# Exploring, improving, and evaluating anchored hybrid enrichment data support for relationships within the Family Cicadidae

Chris Owen
Computational Biology Institute
The George Washington University

Coauthors:
Beth Wade, Dave Marshall, Kathy Hill, Geert Goemans,
Russ Meister, Alan Lemmon, Emily Lemmon, Krushnamegh Kunte,
and Chris Simon

# Outline

1. Lessons learned from transcriptome phylogenomic datasets

2. Dealing with the unknown: what should the phylogeny look like and what are these sequences?

3. Model violation and gene filtering

4. Gene tree / species tree estimation

5. Exploring tree space

# Lessons learned from Hemiptera transcriptome phylogenomic datasets

# Taxon & Ortholog Sampling

Underlying bar chart (Single-copy orthologs vs Taxa):

| Taxon | Single-copy orthologs |
| --- | --- |
| Pediculus humanus | 4099 |
| Rhodnius prolixus | 4099 |
| Acyrthosiphon_pisum | 4099 |
| ★Triatoma_rubida | 89 |
| Triatoma_infestans | 3860 |
| Trialeurodes_vaporariorum | 3425 |
| Tibicen_tibicen | 999 |
| Thrips_tabaci | 3889 |
| Sogatella_fucifera | 3948 |
| Sitobion_avenae | 3825 |
| Schizaphis_graminum | 3210 |
| ★Riptortus_pedestris | 232 |
| Rhopalosiphum_padi | 602 |
| Popplepsalta_sp | 2666 |
| Platypedia_putnami | 1424 |
| Planococcus_citri | 3703 |
| Philaenus_spumarius_strict | 767 |
| Peregrinus_maidis | 688 |
| Pemphigus_spyrothecae | 320 |
| Pachypsylla_venusta | 3620 |
| Oncopeltus_fasciatus | 2846 |
| Oncometopia_nigricans | 471 |
| Nilaparvata_lugens | 3283 |
| Myzus_persicae | 311 |
| Megalurothrips_sjostedti | 1492 |
| ★Magicicada_septendecim | 124 |
| ★Macrosiphum_euphorbiae | 86 |
| Maconellicoccus_hirsutus | 281 |
| Lygus_lineolaris | 3416 |
| Lygus_hesperus | 3670 |
| Lygaeus_kalamii | 3757 |
| Liposcelis_bostrychophila | 3967 |
| Laodelphax_striatella | 773 |
| Kerria_lacca | 3890 |
| Homalodisca_vitripennis | 409 |
| ★Graphocephala_atropunctata | 109 |
| Graminella_nigrifrons | 3867 |
| Gerris_buenoi | 2445 |
| Frankliniella_tritici | 3921 |
| Frankliniella_occidentalis | 3502 |
| Ericerus_pela | 3374 |
| Echinothrips_americanus | 3960 |
| Diaphorina_citri | 1575 |
| Dialeurodes_citri | 2958 |
| Clavigralla_tomentosicollis | 3224 |
| Cimex_lectularis | 2718 |
| Cacopsylla_pruni | 1026 |
| Boisea_trivittata | 3319 |
| Bemisia_tabaci | 3843 |
| Bactericera_cockerelli | 3728 |
| Atrapsalta_sp | 3704 |
| Arma_chinensis | 3945 |
| Apolygus_lucorum | 3931 |
| Aphis_nerii | 3958 |
| Aphis_gossypii | 2508 |
| Aphis_glycines | 2815 |
| Anoplocnemis_curvipes | 1947 |

X-axis: **Single-copy orthologs**
Y-axis: **Taxa**

Outgroups
- Thysanoptera: 5 species
- Phthiraptera: 1 species
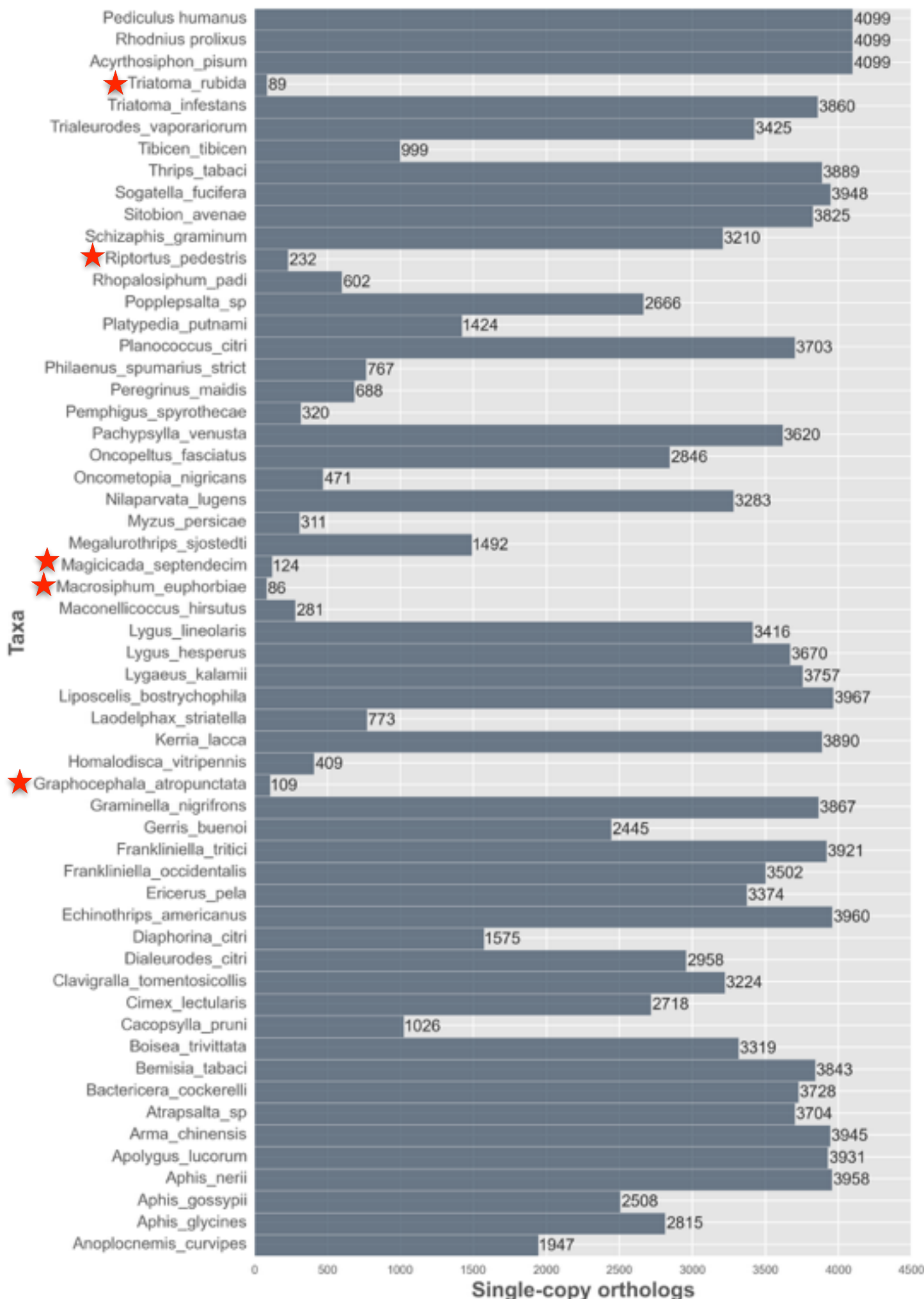- Psocoptera: 1 species

