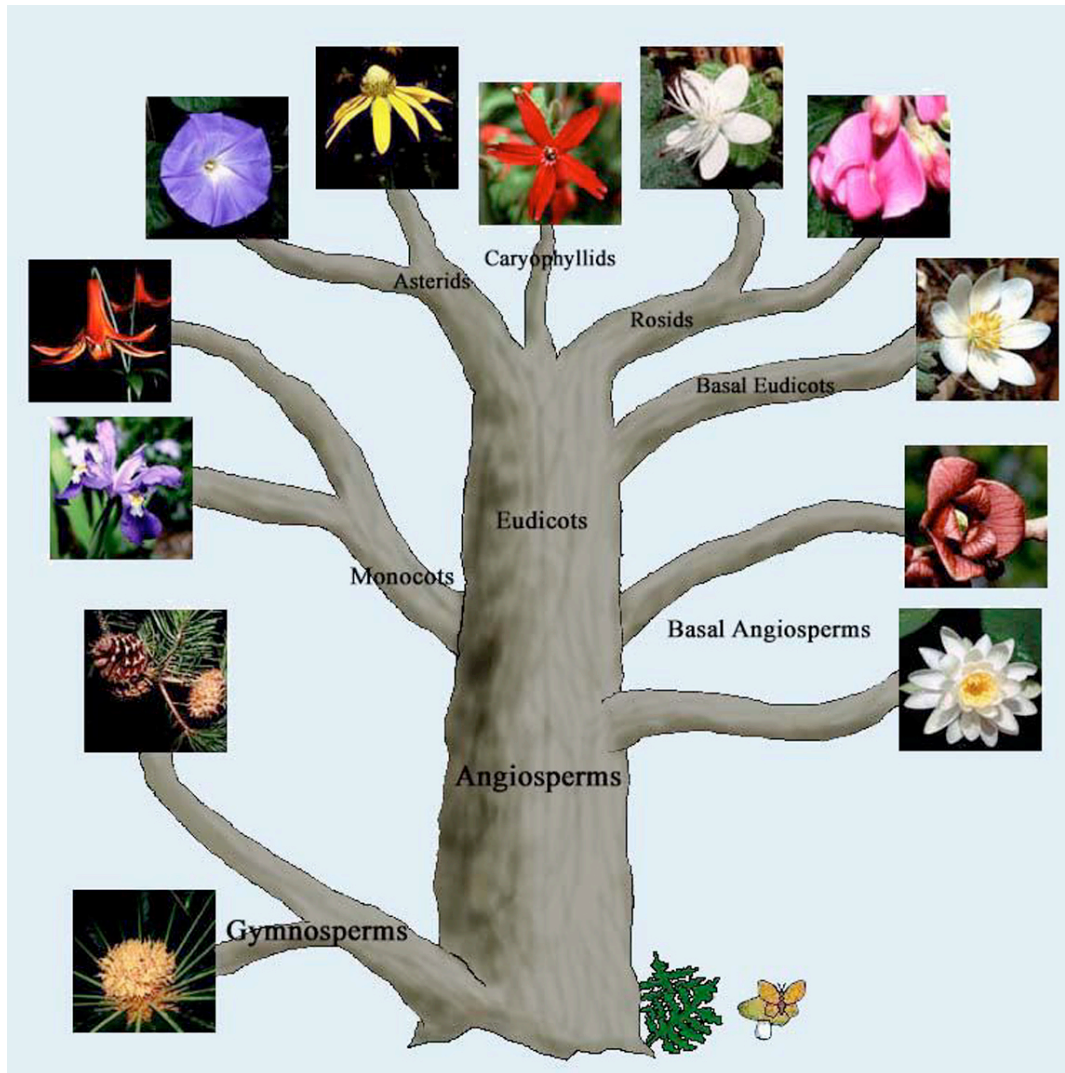


Plant Phylogenomics: Lessons from the 1KP Project



Jim Leebens-Mack

Department of Plant
Biology
University of Georgia

New Methods for
Phylogenomics and
Metagenomics Symposium
Feb, 2013

Ancestral Angiosperm/
Amborella Genome

Vic Albert
Raj Ayyampalayam
Brad Barbazuk
John Bowers
Jim Burnette
Srikar Chamala
Andre Chanderbali
Josh Der
Claude dePamphils
Jamie Estill
Hong Ma
Doug & Pam Soltis
Stephan Schuster
Sue Wessler
Rod Wing
Kerr Wall
Norm Wickett
Eric Wafula

Collaborators

MonAToL

Claude dePamphilis
Tom Givnish
Cecile Ané
Raj Ayyampalayam
Sean Graham
Dennis Stevenson
Jerry Davis
Alejandra Gandolfo
Chris Pires
Norm Wickett
Wendy Zomlefer
Michael McKain
Jill Duarte

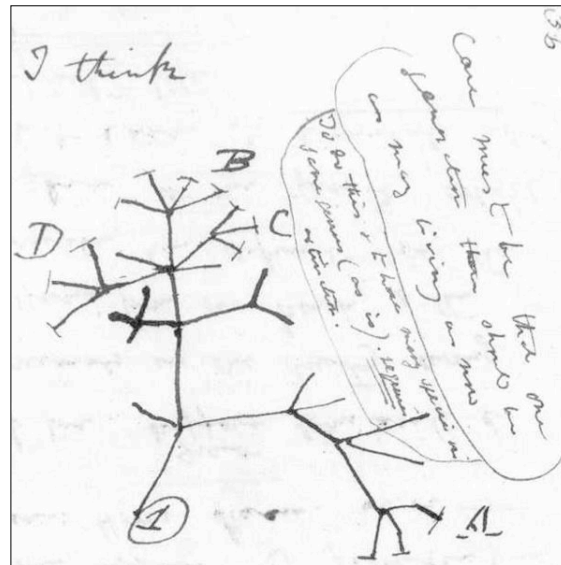
Funding: NSF, iPlant, University
of Georgia, OneKP

OneKP/MSA AToL

Norman Wickett
Nam Nguyen
Siavash Mirarab
Naim Mataci
Gane Ka-Shu Wong
BGI
Eric Carpenter
Brad Ruhfel
Herve Philippe.
Gordon Burleigh
Matt Barker
Claude dePamphilis
Tandy Warnow
Jamie Estill
Raj Ayyampalayam
Doug & Pam Soltis
Sean Graham
Dennis Stevenson
Michael Melkonian
.....OneKP Consortium

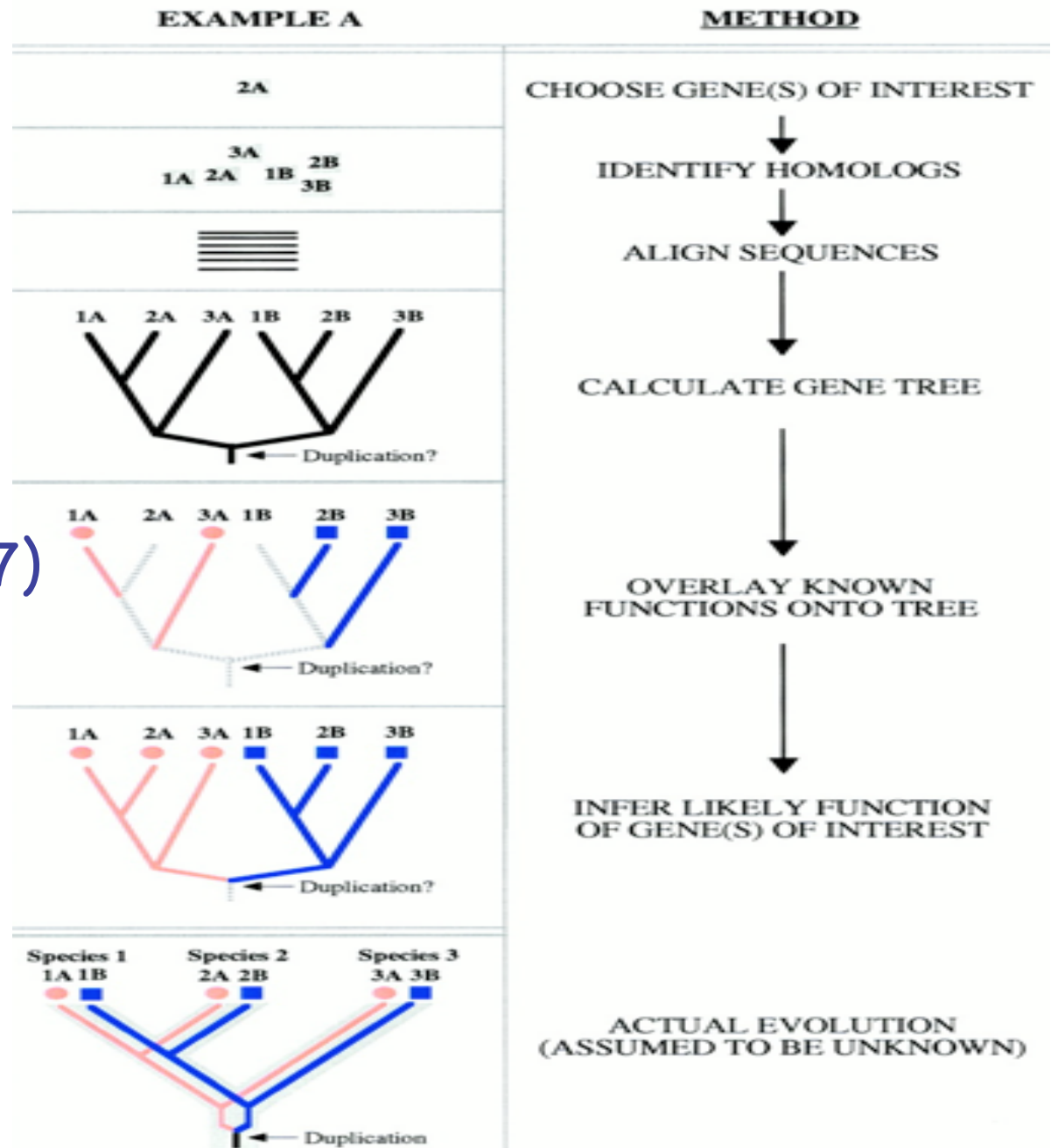
"Nothing in biology makes sense except in the light of evolution" (Theodosius Dobzhansky, 1973)

