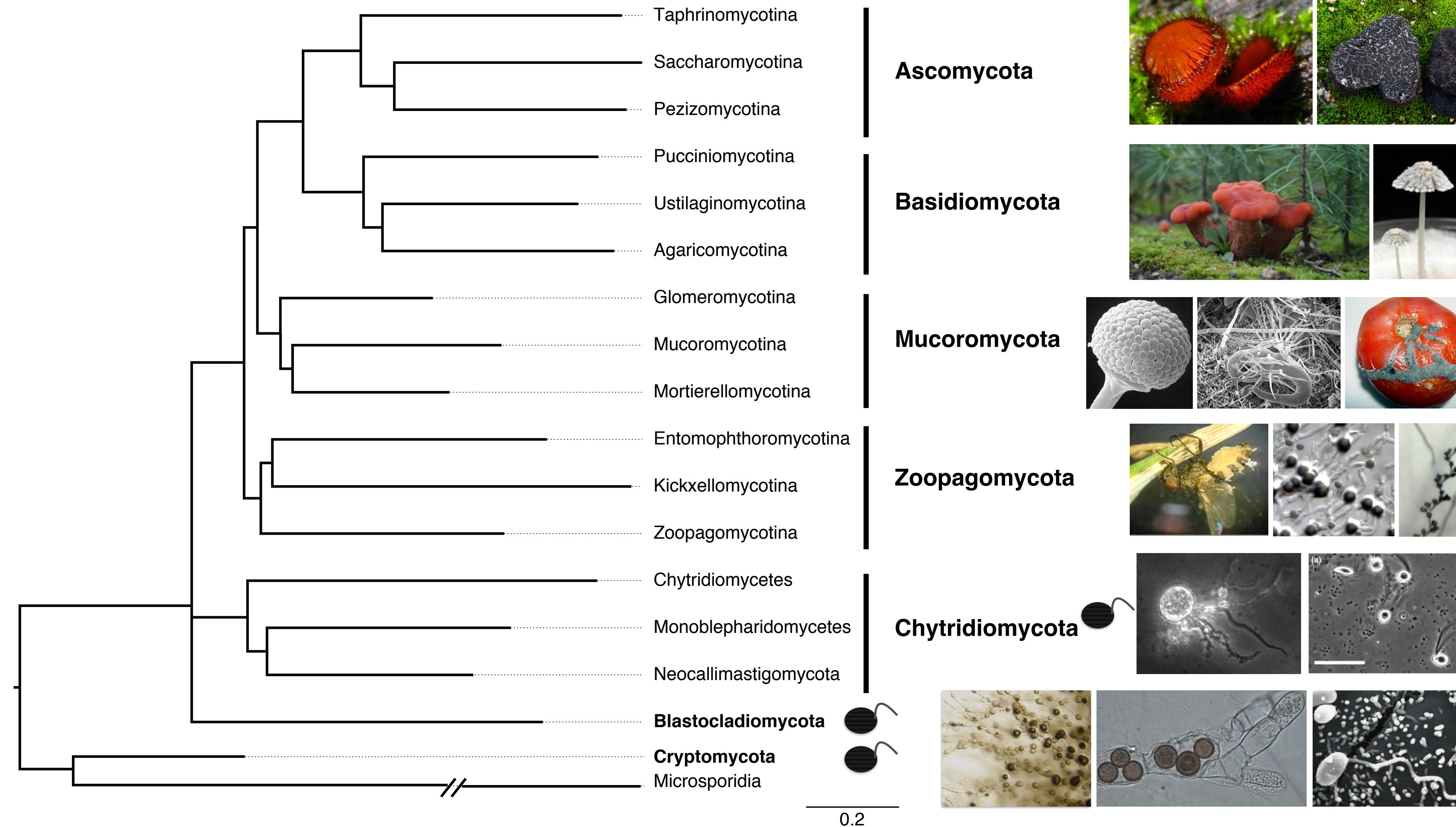


EVOLUTIONARY GENOMICS

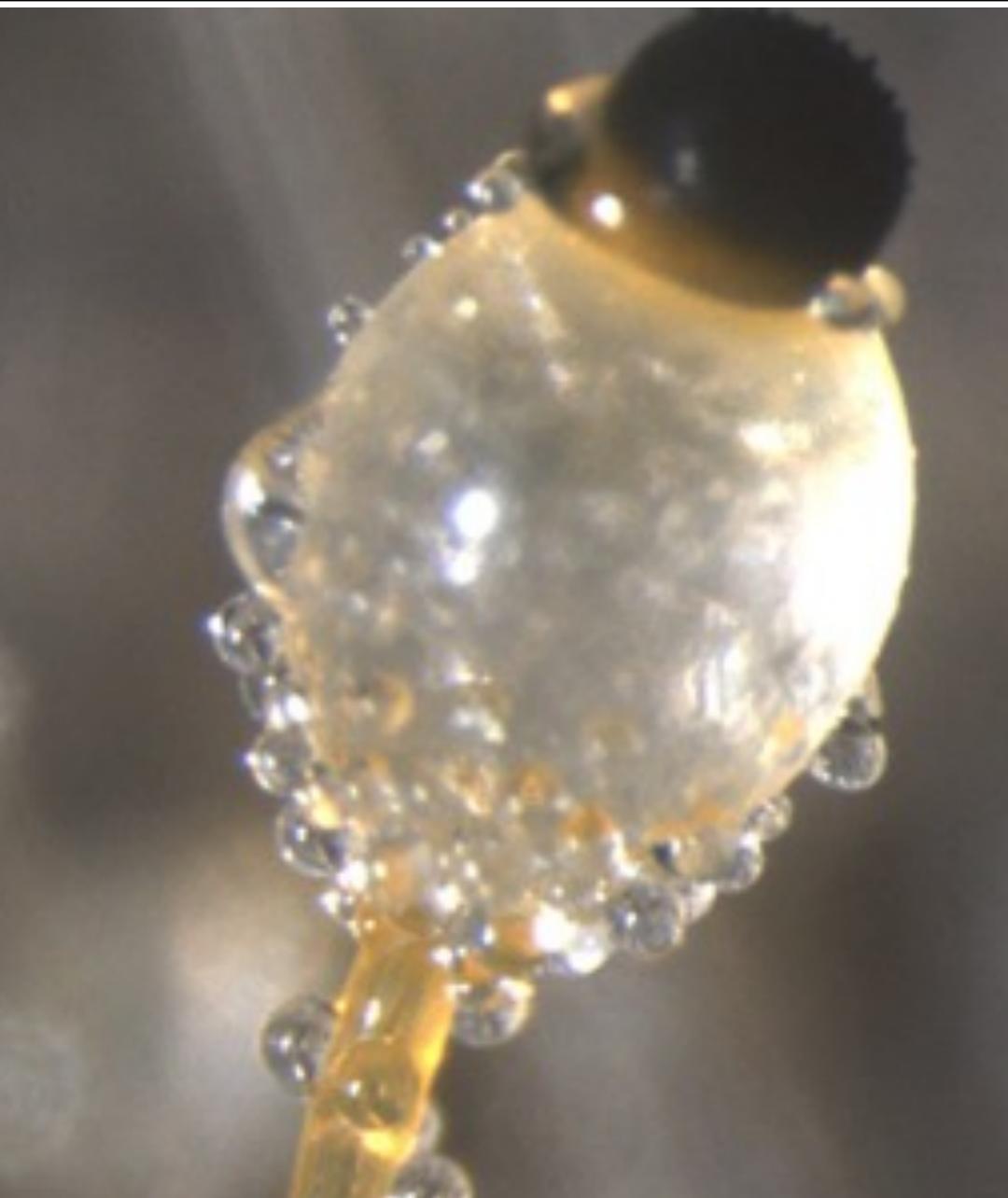
*Jason Stajich
Univ of California-Riverside*

KINGDOM FUNGI





ALESSANDRO DESIRO
ROBBIE ROBESON

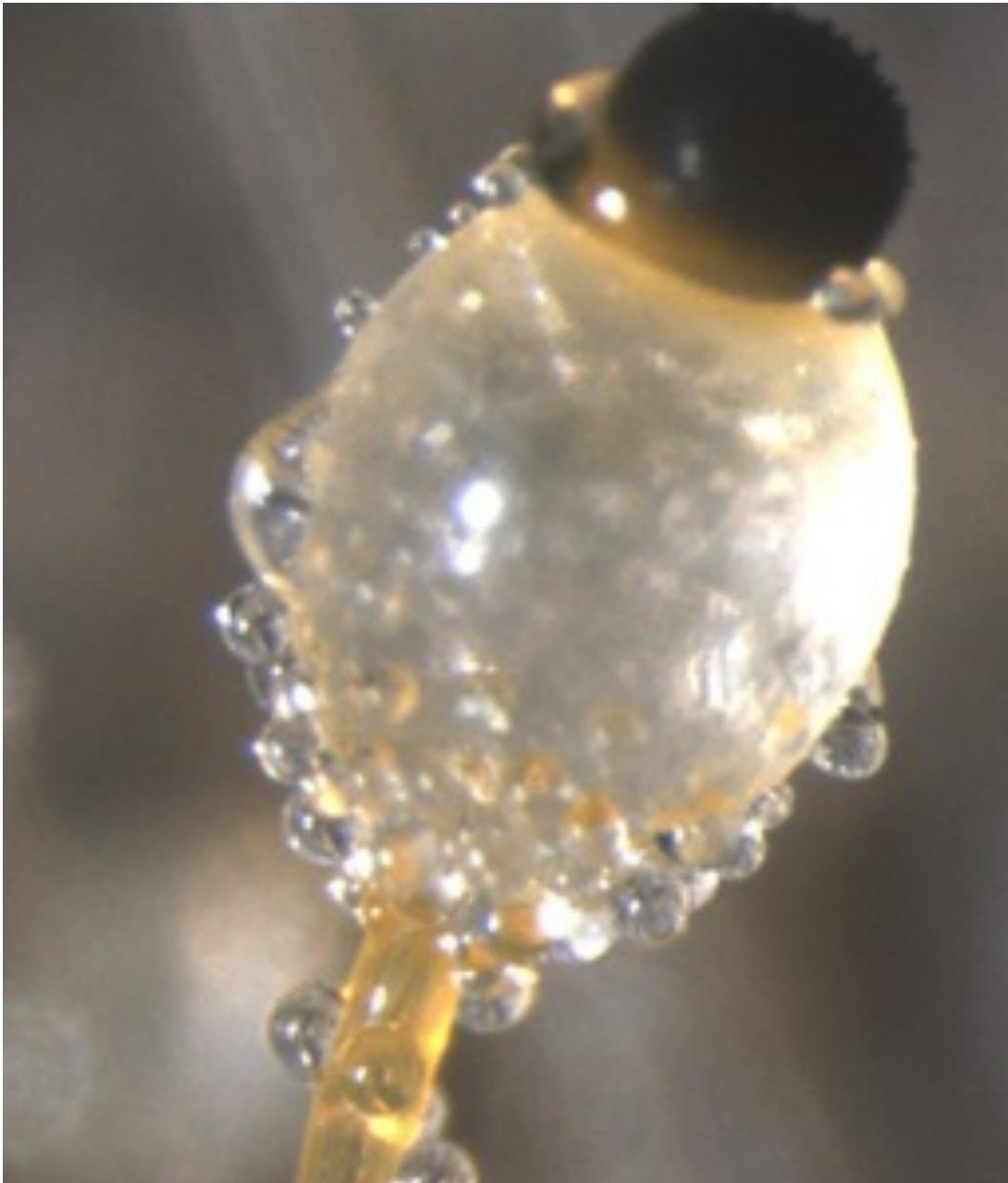
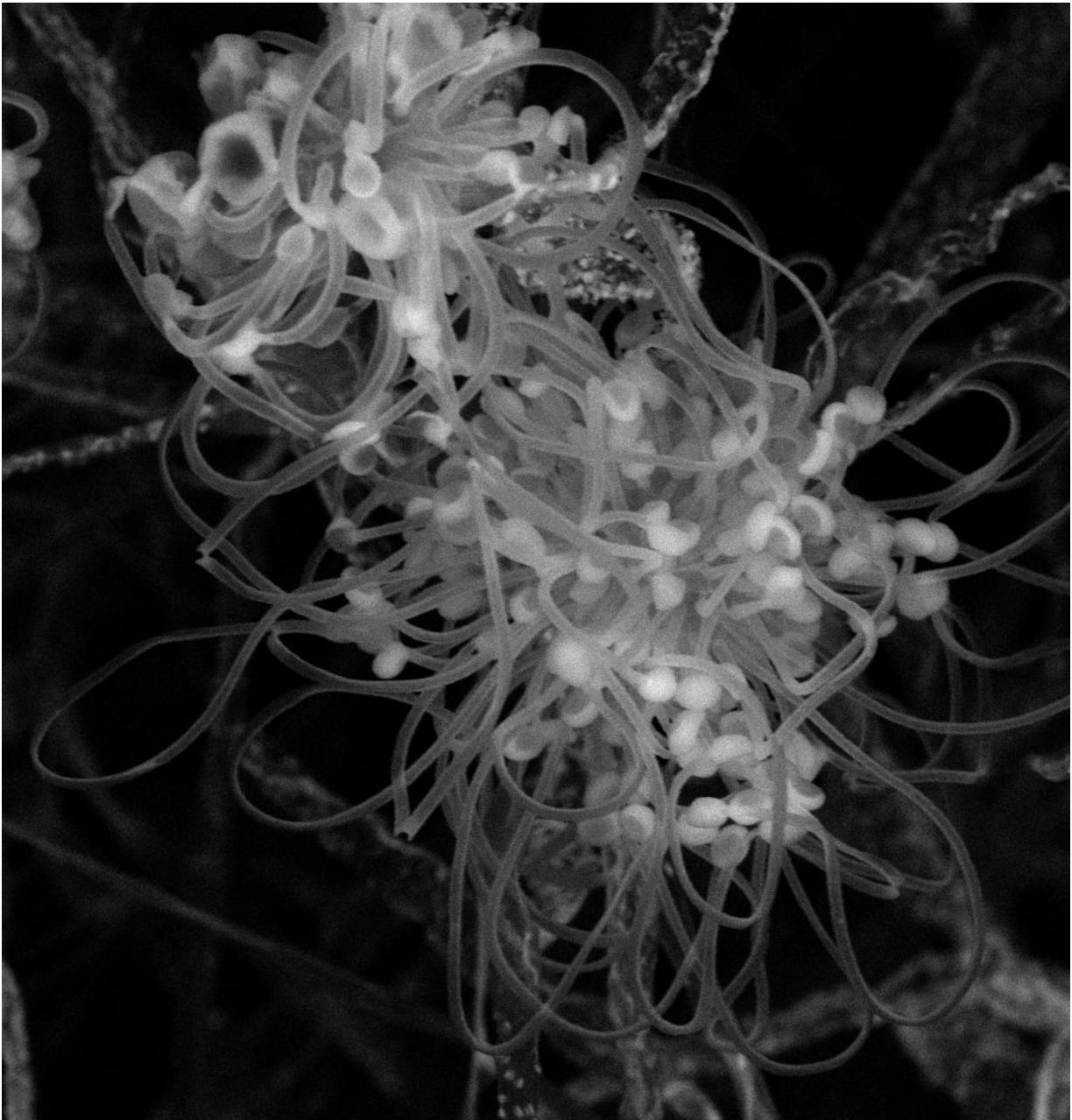


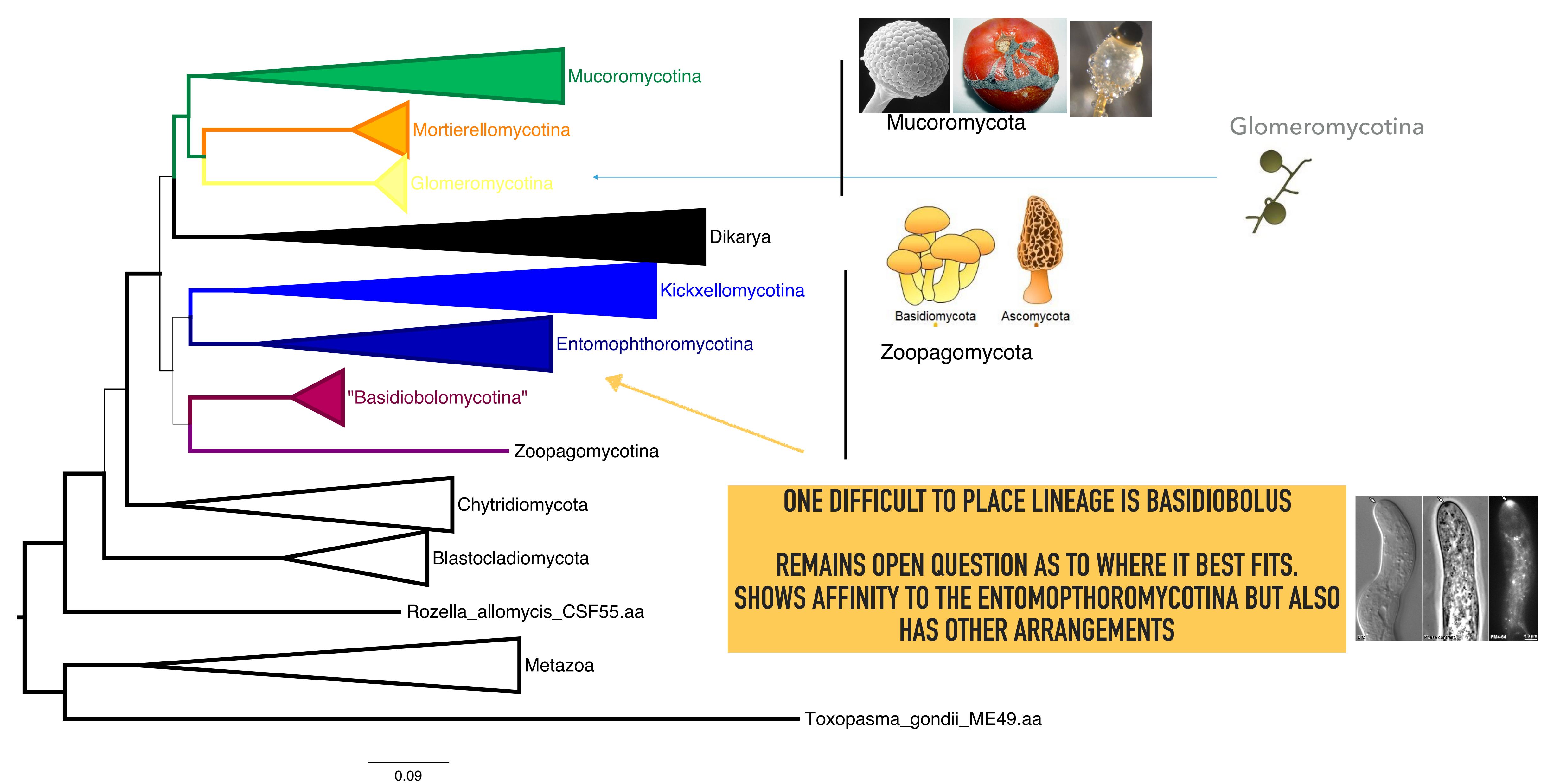
ZYGOMYCETE GENEALOGY OF LIFE

- Revisit phylogenetic relationships with whole genome data. Reference genomes and light coverage.
- Extensive genome sampling of the zygomycete phyla
- Estimating divergence time incorporating fossil data
- Examine evolution of sub-cellular characters on phylogeny (Spitzenkörper; Septa)

ZYgomycete Genealogy of Life

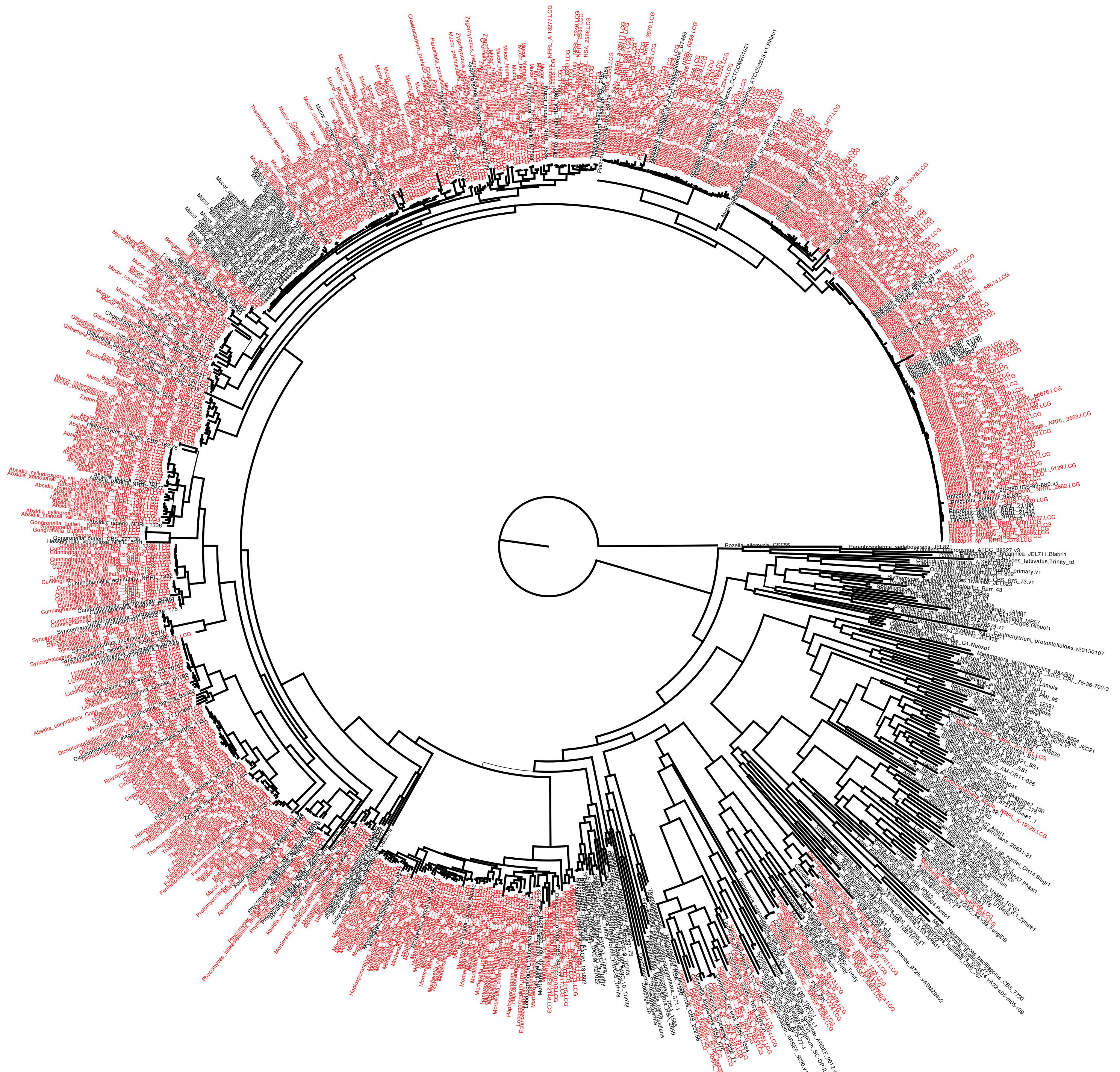
- Revisit phylogenetic relationships with whole genome data. Reference genomes and light coverage.
- Extensive genome sampling of the zygomycete phyla
- Estimating divergence time incorporating fossil data
- Examine evolution of sub-cellular characters on phylogeny (Spitzenkörper; Septa)

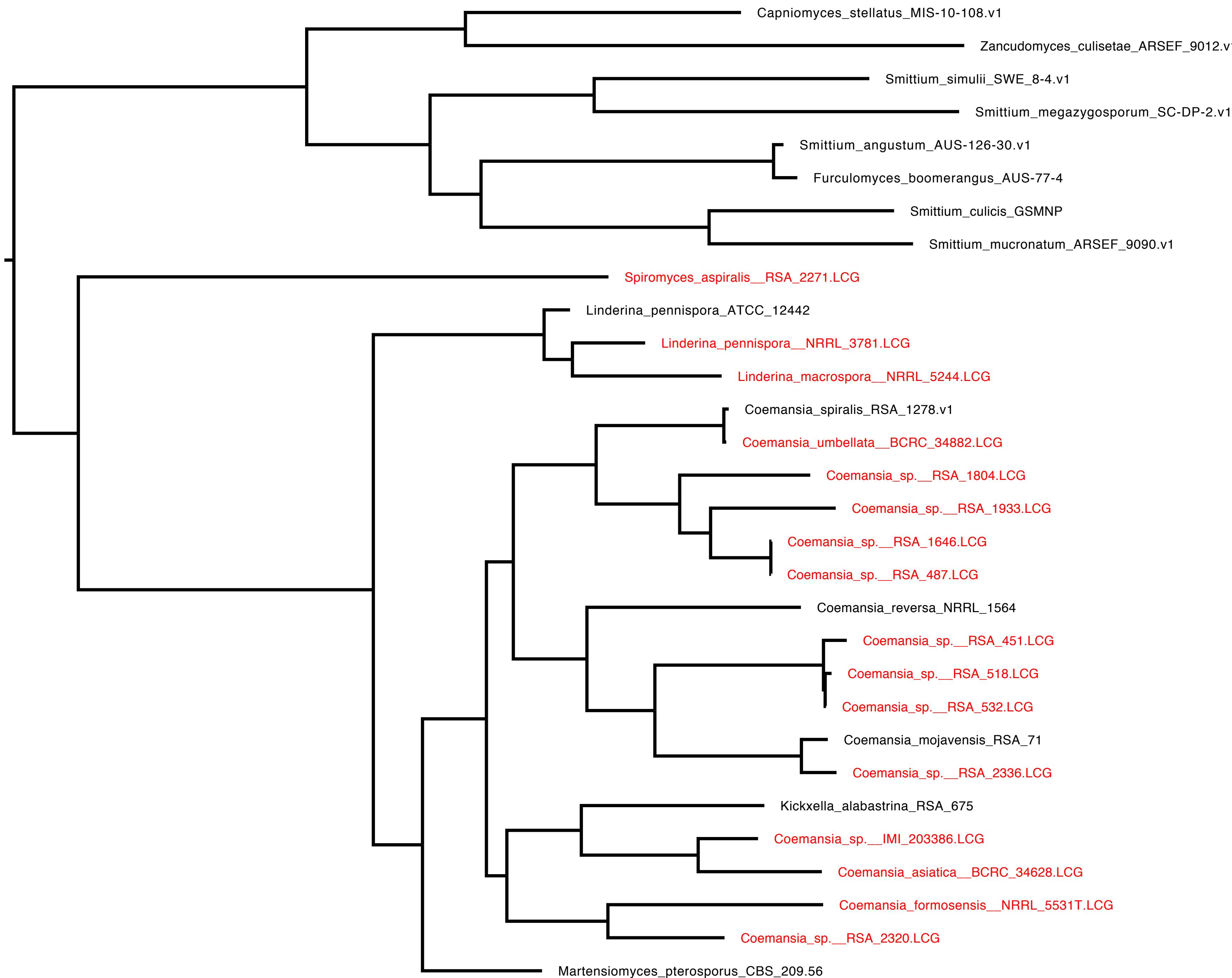




857 Taxa
520 “Light Coverage” Genomes

434 Protein coding Marker genes





GENOME EVOLUTION AND COMPARATIVE ANALYSES

Within Single Individual Genome

Summary Statistics

Genome & Chromosome Structure

Total genome size

Organization and compactness

Gene/Feature Content (domains)

Patterns of gene duplication

Genomes of multiple individuals, species

Comparative analysis

Gene Families content

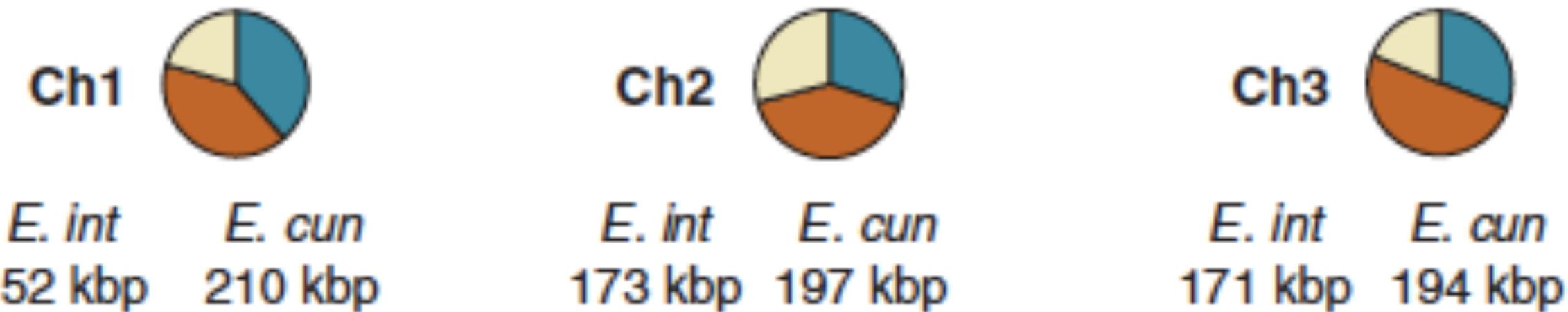
Orthologs & Paralogs

Shared vs Unique genes

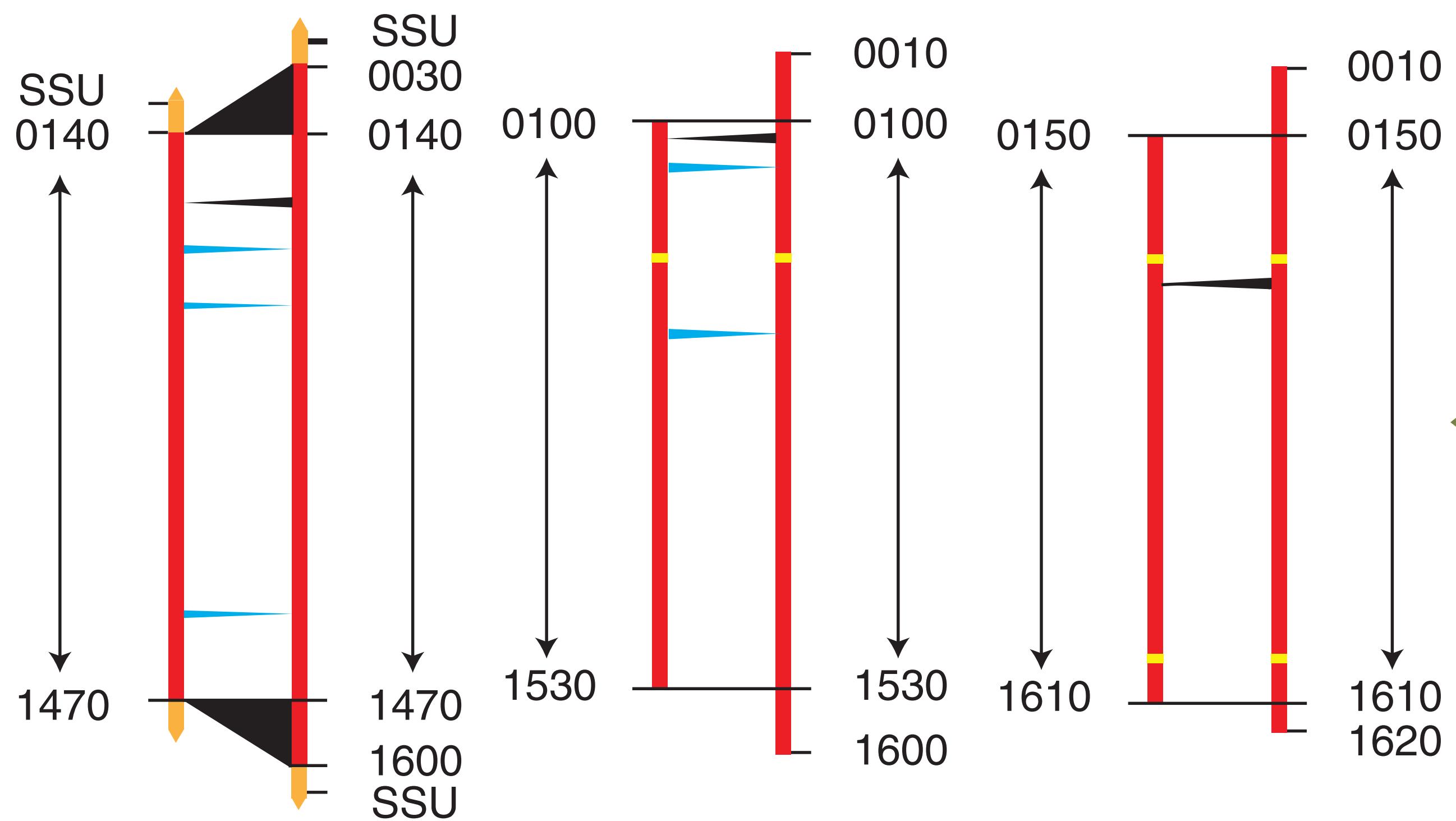
Tests for horizontal transfer

Gene Trees

Tests for Selection



Blue, brown and beige colours represent the portion of proteins that are, respectively, **shorter**, **identical** or **longer** in *E. intestinalis* compared with *E. cuniculi* orthologues



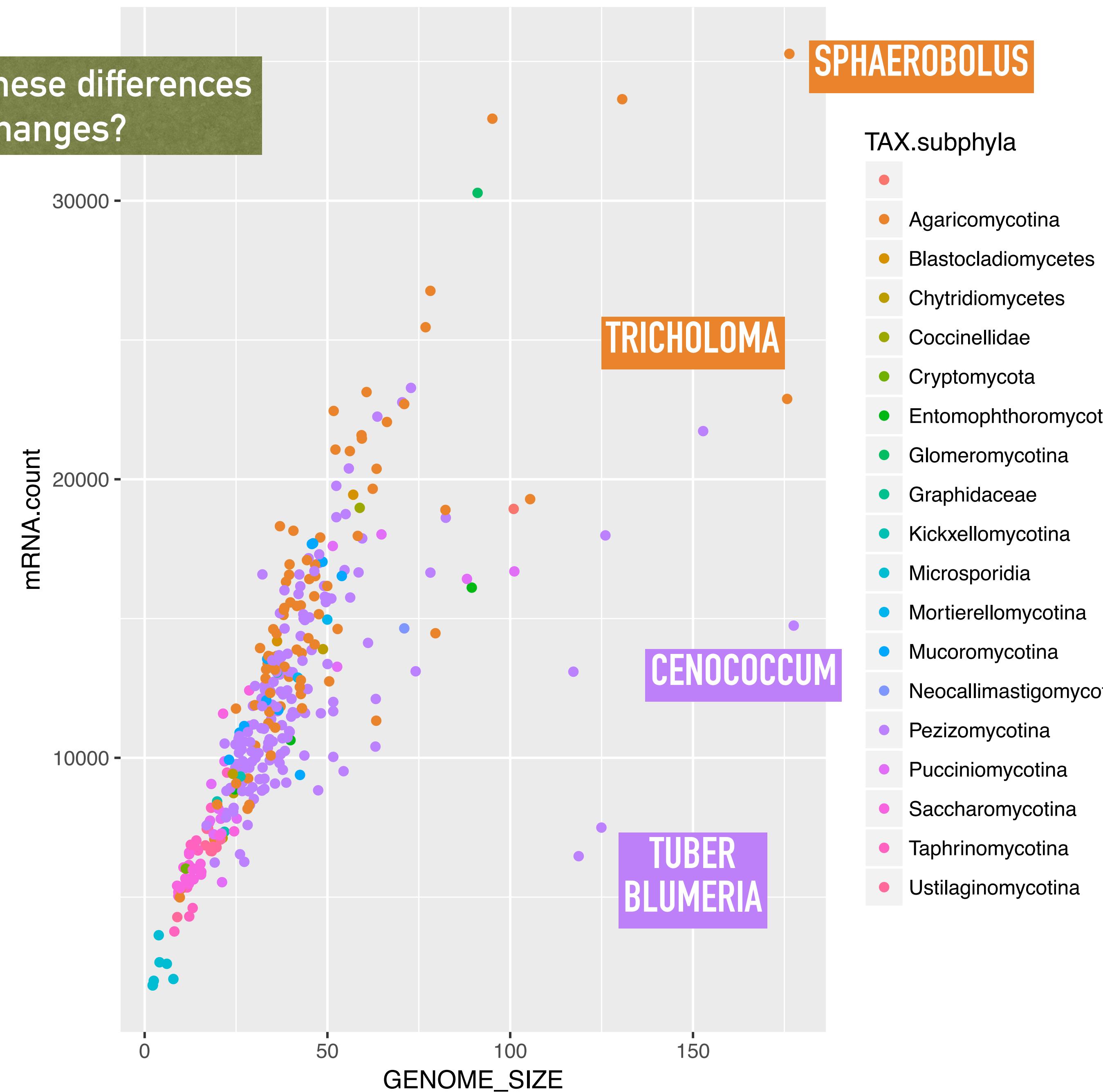
METHODS

- Genome feature investigation
 - Basic gene or genome descriptions
 - Gene content and gene

The Microsporidian *E. intestinalis* has a more compact and reduced genome as compared to its sister species *E. cuniculi*

Fungal Genome size vs Gene number

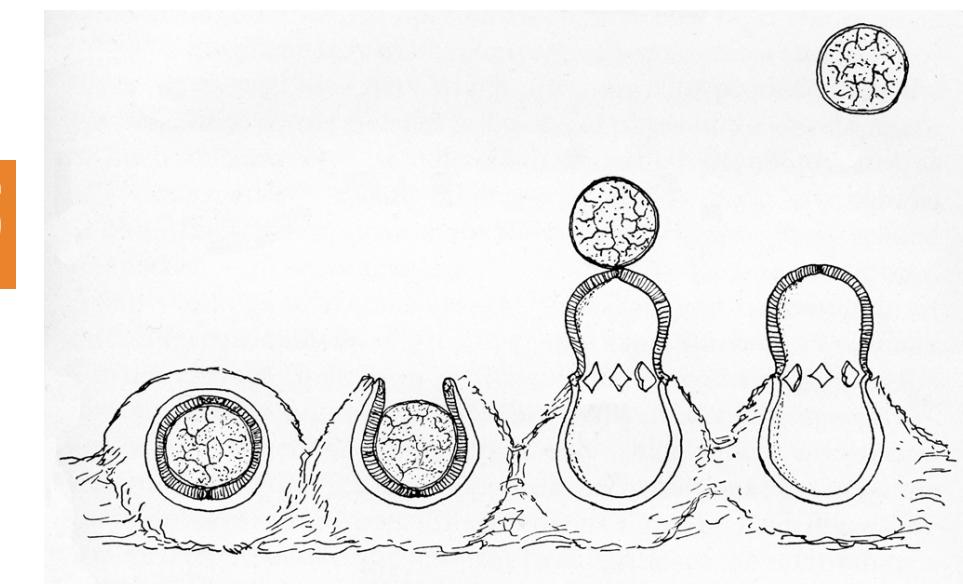
What drives these differences and changes?



TAX.subphyla

- Agaricomycotina
- Blastocladiomycetes
- Chytridiomycetes
- Coccinellidae
- Cryptomycota
- Entomophthoromycotina
- Glomeromycotina
- Graphidaceae
- Kickxellomycotina
- Microsporidia
- Mortierellomycotina
- Mucoromycotina
- Neocallimastigomycota
- Pezizomycotina
- Pucciniomycotina
- Saccharomycotina
- Taphrinomycotina
- Ustilaginomycotina

SPHAEROBOLUS



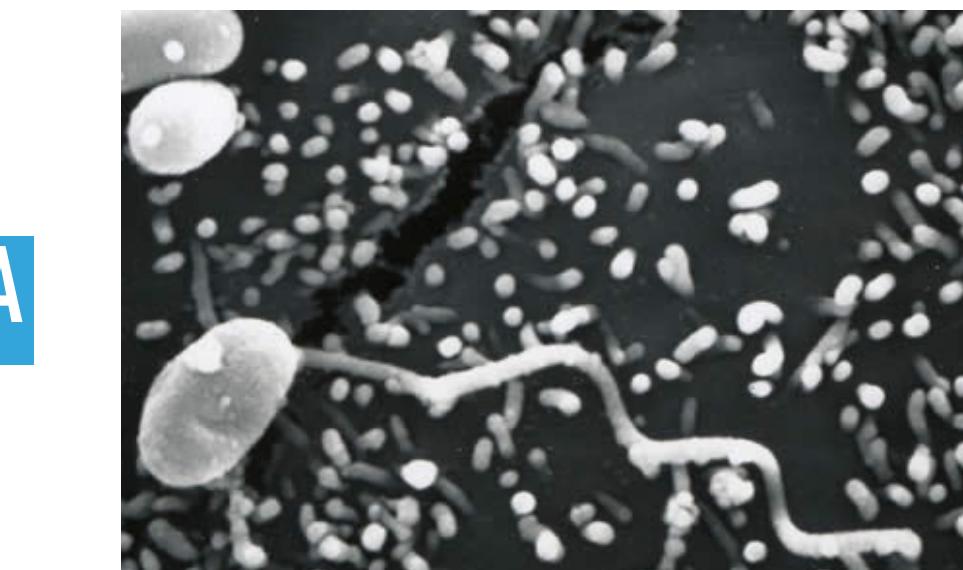
CENOCOCCUM



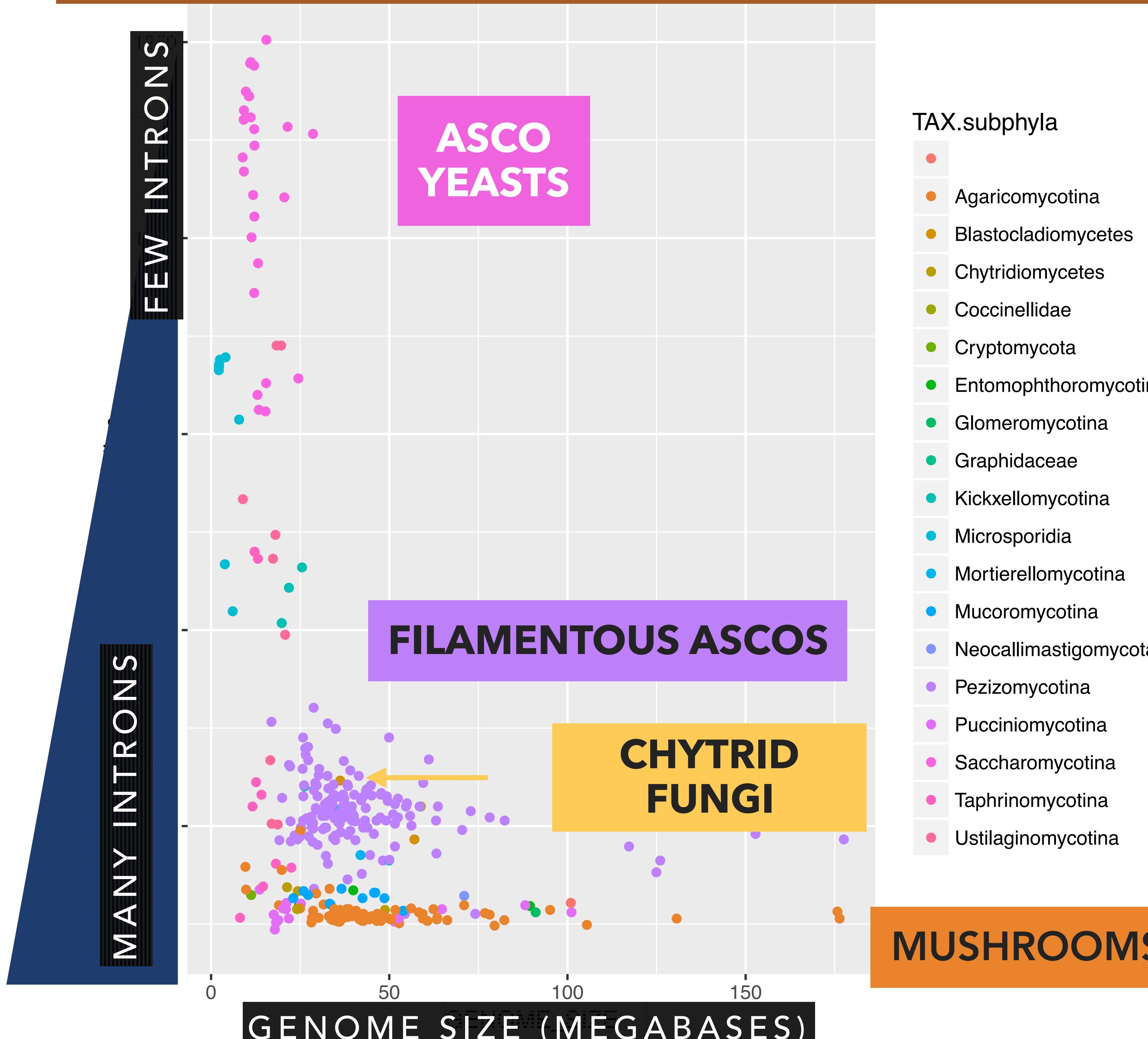
TUBER



MICROSPORIDIA



Genome size vs Intron density



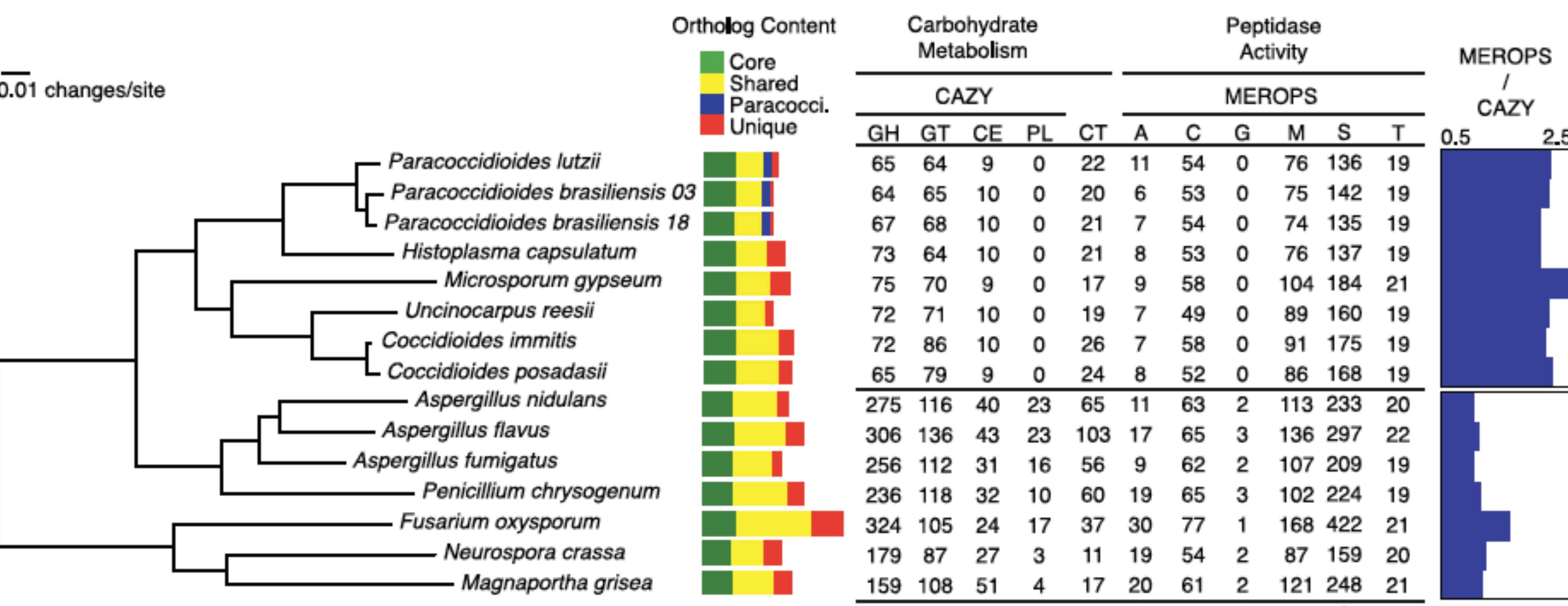
ANCESTRAL FUNGI HAD GENES RICH WITH INTRONS BASED ON OBSERVATIONS OF MANY SHARED INTRON POSITIONS AMONG ORTHOLOGOUS GENES FROM PHYLOGENETICALLY DIVERSE SPECIES

INDEPENDENT INTRON LOSS PRESSURE IN SEVERAL YEAST LINEAGES (SACCHAROMYCOTINA AND TAPHRINOMYCOTINA)

SUMMARY STATISTICS ABOUT GENOME CONTENT WITH SOME COMPARISONS

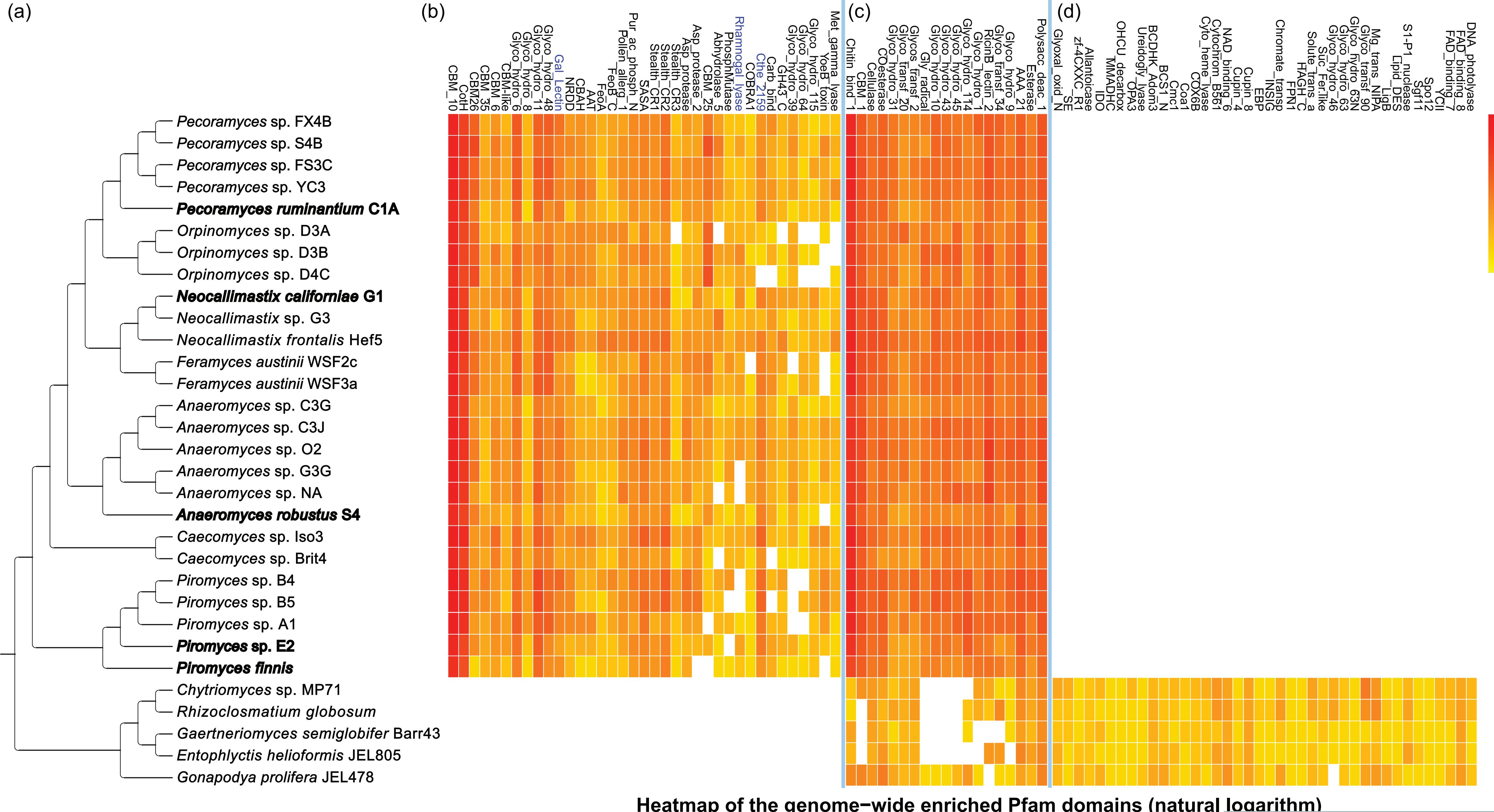
Table 1. Assembly and gene statistics.