Auchenorrhyncha
- Cicadoidea: 5 species
- Cercopoidea: 1 species
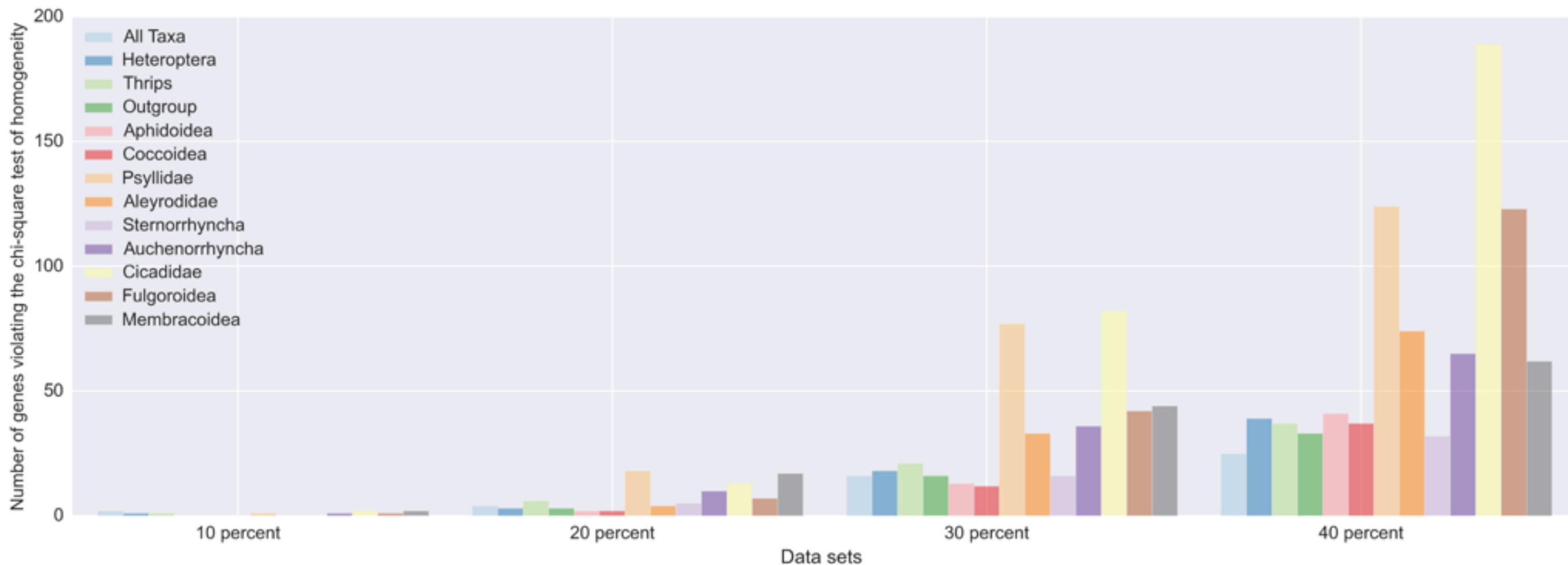- Membracoidea: 4 species
- Fulgoroidea: 4 species

Heteroptera
- Pentatomoidea: 1 species
- Coreoidea: 4 species
- Lygaeoidea: 2 species
- Cimicoidea: 1 species
- Miroidea: 3 species
- Reduvioidea: 3 species
- Gerroidea: 1 species

