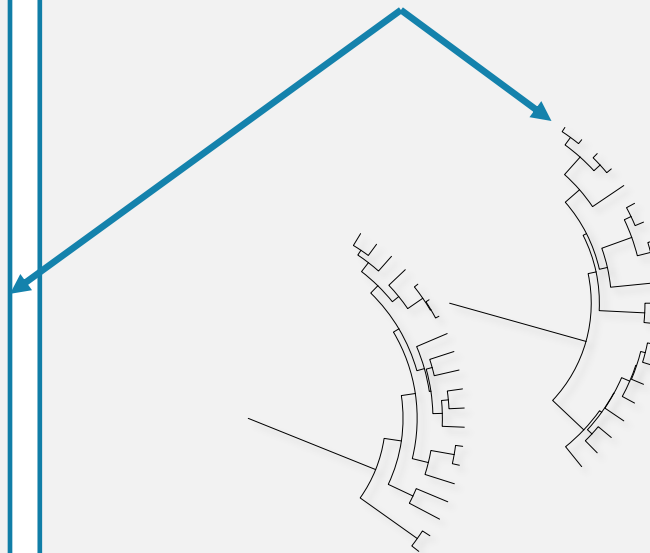
*Nymphaeales* sister



*Amborella*+  
*Nymphaeales*



## Use constraint for ML gene runs



$\#gt * \#constraints$

Which constraint had the highest likelihood (and  $> 2Lnl$ )?

Check every gene

Count genes

Sum  $\Delta Lnl$

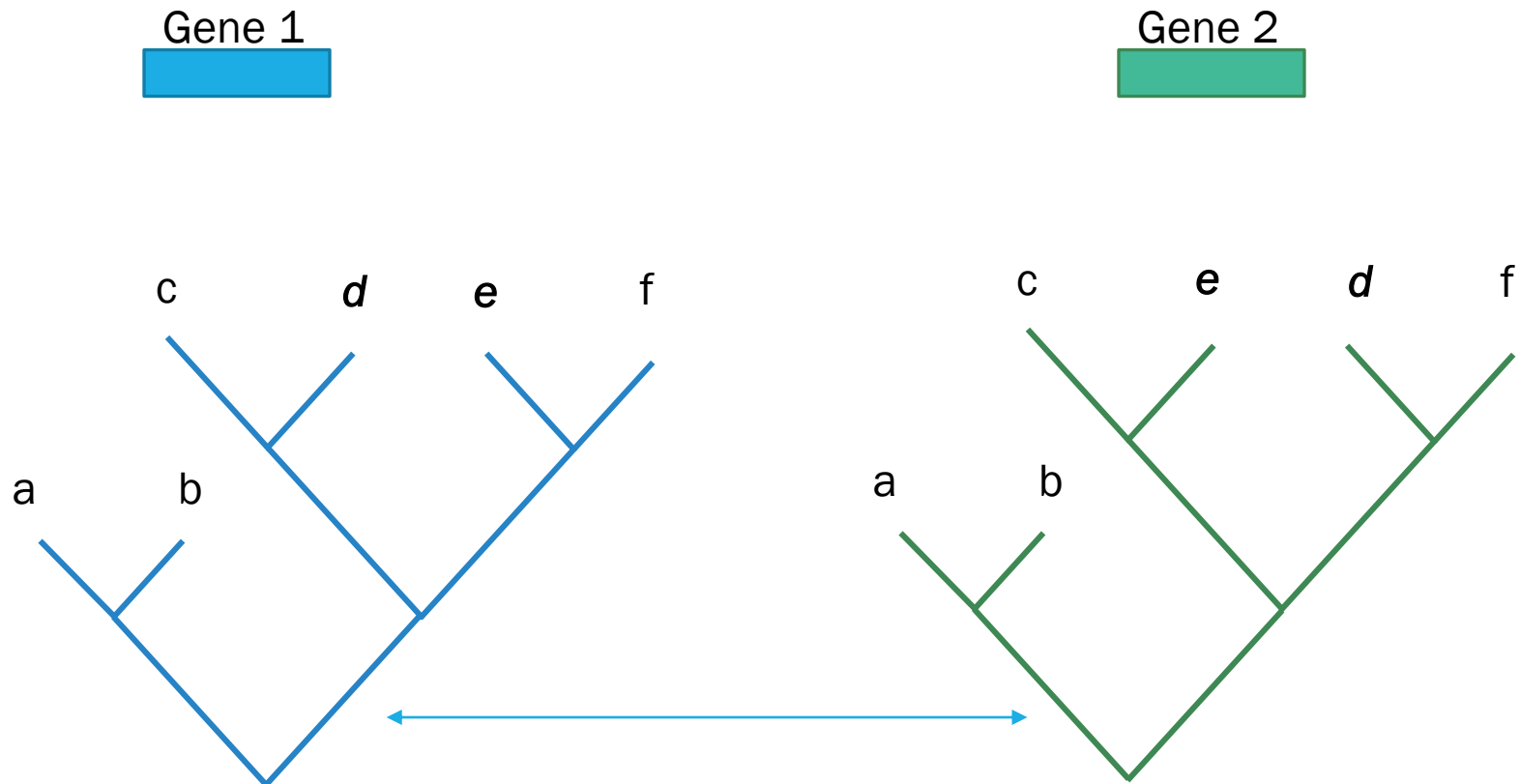
Must have necessary taxa  
Can't have crazy outgroup/ingroup