	<i>P. lutzii</i>	<i>P. brasiliensis</i>
	Pb03	Pb18
Coverage	8.0X	8.9X
Assembly size (Mb)	32.9	29.1
Total contig length (Mb)	32.6	28.8
Scaffolds	111	65
Scaffold N50 (Mb)	1.02	1.97
Contigs	885	552
Contig N50 (kb)	84.3	114.9
Quality \geq Q40 (%)	98.9	98.9
GC (%)	42.8	44.5
Predicted protein-coding genes	9,132	7,875
Dubious genes	1,002	265
High-confidence genes	8,130	7,610
Mean gene length (nt)	1,814	1,833
Mean coding sequence length (nt)	1,330	1,433
Mean intron length (nt)	126	140
Mean intron number per gene	3.1	2.5
Mean exon number per gene	4.1	3.5
GC exonic (%)	49.8	50.4
GC intronic (%)	41.7	42.4
Mean intergenic length (nt)	1,799	1,848
tRNAs	118	103
Transmembrane proteins	1121	1057
Secreted proteins	297	291
GPI-anchored proteins	61	63

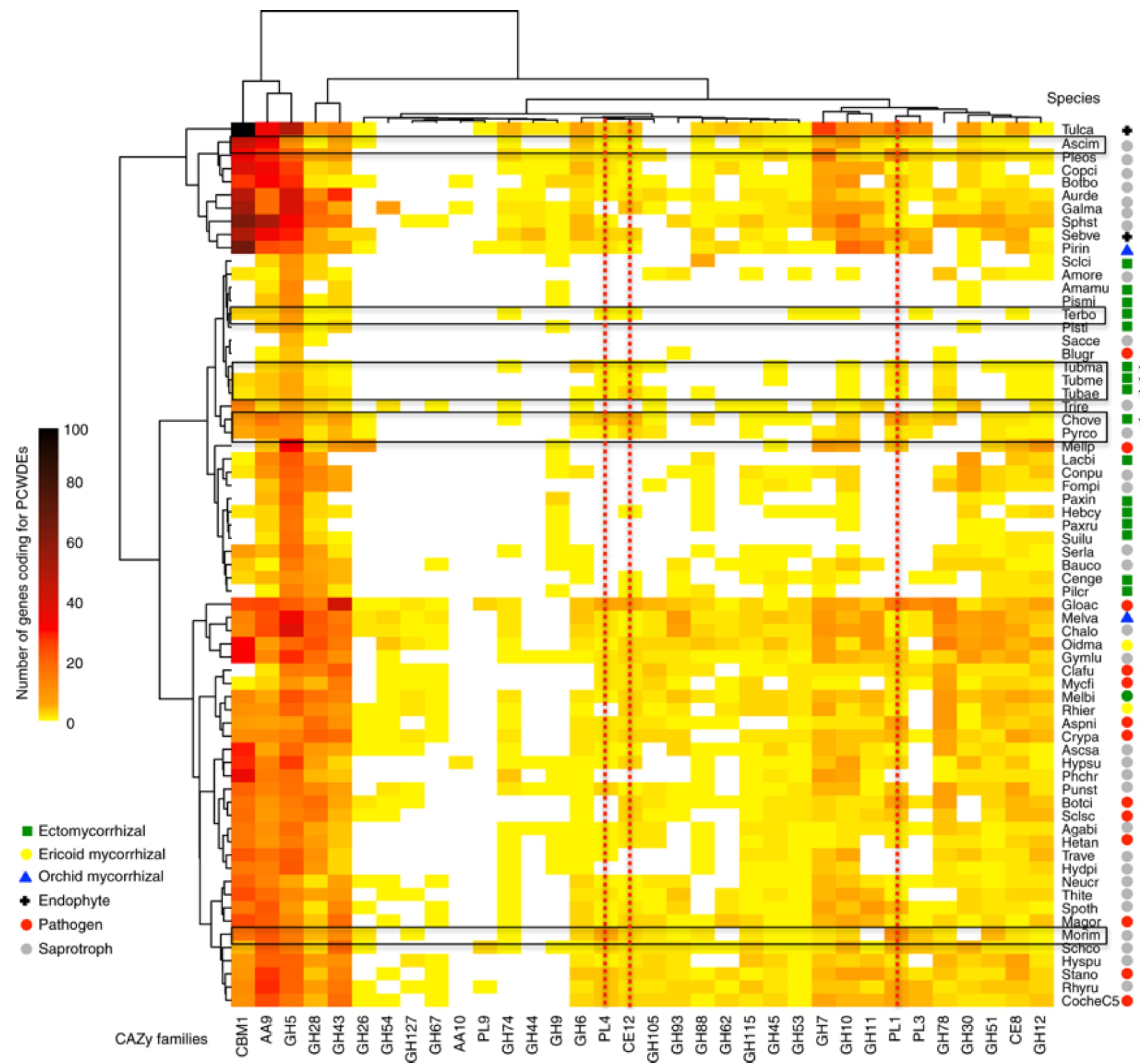


Summary / aggregated Counts by species

Phylum	Subphylum	Species	Ecology 1	Ecology 2	Oxidoreductases	CAZymes	CBM1	
Basidiomycota	Agaricomycotina	<i>Amanita muscaria</i>	Biotroph	Ectomycorrhizal	29	28	0	
		<i>Hebeloma cylindrosporum</i>			22	39	2	
		<i>Laccaria bicolor</i>			34	47	1	
		<i>Paxillus involutus</i>			23	39	0	
		<i>Paxillus rubicundulus</i>			24	35	0	
		<i>Piloderma croceum</i>			41	50	0	
		<i>Pisolithus microcarpus</i>			17	29	0	
		<i>Pisolithus tinctorius</i>			22	26	0	
		<i>Scleroderma citrinum</i>			24	25	0	
		<i>Suillus luteus</i>			29	40	0	
		<i>Tremella mesenterica</i>			Mycoparasite	5	0	
		<i>Sebacina vermifera</i>			Orchid symbiont	21	121	50
		<i>Tulasnella calospora</i>			Root endophyte	9	178	110
		<i>Piriformospora indica</i>				15	103	64
		<i>Auricularia delicata</i>	Saproth	White rot		60	140	42
		<i>Botryobasidium botryosum</i>				16	97	29
		<i>Fomitiporia mediterranea</i>				40	81	10
		<i>Galerina marginata</i>				77	128	52
		<i>Heterobasidion annosum</i>				35	65	17
		<i>Hypholoma sublatentium</i>				52	83	28
		<i>Jaapia argillacea</i>				17	102	24
		<i>Phanerochaete chrysosporium</i>				31	75	28
		<i>Pleurotus ostreatus</i>				43	111	31
		<i>Plicaturopsis crispa</i>				21	70	9
		<i>Punctularia strigosozonata</i>				46	93	26
		<i>Schizophyllum commune</i>				8	103	5
		<i>Sphaerobolus stellatus</i>				(284)	(195)	0
		<i>Trametes versicolor</i>				50	94	23
		<i>Coniophora puteana</i>	Brown rot	Brown rot		16	80	3
		<i>Dacryopinax</i> sp.				14	50	1
		<i>Fomitopsis pinicola</i>				16	65	0
		<i>Gloeophyllum trabeum</i>				13	61	1
		<i>Hydnomyces pinastri</i>				21	86	16
		<i>Serpula lacrymans</i>				12	52	8
	Ascomycota	<i>Agaricus bisporus</i>	Saproth	Soil, litter or other Saproth		47	65	13
		<i>Amanita thiersii</i>				27	71	10
		<i>Coprinopsis cinerea</i>				41	96	44
		<i>Gymnopus luxurians</i>				66	118	32
		<i>Ustilago maydis</i>				Plant pathogen	6	22
	Chytridiomycota	<i>Melampsora larici-populina</i>	Biotroph	Ectomycorrhizal		31	71	0
		<i>Tuber melanosporum</i>				5	19	3
		<i>Oidiodendron maius</i>				Ericoid symbiont	31	119
		<i>Cryphonectria parasitica</i>				18	102	12
		<i>Stagonospora nodorum</i>				16	104	10
		<i>Aspergillus nidulans</i>	Saproth	Animal pathogen		9	83	6
		<i>Trichoderma reesei</i>				10	39	14
	NA	<i>Pichia stipitis</i>	Saproth	Soil, litter or other Saproth		2	10	0
		<i>Phycomyces blakesleeanaus</i>				2	18	1
	Chytridiomycota	<i>Batrachochytrium dendrobatidis</i>	Biotroph	Animal pathogen		3	5	0

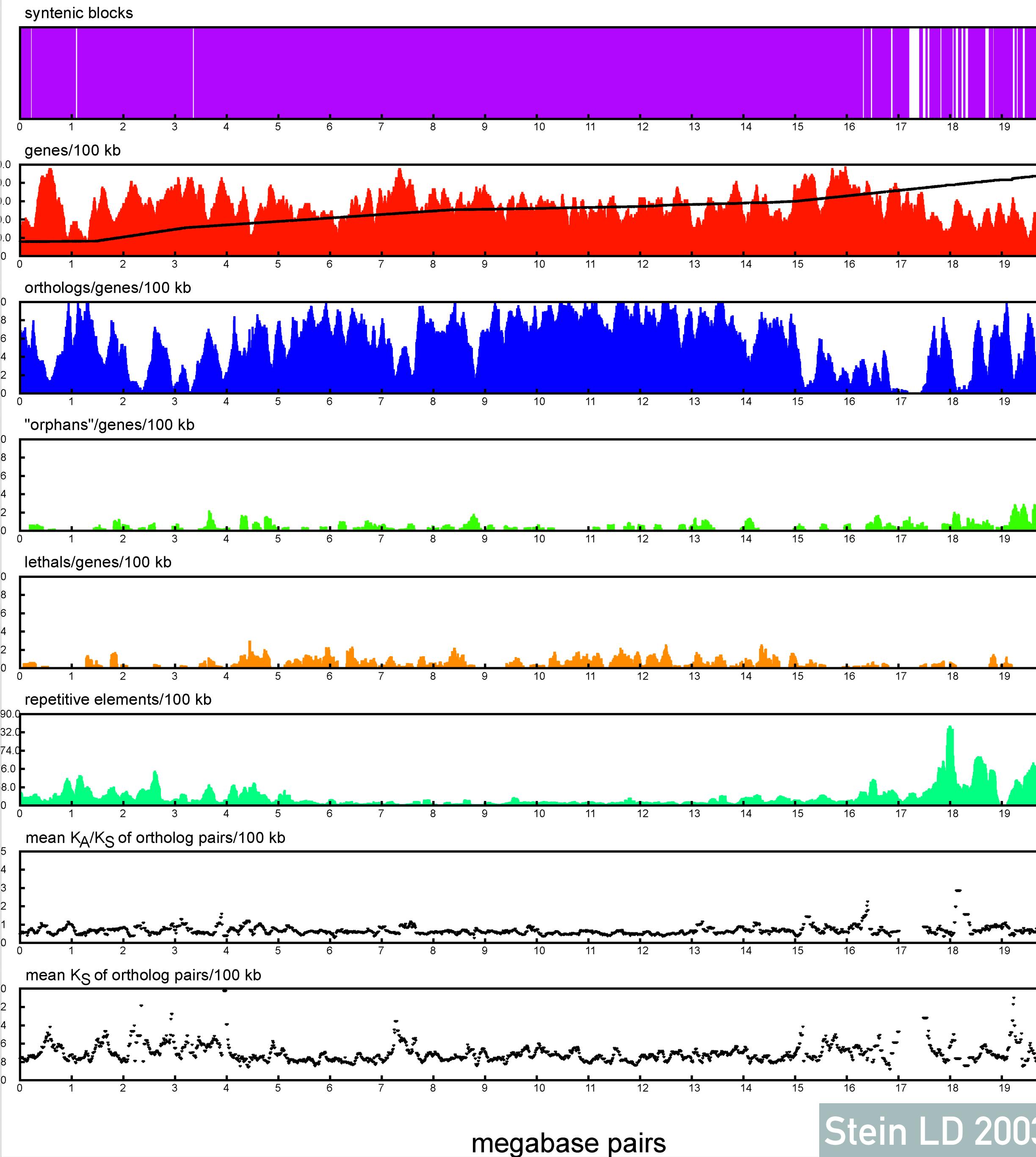


Distribution of secreted plant cell wall degrading enzyme (PCWDE) across collection fungi to test how Tuber group interacts



Black frames highlight Pezizomycetes taxa, whereas black arrows indicate Tuberaceae taxa and red dotted lines highlight PL1, PL4 and CE12 families. The ecology of each species is indicated at the right of the species abbreviations.

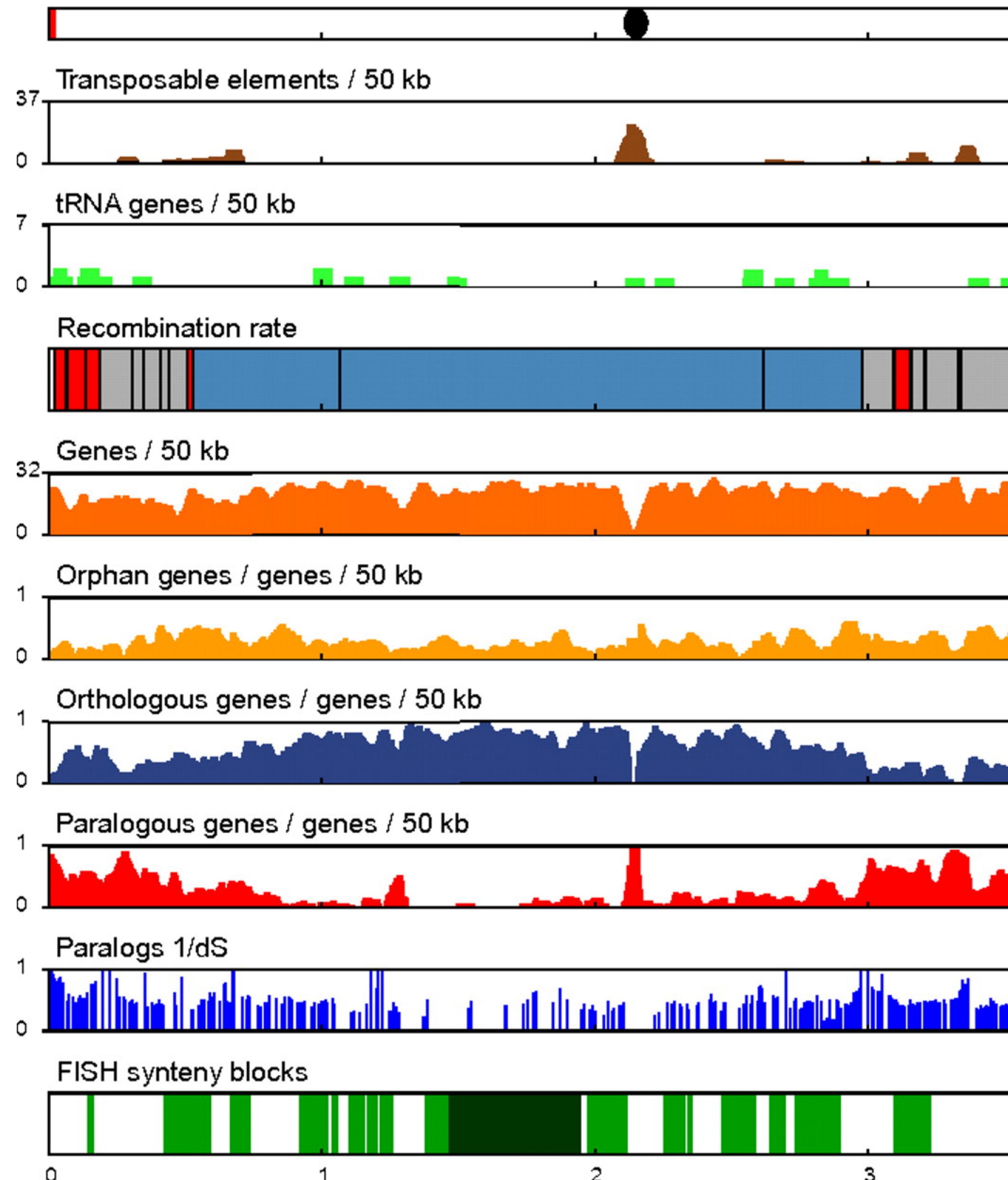
Chromosome V



GENOME CONTENT WITH TRACKS

- *C. briggsae* gene content and features
- Show information about gene content, which genes are shared/unique
- Inverse relationship between Transposons and called genes.

Telomere and Centromere location

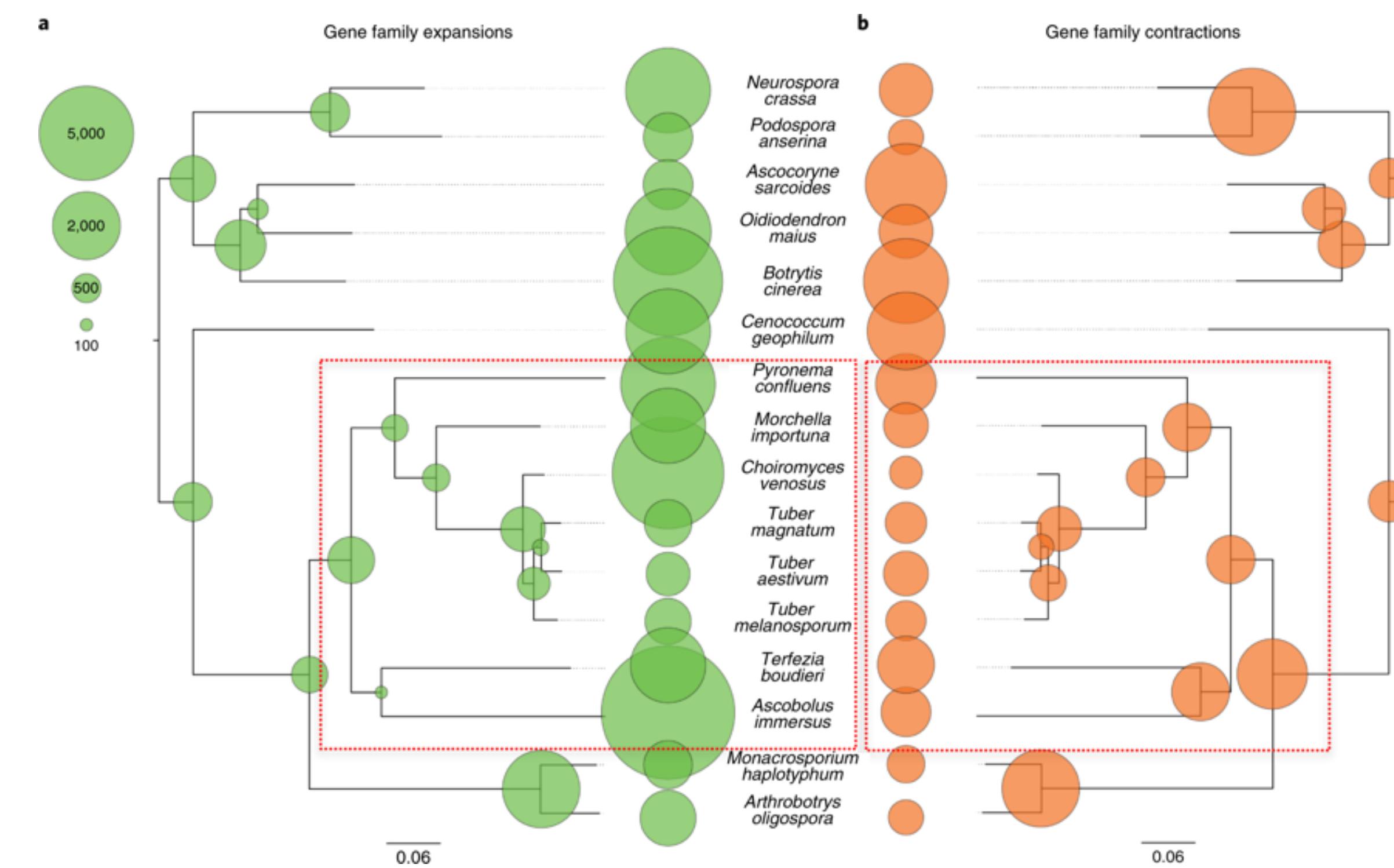
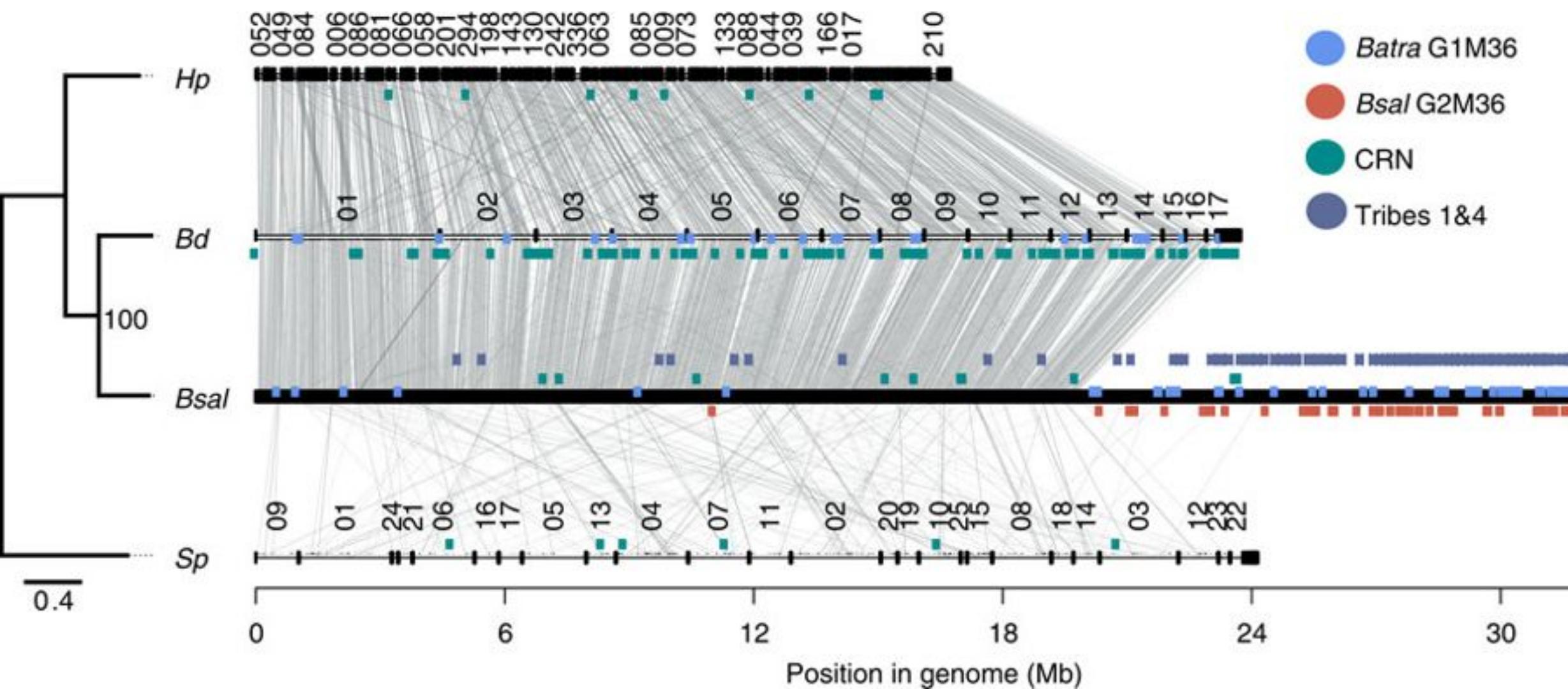


SIMILAR CHROMOSOME CONTENT PLOT

- Comparison of chromosome II for Mushroom *Coprinopsis*.
- Demonstrating other ways to compare content
- Notice the regions of high synteny and having orthologous genes are located in center
- While gene duplications (paralogs) are enriched at telomere

METHODS: COMPARATIVE GENOMICS

- Identifying Shared and Unique Genome content
- Synteny
- Whole genome alignment
- Gene family
 - Gene clustering - examine orthology and paralogy
 - Functional content comparison



GENOME ALIGNMENT AND SYNTENY INFERENCE

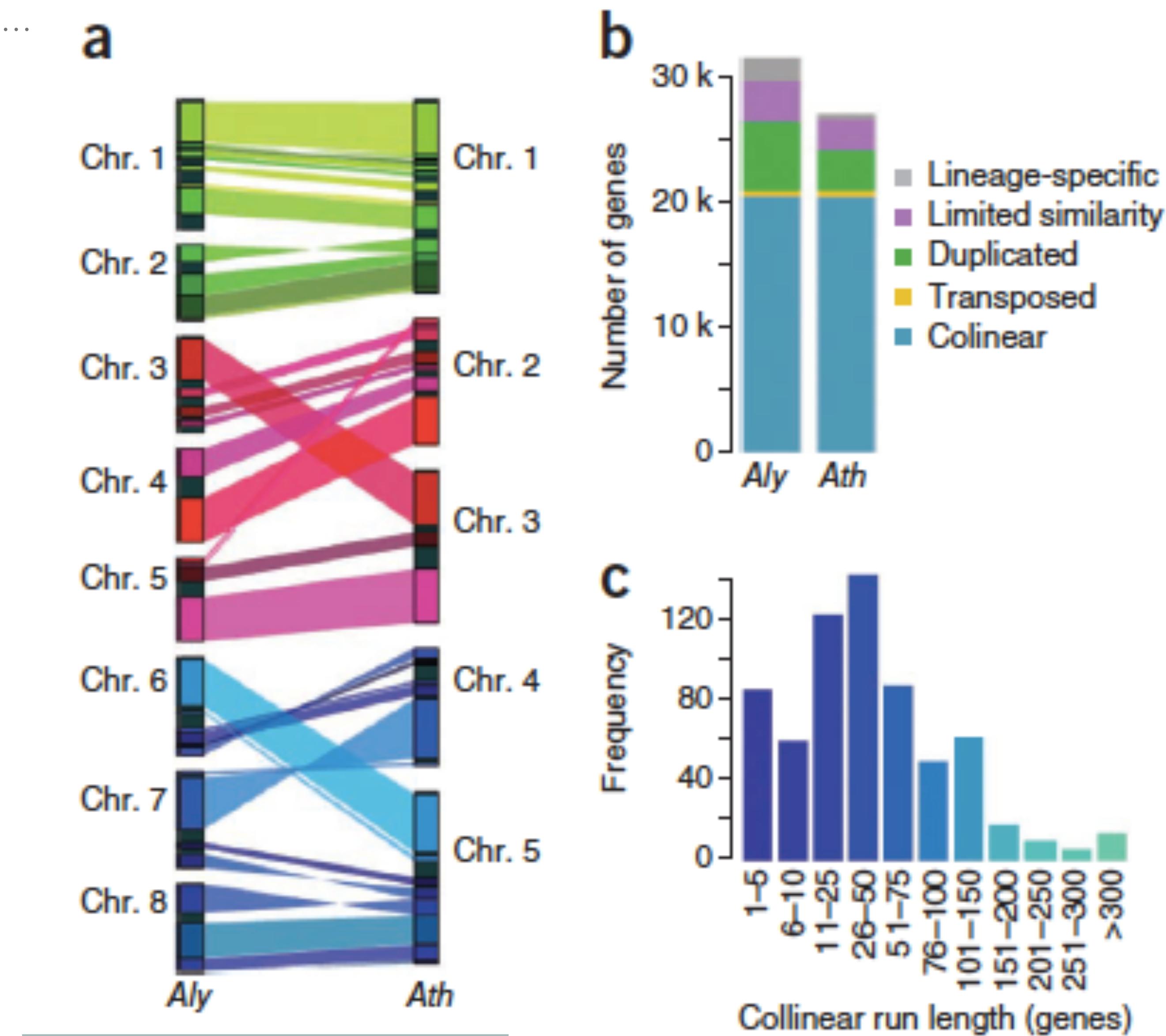
Goals

- Identify homologous segments between regions of genome
- Reconstruct evolutionary history for every nucleotide in the genome

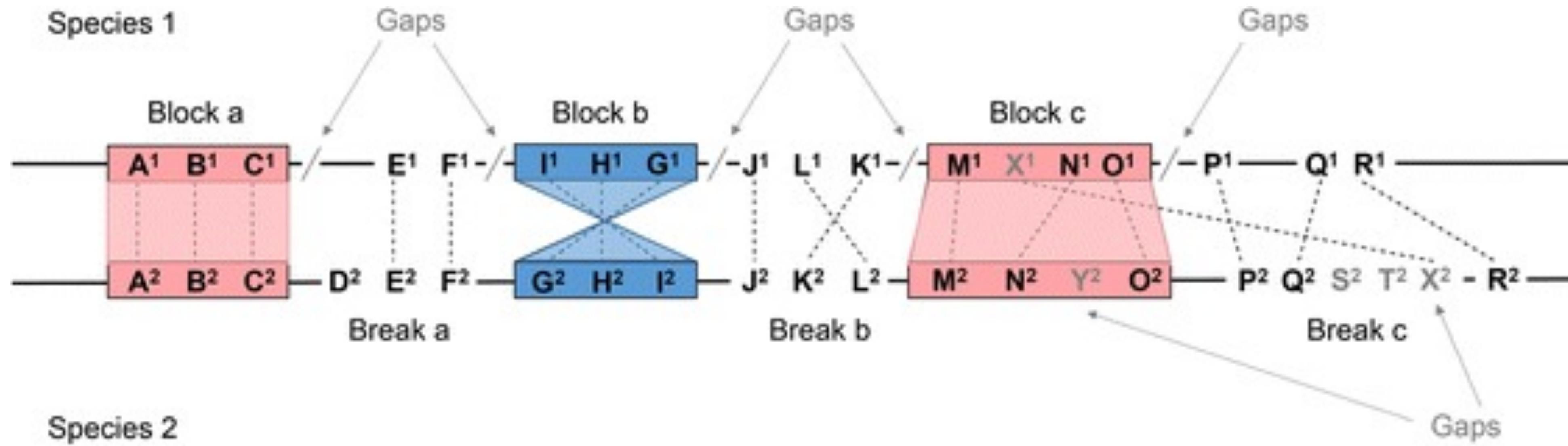
Challenges: Genomes are not necessarily (rarely!) co-linear between species

Identifying Insertion / Deletion / Rearrangements

Not feasible to just run a giant multiple alignment for all chromosomes



SYNTENY



Synteny was originally defined as two more pairs of homologous genes on same chromosomal segment

The stricter definition require collinearity of the genes and orientation

Requires detection of similar regions and chaining these into “runs” of genes

SYNTENY TO EXAMINE GENOME ORGANIZATION AND CONSTRAINT

C. briggsae 100kb
scaffolded with ALLMAPS
(L90 = 1,063 → L90 = 6)

C. elegans
(L90 = 6)

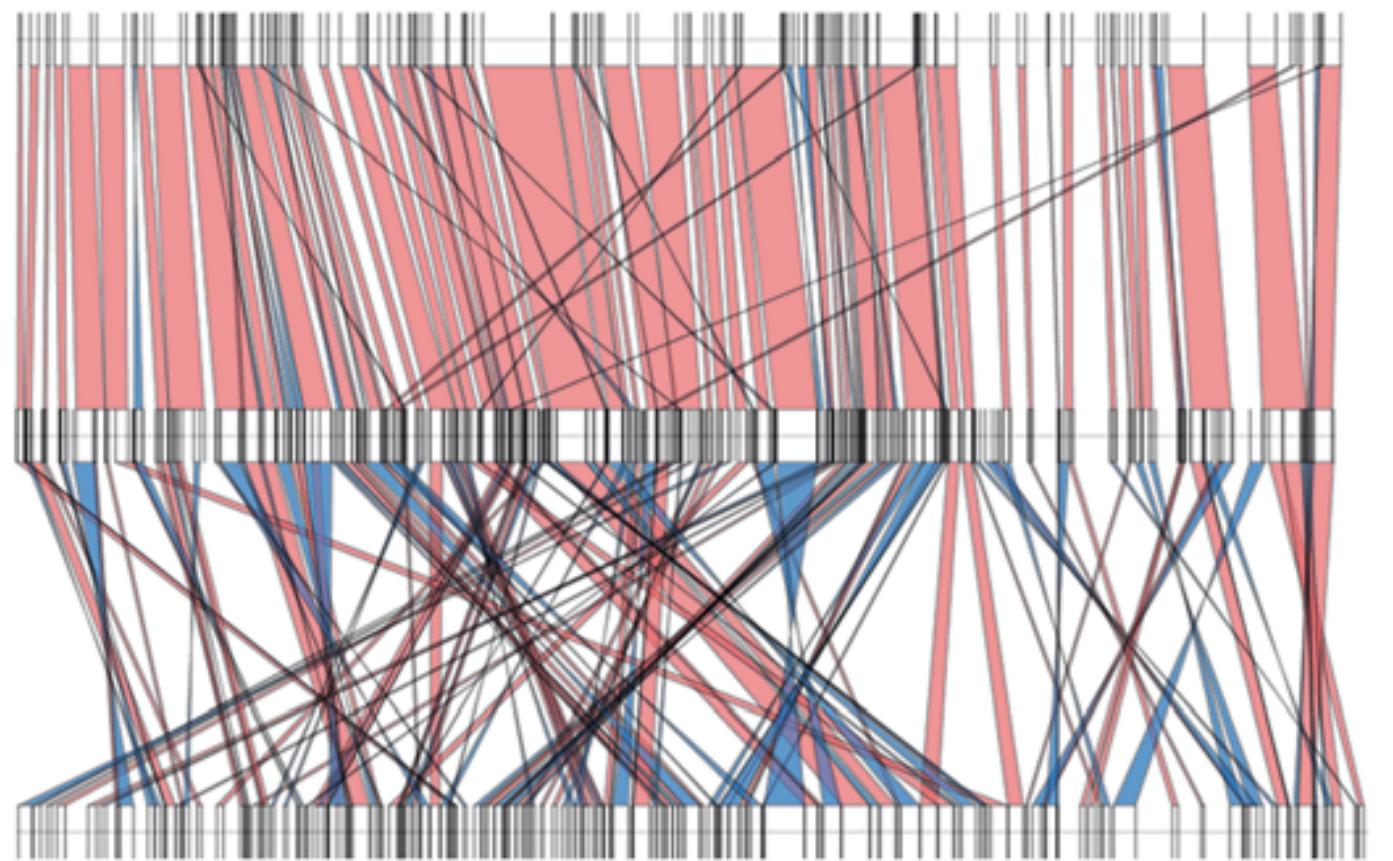
C. briggsae
(L90 = 6)

C. briggsae 100kb
scaffolded with ALLMAPS
(L90 = 1,063 → L90 = 6)

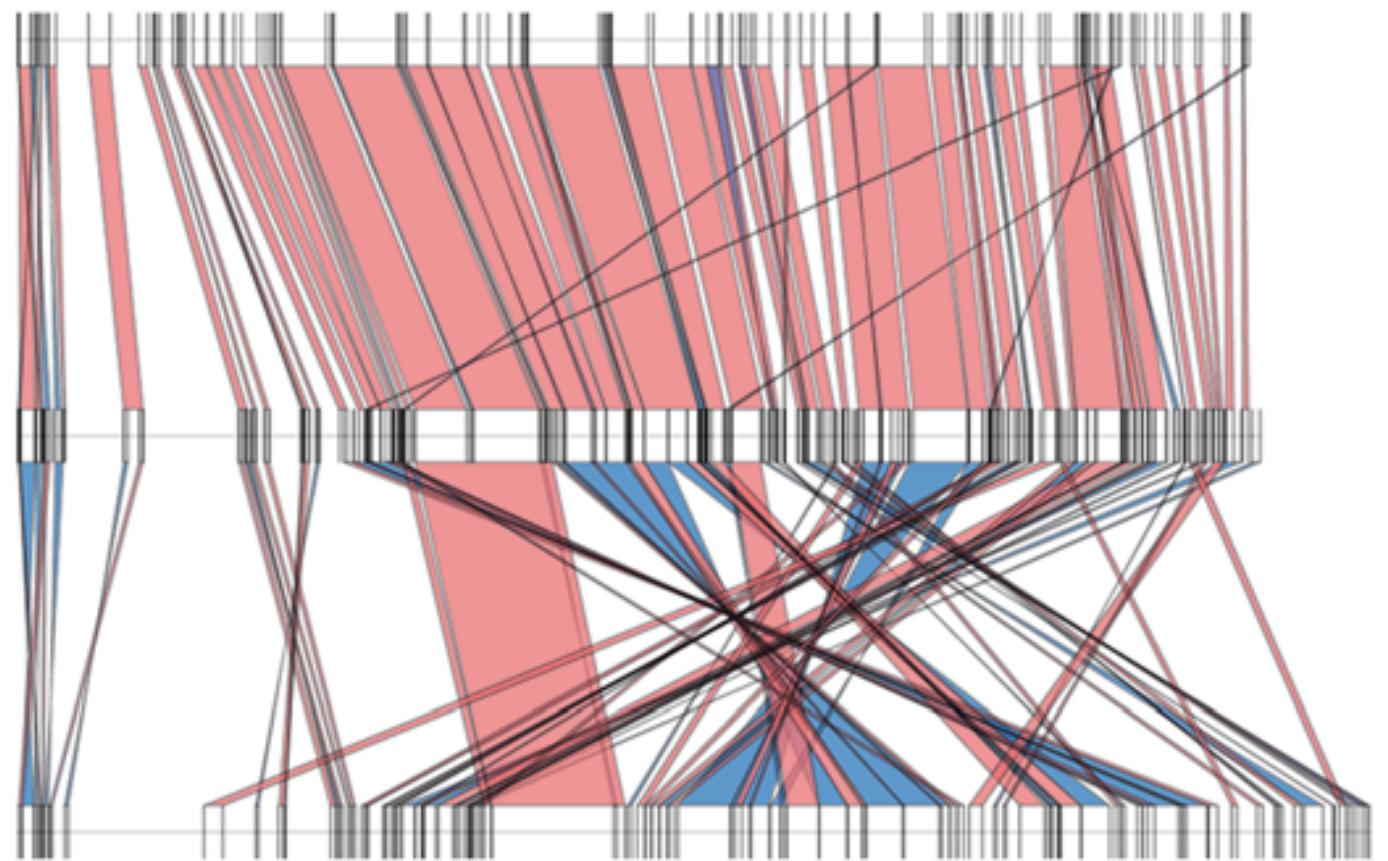
C. elegans
(L90 = 6)