Sternorrhyncha
- Coccoidea: 4 species
- Aphidoidea: 10 species
- Psylloidea: 4 species
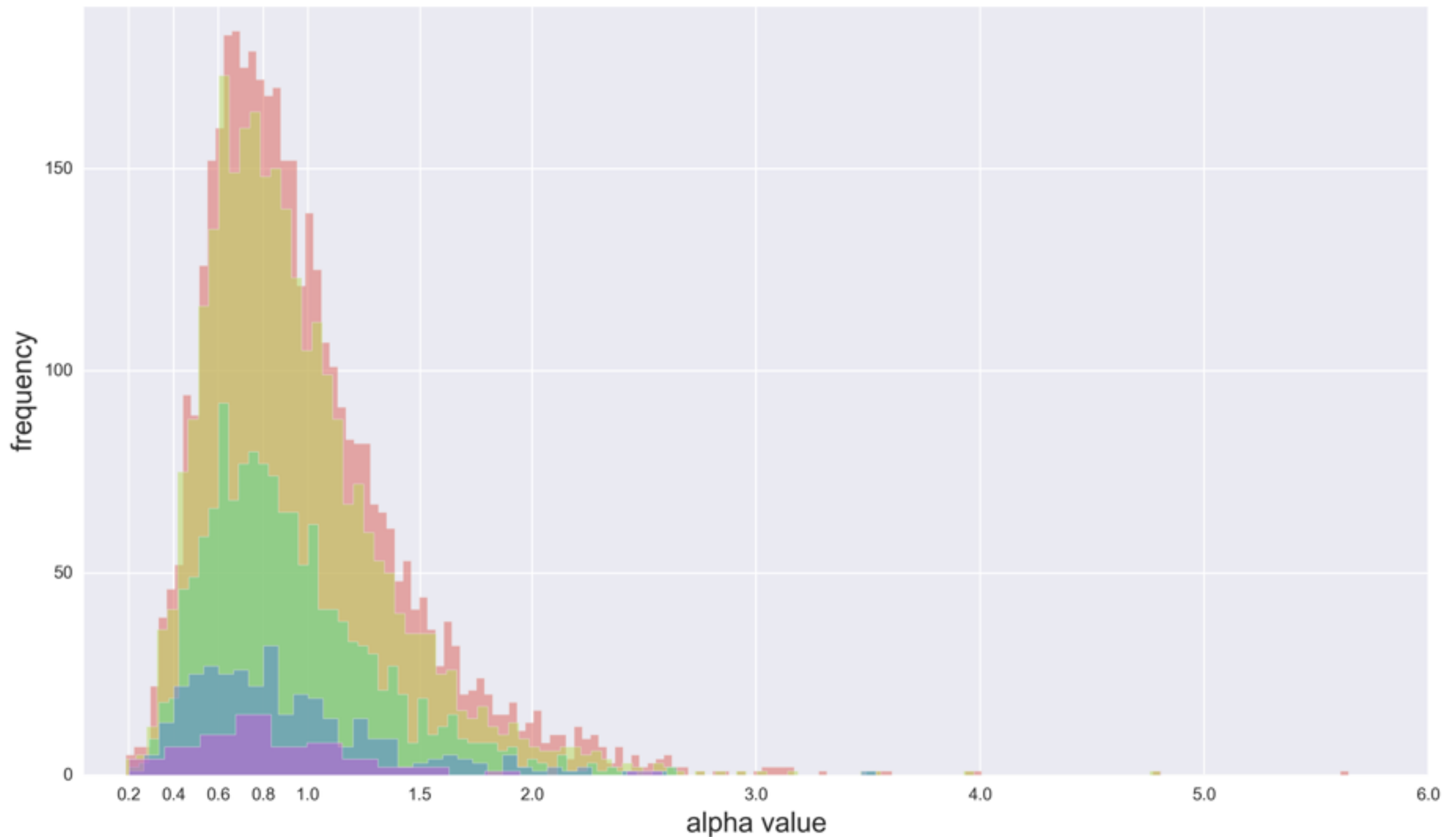- Aleyrodidae: 3 species

# Choosing the best genes: compositional bias



- Chi-square test among taxa and higher taxa (p<0.5)
- Suffers from type 2 error (Foster 2004)
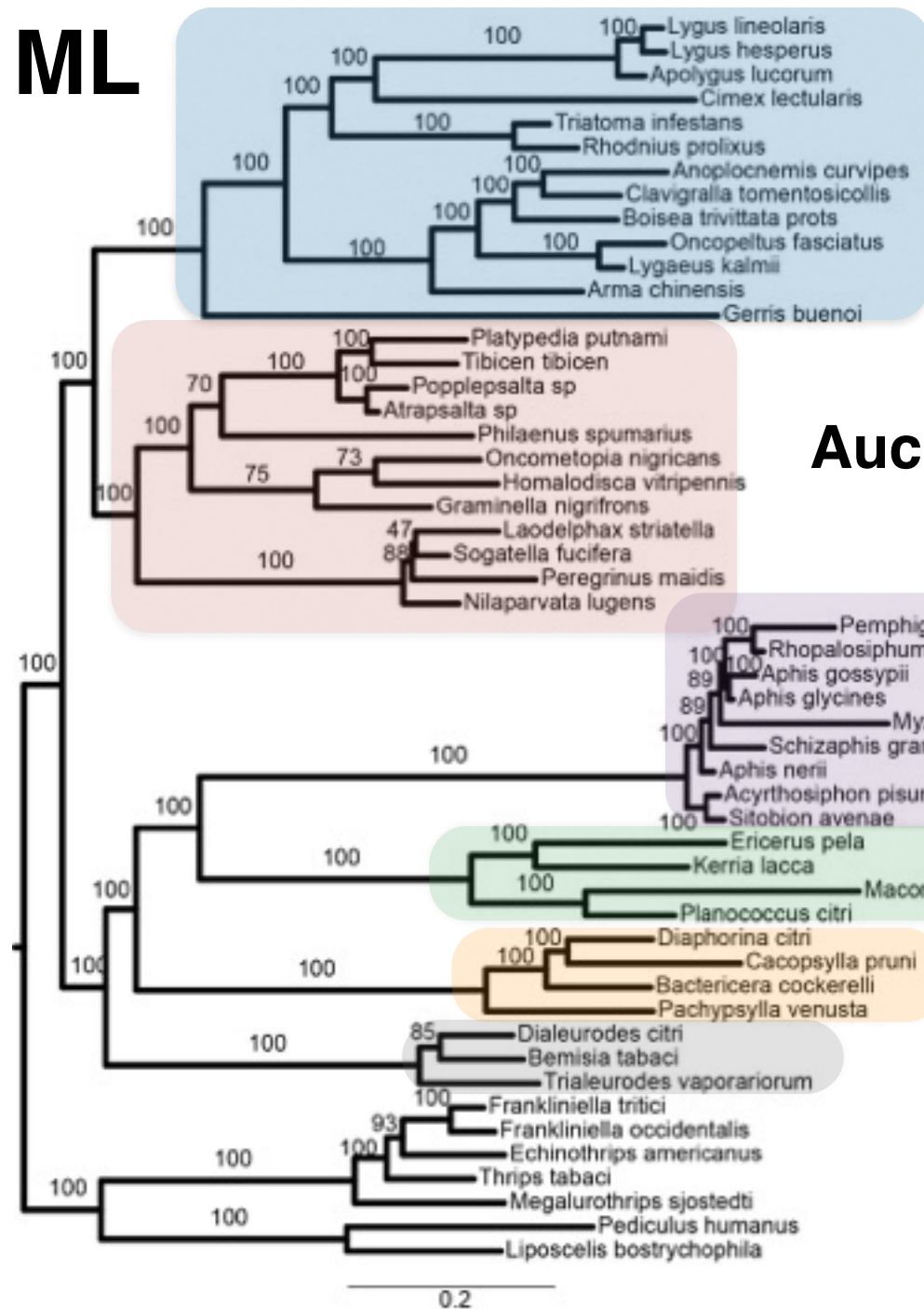
# Choosing the best genes: among site rate variation



