

# Workshop Part III: EvoGeno Methods

## Comparative Genomics and Phylogenomics

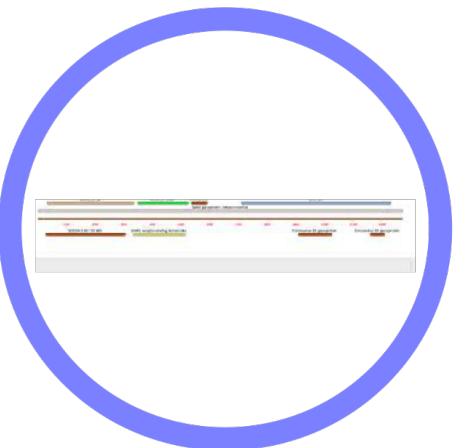
Yocelyn Gutiérrez Guerrero

[ygutierrezgro@berkeley.edu](mailto:ygutierrezgro@berkeley.edu)

May 05, 2023

# Overview

First Part



After assembly



Databases

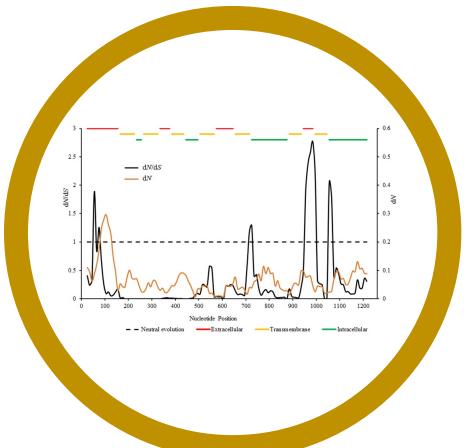


Orthologous identification

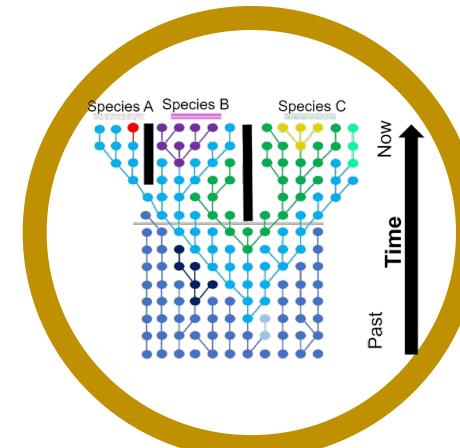
Second Part



Phylogenomics

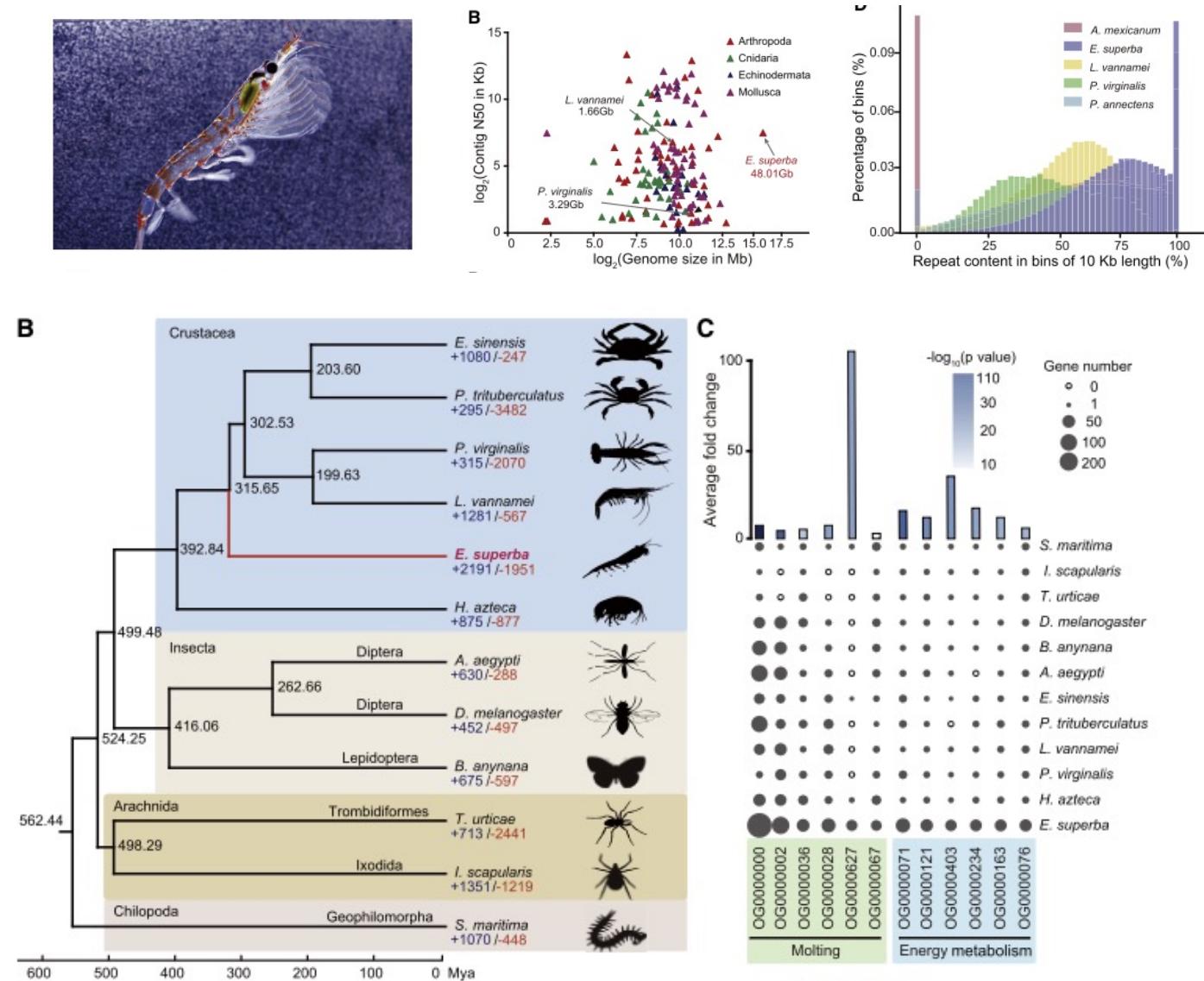
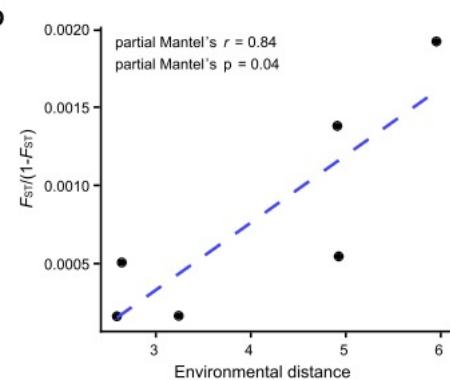
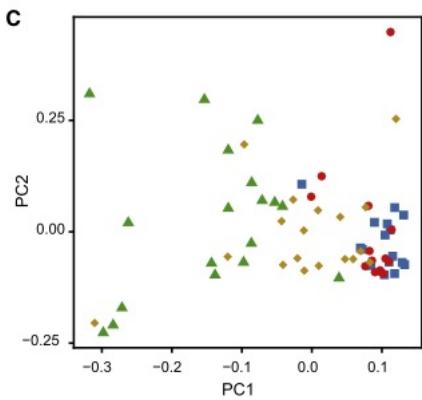
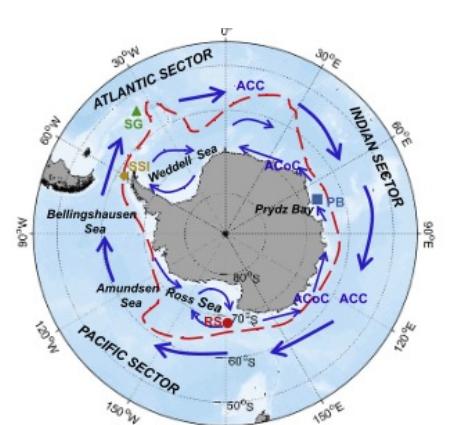


Selection Test



Incomplete lineage  
Sorting

# From Population to Comparative Genomics



# Comparative Genomics: Genomic features of different organism are compared

-Genome architecture

-Common and divergent features

-Gene evolutionary rates

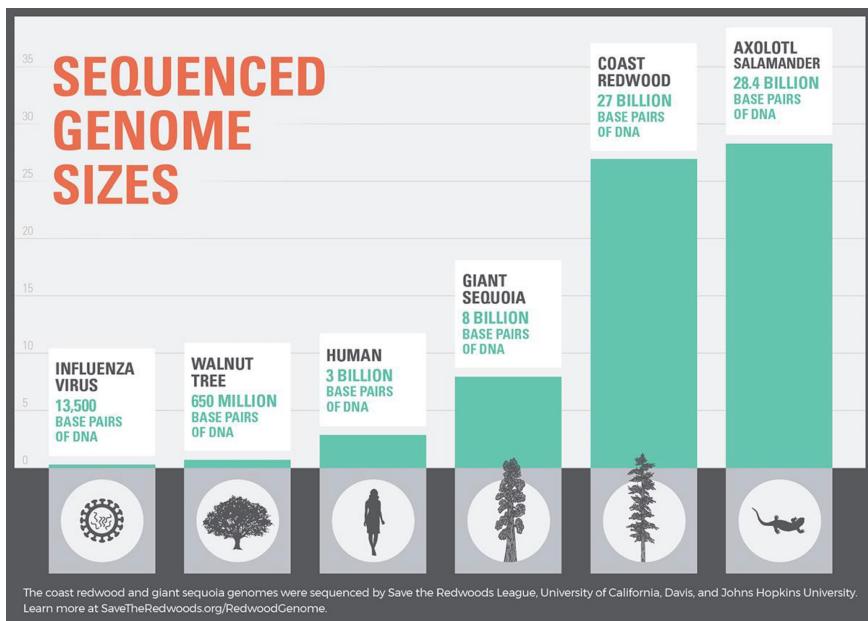
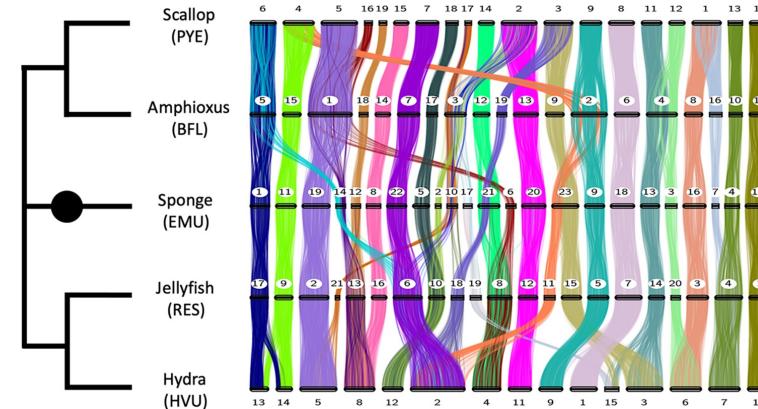
-Selection process

-Non coding elements evolution

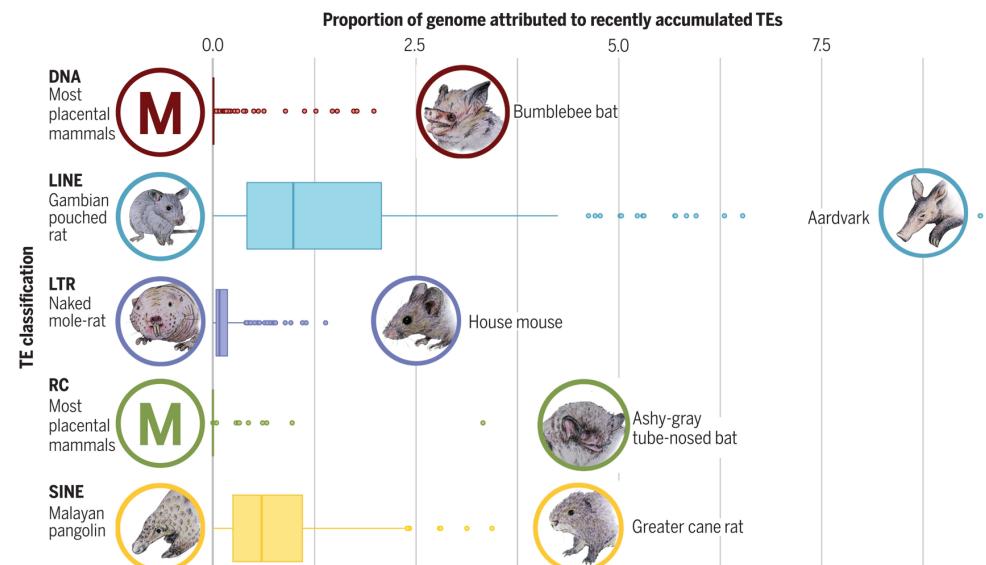
SCIENCE ADVANCES | RESEARCH ARTICLE

GENETICS

## Deeply conserved synteny and the evolution of metazoan chromosomes



## Insights into mammalian TE diversity through the curation of 248 genome assemblies



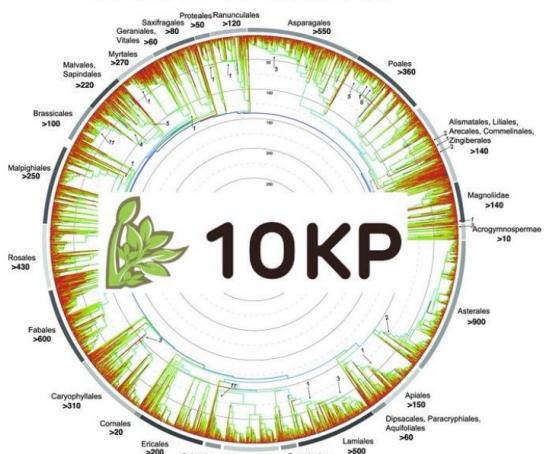


# Avian Phylogenomics Project



VERTEBRATE  
GENOMES  
PROJECT

A PROJECT OF THE G10K CONSORTIUM

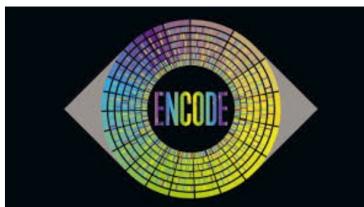


# The 1KP Project

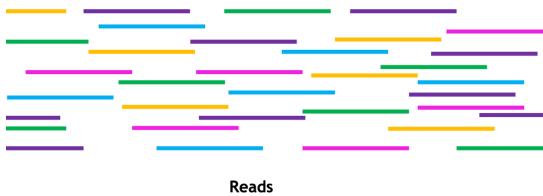
## Plants transcriptome

# 1001 Genomes

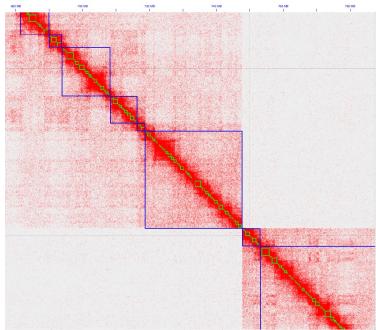
## A Catalog of *Arabidopsis thaliana* Genetic Variation