C. briggsae
(L90 = 6)

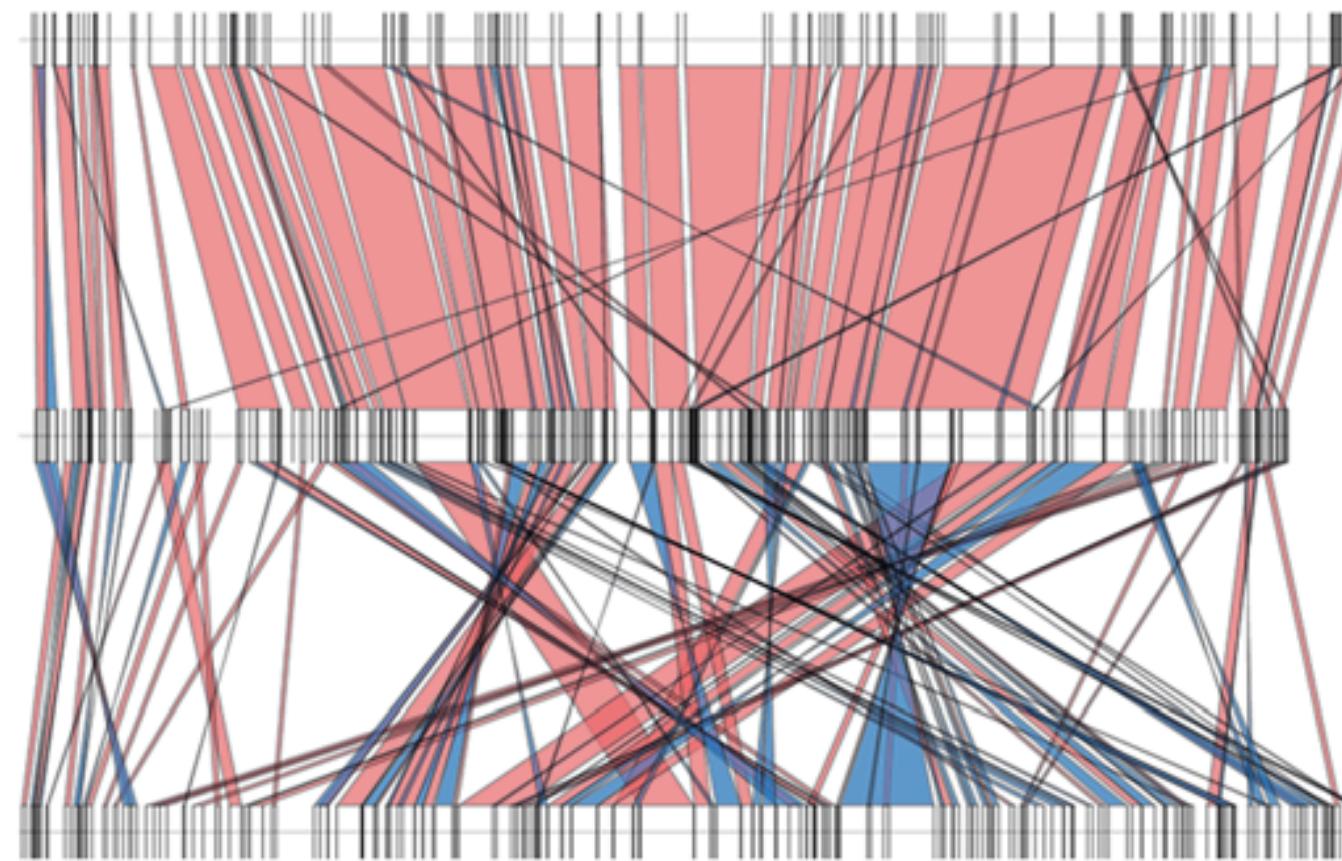
Chr. I



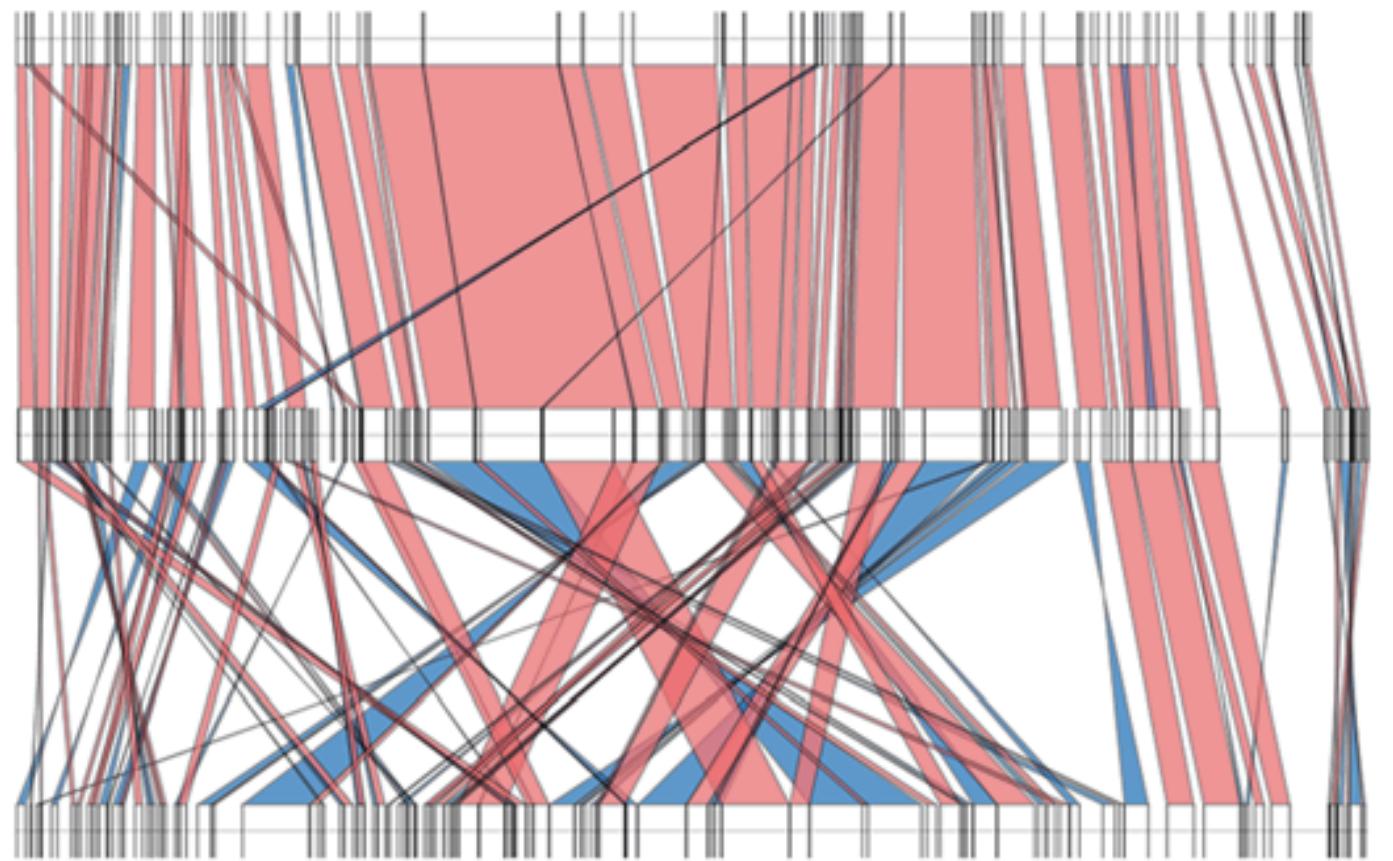
Chr. II



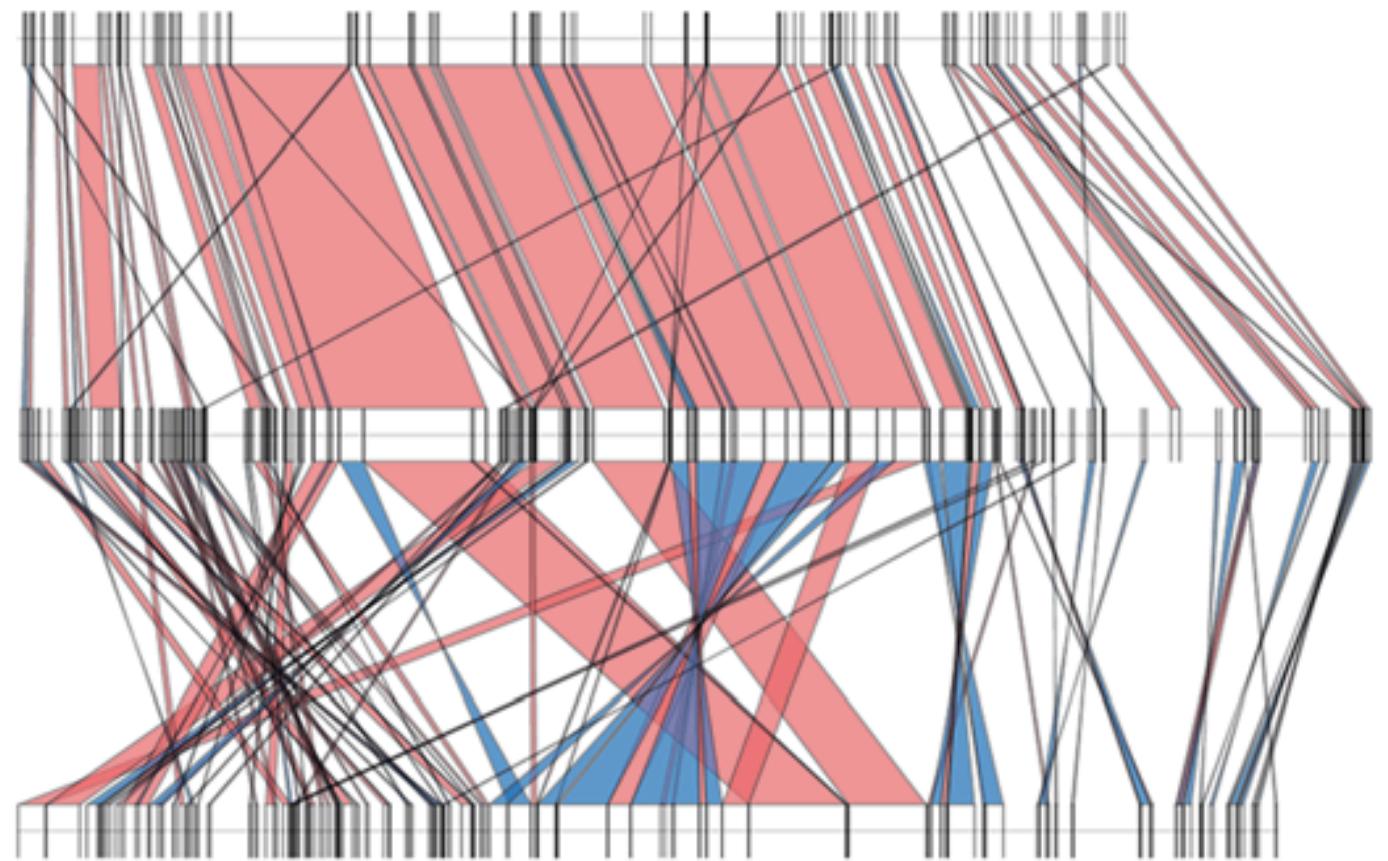
Chr. III



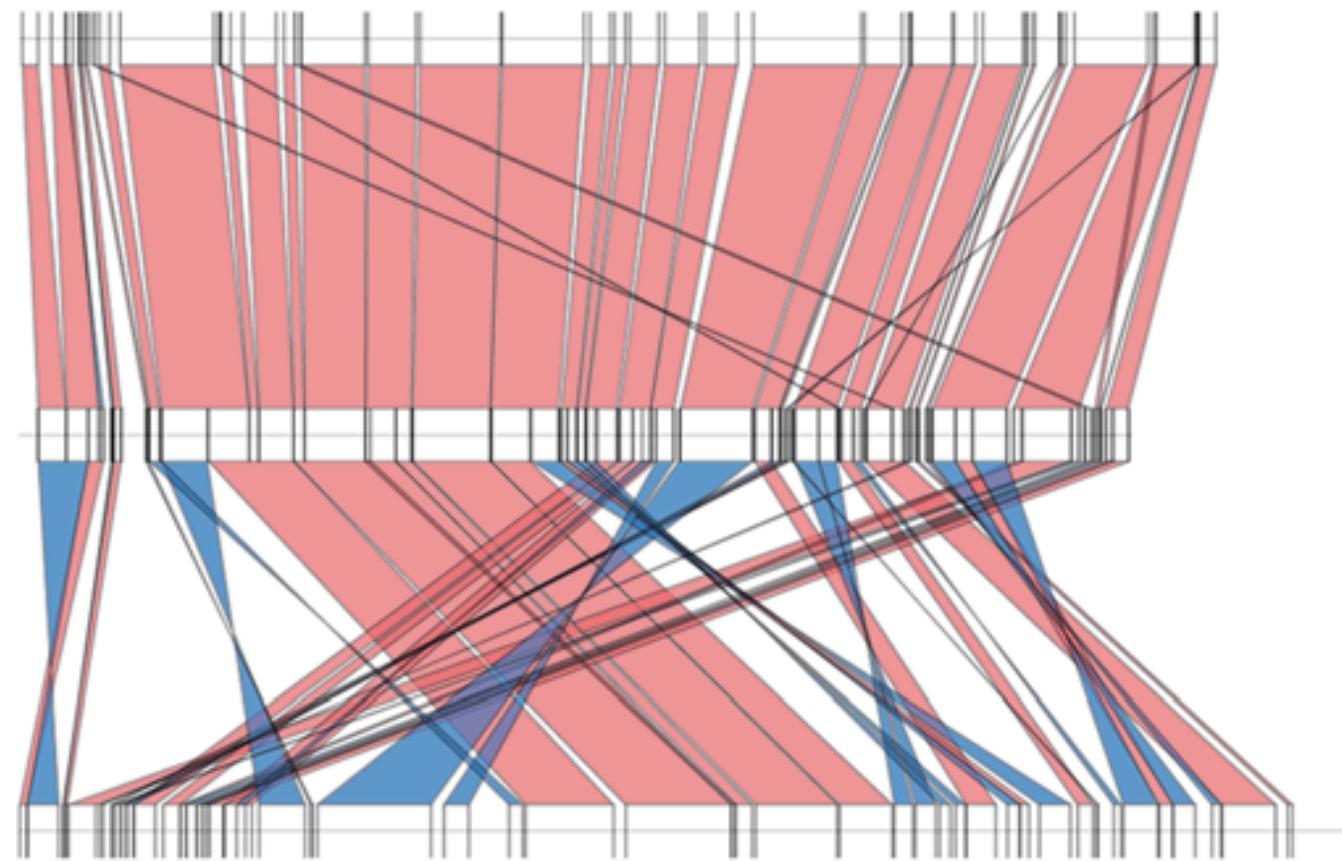
Chr. IV

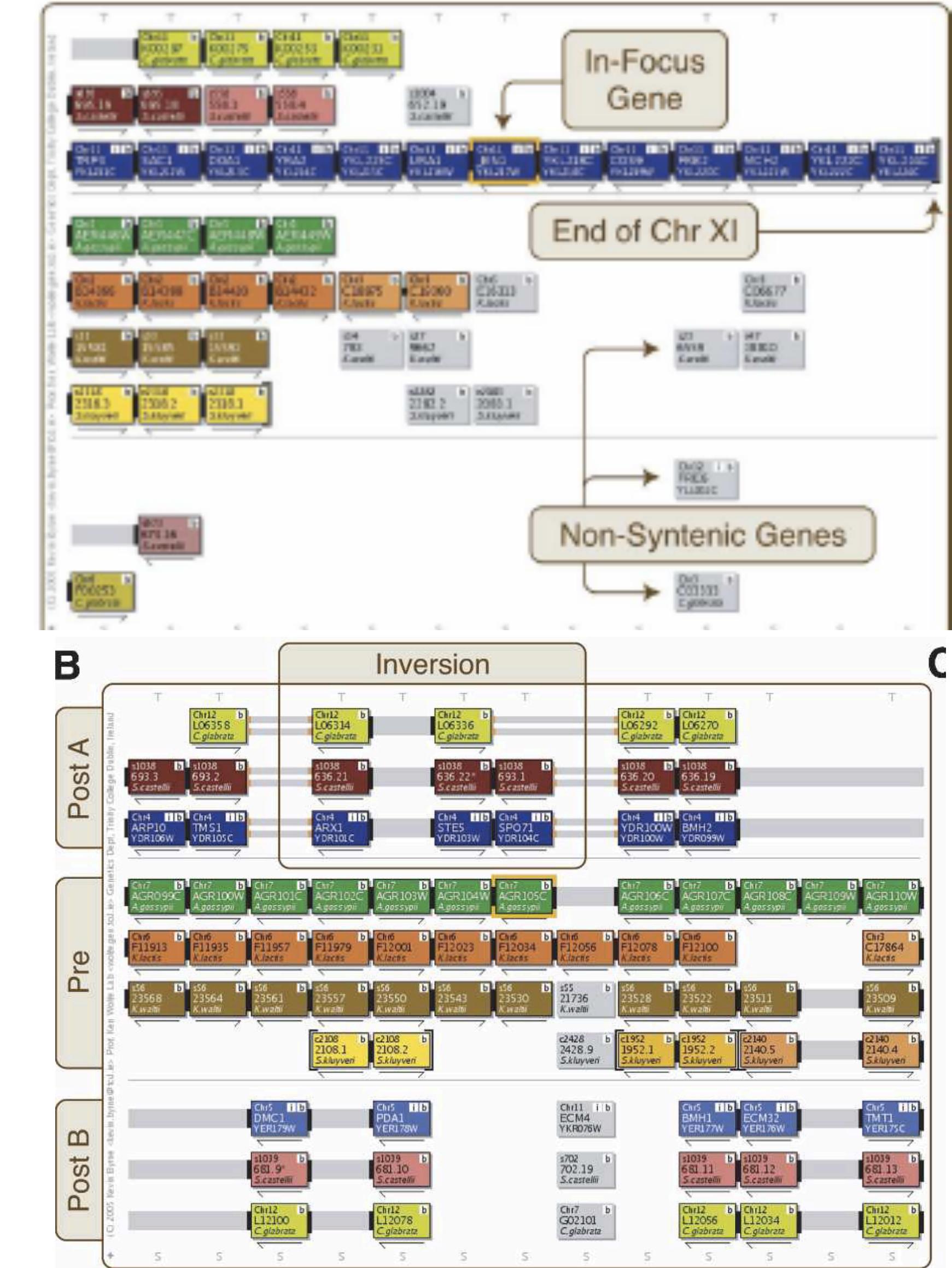
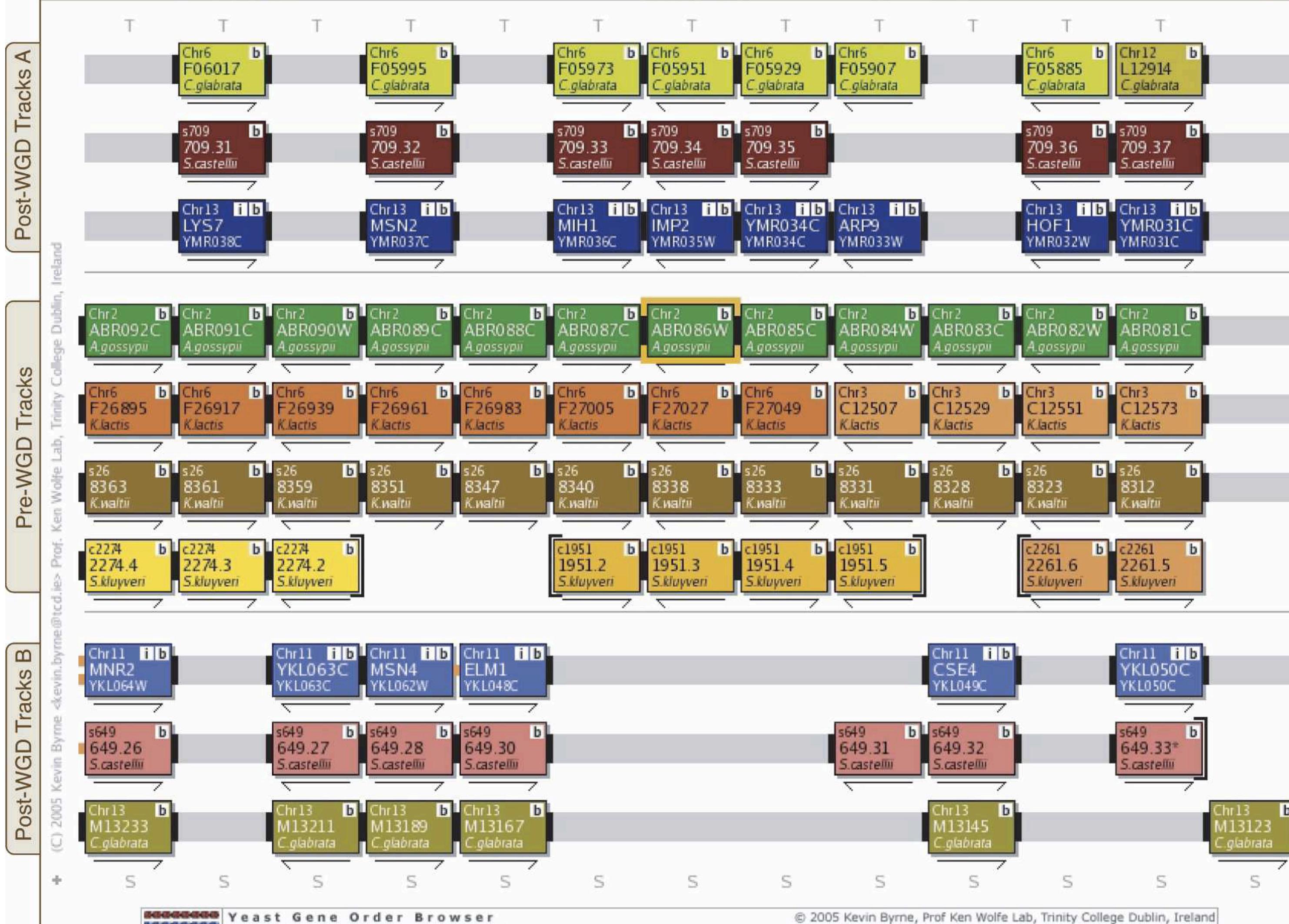


Chr. V



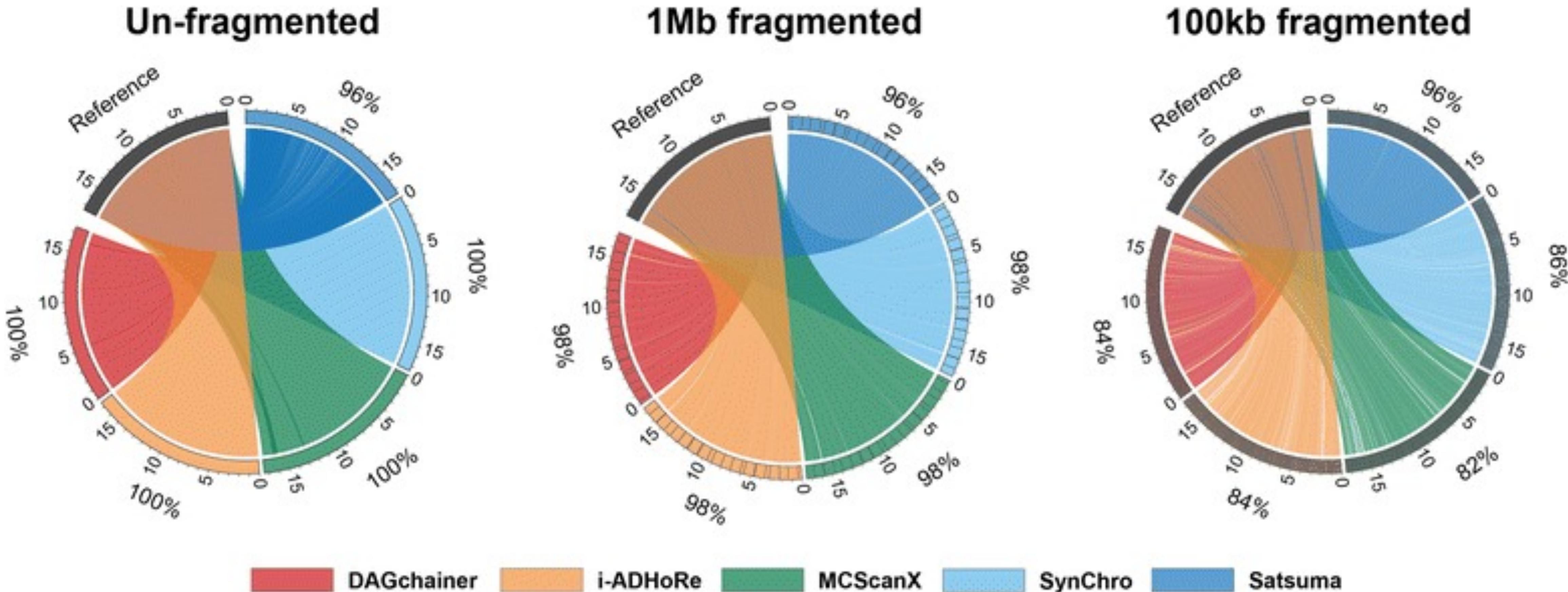
Chr. X





*Yeast Gene Order Browser: Automated and curated synteny Gold standard in synteny and orthology comparisons in *Saccharomycotina* yeasts - <http://ygob.ucd.ie/>*

RECOVERY OF SYNTENIC REGIONS BY DIFFERENT TOOLS ON SIMULATED FRAGMENT C.ELEGANS CHROMOSOME



TOOLS FOR SYNTENY COMPUTATION

Gene or gene anchor based

i-ADHoRe - <http://bioinformatics.psb.ugent.be/beg/tools/i-adhore30>

Multi-way species compare

“highly sensitive software tool to detect degenerated homology relations within and between different genomes.”

DAGchainer- <http://dagchainer.sourceforge.net/>

Multi-way species compare

“computes chains of syntenic genes found within complete genome sequence”

Mercator- <https://www.biostat.wisc.edu/~cdewey/mercator/>

Multi-way species compare

constructs orthology maps from protein anchored maps followed by refinement from multiple alignment of DNA

TOOLS FOR SYNTENY COMPUTATION AND ALIGNMENT

Whole genome alignment based

Satsuma2 - <https://github.com/bioinfologics/satsuma2> Pairwise

“FastFourierTransform cross-correlation based synteny aligner, to reliably align large and complex DNA sequences providing maximum sensitivity (to find all there is to find), specificity (to only find real homology) and speed”

MUMMER - <https://github.com/mummer4/mummer> (v4) and <http://mummer.sourceforge.net/> (v3) Pairwise

“a system for rapidly aligning entire genomes, whether in complete or draft form”

LAST - <http://last.cbrc.jp/> Pairwise

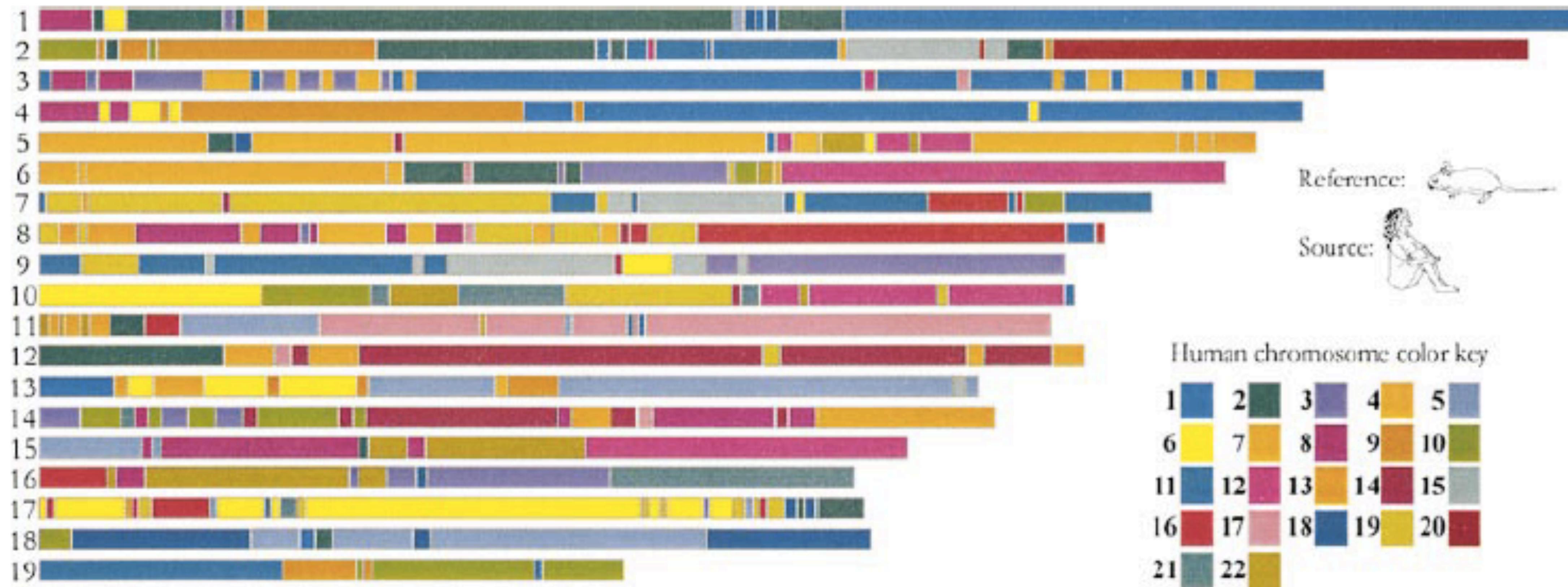
“LAST finds similar regions between sequences, and aligns them.”

Minimap2 - <https://lh3.github.io/minimap2/> Pairwise

A versatile pairwise aligner for genomic and spliced nucleotide sequences

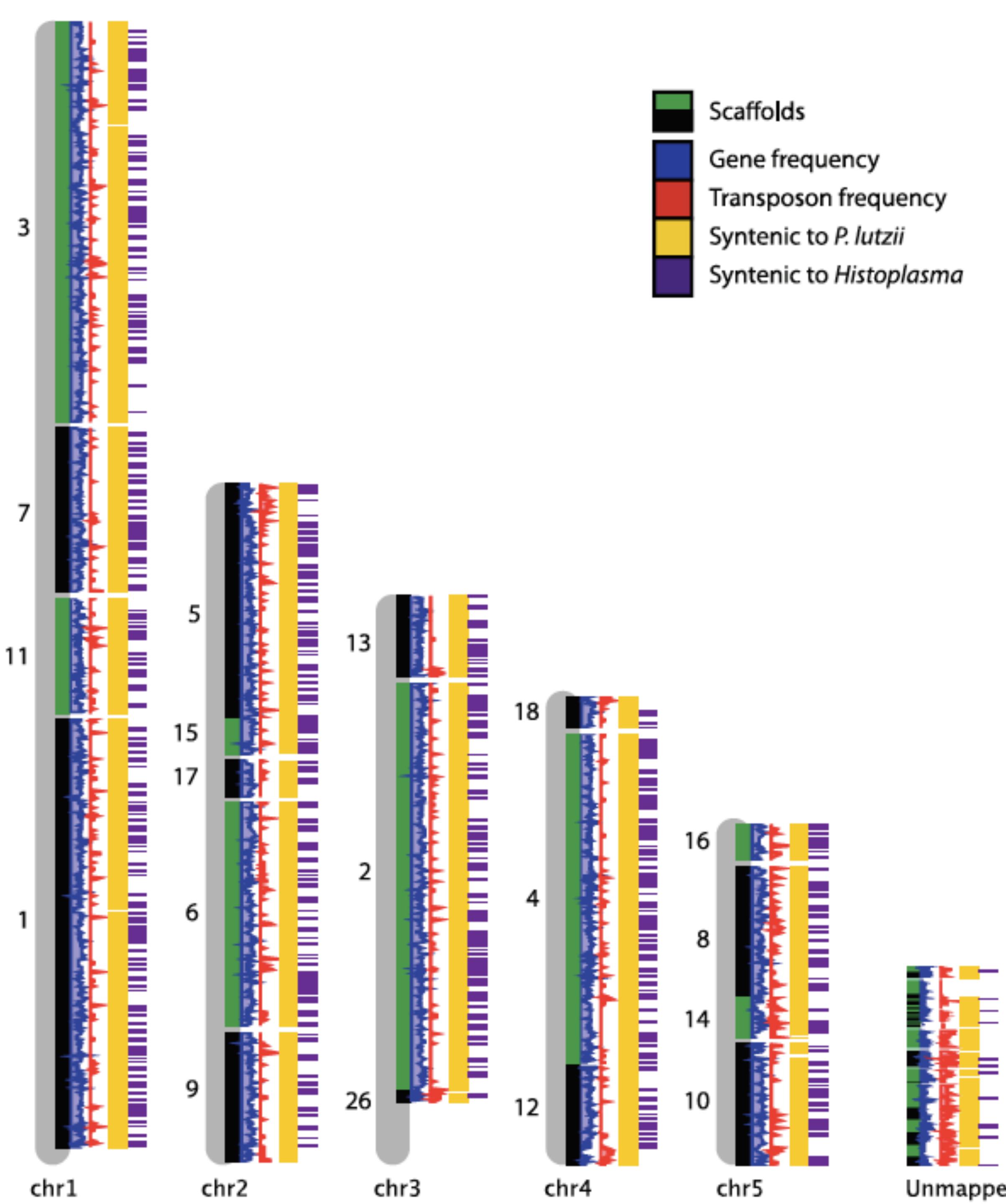
SYNTENY AND HOMOLOGY: MOUSE HUMAN SHARED GENE CONTENT & SYNTENY

Chromosomal Segments Between Species

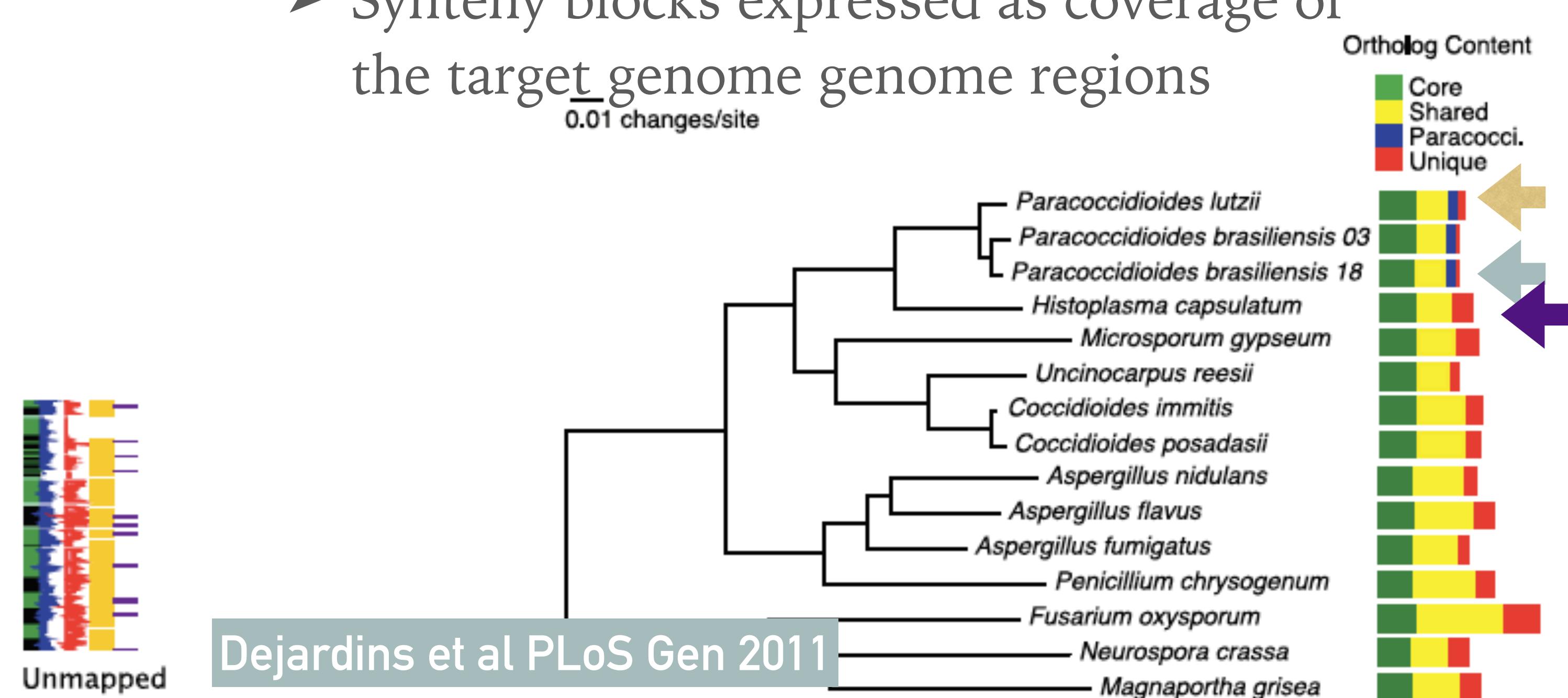


SUMMARIZING SHARED GENOME CONTENT

- 3 fungi are compared here: Pb18 supercontigs (grey) to the five chromosomes from the optical map (chr1-5).
- Summary statistics about gene and transposon density
- Synteny blocks expressed as coverage of the target genome genome regions



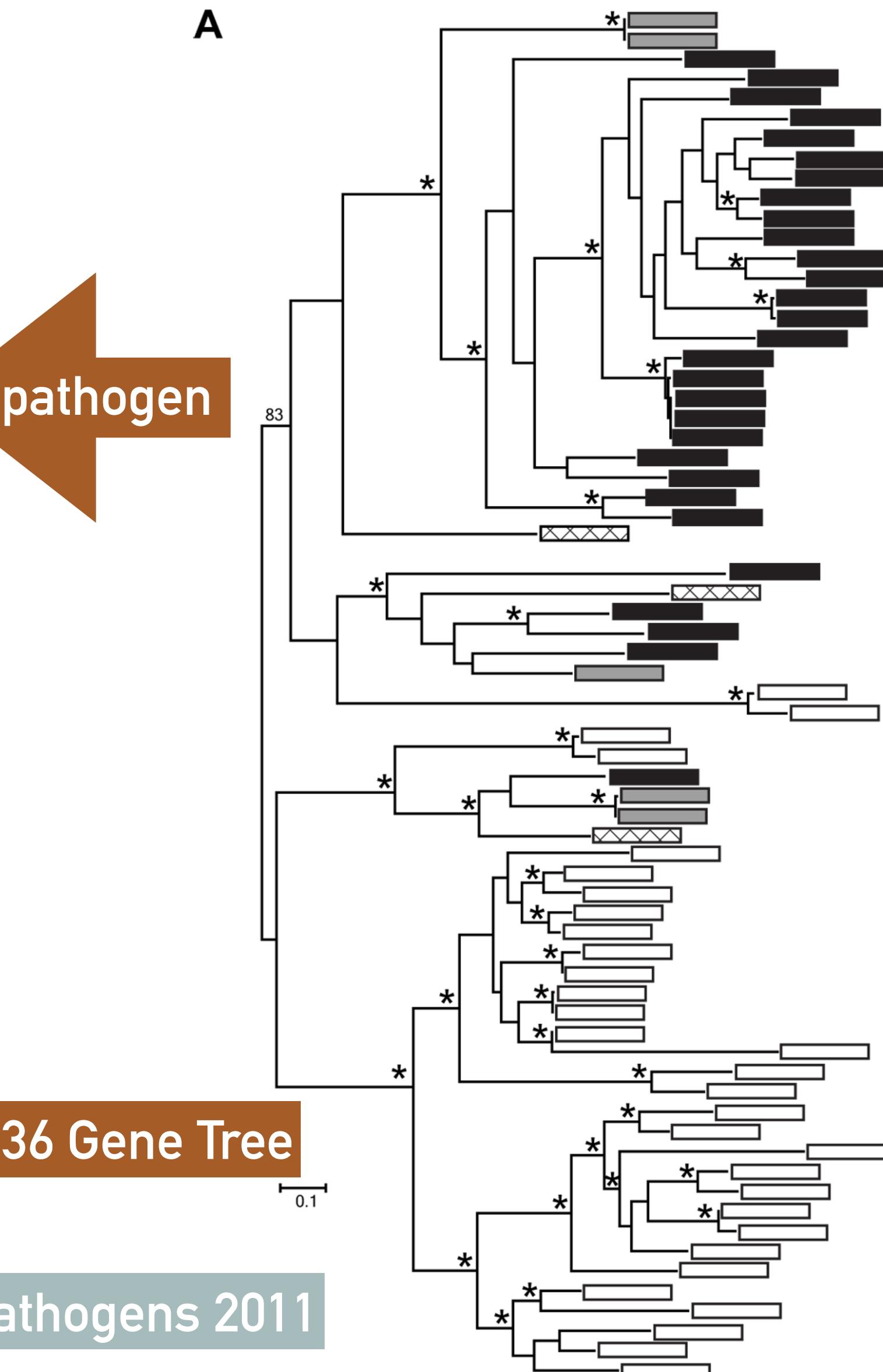
Dejardins et al PLoS Gen 2011



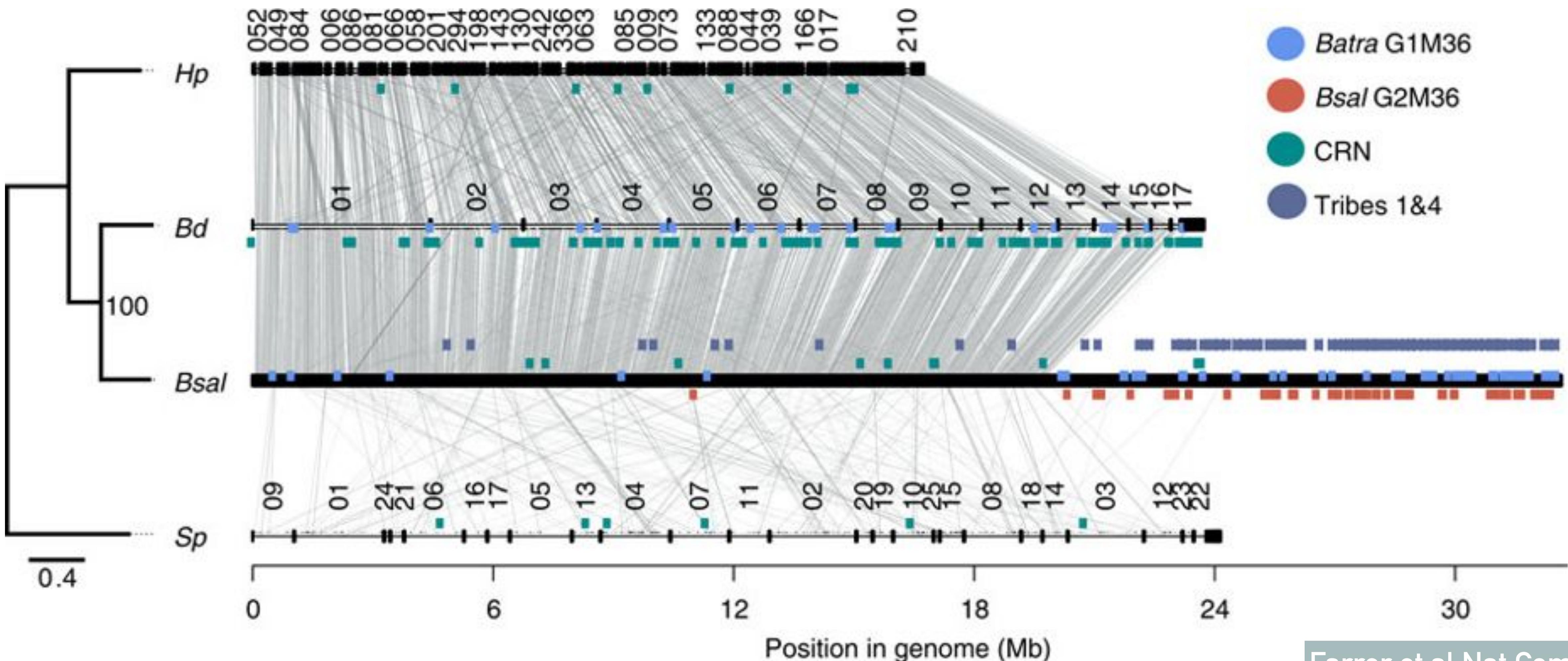
EXAMINATION OF AN EXPANDED GENE FAMILY WITH FURTHER FOLLOWUP OF GENE TREES

	M36	S41	Asp	CRN
<i>Allomyces macrogynus</i>	31	0	6	0
<i>Spizellomyces punctatus</i>	3	3	10	0
<i>Homolaphlyctis polyrhiza</i>	5	3	22	0
<i>Batrachochytrium dendrobatidis</i>	38	32	99	62

Examination of all Pfam and metalloprotease domains among the available chytrids found only a few which were expanded in the pathogenic lineage



SYNTENY, NOVEL EXPANSION OF M36 FAMILY IN *B. SALAMANDRIVORANS* - AMPHIBIAN PATHOGEN

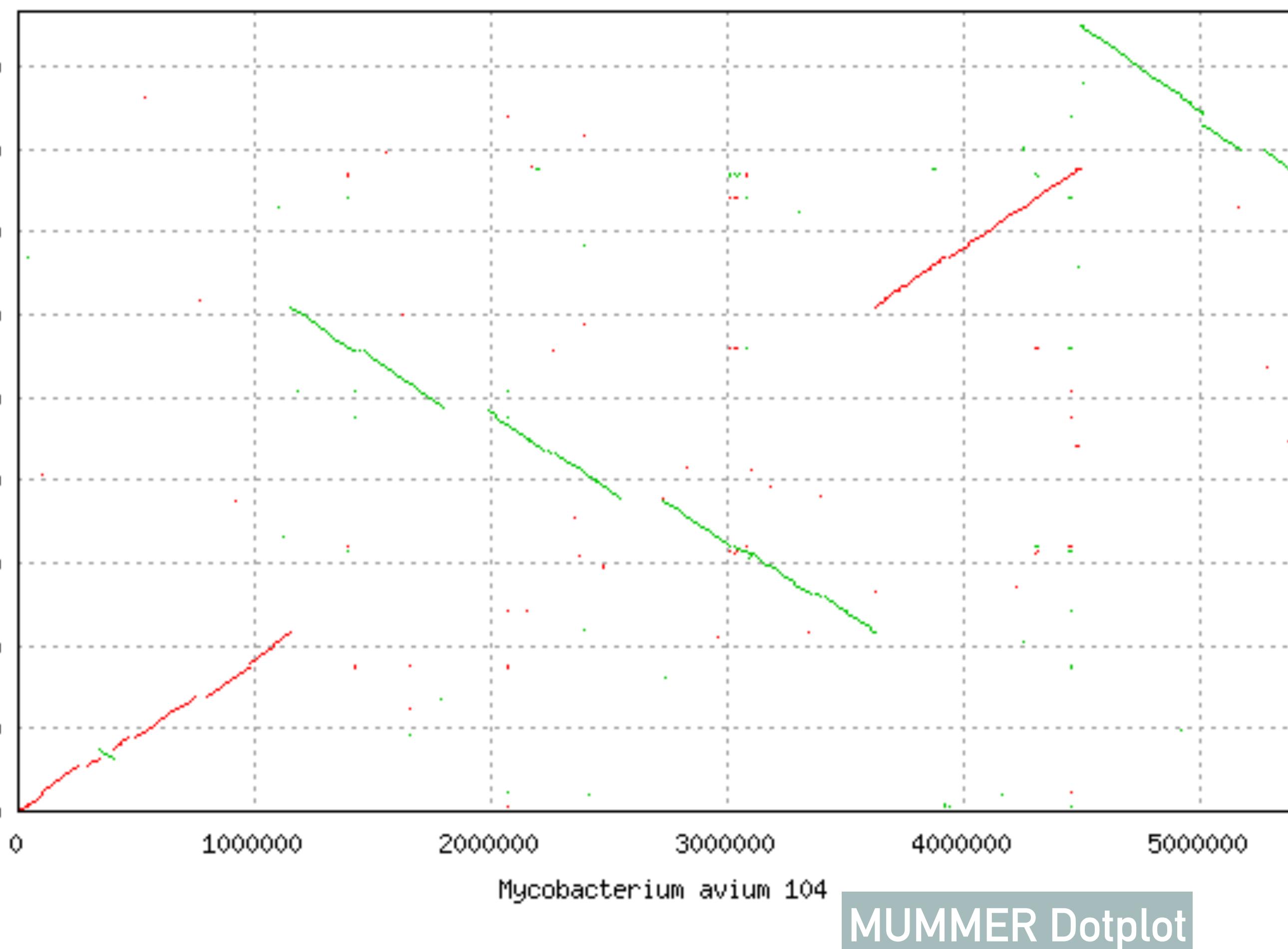


WHOLE GENOME ALIGNMENT

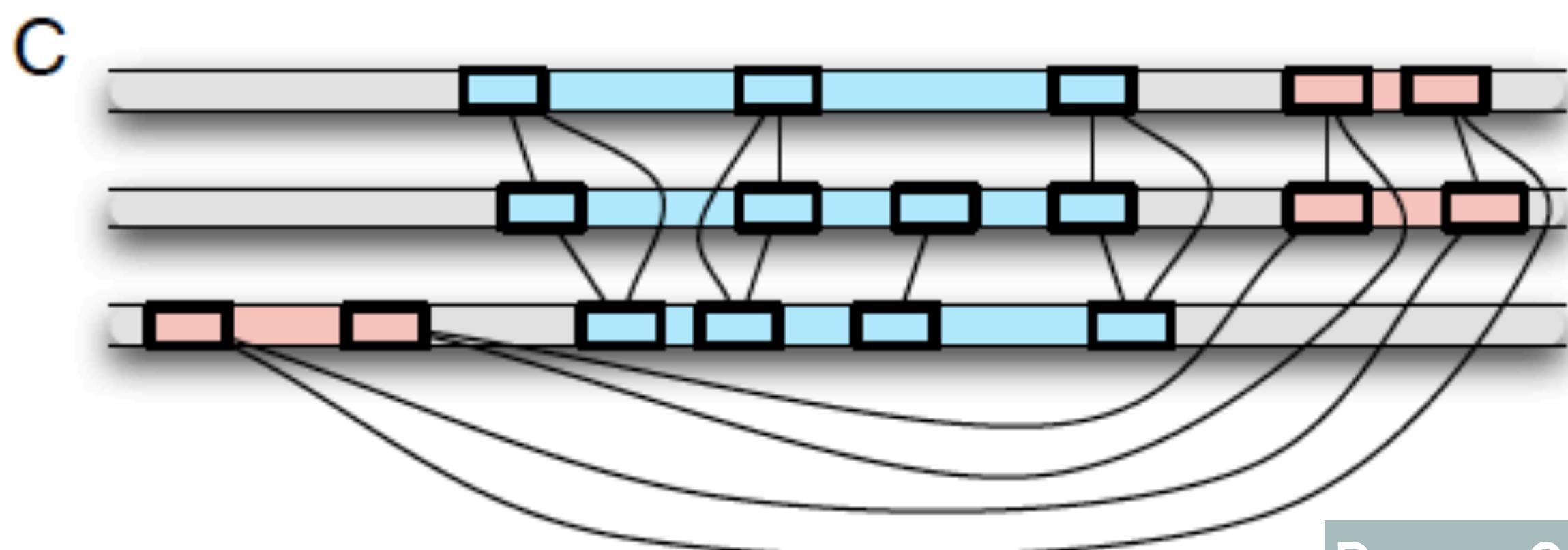
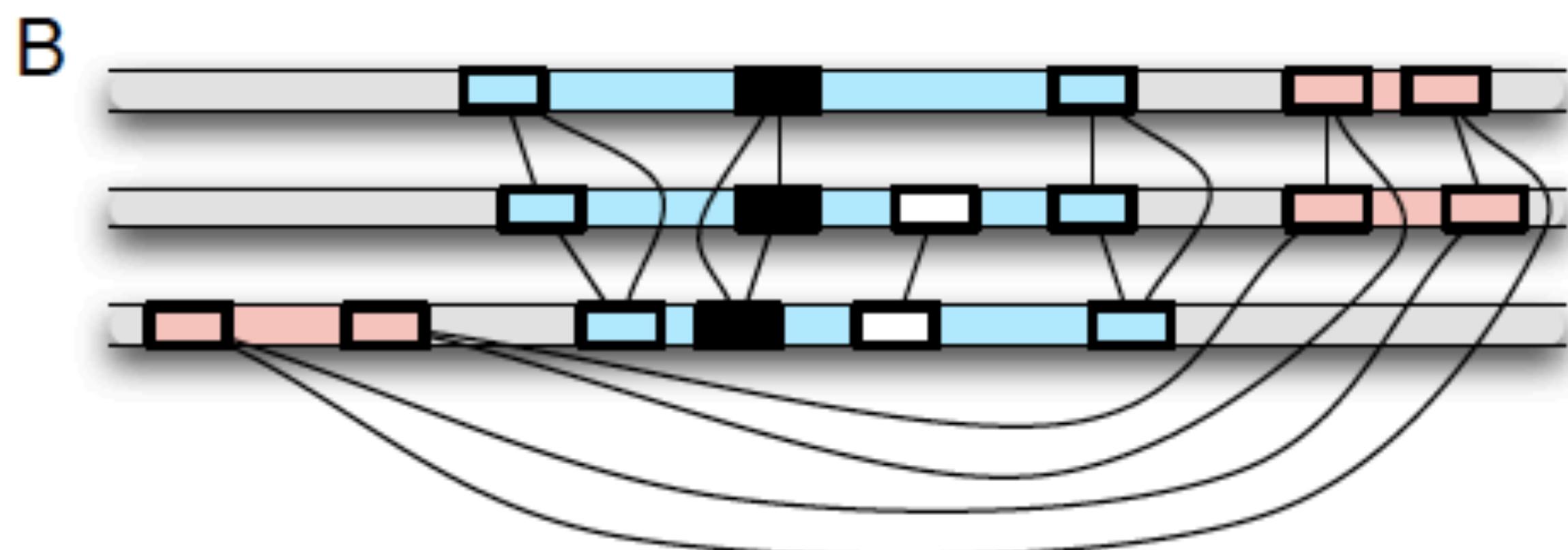
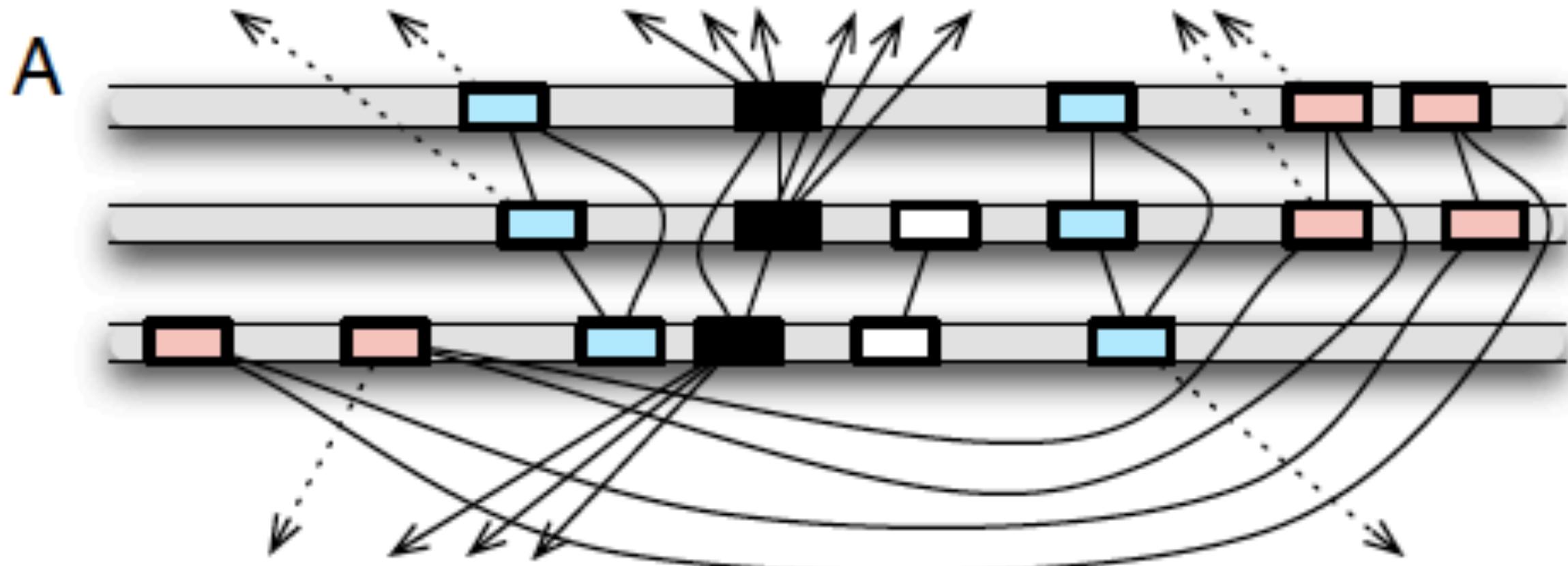
WORD MATCHING & ALIGNMENTS