- removed genes alpha value < 1 (Yang 1994)
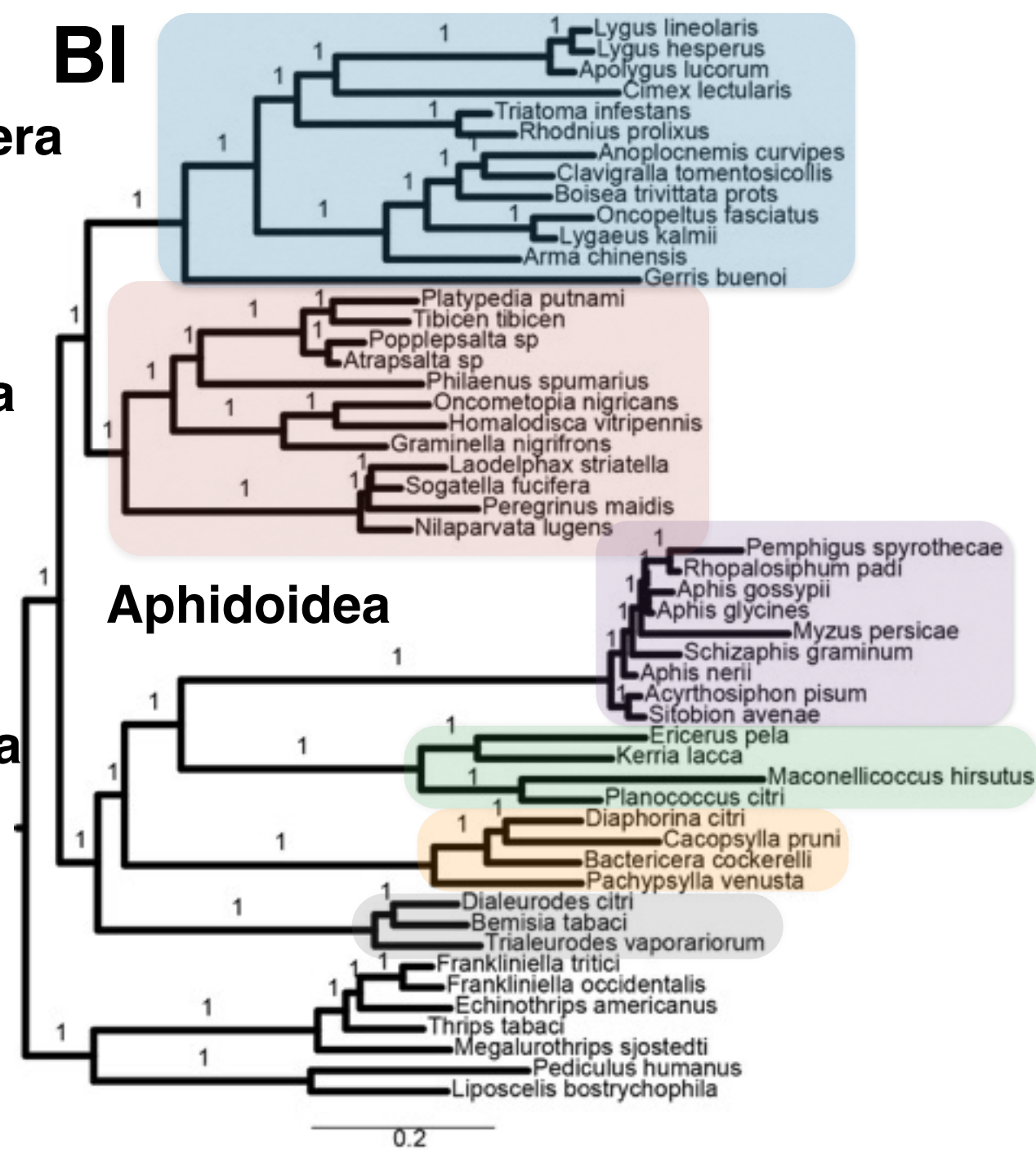
# Hemiptera phylogeny: most curated dataset

Alignment: 463 genes; 106,740 AA; 52 taxa

# Comparison of monophyly among datasets

# Cicadidae hybrid capture phylogenomics

# Dealing with the unknown: what should the phylogeny look like?

- No molecular hypotheses of subfamily and tribal relationships

- Moulds (2005) proposed subfamily relationships and Australian tribal relationships (117 morphological characters)



Tettigarctidae

Tettigadinae

Cicadinae

Cicadidae

Cicadettinae

## Gene matrix occupancy

- 150 loci; 87 taxa; ~50k bp
- 2 taxa < 90 loci
- 2 taxa 90 < loci < 100
- 1 taxon with 150 loci