Smith et al. in prep

# Why not fixed topologies?

Several processes underlie these trees

As we incorporate more taxa or go deeper in the tree, we are likely to include even more complexity



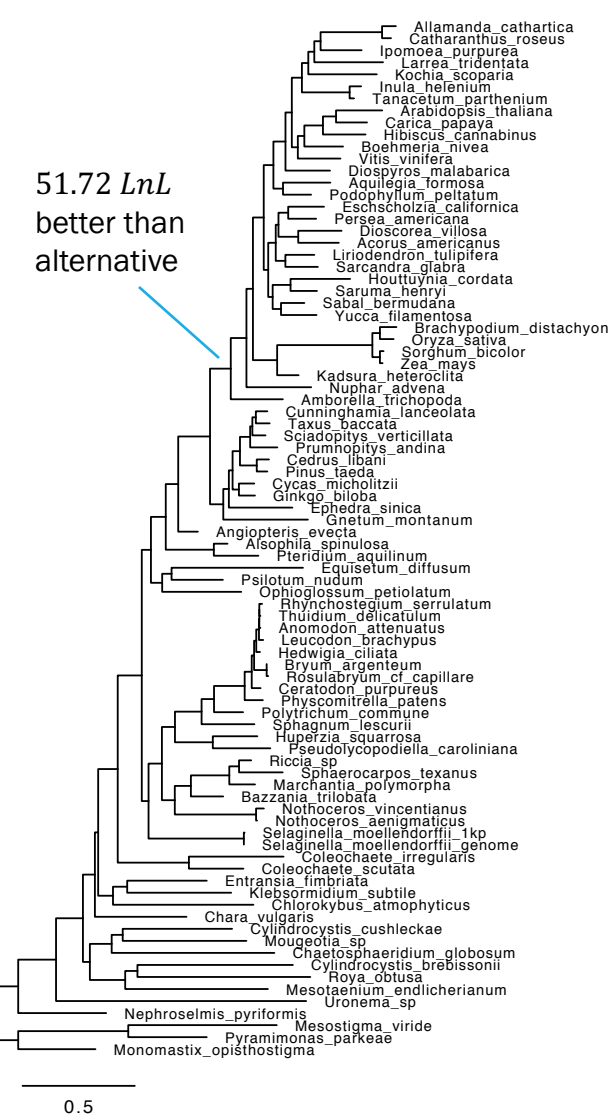
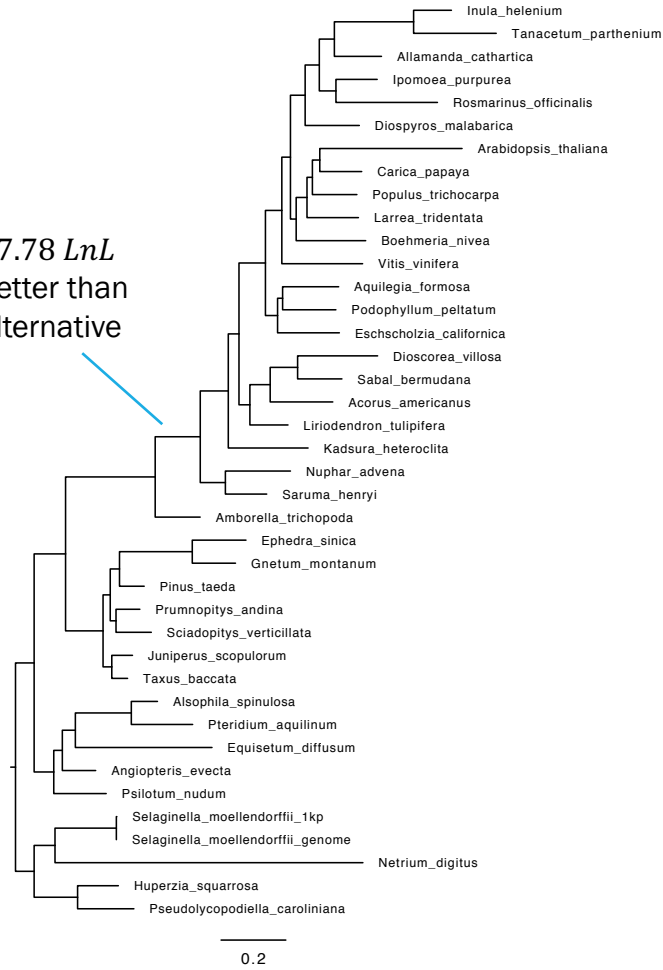
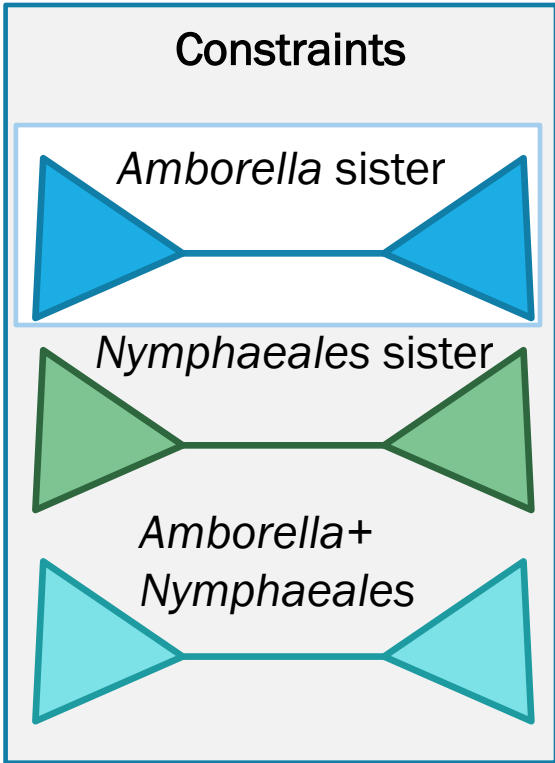