NUCmer Plot: *Mycobacterium avium* 104 vs. *Mycobacterium avium* K-10

(In a Mummer dot-plot, the reference genome is along the x-axis, while the query genome is on the y-axis. Wherever the two sequences match, a colored dot is plotted. The forward matches are displayed in red, while the reverse matches are displayed in green)



- Word-based, fixed or variable size matching
- **MUMMER, LAST, LASTZ, BLASTZ, YASS**
 - MUMMER: Fast lookup with suffix tables by building an index query and target sequences and apply fast suffix-table matching “ultra fast”
 - LAST finds initial matches based on their multiplicity, instead of using a fixed length (e.g. BLAST uses 11-mers). To find these variable-length matches, it also uses a suffix array
 - Generally implemented as a pairwise analysis but can be hierarchically chained to support multiple taxa
 - Annotation free - use sequence matches to define regions to anchor alignments and extend match region



ANCHORING AND ALIGNING - MERCATOR

- Use other features to draw connections between
- Exons - protein coding regions in particular - can be matched between species
- Extend blocks based on shared gene order
- Targeted sequence alignment in the regions that are considered matching based on these features
- Mercator (Dewey and Pachter)

WHOLE GENOME ALIGNMENT

Examine evolutionary history of every base in the genome

Questions that can be answered

Rates of evolution across the genome - using a sliding windows or evaluating each gene.
Tests for directional selection

Individual gene trees from these alignments using subsampled windows

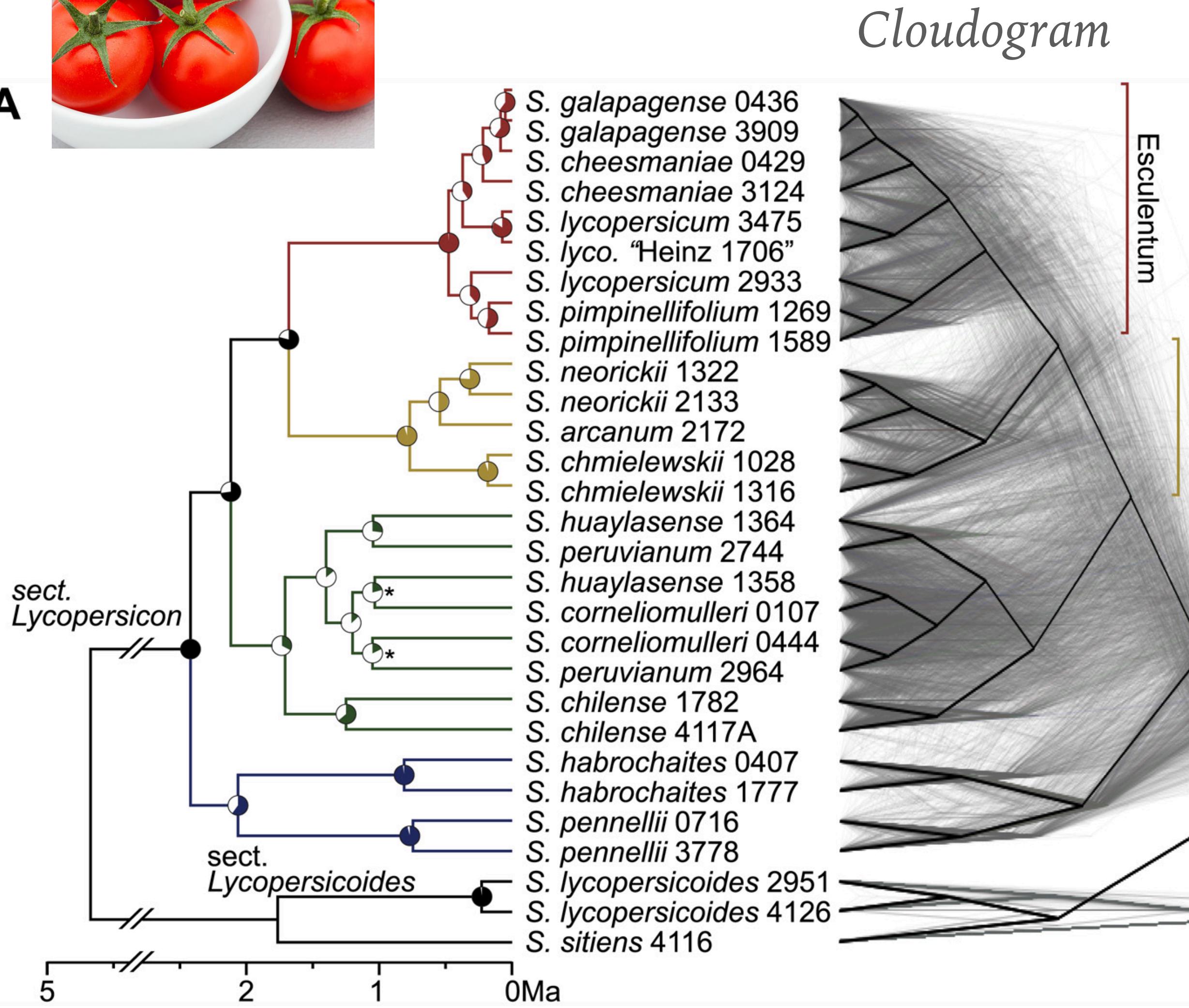
Identify unique regions and shared insertion/deletions

Testing for constraint / co-evolution

EXAMINING SPECIES TREE AND GENE/ WINDOW PHYLOGENY FROM WGA

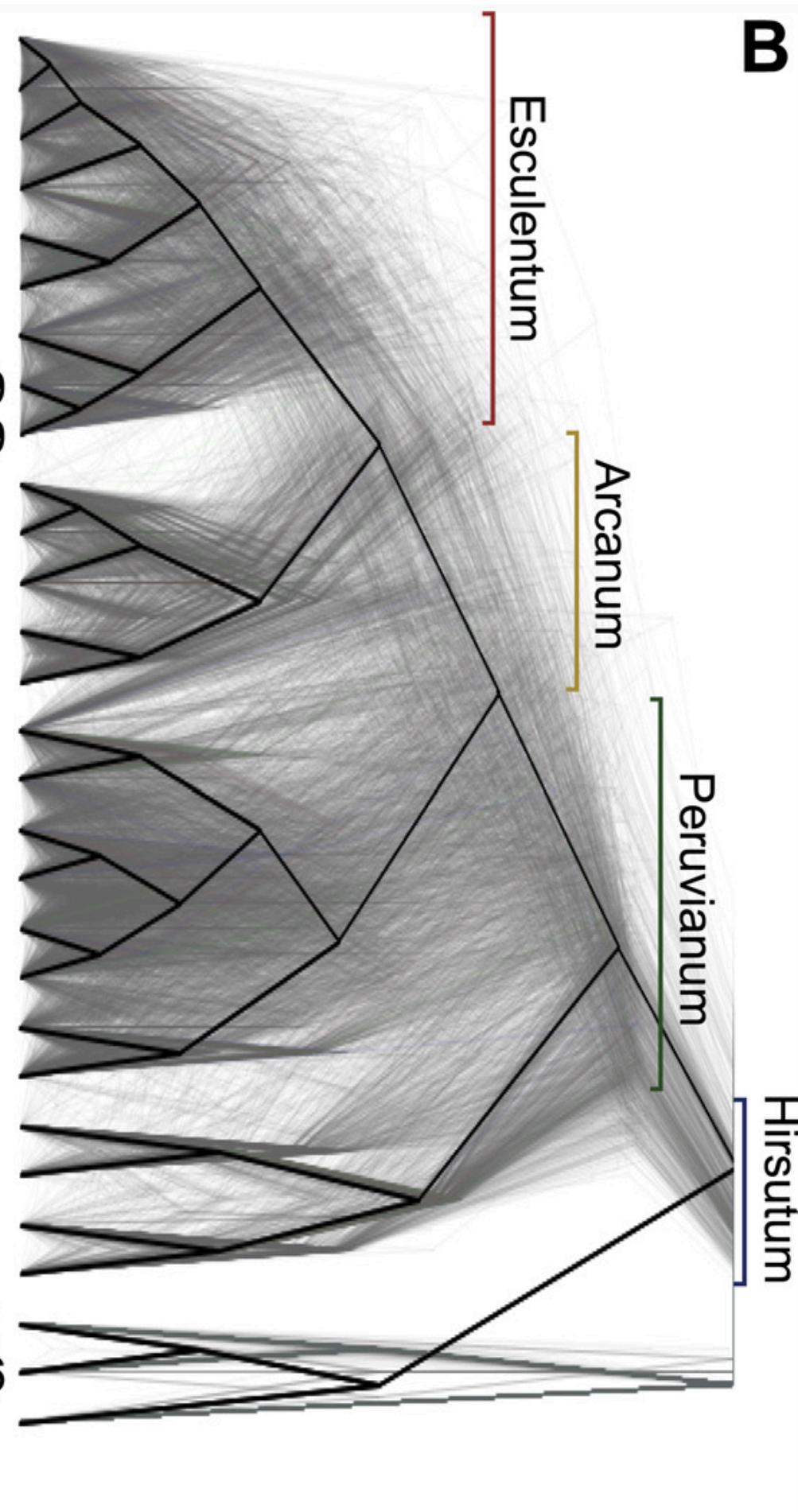


A



Solanum sect. Lycopersicon)

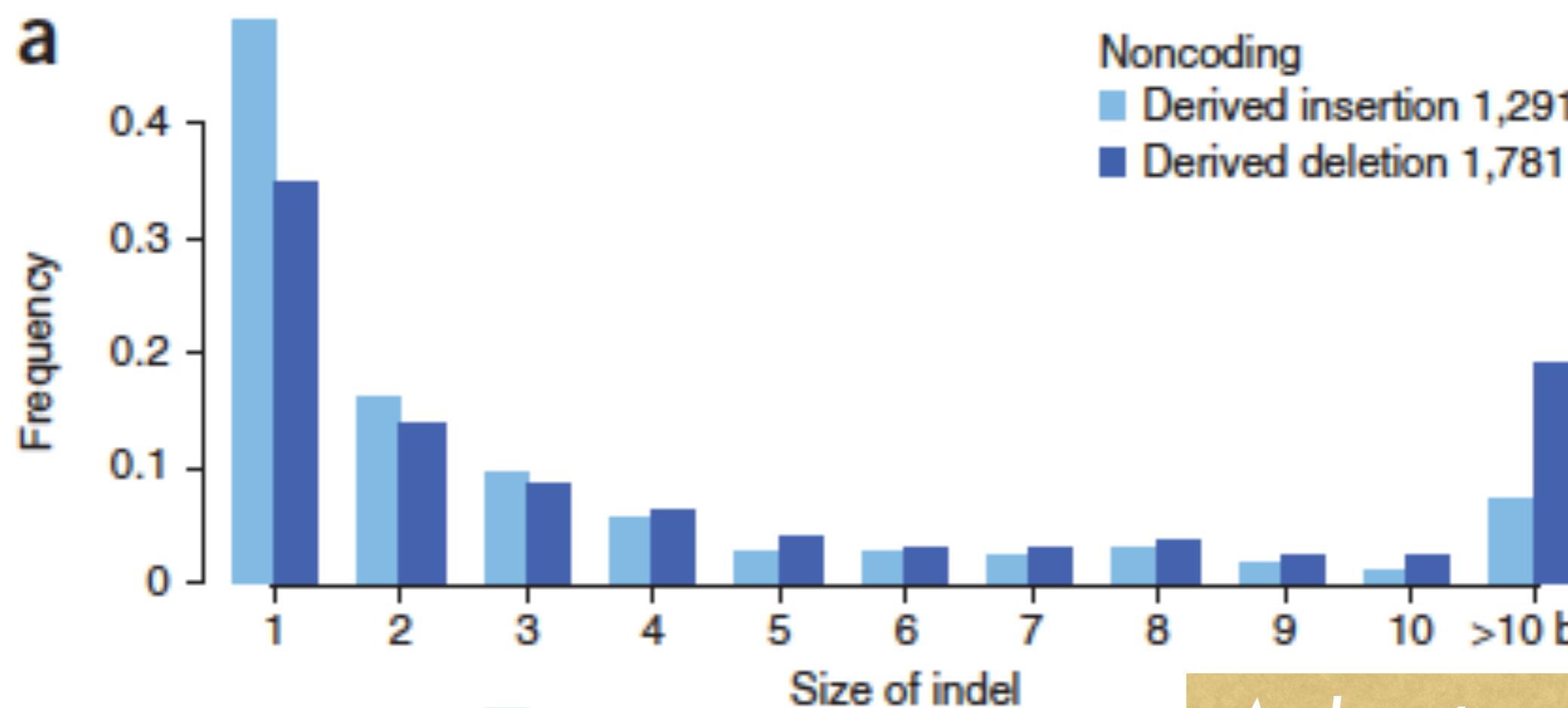
Cloudogram



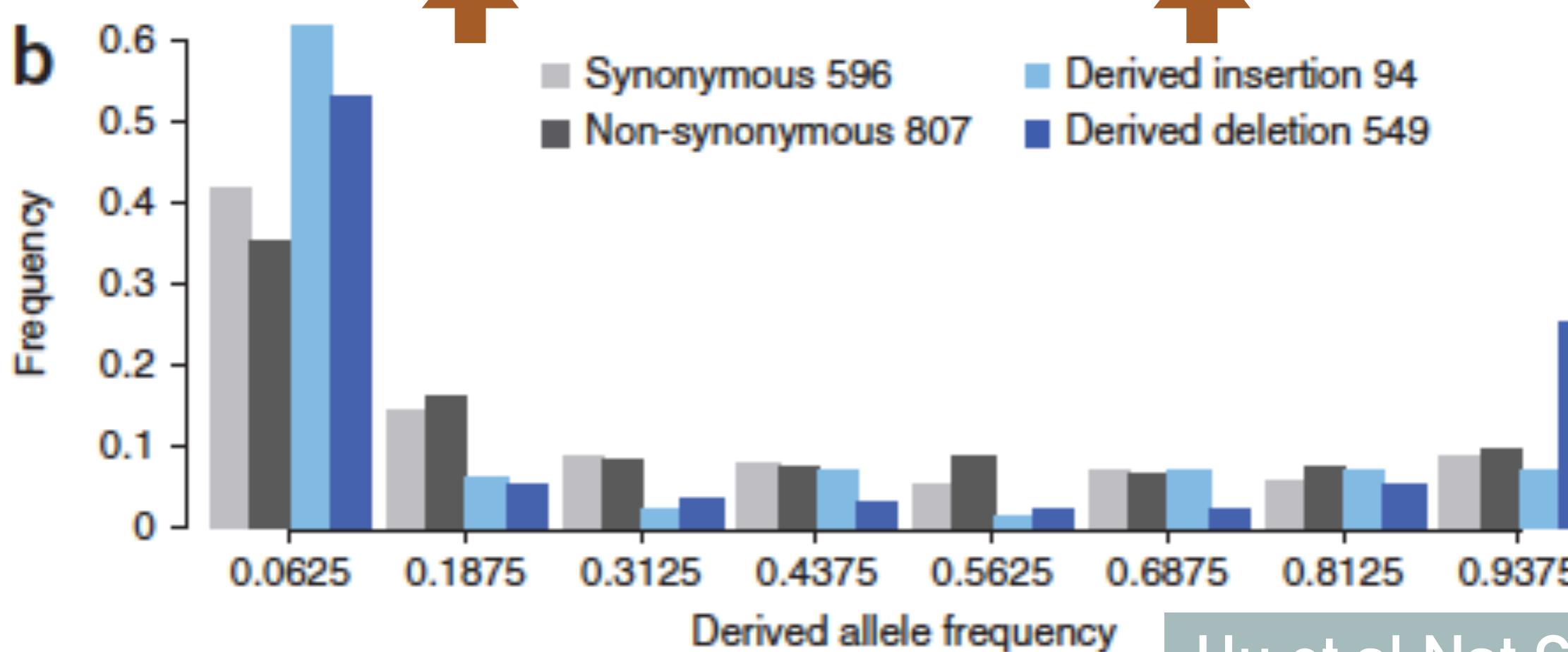
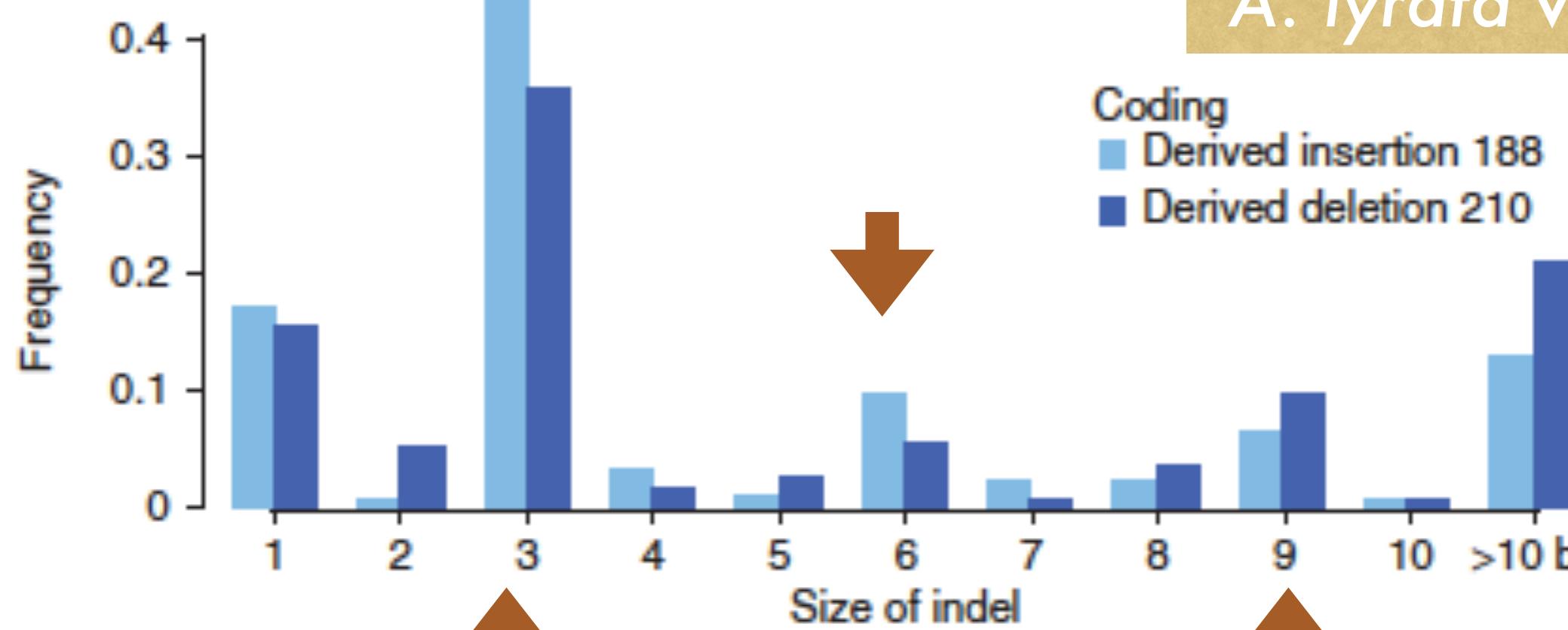
B

- Using a whole genome alignment (or in this case of transcriptome alignments) to examine consistency of phylogenetic signal across the genome.
- Individual gene trees can be inferred from sliding windows across the whole genome/transcriptome alignment
- Examine the consistency of these gene trees in a Cloudogram - drawing these as patterns impose on species tree.

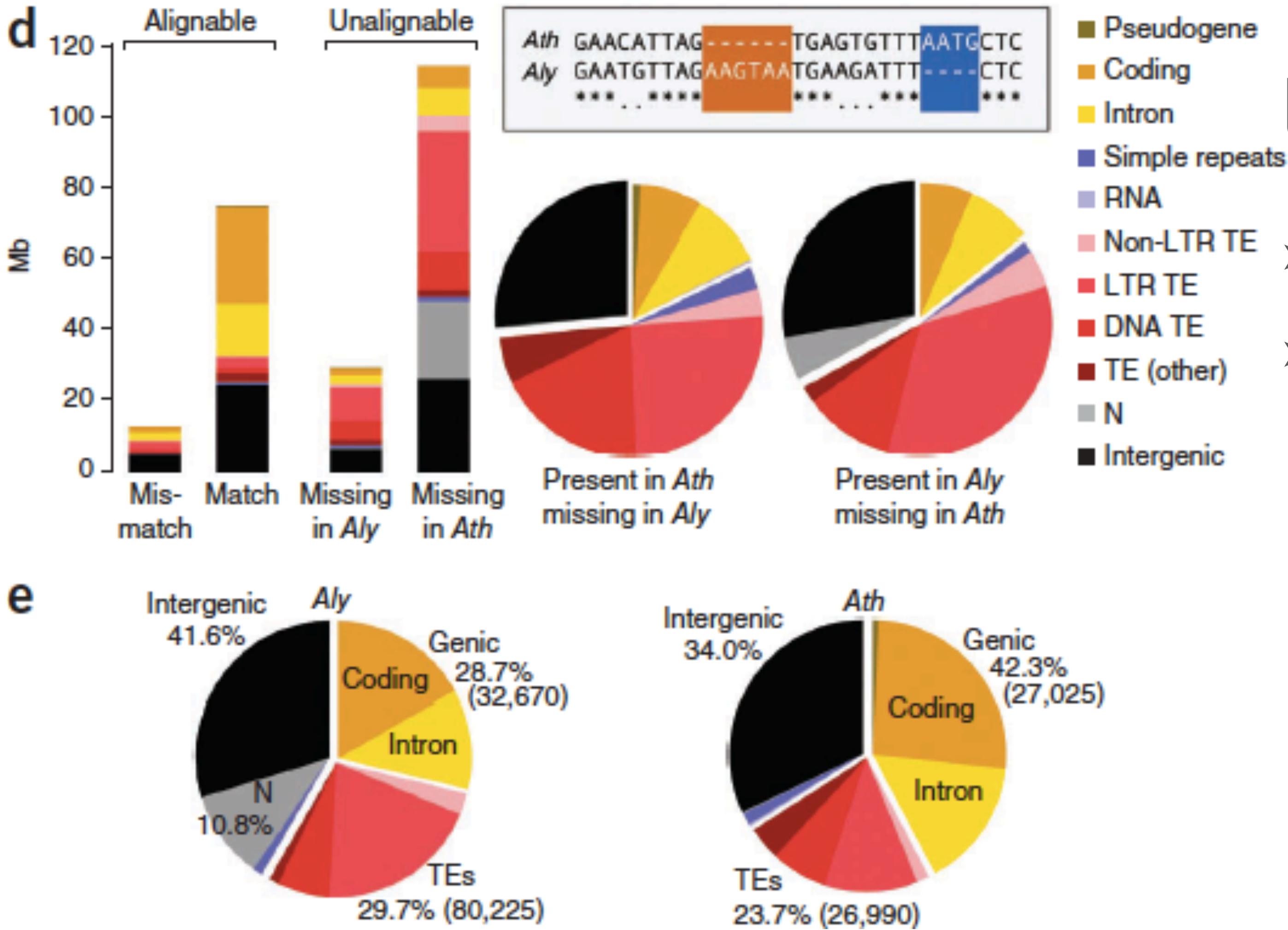
GENOME WIDE DIVERGENCE PATTERNS



A. lyrata vs A. thaliana



- Examination of insertion/deletion differences between two *Arabidopsis* (plant) genomes.
- Comparing the sizes of INDELs reveals constraint on the size of changes in coding regions
- “Derived” status was inferred by using the 95 *A. thaliana* individuals and only considering a site that is fixed in the *A. thaliana* population. *A. lyrata* allele assumed to be ancestral allele when these are different.



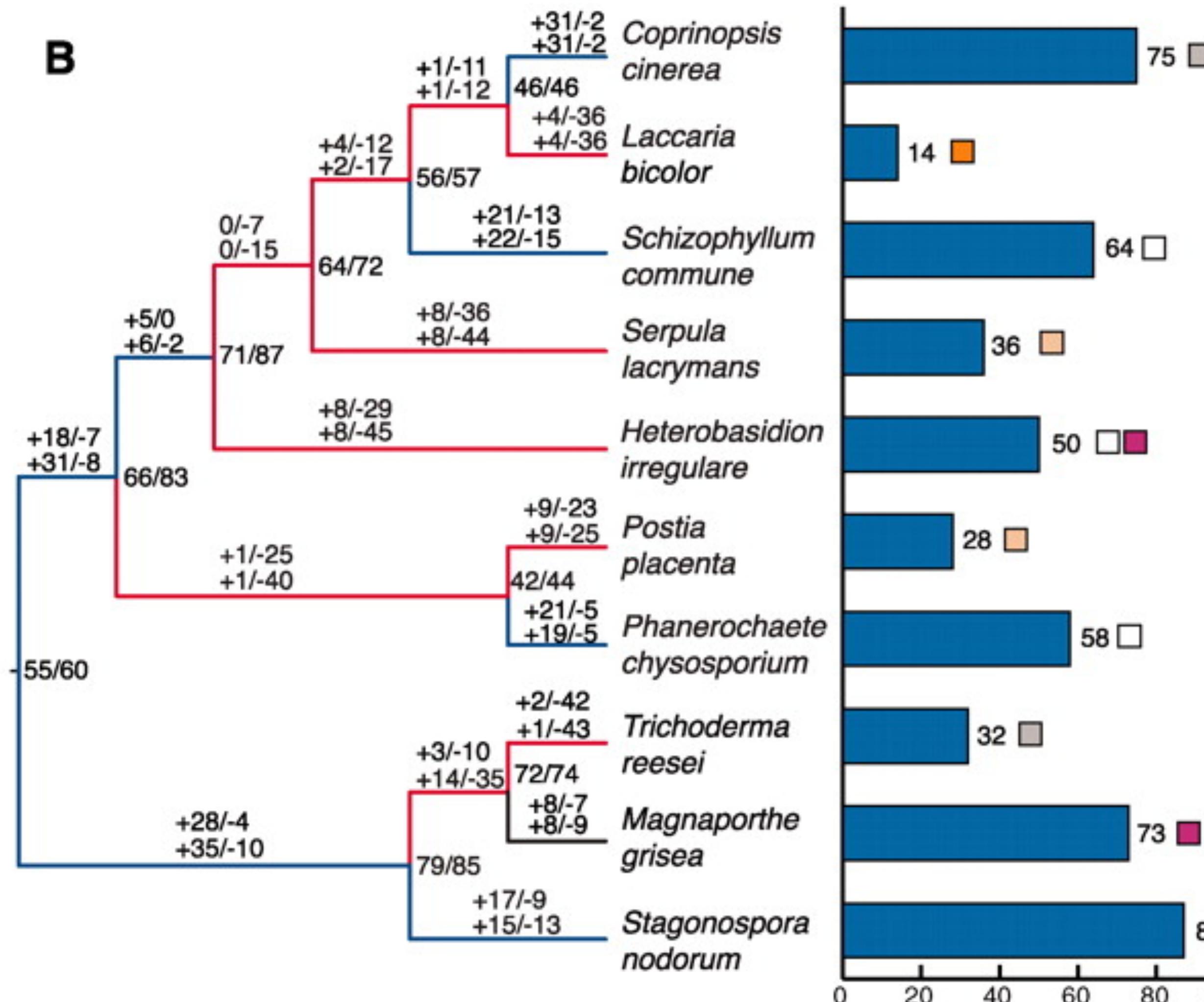
Arabidopsis lyrata genome provides comparative context to how genome size can change

UNALIGNABLE REGIONS

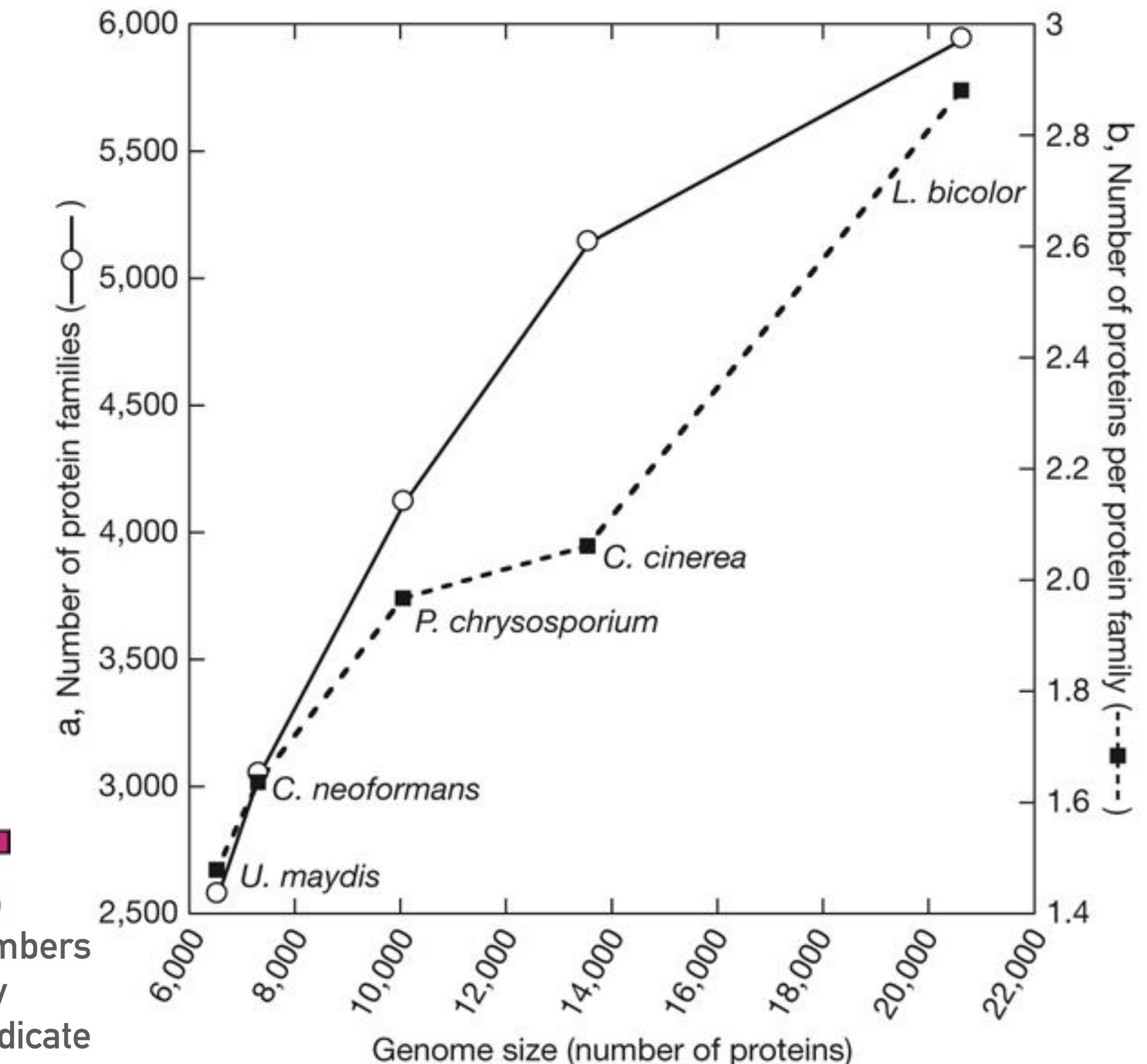
- Finding unalienable regions
- Where did they come from?
“Unalignable sites can be considered as present in one species and absent in the other, as shown in the boxed sequence diagram; matches are indicated by asterisks, and mismatches by periods.
- The histogram on the left indicates the absolute number of unalignable sites, and the pie charts in the middle compare their relative distribution over different genomic features.”

GENE FAMILY SIZE EVOLUTION

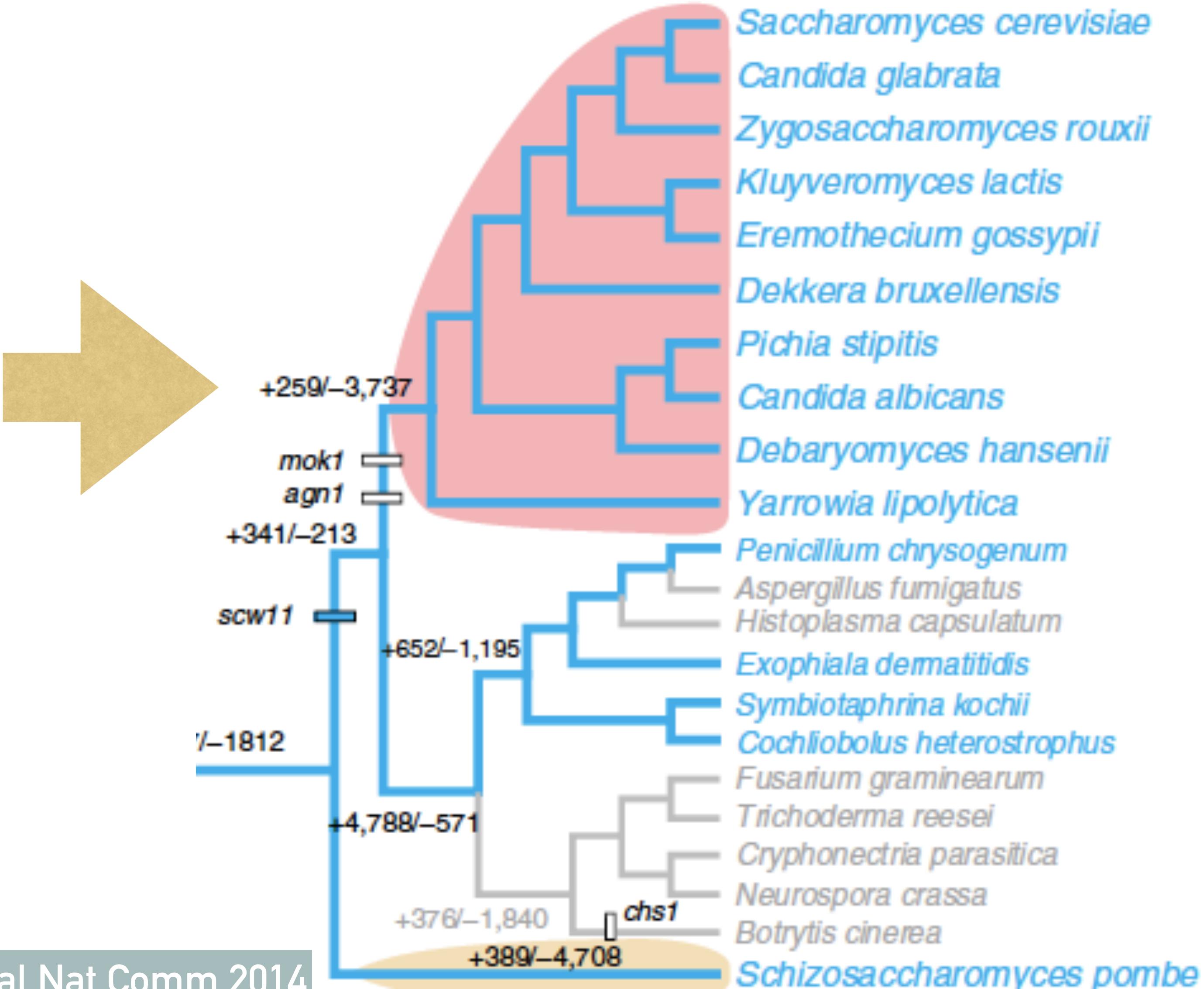
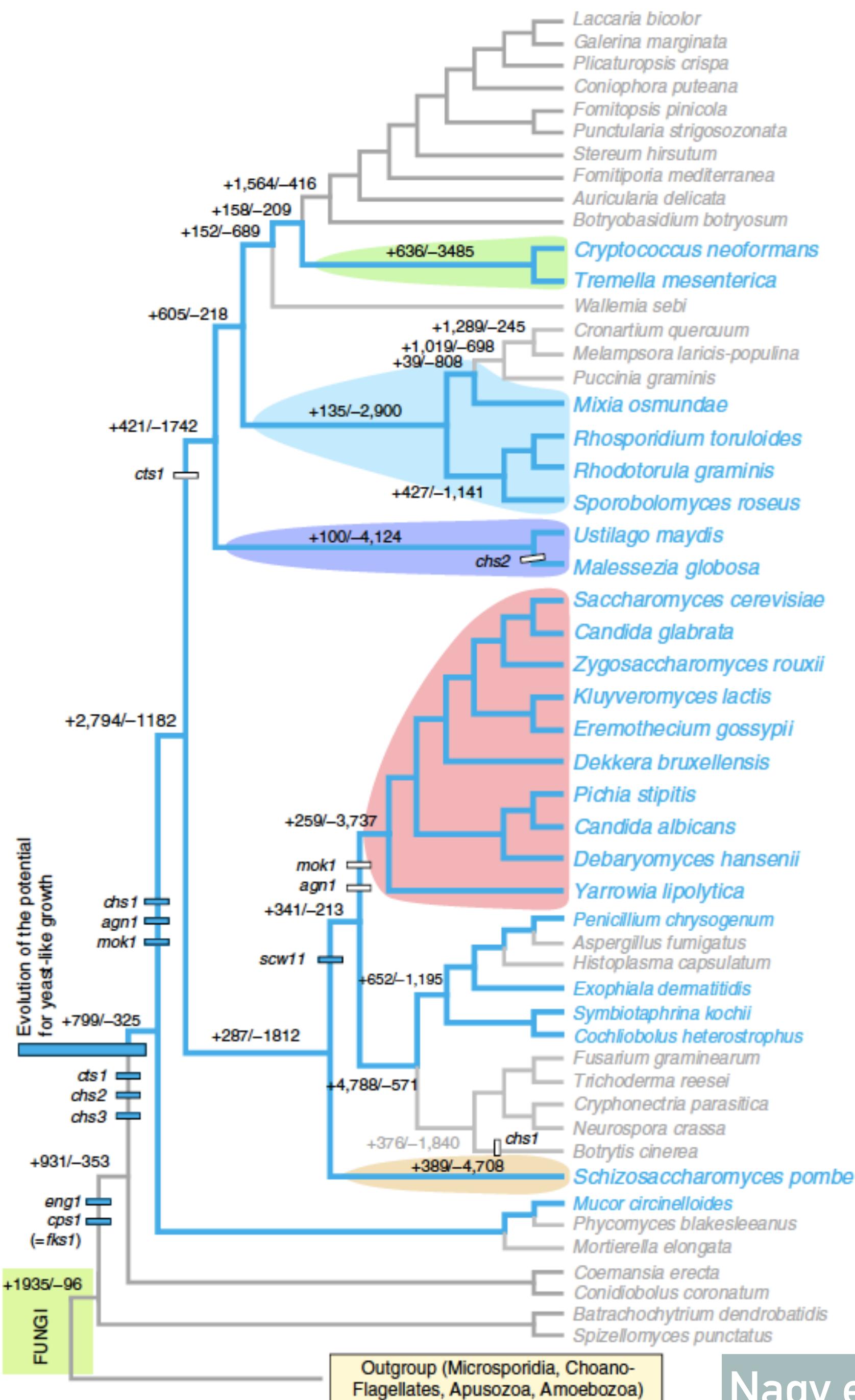
B



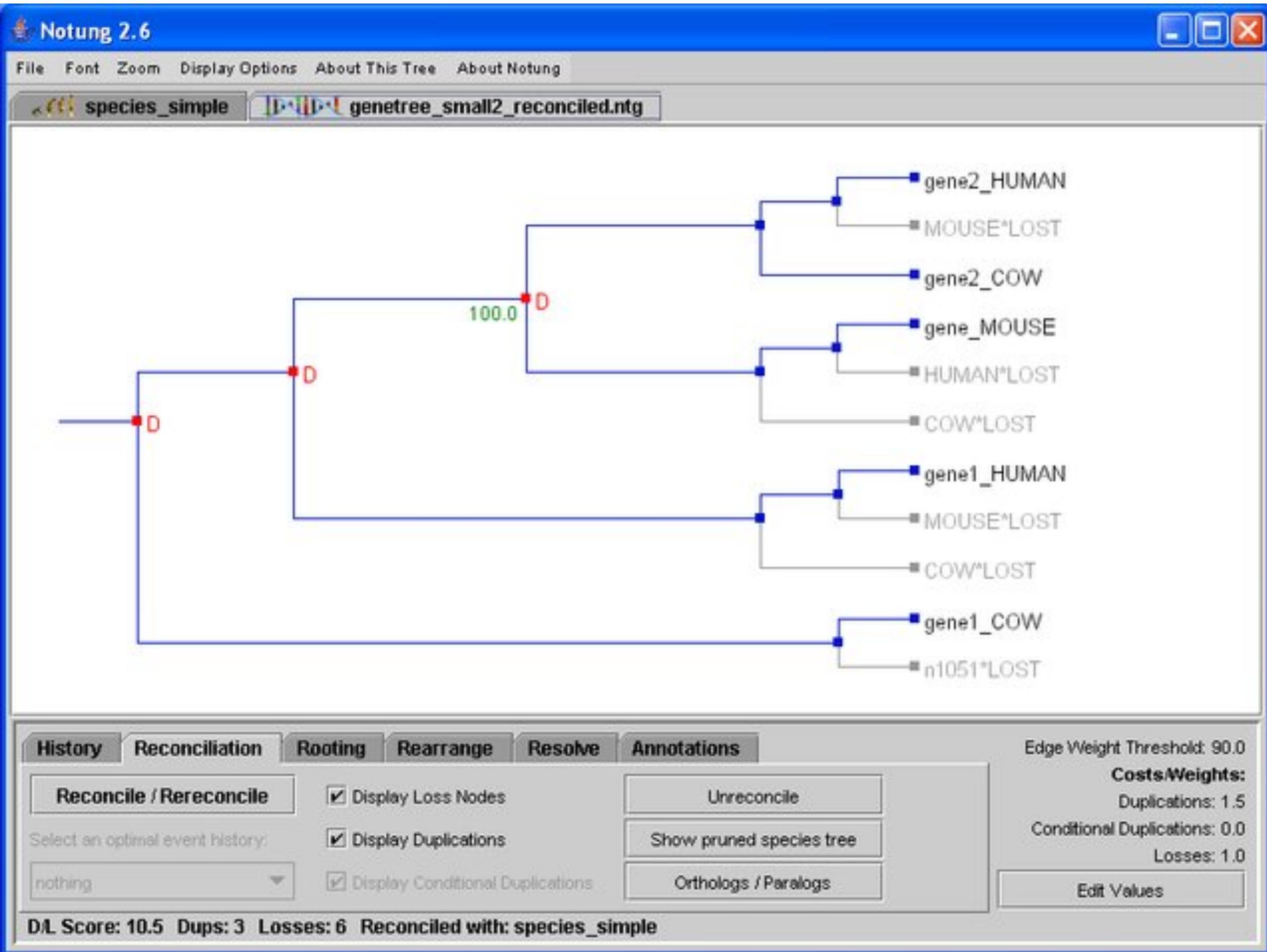
Gain / Loss Numbers at nodes and along branches indicate estimated copy numbers for ancestral species and ranges of gains and losses, respectively, estimated by using 90 and 75% bootstrap thresholds for gene trees in reconciliations. Bars indicate copy numbers in sampled genomes.