"Nothing in evolution makes sense except in the light of phylogeny"



Darwin (1837) First Notebook on Transmutation of Species

“Phylogenomics” -
Jonathan Eisen
(1998; Genome
Research 8:163-167)



Current Usages

1. Using genome-scale data to resolve phylogenetic relationships
2. Genome-Scale comparisons placed within a phylogenetic context

What About Nuclear Gene Histories?

Grouped by Phylogeny

[Angiosperms](#)

[Non Flowering](#)

[Green Algae](#)

Grouped by Application

[Agricultural](#)

[Biochemical](#)

[Medicinal](#)

[Extremophytes](#)



www.onekp.com

Gane Ka-Shu Wong (Alberta)



musea ventures



AAGP Ancestral Angiosperm Genome Project

ancangio.uga.edu



MonATol

monatol.uga.edu

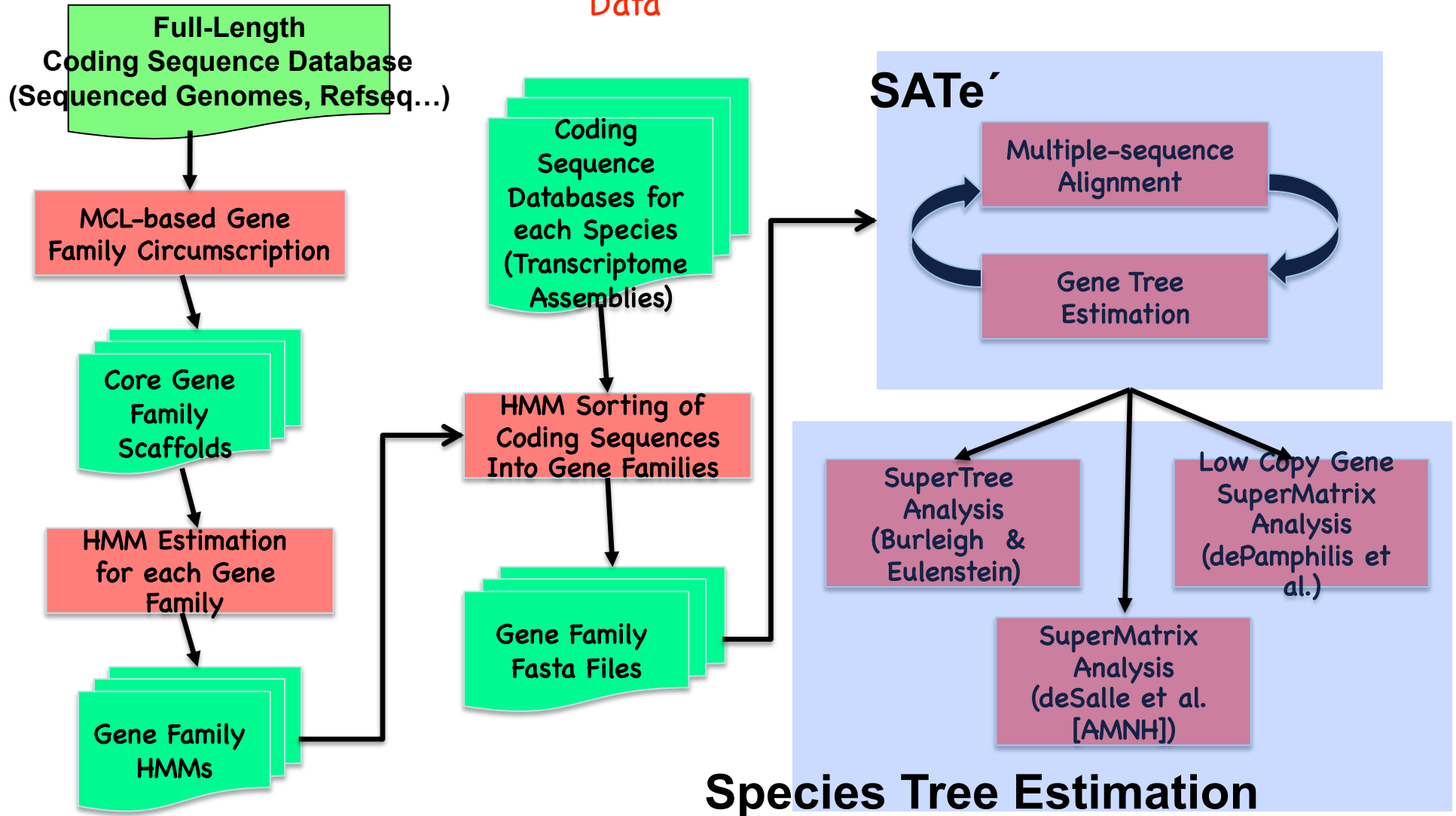


iPlant Tree of Life

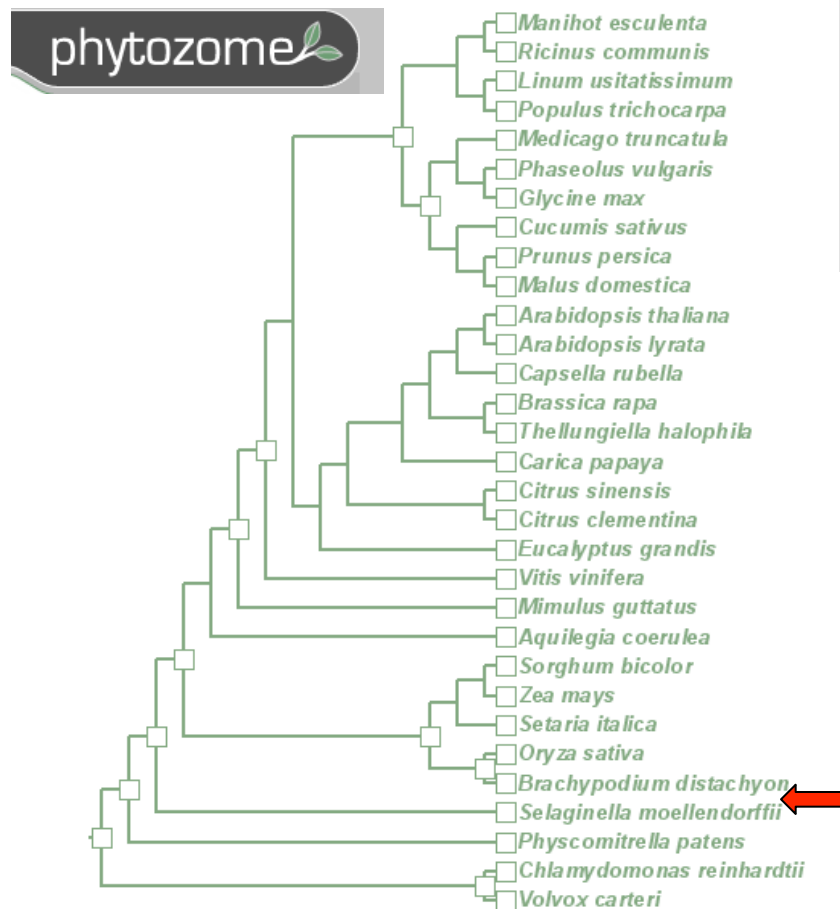


Gene Family Circumscription, Sequence Alignment, Gene Tree Estimation and Species Tree Estimation