# After genome assembly



Relative order



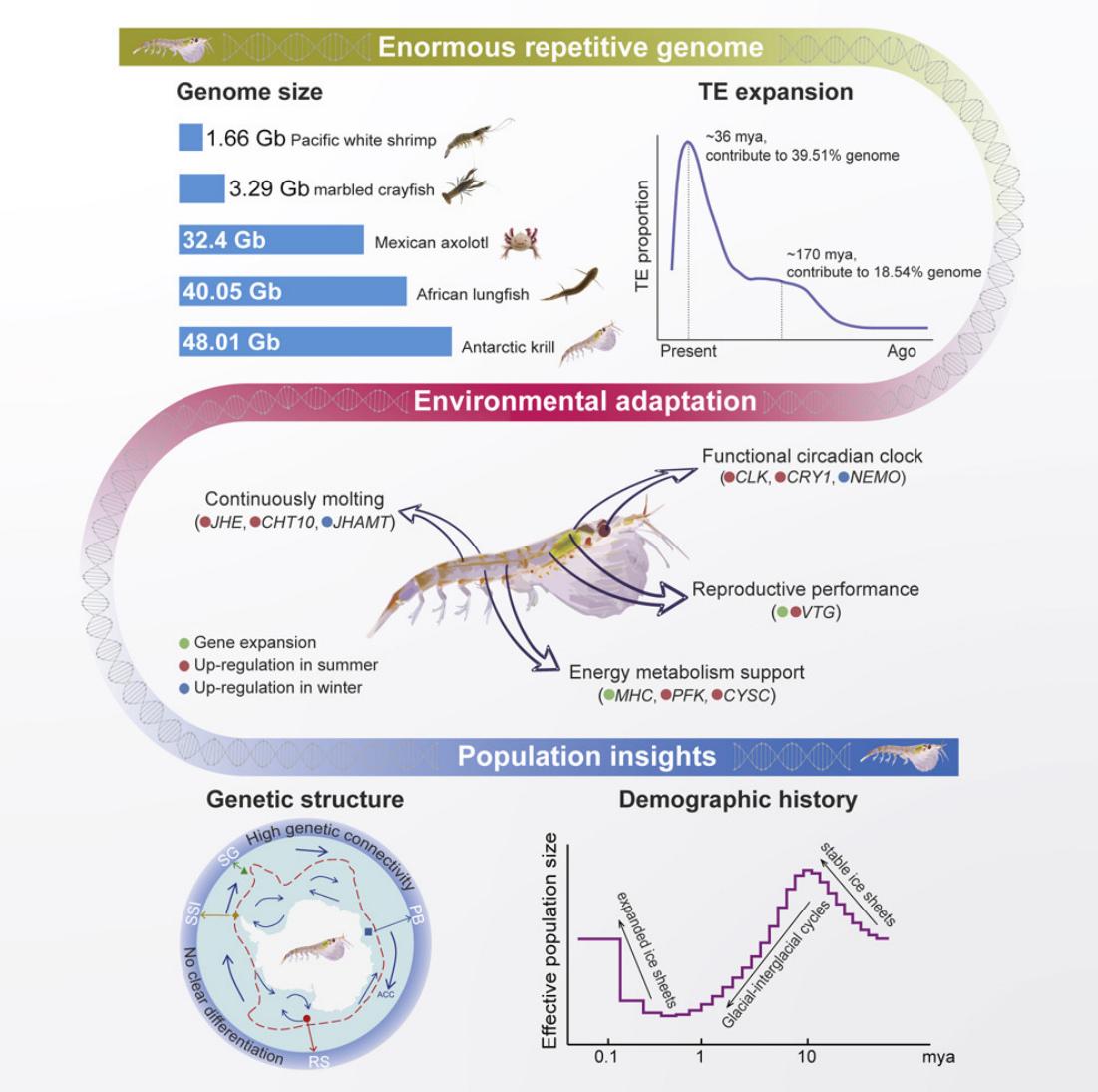
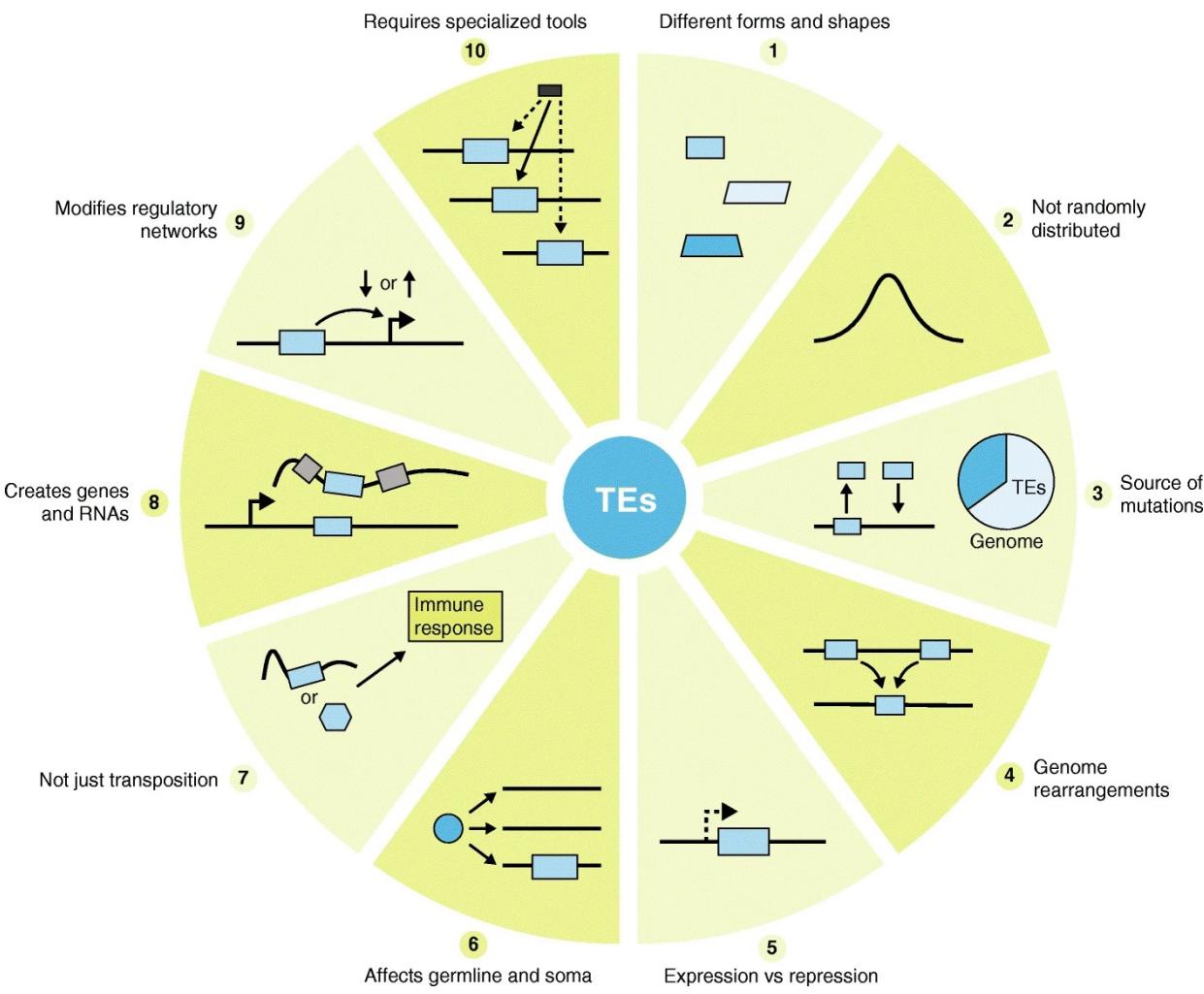
Repetitive elements  
identification



Mask genome

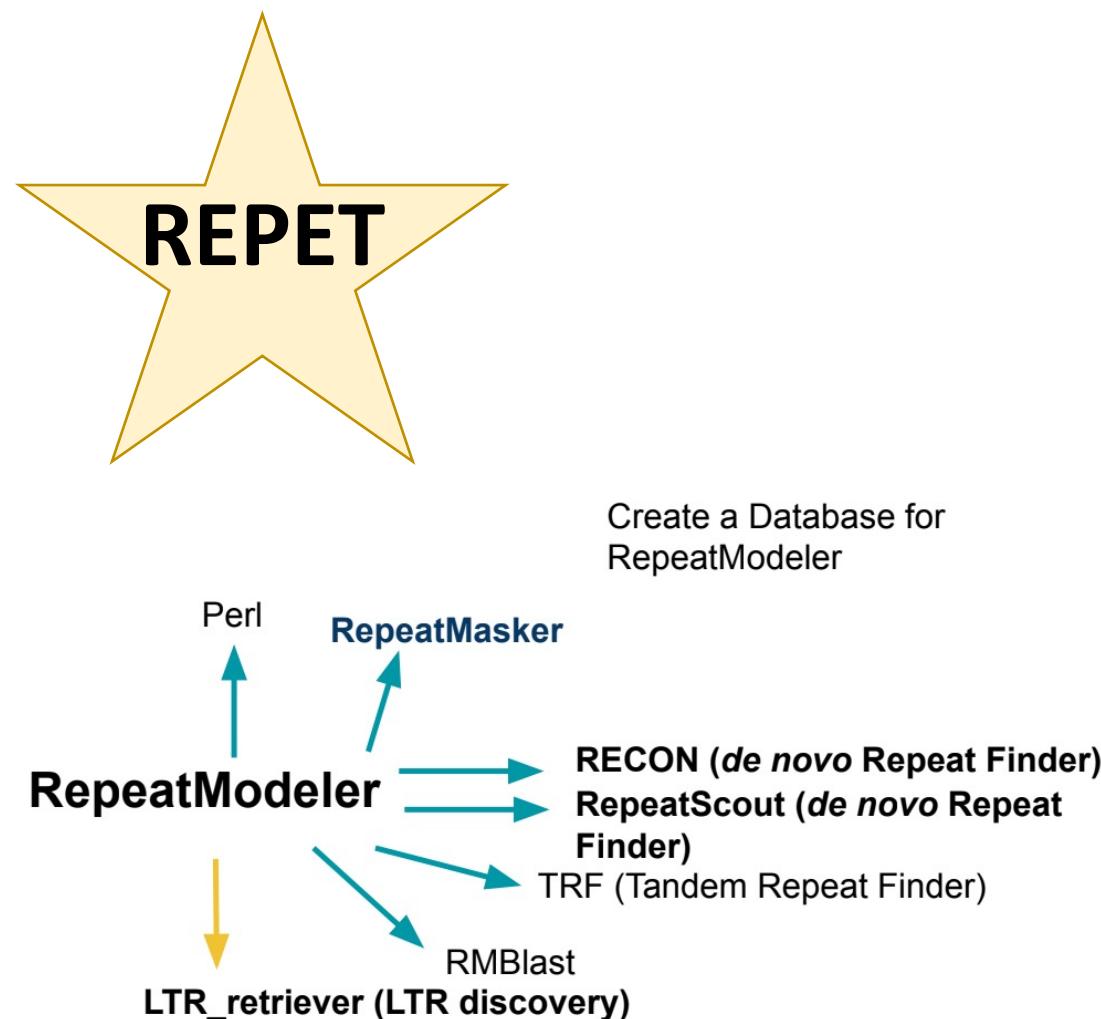
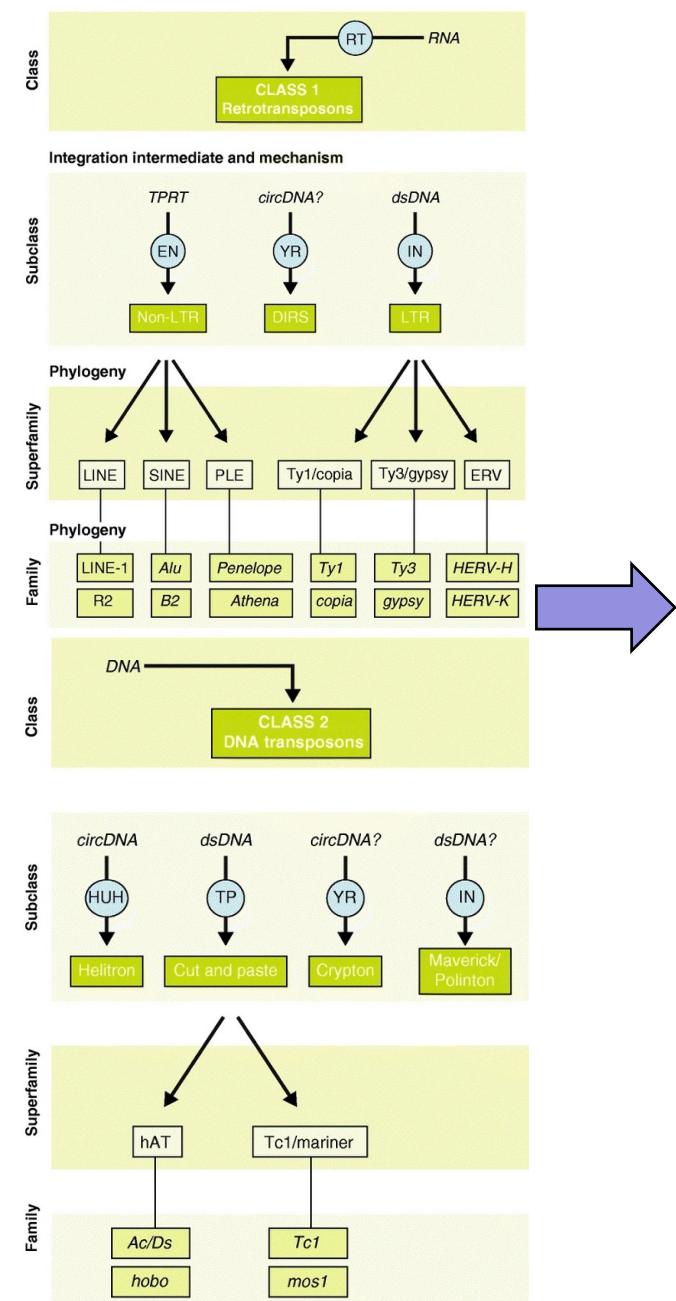
Gene annotation and  
genome features