IDENTIFYING AND COMPARING GENE SIZE CHANGE FAMILIES



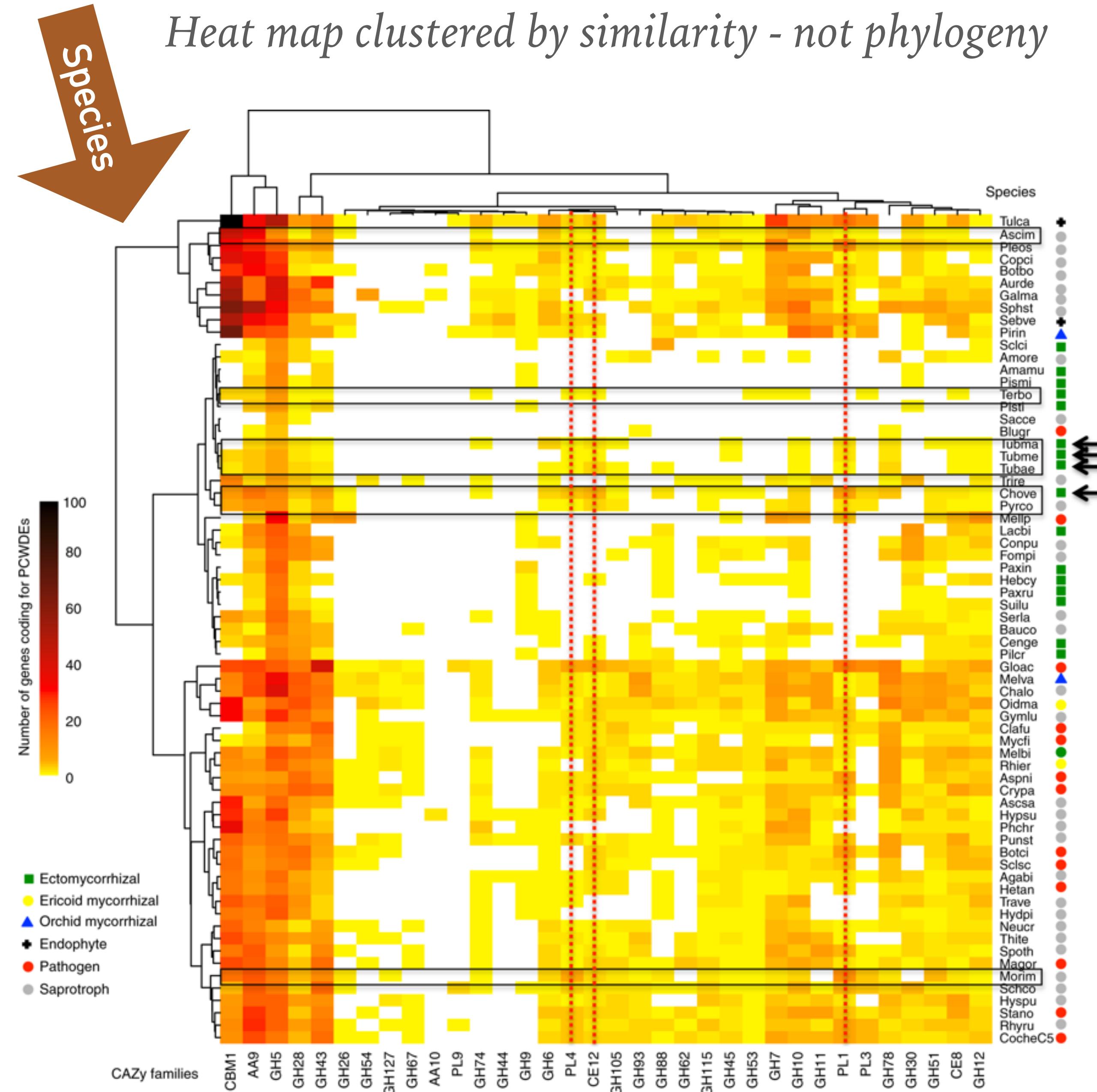
HOW TO INFER LOSSES OR GAINS?



- Gene / Tree species tree reconciliation
- NOTUNG can classify nodes and lineages on the gene tree as gains or losses
- Requires a known species tree and then input gene trees, each is reconciled with the species tree to determine

COMPARING GENE CONTENT

Heat map clustered by similarity - not phylogeny

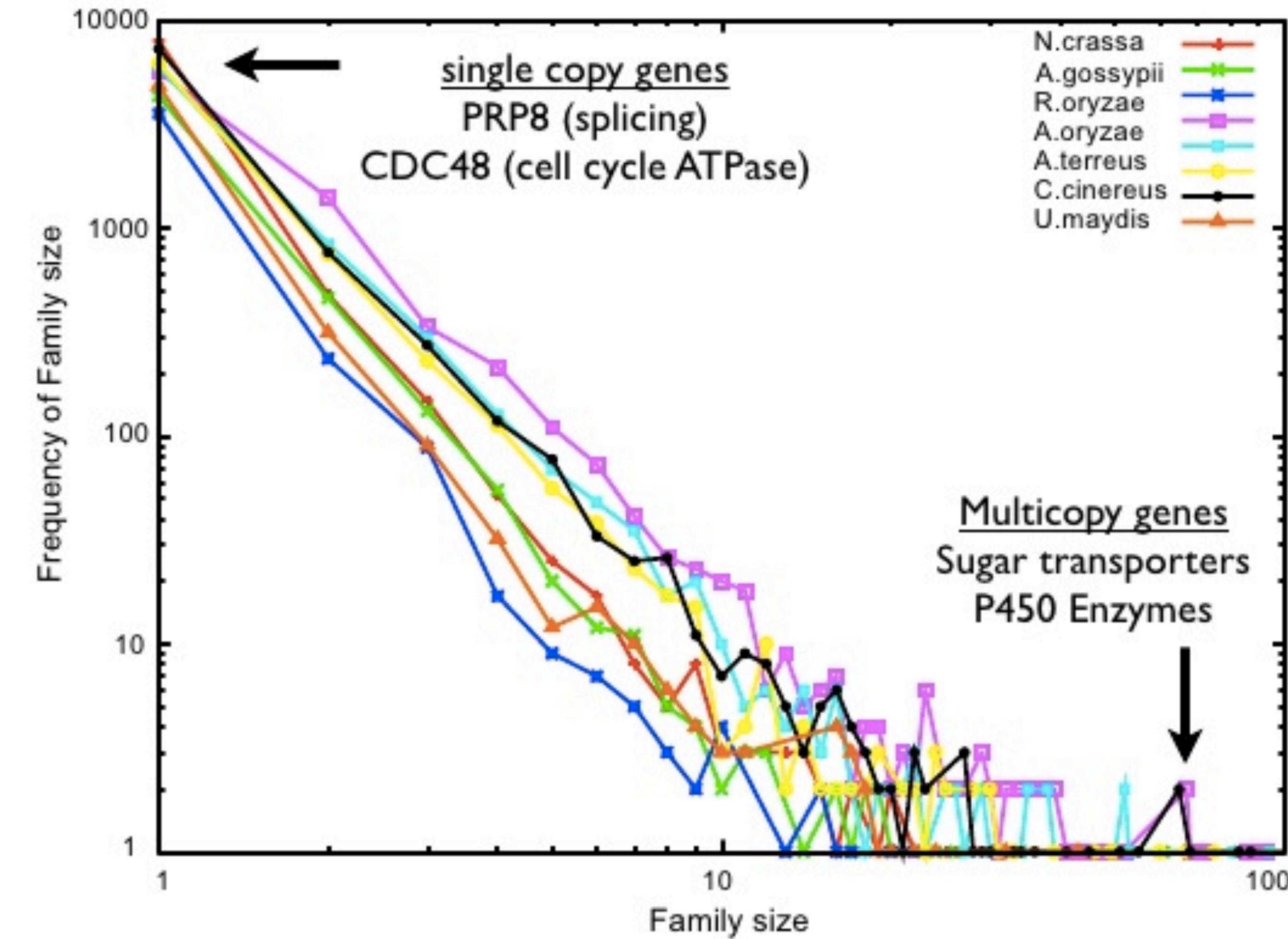
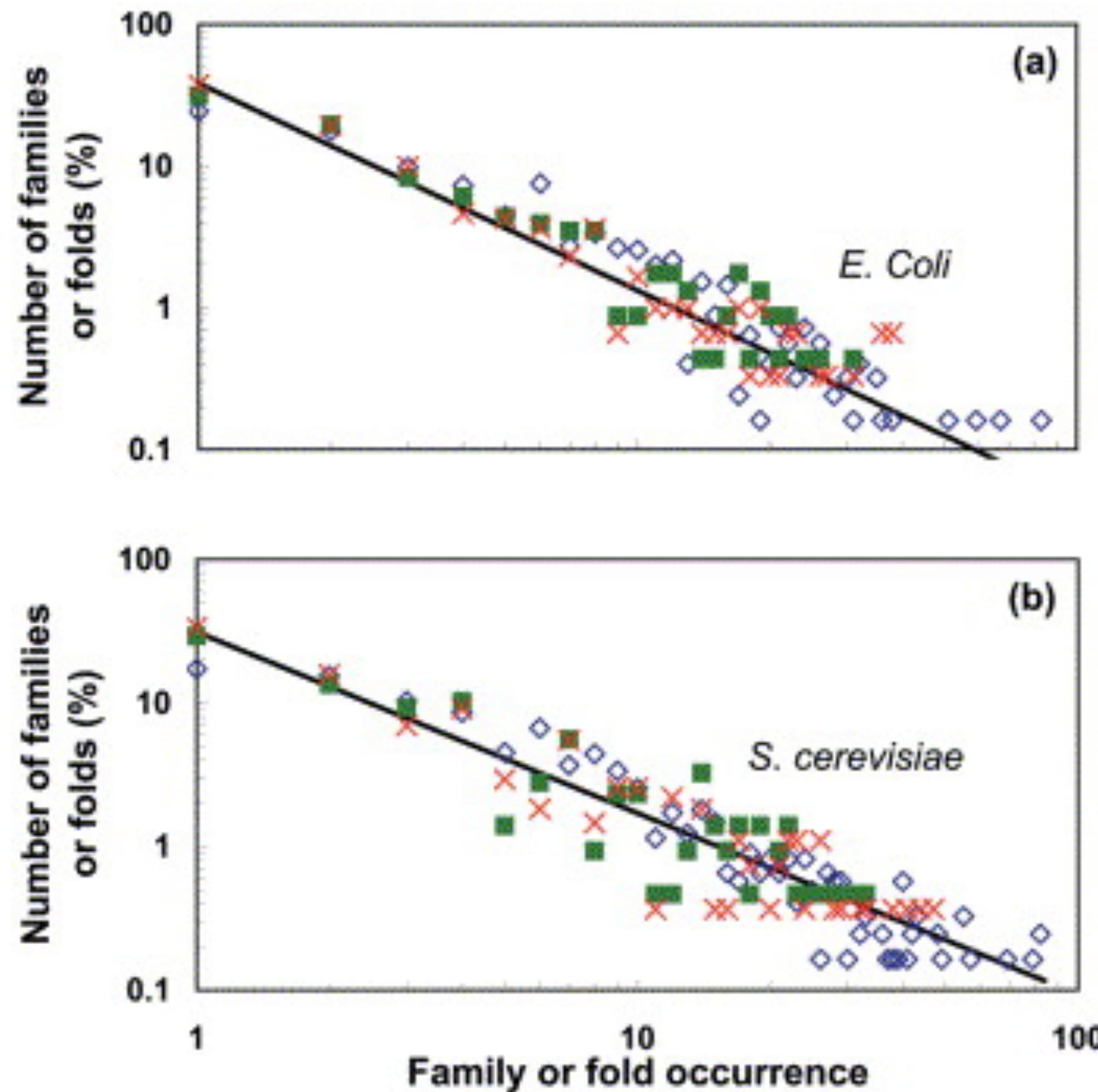


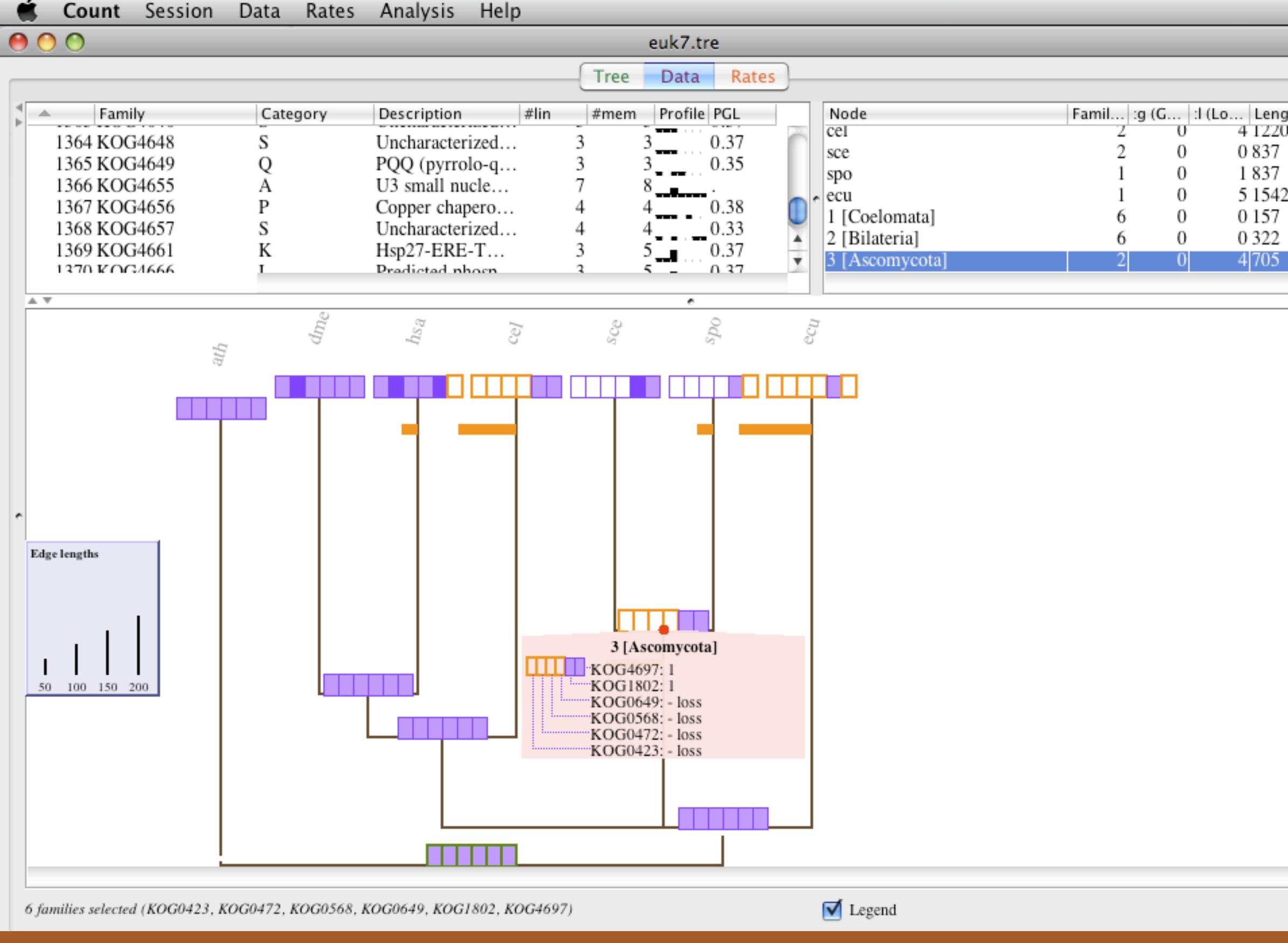
Heat map columns grouped by ecological status

A heatmap showing the correlation of various enzymes with phenotypes across different species. The columns represent species: Trave, Galma, Aurdre, Formme, Phchr, Cersu, Dicsq, Punst, Phaca, Heta, Steni, Schco, Jaajar, Bolbo, Fompi, Paspl, Wolco, Dacsp, Glocr, Conpu, and Serla. The rows represent enzyme families: CBM1, GH6, GH7, AA9, POD, AA3_2, AA5_1, AA1_1, AA3_3, AA7, AA1_2, AA3_1, AA3_4, AA1_dist, AA6, AA8, AA4, FAS, NRPS, NR-PKS, R-PKS, and TS. A color scale at the top indicates the correlation with phenotype, ranging from -1.0 (blue) to 1.0 (orange). An orange arrow points down the right side of the heatmap. A legend on the right lists the enzyme families and their descriptions.

	Trave	Galma	Aurdre	Formme	Phchr	Cersu	Dicsq	Punst	Phaca	Heta	Steni	Schco	Jaajar	Bolbo	Fompi	Paspl	Wolco	Dacsp	Glocr	Conpu	Serla	
CBM1	23	54	43	6	30	17	17	21	27	33	17	17	5	24	28	0	0	0	1	2	8	
GH6	1	3	2	2	1	1	1	1	1	3	1	1	1	3	3	0	0	0	0	2	1	
GH7	4	8	8	2	9	3	4	5	5	16	1	3	2	5	7	0	0	0	0	2	0	
AA9	18	19	20	13	15	9	15	14	11	29	10	16	22	15	32	4	2	2	0	4	10	5
POD	25	10	5	16	15	15	12	10	8	9	8	5	0	0	0	0	0	0	0	0	0	
AA3_2	17	32	30	24	27	18	30	19	37	36	29	40	18	16	21	15	21	9	8	20	14	8
AA5_1	9	15	8	4	7	3	9	9	6	16	5	8	2	4	5	4	3	4	3	2	6	3
AA1_1	7	8	0	10	0	7	11	12	0	11	14	15	2	1	0	5	2	3	0	4	6	4
AA3_3	4	6	6	3	3	3	4	4	4	4	3	7	4	2	3	5	4	5	1	2	5	5
AA7	0	9	2	0	0	0	4	1	0	3	3	3	4	1	4	5	0	0	3	0	0	0
AA1_2	2	1	1	1	1	1	1	1	1	1	1	2	0	1	1	1	1	1	1	1	1	1
AA3_1	1	1	1	1	1	1	1	1	1	1	1	1	3	0	0	0	0	0	1	2	2	0
AA3_4	1	0	3	0	1	0	0	1	0	0	0	0	1	2	0	0	0	0	1	0	0	0
AA1_dist	1	0	7	1	4	1	1	0	8	0	2	3	4	1	4	1	1	1	2	0	1	1
AA6	1	3	4	3	4	0	1	2	3	2	2	1	4	3	1	1	1	1	1	3	2	2
AA8	2	1	2	1	2	2	2	1	2	2	3	2	2	0	0	0	0	0	0	4	4	4
AA4	0	0	0	0	0	0	0	0	0	0	0	0	2	0	0	0	0	2	3	0	0	0
FAS	1	2	1	1	1	1	1	1	1	1	2	1	1	1	1	1	1	1	1	1	1	1
NRPS	2	1	1	1	1	1	2	1	9	1	0	1	7	0	3	0	0	0	0	0	1	8
NR-PKS	1	1	0	1	0	1	1	1	1	1	2	1	1	5	1	1	1	0	1	1	1	1
R-PKS	1	4	1	4	0	2	2	4	0	0	2	5	0	2	1	6	5	6	1	10	4	9
TS	5	4	3	13	1	6	9	3	1	7	6	6	2	11	5	8	6	6	2	6	7	5

GENE FAMILY SIZES FOLLOW A POWER LAW AND CAN BE MODELED WITH A BIRTH-DEATH MODEL





http://www.iro.umontreal.ca/~csuros/gene_content/count.html

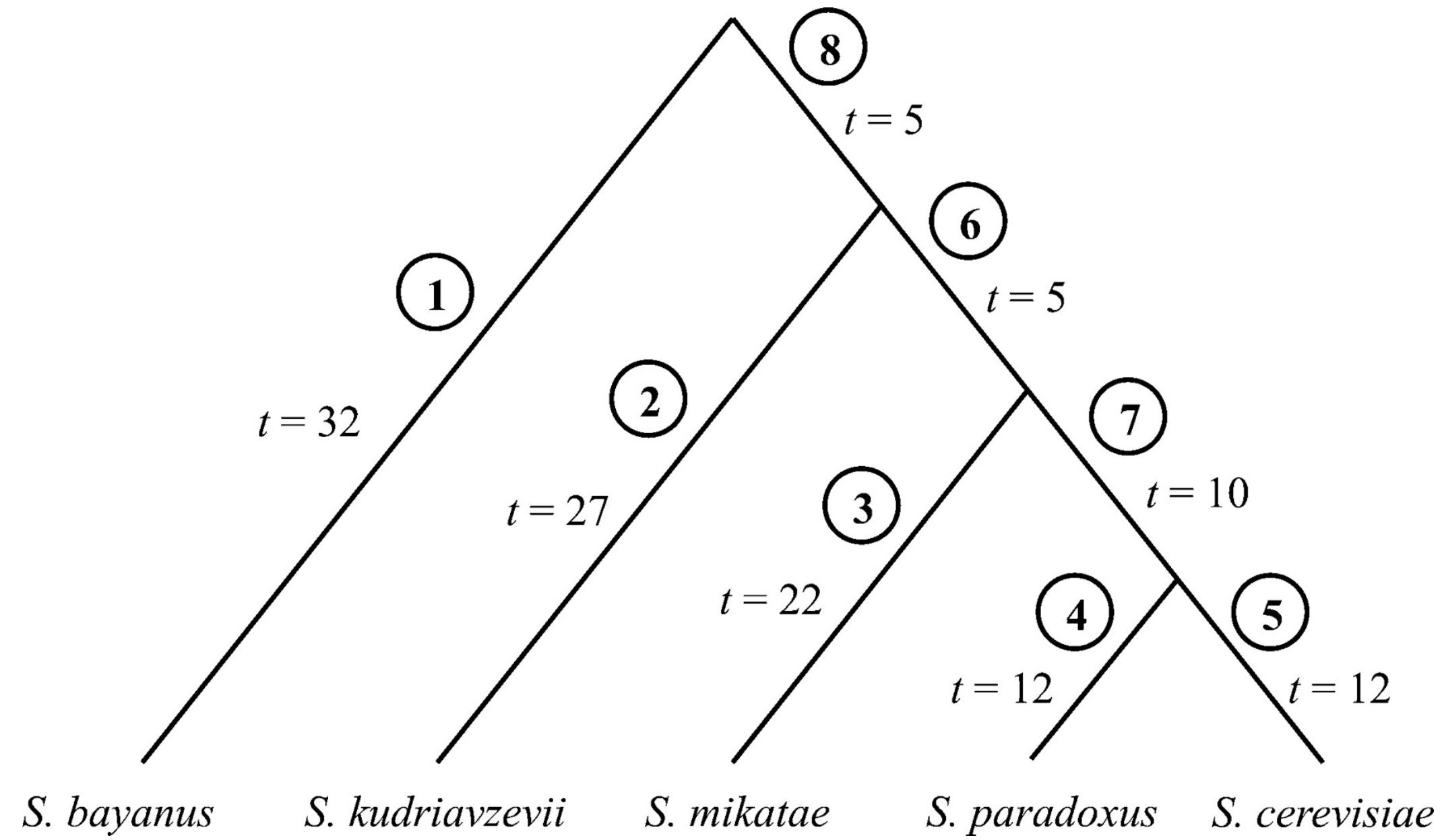
Species ->	A	B	C
OG1	1	3	1
OG2	10	2	2
OG3	5	5	2

GENE FAMILY SIZE CHANGE USING COUNTS

- Using counts of copy number of genes in a gene family can examine the patterns of gains and losses to:
- Find branches with excess of gains/losses
- Identify gene families with extreme changes

- Count is a tool to evaluate homolog family sizes (phylogenetic profiles), or other numerical census-type characters along a phylogeny.

CAFE - COMPUTATION ANALYSIS OF (GENE) FAMILY EVOLUTION



Branch #	Expansions	No change	Contractions	Average expansion
1 ($t = 32$)	97	3181	239	-0.050
2 ($t = 27$)	383	3032	102	0.095
3 ($t = 22$)	509	2922	86	0.147
4 ($t = 12$)	96	3383	38	0.019
5 ($t = 12$)	44	3426	47	0.021
6 ($t = 5$)	3	3491	23	-0.005
7 ($t = 10$)	10	3313	194	-0.052
8 ($t = 5$)	2	3515	0	0.001

Family	Family sizes	Branch	Score
Flavodoxin	(2 (3 (5 (1 1))))	3	0.11

EVOLUTIONARY PATTERNS GENOME DUPLICATION

GBNTN_000002-T1 GACCTGCTTGA**CTATGAAGAAGCCTCGCATGGTATTCTTGTGAACCTCCTCCACCCAA**T
GBNWO_006121-T1 GACCTGCTTGA**CTATGAAGAAGCCTCGCATGGTATTCTTGTGAACCACCTTACCCAA**T

GBNTN_000002-T1 ACACTGTTACTAGGATCTGAGAATGATGGTGGTACTTCTGTTATATCGATTCCCTTATTG

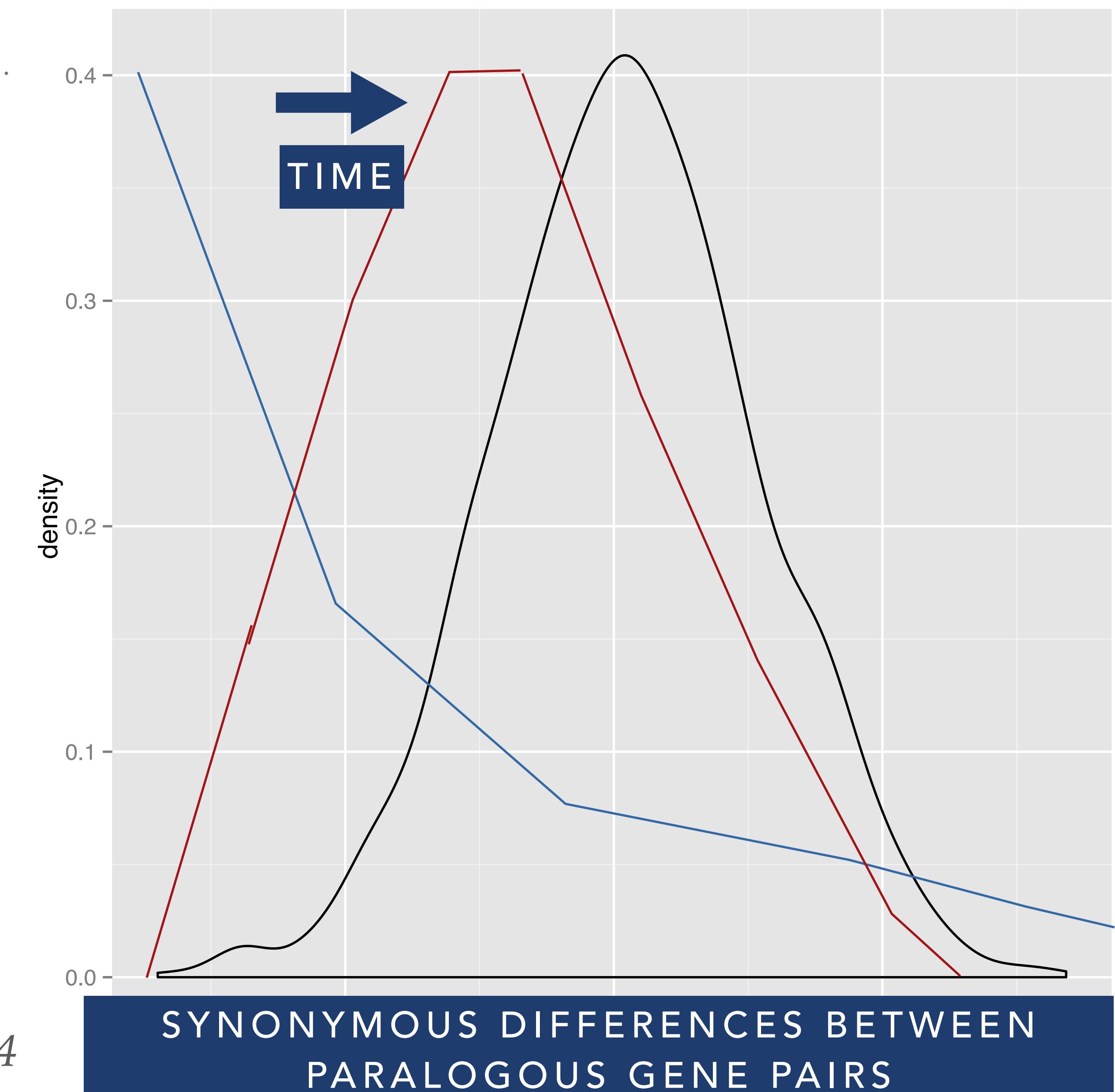
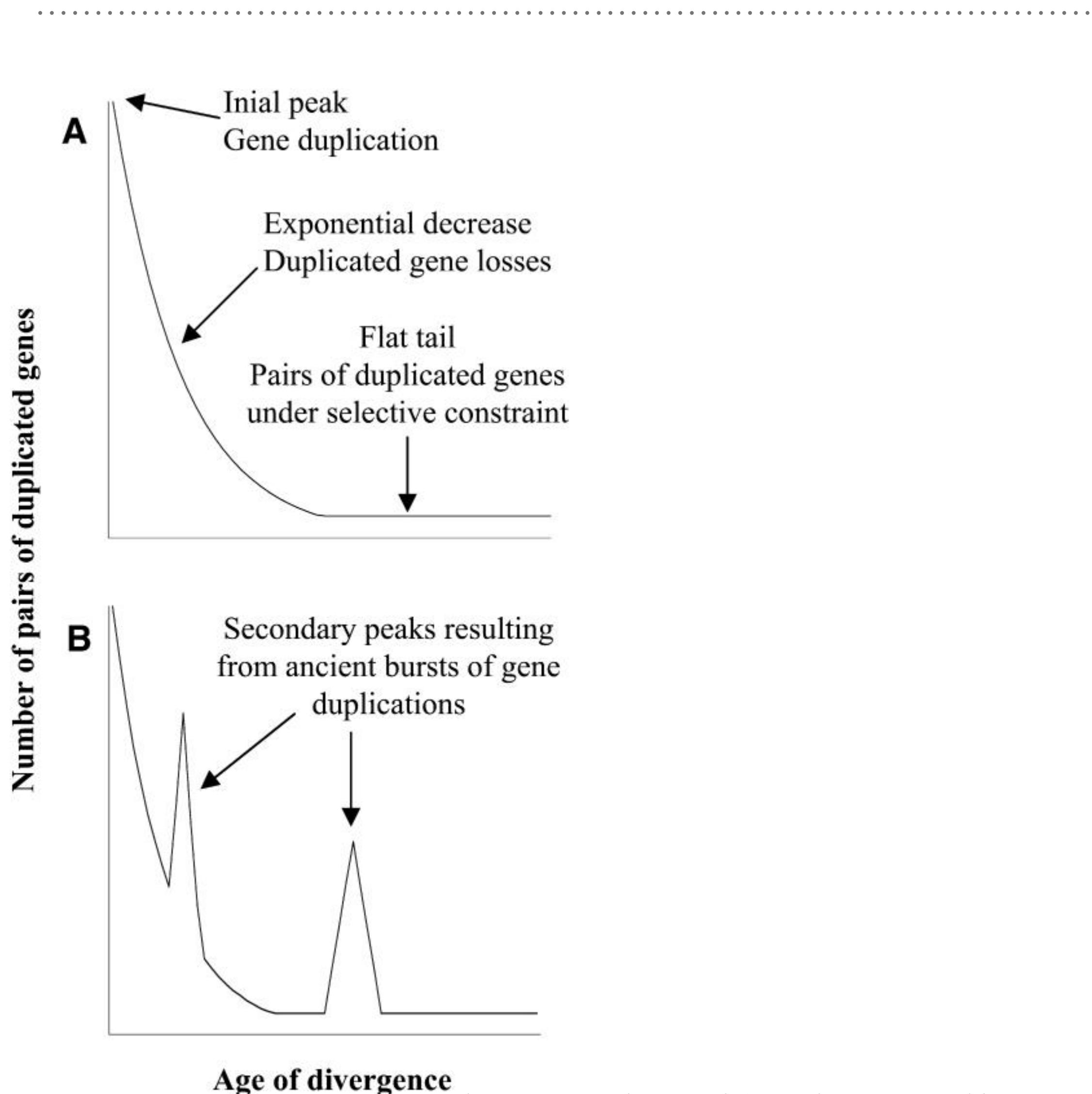
GBNW0_006121-T1 ACGCTGCTATTAGGATCTGAGAATGATGGTGGTACTTCTGTTATATTGATTCCCTTATTG

GBNTN_00002-T1 EAKWNGQDEAGHWTQEDASELFITE-----TFDLPYLPFQIRLFHGANKDSDDDR
GBNWO_006121-T1 EAKWNGQDEAGHWTQEDASELFITETFDLPYLPTFDLPYLPFQIRLFHGANKDSDDDR

MOLECULAR EVOLUTION

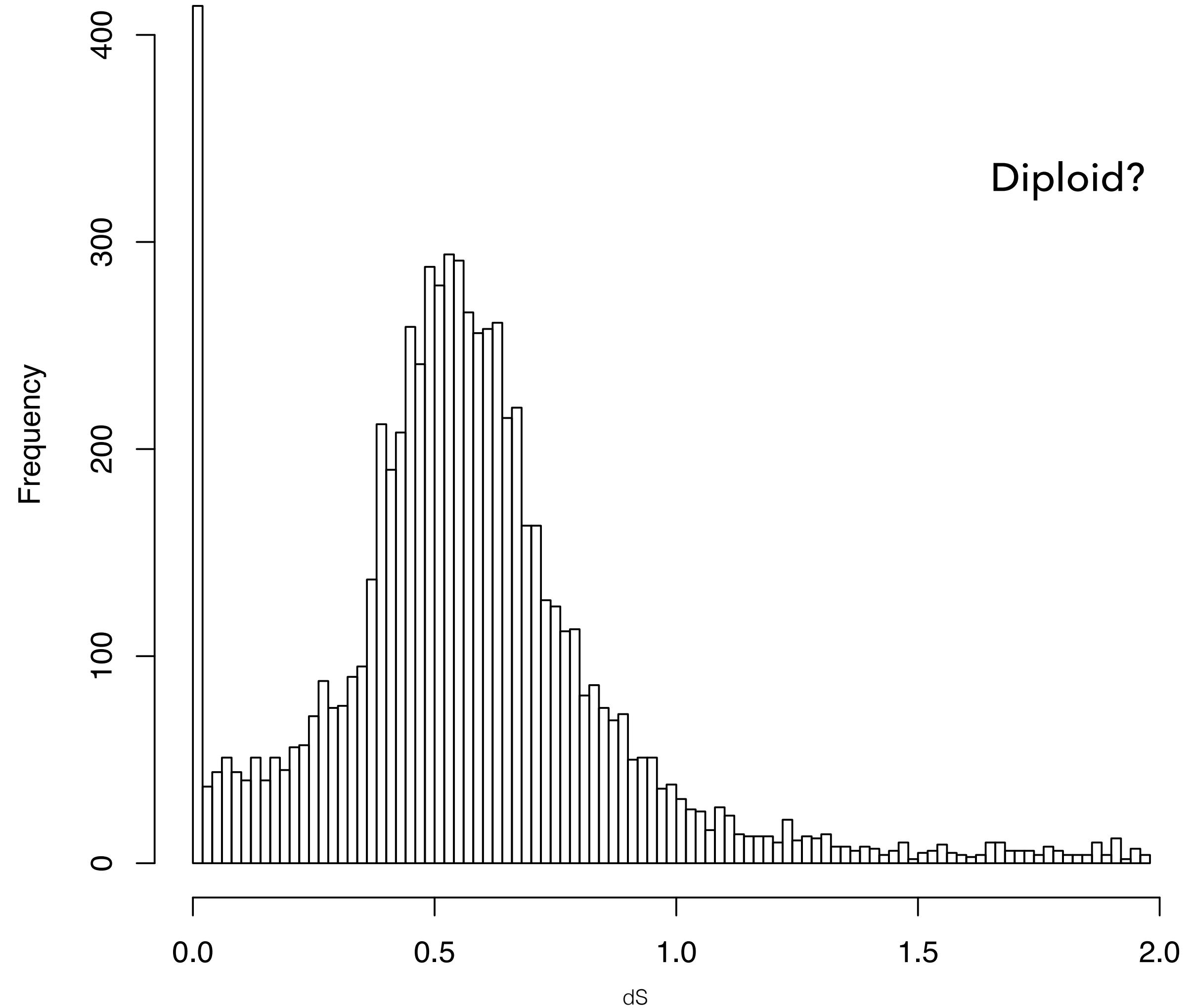
- Phylogenetics is about inferring evolutionary history of group of sequences, characters, taxa, etc
 - Molecular evolution is specifically focusing on changes in molecular sequence (protein, DNA/RNA). One aspect is examining rates of change of sequences.
 - changes in sequence calculated under different substitution models and calculate the distance.
 - K_s = synonymous site change
 dS = synonymous site rate of change
 - K_a = non-synonymous site change
 dN = non-synonymous site rate of change

KS PLOTS AND DUPLICATION



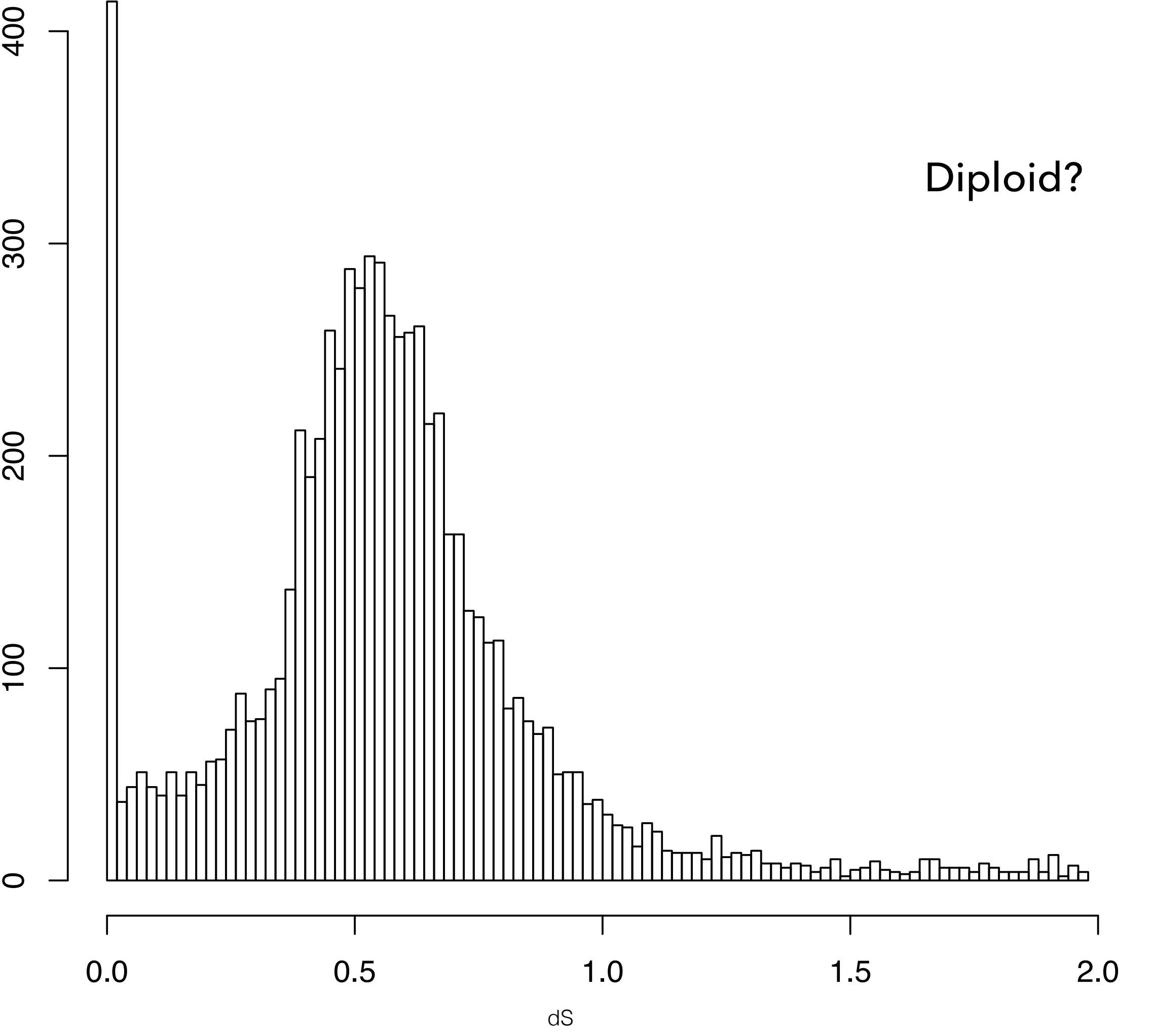
Allomyces macrogynus
(Blastocladiomycota Chytrid)

Diploid?