# Analyses again...

These are the better vs the alternatives and are concordant with the ML tree (so there isn't some even better alternative)

97.78 LnL better than alternative

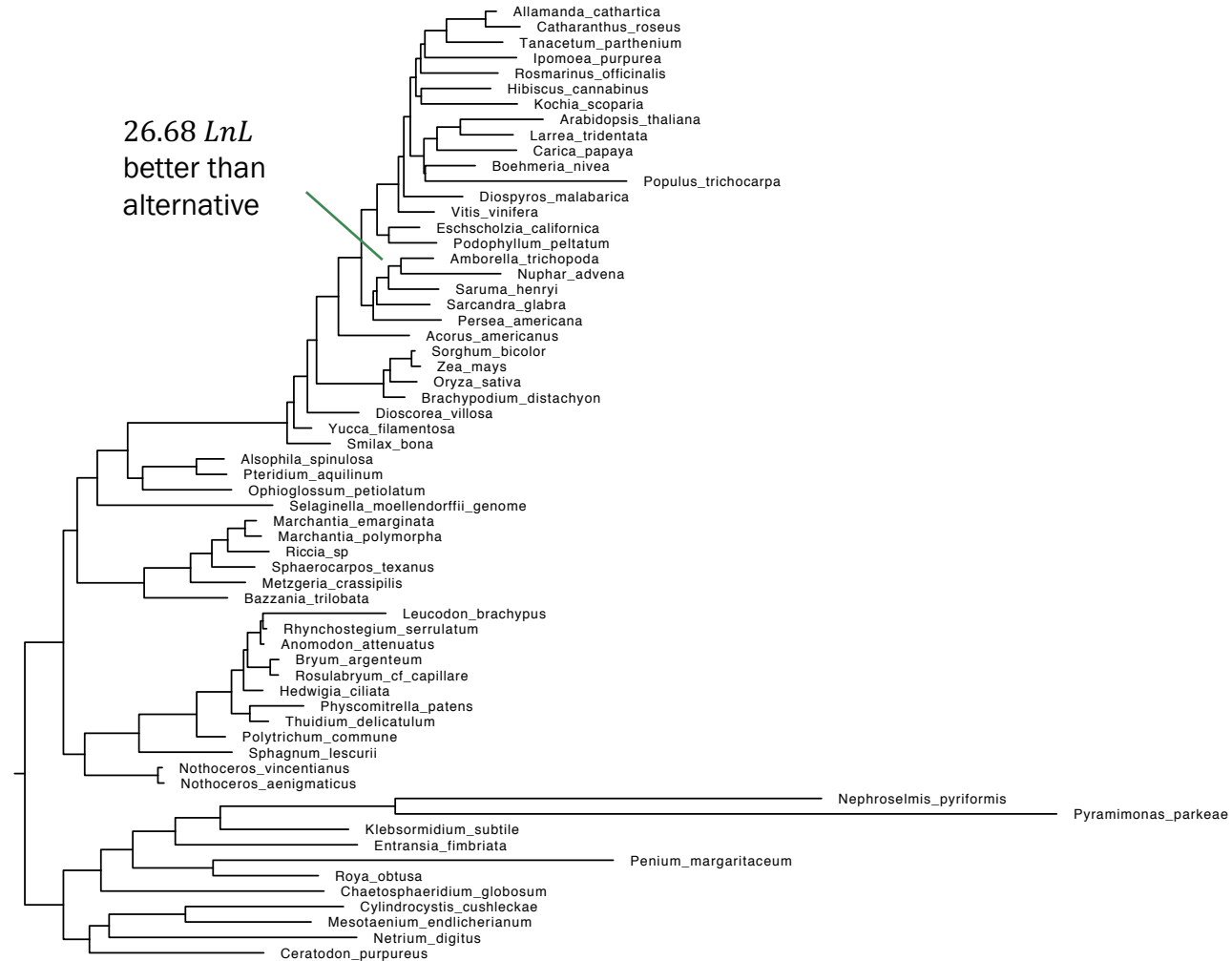
51.72 LnL better than alternative



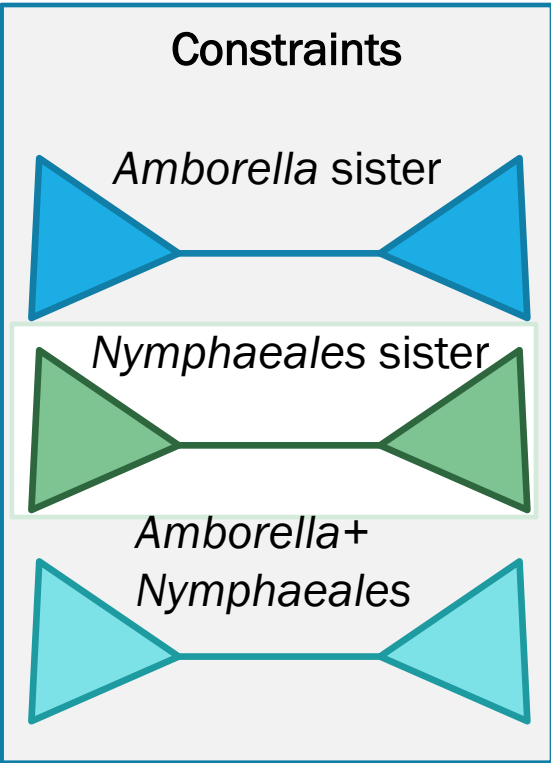
# Analyses again...

These are the better vs the alternatives and are concordant with the ML tree (so there isn't some even better alternative)

26.68 LnL better than alternative

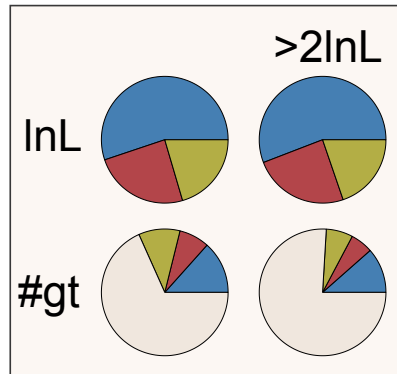
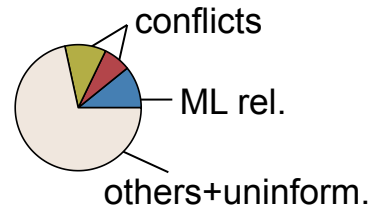
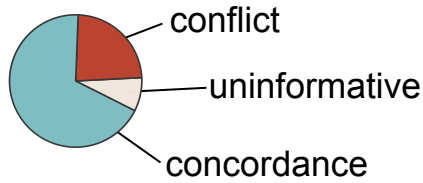