# Repetitive elements content



Tandem Repeats (TE) is over 92%, two expansion events

# Repetitive elements identification



Dfam

giri  
REPBASE

### Fragments identified in sequence Possible families

MER58 or MER58B

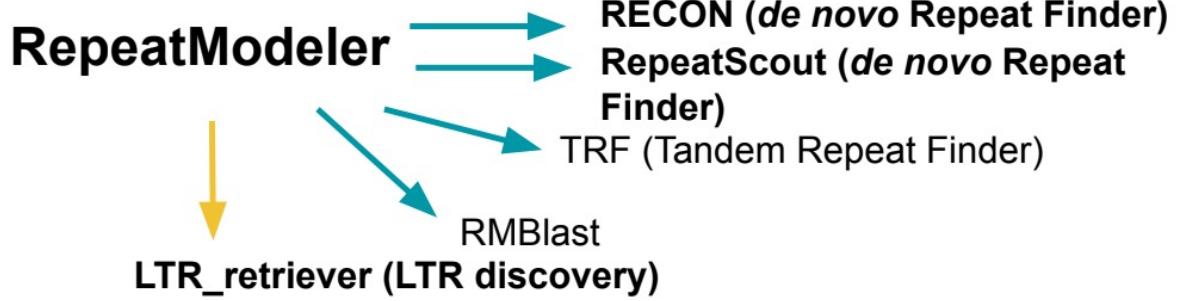
MER58 only

MER58, MER58B or MER58A

```
# STOCKHOLM 1.0
#=GF ID      rnd-1_family-20
#=GF DE      RepeatModeler Generated - rnd-1_family-20, RepeatScout: [ Index = R=13,
RS Size = 240, Refiner Input Size = 100, Final Multiple Alignment Size = 100 ]
#=GF SQ      100
#=GC RF      XXXXX..XX.XXXXXX..XX.XXXXXXX..XX..XXXX..XXX.....XXXX.....XX..XX.XXXXXXX.
.XXXXXXX.....XXX..X..X..XXX..XX...XXXXXX..XXXXXXXXXXXXXXXXXX..XXXXX..X..XXX..X.
XXXX..XXXXX..XXXXX..X..XXXXXXX..XXXX
NC_044211.1_RaG00:30307278-30307399    TGAGG.GG.CTGA.GC.GTGTCC.A..GAGA.AGG.....
.GCAA.....CA..GA.GCTGGGG..AAGGGTC.....TGG..T.G.CAC.AA..GTCTTAT.GAGGAGTGG
CTGAGGGAG..CTGGG...G.GTA..T.TTAT..CTTGG.AGAAA.A.GGAGGATC.AGGG
NC_044211.1_RaG00:69032180-69032300    TGAGG.AG.CTGA.GT.GTGTCC.A..GAGA.AGG.....
.GCAA.....TG..GA.ACGGGTG..AAGGGTC.....TGG..A.G.CAC.AA..GTCCTGT.GAGGAGCAG
CTCAGGGAG..CTGAA..G.TTC..T.TTAG..CC.AG.AGAAA.A.GGAAGCTC.AGGG
NC_044213.1_RaG00:139863354-139863474    TGAGG.TG.CTGA.GT.GTGTCC.A..GAGA.AGG.....
..A.AA.....CA..GA.GCTGGGG..AAGGGC.....TGG..A.G.CAC.AG..GGCTGT.GAGGAGC
AGCTGAGGGAG..CTGGG...G.GTG..C.TCAG..CCTGG.AGAA.A.GGAGGCTC.AGGG
NC_044214.1_RaG00:54951091-54951212    TGTGG.TG.CTGA.TT.GTGCC.A..GAGA.AGG.....
.GCAA.....TG..GA.TGTGGT..AAGGGTC.....TGG..A.G.CAC.AA..GCCCTAC.GAGGAGCAC
TTGAGGGAG..CTGGG...G.GTG..T.TTAC..TCTGC.AGAAA.A.GGAGGCC.TGGG
NC_044213.1_RaG00:64961061-64961182    TGAGG.GG.CTGA.GC.TTGTCC.A..GAGA.ATG.....
families.stk]
```

## Consensus sequences identified

# *de novo* Identification with Repeat Modeler



# Genome assembly

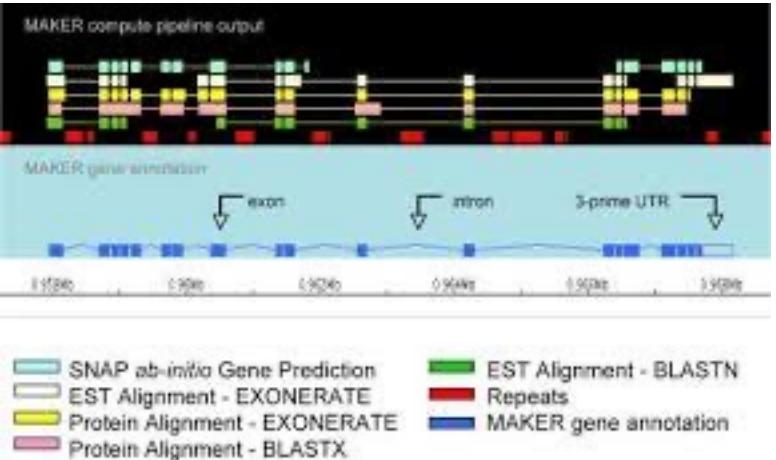
GGAAAGGAATACCTTGCAAAAAAAAGTAAAGGAAGAACTTGAAAAGCTG  
 AAAATGAGAAGAAAGAGGAAGCTGAAAAGGAGGAAGGAAGAAAATTGCT  
 GCAGAAGAAATTCTAAGTTCAAGAGCATGTGCGTAATTCACCTCAG  
 TTTGCATGTCGTTGTTGAATTGATATTCCACTACAGTTATAATC  
 AGTTAACAGGTTTAAAGTCCAAGTACTATTCTAAATCCATTATTG  
 ATTCAATTACAATAACAGCATCTATTGCAGATAAGAGTATGAAACTTGA  
 TTTTAGACAAAGATGCTCCTTCTTCATAGATATTGAATTCTACAAGA  
 TTACTCAAGATCATATTCTTATCCATATGAACTGCTTTATTCTTC  
 AGTAAATAATTNTTTGAGNNNNNNNNNNNNNNNNNNNNNNNNNNNN  
 NNN  
 NNN  
 CCTGGTGCAGAAAACCAAGGGATAATTCTGGCAGGACTGCACTTAA  
 TGAGAGCCCTTCCAATCCTGCTGCTTATTAAACAGAGAAATTGCTCTGG  
 ATTCTCCTCCTCCAATCCAGCCCTCTCCATGGTCAGCTCC  
 TTCCTCCCCCAGGCCAGGCTGCCCCATAAACACTTGTACATCAA  
 ATAGGATTAATTACAGAGCCAGAAATATTCTGCCACATTCACCCCT  
 GCACAGCCCACAAAGCCCTGGCAATTCAAATATGTGGATTATGGTTG  
 CTTCTCAGGGTCACTGGAGCTAAGGACAGATGCCTGGAGTCCTCTGG  
 CTGTCAGGAATGCTGCTCTTTAGACATTAGTCCAAATAAATAAT  
 AATCTAAAGCCAAGCTGGTTTACTCTATTGAAACCTCACATTGATG

Masked genome

##date 2021-02-22	##sequence-region RaGOO_Dbar-Tgut_500.fasta	RepeatMasker	similarity	860	1125	26.2	+	.	Target "Motif:rnd-5_family-980" 6
269	Chr0_RaGOO	RepeatMasker	similarity	1482	1584	22.6	+	.	Target "Motif:rnd-5_family-980" 65
9 761	Chr0_RaGOO	RepeatMasker	similarity	1516	1788	28.6	-	.	Target "Motif:rnd-2_family-34" 235
868	Chr0_RaGOO	RepeatMasker	similarity	1789	2276	21.1	+	.	Target "Motif:rnd-5_family-2770" 4
6 511	Chr0_RaGOO	RepeatMasker	similarity	2277	2365	28.6	-	.	Target "Motif:rnd-2_family-34" 23
234	Chr0_RaGOO	RepeatMasker	similarity	2820	2884	20.2	+	.	Target "Motif:A-rich" 1 62
Chr0_RaGOO	RepeatMasker	similarity	3358	3402	19.4	+	.	Target "Motif:A-rich" 1 43	
Chr0_RaGOO	RepeatMasker	similarity	4162	4224	26.5	+	.	Target "Motif:(GCA)n" 1 67	
Chr0_RaGOO	RepeatMasker	similarity	4534	4576	9.3	-	.	Target "Motif:rnd-6_family-3719" 3	
45	Chr0_RaGOO	RepeatMasker	similarity	4700	4854	20.6	-	.	Target "Motif:rnd-2_family-7" 142
336	Chr0_RaGOO	RepeatMasker	similarity	4868	4907	15.0	-	.	Target "Motif:rnd-6_family-5569" 1
99 238	Chr0_RaGOO	RepeatMasker	similarity	5171	5259	27.0	-	.	Target "Motif:rnd-5_family-247" 15
4 242	Chr0_RaGOO	RepeatMasker	similarity	6335	6365	15.5	+	.	Target "Motif:(TA)n" 1 30
508	Chr0_RaGOO	RepeatMasker	similarity	6390	6529	17.4	+	.	Target "Motif:rnd-3_family-51" 367
Chr0_RaGOO	RepeatMasker	similarity	9674	9708	9.1	+	.	Target "Motif:(A)n" 1 35	
Chr0_RaGOO	RepeatMasker	similarity	11698	11733	5.9	+	.	Target "Motif:(AGTCT)n" 1 36	
Chr0_RaGOO	RepeatMasker	similarity	13317	13472	23.1	-	.	Target "Motif:rnd-5_family-288" 1	

# Genome Annotation

TATATATAT-  
NNNNNNN-



## RNA/Transcript Evidence (the options are called EST for historic reasons)

```
est= #set of ESTs or assembled mRNA-seq in fasta format
altest= #EST/cDNA sequence file in fasta format from an alternate organism
est_gff= #aligned ESTs or mRNA-seq from an external GFF3 file
altest_gff= #aligned ESTs from a closely related species in GFF3 format
```

## Protein Homology Evidence

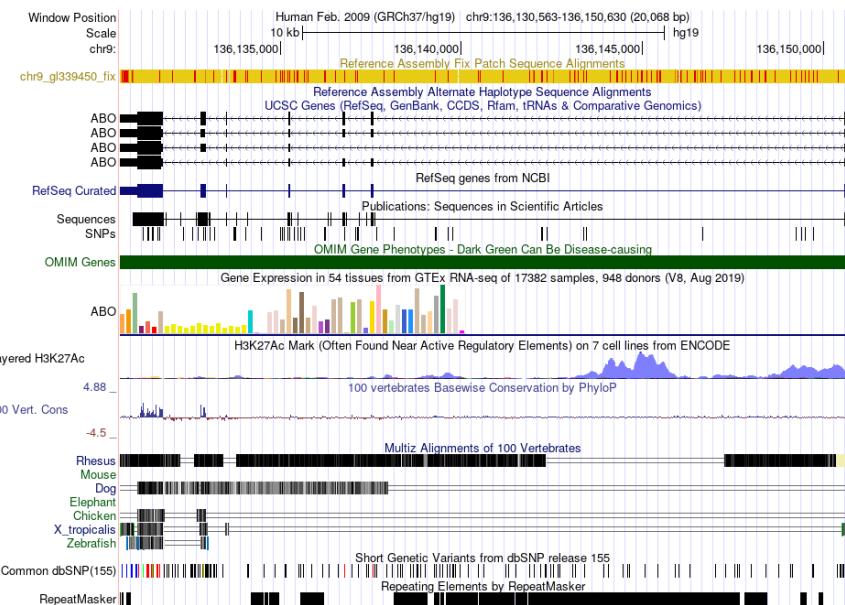
```
protein= #protein sequence file in fasta format (i.e. from multiple organisms)
protein_gff= #aligned protein homology evidence from an external GFF3 file
```

## Repeat Masking

```
model_org=all #select a model organism for RepBase masking in RepeatMasker
rmlib= #provide an organism specific repeat library in fasta format for RepeatMasker
repeat_protein= #provide a fasta file of transposable element proteins for RepeatRunner
rm_gff= #pre-identified repeat elements from an external GFF3 file
prok_rm=0 #forces MAKER to repeatmask prokaryotes (no reason to change this), 1 = yes, 0 = no
softmask=1 #use soft-masking rather than hard-masking in BLAST (i.e. seg and dust filtering)
```