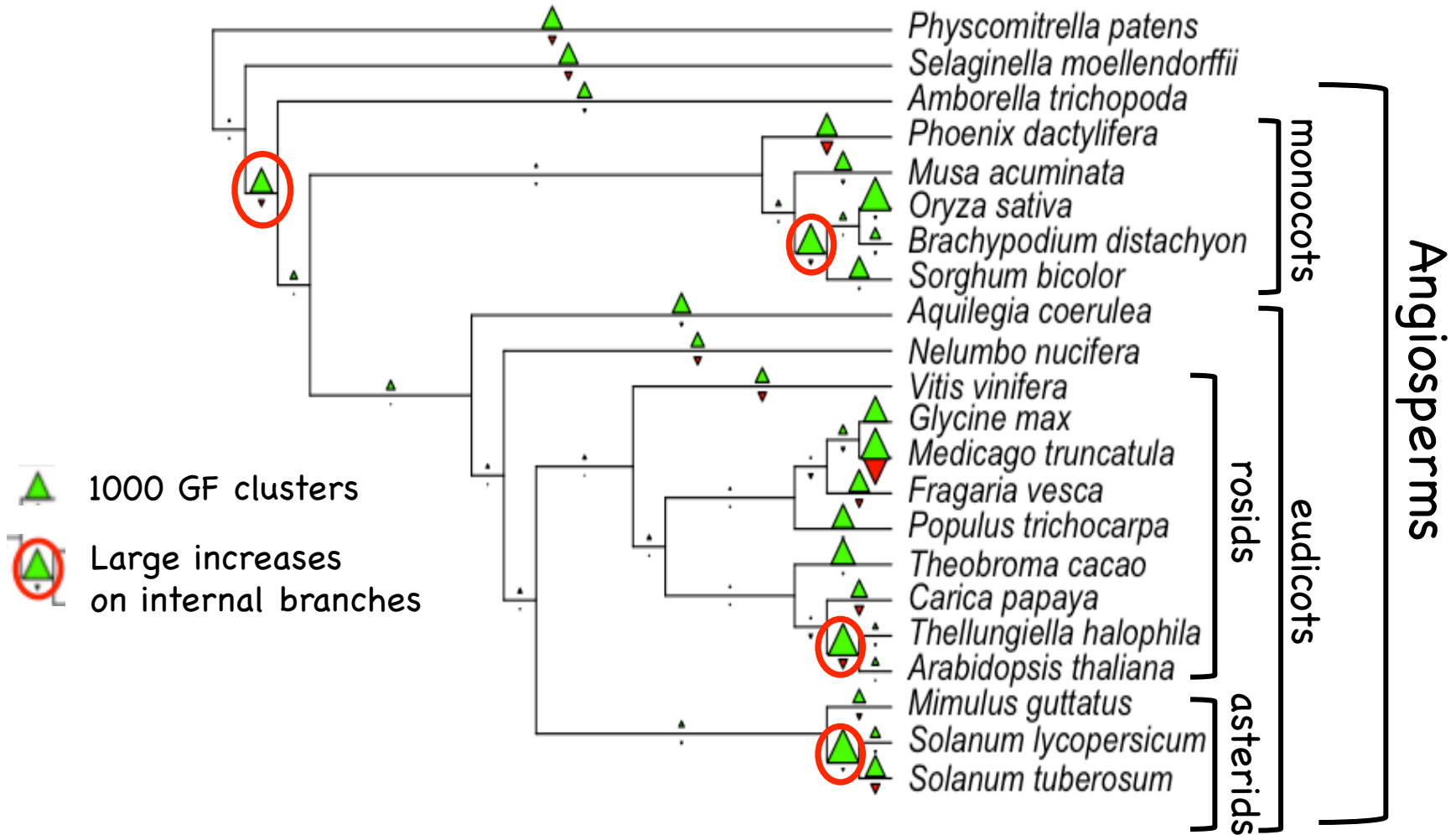
Reference Proteomes + New Transcriptome Data

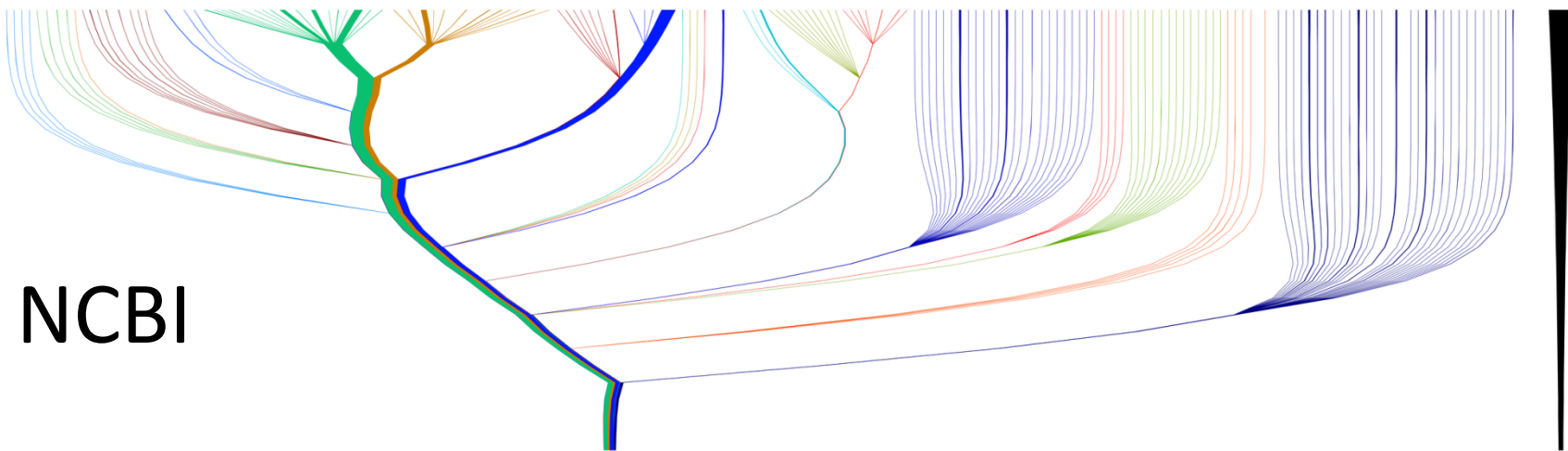
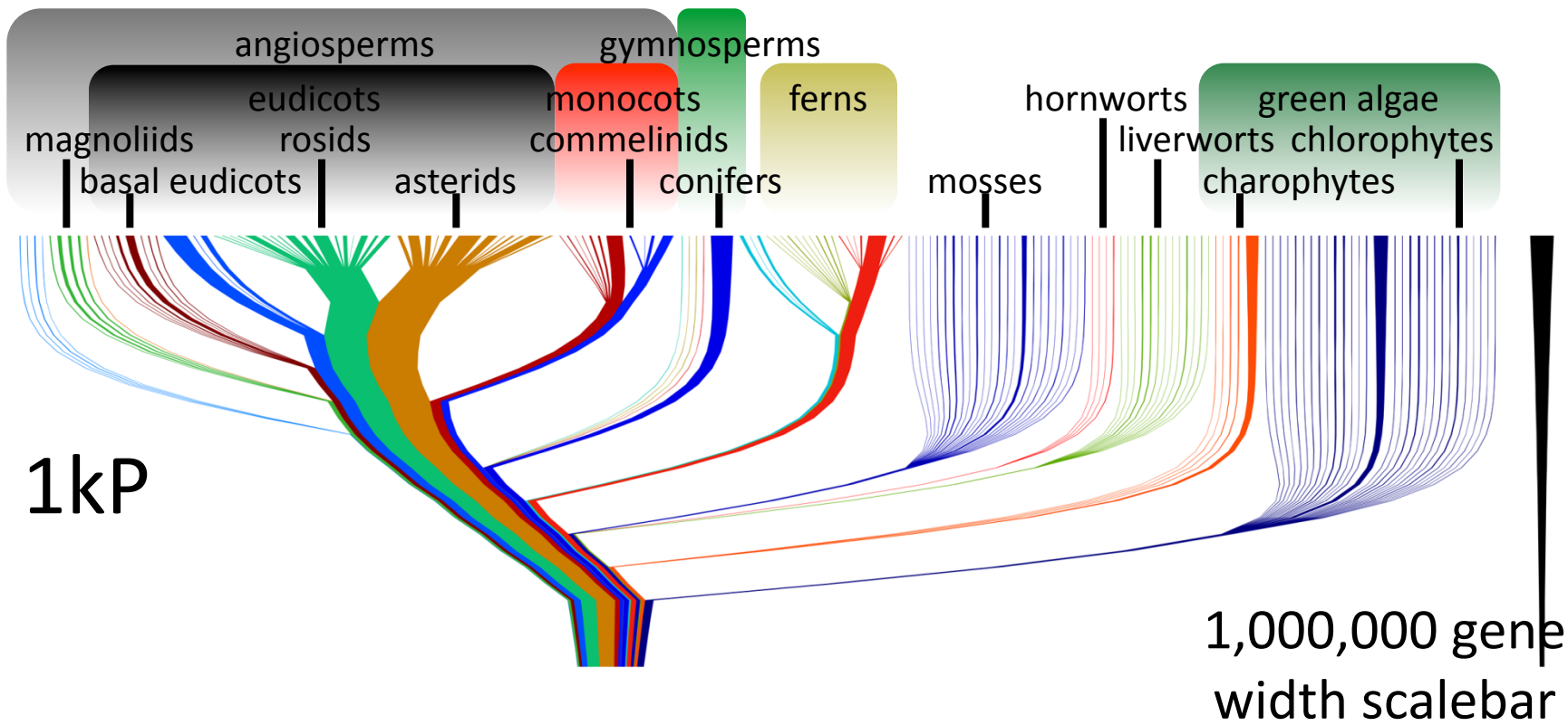