0.3



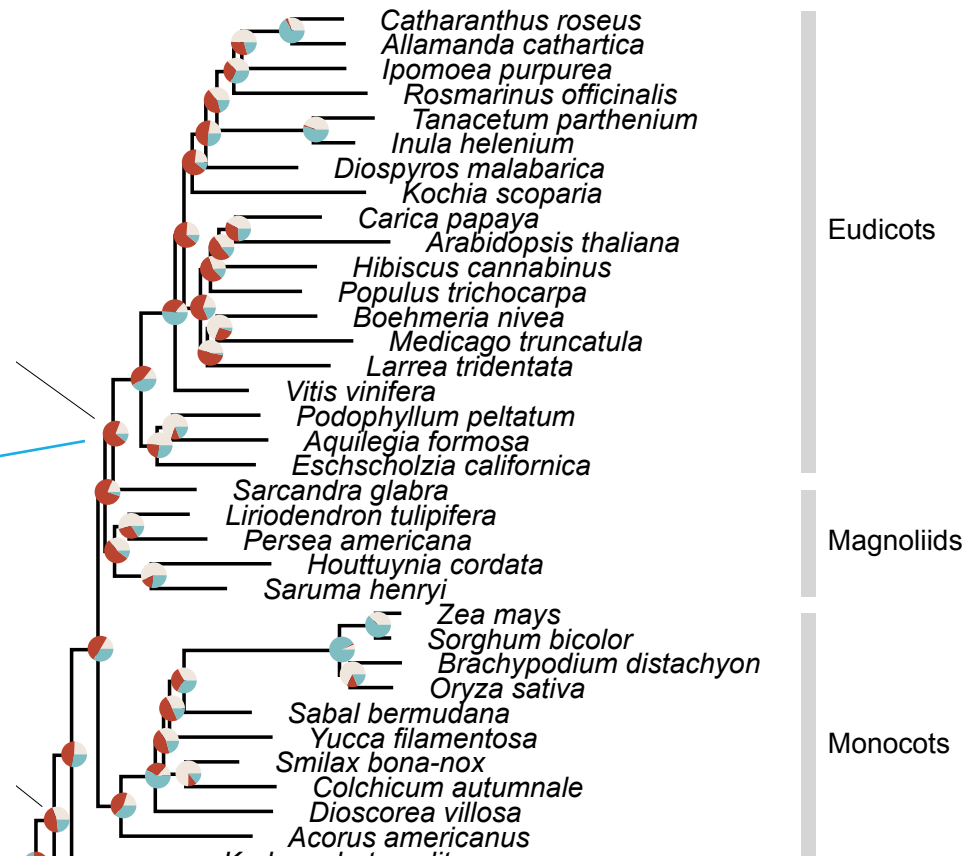
# Results

Major clade	Resolutions	Genes	Genes (> 2lnL)	DlnL	DlnL > 2
Horworts	Hornworts sister*	110	83	677.6	654.1
	Liverworts sister	56	41	294.1	280.8
	Mosses+liverworts	81	40	228.9	190.2
	All monophyly	81	37	185.3	148.5
Gymnosperms	monophyly*	288	264	7259.0	7233.8
	<i>Gnetum</i> sister	45	31	229.8	216.0
	<i>Cycas</i> sister	39	18	120.3	105.2
Gymno relat.	Gnepine*	107	85	1017.2	994.4
	conifers	93	79	800.0	787.2
	Gnetifers	134	55	288.1	217.8
	Gnetales sister	76	40	211.2	176.3
<i>Amborella</i>	<i>Amborella</i> sister*	184	152	1501.1	1470.0
	<i>Amborella</i> + <i>Nuphar</i>	118	75	564.2	526.3
	<i>Nuphar</i> sister	111	62	392.2	345.2
Eudicots	Magnoliids+eudicots*	114	98	1223.4	1204.3
	Monocots+eudicots	66	49	541.5	526.5
	Monocots+magnoliids	90	58	453.3	425.5



Many that were in conflict are uninformative (weakly in conflict, weird relationships, etc) or represent very minority relationships

Major clade	Resolutions	Genes	Genes (> 2lnL)	DlnL	DlnL > 2
Eudicots	Magnoliids+eudicots*	114	98	1223.4	1204.3
	Monocots+eudicots	66	49	541.5	526.5
	Monocots+magnoliids	90	58	453.3	425.5



In this dataset...