## Gene Prediction

```
snapHmm= #SNAP HMM file
gmHmm= #GeneMark HMM file
augustus_species= #Augustus gene prediction species model
fgenesh_par_file= #FGENESH parameter file
pred_gff= #ab-initio predictions from an external GFF3 file
model_gff= #annotated gene models from an external GFF3 file (annotation pass-through)
est2genome=1 #infer gene predictions directly from ESTs, 1 = yes, 0 = no
protein2genome=0 #infer predictions from protein homology, 1 = yes, 0 = no
trna=0 #find tRNAs with tRNAscan, 1 = yes, 0 = no
snoscan_rrna= #rRNA file to have Snoscan find snoRNAs
unmask=0 #also run ab-initio prediction programs on unmasked sequence, 1 = yes, 0 = no
```



# Estimating phylogenies from genomes

Phylogenomics focuses largely on analyzing evolutionary histories to reconstruct relationships between taxa

A persistent problem in phylogenomics is the selection of appropriate genetic data or markers for a given taxonomic groups

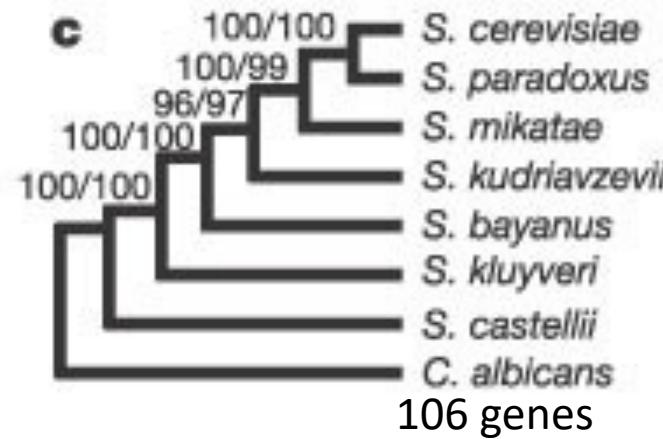
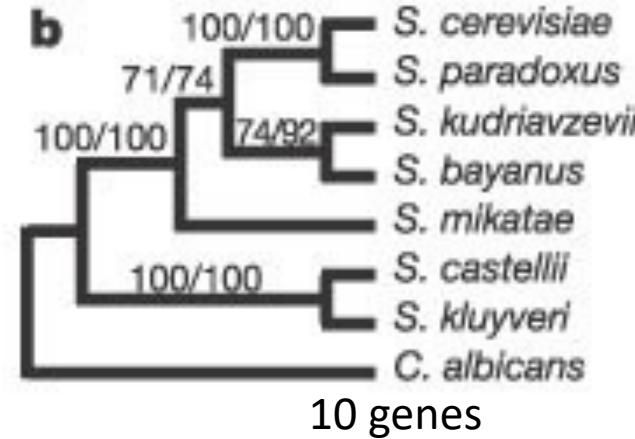
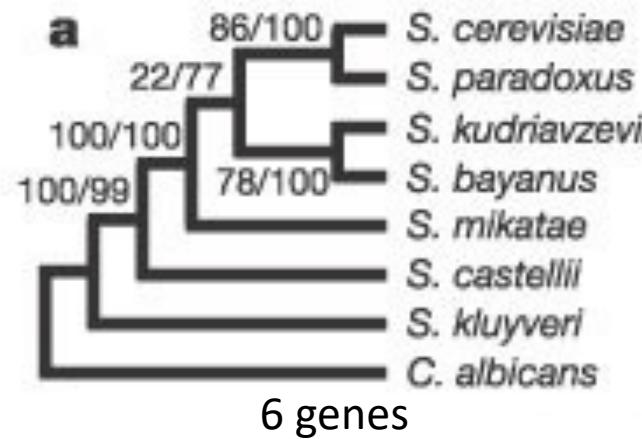
Practical limitations or considerations:

- Genome quality

- Genome species sequenced  
(Transcriptome alternative)

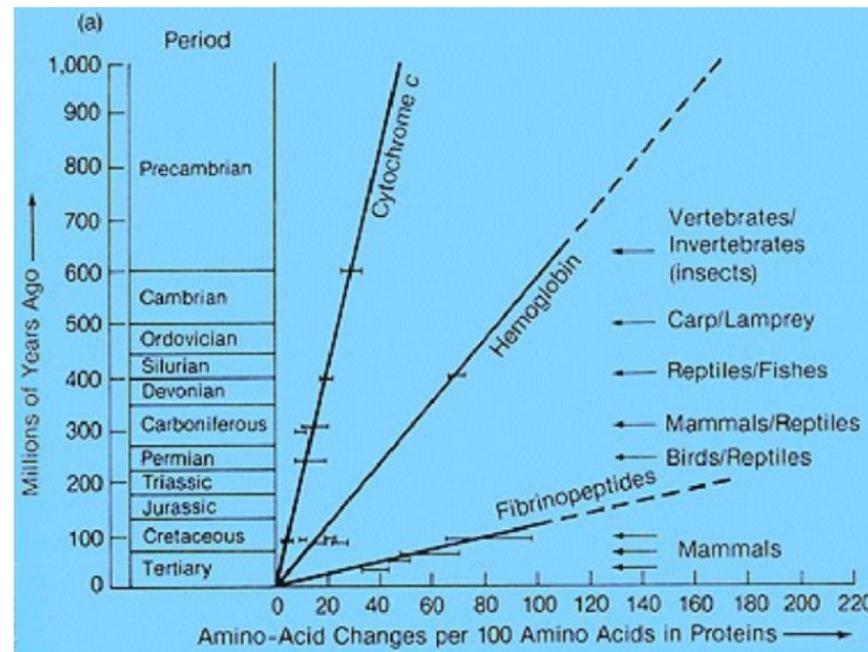
- Computational time

# More genes with different evolutionary rates increase and fully resolve species tree with maximum support



Gene	Description	Reference
<i>EF-1<math>\alpha</math></i>	Elongation factor-1 $\alpha$ , Role in protein synthesis.	[52]
<i>rpoA gene</i>	Encoding the alpha subunit of RNA polymerase	[53]
<i>atpB</i>	Encode the beta subunit of ATP synthase	[54]
<i>dnaA</i>	involved in DNA synthesis initiation	[55]
<i>ftsZ</i>	Role in cell division	[56]
<i>gapA</i>	Codes for glyceraldehyde phosphate dehydrogenase	[57]
<i>groEL</i>	Encodes bacterial heat shock protein.	[58]
<i>gltA</i>	Encoding citrate synthase	[59]
<i>ITS</i>	Piece of non-functional RNA situated between structural ribosomal RNAs precursor transcript.	[60]
<i>lux Gene</i>	encode proteins involved in luminescence	[61]
<i>PEPCK</i>	Codes for phosphoenolpyruvate carboxykinase	[62]
<i>pyrH genes</i>	Codes for uridine monophosphate (UMP) kinases	[63]
<i>recA</i>	Role in recombination	[64]
<i>U2 snRNA</i>	Component of the spliceosome	[65]
<i>Wsp gene</i>	Encodes a major cell surface coat protein	[66]
<i>Nuclear H3</i>	Codes for protein which is associated with DNA	[67]
<i>trnH-psbA</i>	Non-coding intergenic spacer region located in plastid genome	[68]
<i>rpoB, rpoC1</i>	Coding region located in plastid genome	[69]

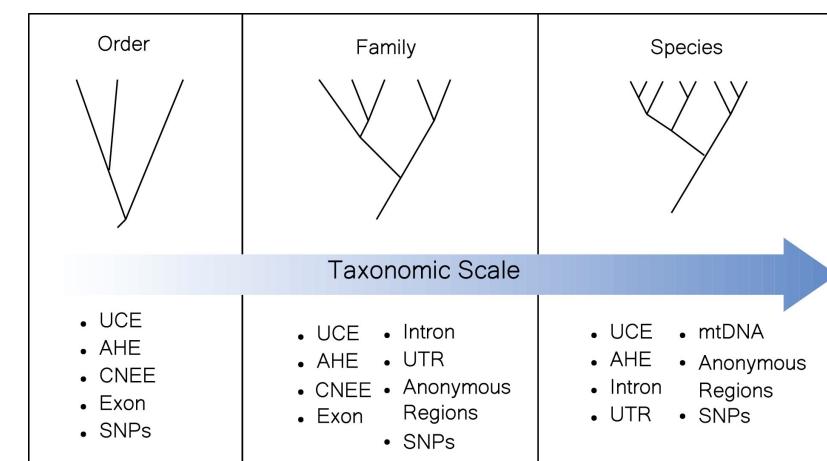
Table 1: List of some other molecular markers used in phylogeny research.

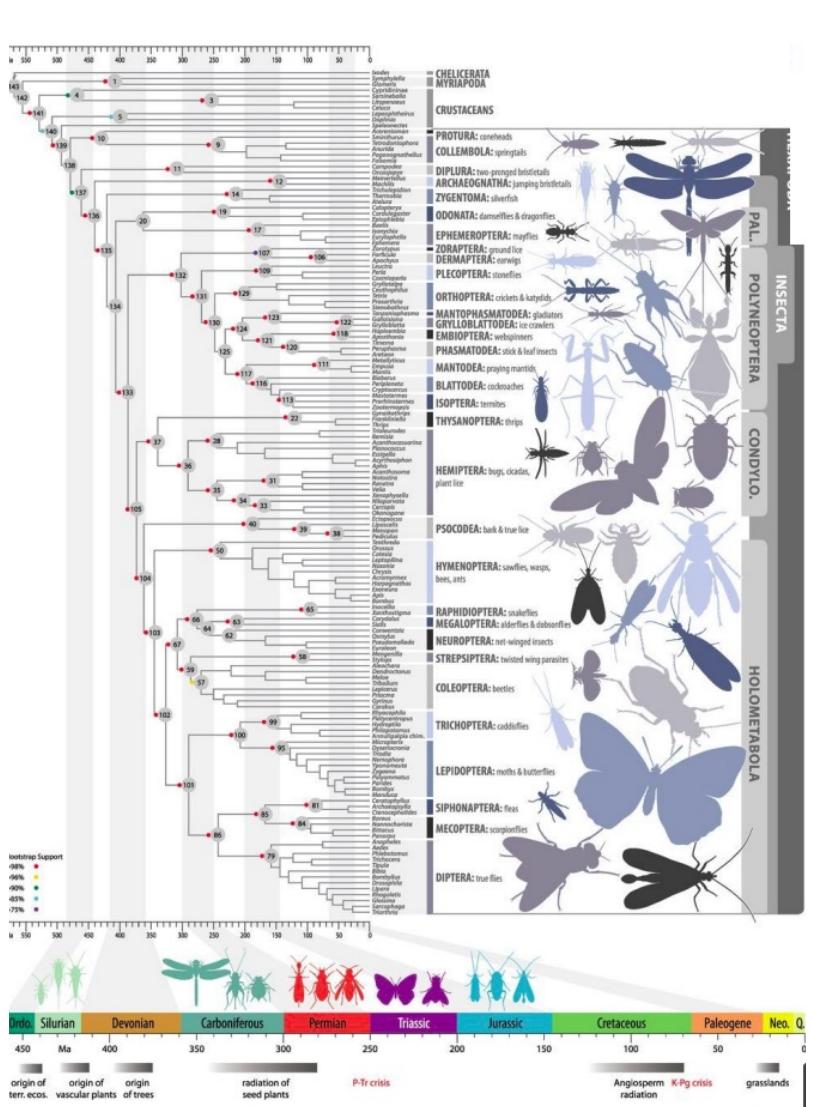


Marker types	Specimen samples needed/minimum quality needed	Evolutionary history estimate (deep-moderate-shallow)	Types of genetic and genomic method needed/reference genome requirement	Relative cost based on genomic method needed (see previous column)
UCEs = ultraconserved element flanking regions	Historic specimens, fresh tissue, blood samples/low quality	Deep to shallow time estimates	Target method, WGS data with computational target method/reference genome needed if designing probes	Low to moderate
AHE = anchor hybrid enrichment	Historic specimens, fresh tissue, blood samples/low quality	Deep to shallow time estimates	Target method, WGS data with computational target method/reference genome needed if designing probes	Low to moderate
CNEEs = conserved nonexonic elements	Historic specimens, fresh tissue, blood tissue/moderate quality	Moderate to deep time estimates	Target method or WGS data (with computational target method (preferred)/reference genome preferred	Low to moderate
Exon	Historic specimens, fresh tissue, blood sample, RNA samples/low to high quality	Mostly deep time estimates	WGS data, RNA/transcriptomic data, exonic target capture/reference genome preferred	Moderate to high
Introns	Historic specimens, fresh tissue, blood sample/moderate quality	Deep to moderate time estimates	WGS data/transcriptomic data, intronic target capture/reference genome preferred	Low to moderate
UTRs = untranslated regions	Fresh tissue, blood sample, RNA samples/high quality	Deep or shallow time estimates	WGS, RNA/transcriptomic data/reference genome preferred	Moderate to high
Mitochondrial DNA	Historic specimens, fresh tissue, blood samples/low quality	Shallow time estimates	mtDNA primers, WGS/no reference genome required	Low
Anonymous loci/regions (RADseq)	Historic specimens (not including RADseq), blood sample, tissue sample/moderate quality	Moderate to shallow time estimates; ultimately depends on the primers used	PCR primers, WGS/no reference genome required for PCR primers	Low

**Estimating phylogenies from genomes: A beginners review of commonly used genomic data in vertebrate phylogenomics** 

Javan K Carter , Rebecca T Kimball, Erik R Funk, Nolan C Kane, Drew R Schield, Garth M Spellman, Rebecca J Safran

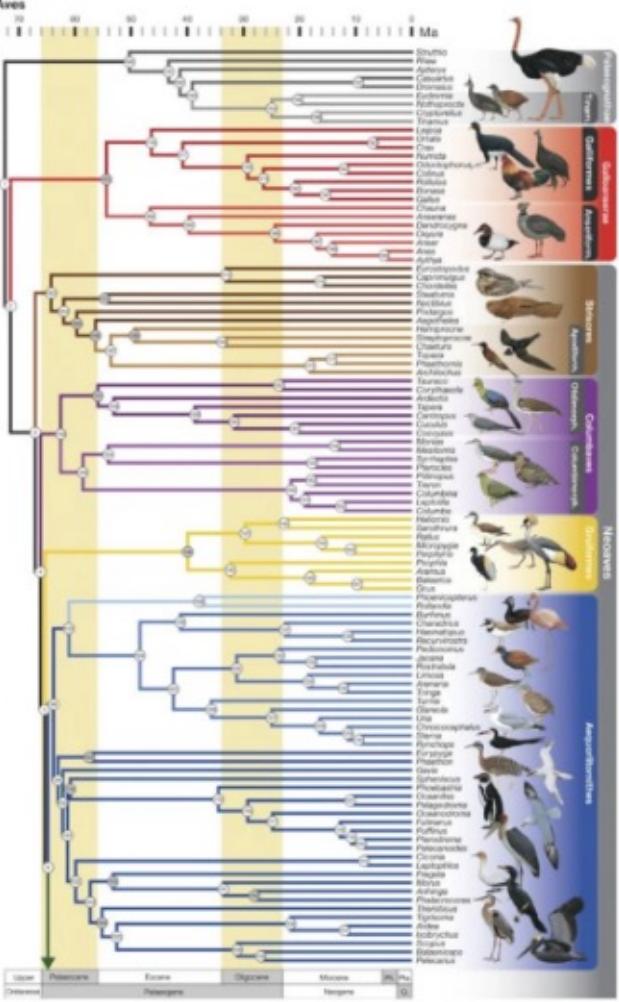




**REPORT**  
Phylogenomics resolves the timing and pattern of insect evolution

Bernhard Misof<sup>1,\*†</sup>, Shanlin Liu<sup>2,3,\*</sup>, Karen Meusemann<sup>1,4,\*</sup>, Ralph S. Peters<sup>5,\*</sup>, Alexander Donath<sup>1,\*</sup>, Christoph Mayer<sup>1,\*</sup>, Pau...