Gene Family Classification for Full-Length Genes from Sequenced Genomes – orthoMCL clustering



Gain and Loss of "Gene Families"





1kp – Questions to be Addressed through Estimation of Gene Trees and Species Relationships

- What are the relationships among lineages across the green plant tree of life?
- What was the nature and gene composition of the likely ancestor of the Viridiplantae?
- What was the impact of lateral gene transfer (from bacteria) on the early evolution of the Viridiplantae and the Chlorophyta/Streptophyta?
- Were gene and/or genome duplication events associated with innovations in green plant evolution including the origin of the flower, the seed, the vascular cambium (wood), alternation of generations, the shift from a life history dominated by the haploid phase to a dominant diploid life stage, colonization of land and shifts from single cells to multicellularity,
- Has polyploidy played a role in the diversification of angiosperms?

Check out OneKP data



Blast for OneKP Project

[Document](#)

Terms of use: These sequences are being released in advance of publication as a service to the community. We only ask that you follow the spirit of the [Fort Lauderdale agreement](#) and refrain from doing the kinds of analyses for which these data were generated, as described on the [1KP project website](#). More generally we ask that you refrain from using these data in any studies that involve multiple genes across multiple species (e.g. studies of a biological process in a particular taxon). Analyses of one or two genes across multiple species, or multiple genes in one species, are generally not in conflict. If you wish to publish your findings, we ask that you credit 1KP with minor authorship. The individuals who must be credited will depend on the particular data that you use. Please contact gane@ualberta.ca or ejc@ualberta.ca for additional clarification.



Enter Query sequence

Enter FASTA sequence

clear

Or, upload file: no file selected



Basic Param

Blast Tools

blastn

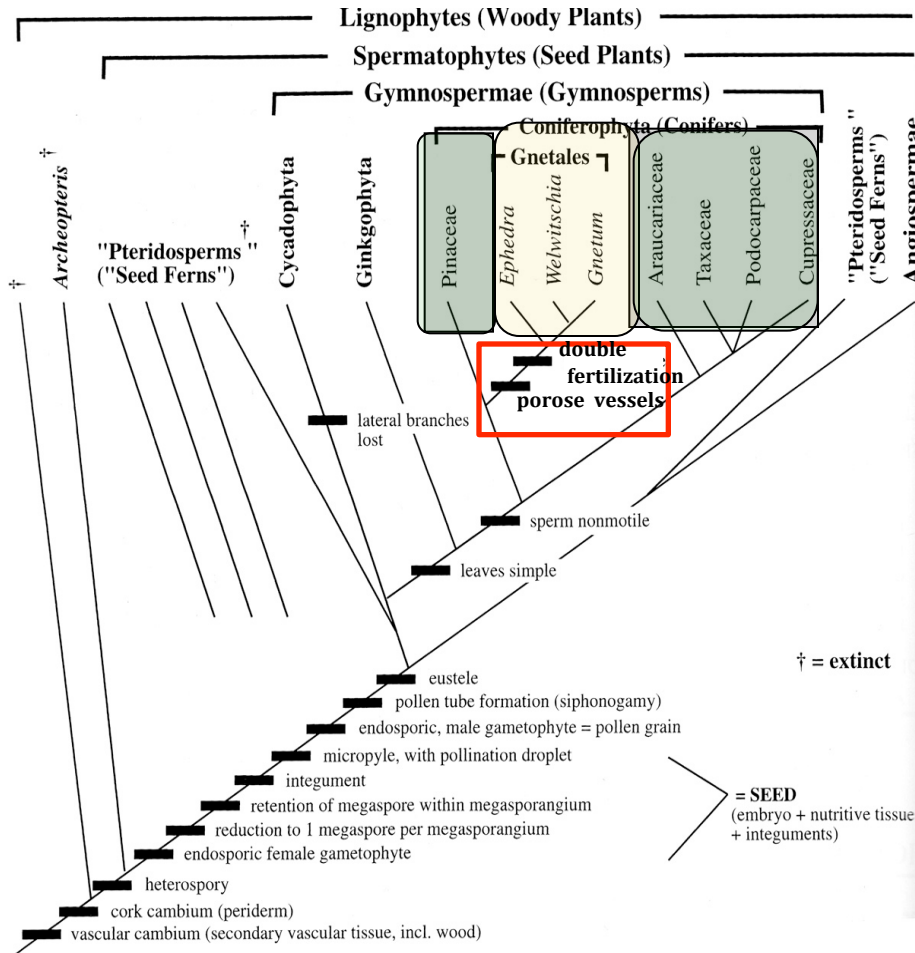
DataBase

Class: 1059-All_Sample

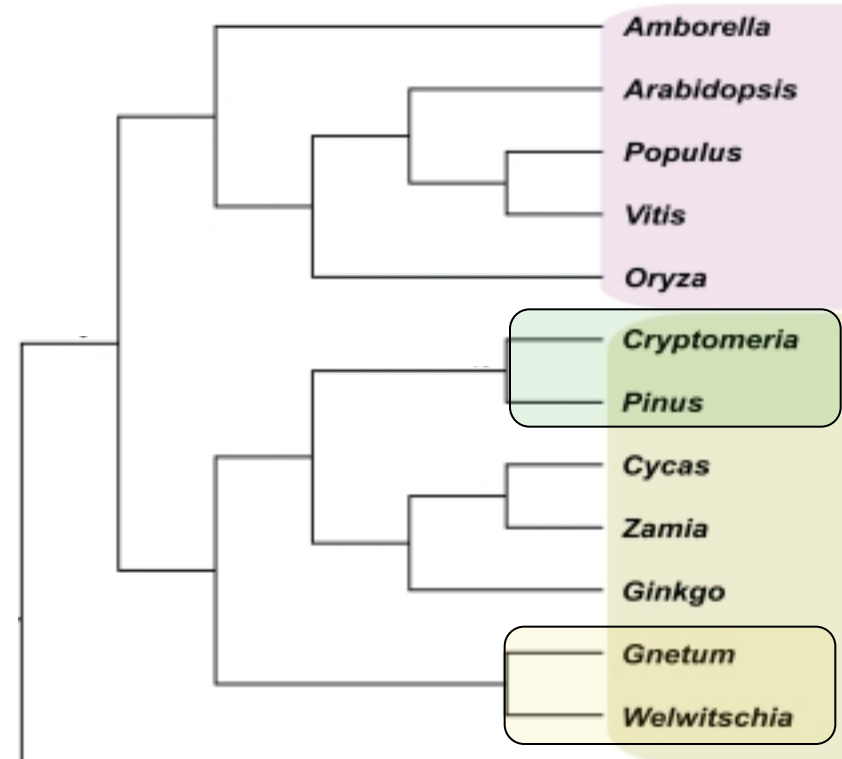
Sample: All_Sample

Surprises from Phylogenomics: Monophyly of Conifers?