# Dealing with the unknown: what are these loci?

**Pipeline**

1. Blast cds transcripts to loci or map transcriptome short reads to loci using bowtie2

2. Use MACSE (Ranwez et al. 2011) for first alignment

3. Use Muscle for refinement (-*refine*) to clean it up

**Results**

147 loci are a combination of coding and non-coding

3 loci are non-coding UTR's

# Sequence partitioning, modeling, and phylogeny estimation

1. Each gene partitioned by coding PartitionFinder and non-coding

2. Models estimated in PartitionFinder

3. Gblocks / no Gblocks

4. Phylogenies:

    • Partitioned gene trees: Garli

    • Concatenated matrices (including k-means cluster modeling (Frandsen et al. 2015)): RAxML

    •Concatenated partitioned matrices: RAxML

    •Gene tree / species tree methods: Astral

# Gene filtering: base composition



- Chi-square test in BaCoCa (Kuck and Struck 2014)
- Cicadidae hybrid capture dataset: No significant differences at $p = 0.05$ -> $p = 0.20$

# Gene tree filtering: removing loci with long branches



## Long Branch Score

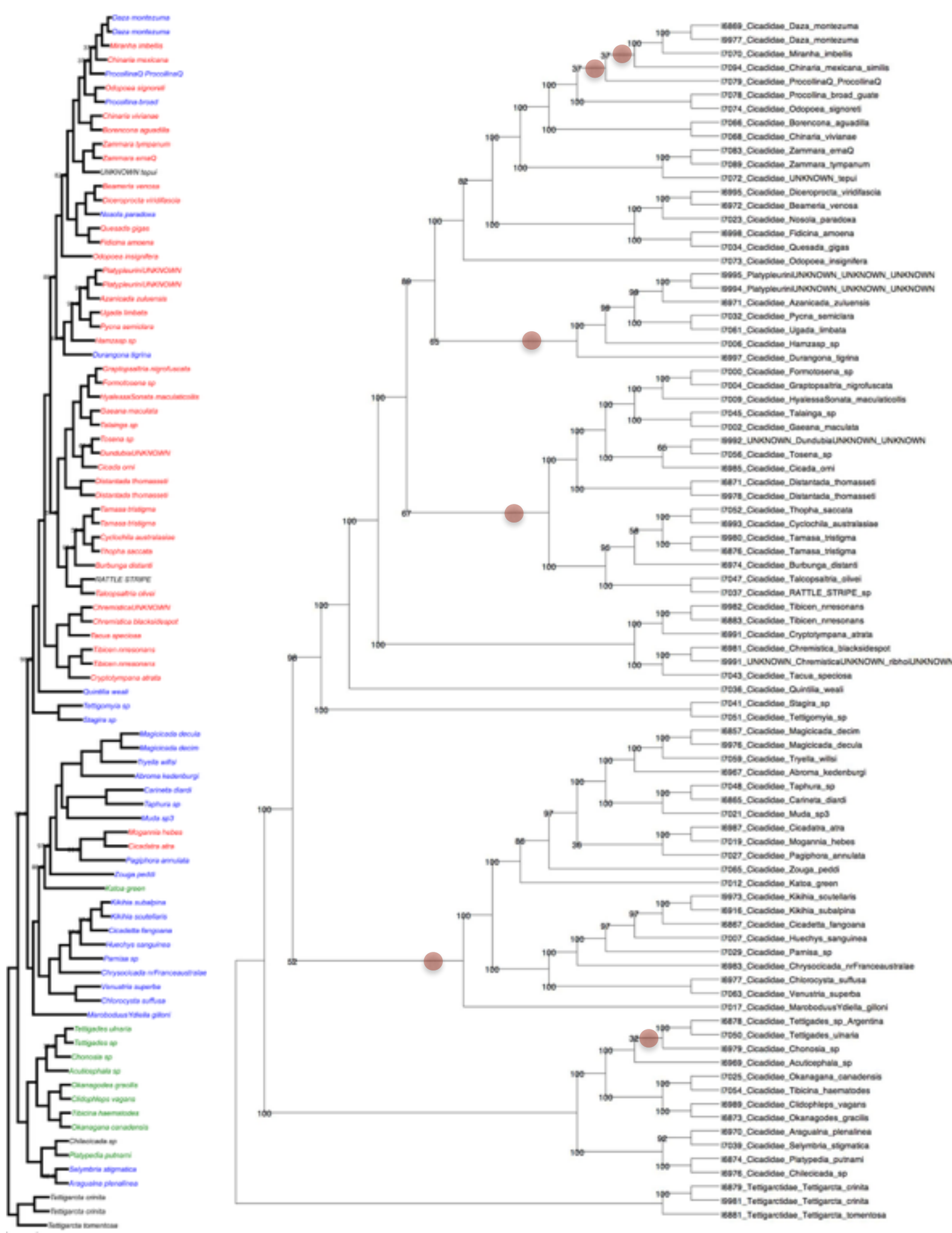$$LB_i = \left( \frac{\overline{PD_i}}{\overline{PD_a}} - 1 \right) * 100$$

- mean pairwise patristic distance (PD) of a taxon to all other taxa in the tree relative to the avg. pairwise PD over all taxa

