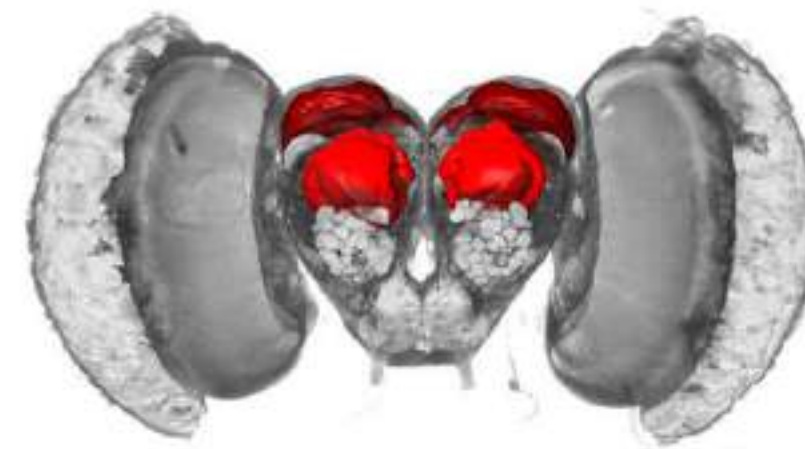# Day 9 - Comparative Genomics
## Attempting to make sense of how evolution works

*F. Cicconardi, PhD*

**EBaB lab**

2025 WORKSHOP ON GENOMICS, CESKY KRUMLOV
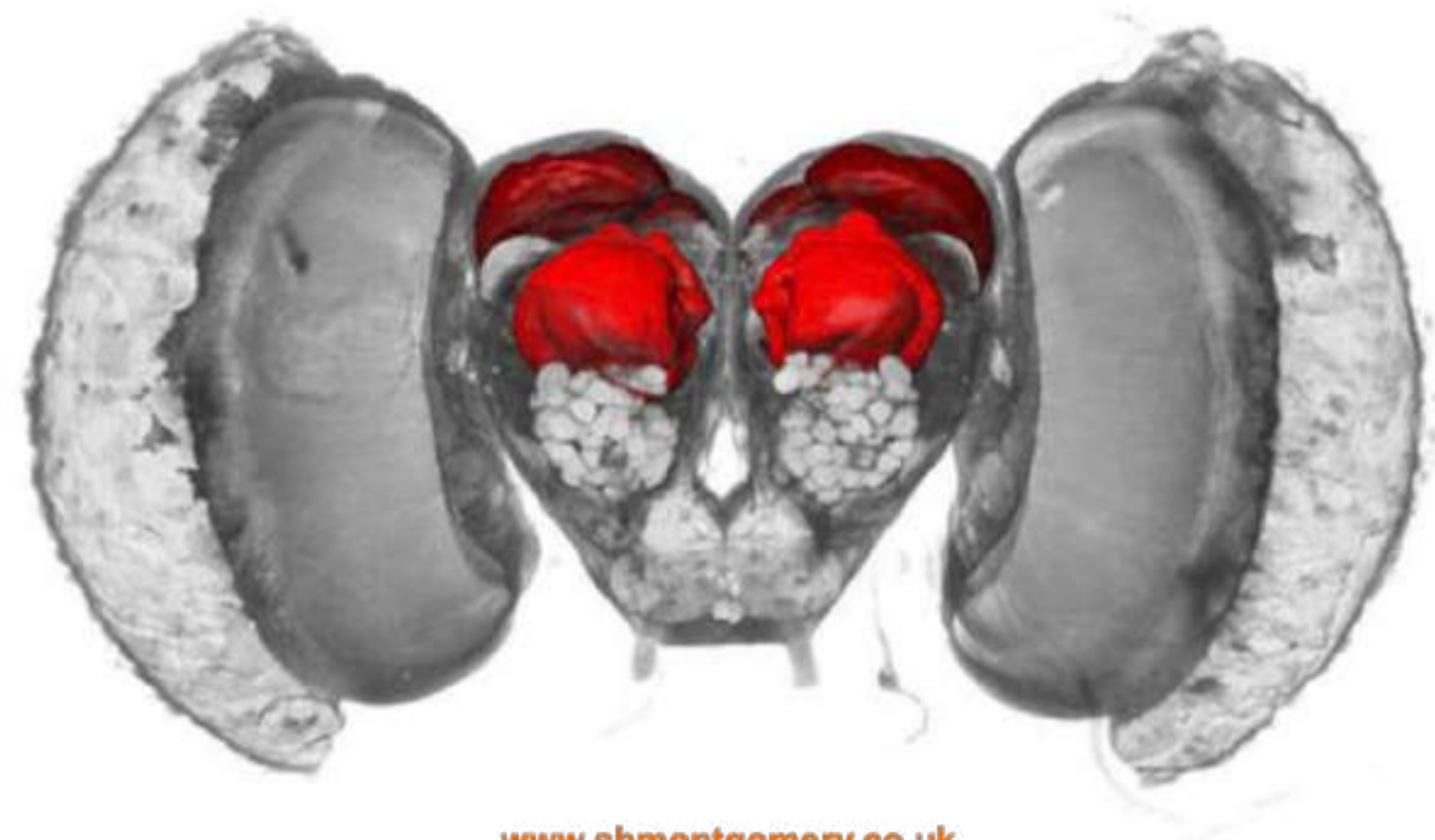
# Evolution of Brain and Behaviour Lab

EBAB LAB

*Dr. Stephen Montgomery*

www.shmontgomery.co.uk

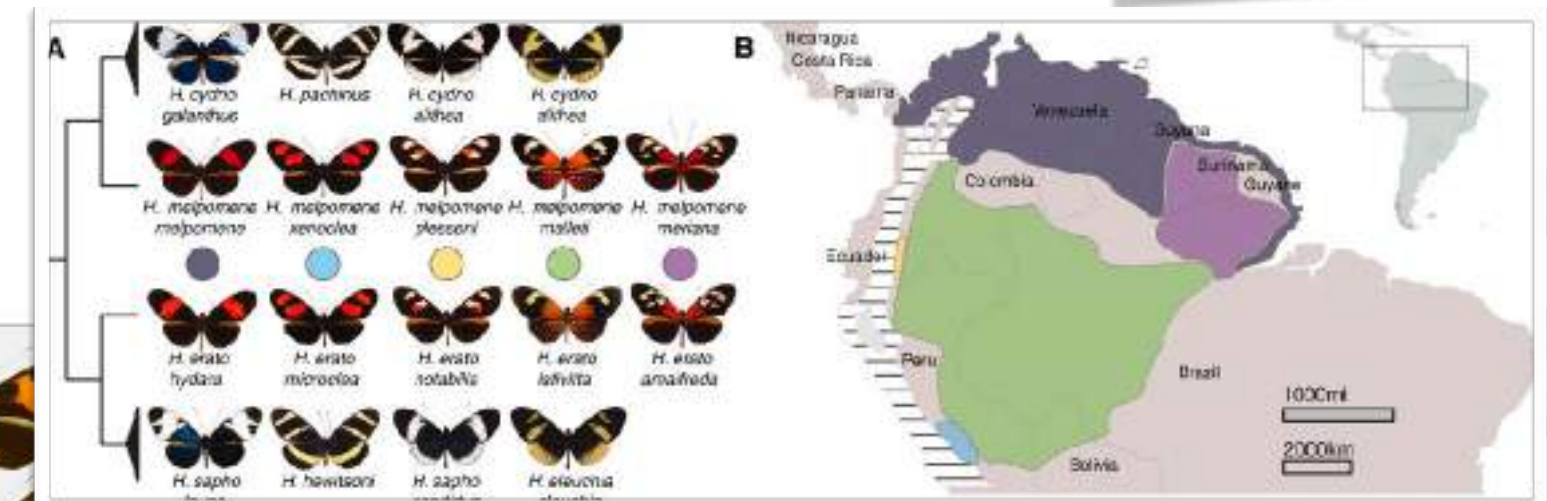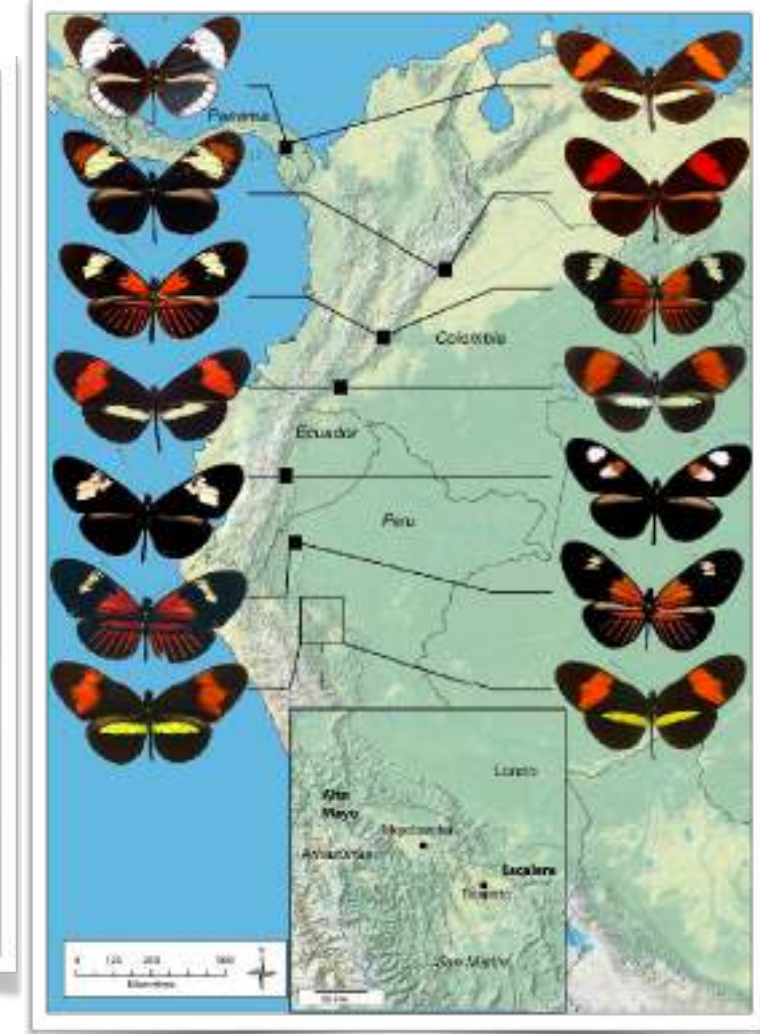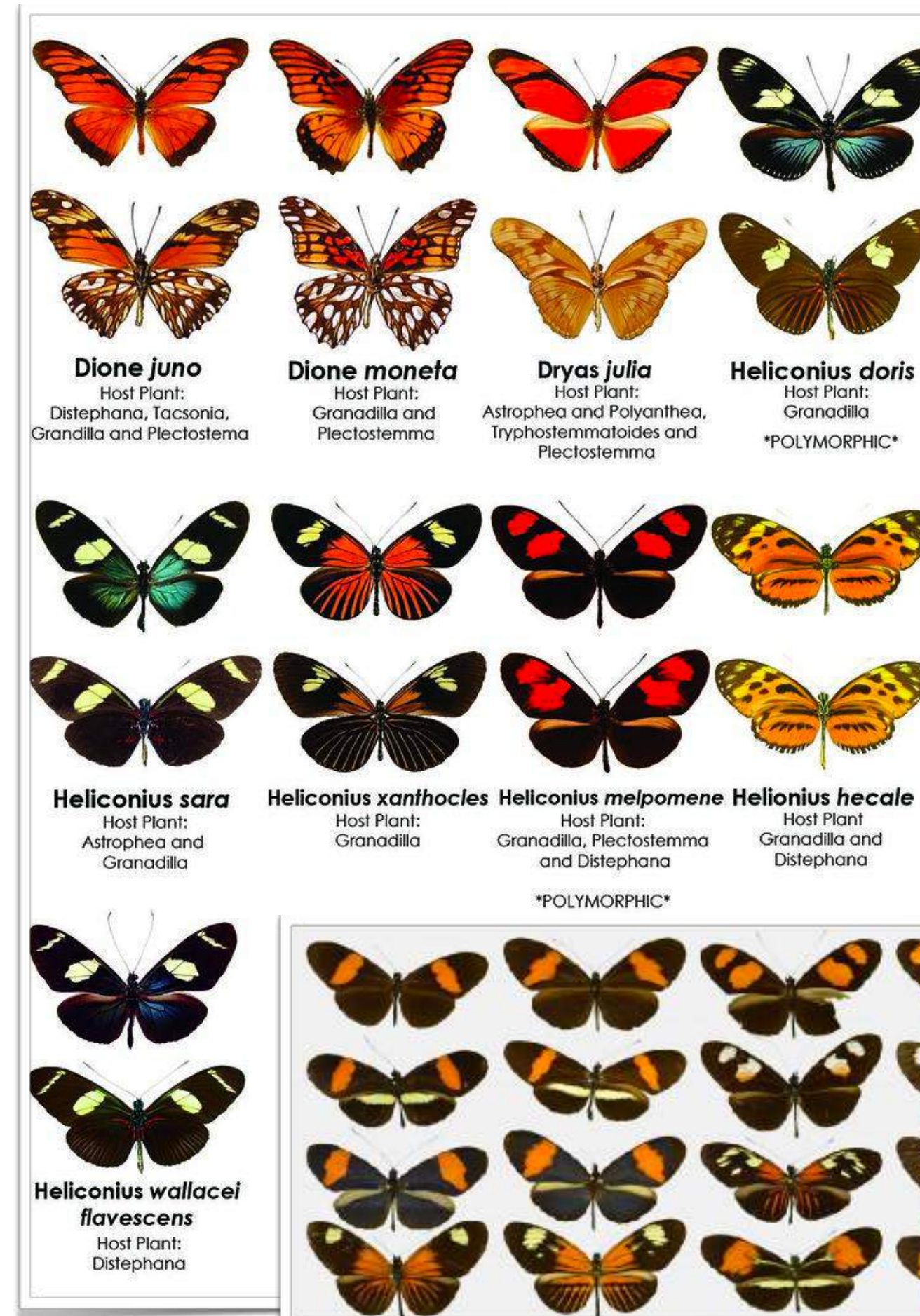# Adaptive Radiation of Heliconiini (Family: Nymphalidae)
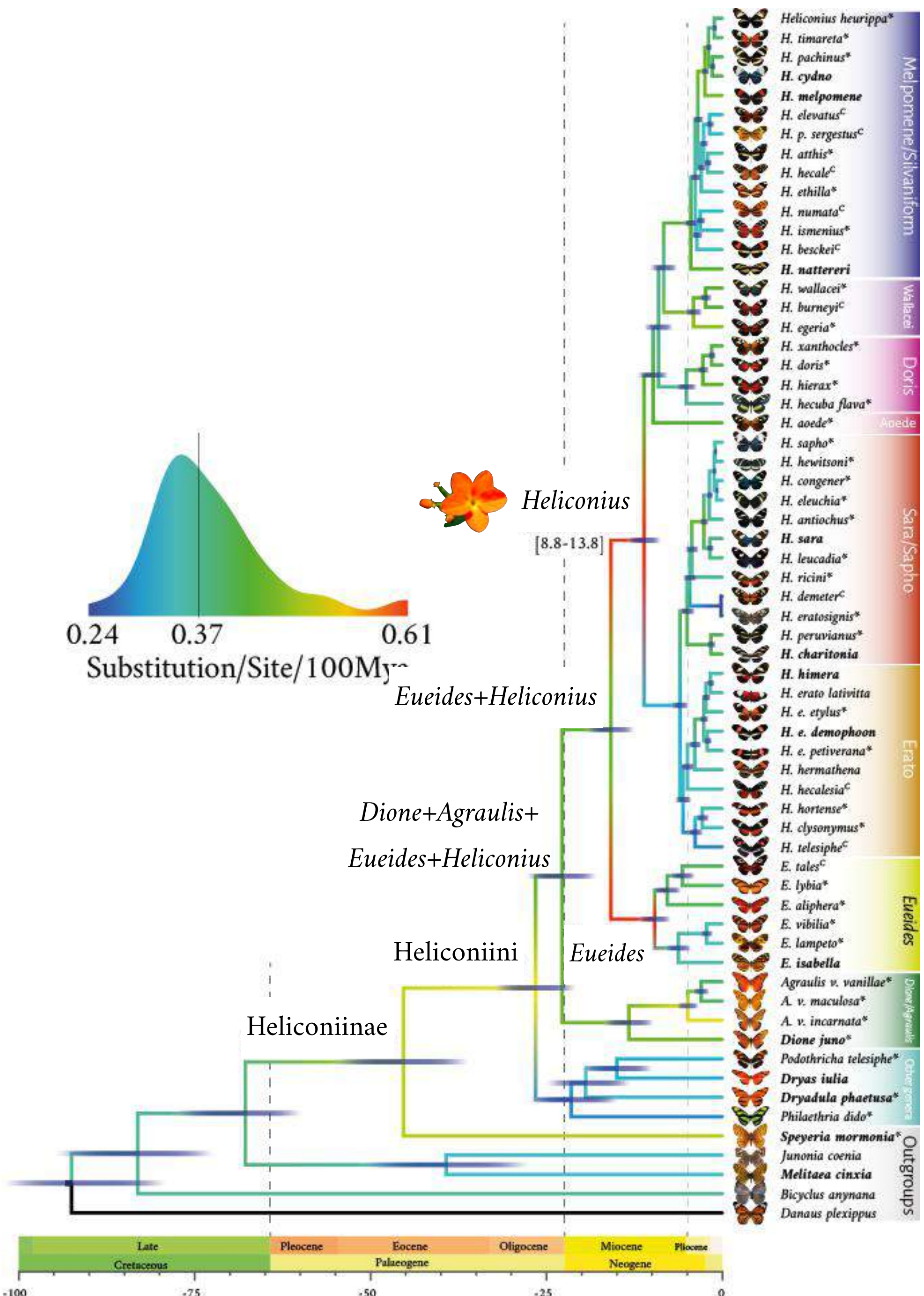