Plastid Genome Analyses

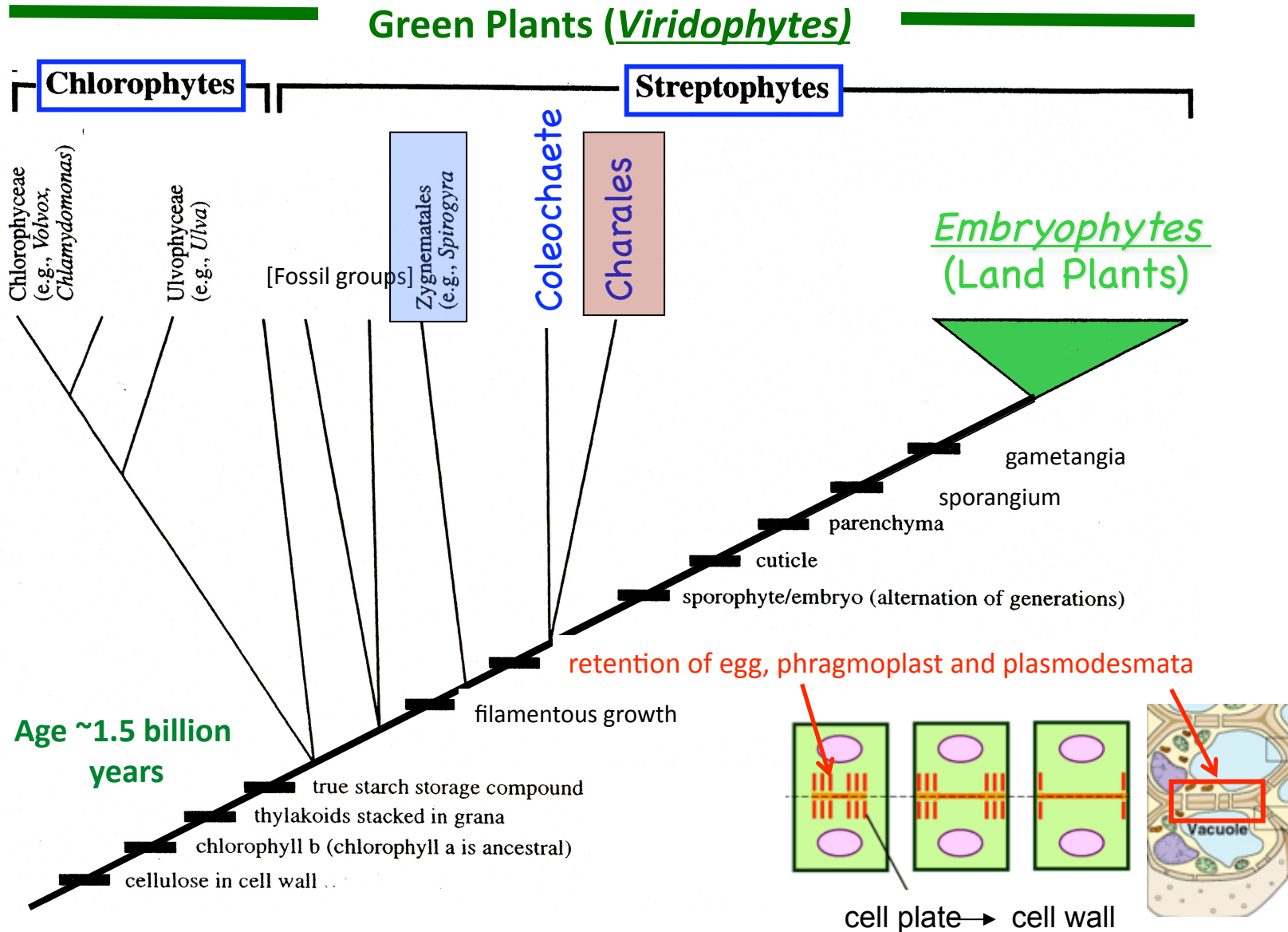


Large Nuclear Gene Analyses

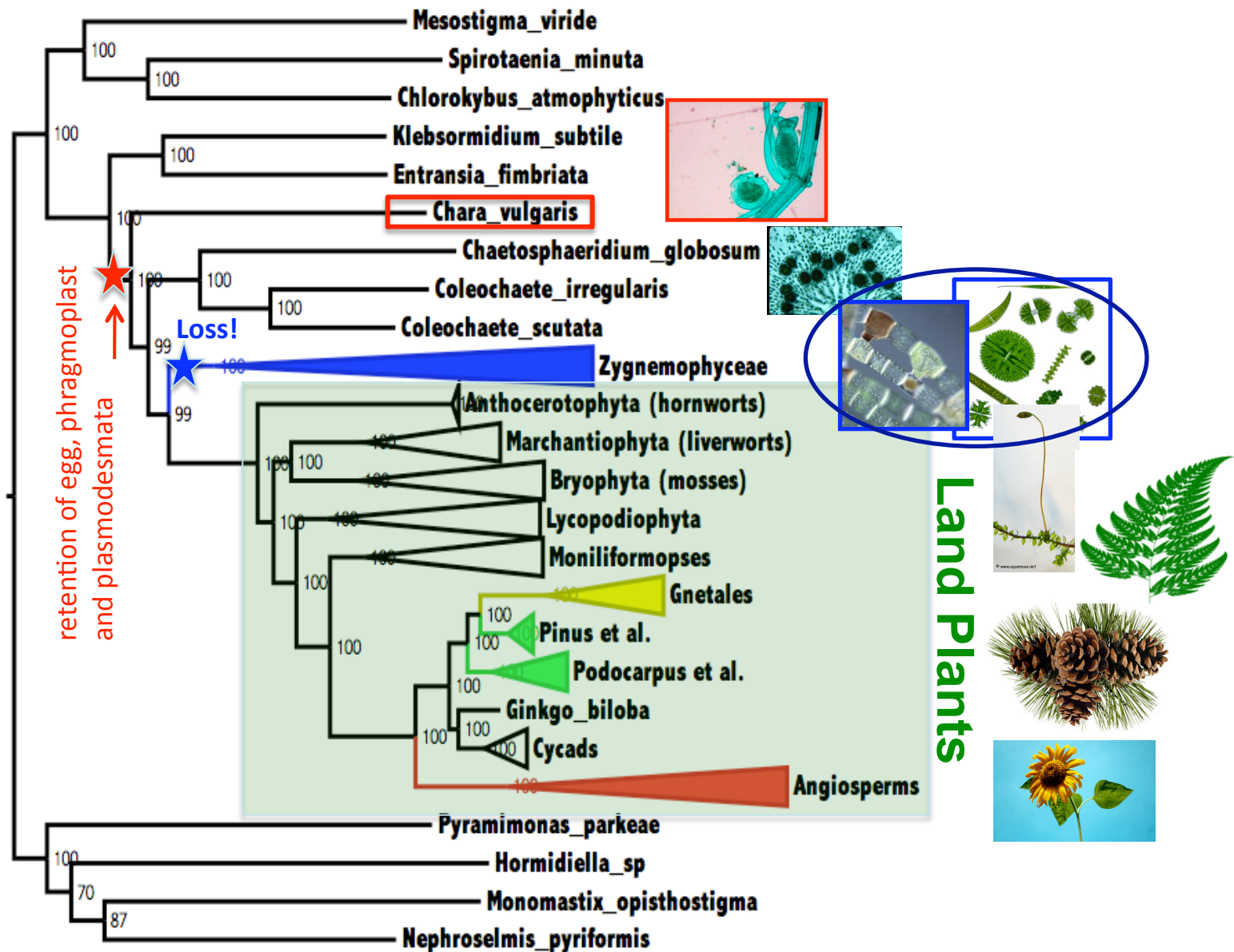


Cibrián et al. 2010, Lee et al 2011, OneKp unpublished

Origin of Land Plants: Retention of egg, phragmoplast and plasmodesmata preadaptations for colonization of land