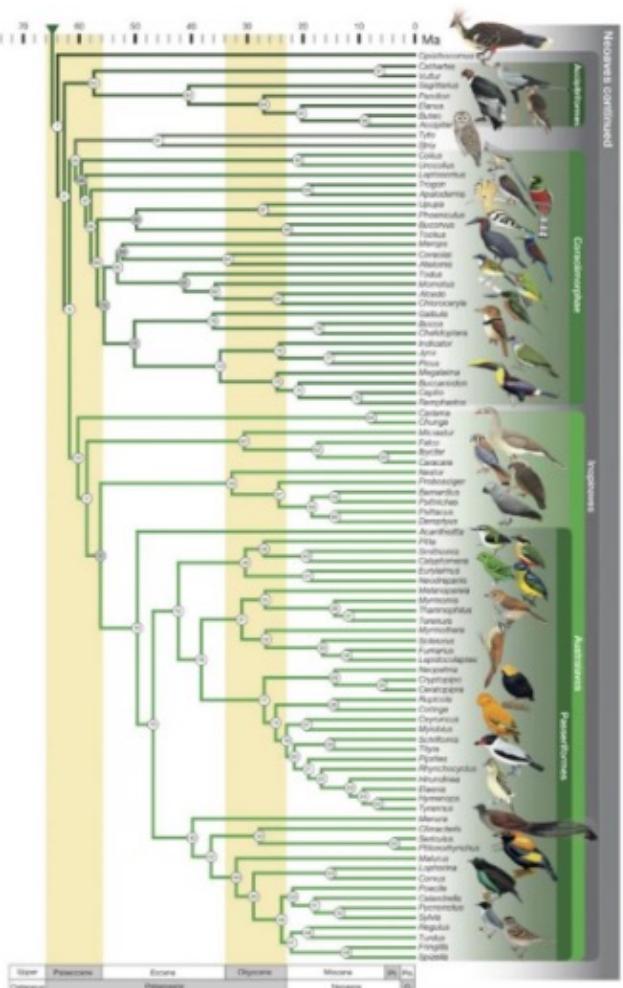
1,478 genes across species

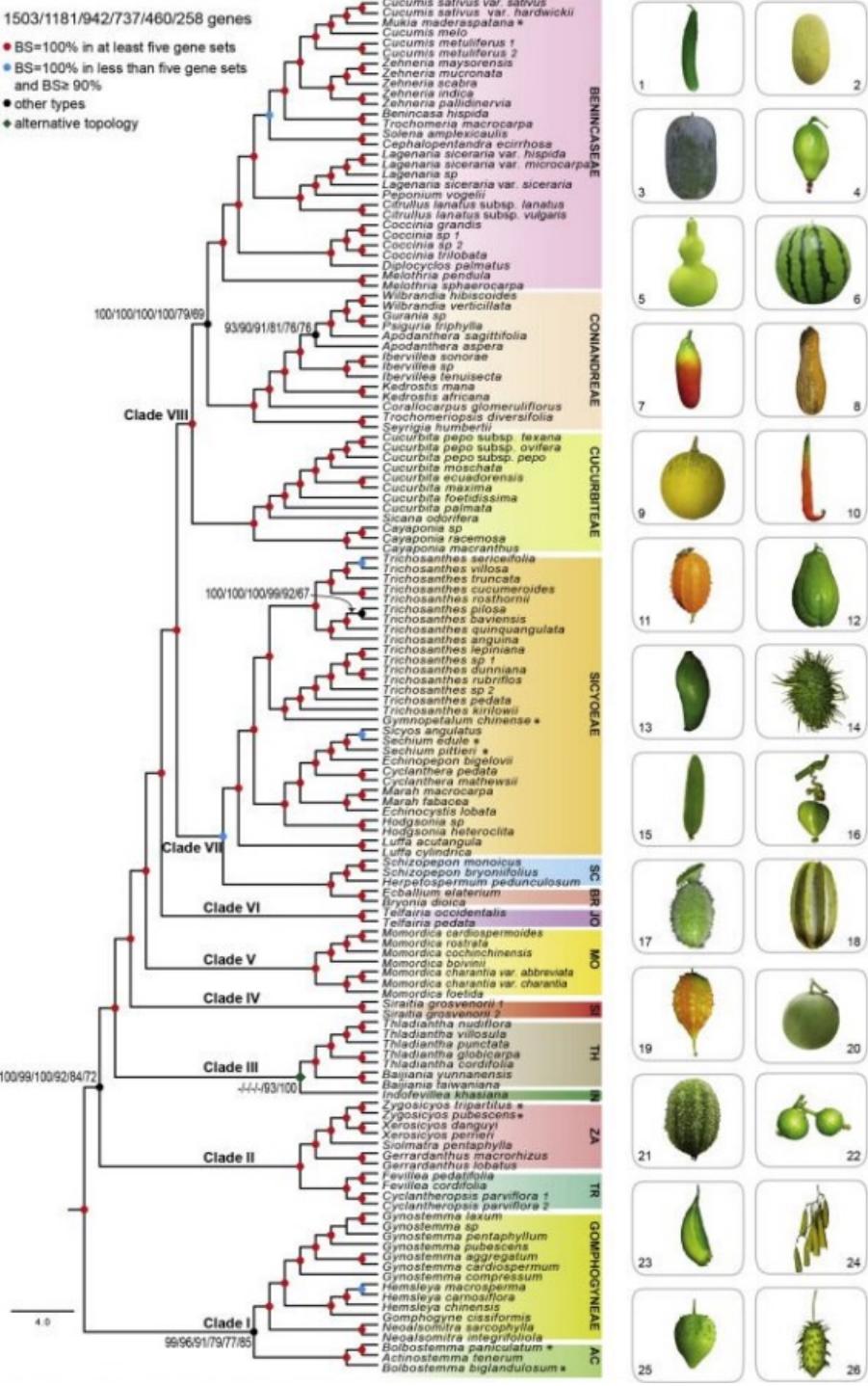


259 nuclear loci (1,523 bp)

# A comprehensive phylogeny of birds (Aves) using targeted next-generation DNA sequencing

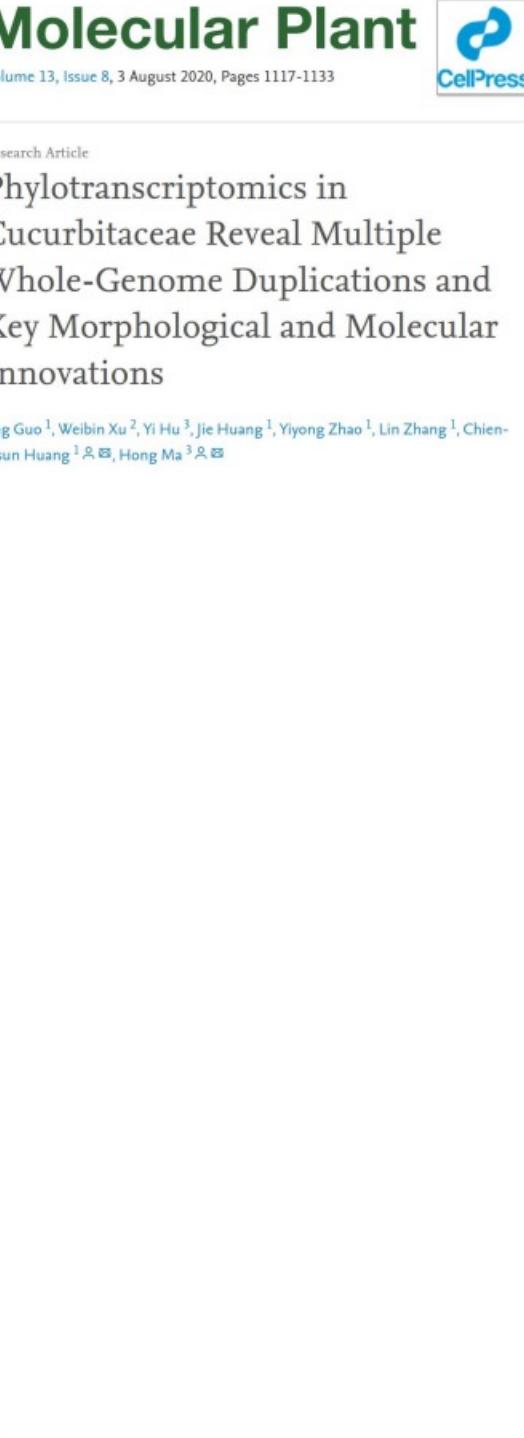
Richard O. Prum , Jacob S. Berv , Alex Dornburg, Daniel J. Field, Jeffrey P. Townsend, Emily Moriarty Lemmon & Alan R. Lemmon





# Phylogenetic transcriptomics in Cucurbitaceae reveal multiple whole-genome duplications and key morphological and molecular innovations

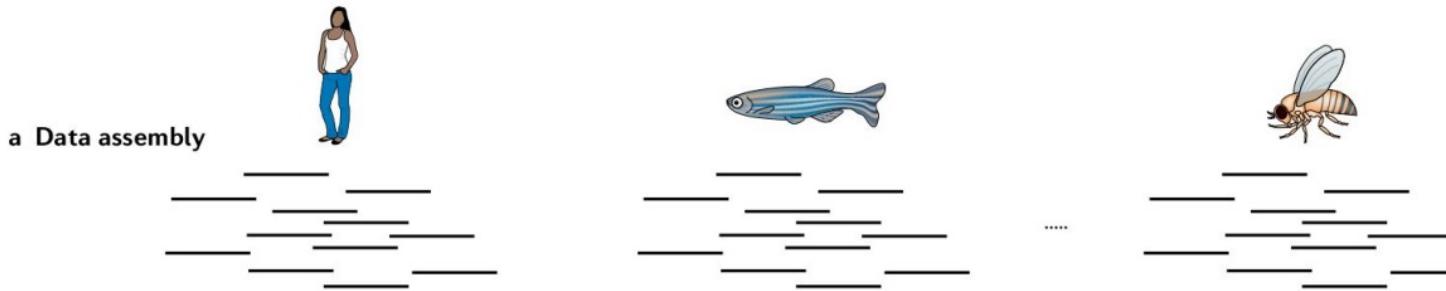
Guo<sup>1</sup>, Weibin Xu<sup>2</sup>, Yi Hu<sup>3</sup>, Jie Huang<sup>1</sup>, Yiyong Zhao<sup>1</sup>, Lin Zhang<sup>1</sup>, Chien-sun Huang<sup>1</sup>✉, Hong Ma<sup>3</sup>✉



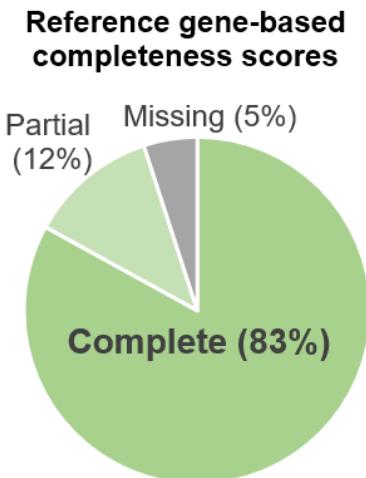
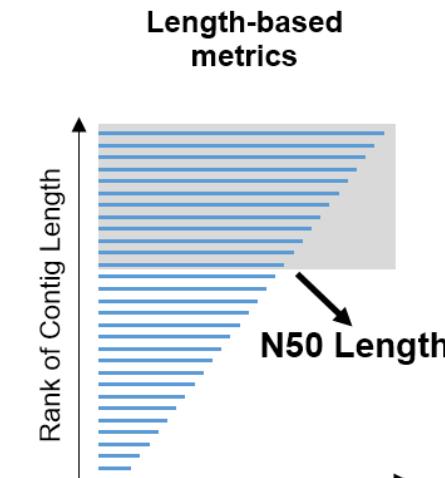
# Transcriptomic data for phylogenetic reconstruction using only coding genes as an alternative to explore complex genomes