Photo credit: @mena_sebas

- *8 Genera*

- *87 Species*

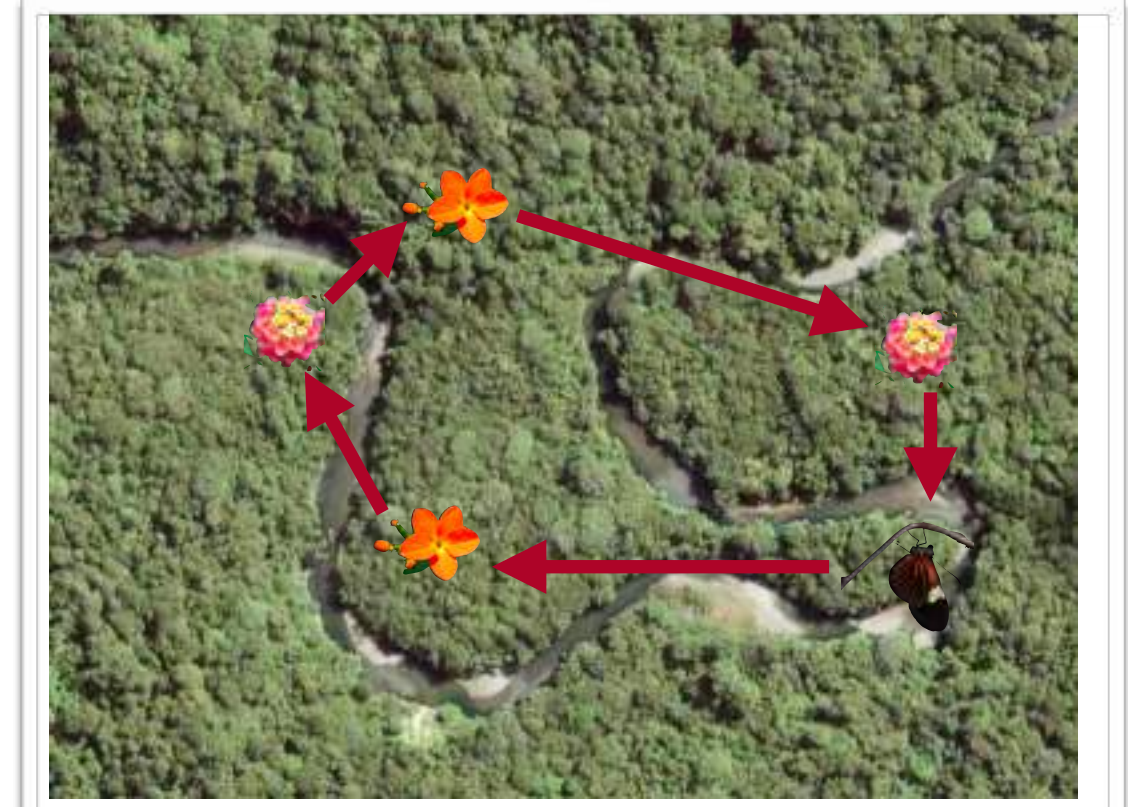- *~440 sub-species*
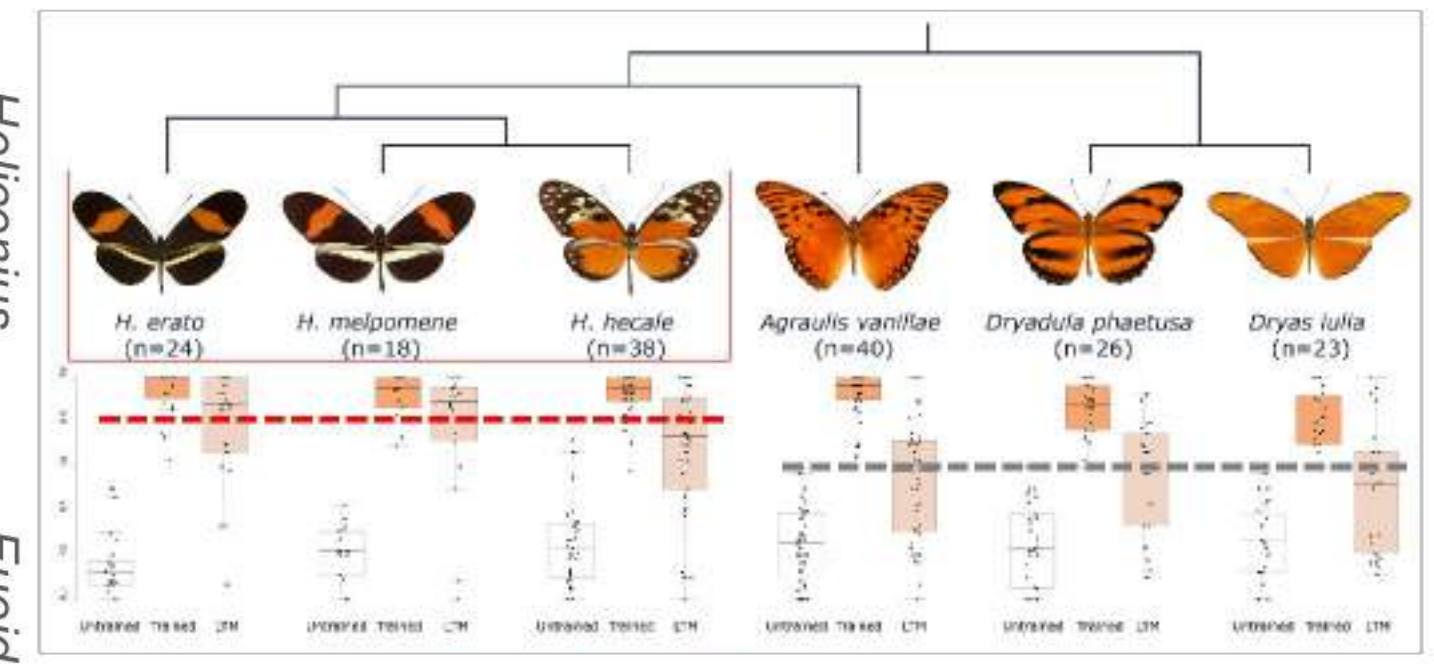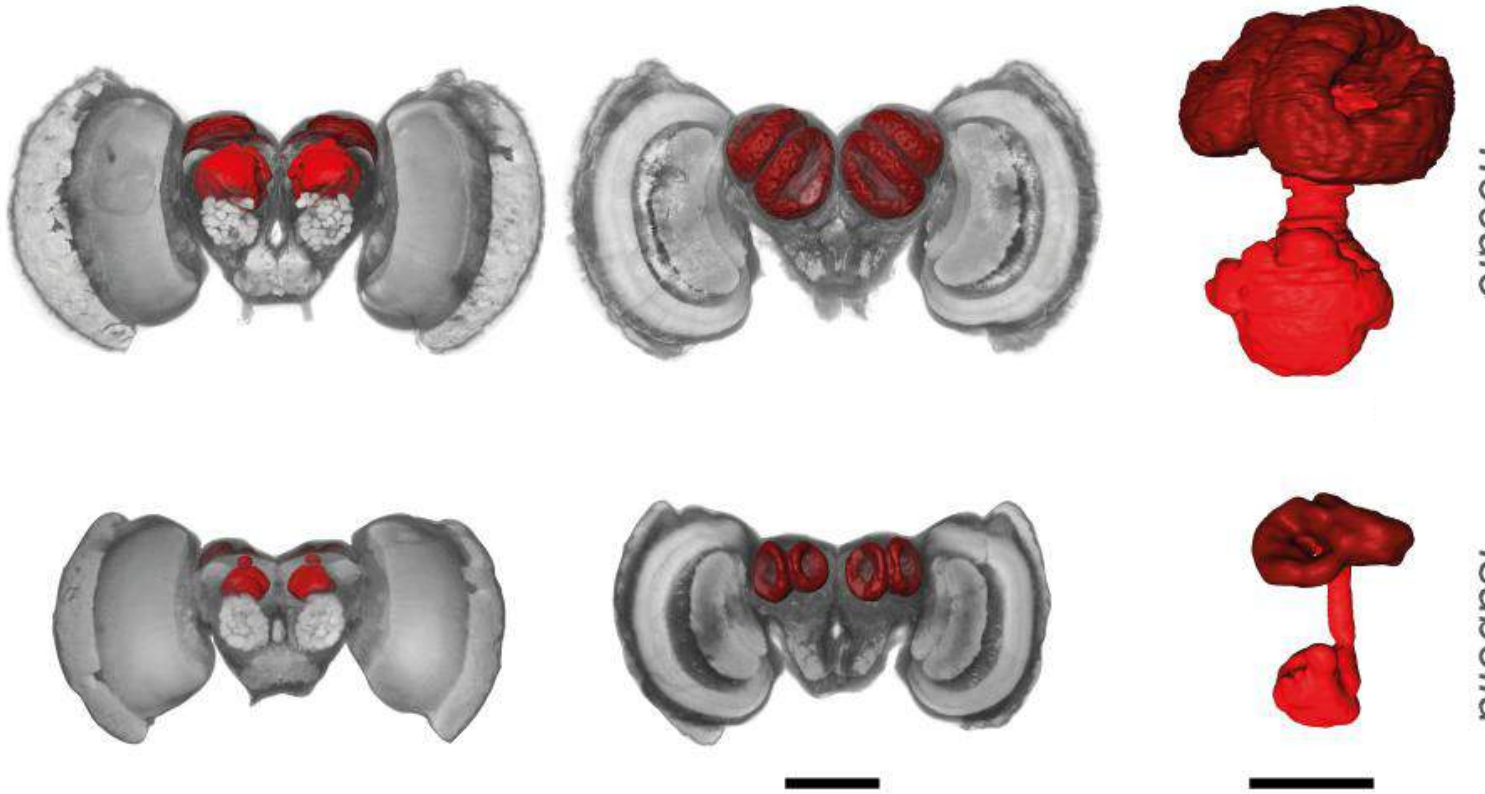
» *Pollen-feeding*

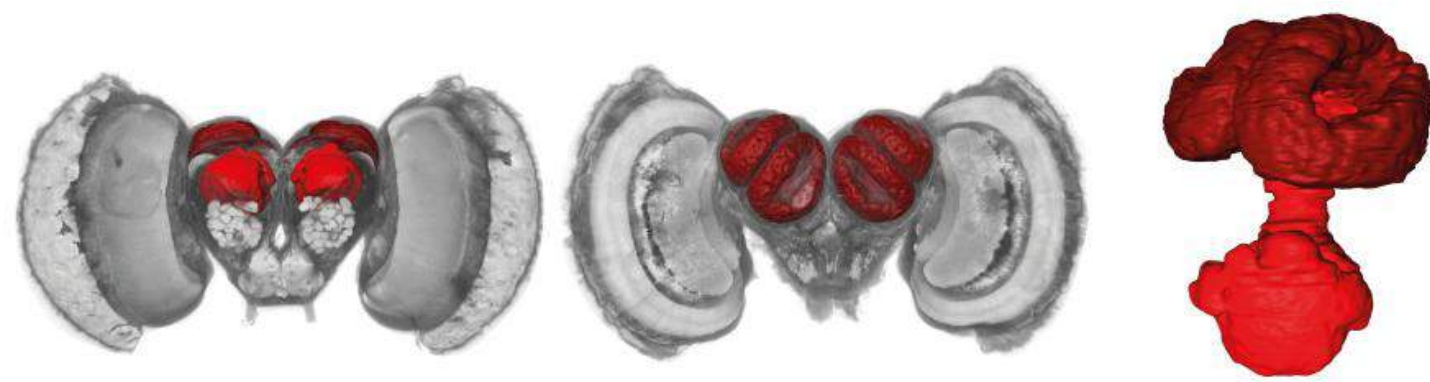Photo credit: @mena_sebas

» *Trap-lining behaviour*

Allocentric trap-lining in home ranges up to 1 km²

Moura et al.. **2021**. *Functional Ecoloav.*

» *Mushroom body expansion (Brain)*

Interspecific variation: $X_2$ = 76.2, p < 0.001
Small vs. big MBs: $X_8$ = 182.8, p < 0.001

*Heliconius*

non-*Heliconius*

*H. melpomene*

*H. e. demophoon*

*Agraaulis v. vanillae*

*Dryas iulia*

ATACseq

ATACseq

ATACseq

ATACseq

L3

L4

L5

*Larva*

*Pupa*

*Adult*

A...

RNAseq

PP

P12h

P24h

P36h

P5.5d

A12h

RNAseq

RNAseq

RNAseq

RNAseq

Stacked Histogram of Conservation Levels by Species

1.00

0.75

0.50

0.25

0.00

Count

Conservation

NotConserved
NotConsPeak(1)
NotConsPeak(2)
NotConsPeak(3)
ConsReag(1)
ConsReag(2)
ConsReag(3)
ConsPeak(1)
ConsPeak(2)
ConsPeak(3)

Hmel

Herd

Avcr

Diul

Species

*Heliconius*

non-*Heliconius*

225,859 nuclei of *H. melpomene* brains

176,804 nuclei of *D. iulia* brains

» But how do you tackle these problems? »

» How do we make sense of how evolution works? »

» But how do you tackle these problems? »

» How do we make sense of how evolution works? »

» **Comparing "things"!** »

# Comparative Genomics

*Researchers choose the appropriate time-scale of evolutionary conservation for the question being addressed.*

**Common features of different organisms** such as humans and fish are often encoded within the DNA evolutionarily conserved between them.

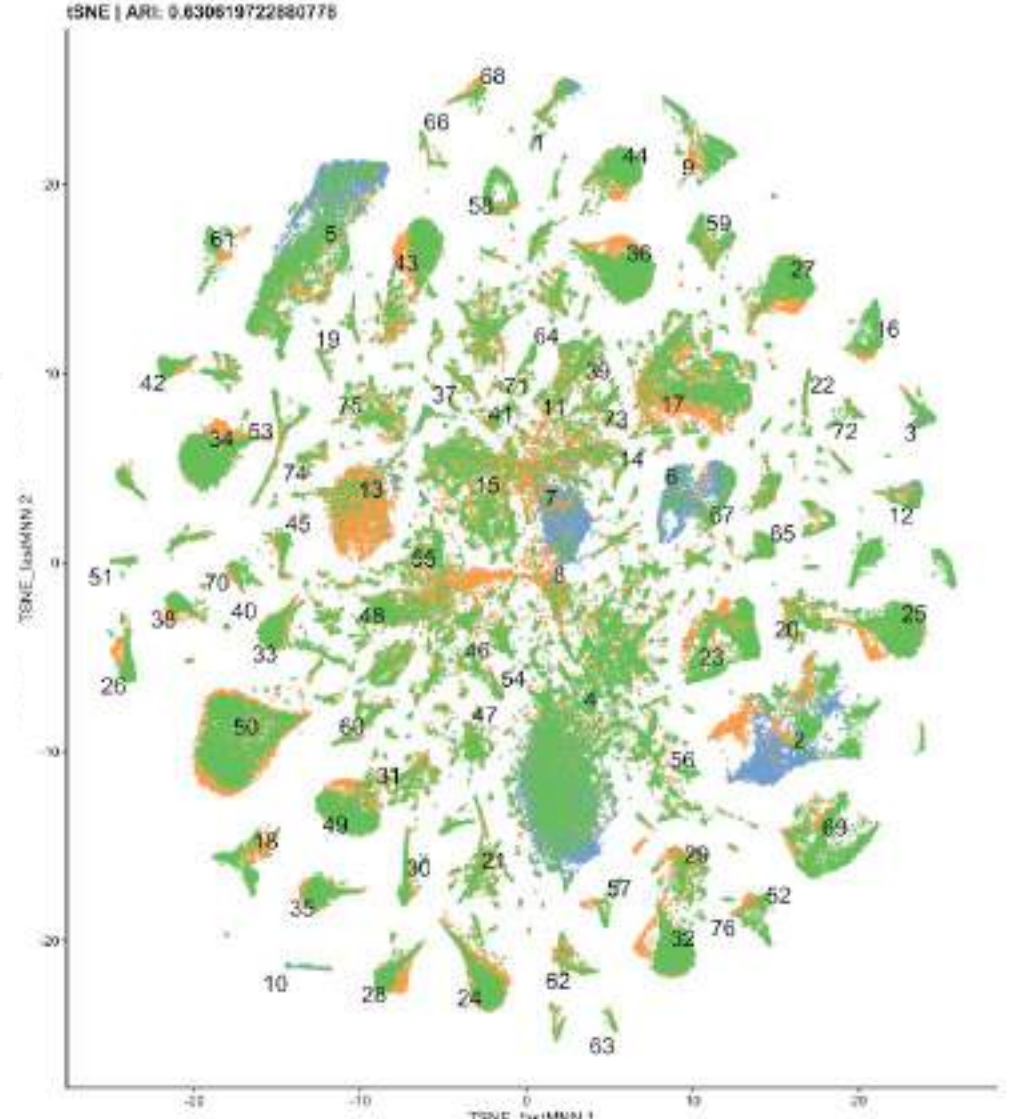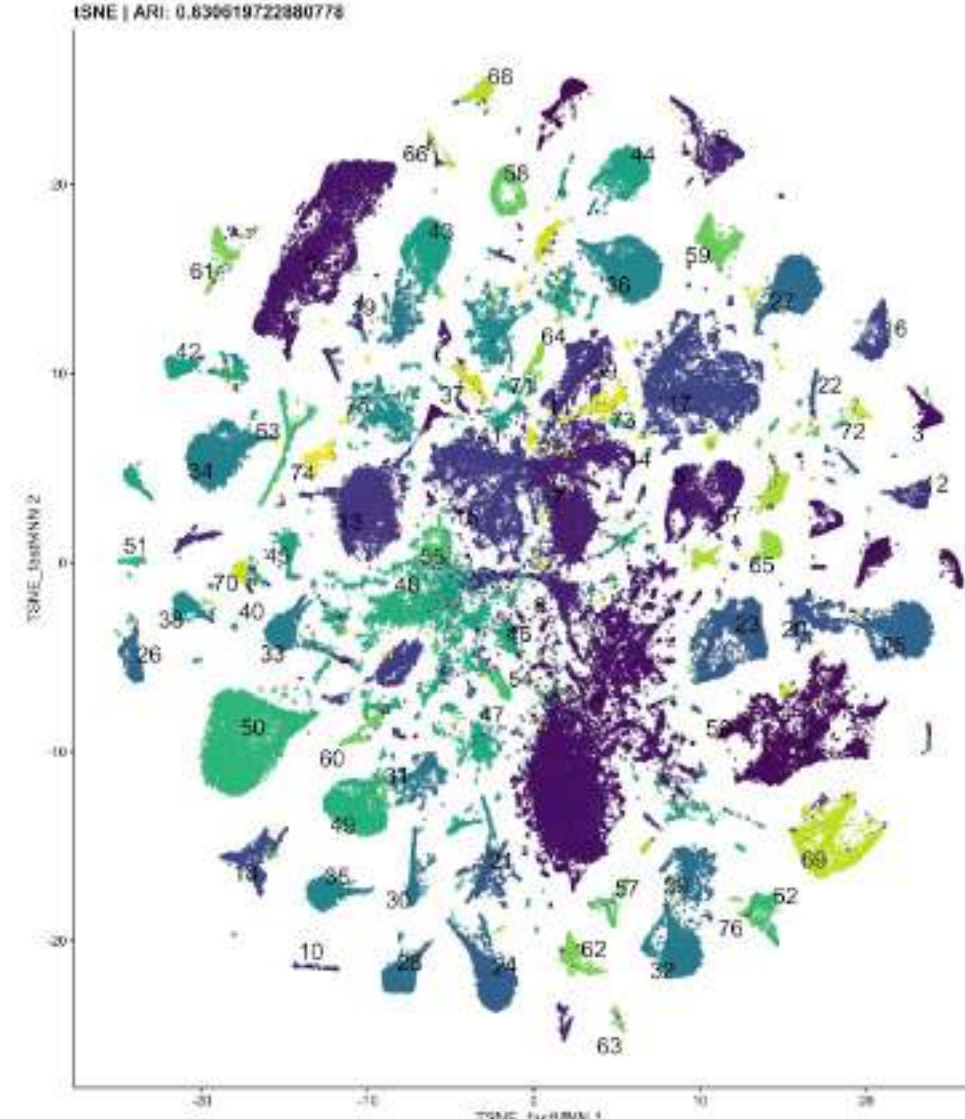Looking at **closely related species** such as humans and chimpanzees shows which genomic elements are unique to each.