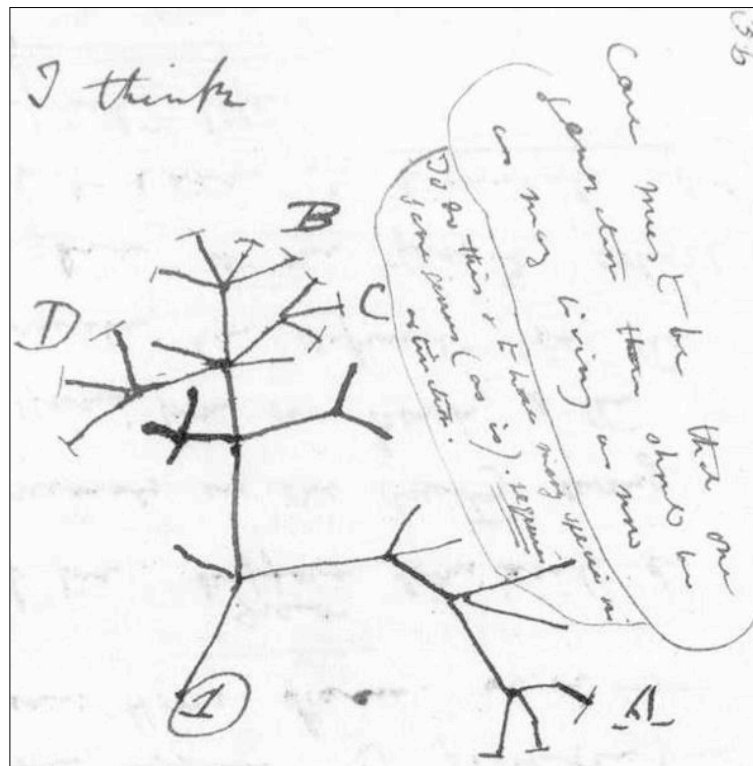


Analysis of 604 single-copy nuclear genes agrees with plastome analysis suggesting that retention of egg, plasmodesmata, and phragmoplast lost in diverse Zygnemophyceae



Why should you believe these results?

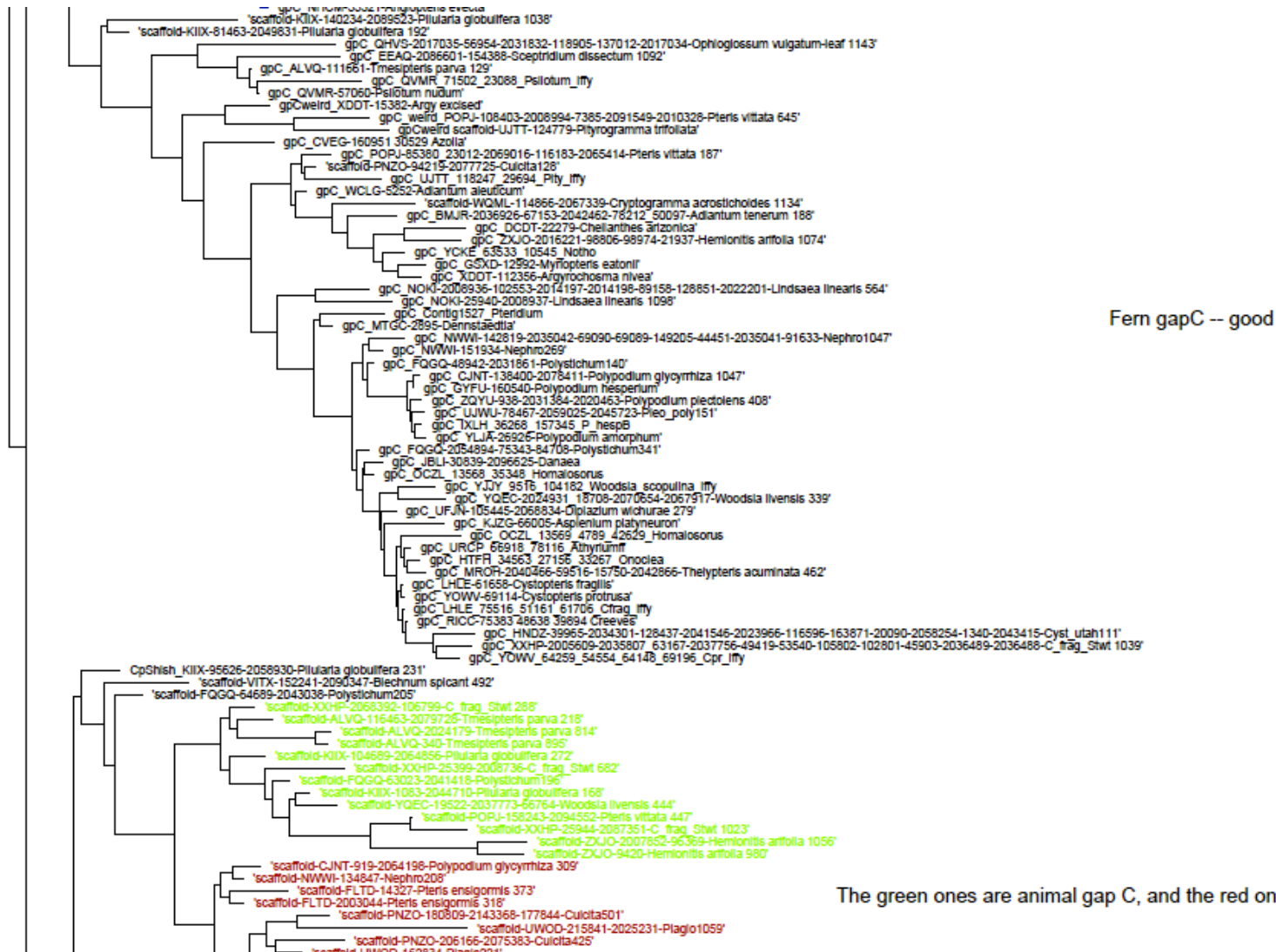
- You shouldn't... Inferred trees are hypotheses!



Consider Possible Complications

- Lots of missing data!
 - Remove genes/taxa with lots of missing data
- Ortholog identification
 - Focus on genes/genomes that tend not to be duplicated
- Contamination
 - BLAST, Sequence placement on reference tree (SEPP?), Long branch trimming
- Model misspecification
 - Model validation – simulations
 - Test robustness of hypothesis among inferences derived from alternative models

Contamination is an issue!