- LB score upper quartile of each partition in TreSpEx (Struck 2014)
- unrooted gene trees
- Brinkman & Philippe 2008; Bergsten 2005
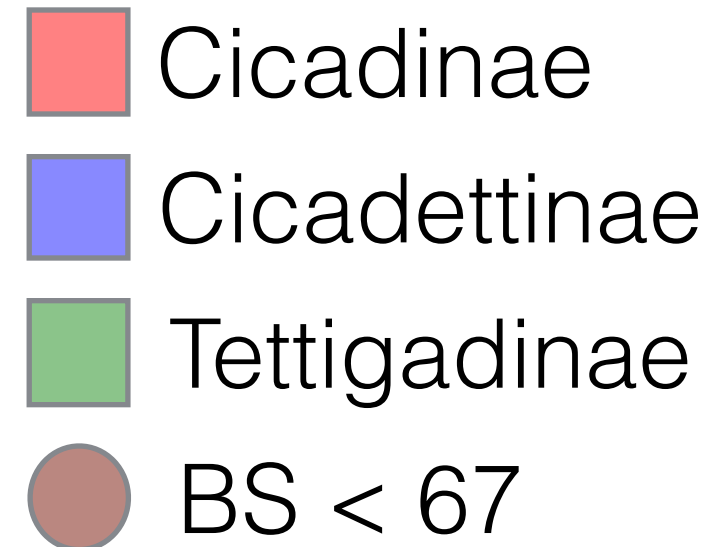
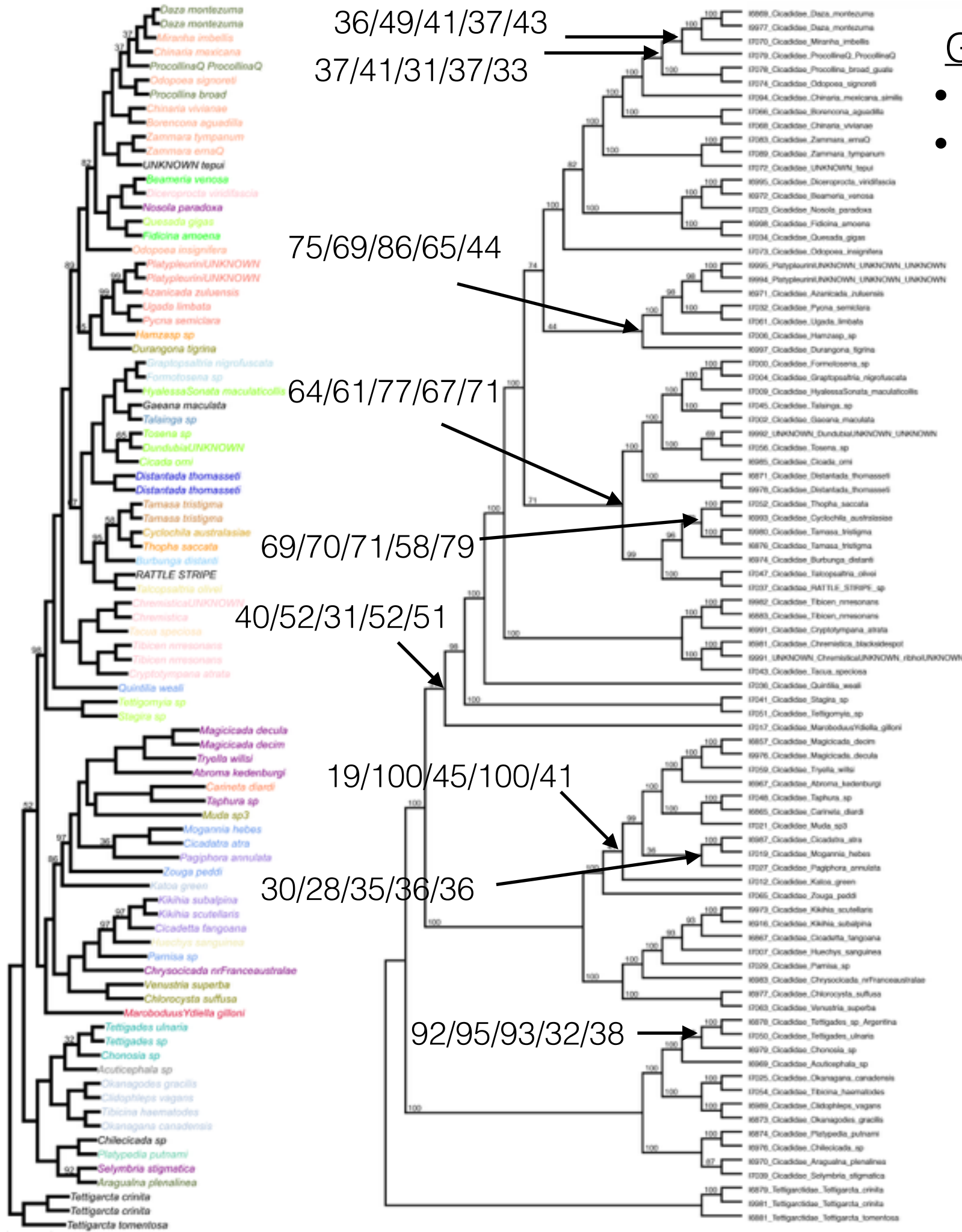# Gene filtering: Robinson-Foulds distances


Robinson-Foulds Distances PCA

- Weighted Robinson-Foulds distances
  - shared bipartitions and sum of square differences in branch lengths
- 10 best trees from each gene tree search
- Species tree from RAxML partitioned concatenated search
- Removed 12 genes

Concatenated Tree

- 6 branches not supported by bootstrap support
- 2 unsupported branches are sister to larger clades
- RAxML gene partitioned