Genetic differences **within one species** such as our own can reveal variants with a role in disease.

**NIH** National Human Genome Research Institute

# Comparative genomics reveals the dynamics of chromosome evolution in Lepidoptera

Charlotte J. Wright [1], Lewis Stevens [1], Alexander Mackintosh [2], Mara Lawniczak [1] & Mark Blaxter [1]
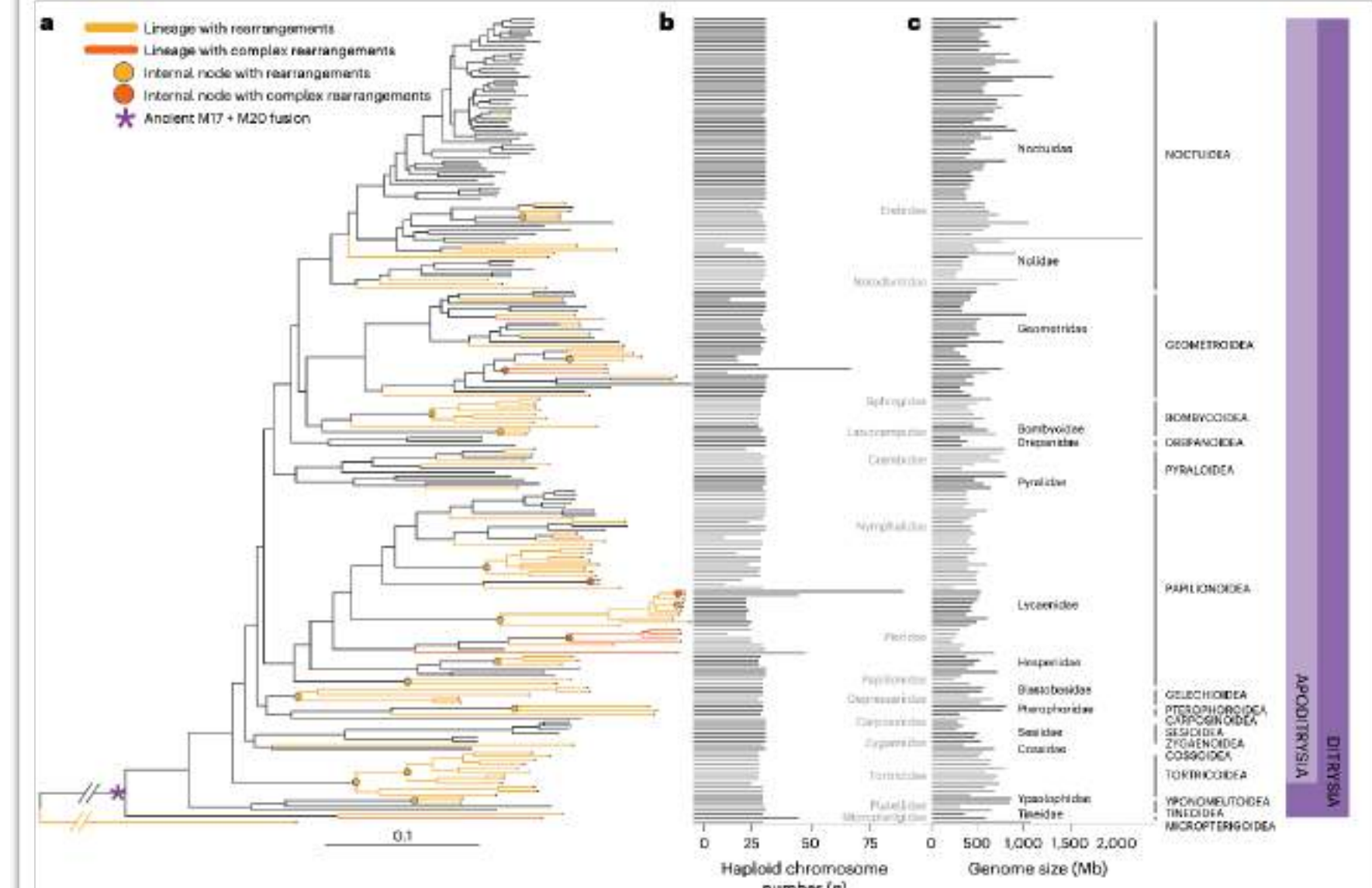


---
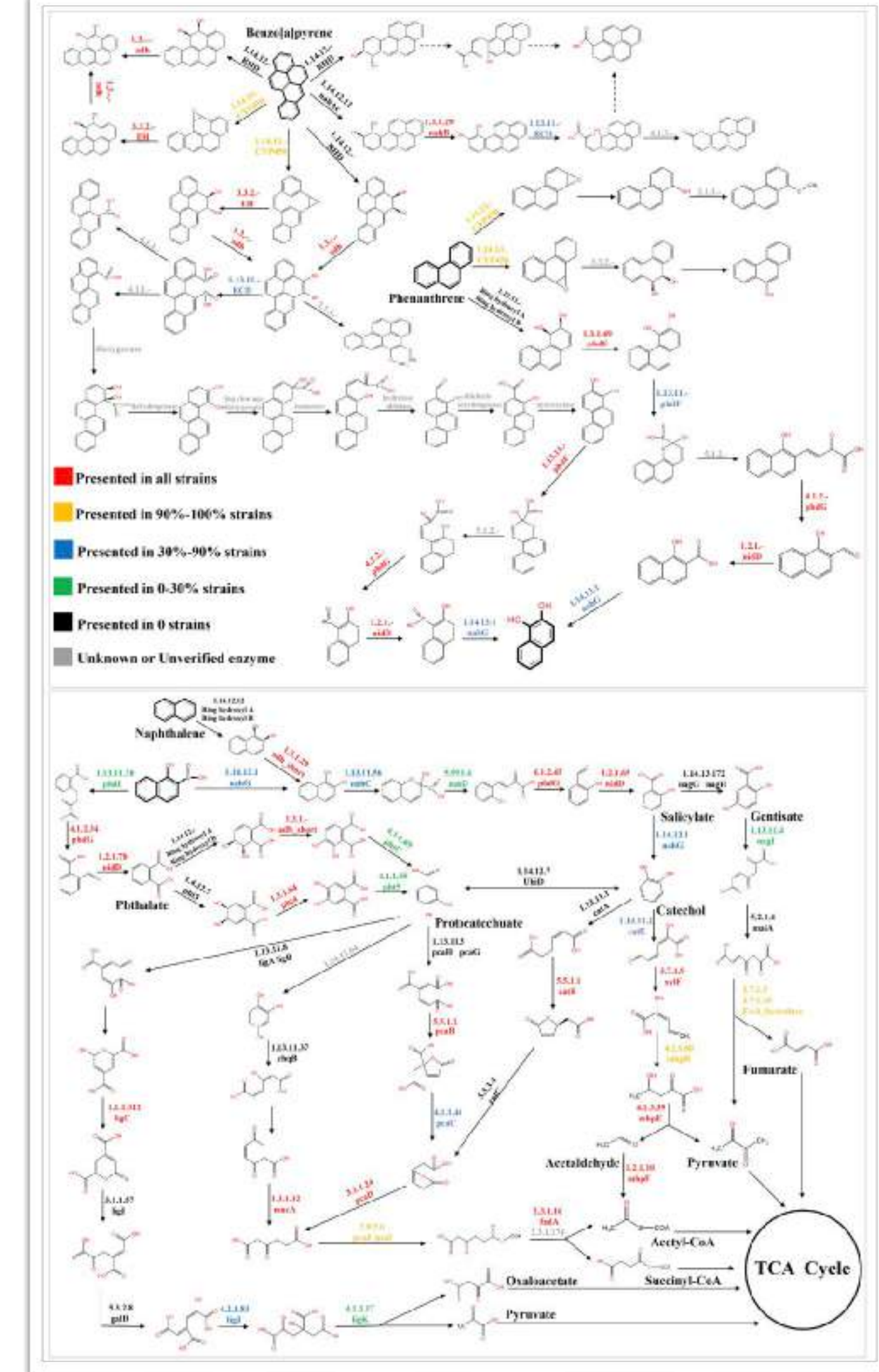
# Comparative genomics reveals evidence of polycyclic aromatic hydrocarbon degradation in the moderately halophilic genus *Pontibacillus*

Haichen Yang [a], Zhihui Qian [a], Yongjin Liu [a], Fei Yu [a], Tongwang Huang [a], Bing Zhang [a], Tao Peng [a,*], Zhong Hu [a,b,*]

[a] Department of Biology, Shantou University, Shantou, Guangdong 515063, PR China
[b] Guangdong Research Center of Offshore Environmental Pollution Control Engineering, Shantou University, Shantou 515063, Guangdong, PR China



---

# Population comparative genomics discovers gene gain and loss during grapevine domestication

Qiming Long [1,†], Shuo Cao [1,2,†], Guizhou Huang [1], Xu Wang [1,3], Zhongjie Liu [1], Wenwen Liu [1], Yiwen Wang [1], Hua Xiao [1], Yanling Peng [1,*] and Yongfeng Zhou [1,4,*]

# » … Some definitions … »

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » Homology/Orthology definition »