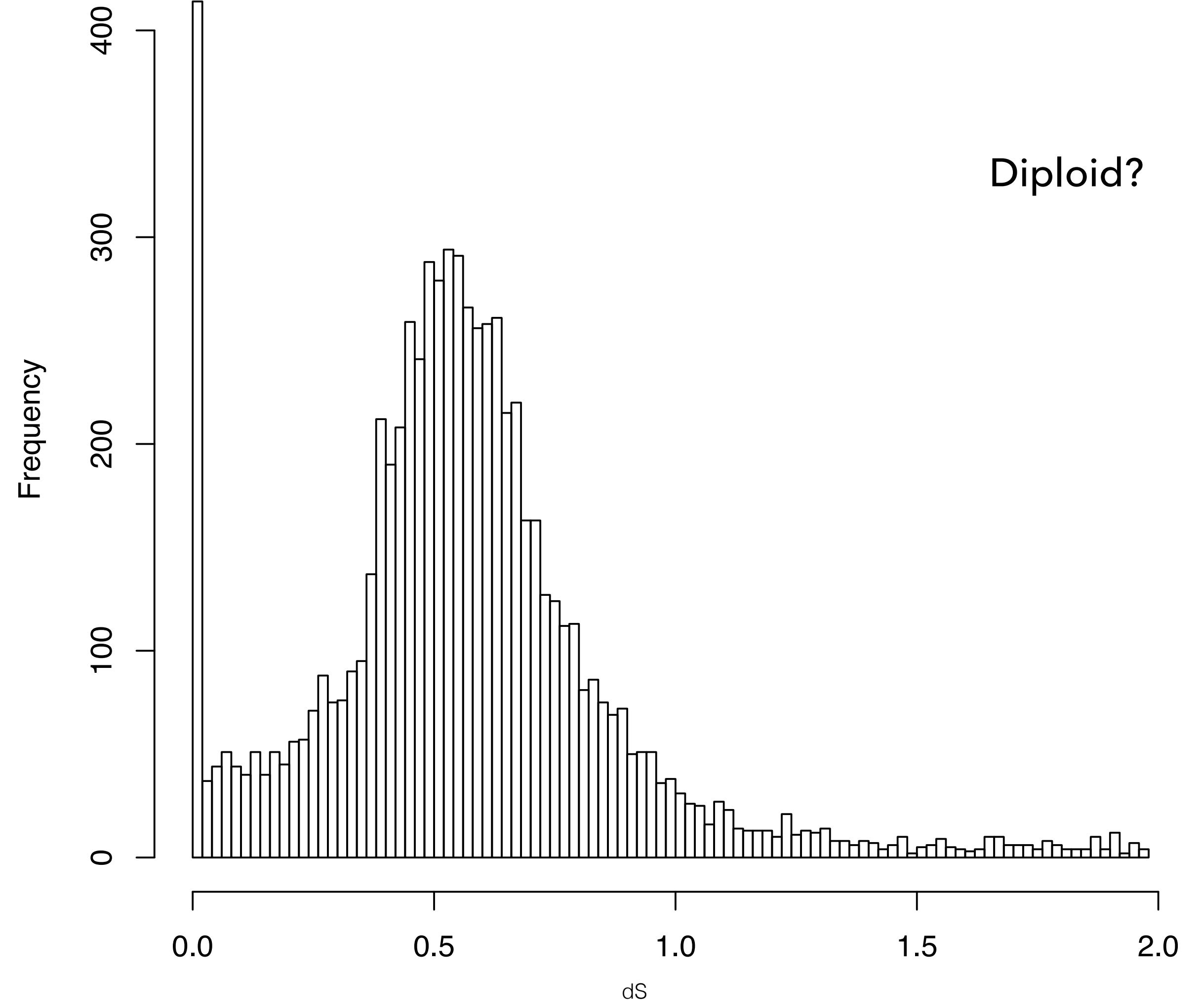
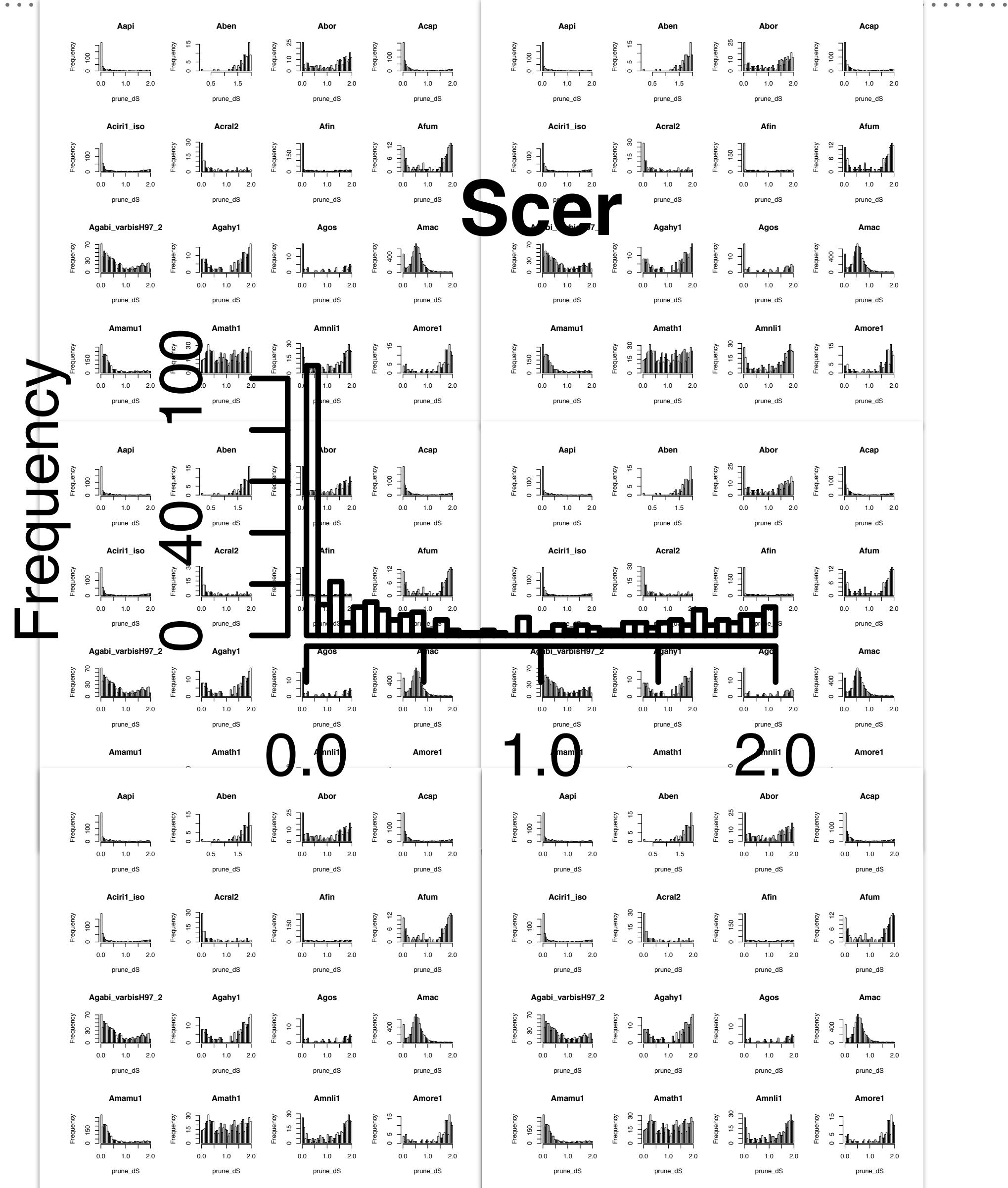


Allomyces macrogynus (Blastocladiomycota Chytrid)

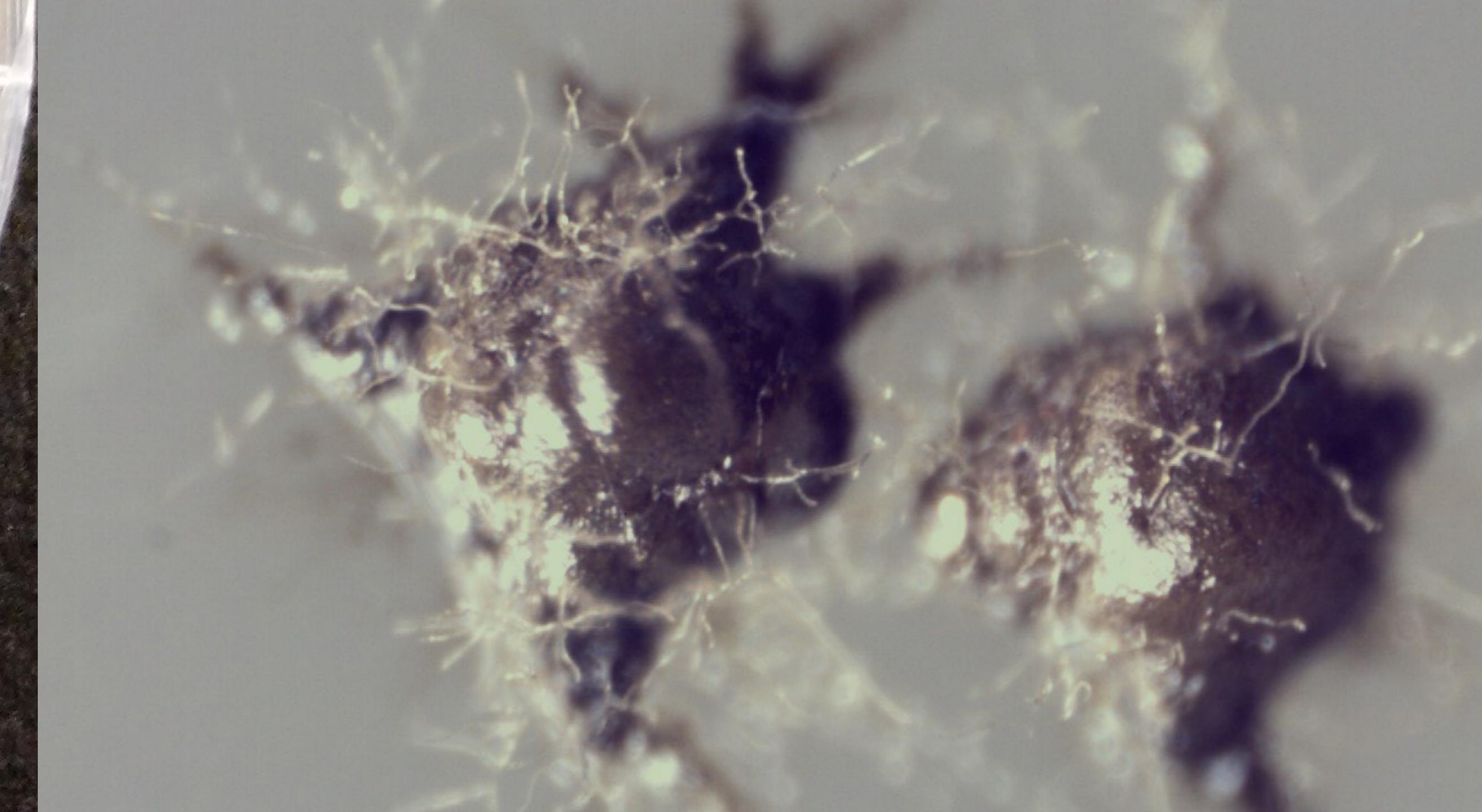
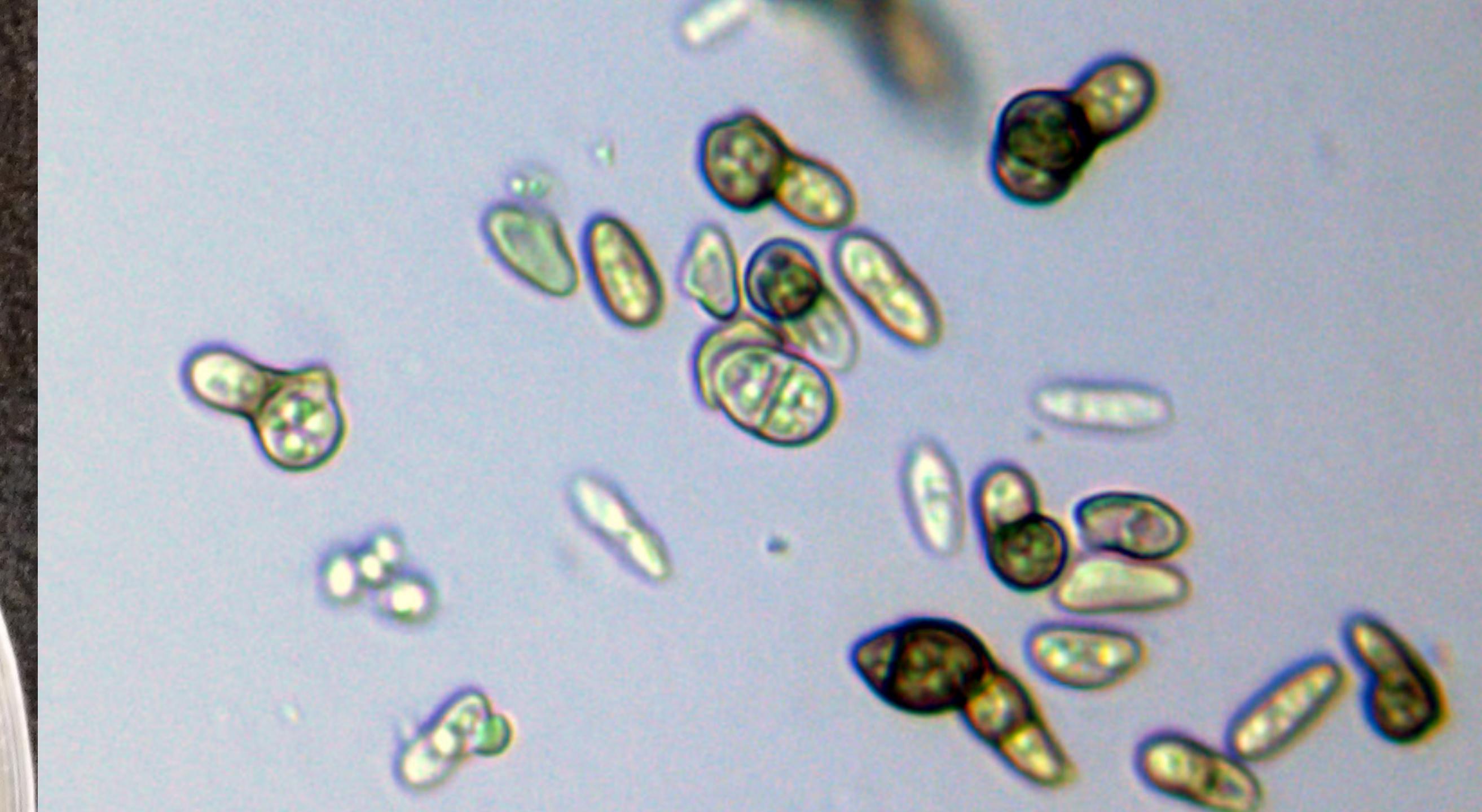
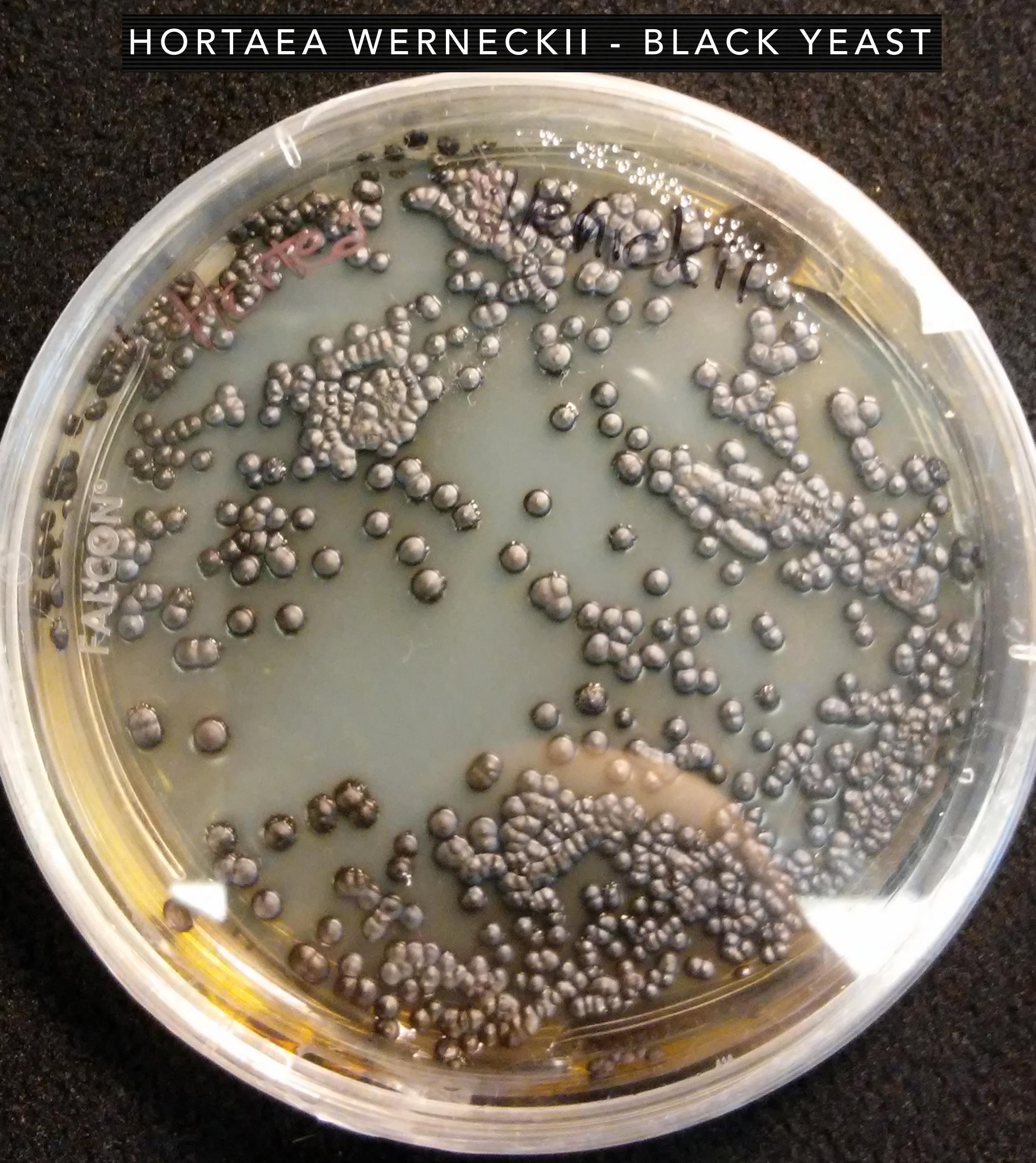
Diploid?



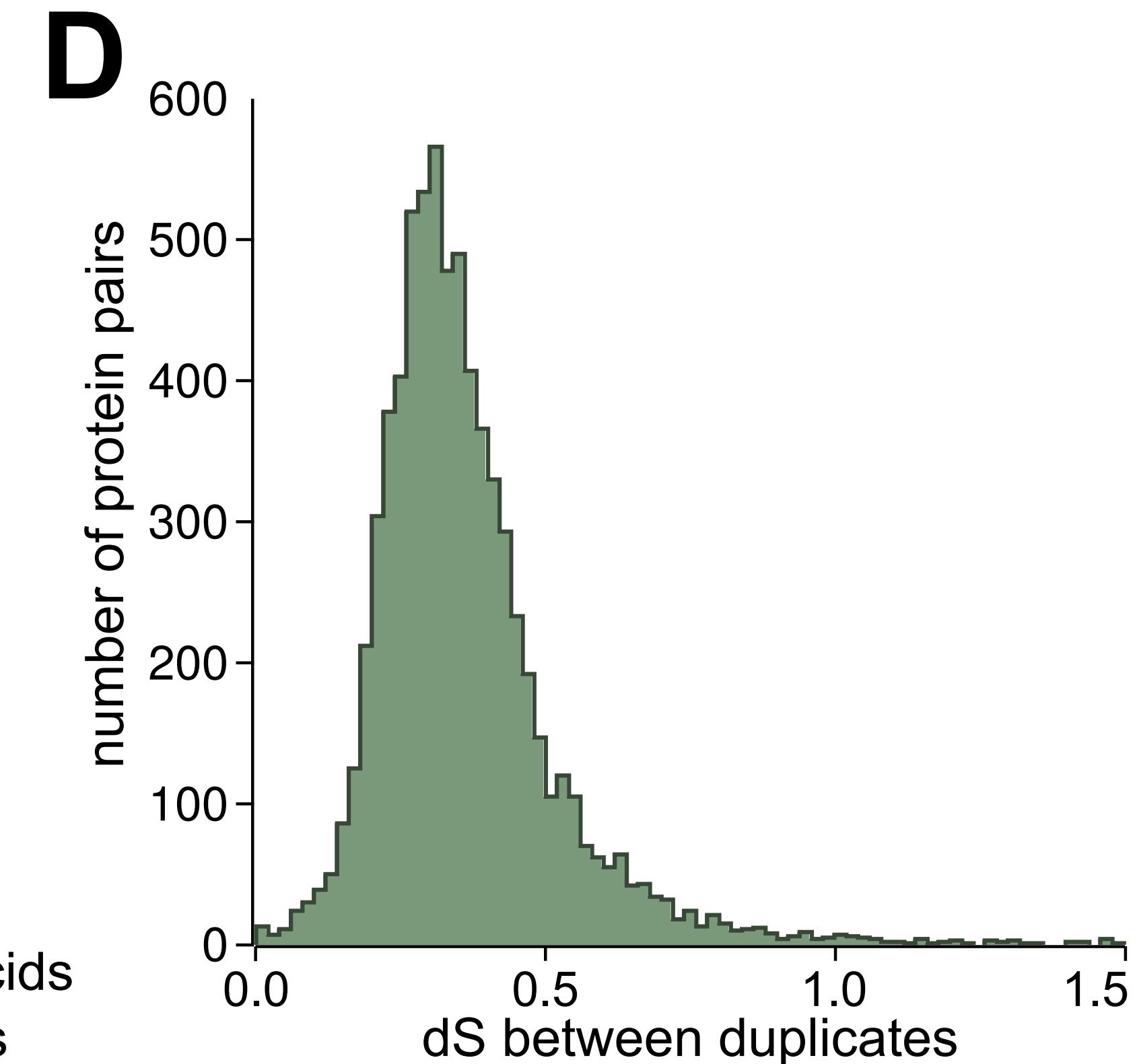
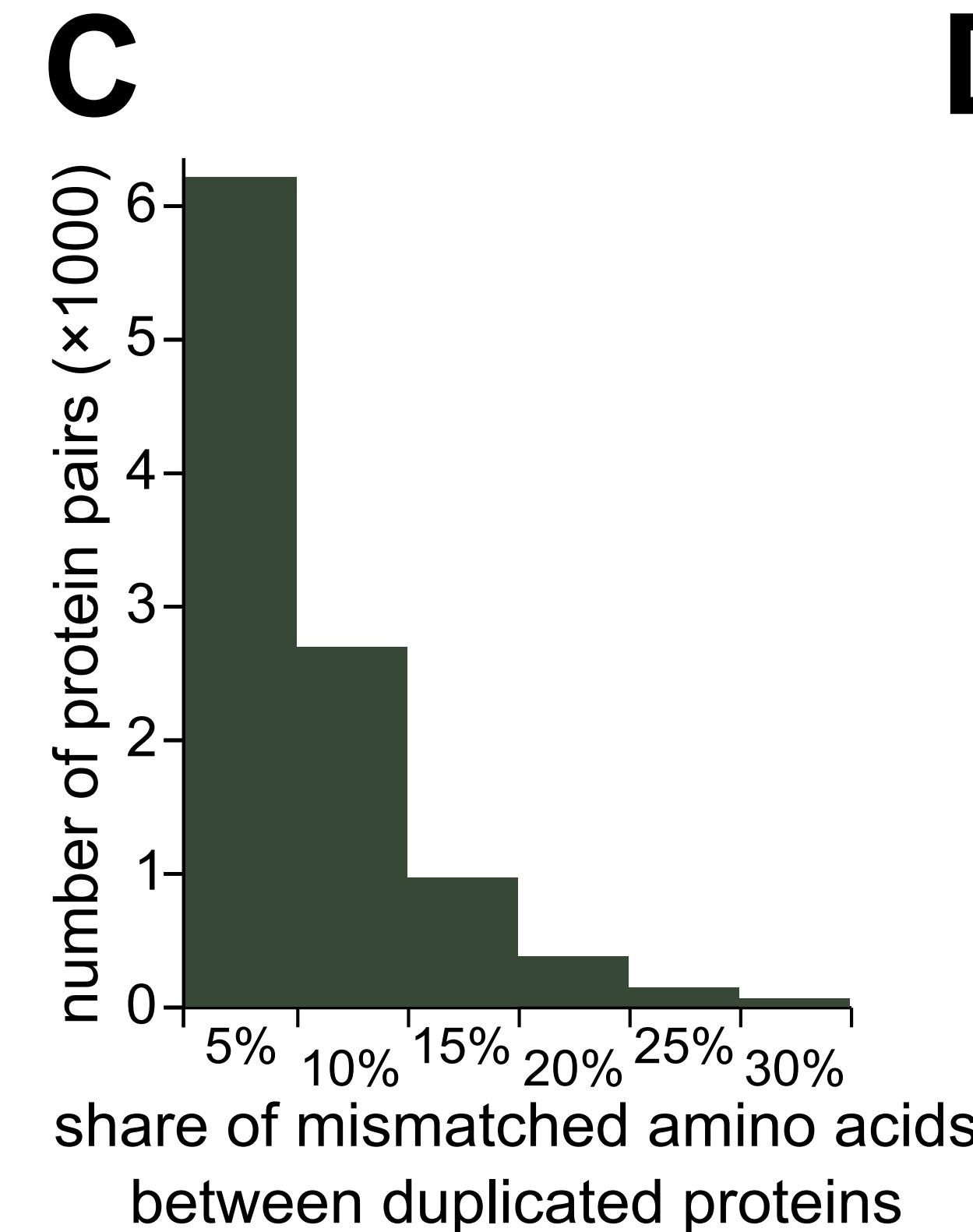
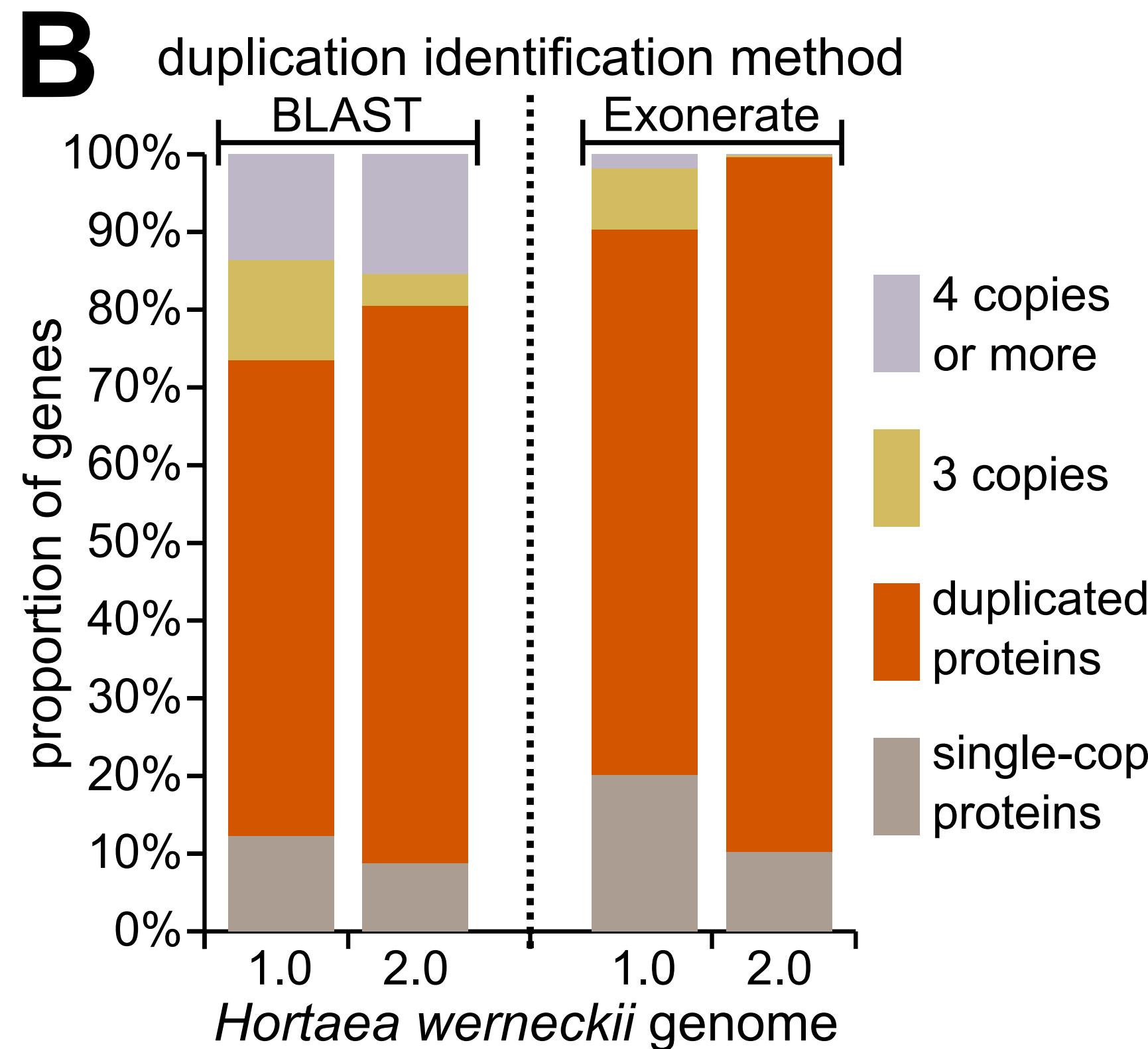
Allomyces macrogynus
(Blastocladiomycota Chytrid)



HORTAEA WERNECKII - BLACK YEAST

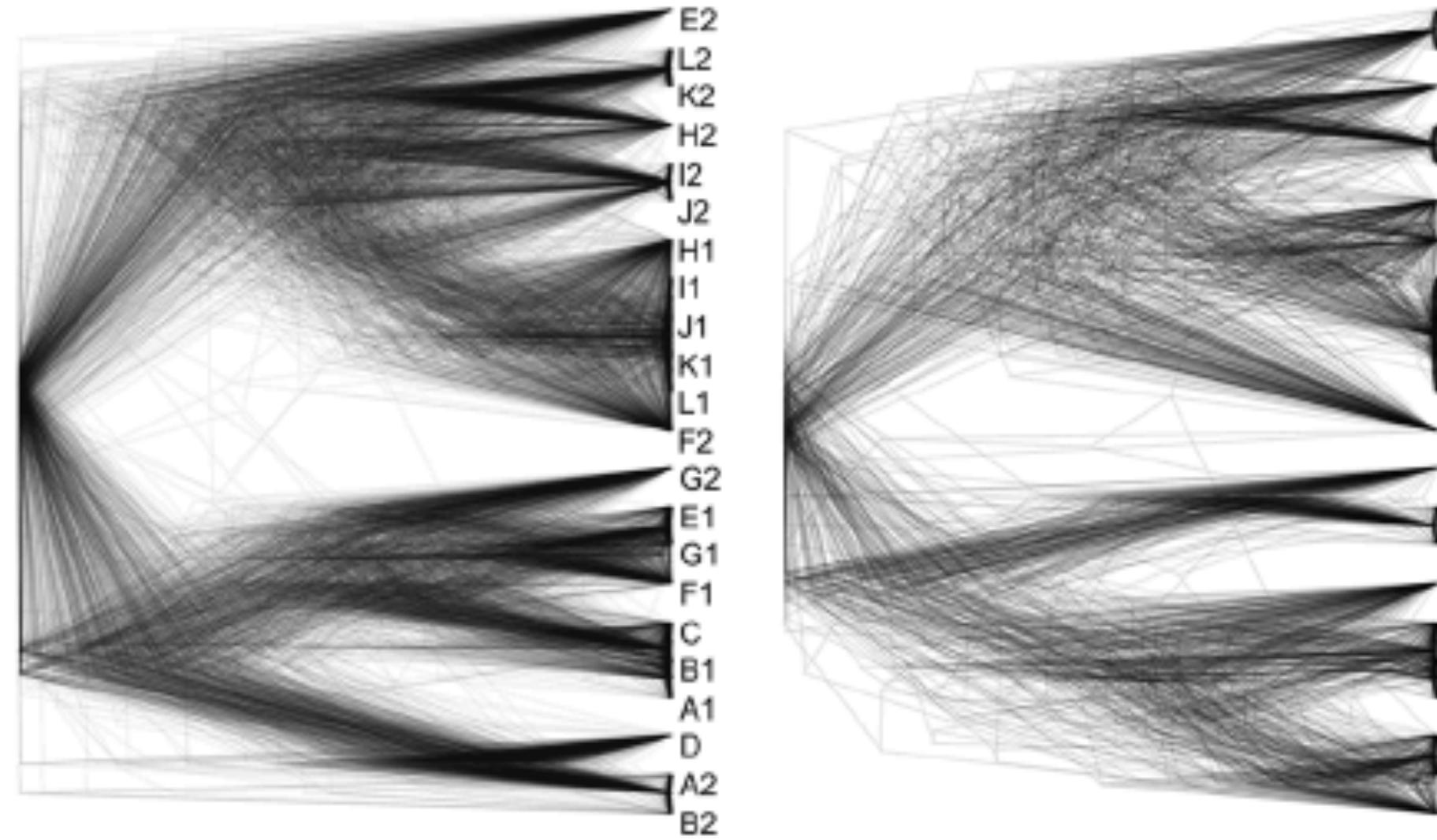


Recent fungal whole genome duplications- *Hortaea werneckii*



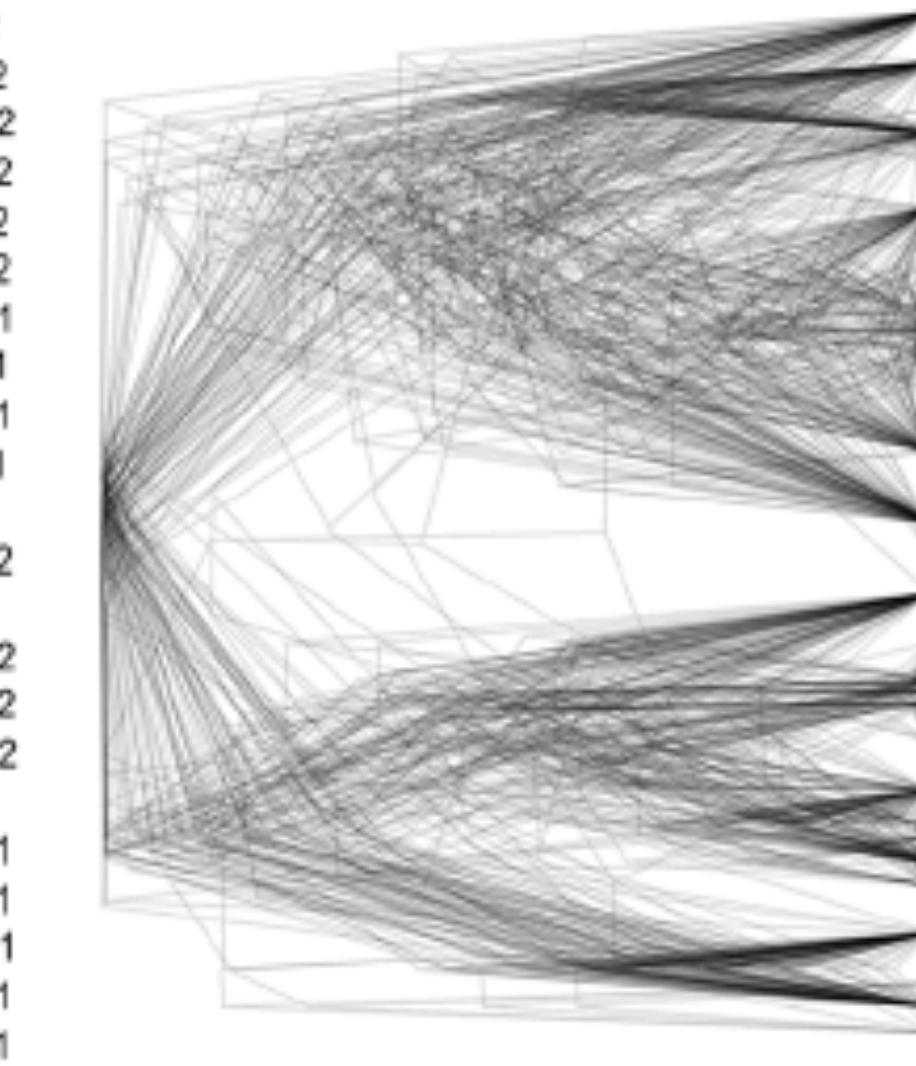
HORTAEA WERNECKII HISTORY IS MORE COMPLICATED THAN A WGD

A 743 trees (55%)



193 trees (14%)

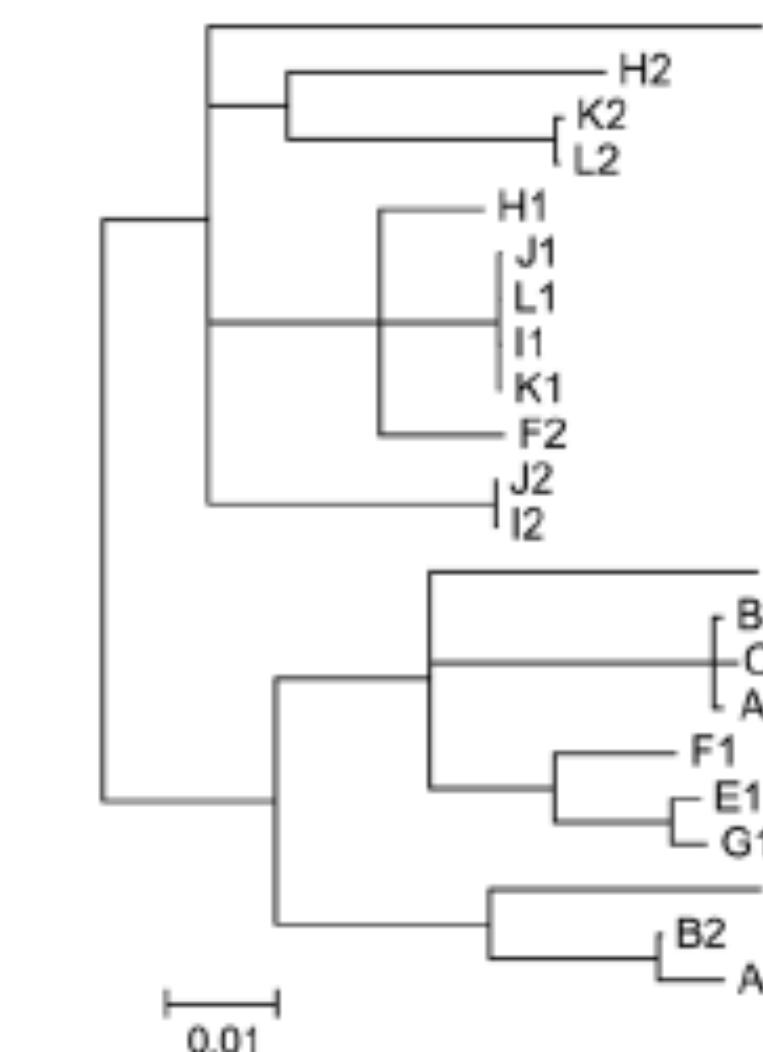
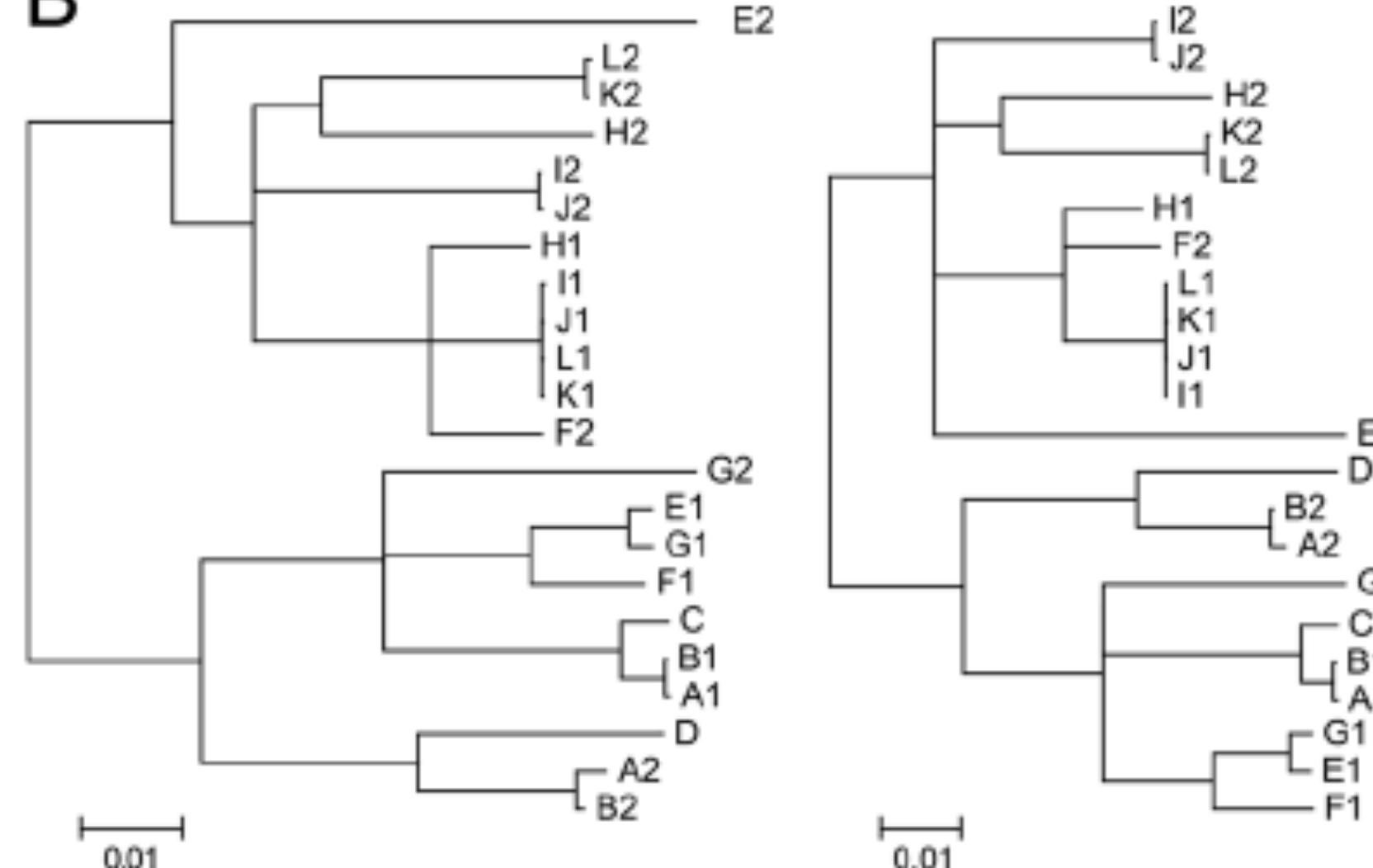
97 trees (7%)



Examining 11 more strains

Multiple evolutionary histories among genomic loci

B

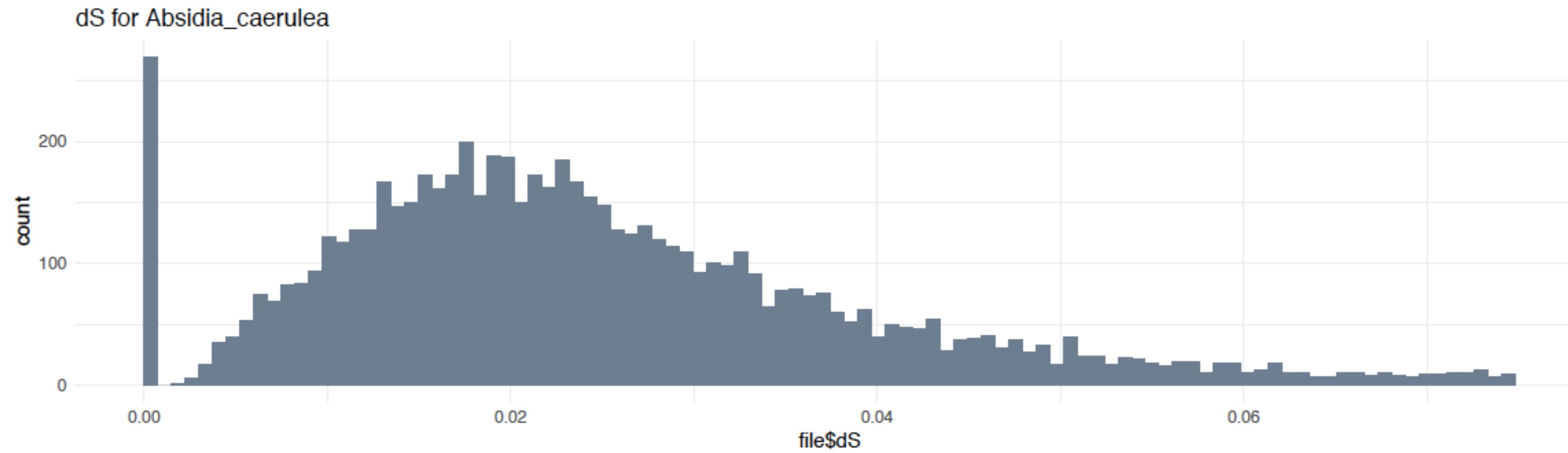


Evidence that some strains are haploid and some are diploid. First genome sequenced was diploid and where haplotypes had sufficient divergence they assembled into separate scaffolds

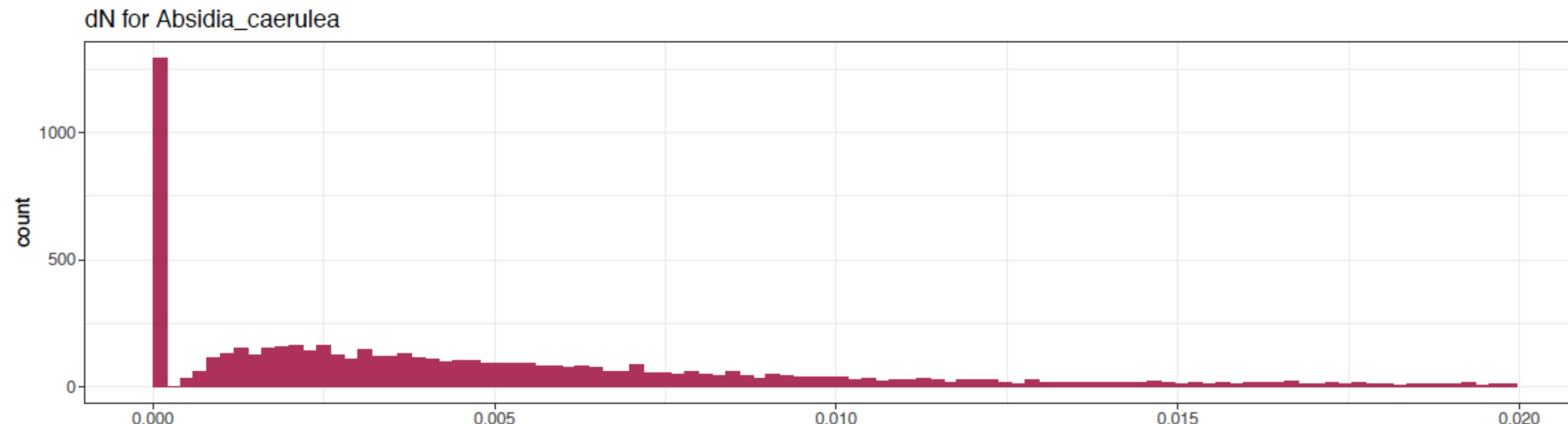
HOW TO COMPUTE KS PLOTS AND TEST FOR DUPLICATIONS?