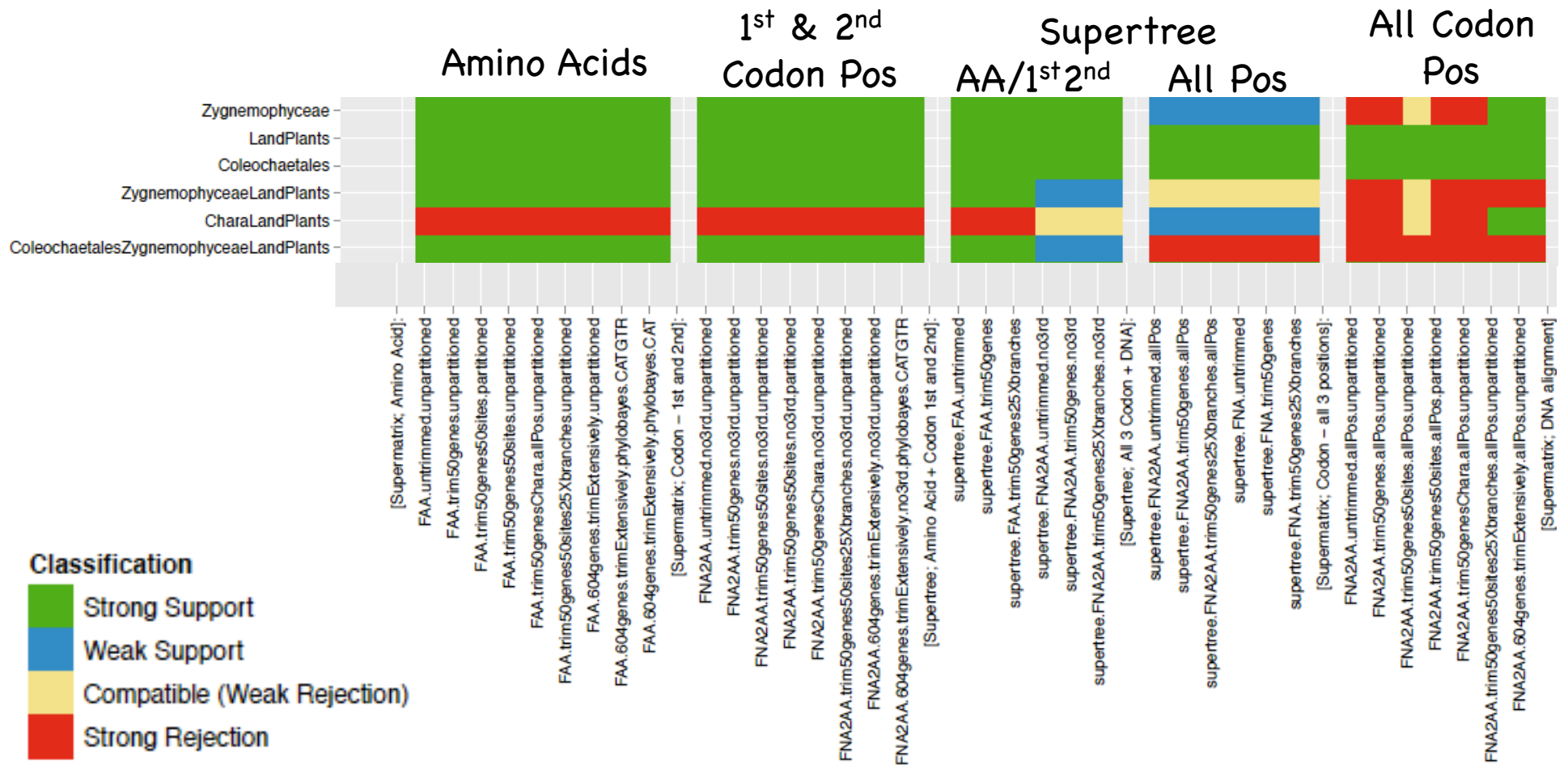
Data Matrix/Analysis Perturbation

- DNA vs Amino Acids
- All codon positions vs remove third position
- Remove gappy genes or sites
- Remove genes on long branches in gene trees
- Remove genes where contamination seems to persist.
- Supermatrix and Supertree estimation

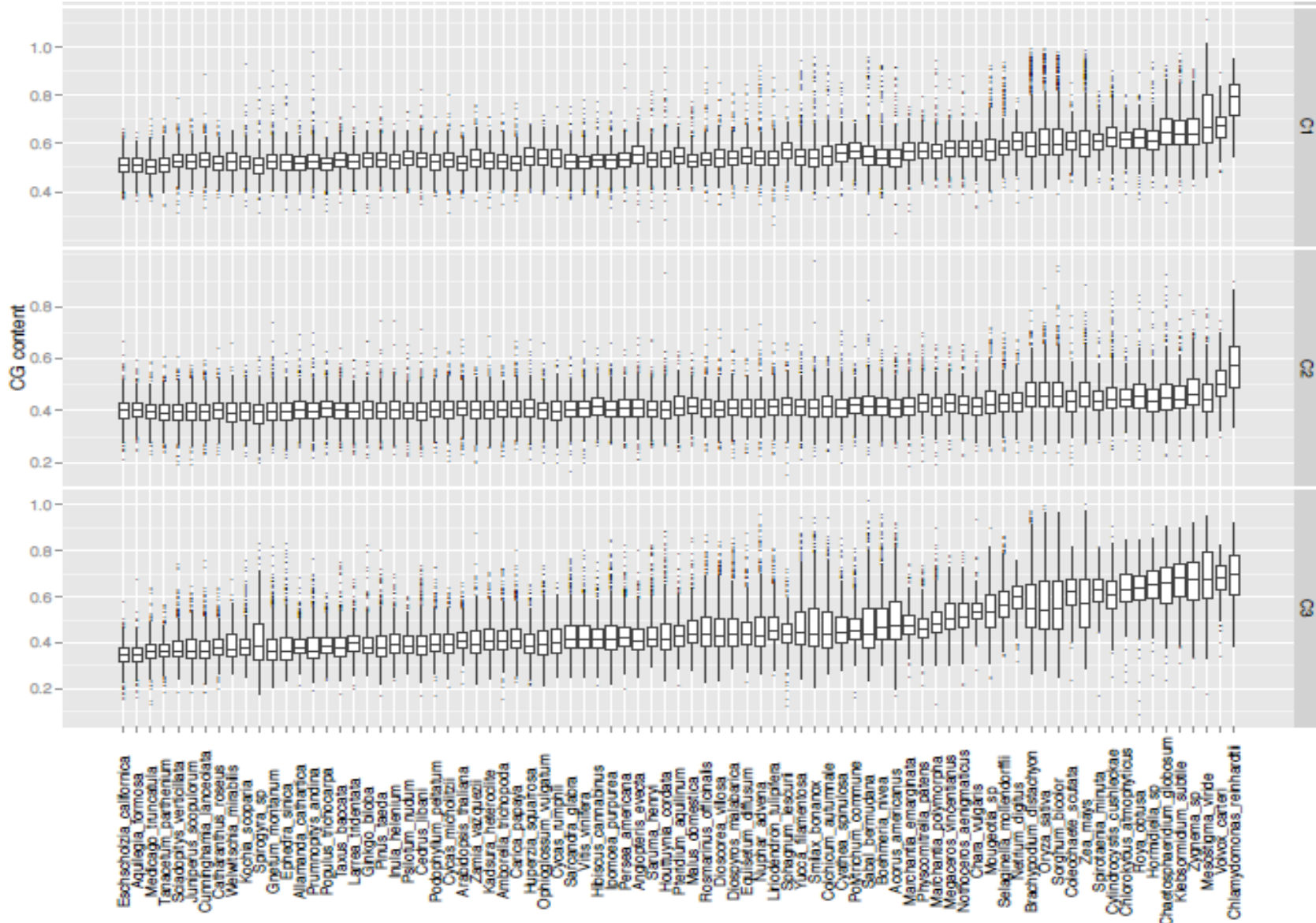
Good News: Some results seem to be very robust

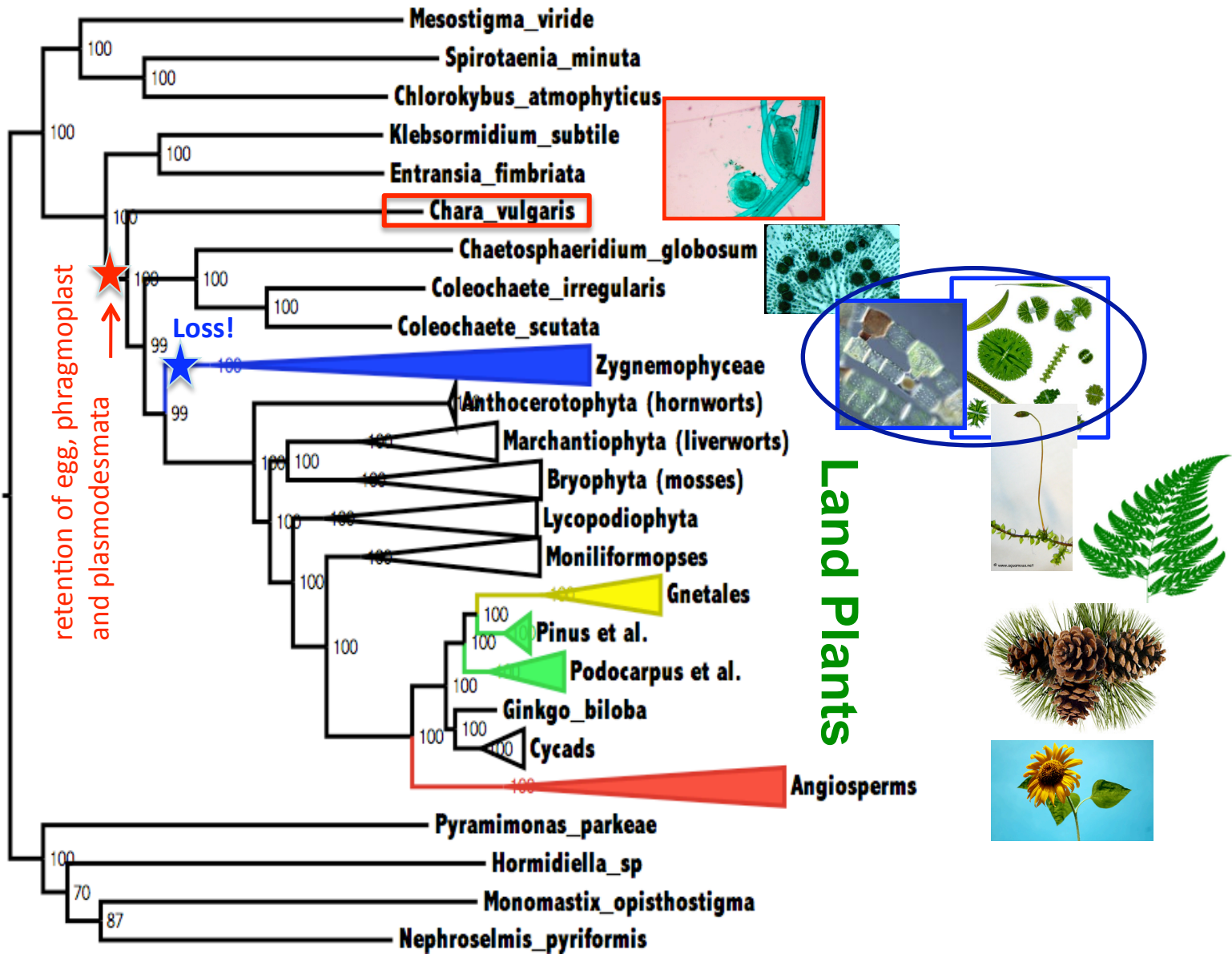


Issues with inclusion of 3rd codon position?