**Homology:** it describes descent from a <u>common evolutionary origin</u>: *two loci (genes) are homologous if they derive from the same ancestral locus (gene).*
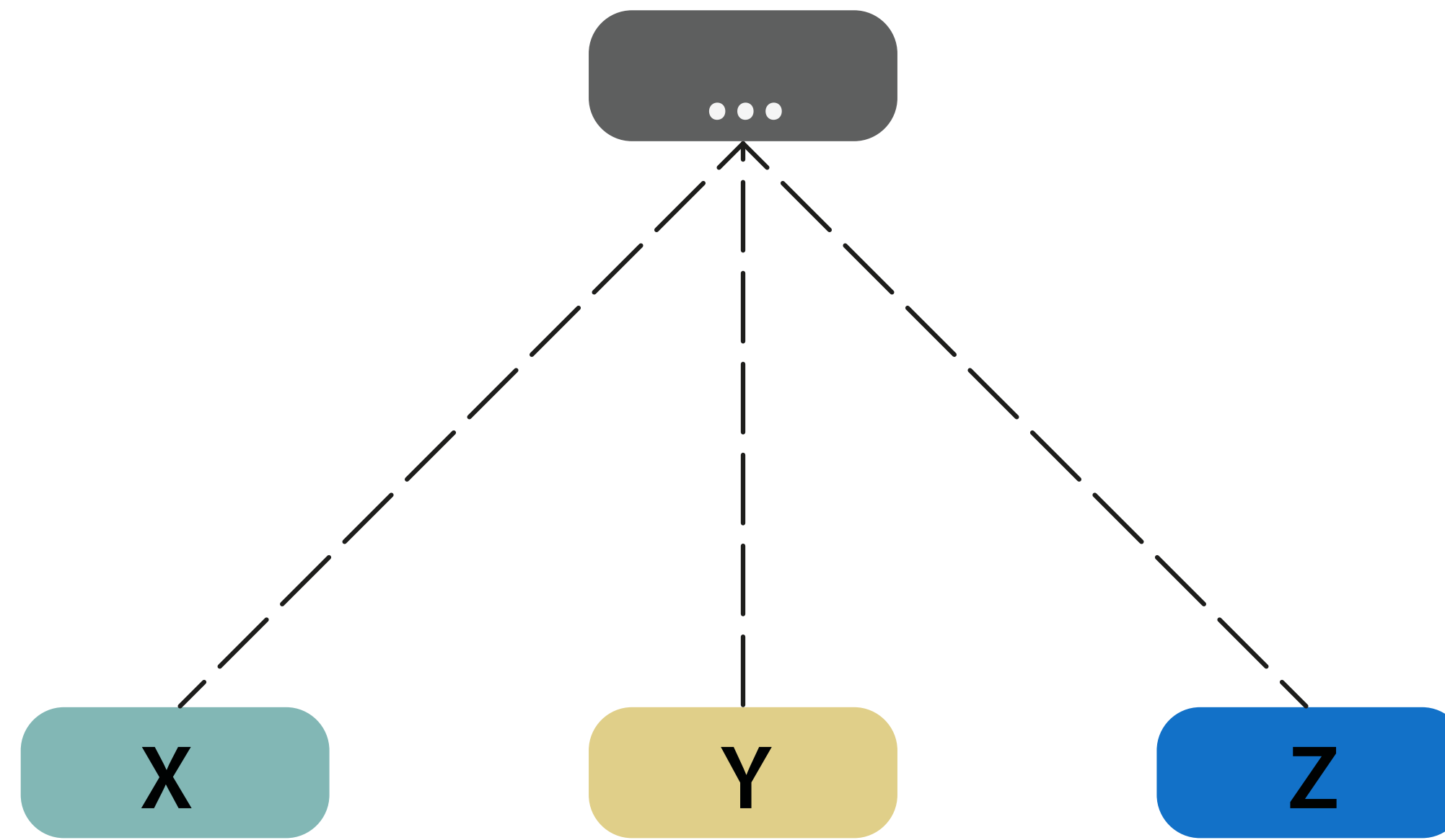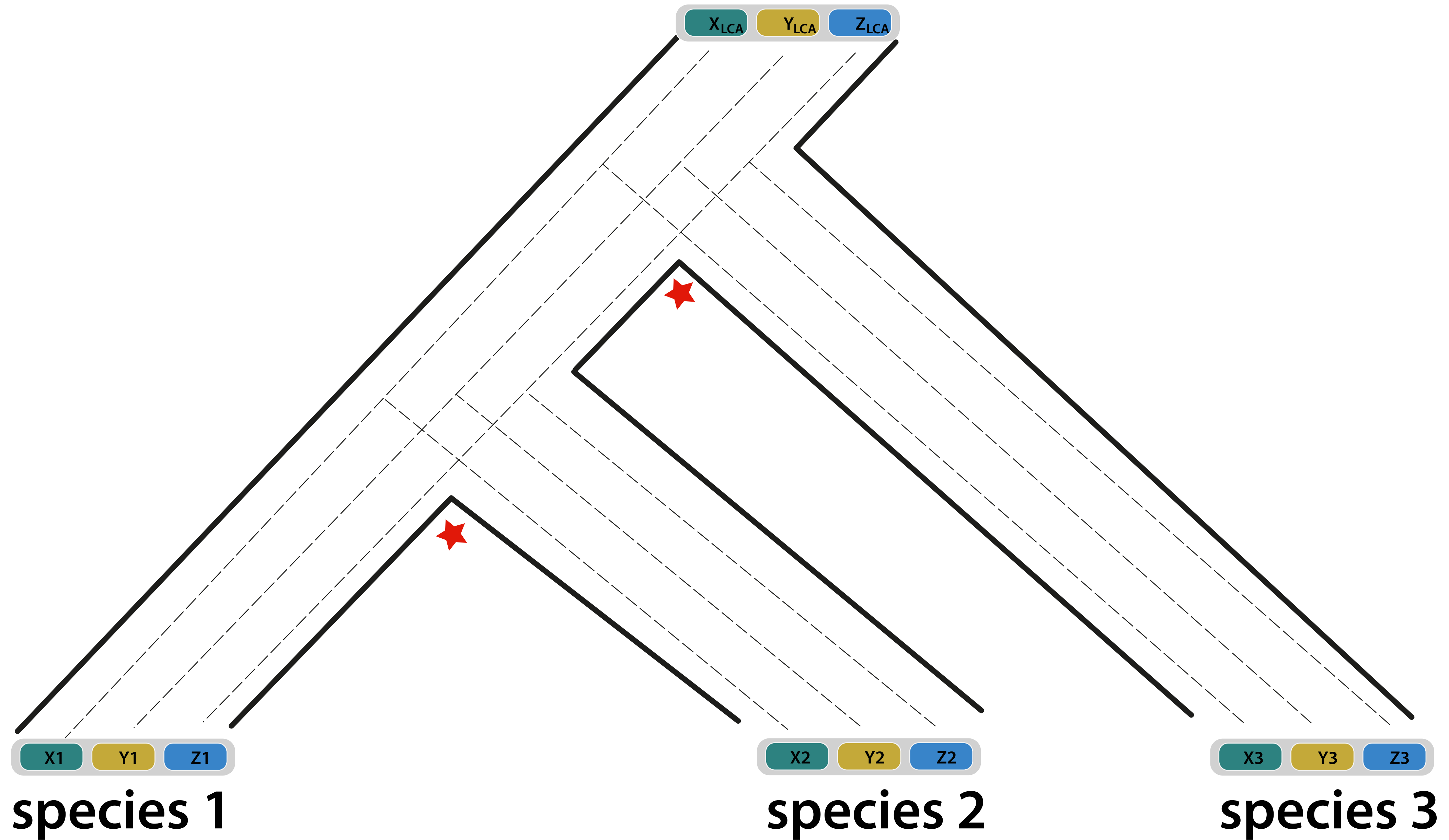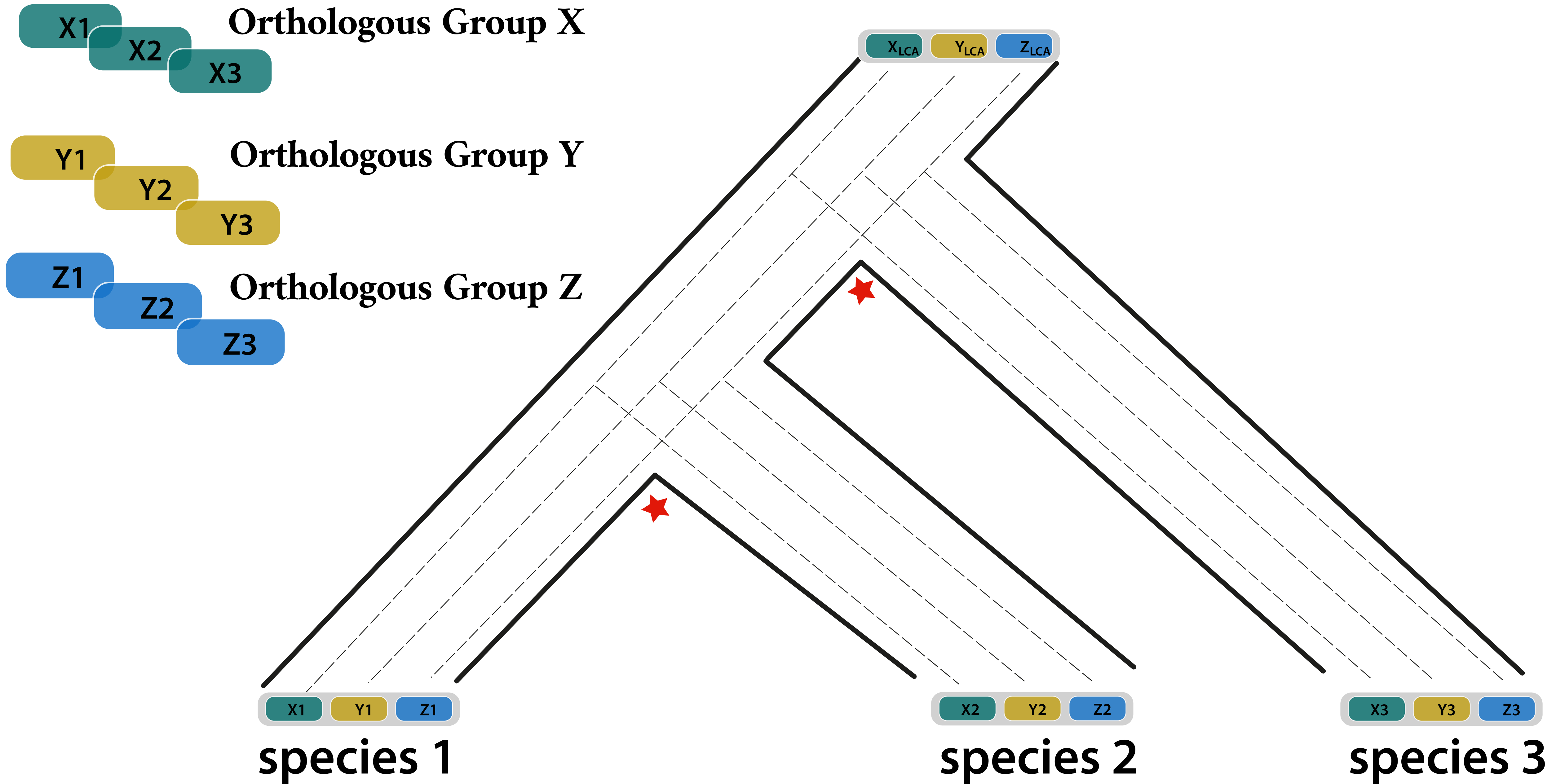
Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » Homology/Orthology definition »
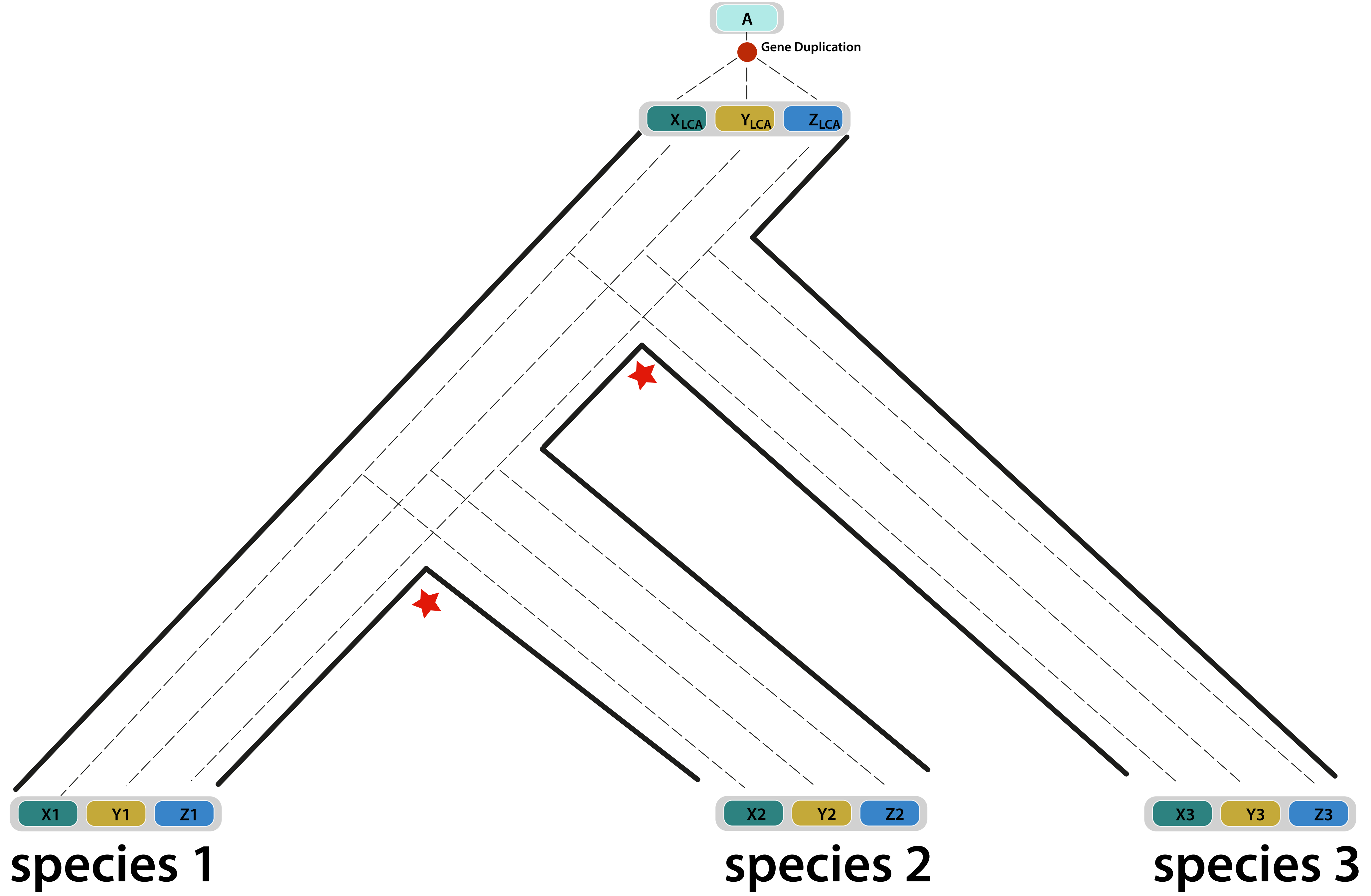
**Orthology:** it describes a relationship due to a <u>speciation event</u>:
*two loci (genes) are orthologous if they derive from speciation.*

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

Orthologous Group X

X1 X2 X3

Orthologous Group Y

Y1 Y2 Y3

Orthologous Group Z

Z1 Z2 Z3

$X_{LCA}$  $Y_{LCA}$  $Z_{LCA}$

X1  Y1  Z1

X2  Y2  Z2

X3  Y3  Z3

species 1

species 2

species 3

A

Gene Duplication

$X_{LCA}$  $Y_{LCA}$  $Z_{LCA}$

X1  Y1  Z1

X2  Y2  Z2

X3  Y3  Z3

species 1

species 2

species 3

Orthologous Group X

X1 X2 X3

Orthologous Group Y

Y1 Y2 Y3

Orthologous Group Z

Z1 Z2 Z3

All Homologous!!!
All Single-copy!!!

A

Gene Duplication

X_LCA  Y_LCA  Z_LCA

X1 Y1 Z1

X2 Y2 Z2

X3 Y3 Z3

species 1

species 2

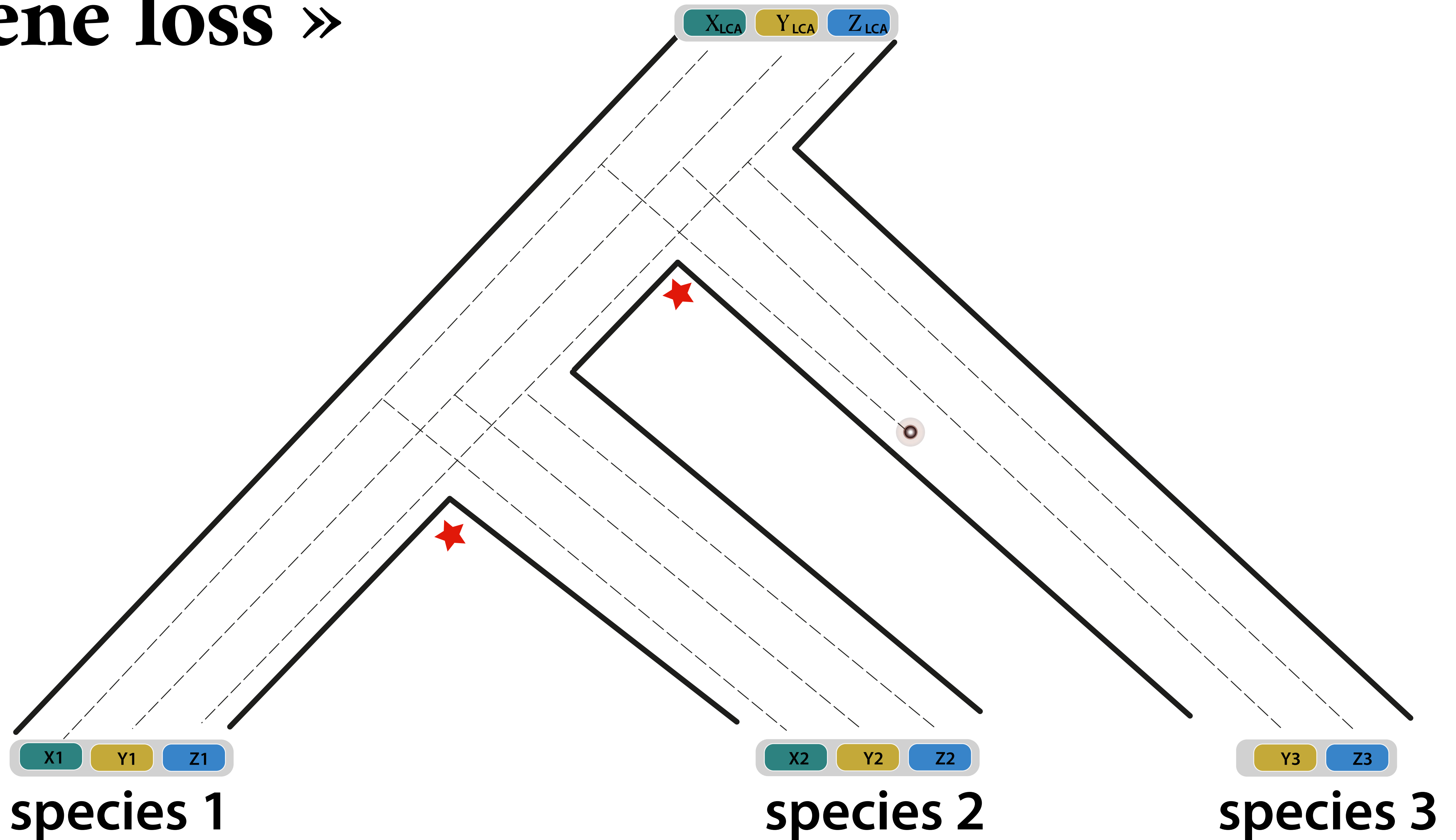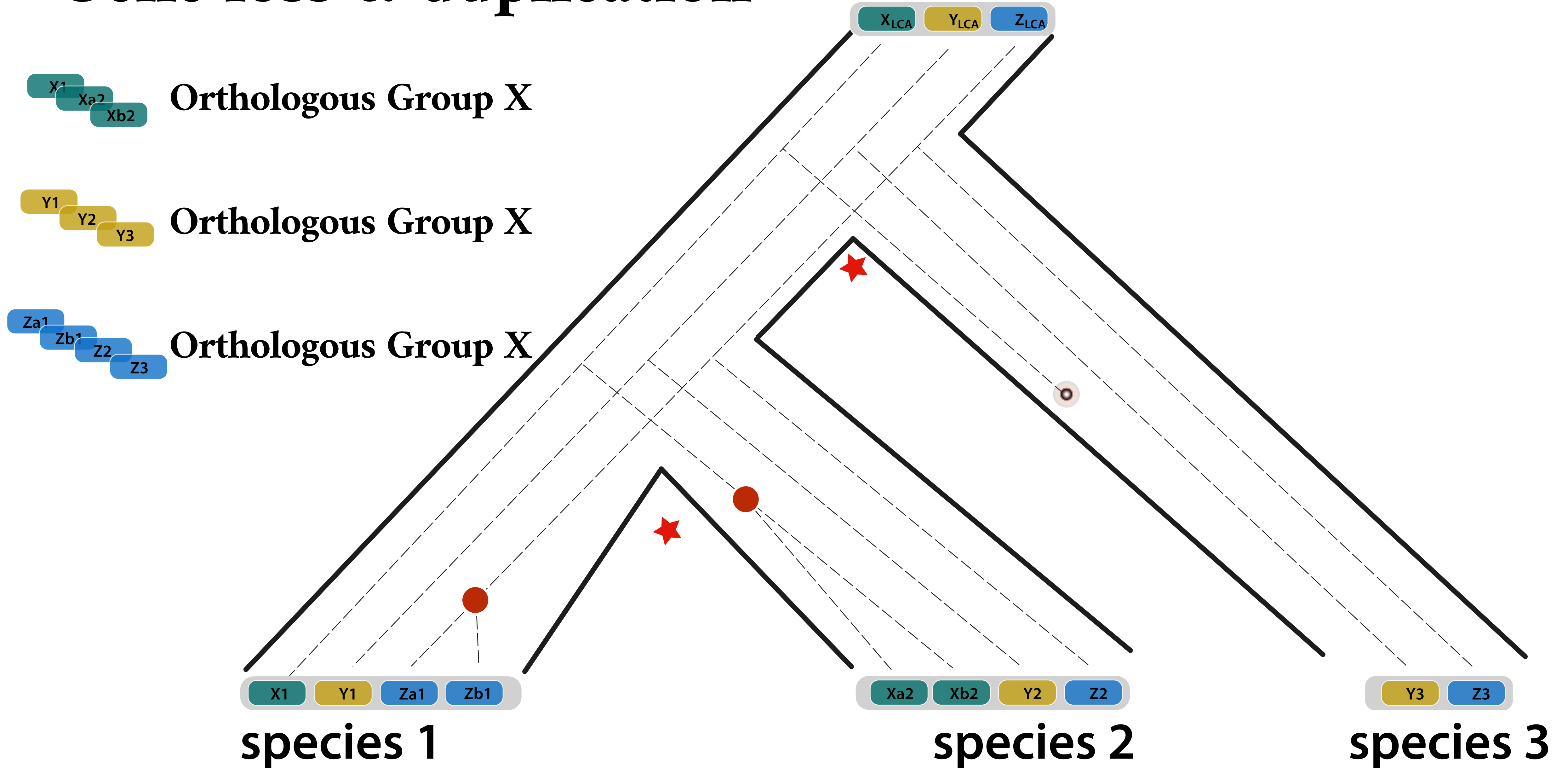species 3

» **Gene loss** »