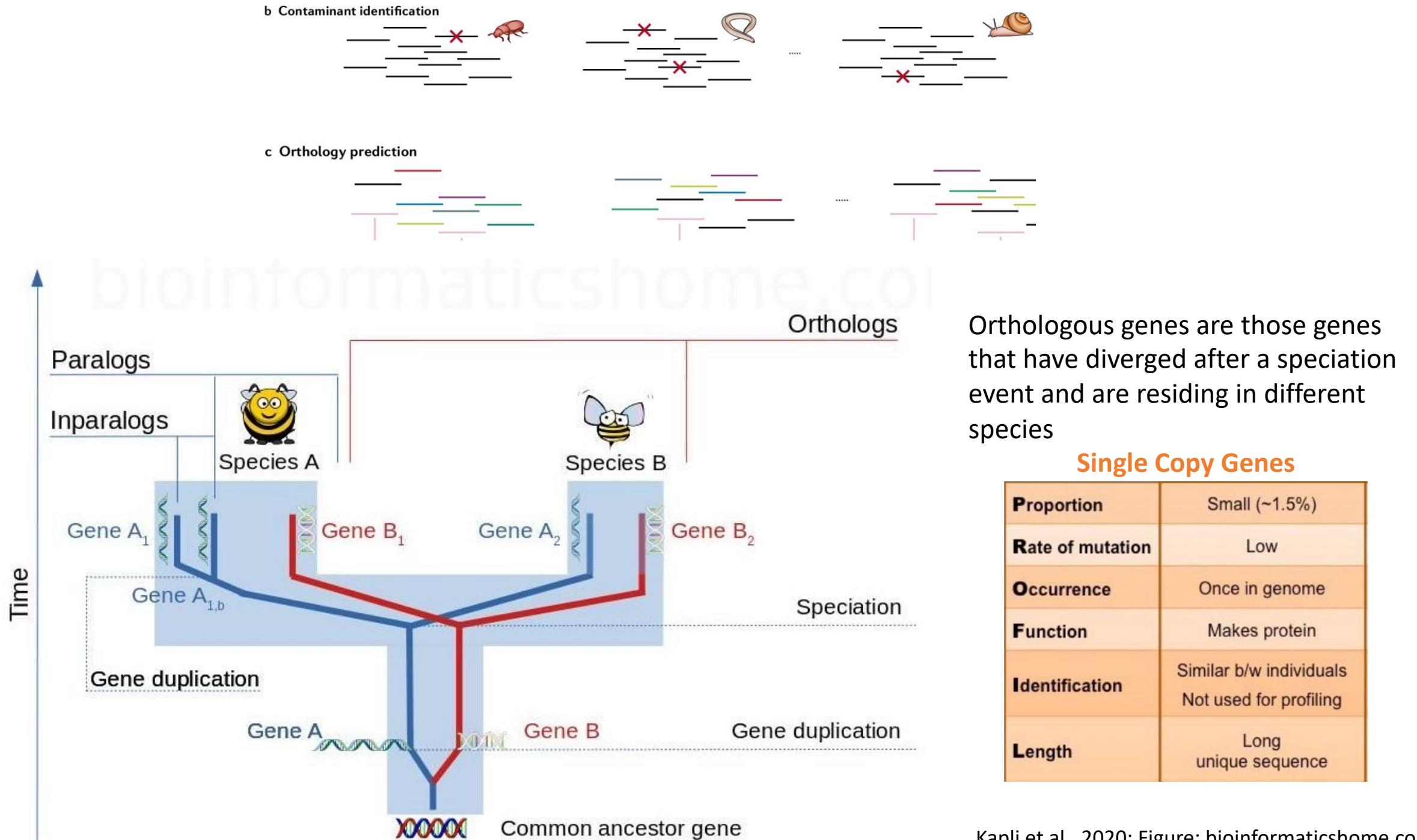
# Phylogenetic tree building in the genomic age

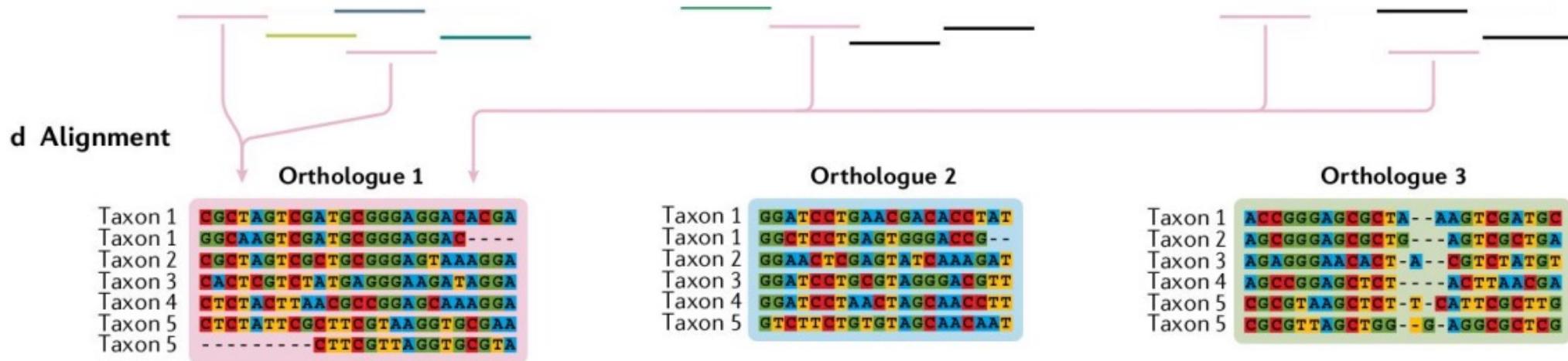
Paschalia Kapli  Ziheng Yang  and Maximilian J. Telford  



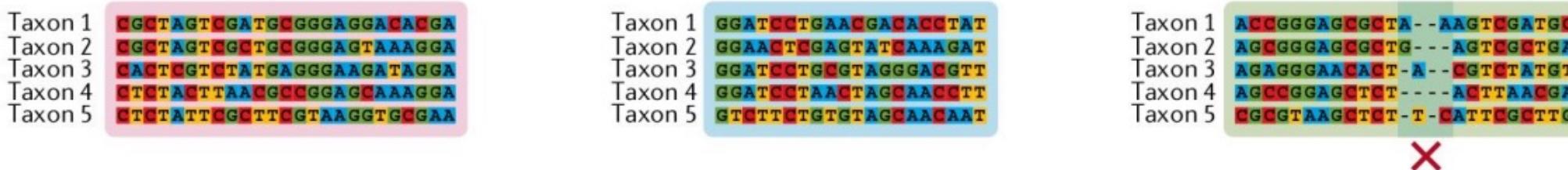
  
[Parent Directory](#)  
[Annotation\\_comparison/](#)  
[Evidence\\_alignments/](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_assembly\\_structure/](#)  
[Gnomon\\_models/](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_assembly\\_report.txt](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_assembly\\_stats.txt](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_cds\\_from\\_genomic.fna.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_feature\\_count.txt.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_feature\\_table.txt.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_genomic.fna.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_genomic.gbff.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_genomic.gff.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_genomic.gtf.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_genomic\\_gaps.txt.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_protein.faa.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_protein.gpff.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_pseudo\\_without\\_product.fna.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_rm.out.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_rm.run](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_rna.fna.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_rna.gbff.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_rna\\_from\\_genomic.fna.gz](#)  
[GCF\\_000325575.1\\_ASM32557v1\\_translated\\_cds.faa.gz](#)  
[GCF\\_000325575.1\\_knownrefseq\\_alns.bam](#)  
[GCF\\_000325575.1\\_knownrefseq\\_alns.bam.bai](#)  
[GCF\\_000325575.1\\_modelrefseq\\_alns.bam](#)  
[GCF\\_000325575.1\\_modelrefseq\\_alns.bam.bai](#)  
[Pteropus\\_alecto\\_AR102\\_annotation\\_report.xml](#)  
[README.txt](#)



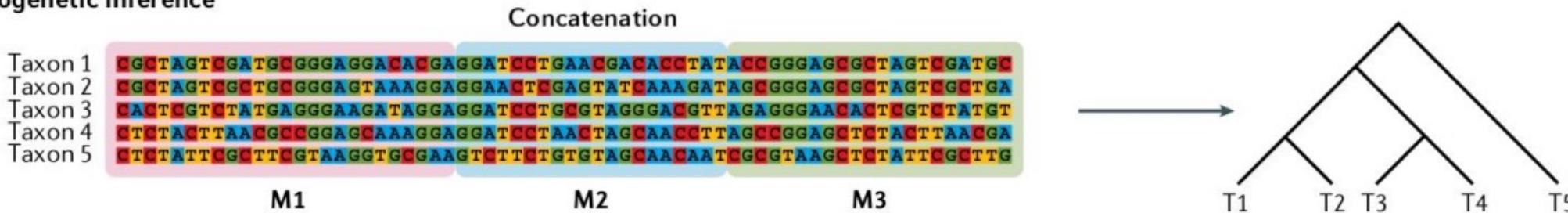




**f Alignment/site filtering**



**g Phylogenetic inference**



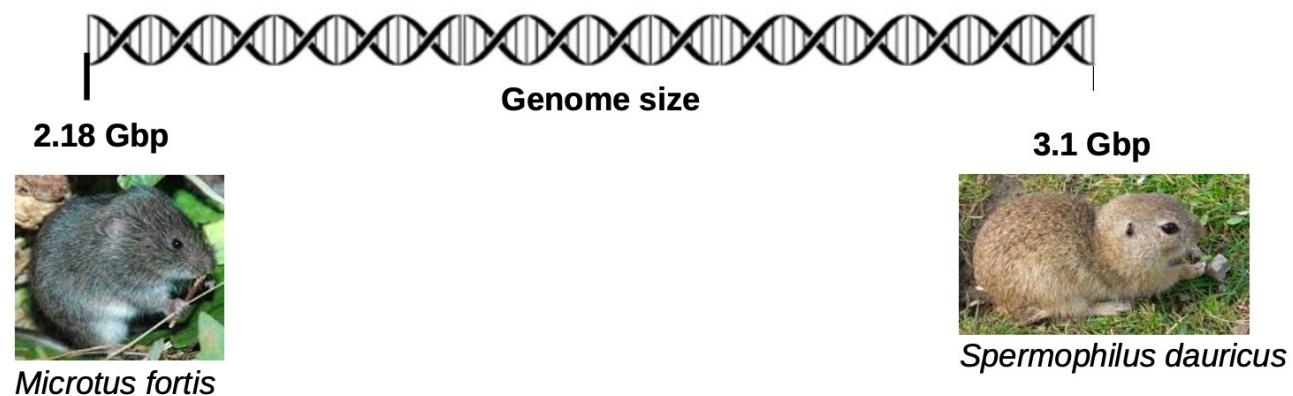
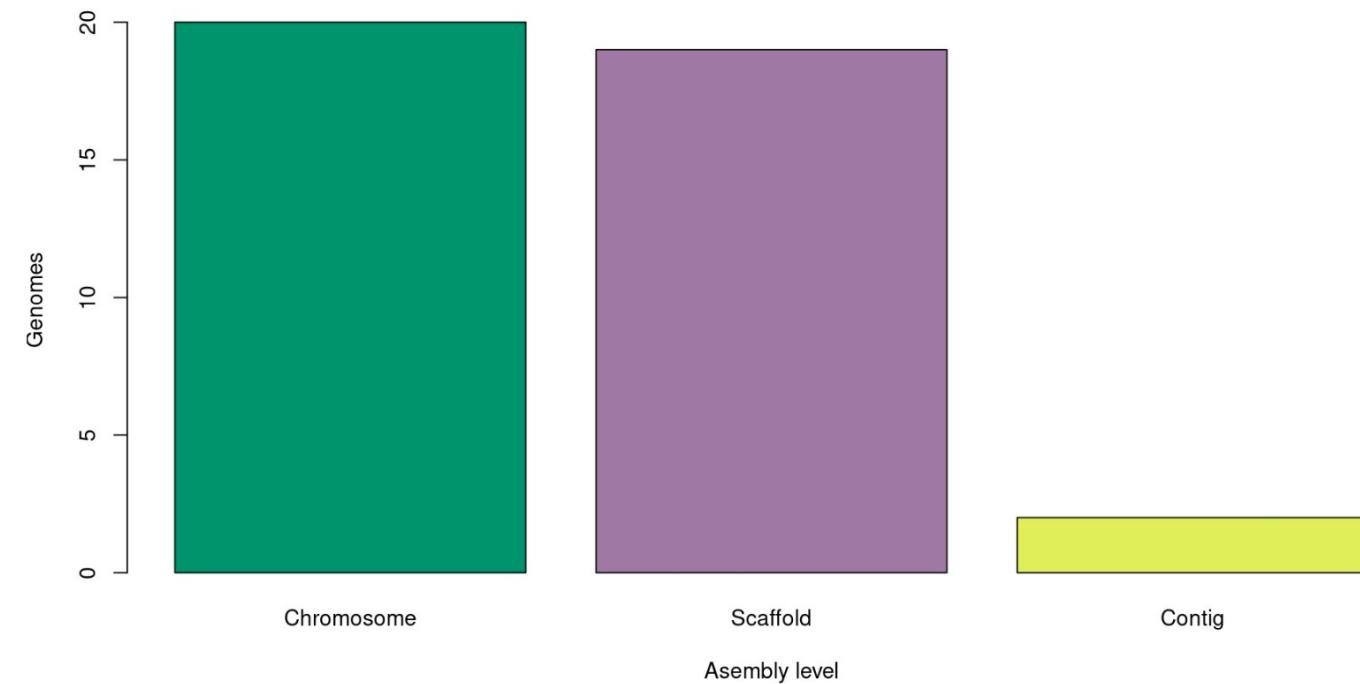
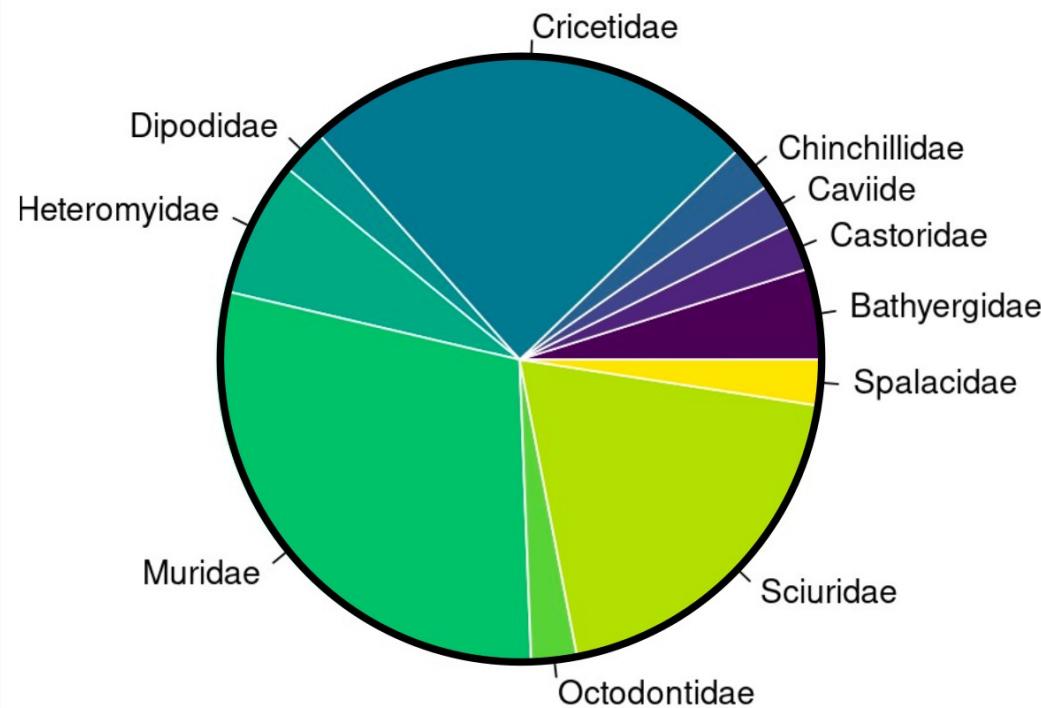
# Phylogenomic exercise