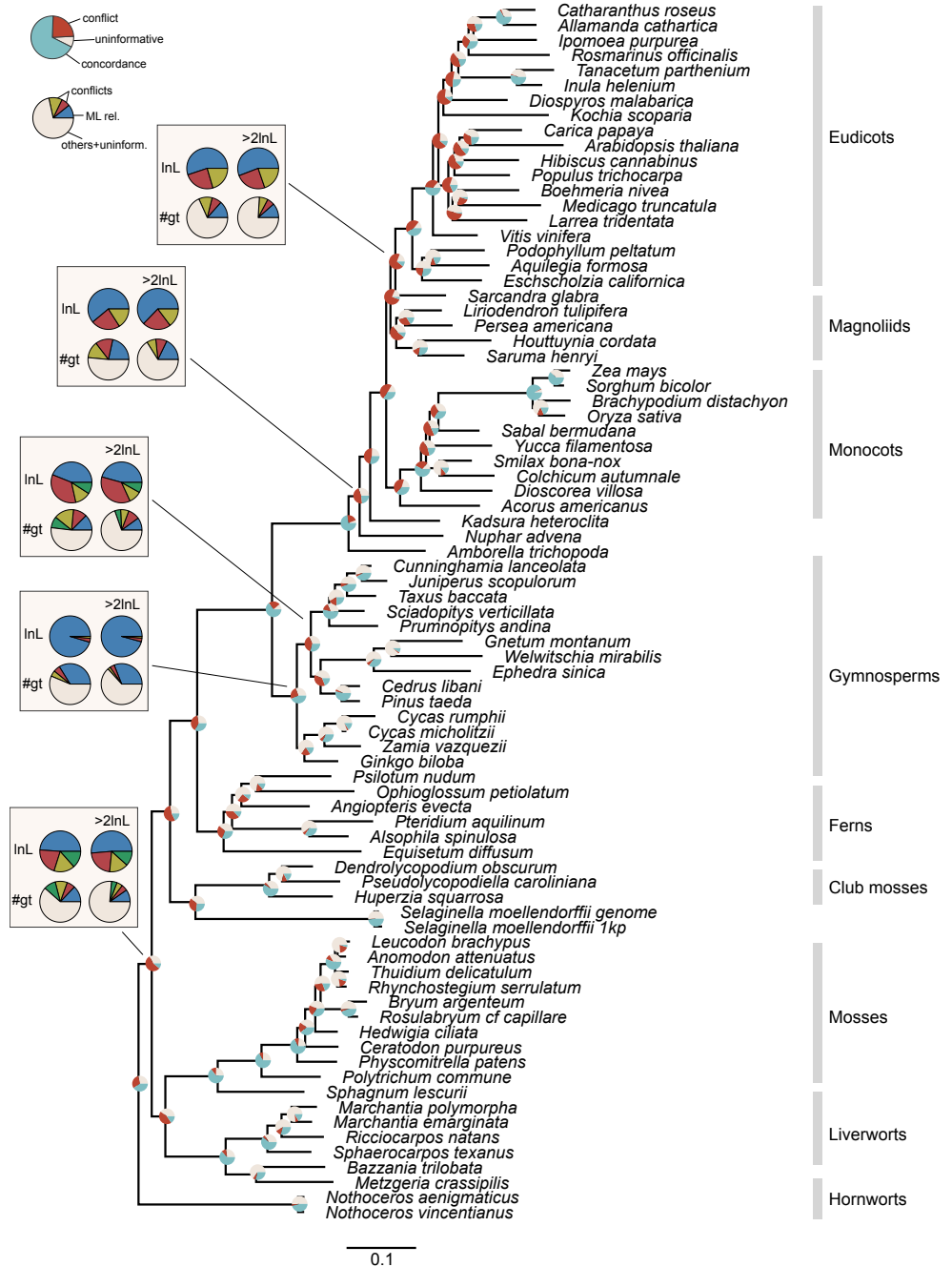
Magnoliids are sister to eudicots

Amborella is sister to Angiosperms

Relationships within Gymnosperms are not clear

Gymnosperms are monophyletic

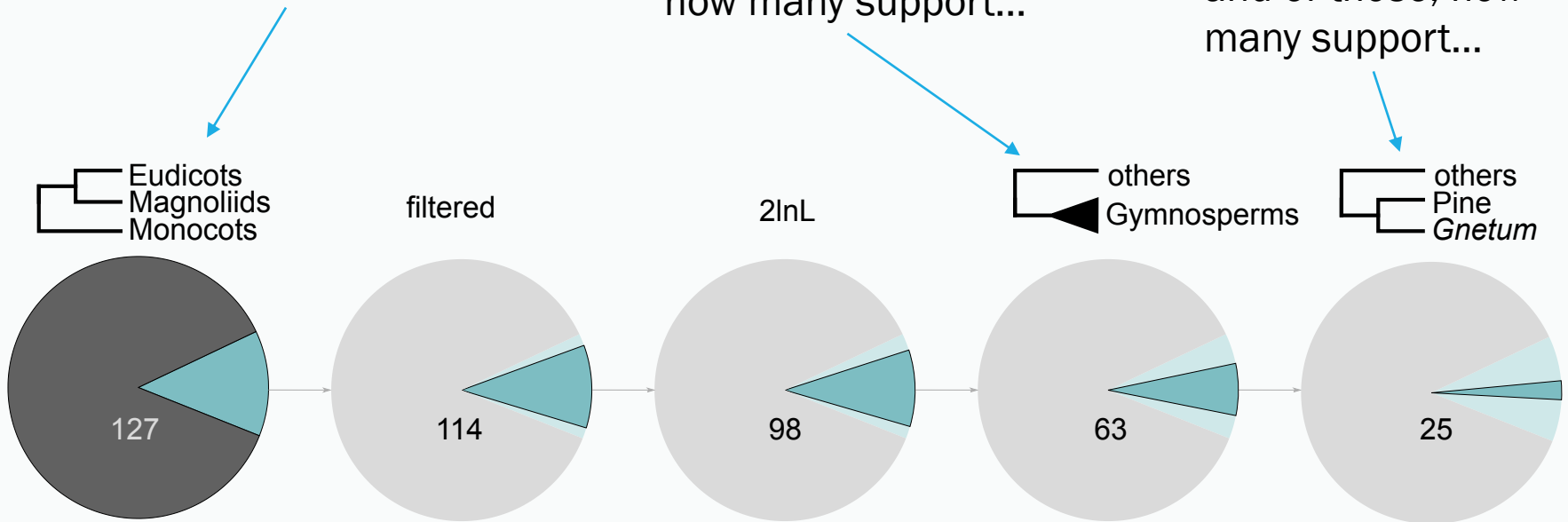
Bryophytes are *not* monophyletic

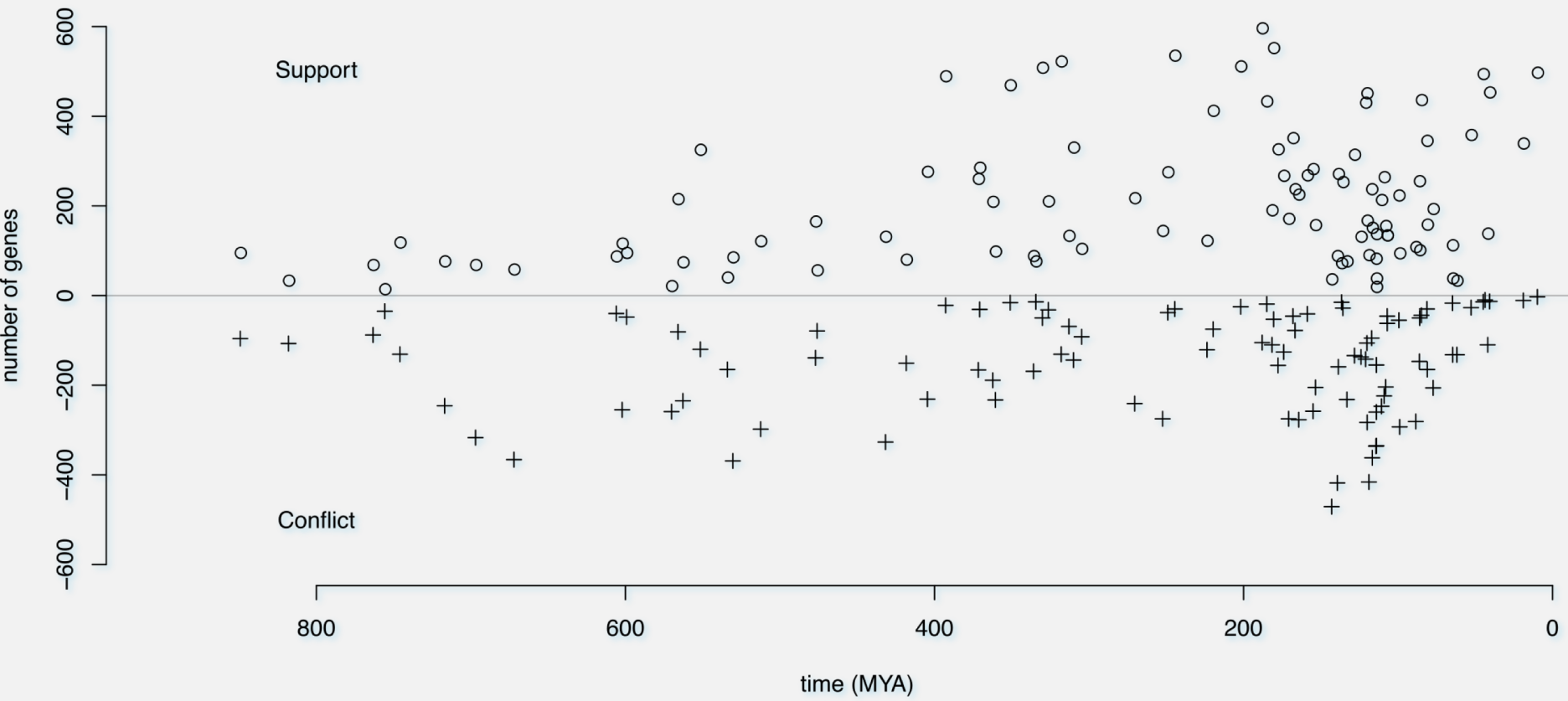


Of those trees that support...

how many support...

and of those, how many support...





# Conclusions

By examining individual relationships we may be able to more confidently make conclusions regarding *datasets* and *relationships*

It is probably unreasonable to assume that all genes will speak to all the edges of a tree

- Rate of evolution
- Gene specific evolutionary processes

These approaches explored here are just the first steps, but