# » Gene loss & duplication »

Orthologs, Paralogs, and Evolutionary Genomics | https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725
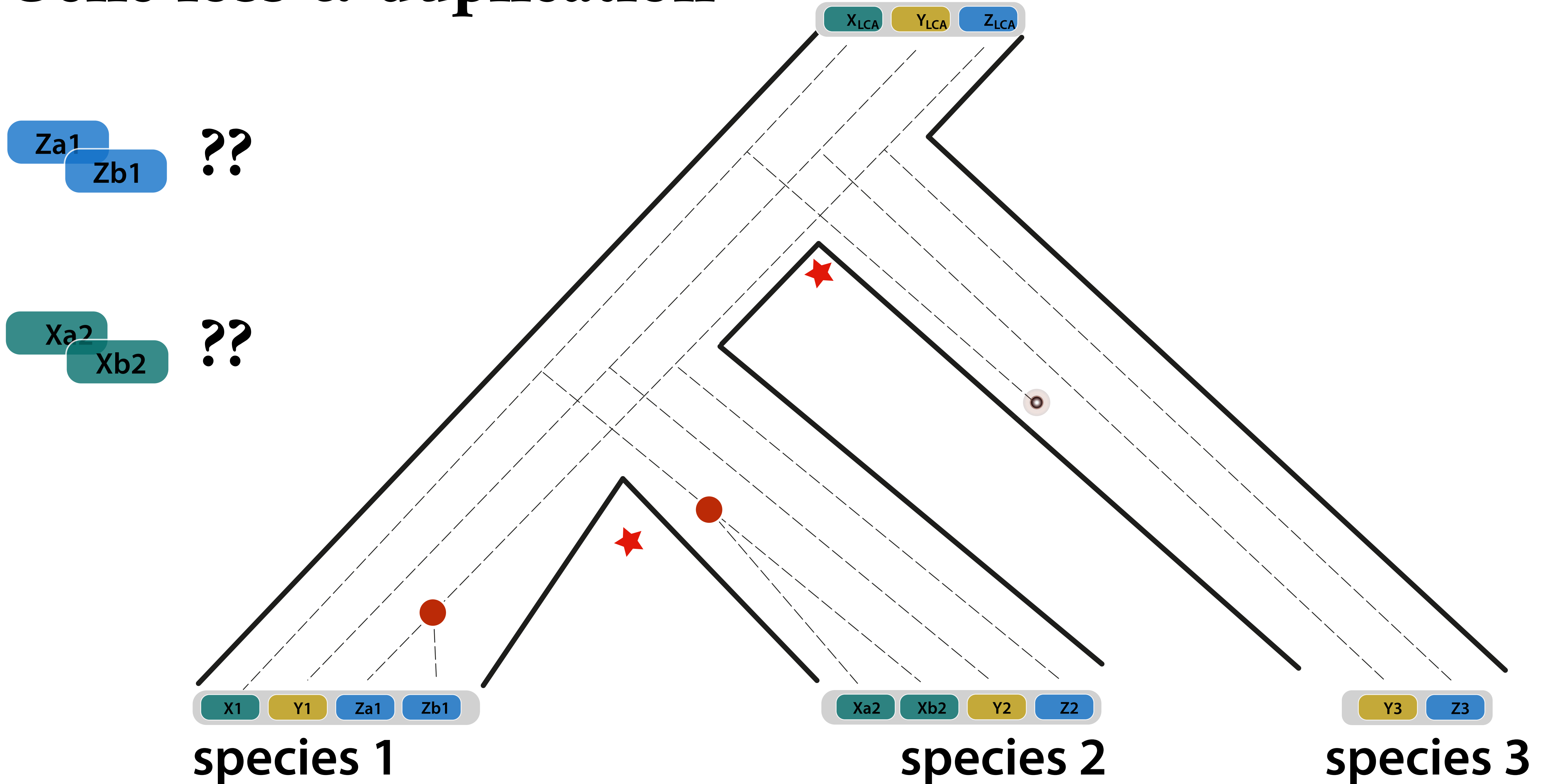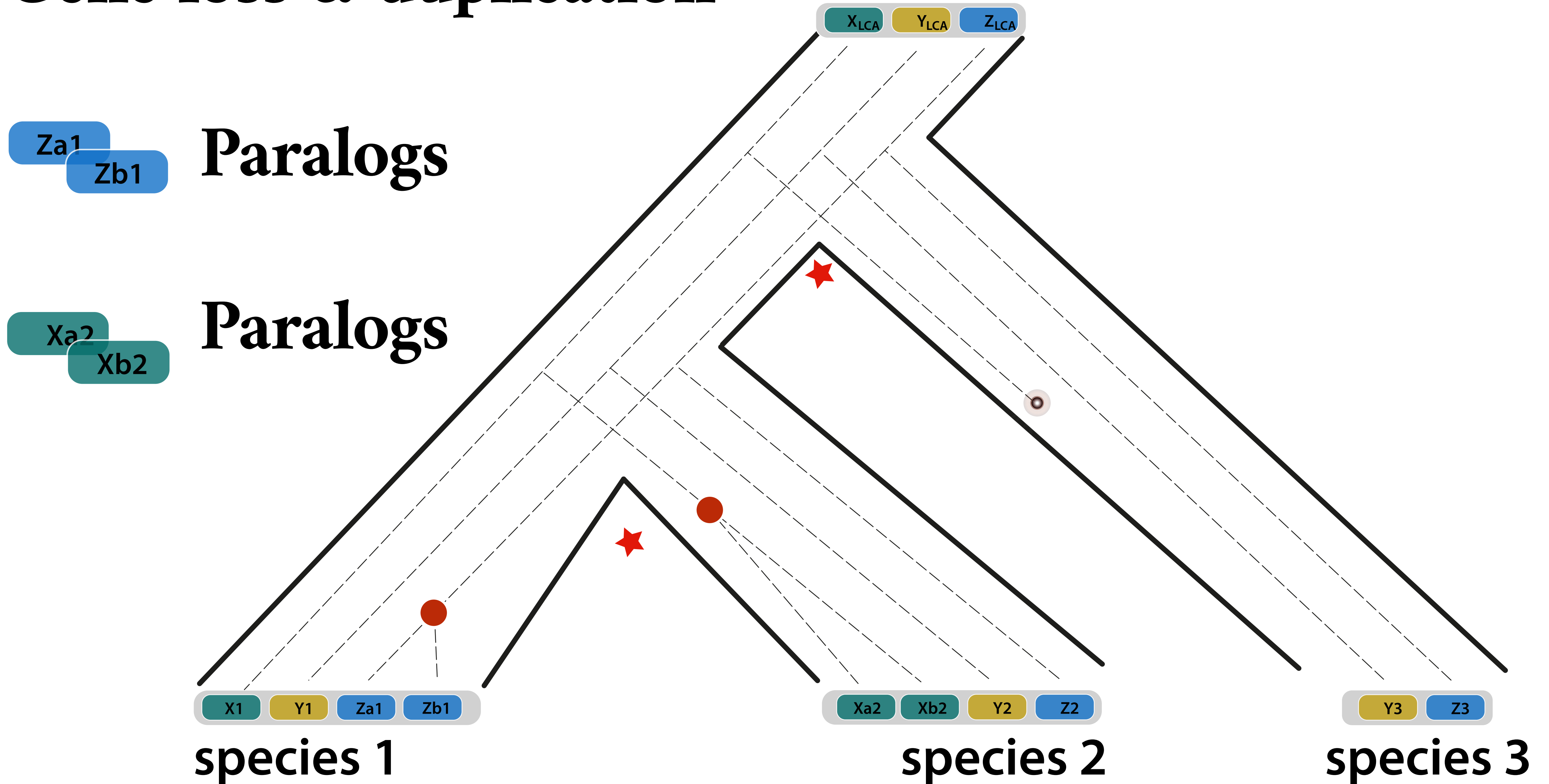Functional and evolutionary implications of gene orthology | https://www.nature.com/articles/nrg3456

» Gene loss & duplication »

» **Gene loss & duplication** »

Paralogs

Paralogs

species 1

species 2

species 3

**Paralogy:** it describes a relationship that involve a duplication:
*if a locus (gene) is generated by an event of tandem duplication.*



**Paralogs**

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
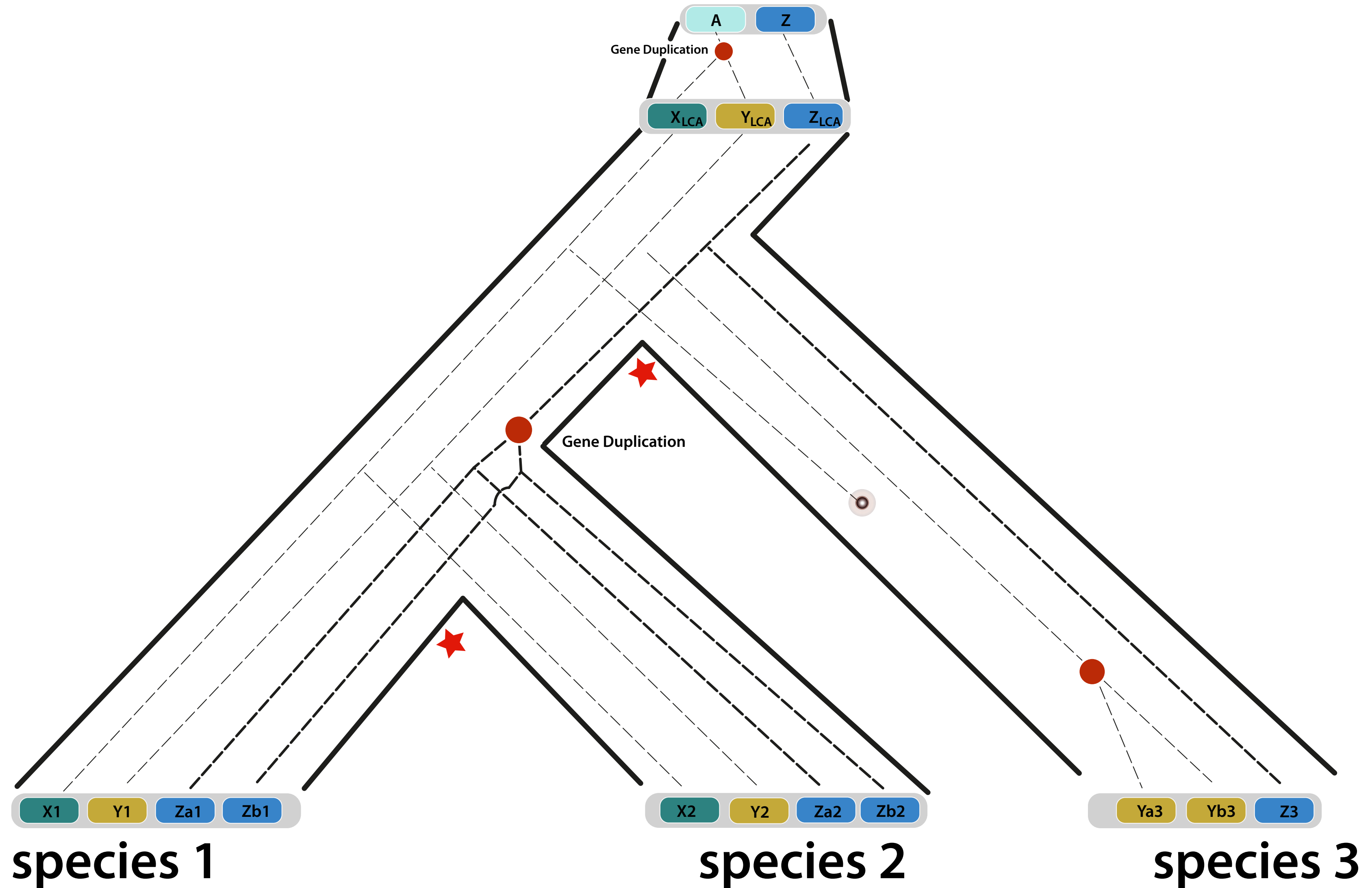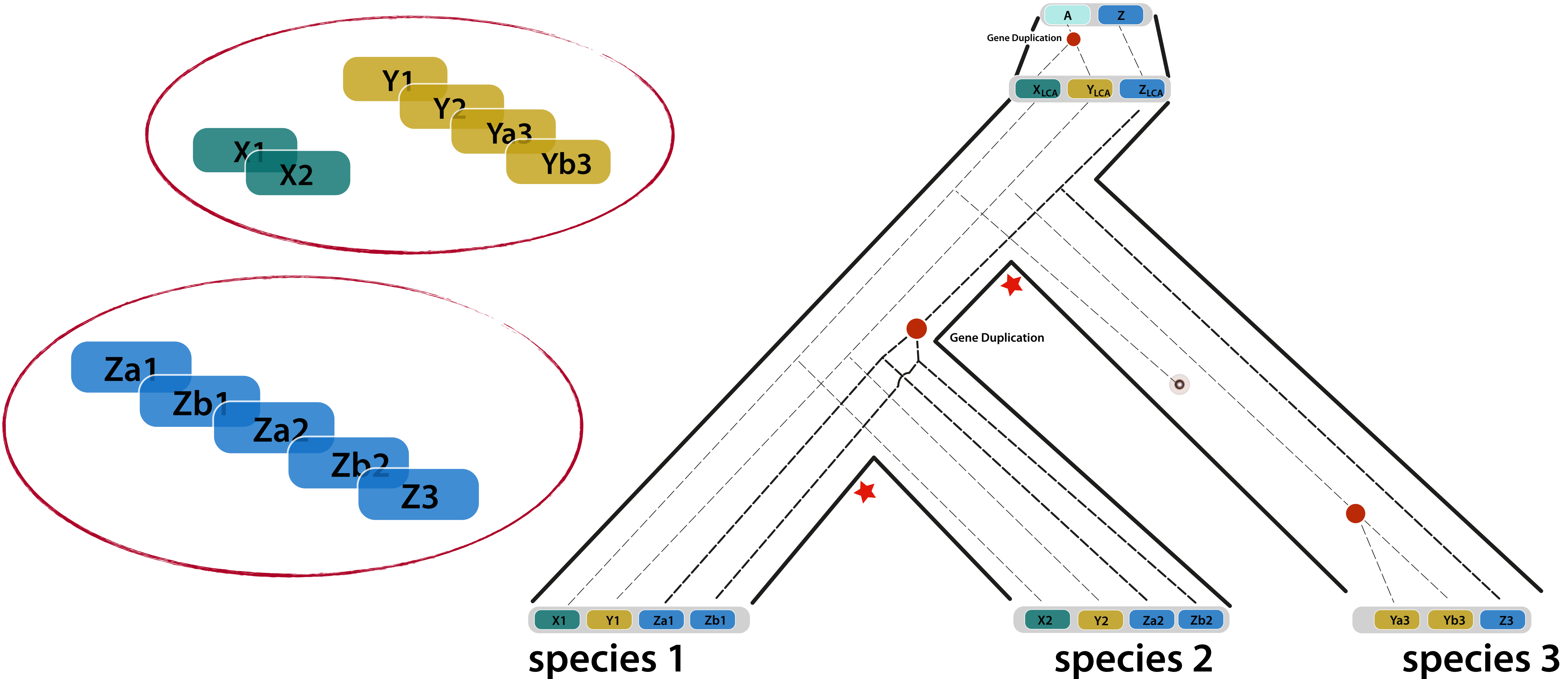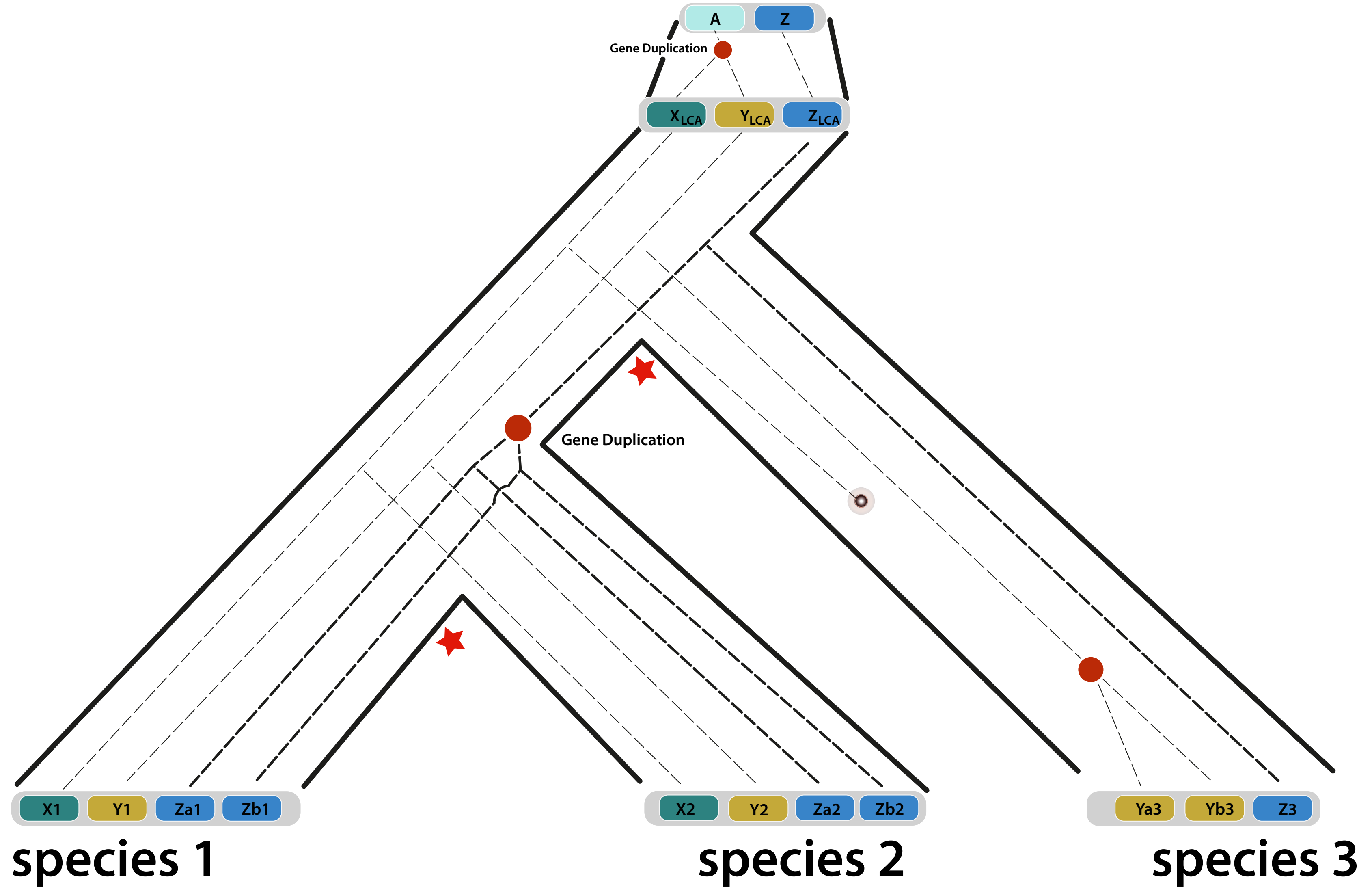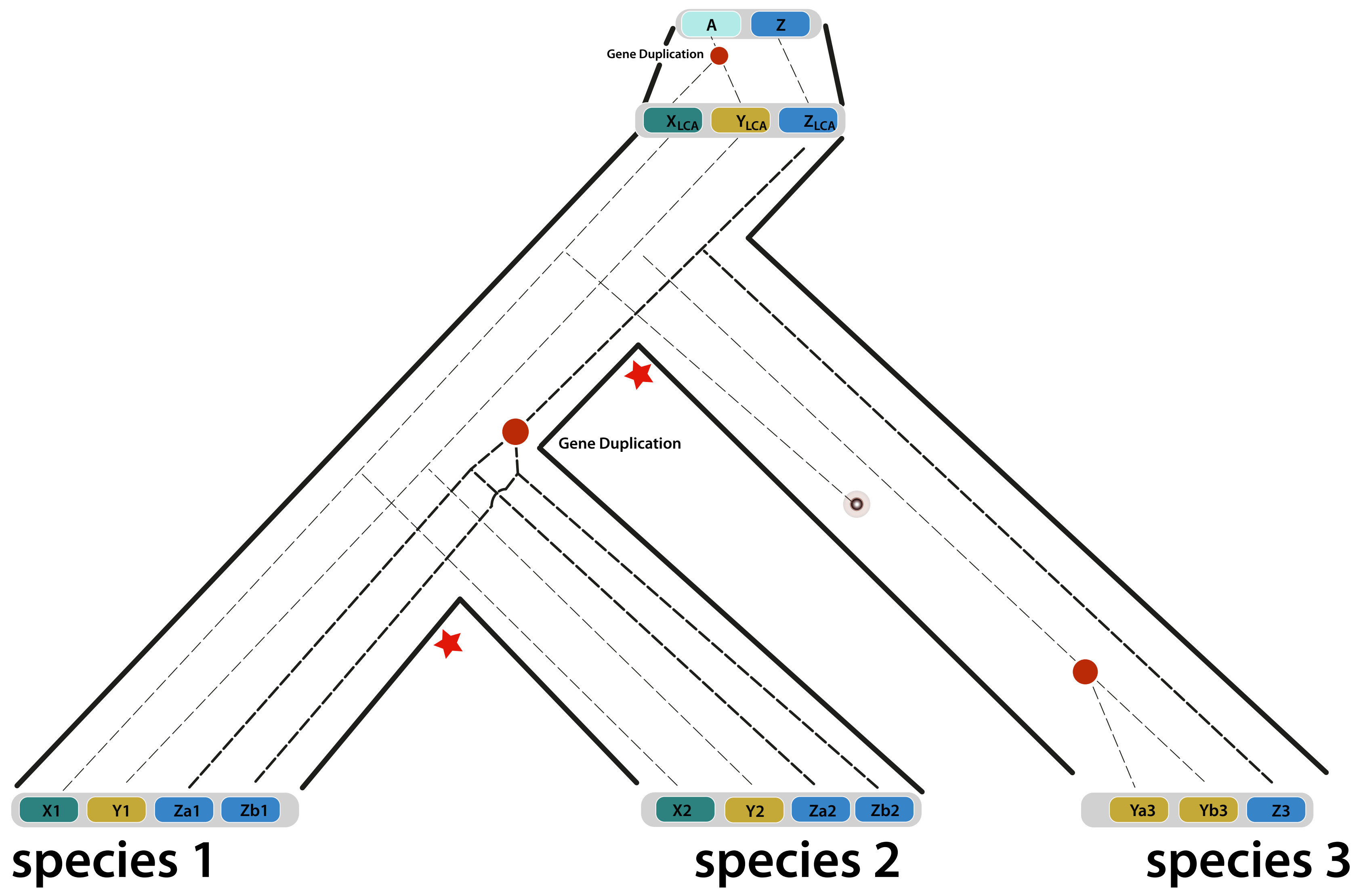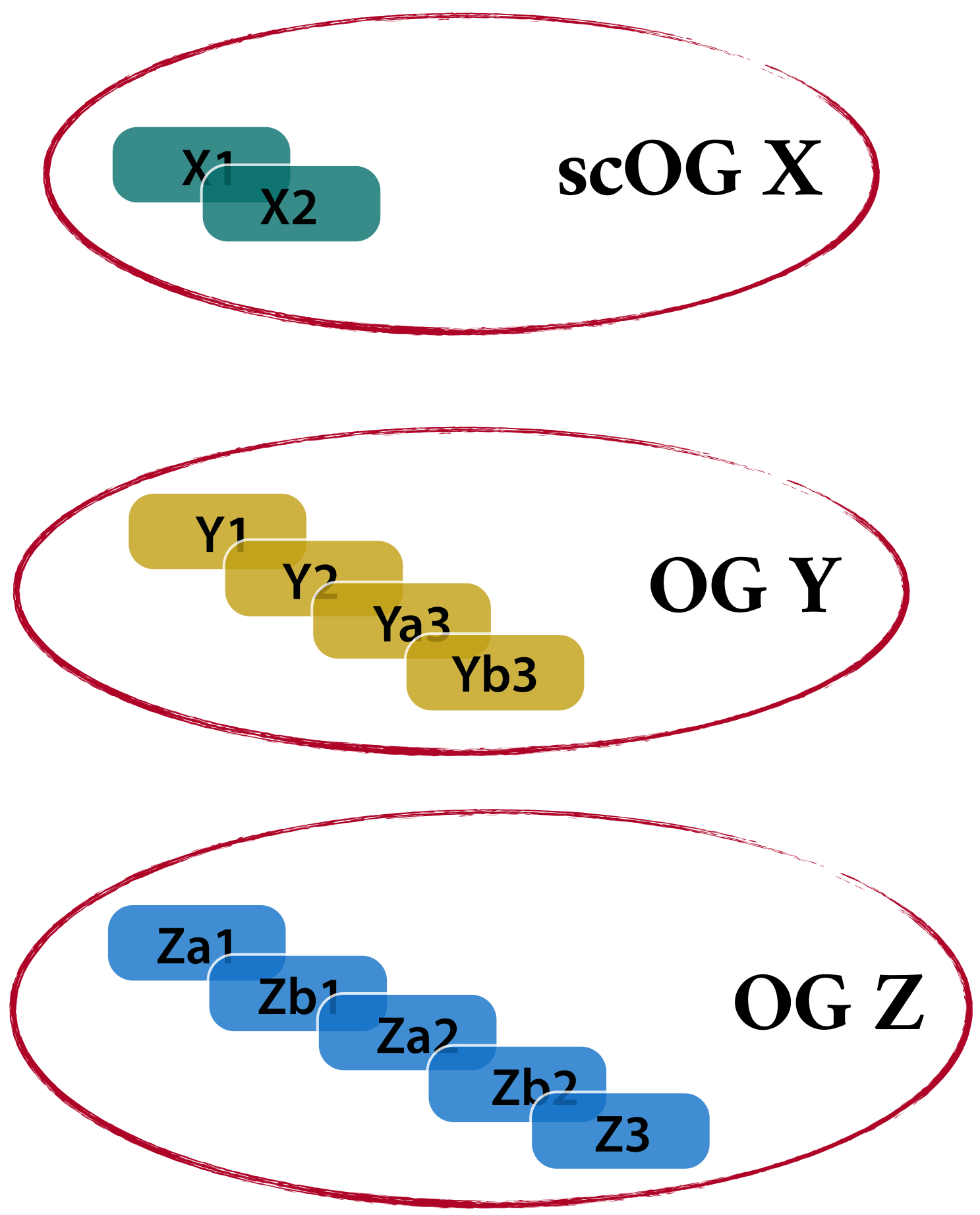Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » Homologs ? »