- All vs All clustering of proteins (or gene DNA sequence)
- Identify 1:1 paralogs, 3 or more members of a gene family are going to be hard to resolve so ignore these as the hypothesis for a WGD is it creates a lot of duplicate pairs
- Align proteins and project in to coding sequence alignment to get codon alignment
- calculate Ks (and Ka if you like) for all pairs.
- PAML, YN00 are suitable. I also wrote a fast tool based on YN00 that generates simple tabular output - <https://github.com/hyphaltip/subopt-kaks/>
- yn00_cds_prealigned or yn00_cds_optimal (let it align coding sequences for you)

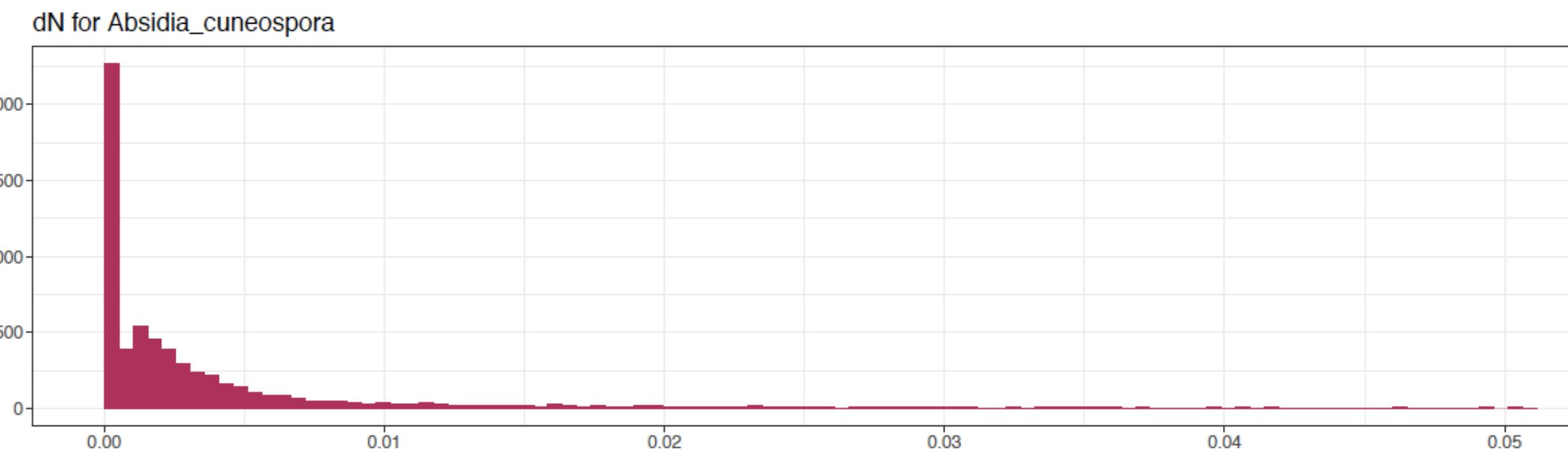
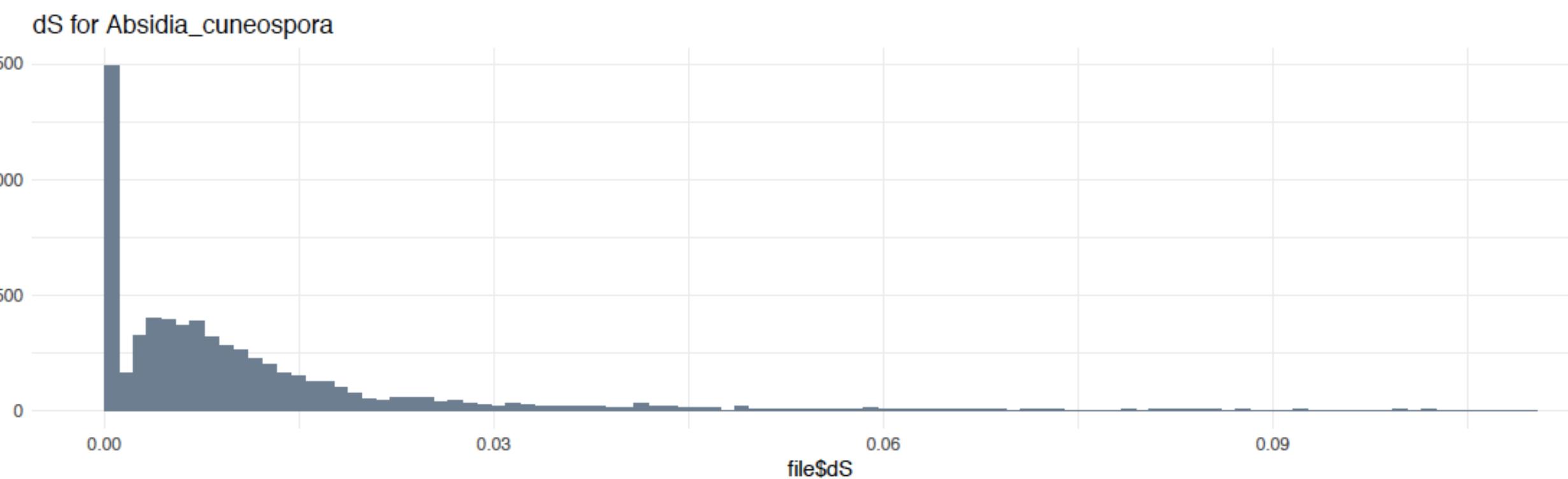
USEFUL INVESTIGATING BETWEEN SPECIES DIVERGENCES



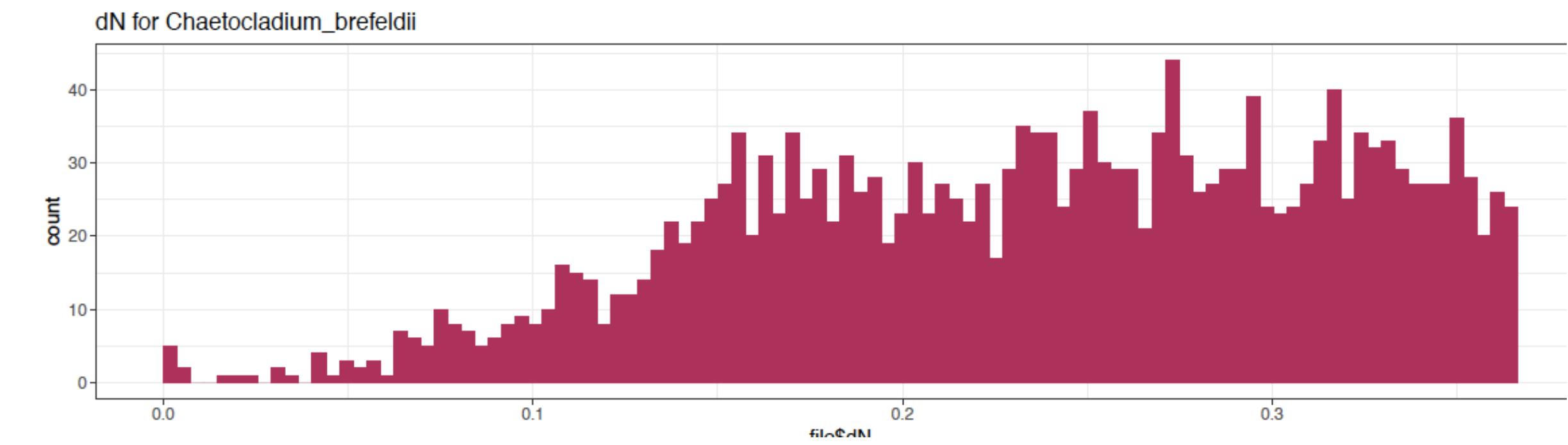
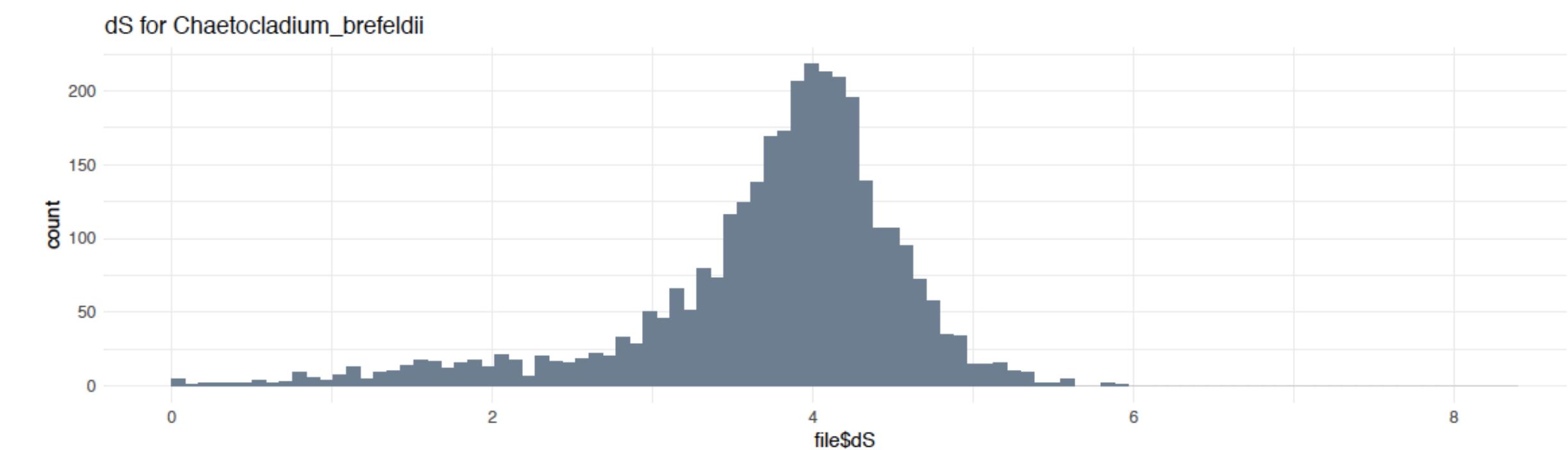
Comparing 2 strains



USEFUL INVESTIGATING BETWEEN SPECIES DIVERGENCES



Comparing 2 strains



Comparing 2 species

DATA SOURCES

- Primary sequence data
 - Genome assembly, annotation, raw sequencing reads - GenBank (NCBI,EMBL,DDBJ)
- Curated, organized, managed
 - Model organism databases (flybase.org, yeastgenome.org, wormbase.org)
 - Collected databases of genomes EnSEML, EuPathDB/FungiDB (ensembl.org, fungidb.org, eupathdb.org)
- Data providers and integrated tools
 - US Joint Genome Institute Mycocosm, Phytozome, IMG (Bacteria)
 - Seq center websites/databases (Genoscope, WUSTL)

DATA FORMATS

Fasta

```
>sequenceID  
AGAGCATAT
```

GFF - Generic feature format

9 columns, tab delimited

chrom, source, type, start, end, score, strand, frame, Group

Genbank: Sequence + annotation

LOCUS MSJE01000001.1 430834 bp DNA linear PLN 03-JAN-2019
DEFINITION Aspergillus terreus strain IMV 01167.
ACCESSION
VERSION
KEYWORDS .
SOURCE Aspergillus terreus
ORGANISM Aspergillus terreus
Eukaryota; Fungi; Dikarya; Ascomycota; Pezizomycotina;
Eurotiomycetes; Eurotiomycetidae; Eurotiales; Aspergillaceae;
Aspergillus.
REFERENCE 1 (bases 1 to 430834)
AUTHORS Stajich,J.E.
TITLE Annotation of genomes of fungal isolates from surfaces in
International Space Station
JOURNAL Unpublished
REFERENCE 2 (bases 1 to 430834)
AUTHORS Stajich,J.E.
TITLE Direct Submission
JOURNAL Submitted (03-JAN-2019) Plant Pathology and Microbiology,
University of California-Riverside, 900 University Ave, Riverside,
CA 92521, USA
COMMENT 'Annotated using funannotate v1.5.1'.
FEATURES Location/Qualifiers
source 1..430834
/organism="Aspergillus terreus"
/mol_type="genomic DNA"
/strain="IMV 01167"
/db_xref="taxon:33178"
gene complement(<14..1332)
/locus_tag="BS087_000011"
mRNA complement(join(<14..151,209..1006,1071..1332))
/product="hypothetical protein"
CDS complement(join(<14..151,209..1006,1071..1332))

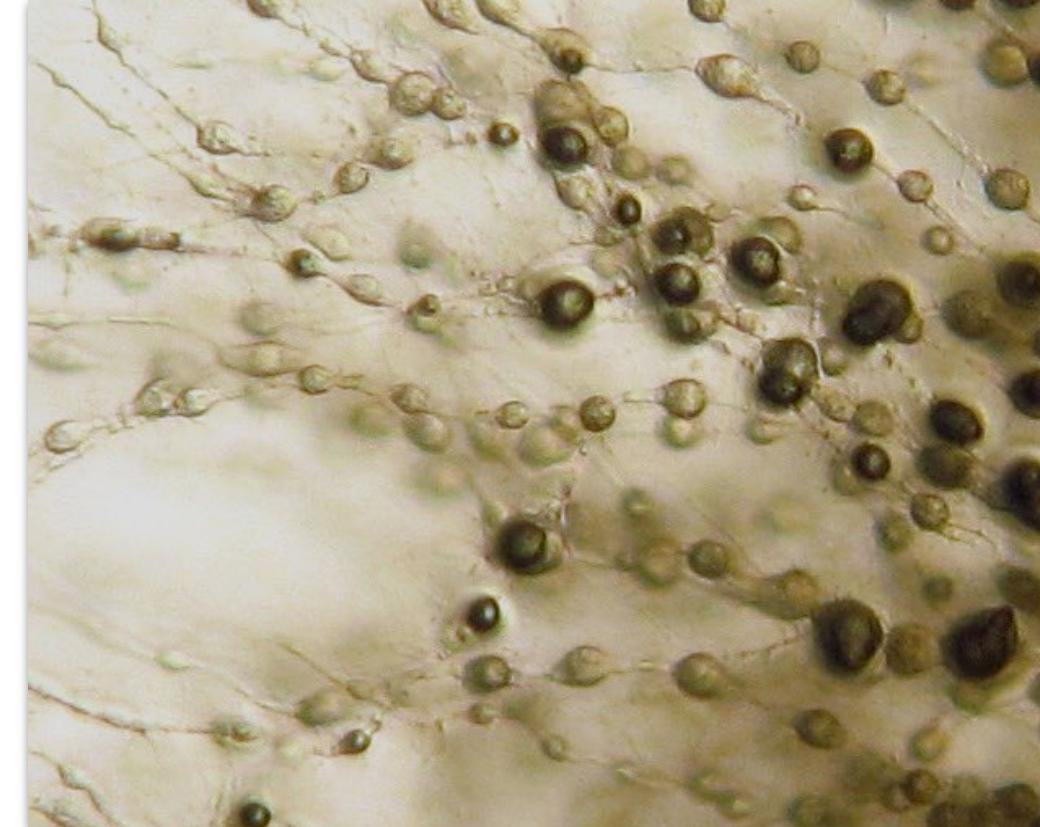
GenBank format

GFF FORMAT

```
MSJE0100001.1  GenBank  gene 14  1332 .  -  .  ID=BS087_000011;
MSJE0100001.1  GenBank  mRNA 14  1332 .  -  .  ID=BS087_000011-T1;Parent=BS087_000011;product=hypothetical protein;
MSJE0100001.1  GenBank  exon 1071 1332 .  -  .  ID=BS087_000011-T1.exon1;Parent=BS087_000011-T1;
MSJE0100001.1  GenBank  exon 209  1006 .  -  .  ID=BS087_000011-T1.exon2;Parent=BS087_000011-T1;
MSJE0100001.1  GenBank  exon 14   151  .  -  .  ID=BS087_000011-T1.exon3;Parent=BS087_000011-T1;
MSJE0100001.1  GenBank  CDS   1071 1332 .  -  0  ID=BS087_000011-T1.cds;Parent=BS087_000011-T1;
MSJE0100001.1  GenBank  CDS   209  1006 .  -  2  ID=BS087_000011-T1.cds;Parent=BS087_000011-T1;
MSJE0100001.1  GenBank  CDS   14    151  .  -  2  ID=BS087_000011-T1.cds;Parent=BS087_000011-T1;
```

GENOME ANNOTATION

- Pipelines to run automatic prediction. Eukaryotic gene annotation harder than Bacteria or Archaea prediction
- expressed RNA (RNASeq, ESTs)
- Comparative data - alignment of proteins and transcripts from other closely related species
- Ab initio gene predictors trained on the organism (Augustus, GeneMark.HMM, etc)
- Data are combined into consensus gene models
- Existing pipelines that I have familiarity with MAKER and Funannotate



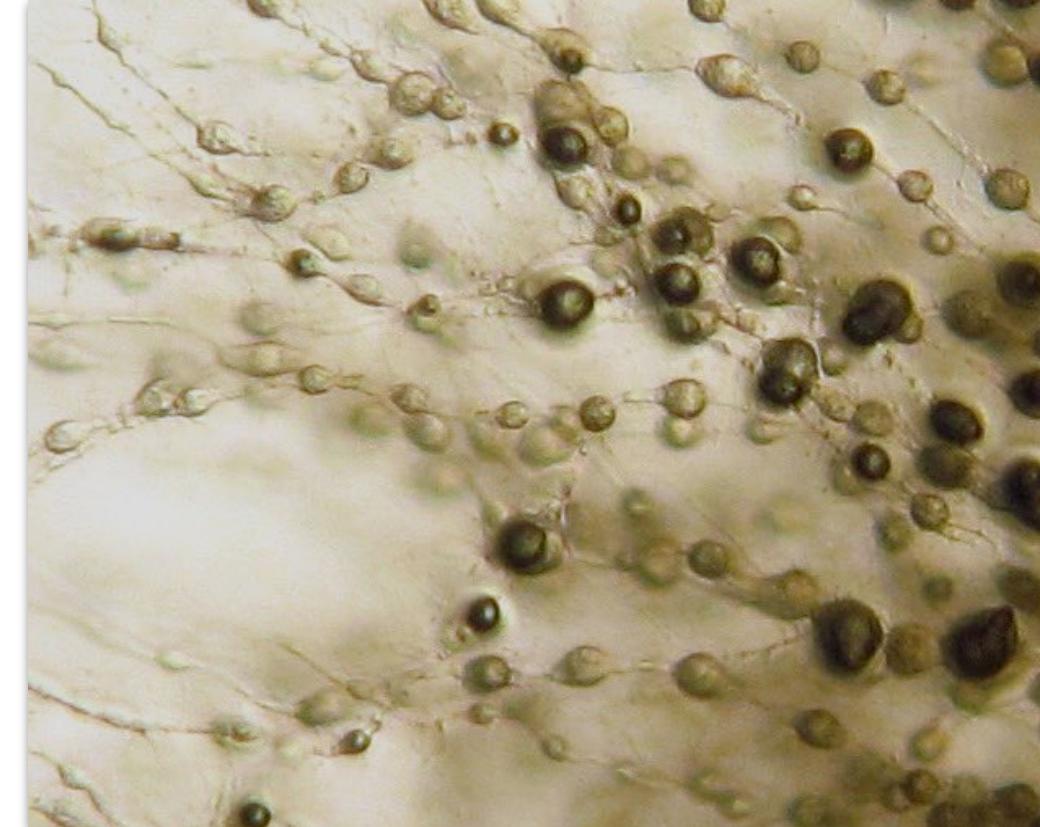
Catenaria (Blastocladiomycota)



Rhizopus (Mucoromycota)

ANN NGUYEN, GREG JEDD
OUSMANE CISSÉ, MINOU NOWROUSIAN, DAVID HEWITT

EVOLUTION OF MULTICELLULARITY IN FUNGI



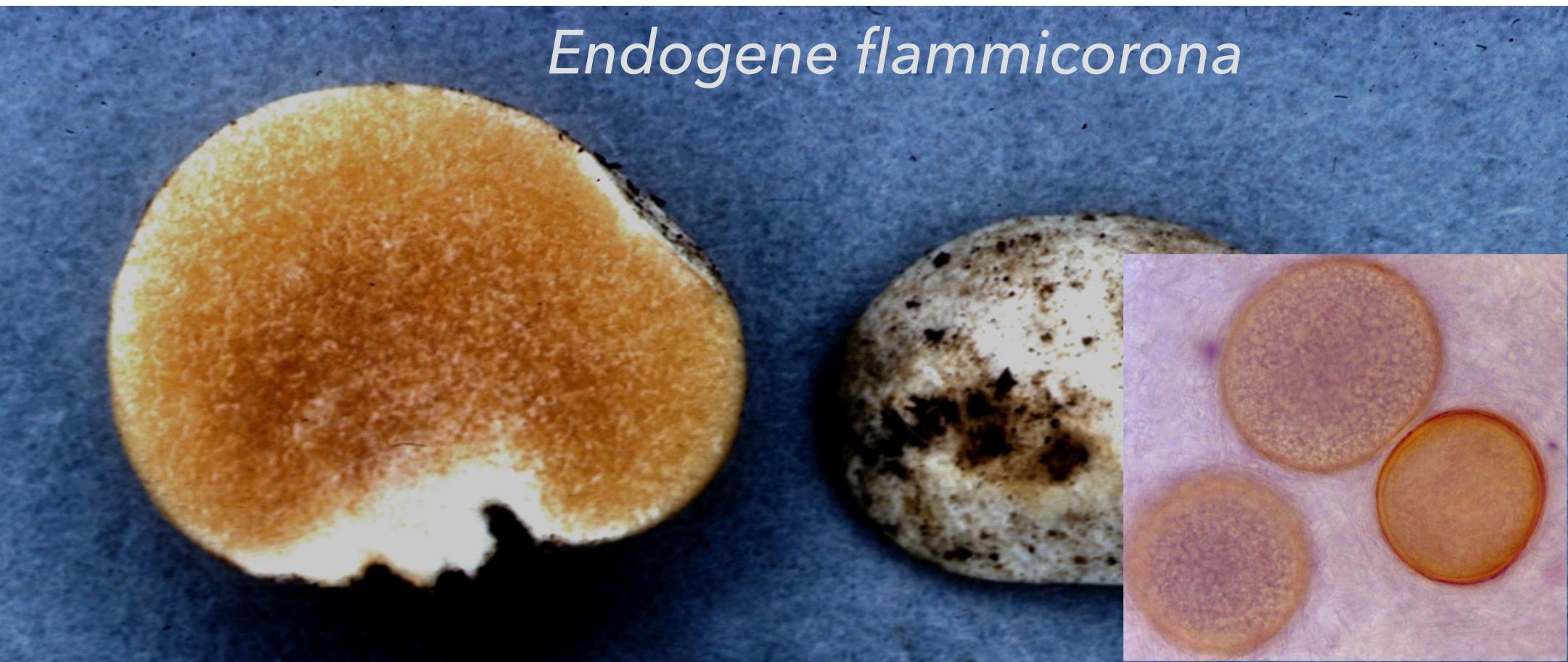
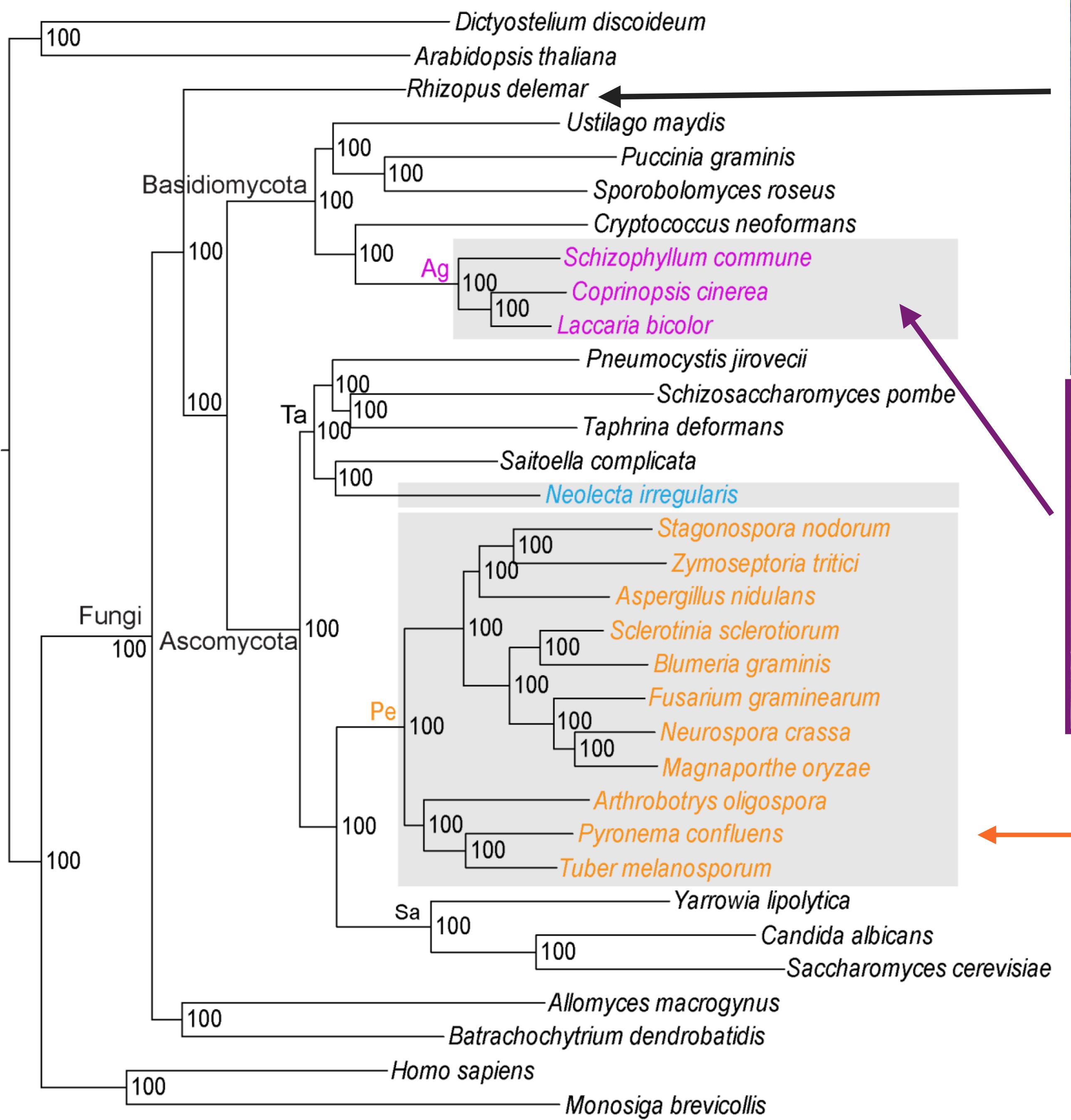
Catenaria (Blastocladiomycota)



Rhizopus (Mucoromycota)

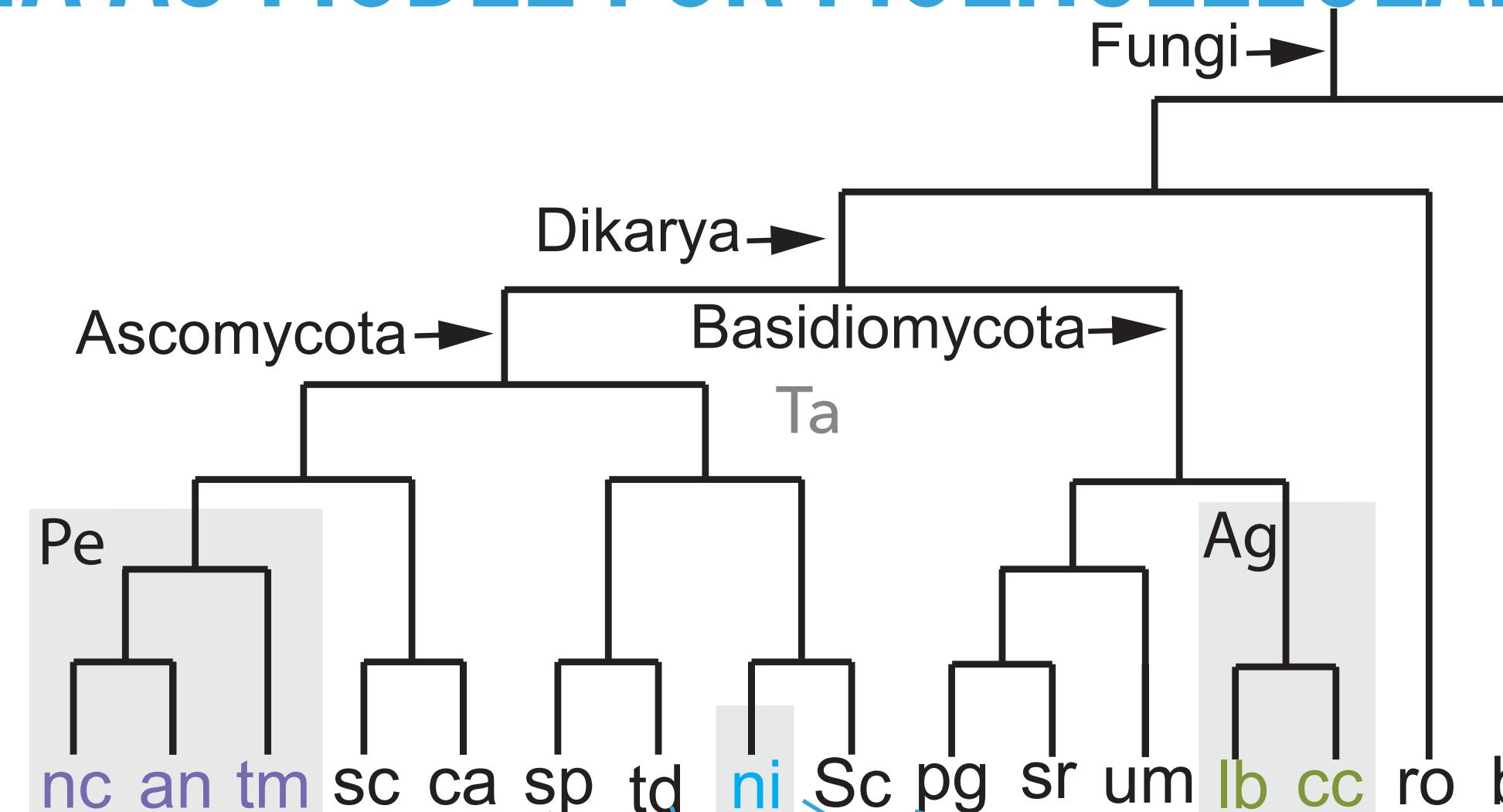
ANN NGUYEN, GREG JEDD
OUSMANE CISSÉ, MINOU NOWROUSIAN, DAVID HEWITT

EVOLUTION OF MULTICELLULARITY IN FUNGI



MULTICELLULARITY

NEOLECTA AS MODEL FOR MULTICELLULAR EVOLUTION

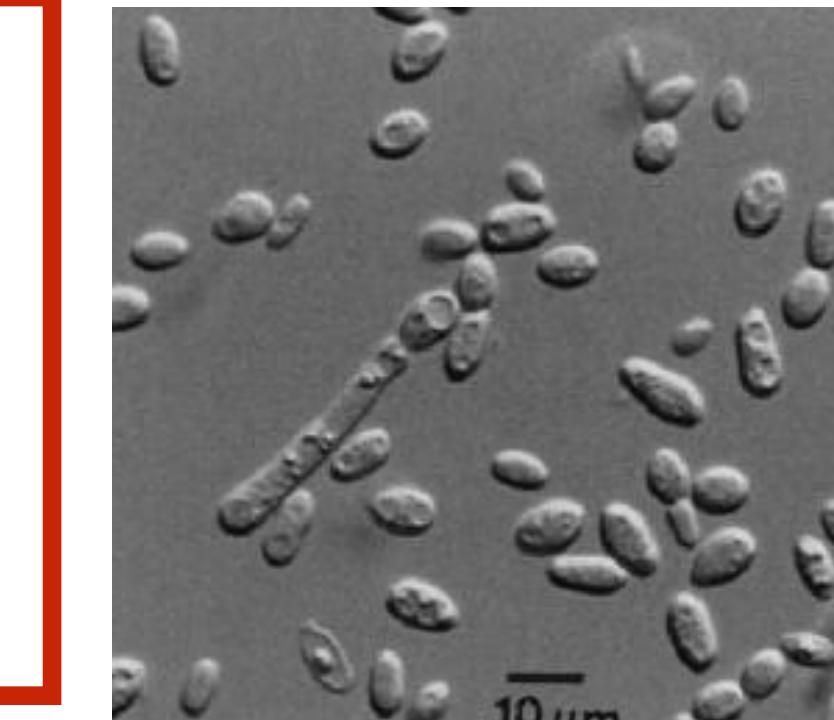
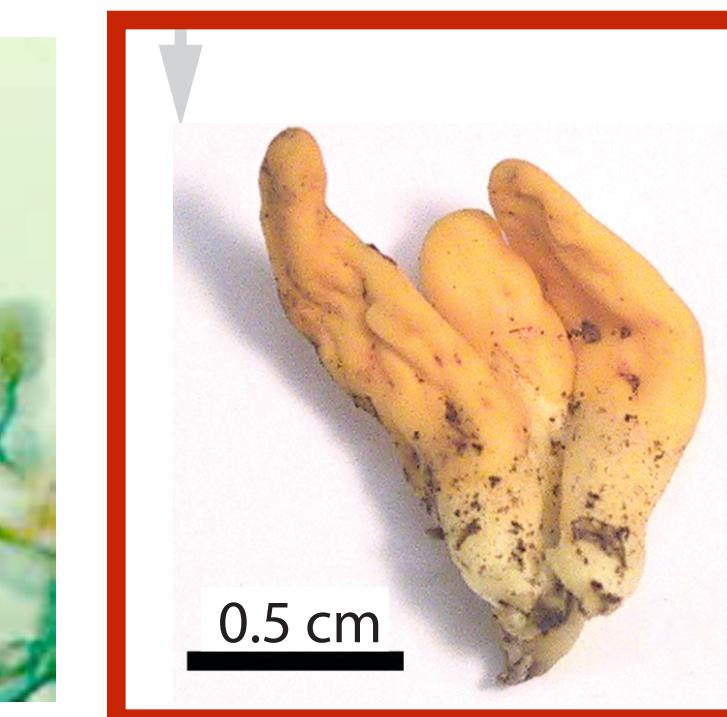
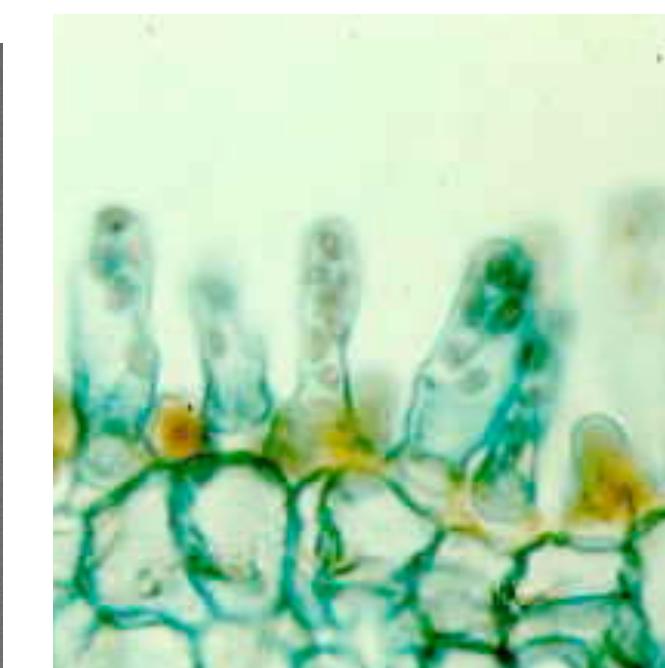
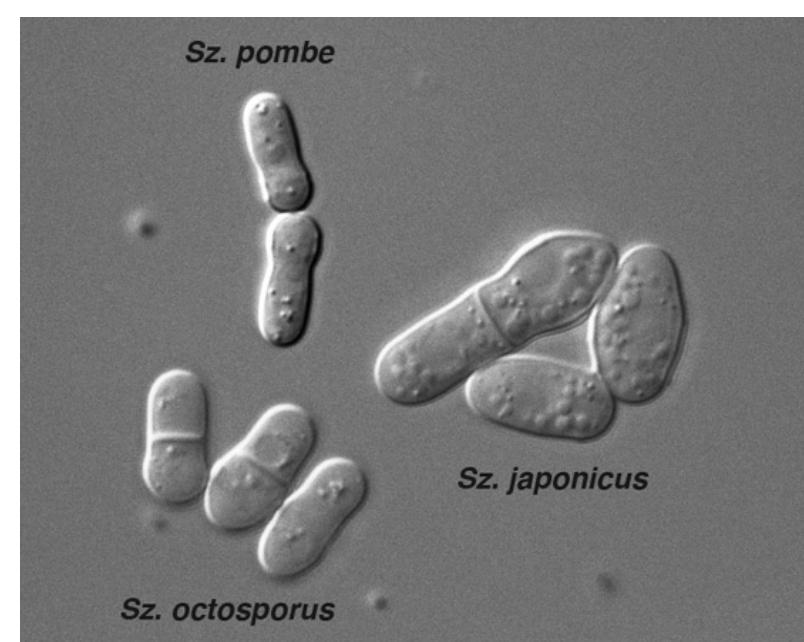


NEOLECTA IRREGULARIS

- * PLANT ROOT ASSOCIATED

- * SO FAR UNCULTURABLE IN THE LAB

- * FORMS SPOROCARPS - FRUITING BODIES



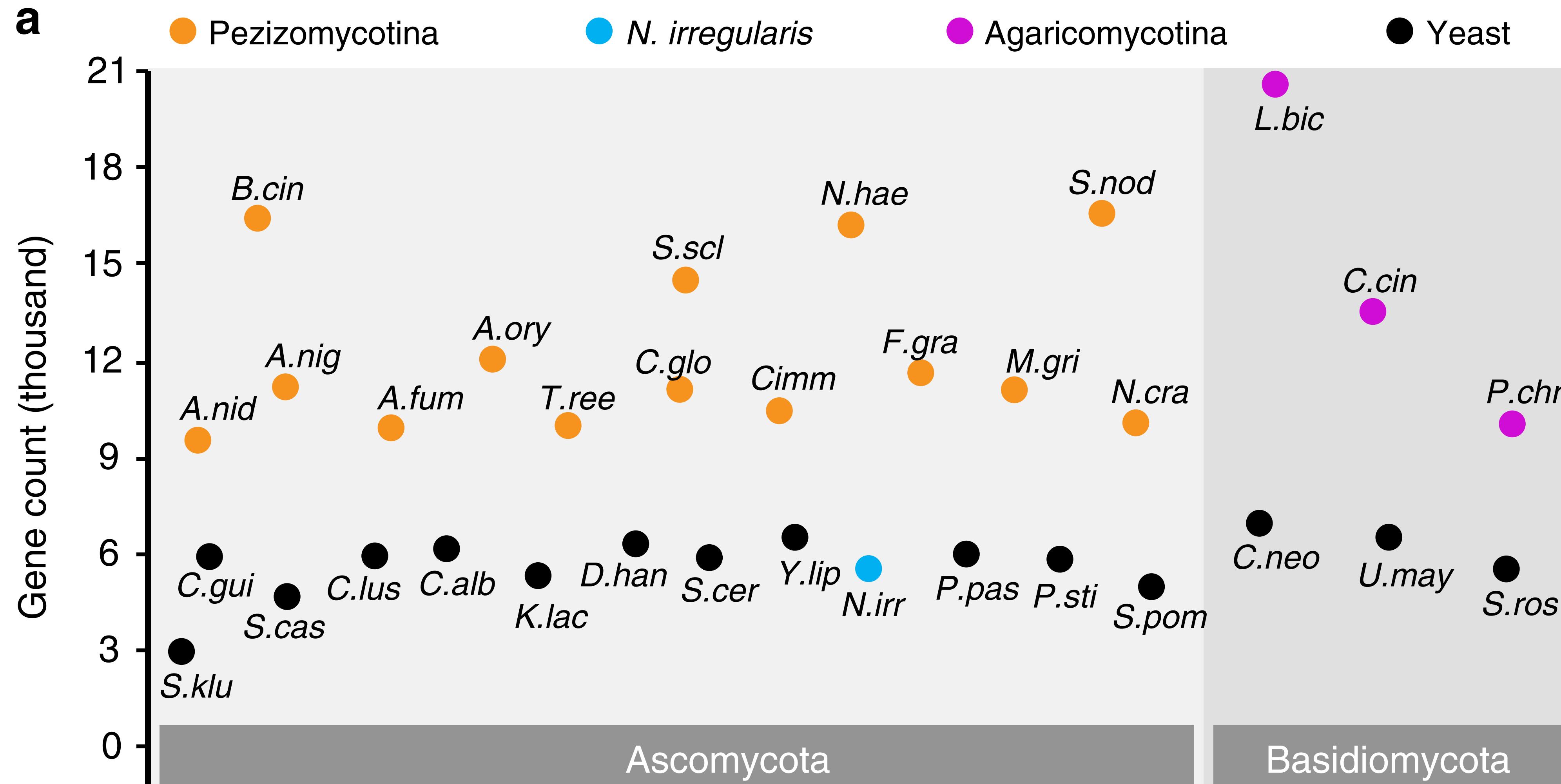
Schizosaccharomyces pombe

TAPHRINA DEFORMANS

NEOLECTA IRREGULARIS

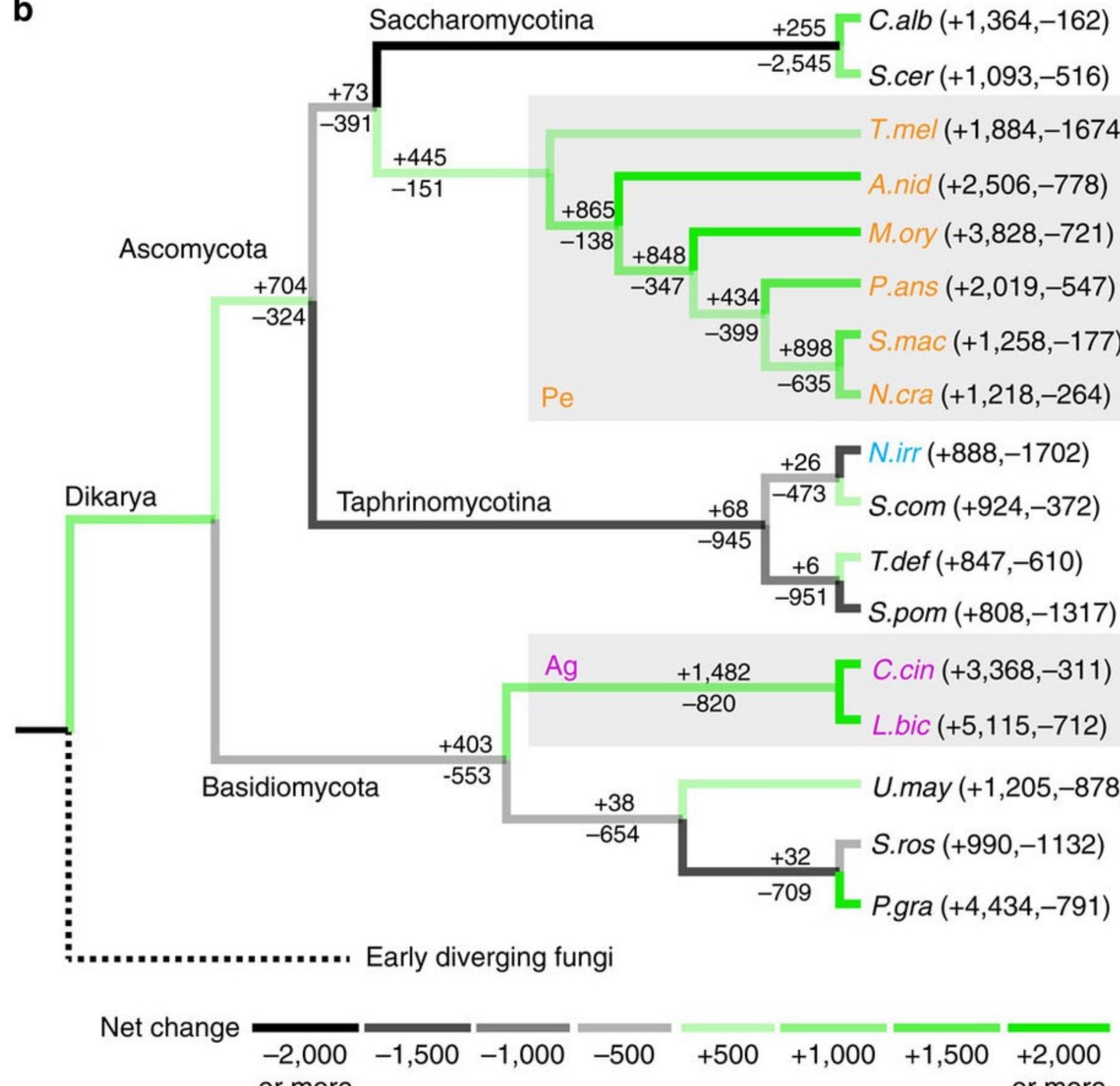
SAITOELLA COMPLICATA

GENOME SIZE DOES NOT PREDICT COMPLEX MULTICELLULARITY



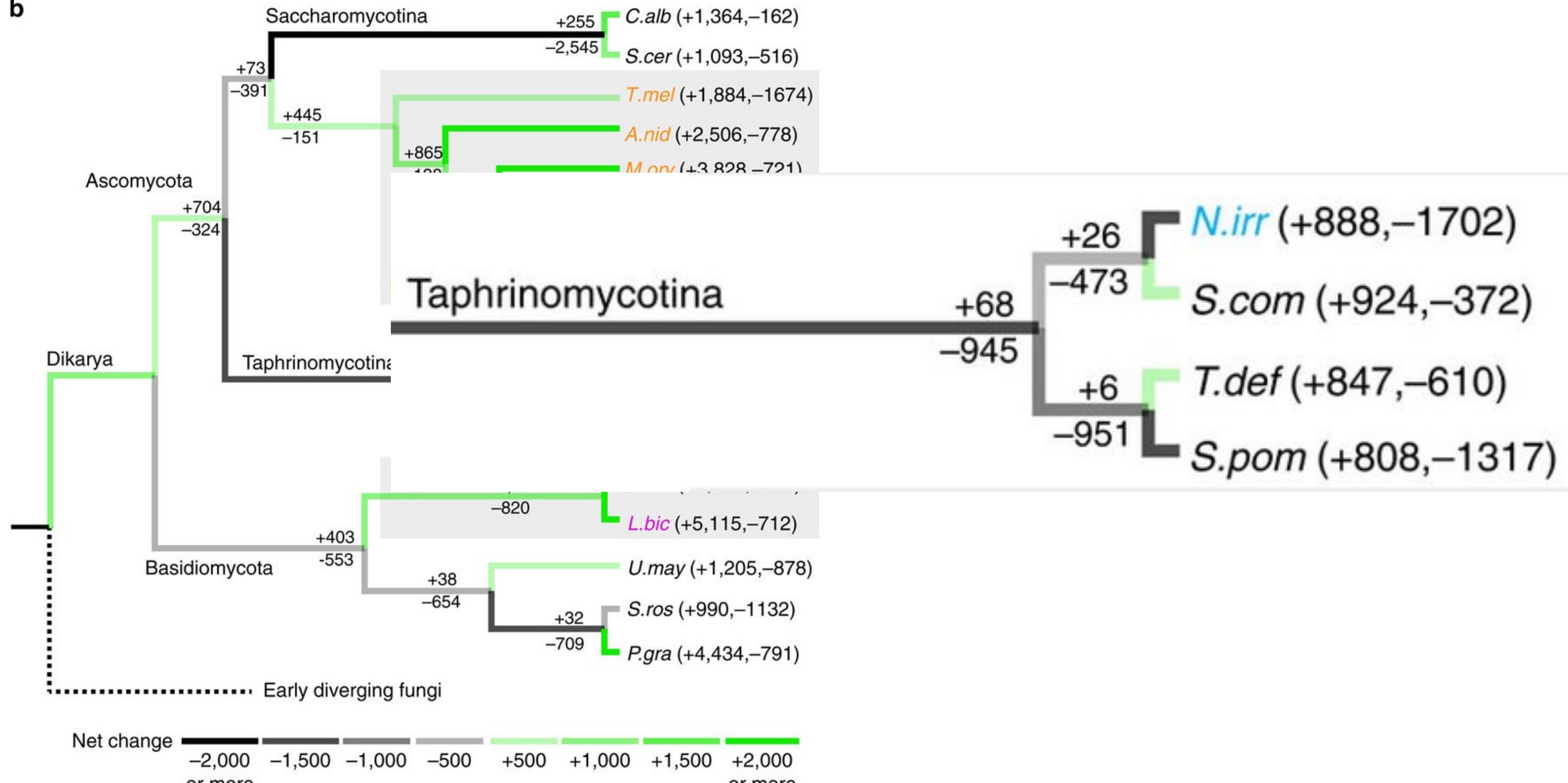
NEOLECTA LINEAGE DID NOT EXPERIENCE LARGE RECENT GAINS OF GENES

b



NEOLECTA LINEAGE DID NOT EXPERIENCE LARGE RECENT GAINS OF GENES

b

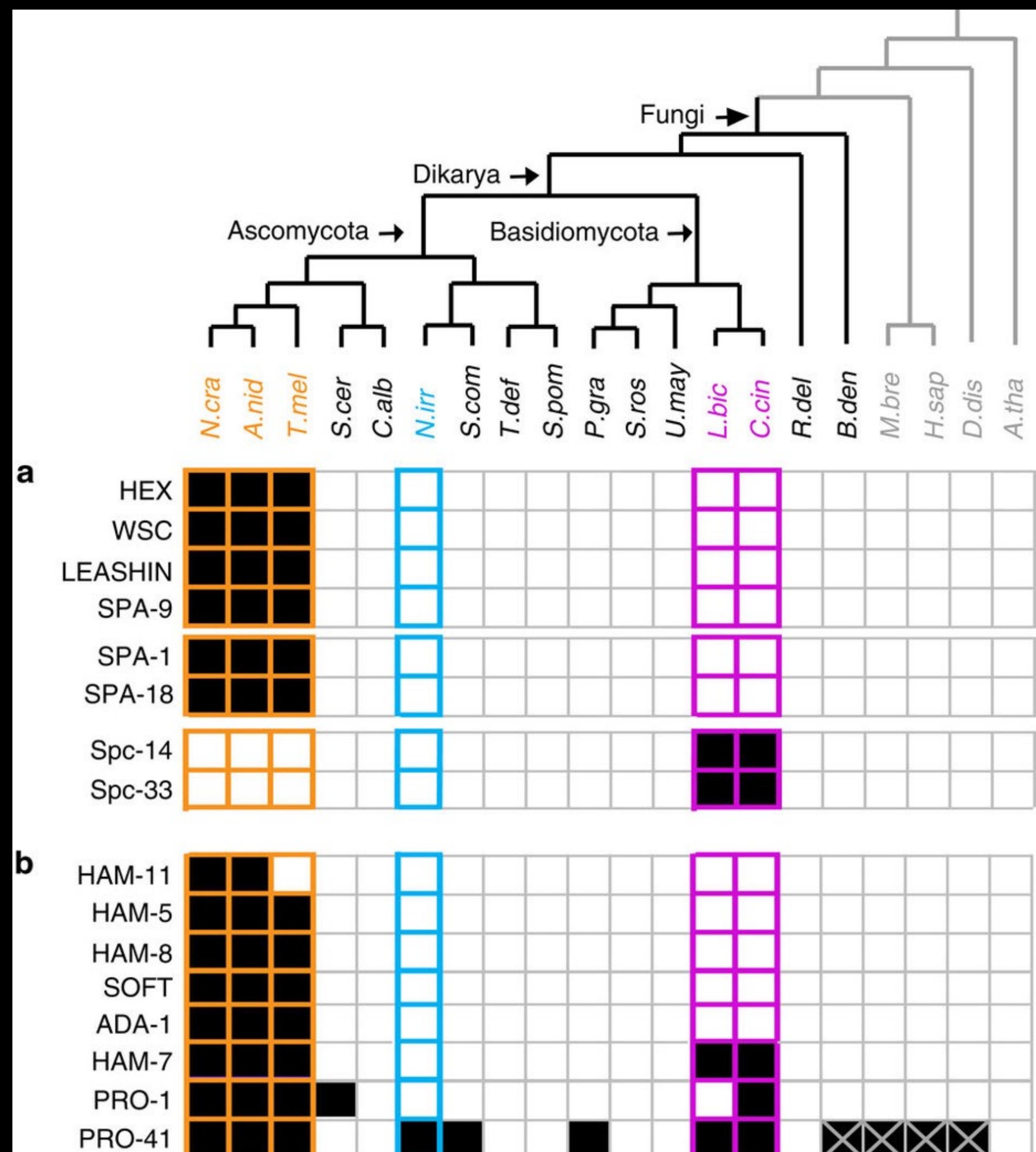


SEARCHING FOR COMPLEX MULTICELLULARITY (CM) SIGNATURES

- Are there genes that fungi that have CM share?
- These fungi make septated hyphae, with gating channels to separate cells.
 - When did these gating mechanisms evolve?
- Hyphal fusion and remodeling needed to make complex multicellular structures (e.g. fruiting bodies)