Cicadinae

Cicadettinae

Tettigadinae

BS < 67

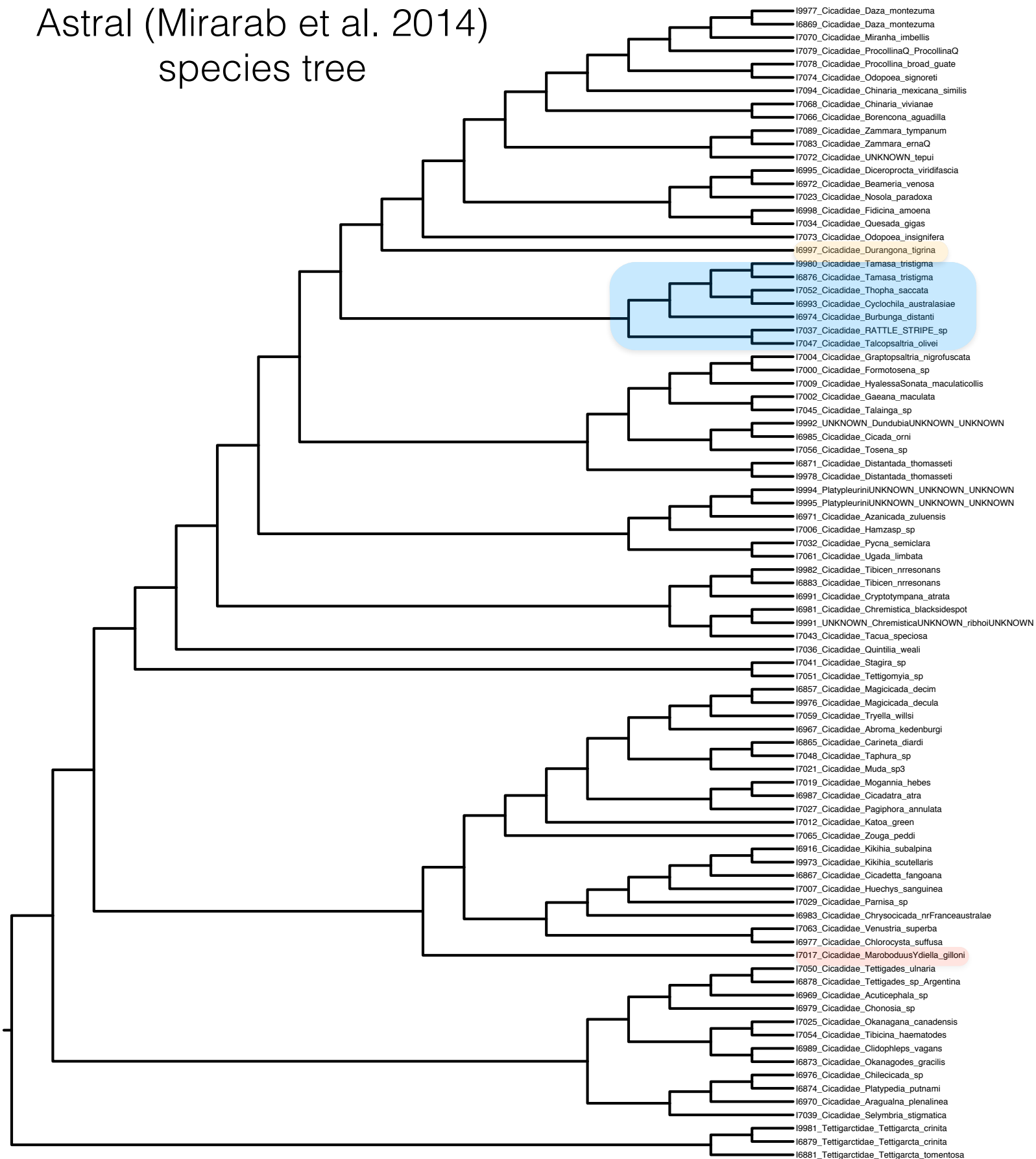Gene filtering and alternative models
- all produce the same concat tree
- six unsupported branches still unsupported to varying degrees in addition to other branches

Bootstrap support values =
no filtering, one model /
Gblocks, partitioned/
Gblocks, K-means/
RF 12 loci rm, gene partitioned/
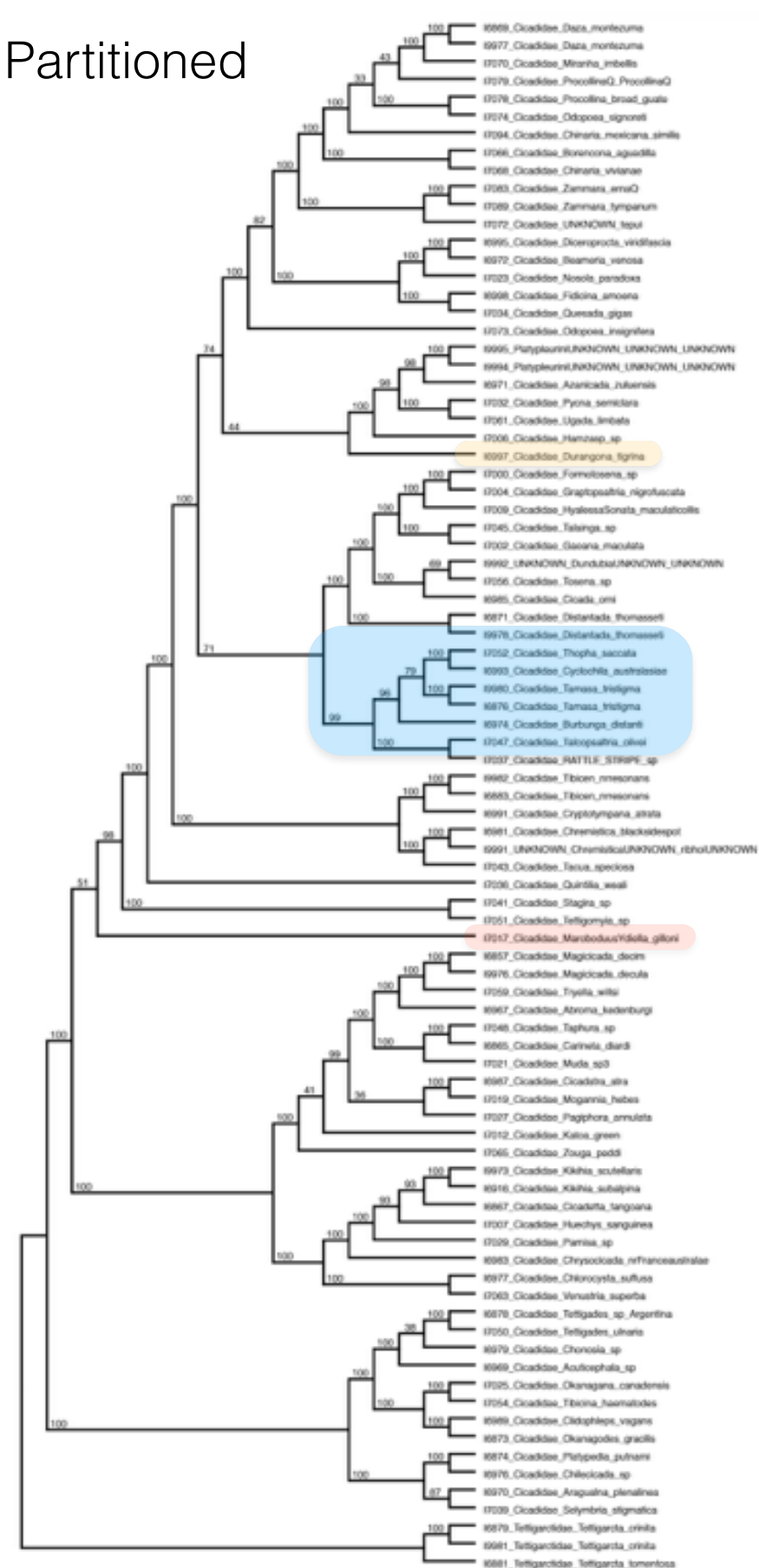RF 12 loci rm, k-means partitioned