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » Homologs ? »

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » Orthologs ? »

# » Orthologs ? »

# » Paralogs ? »

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » Paralogs ? »

# » Paralogs ? »
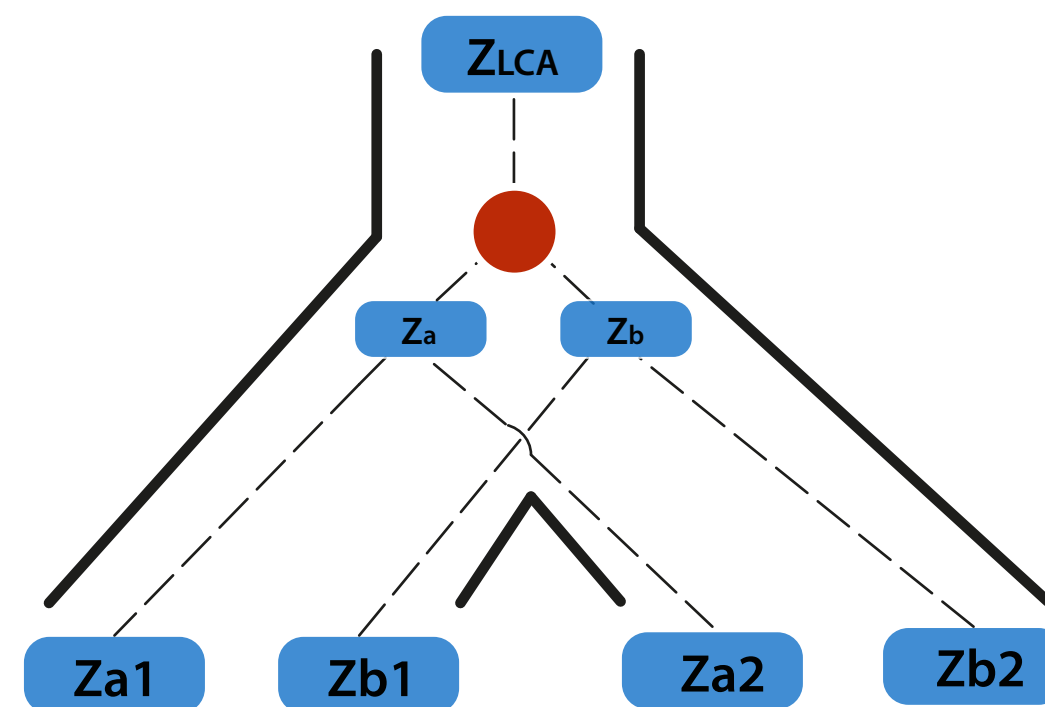
**IN-Paralogy:** it describes a relationship that involve a duplication occurred <u>within a species</u>



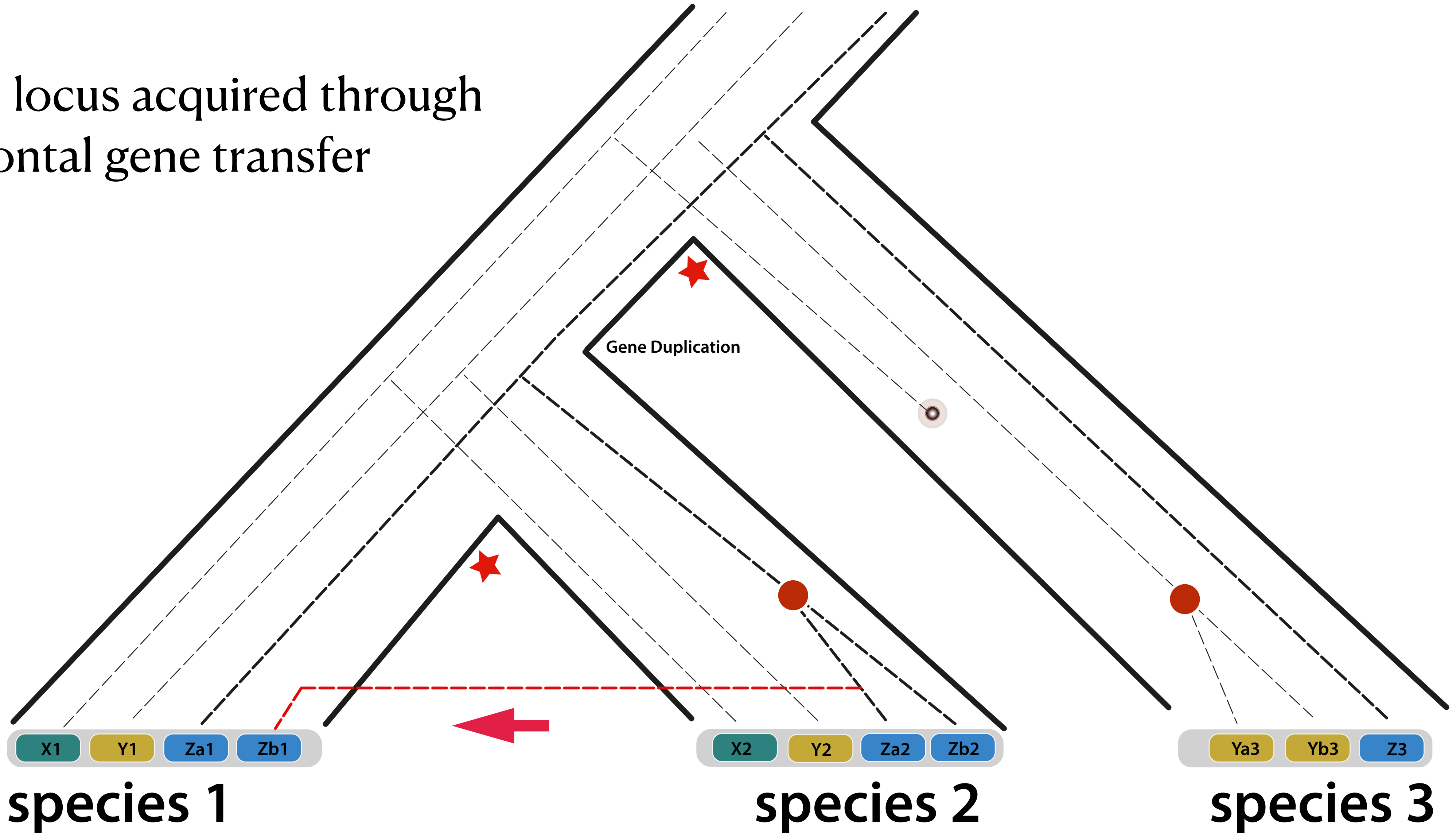**IN-Paralogs**

**OUT-Paralogy:** it describes a relationship that involve a duplication occurred <u>in one of the ancestor</u>

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*
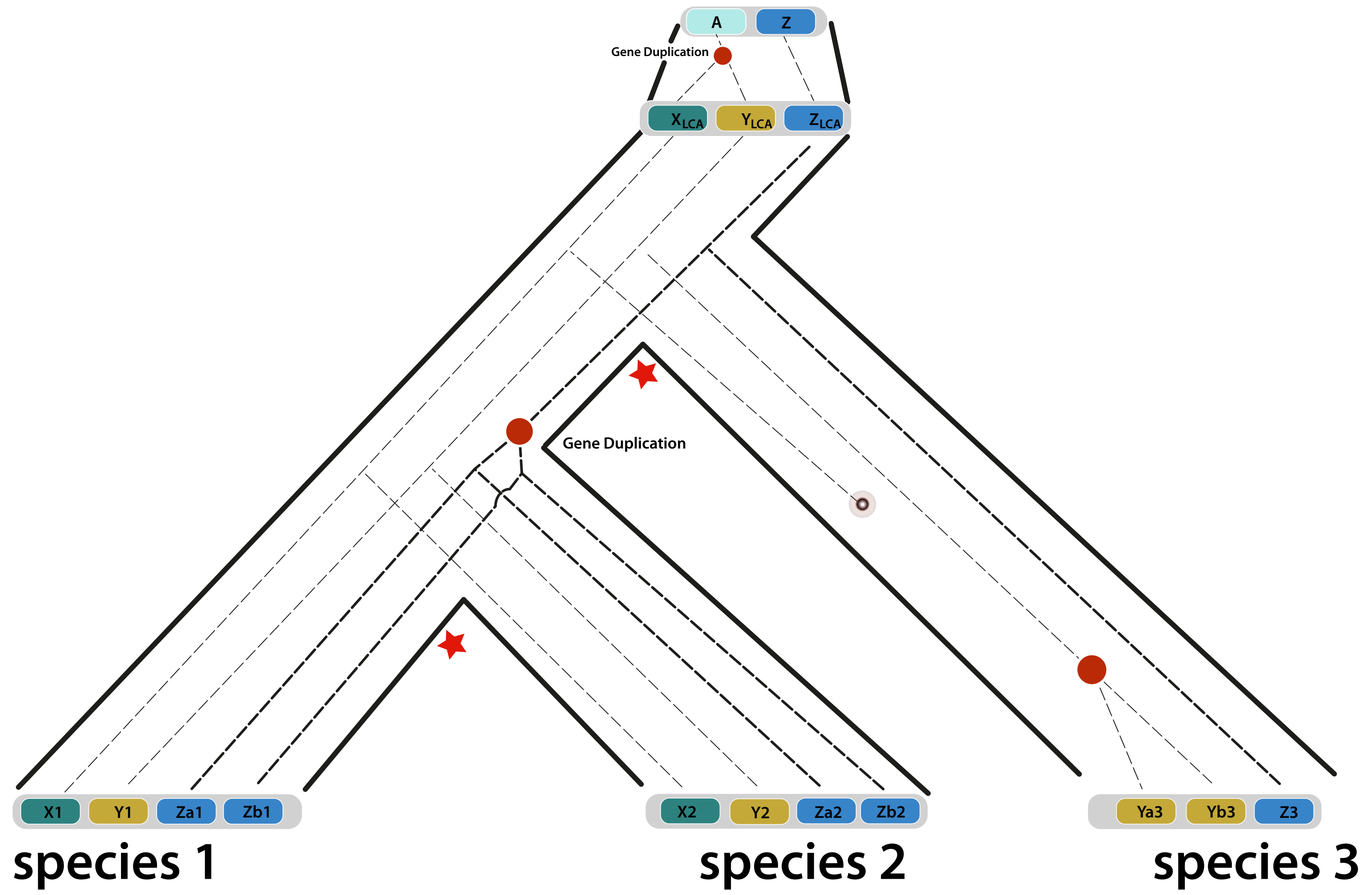
# » Horizontal Gene Transfer (HGT) »