- They are tractable
- They are easily extensible
- They support a great deal of complexity

Questions remain

- Are these consistent with the coalescent?

Look for more coming soon!



# Do you want to do anything that is here?

Check out the github and bitbucket under the

- User: blackrim
- Organization: FePhyFoFum

## phyparts

- General tool for conflict and concordance
- If you have duplications
- <https://bitbucket.org/blackrim/phyparts>

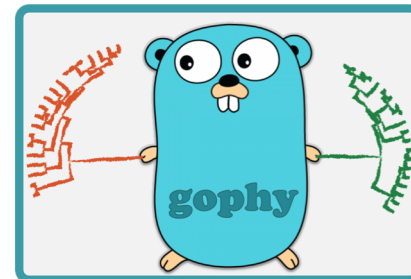
## phyx

- Many tools but pxbp can be useful if you have mostly (all) overlapping taxa
- <https://github.com/FePhyFoFum/phyx>



## gophy

- New tool for intended to be faster and more efficient
- <https://github.com/FePhyFoFum/gophy>



Look for more things coming out of the lab over the next couple months

# Acknowledgements

- Funding sources
  - National Science Foundation
    - DEB
    - ABI
  - University of Michigan
- Collaborators
  - Smith lab
    - Grad students
      - Joseph Walker
      - Drew Larson
      - Lijun Zhao
    - Postdoc
      - Greg Stull
      - Ning Wang
      - Oscar Vargas
      - Diego Serrano
    - Undergrads
      - Sonia Ahluwalia
      - Julia Olivieri
      - EJ Huang



- Former postdocs
  - Ya Yang
  - James Pease
  - Joseph Brown
  - Cody Hinchliff
- Michael Moore
- Sam Brockington
- Douglas Soltis
- Pam Soltis
- All the Open Tree of Life folks



<https://github.com/blackrim/> [blackrim.org](https://blackrim.org)