36/49/41/37/43
37/41/31/37/33
75/69/86/65/44
64/61/77/67/71
69/70/71/58/79
40/52/31/52/51
19/100/45/100/41
30/28/35/36/36
92/95/93/32/38

unique colors = tribes

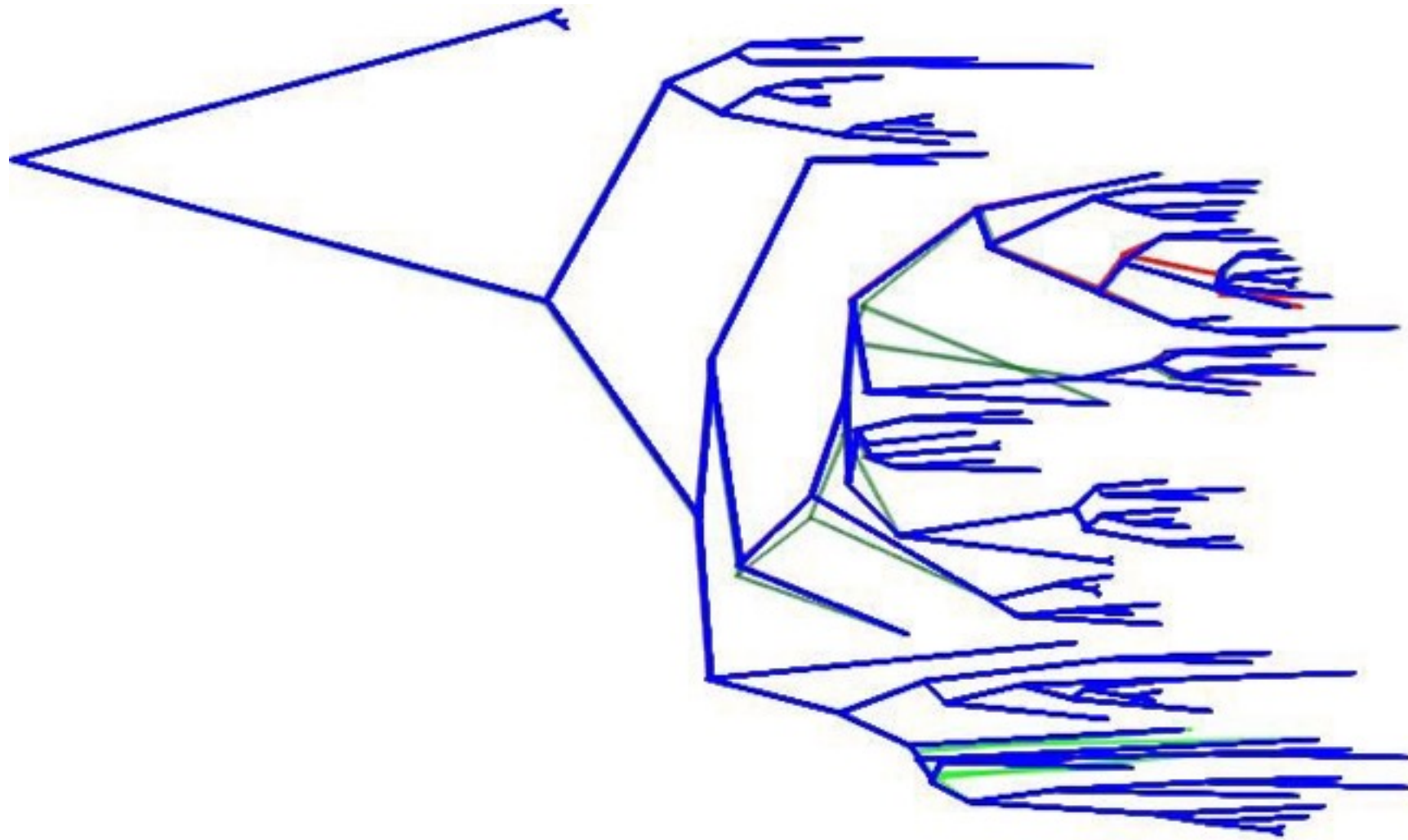# Gene tree / species tree analysis



Astral (Mirarab et al. 2014)
species tree

RAxML Partitioned

# Exploring tree space: highly divergent starting trees



- 100 random starting trees sharing no bipartitions from ML tree
- RAxML gene partitioned analyses

# Conclusions

1. Most of the tree is supported, but some rogue taxa and clades

2. Higher taxonomy needs to be reexamined in light of well-supported relationships

3. Additional sensitivity analyses needed

# Questions?