**Xenologs:** locus acquired through horizontal gene transfer



Gene Duplication

species 1    species 2    species 3

» **Gene Family Tree** »

Gene Duplication

A    Z

$X_{LCA}$    $Y_{LCA}$    $Z_{LCA}$

Gene Duplication

X1    Y1    Za1    Zb1

X2    Y2    Za2    Zb2

Ya3    Yb3    Z3

species 1

species 2

species 3
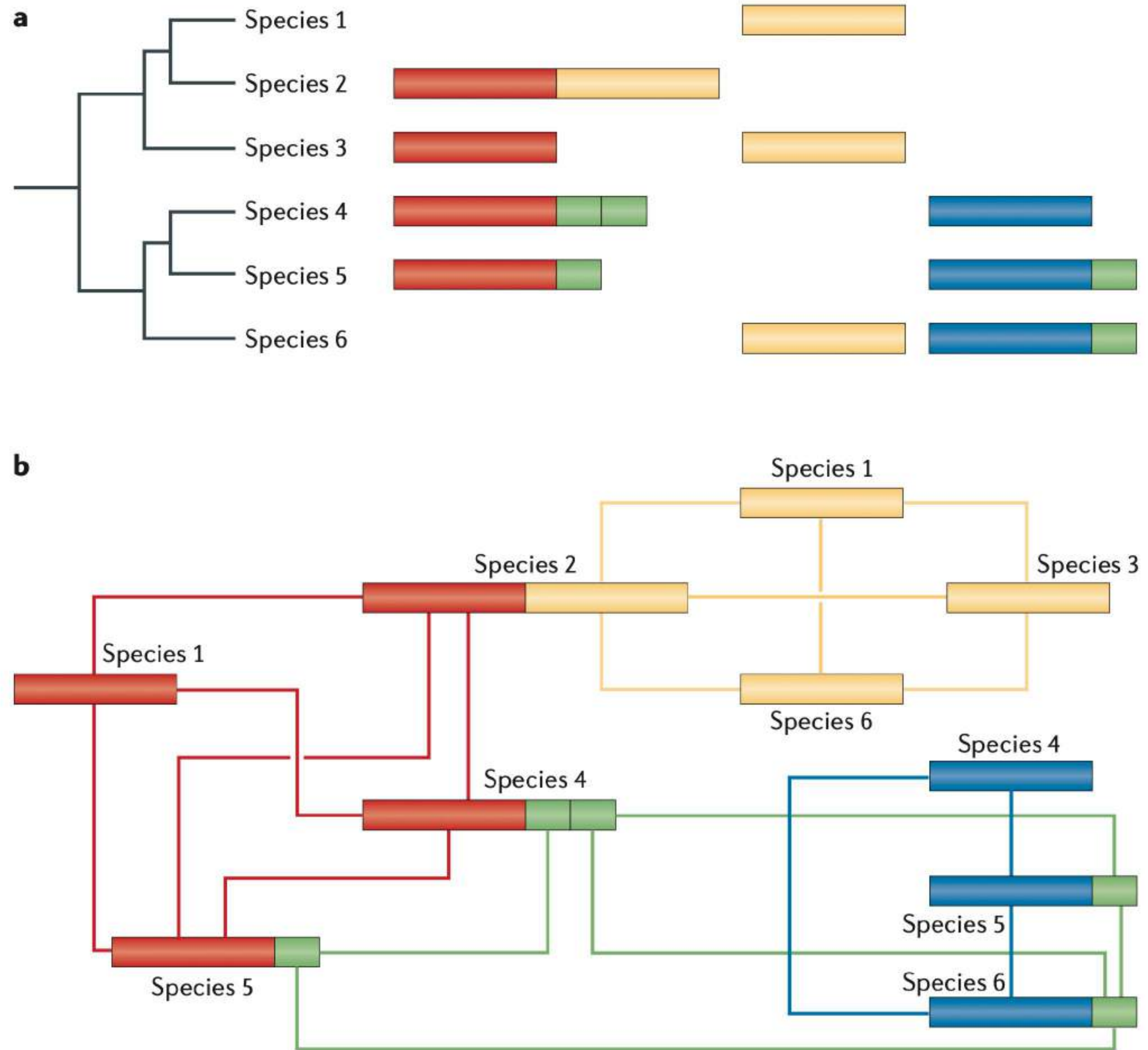
# » Gene Family Trees »



Gene duplication

Speciation

Out-paralogs

In-paralogs

A    Z

Gene Duplication

X_LCA    Y_LCA    Z_LCA

Gene Duplication

X1  Y1  Za1  Zb1          X2  Y2  Za2  Zb2          Ya3  Yb3  Z3

species 1                    species 2                    species 3

X1    X2    Y1    Y2    Ya3  Yb3    Za1    Za2    Zb1    Zb2    Z3

scOG X                           OG Y                              OG Z

# » Units of orthology »

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » **Real examples** »

*Transferrin family*

Orthologs, Paralogs, and Evolutionary Genomics | *https://www.annualreviews.org/content/journals/10.1146/annurev.genet.39.073003.114725*
Functional and evolutionary implications of gene orthology | *https://www.nature.com/articles/nrg3456*

# » Real examples »

*Gustatory receptors*

**8,320** *tips*

**16,638** *nodes*

# » Real examples »



**Heliconinae**
(subfamily)
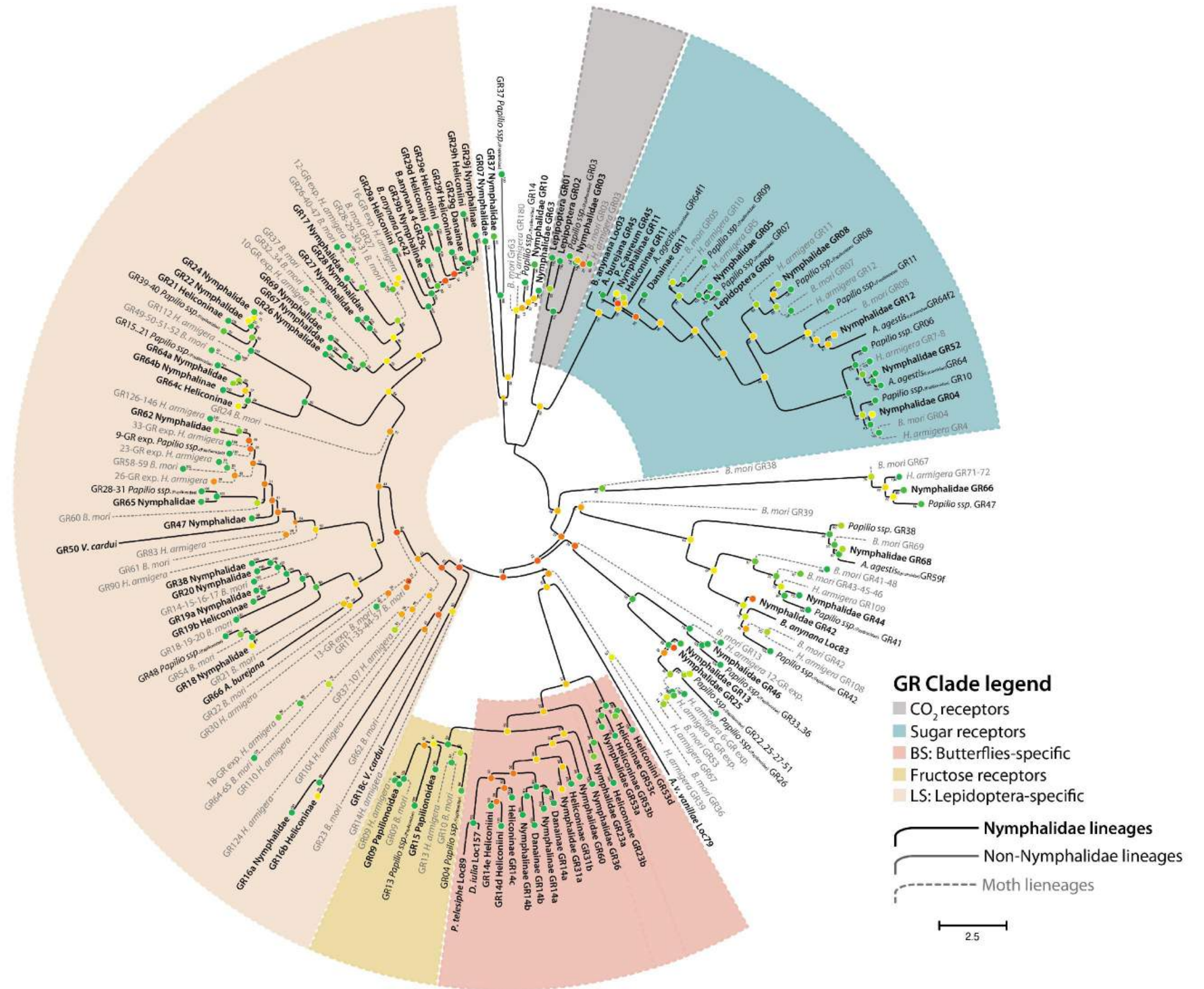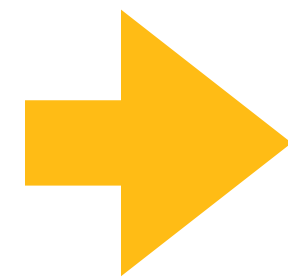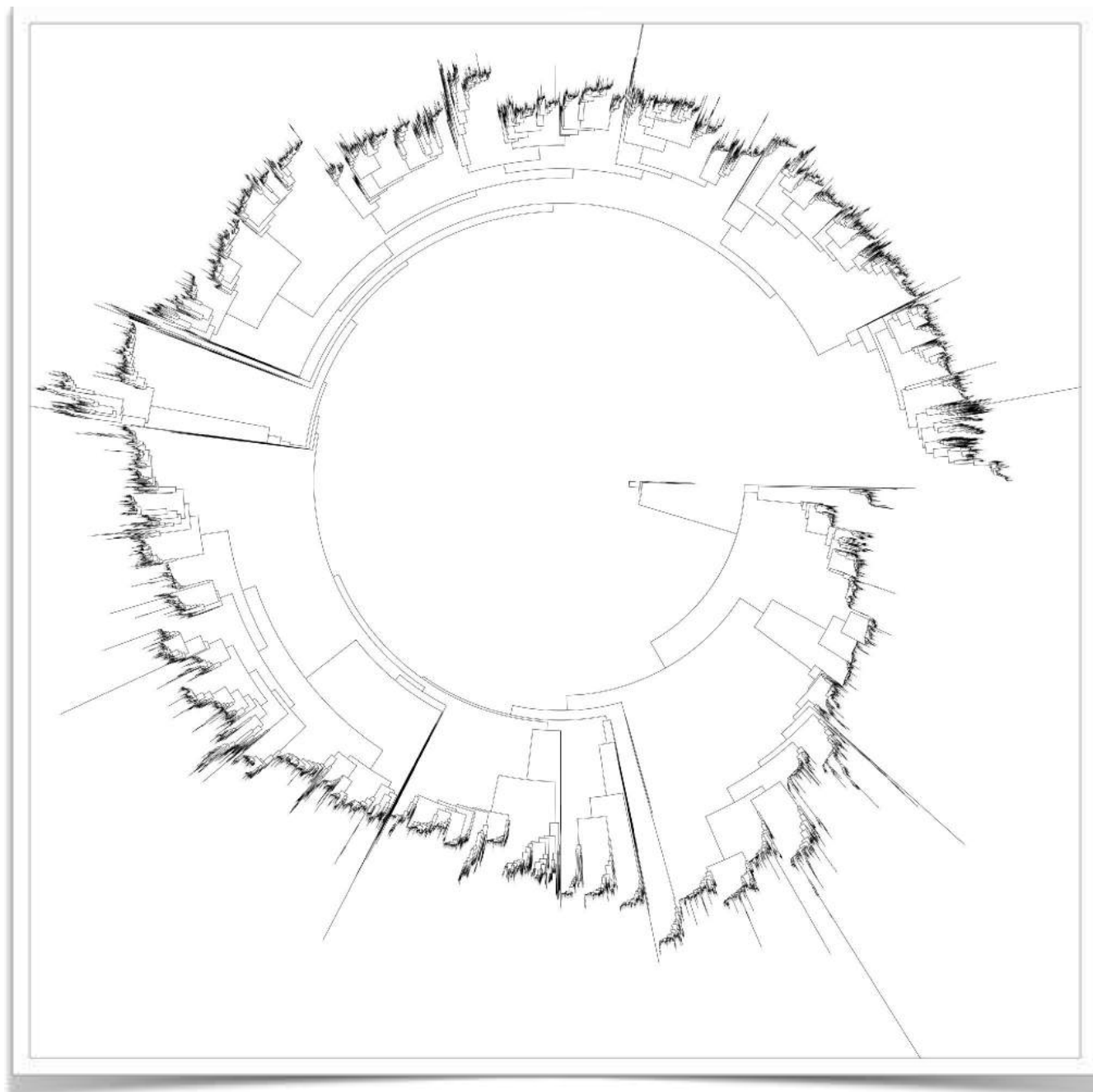
**Heliconiini**
(tribe)

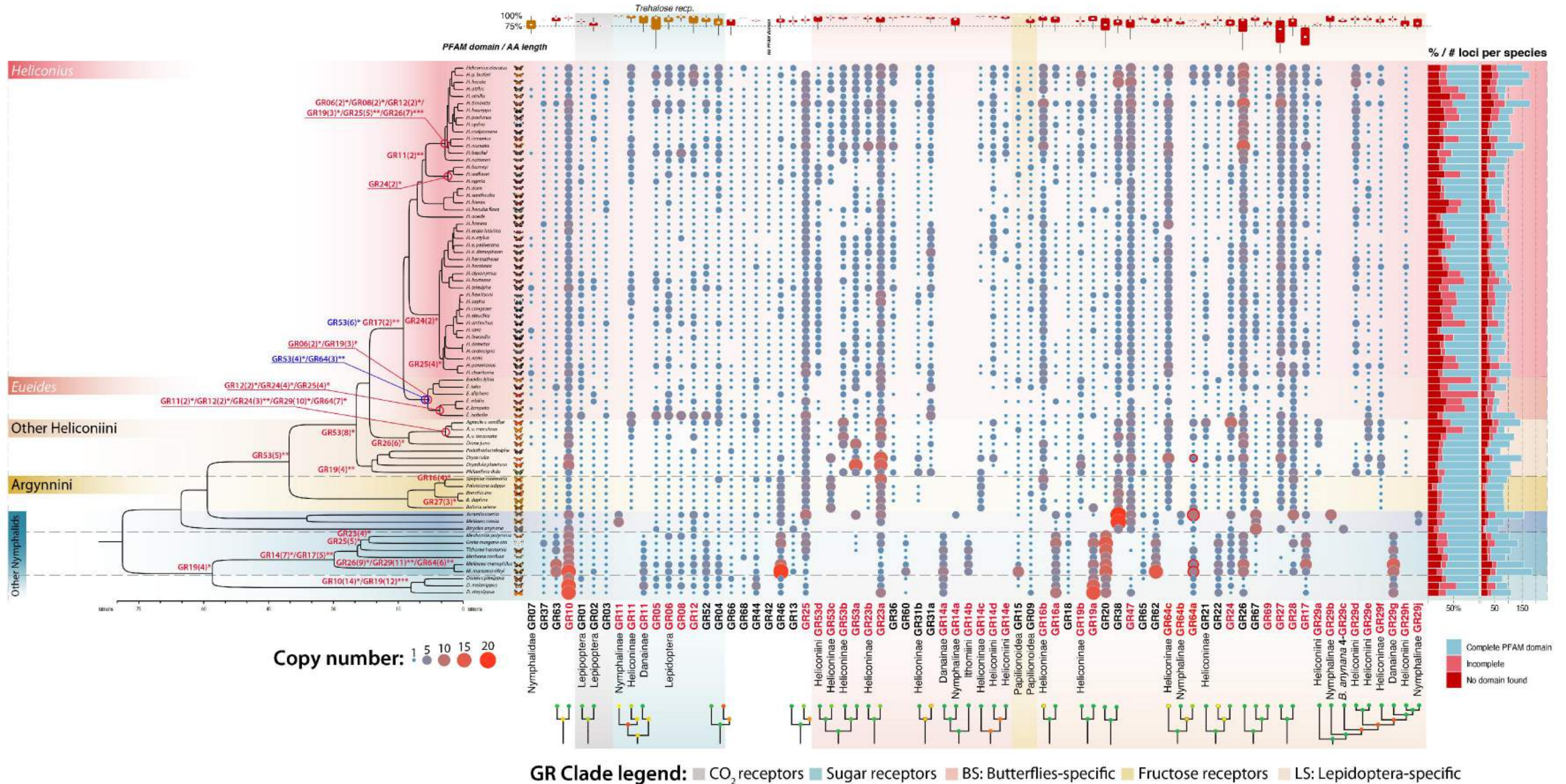# » Why this is relevant ? »

# » Species Tree Estimation »

*Dedicated to Joan*

Phylogenetic tree building in the genomic age | *https://www.nature.com/articles/s41576-020-0233-0*

# » Gene family reconstruction | expansions/contractions »



GR Clade legend
- CO₂ receptors
- Sugar receptors
- BS: Butterflies-specific
- Fructose receptors
- LS: Lepidoptera-specific