## First Part



Photo credits: inaturalista; mammal.org; wildlife; Wikispecies

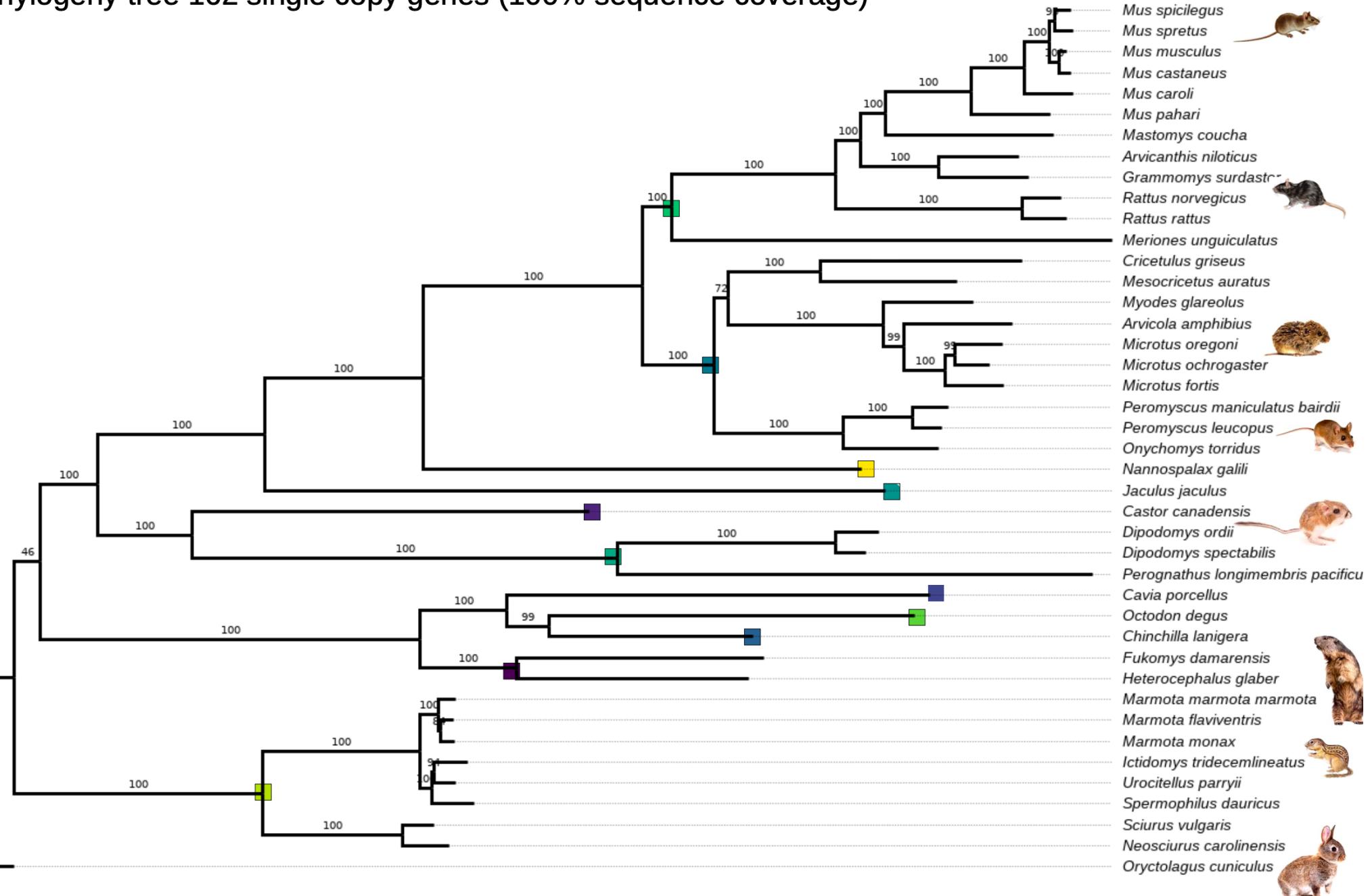
41 rodent species  
11 families



# 12,486 gene families; Phylogeny tree 162 single copy genes (100% sequence coverage)

## Family

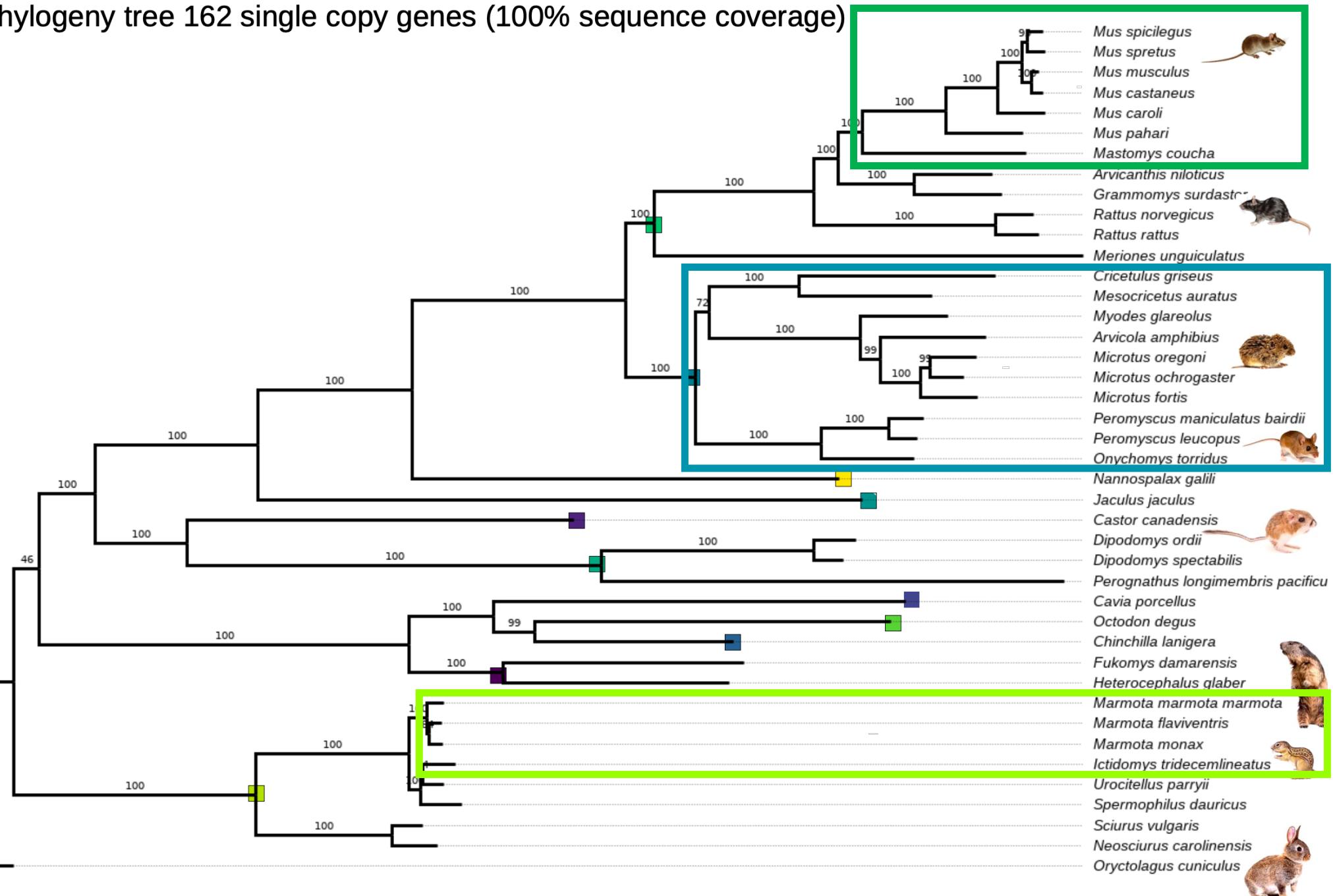
- █ Muridae
- █ Cricetidae
- █ Spalacidae
- █ Dipodidae
- █ Castoridae
- █ Heteromyidae
- █ Octodontidae
- █ Chinchillidae
- █ Caviide
- █ Bathyergidae
- █ Sciuridae



# 12,486 gene families; Phylogeny tree 162 single copy genes (100% sequence coverage)

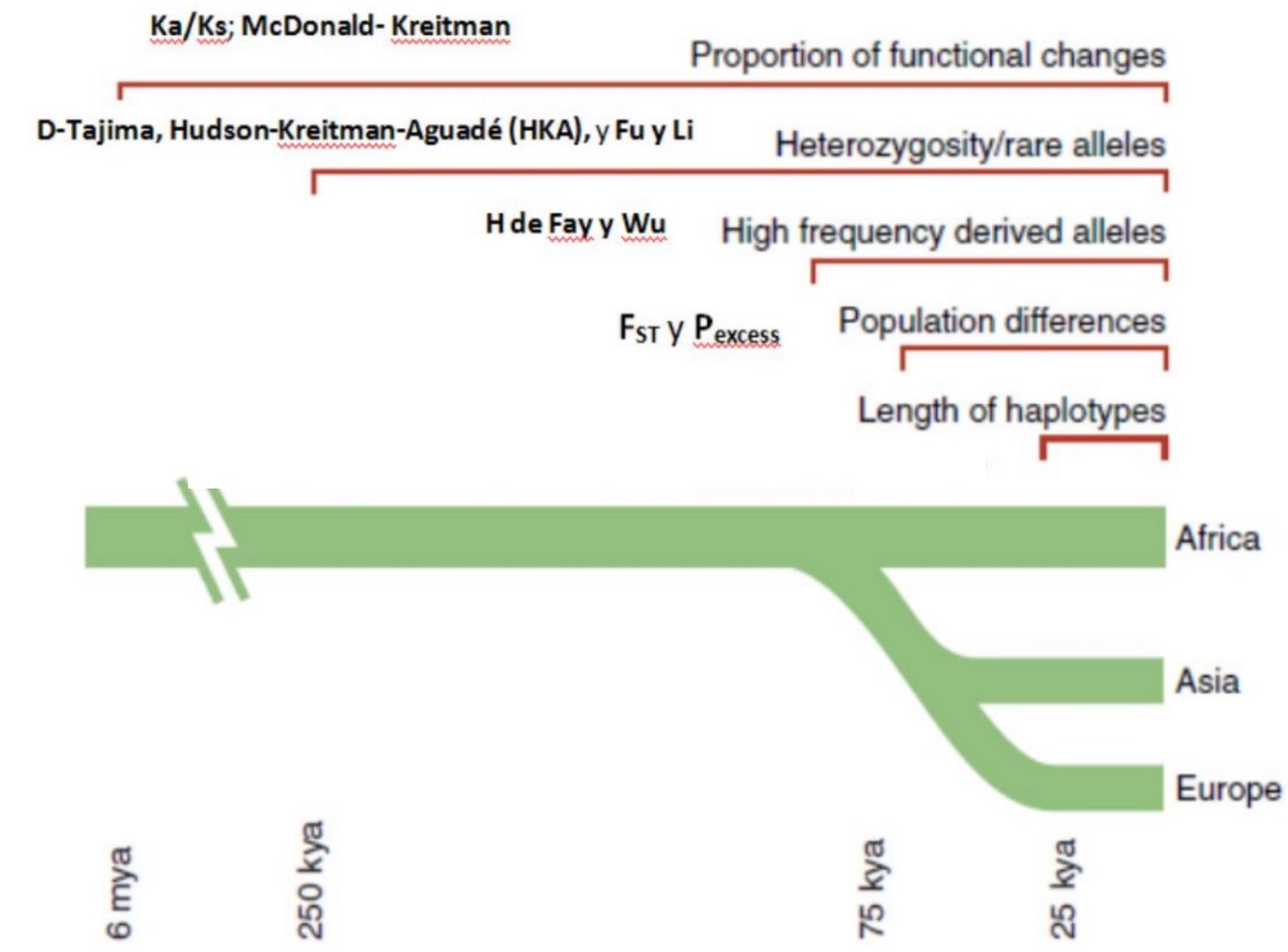
## Family

- Muridae
- Cricetidae
- Spalacidae
- Dipodidae
- Castoridae
- Heteromyidae
- Octodontidae
- Chinchillidae
- Caviide
- Bathyergidae
- Sciuridae



# Selection test over long evolutionary time scales ( $dN/dS$ )

- Different hypothesis and resolution
- Depend on the questions
- Timing and divergence



The ratio of non-synonymous to synonymous substitutions ( $dN/dS$ ) is a statistic measure of the strength and mode of natural selection acting on protein-coding genes in species that have diverged