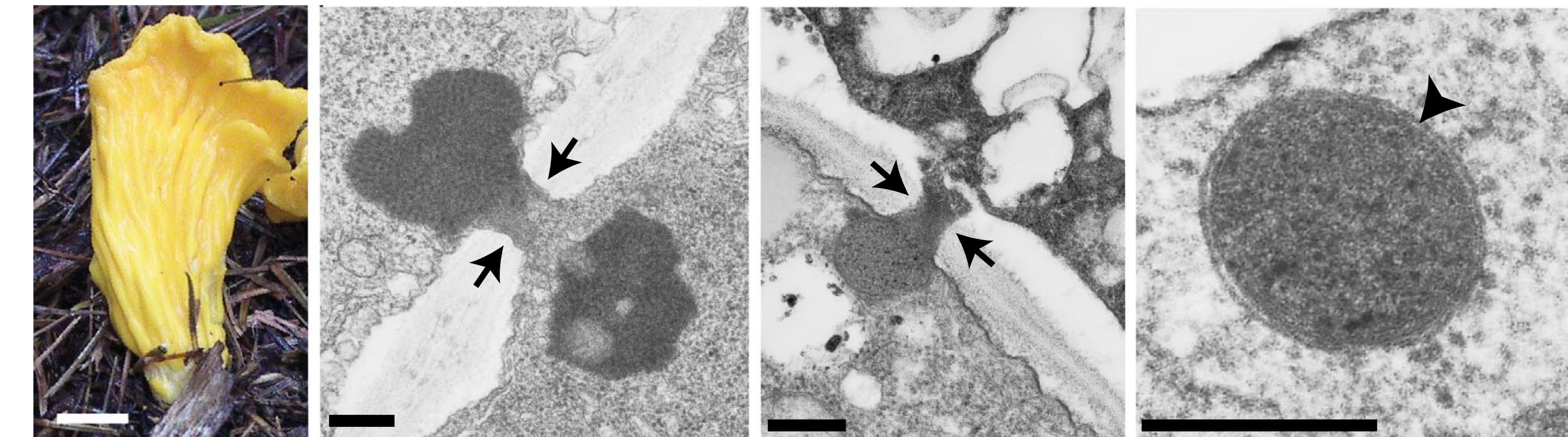
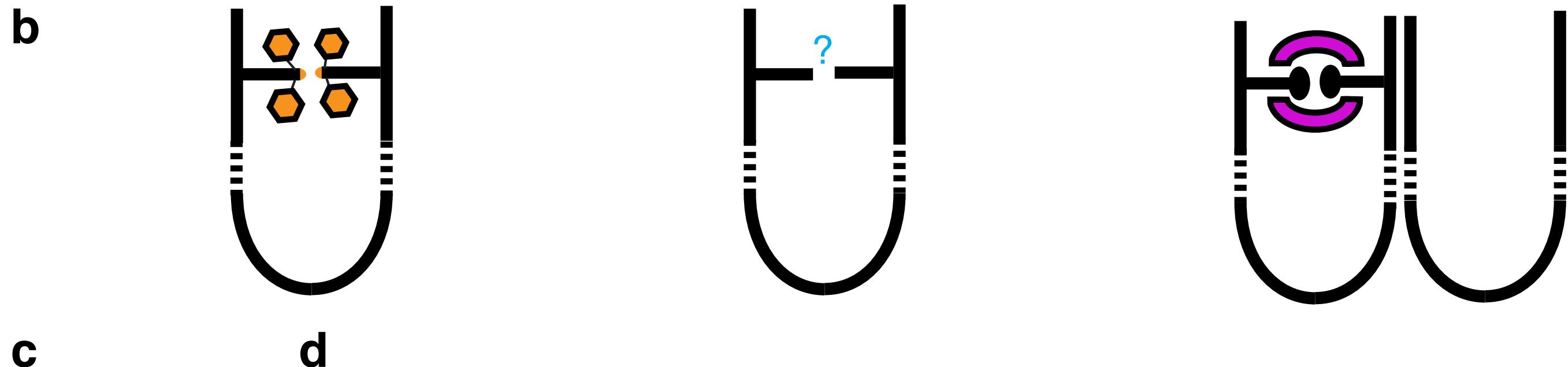
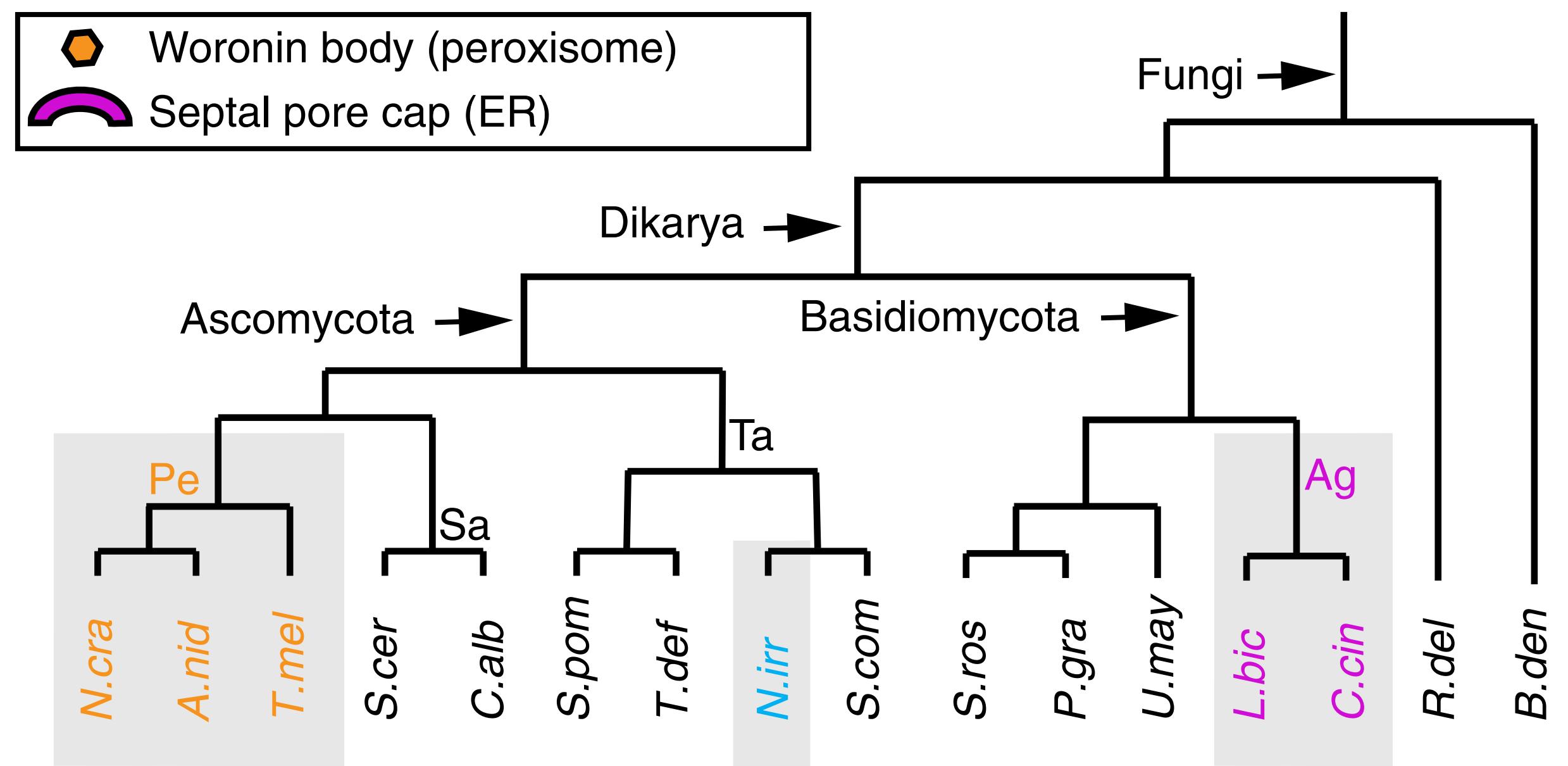
SEARCHING FOR CONSERVED GENES AMONG FUNGI WITH CM

- Search *N. crassa* proteins against collection of fungi
 - Identify those shared among CM fungi or lost in yeasts.
- 1,050 genes found in a CM associated
- 37% are absent and 47% highly divergent in yeasts
- over-represented functional categories:
endomembrane transport and organization (transport routes, substrate transport, peroxisome) and aerobic respiration (electron transport and redox-related enzymes)

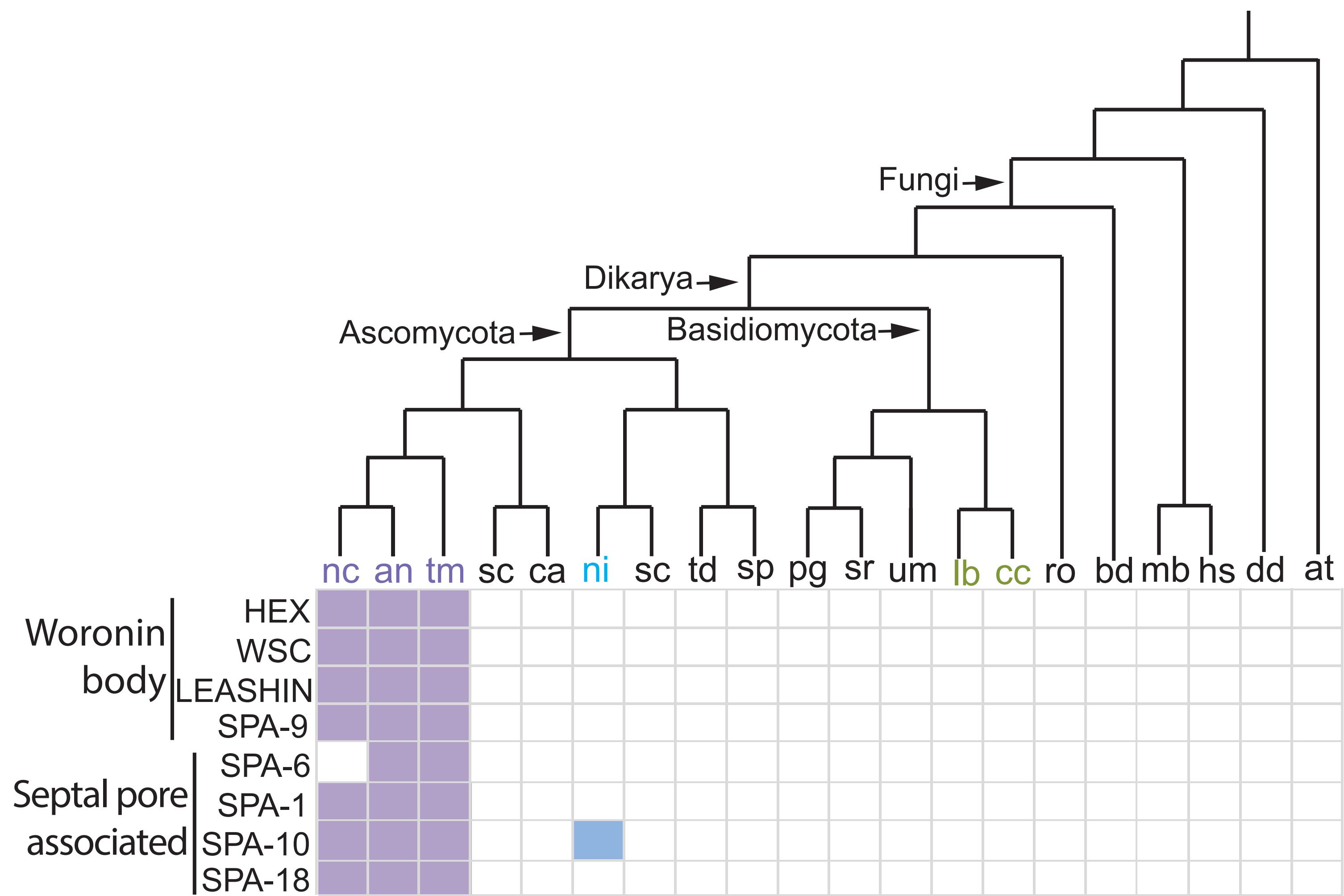


MECHANISMS FOR SEPTAL PLUGGING VARIES ACROSS THE FUNGI

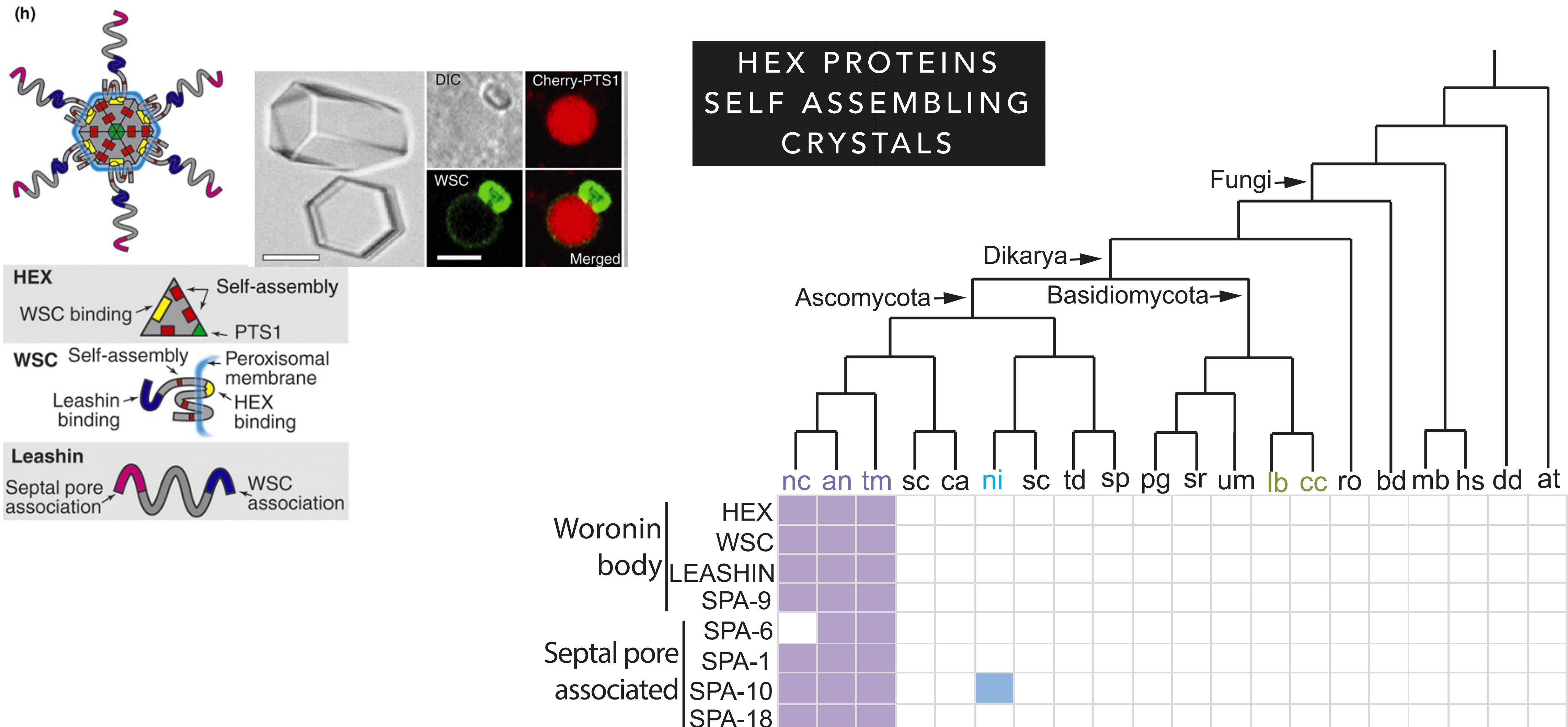
ARE MOLECULES SHARED THAT
PERFORM THIS PLUGGING?



NO WORONIN BODYGENES IN NEOLECTA: RESTRICTED TO PEZIZOMYCOTINA

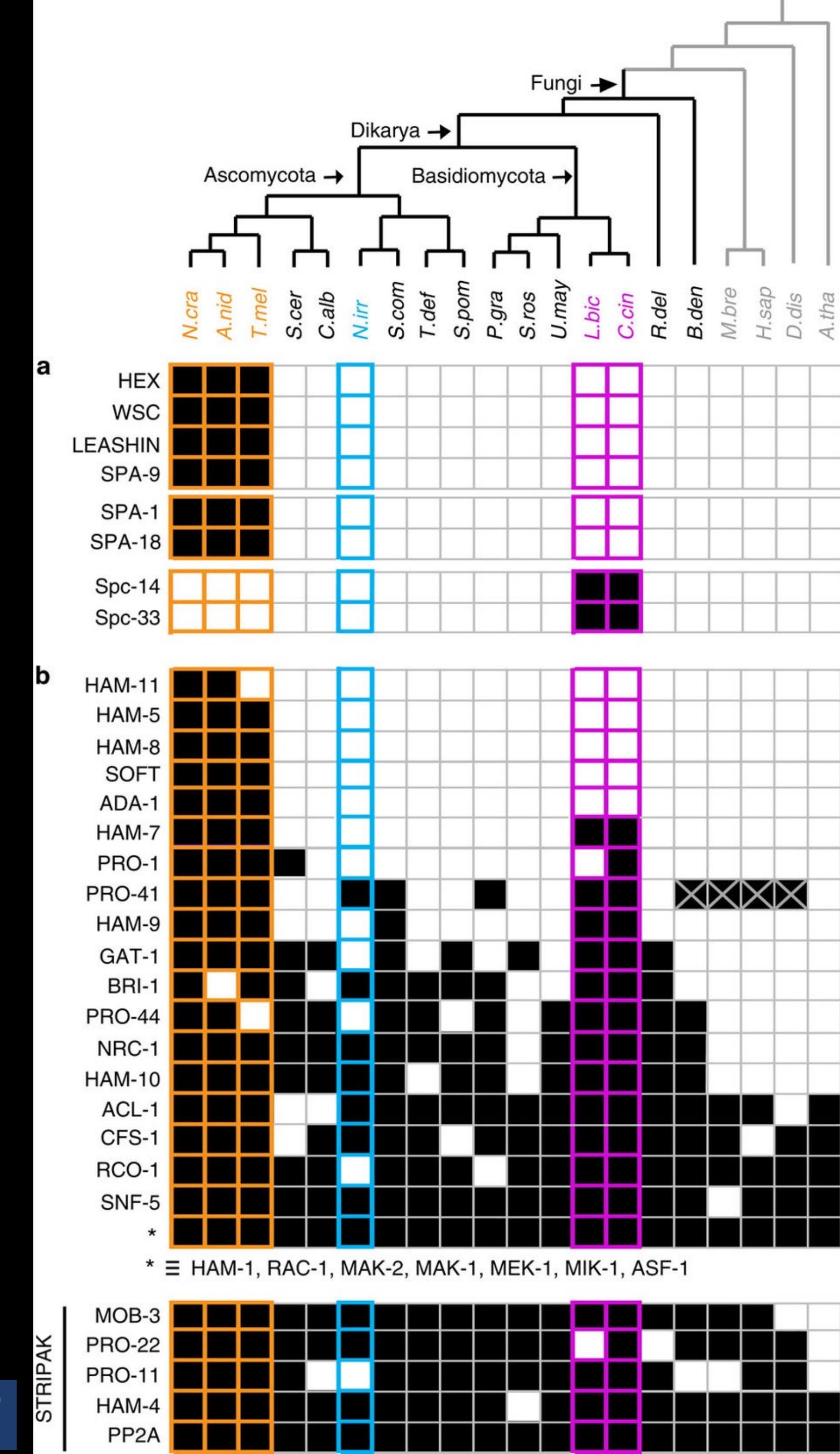


NO WORONIN BODY GENES IN NEOLECTA: RESTRICTED TO PEZIZOMYCOTINA



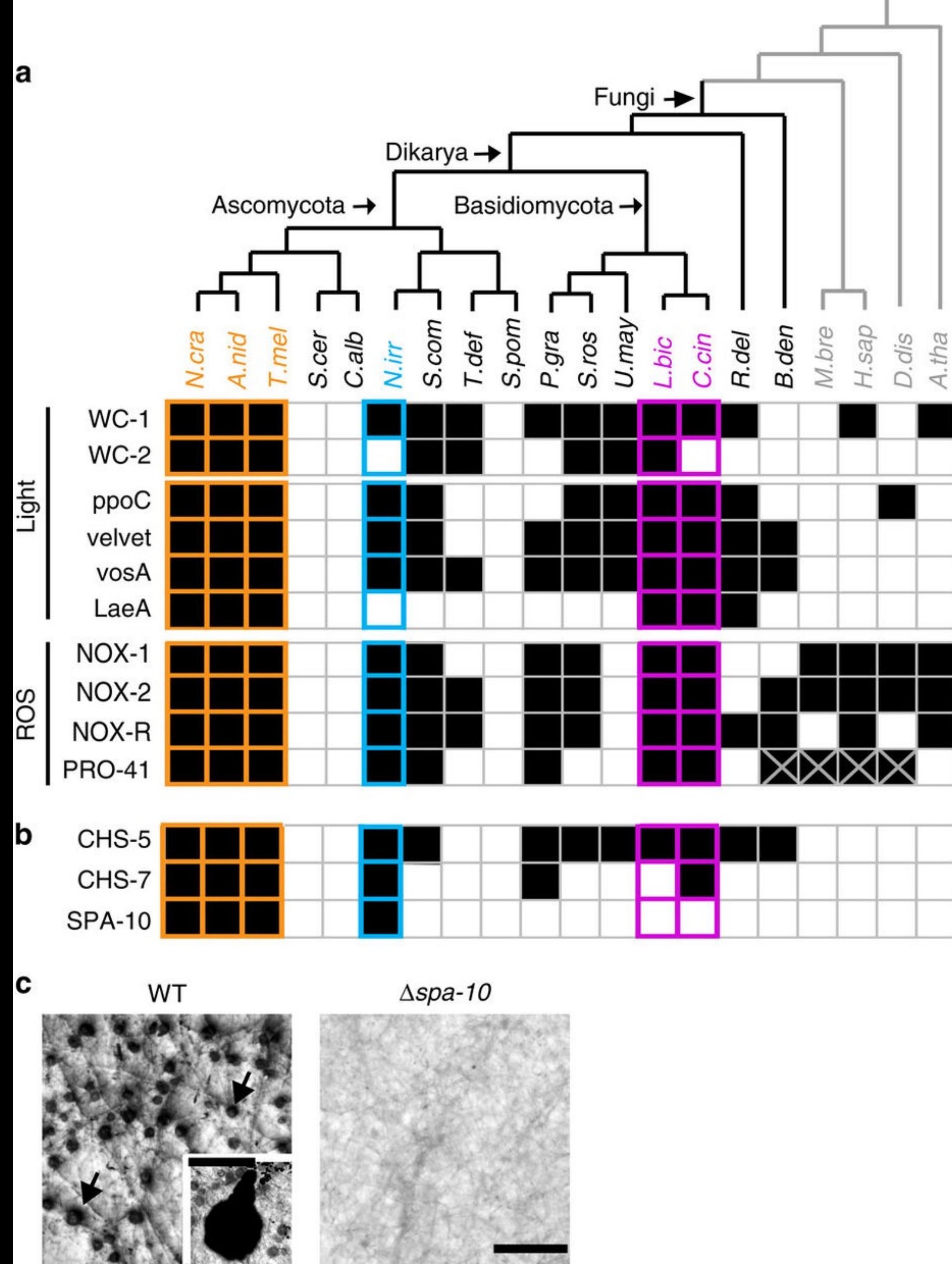
GENE SHARING PATTERNS ACROSS MULTICELLULAR FUNGI

- Septal pore gating and hyphal fusion proteins (a)
 - Pezizomycotina & Agaricomycotina specific proteins for pore gating are missing in Neolecta
 - Some hyphal fusion proteins are missing Neolecta (b)
 - Striatin-interacting phosphatases and kinases (STRIPAK) complex is ancient and conserved throughout. Important in development (c)

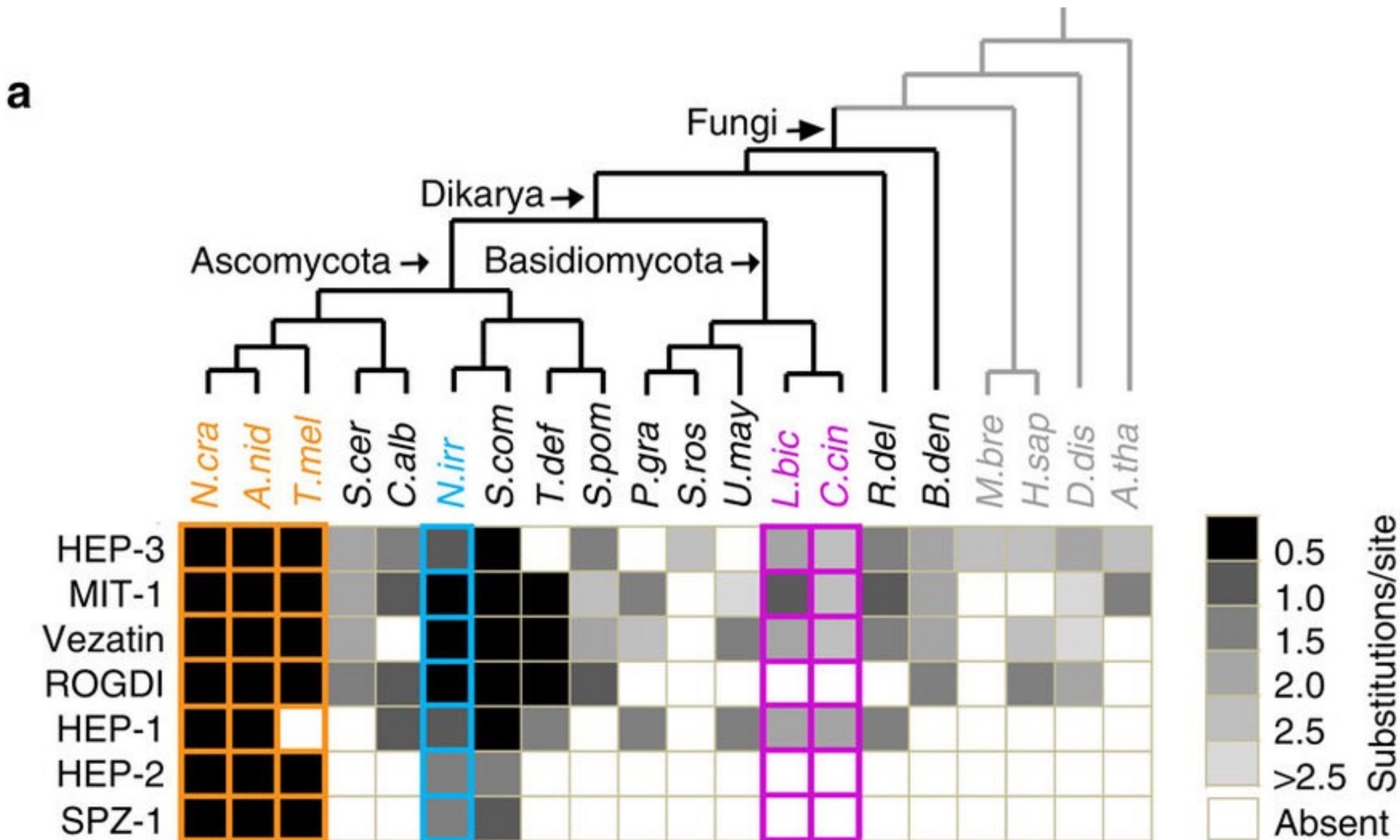


GENES SHARED AMONG SPECIES WITH COMPLEX MORPHOLOGY

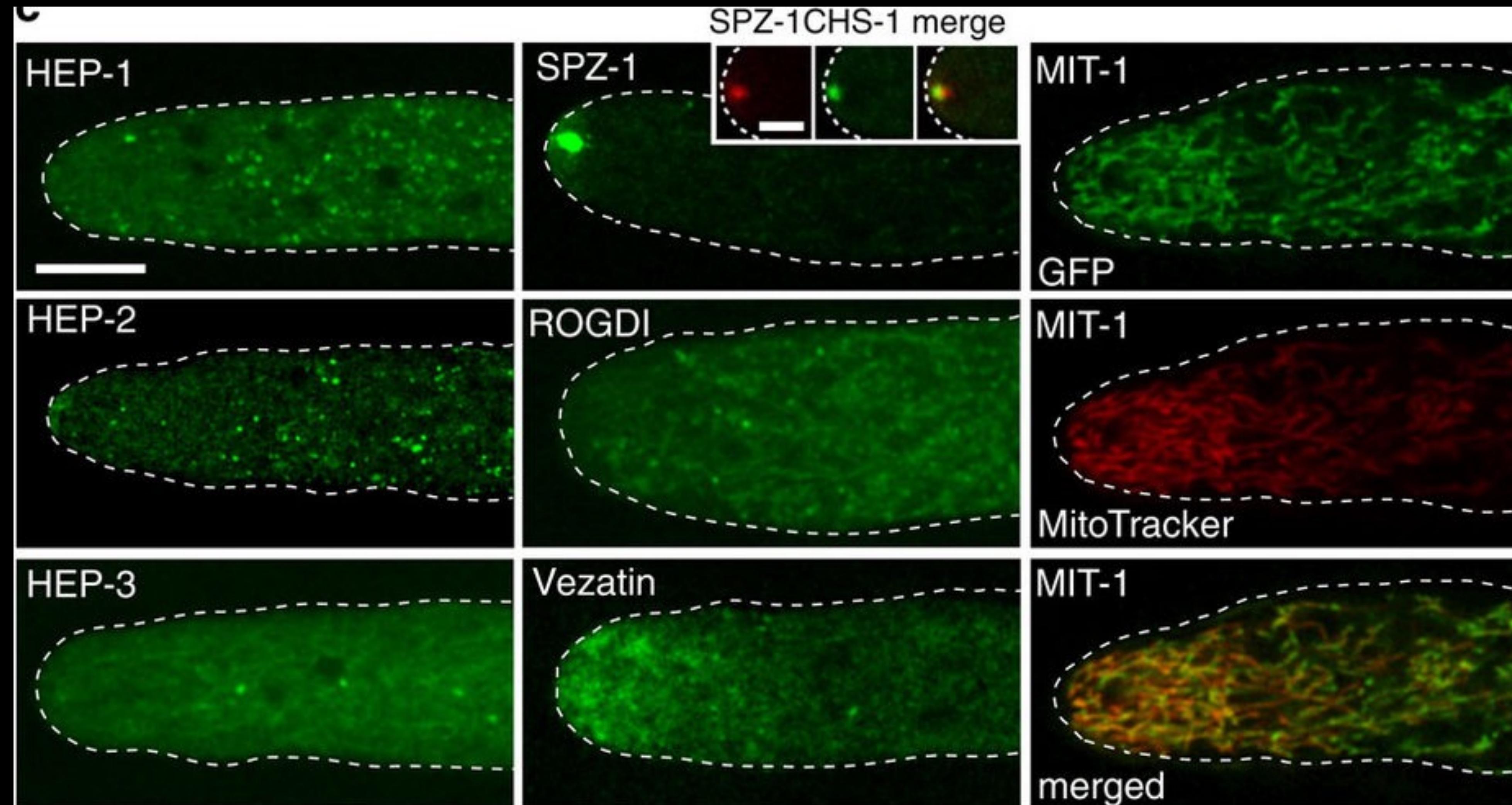
- Proteins important for signaling and morphogenesis
- Reactive Oxygen Species Signaling
- Remodeling and construction
- $\Delta spa-10$ are female sterile



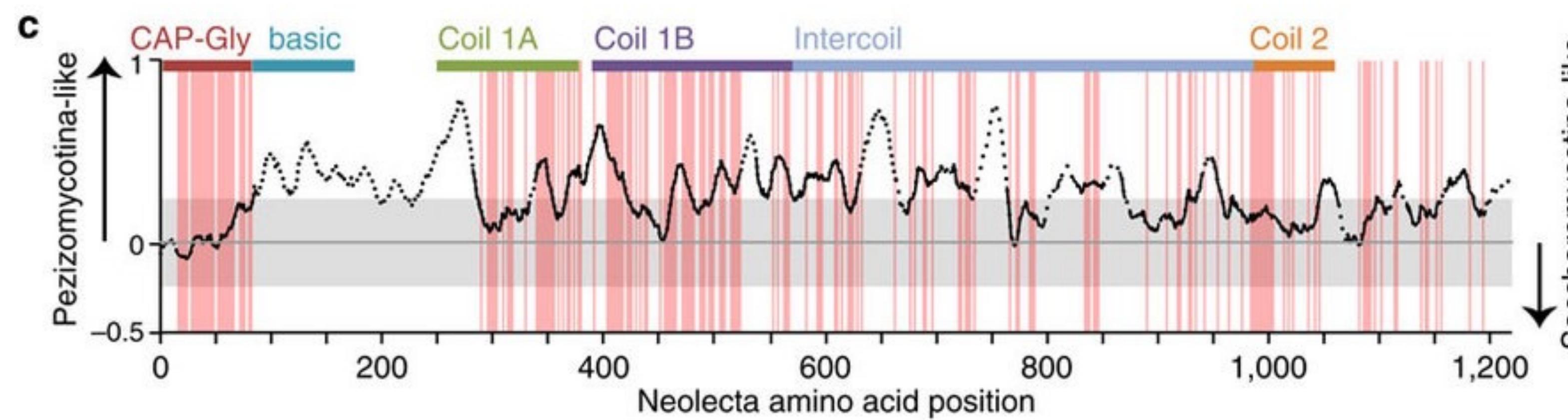
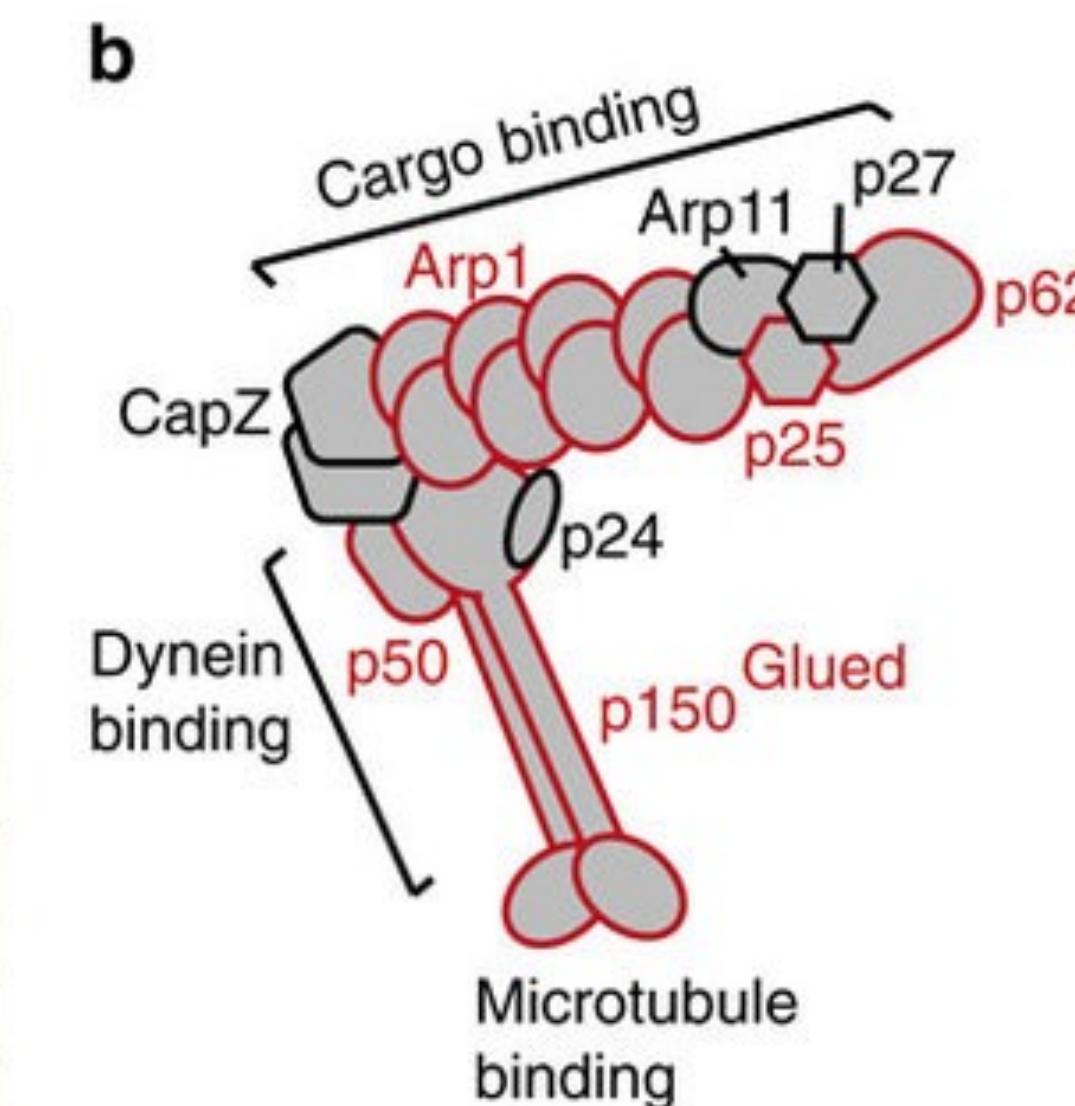
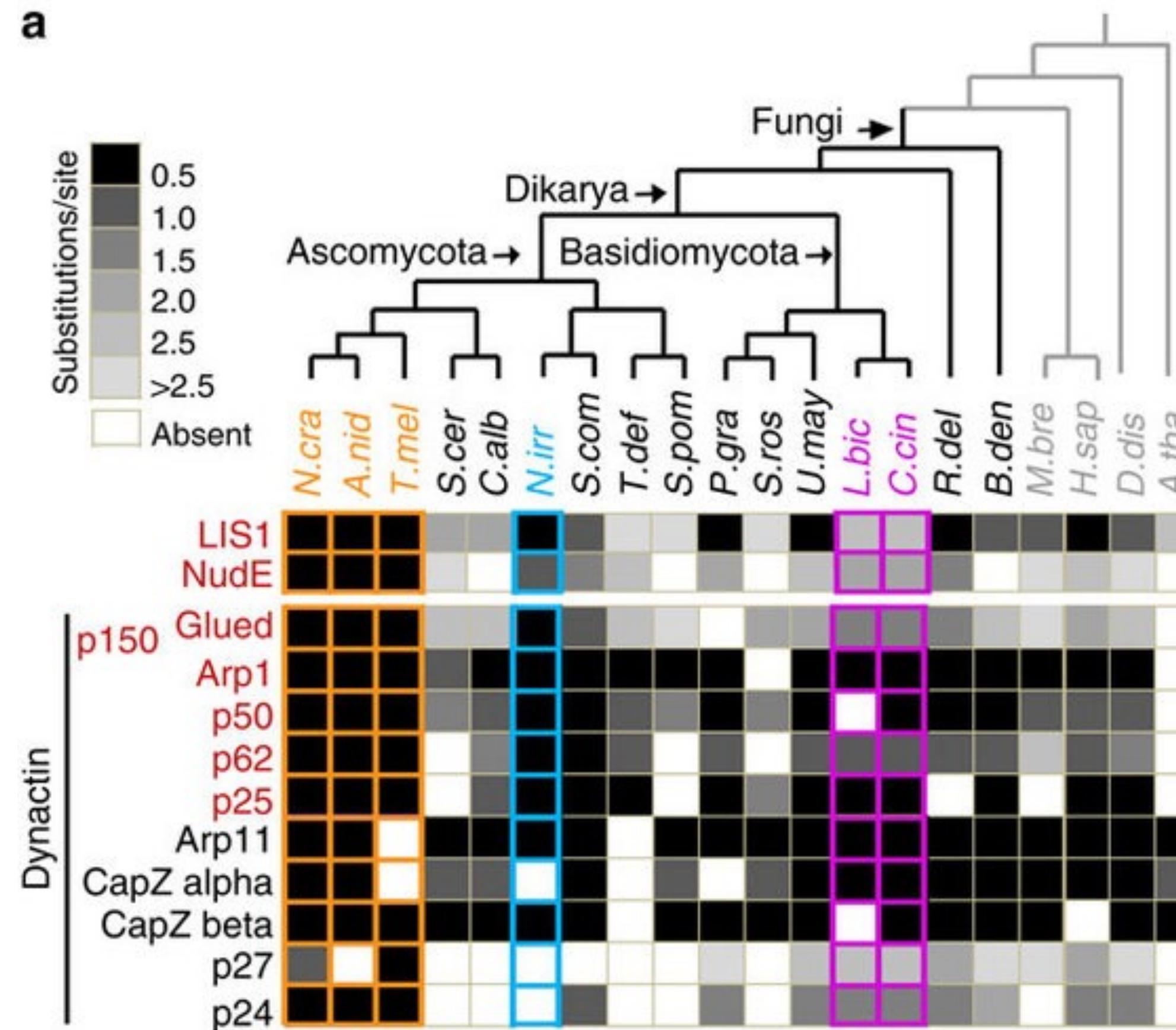
THERE ARE ALSO NOVEL PROTEINS IMPLICATED IN COMPLEX MULTICELLULARITY



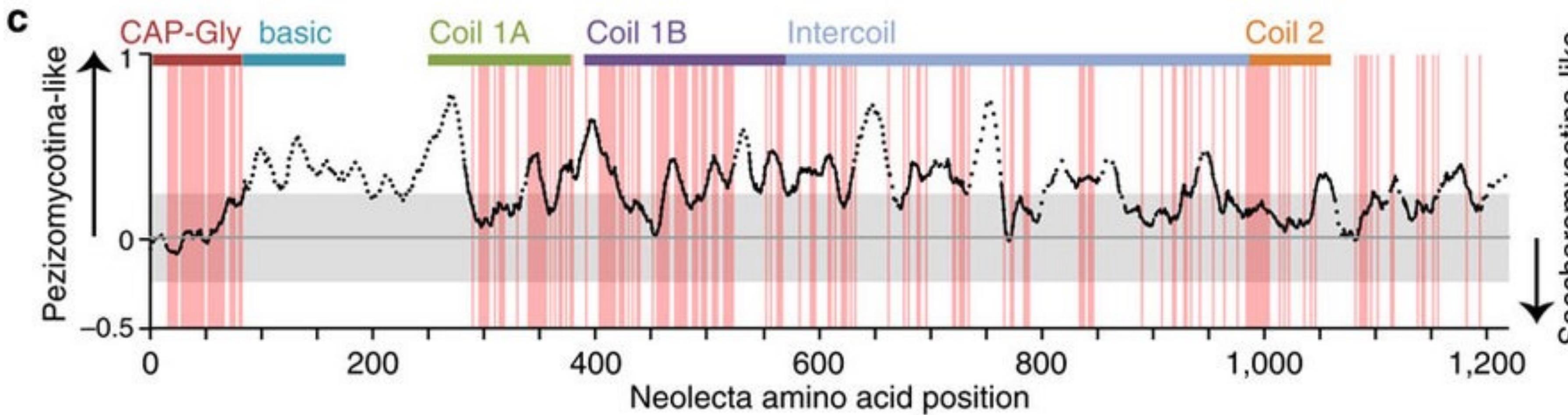
Novel proteins' localization
Enriched for transmembrane domains
MIT-1 is novel mitochondrial localized protein



DYNEIN AND ITS REGULATORS HAVE A COMPLEX MULTICELLULARITY SIGNATURE OF CONSERVATION



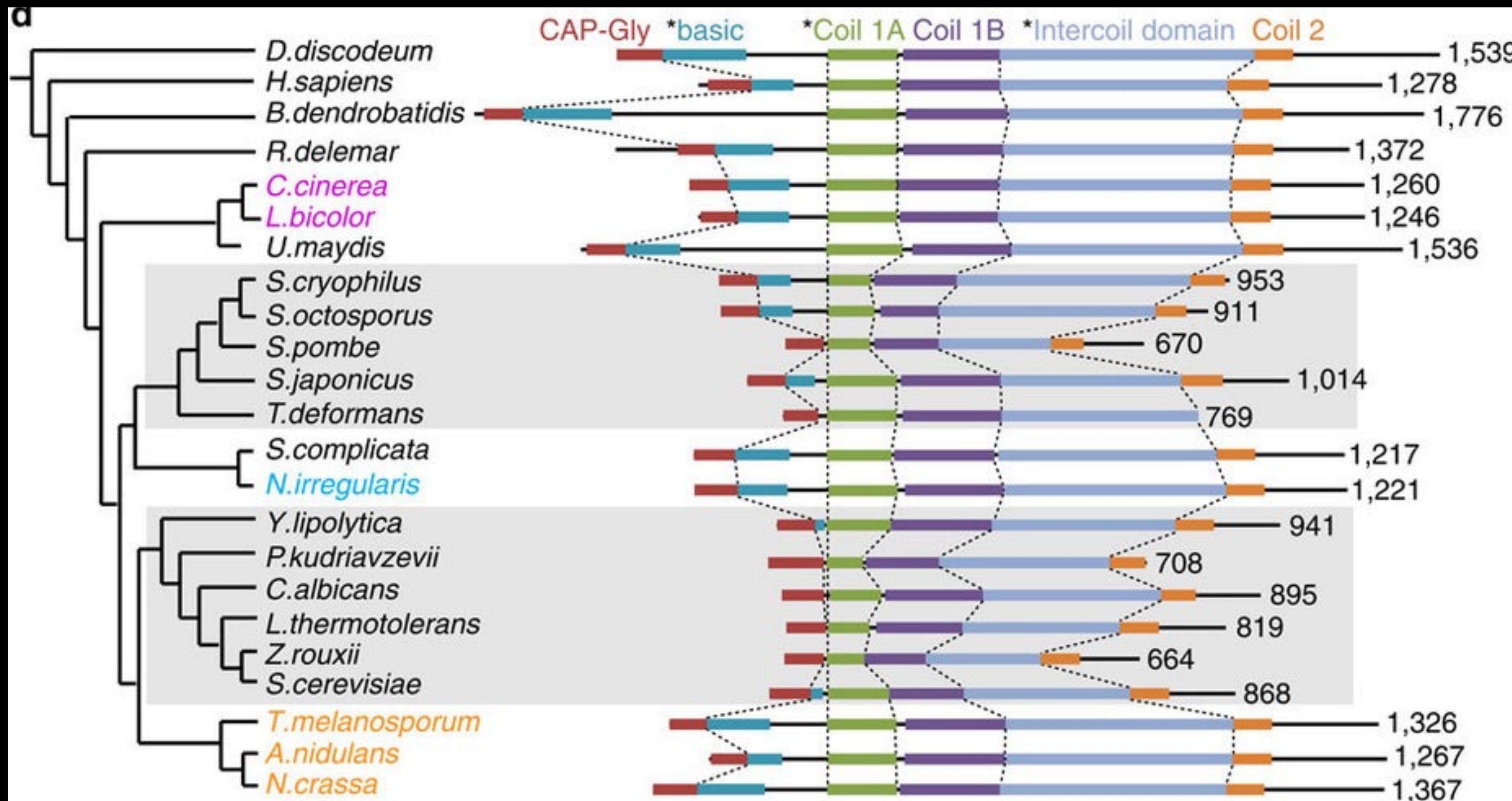
REGIONS OF CONSERVATION IN P150^{GLUED} PROTEIN



DOTTED LINES ARE MISSING IN YEAST - LOSING REGIONS RELATED TO MOTOR PROCESSIVITY

NEOLECTA P150^{GLUED} IS MORE SIMILAR TO PEZIZOMYCOTINA ORTHOLOGUES OVER THE ENTIRE LENGTH OF THE PROTEIN

INDEPENDENT DYNEIN PROTEIN CONTRACTION IN YEASTS



SUMMARY

- Patterns of gene loss in fungi from ancient genes: e.g. flagella, also independent losses dominate in formation of many yeast lineages
- Some aspects of complex multicellularity in fungi is likely not a dramatic gain of genes. Losses of Transcription factors and families in yeast lineages (e.g Laszlo Nagy et al Nat Com 2014.)
- Complexity may be result of repurposing or redirection of existing core genes.
- Ancestral fungus had gene more extensive than what we see in models yeasts
- Utility of comparing similarities of retained genes in species with complex morphology

FIN