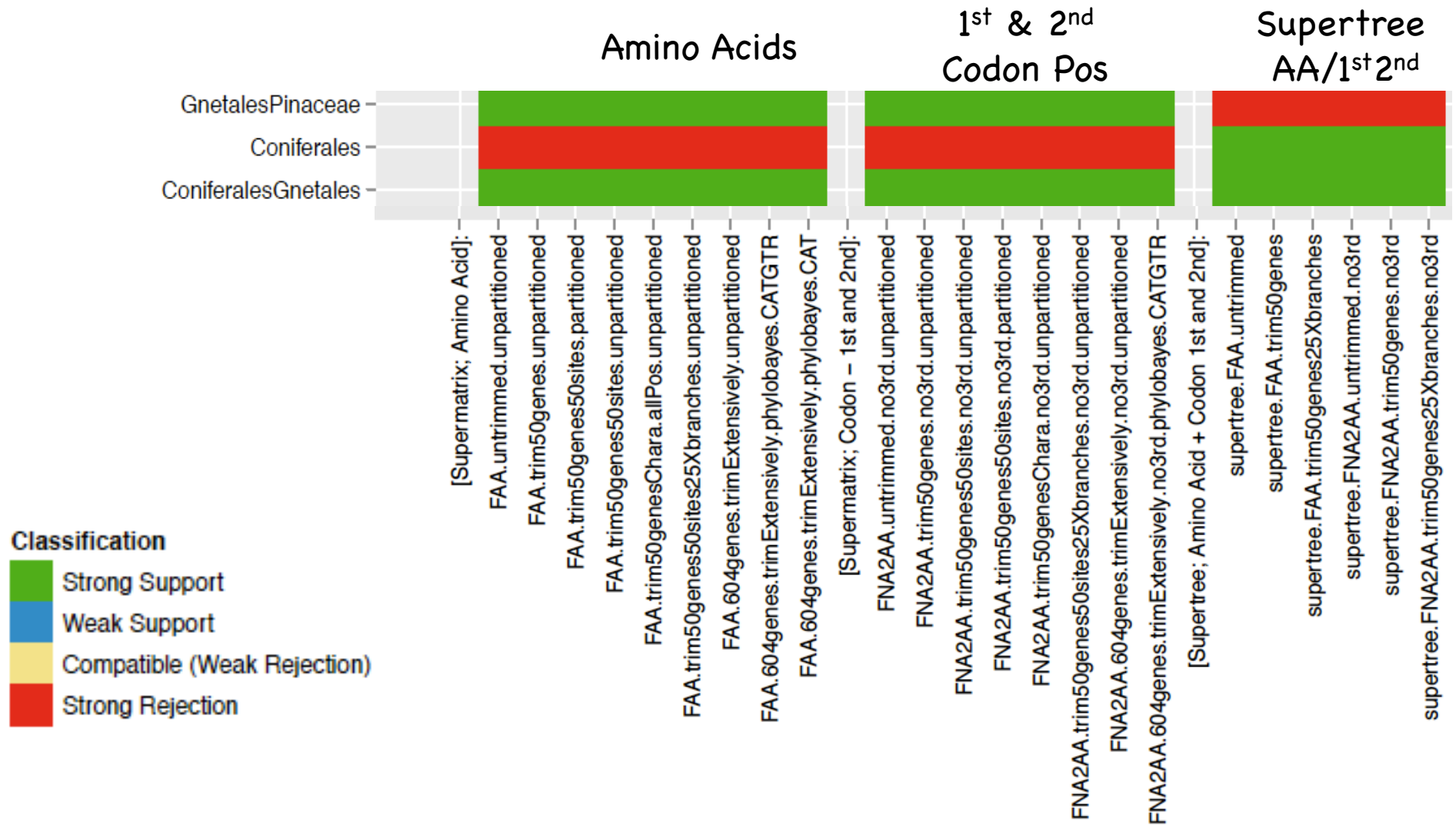


Problems due to heterogeneous substitution process?





Supertree and Supermatrix Analyses Provide Different Inferences Concerning the Monophyly of Conifers



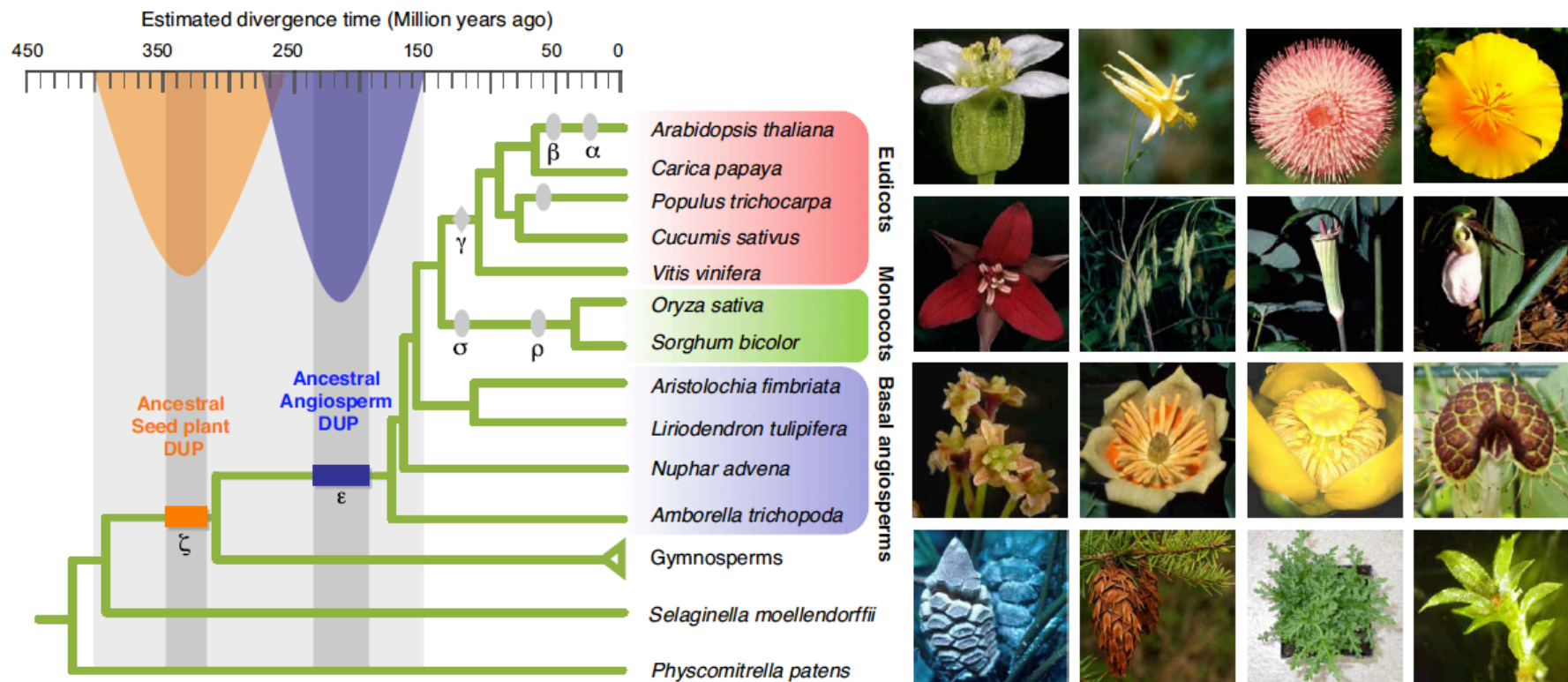
Conclusions

- Genome-scale phylogenetic inference is complex!
- Model mis-specification can result in statistical inconsistency: more data → stronger support for the wrong answer.
- Proposal: Apply multiple analysis strategies and explore/understand basis of conflicts among resulting trees

Phylogenomics???

1. Using genome-scale data to resolve phylogenetic relationships
2. Genome-Scale comparisons placed in a phylogenetic context

Phylogenetic Analysis and Molecular Dating of 100's of Gene Families Implicates Genome Duplications Associated with Origin of Angiosperms and Seed Plants



Jiao et al. 2011

Thank You!