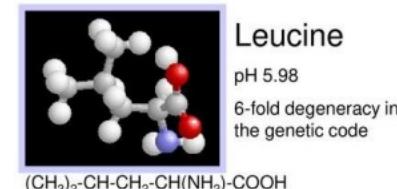
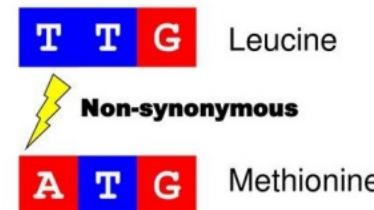
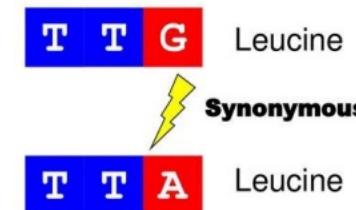
Non-synonymous substitutions are nucleotide changes that alter the protein sequence and might modify protein structure  
Synonymous substitutions do not change the protein sequence

### A Nonsynonymous / Synonymous substitution

TCC	GAT	<u>AT</u> A	TGG	<u>CA</u> A	CCC	<u>GA</u> C	AAA
S	D	I	W	Q	P	D	K

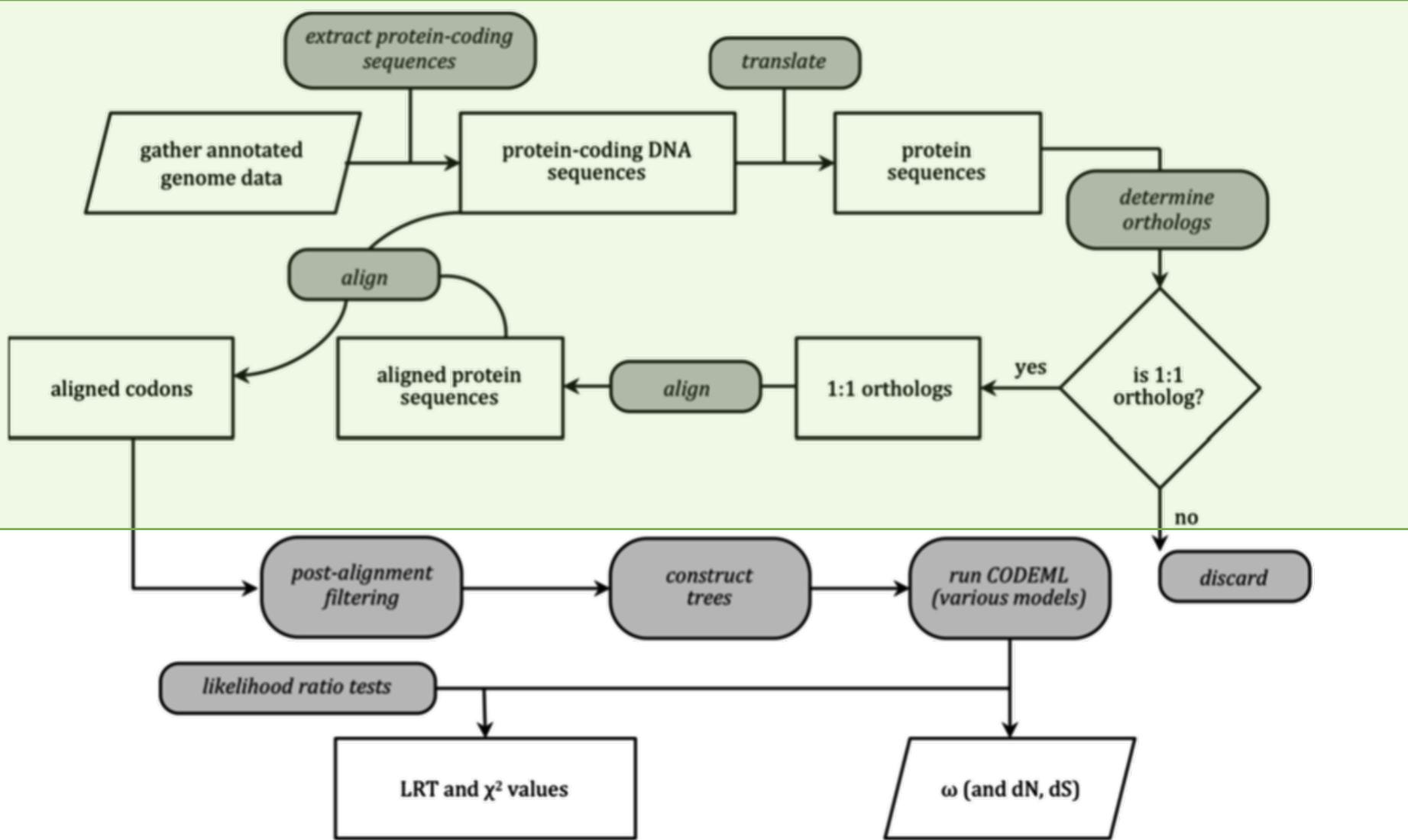
  

TCA	GAT	<u>CT</u> A	TGG	<u>CA</u> G	CCC	<u>CA</u> C	AAA
S	D	L	W	Q	P	R	K



- $dN/dS = 1$  -----> Neutral evolution  
 $dN/dS > 1$  -----> Positive Selection  
 $dN/dS < 1$  -----> Purifying Selection

## Phylogenomic



CODEML calculates  $dN / dS$  using the observed changes present in a multiple alignment of protein coding gene sequences from several species in a phylogeny

Statistical estimation of  $dN / dS$  ratio uses maximum likelihood, and correct for multiple changes, accounting for the different numbers of non-synonymous and synonymous sites

# CODEML

Models	<i>k</i>	Hypothesis tested
Site models:	M2a vs. M1a	$2^a$ Does adding a third class of sites with $\omega > 1$ (adaptive evolution) fit the data better than a model with two classes $\omega < 1$ , $\omega = 1$ ?
	M8 vs. M7	$2^a$ Does adding an extra class of sites with $\omega > 1$ (adaptive evolution) fit the data better than a model with ten classes with flexible normalized non-synonymous ratio distribution?
Branch models:	Free-ratio vs. one-ratio model	$2s-4$ For a tree of $s$ species, is $\omega$ different among lineages?
	Two-ratio vs. one-ratio model	$1$ Are the foreground branches that you specify more likely to have different $\omega$ from background branches?
Branch-site models:	MA( $\omega > 1$ ) vs. MA( $\omega = 1$ )	$1^b$ Is the defined “foreground branch” more likely to contain sites with $\omega > 1$

## HyPhy: hypothesis testing using phylogenies

Sergei L. Kosakovsky Pond ✉, Simon D. W. Frost, Spencer V. Muse



The Site-Wise Log-Likelihood Score is a Good Predictor of Genes under Positive Selection

[Huai-Chun Wang](#) ✉, [Edward Susko](#) & [Andrew J. Roger](#)

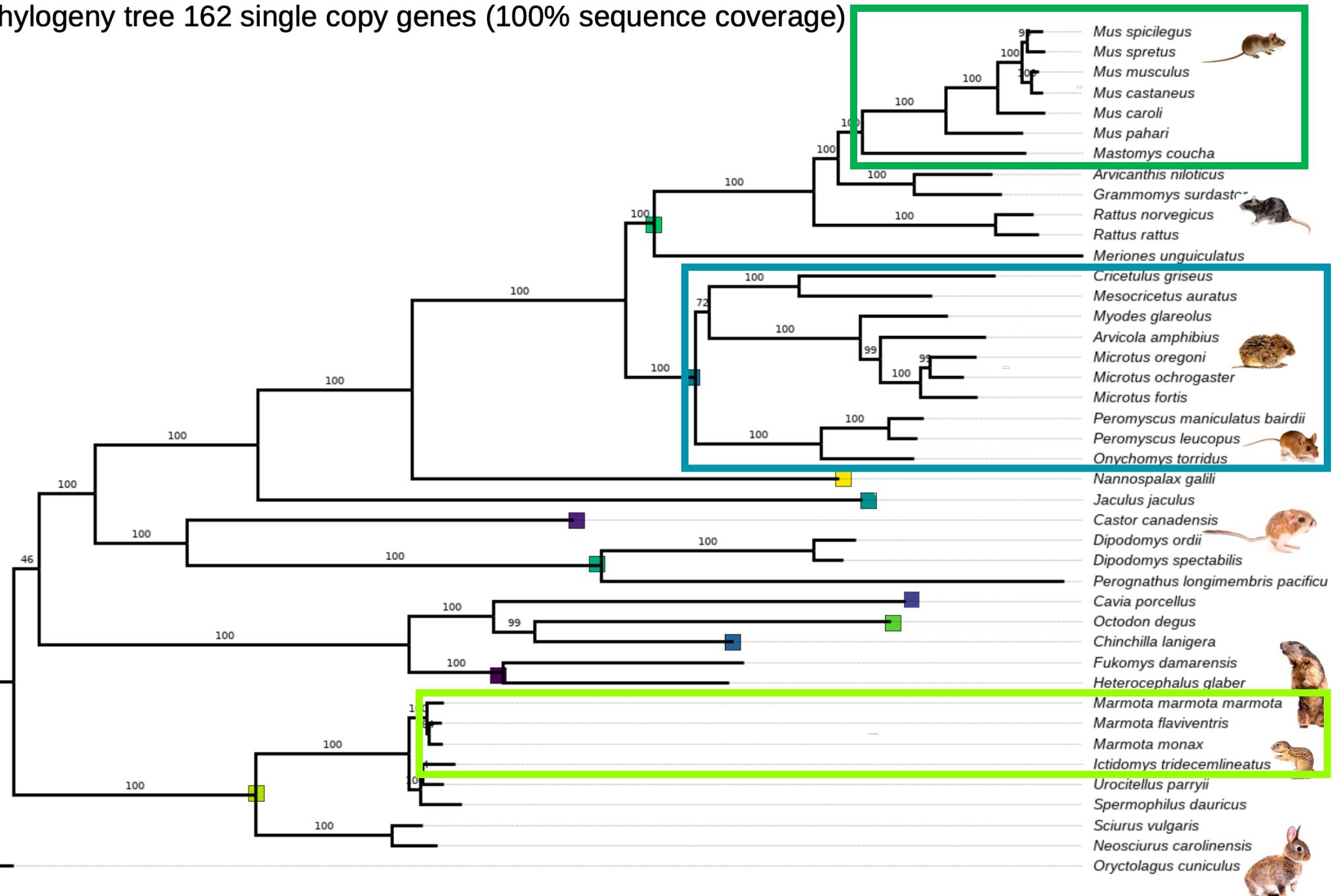
# $dN/dS$ selection test exercise

**Second Part**

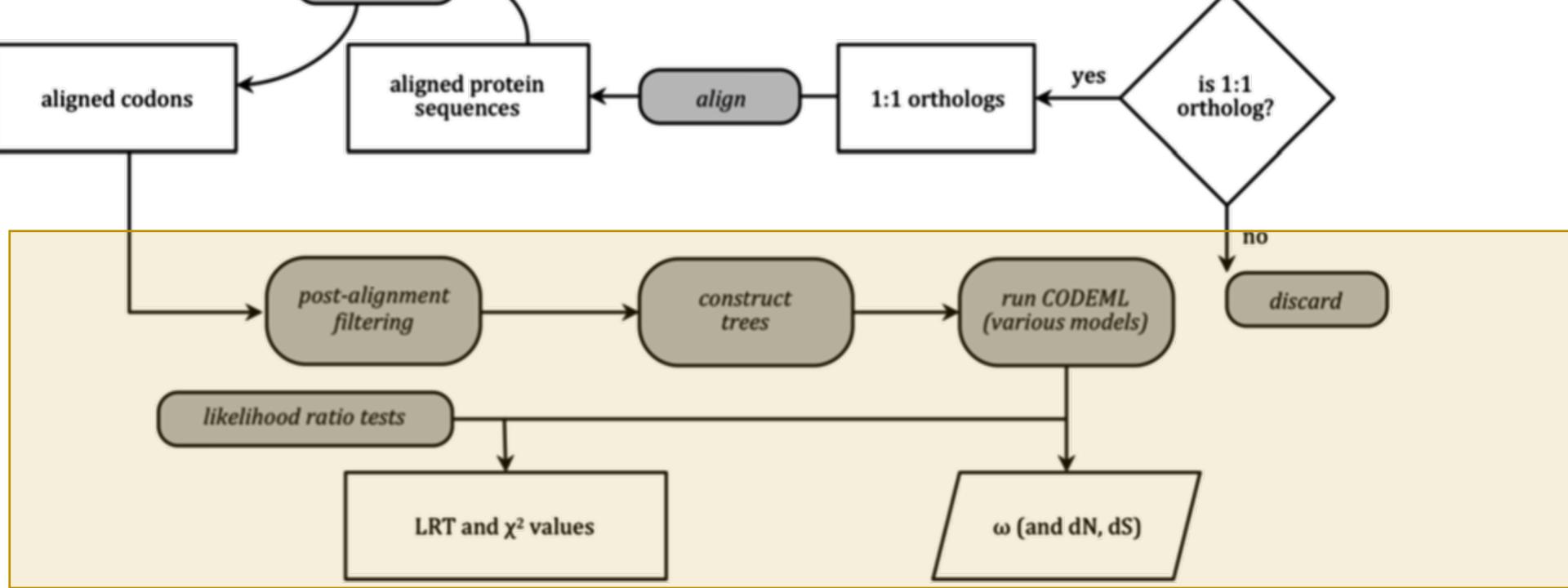
# 12,486 gene families; Phylogeny tree 162 single copy genes (100% sequence coverage)

## Family

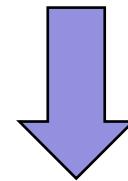
- Muridae
- Cricetidae
- Spalacidae
- Dipodidae
- Castoridae
- Heteromyidae
- Octodontidae
- Chinchillidae
- Caviide
- Bathyergidae
- Sciuridae



## Phylogenomic



## $dN/dS$ test with CODEML



Annotation

and

Gene Enrichment analysis

(topGO, R library)