**Nymphalidae lineages**
Non-Nymphalidae lineages
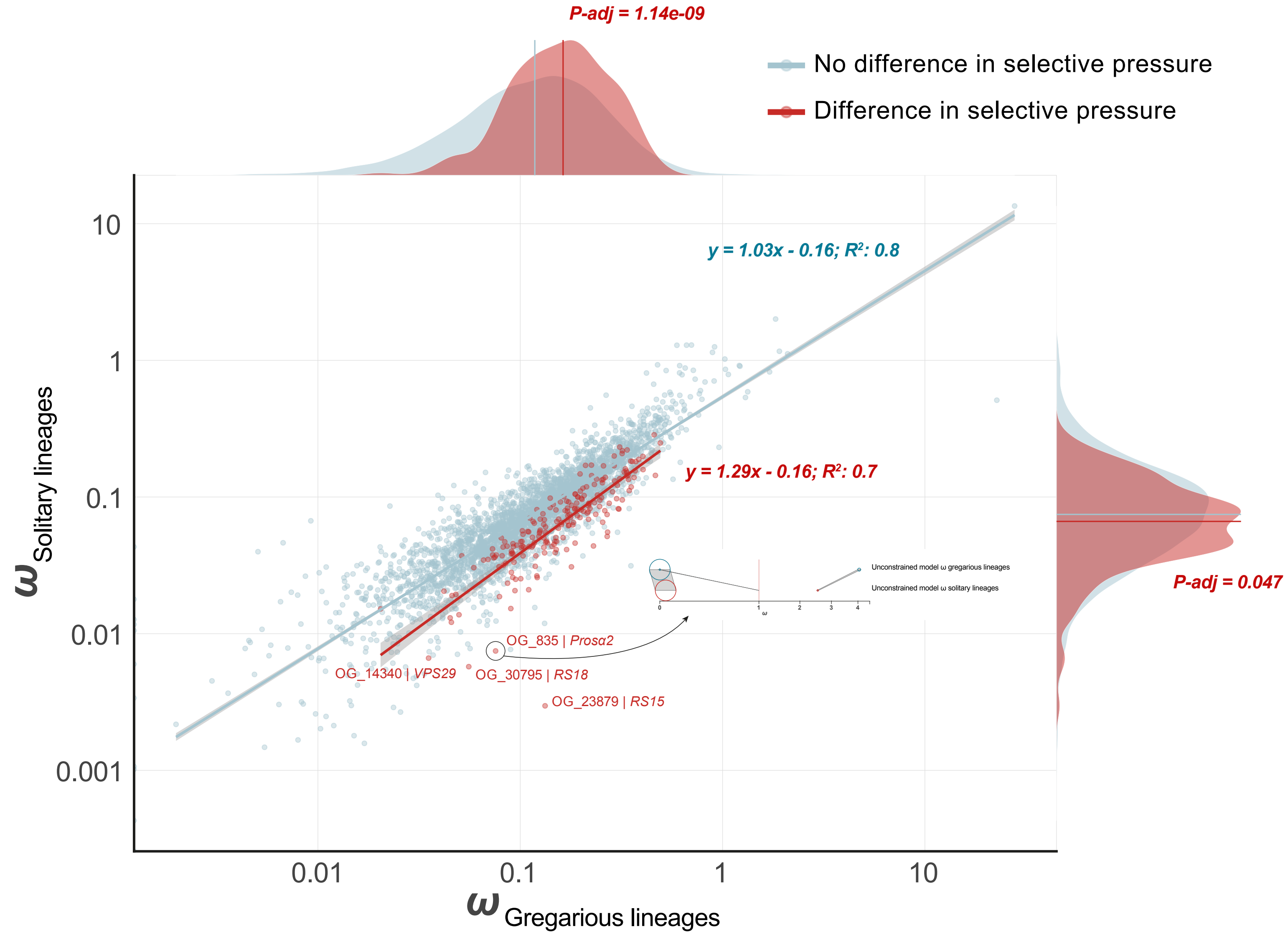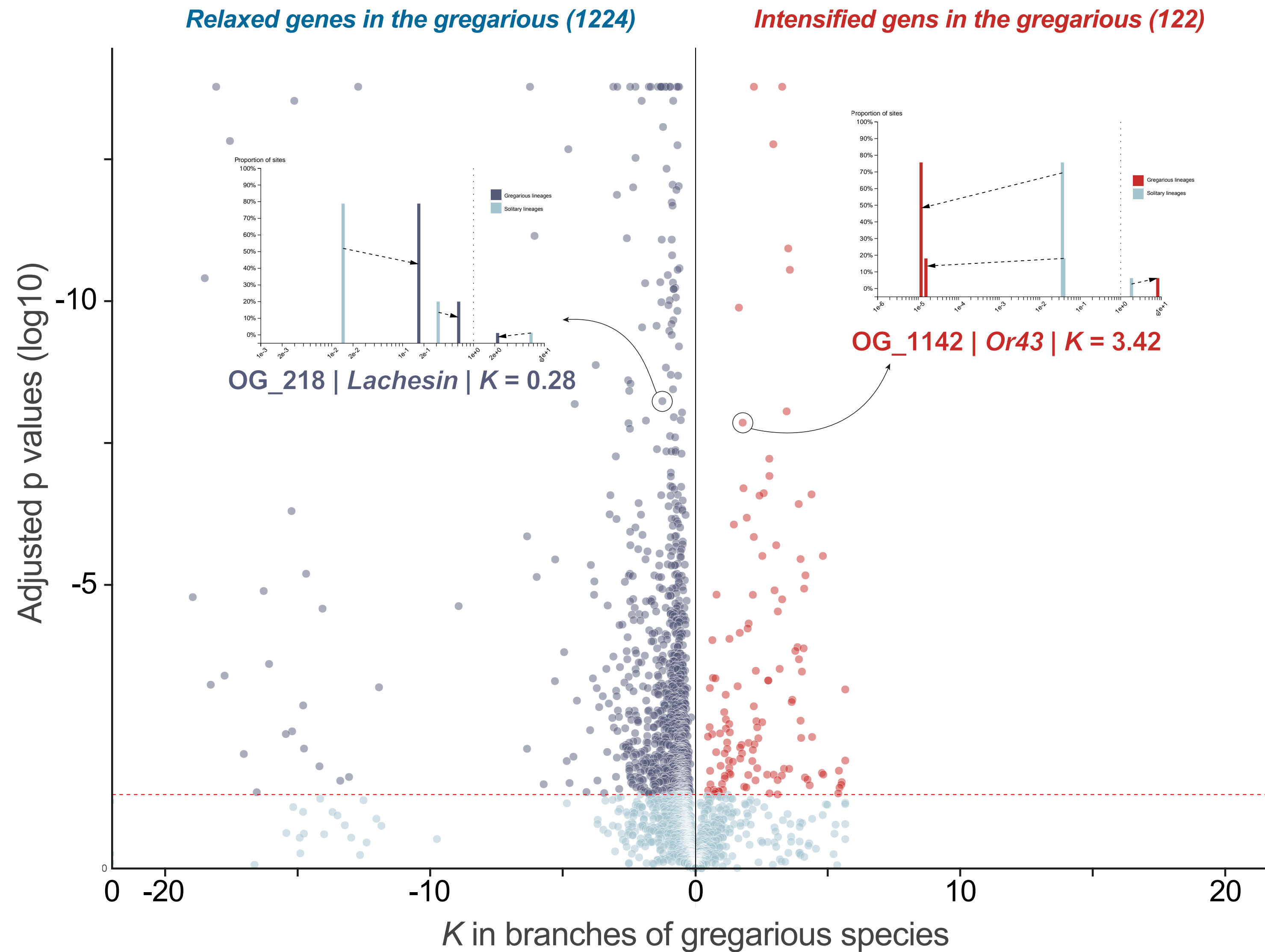Moth lineages

# » Gene family reconstruction | expansions/contractions »

# » Selecting pressures: Positive/Purifying/Convergence »

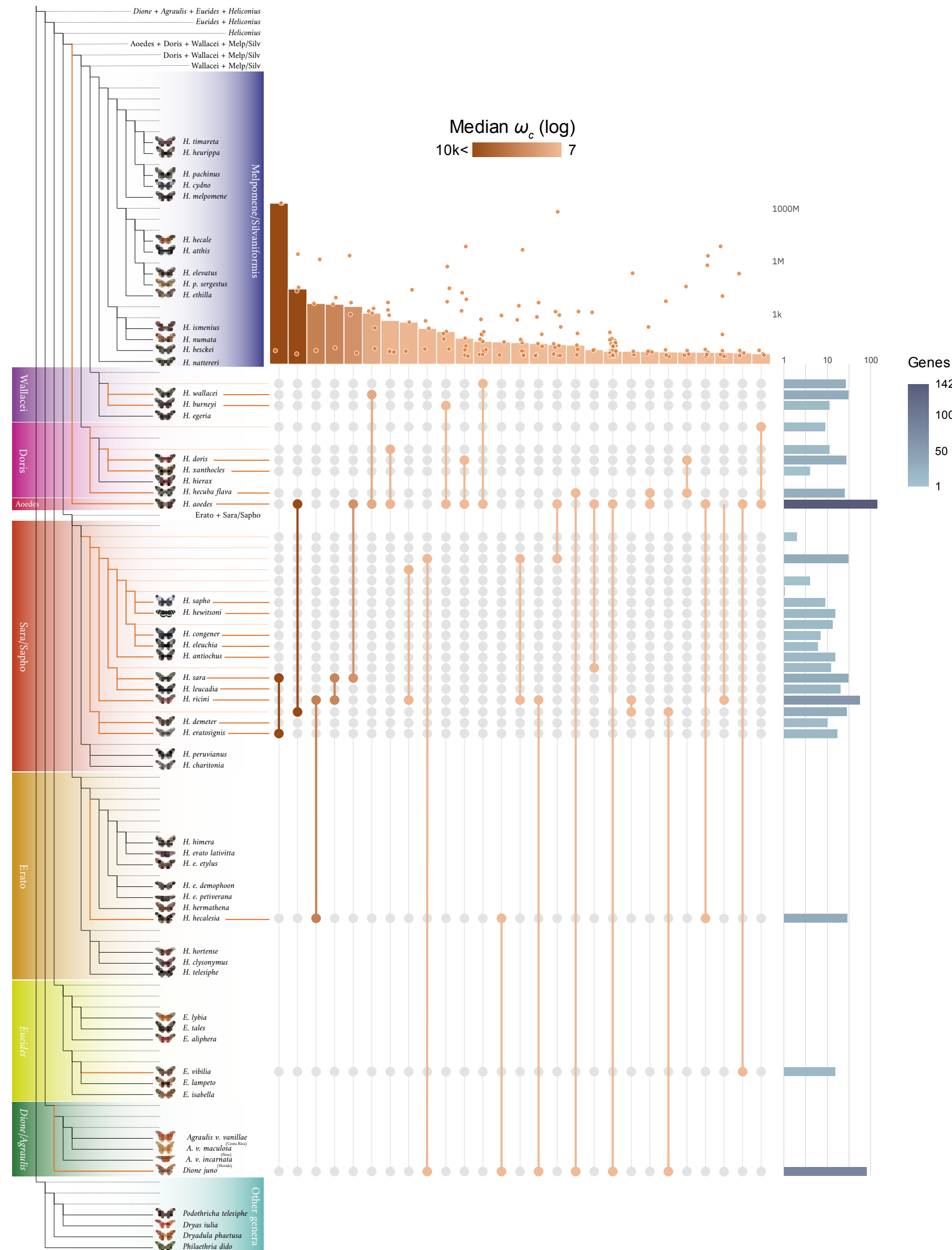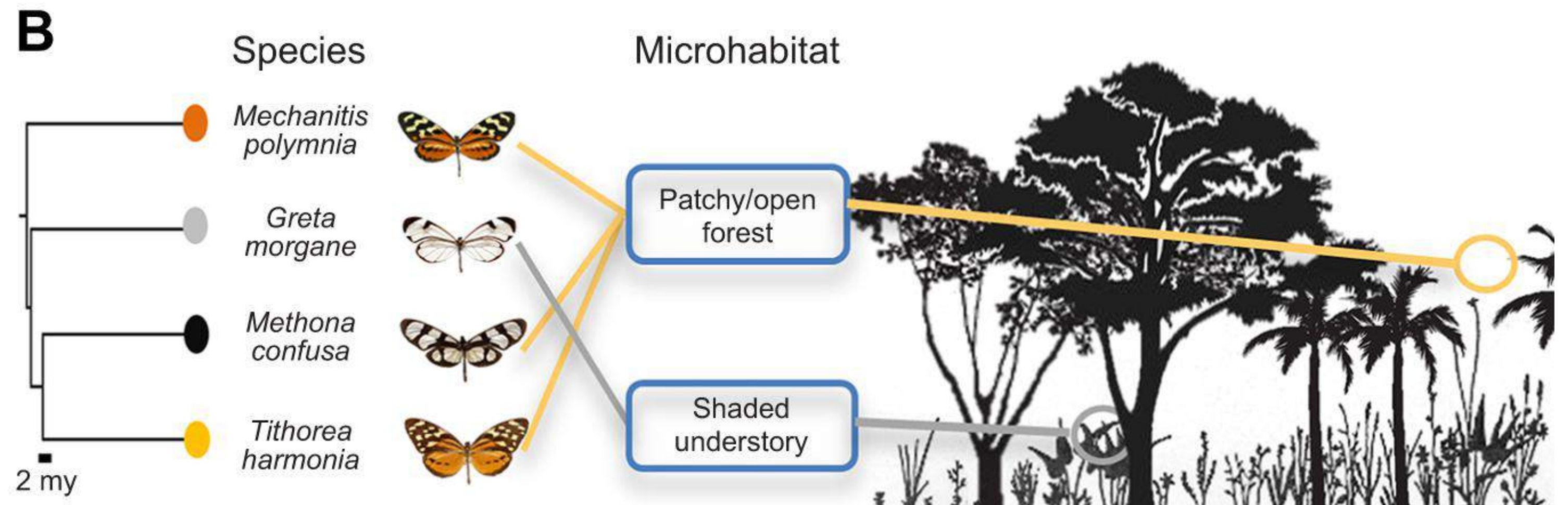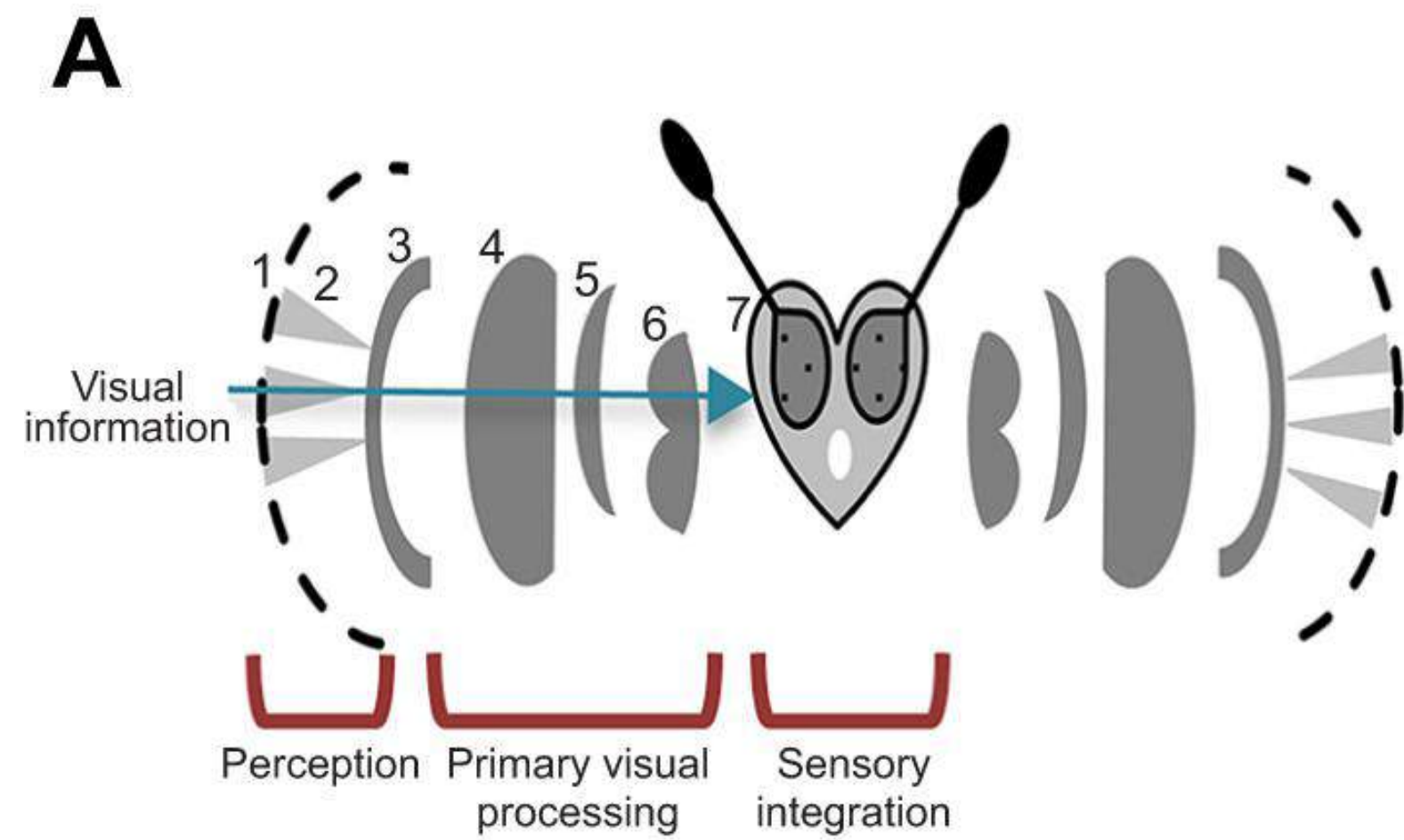# » Selecting pressures: Positive/Purifying/Convergence »



Relaxed genes in the gregarious (1224)

Intensified gens in the gregarious (122)

OG_218 | *Lachesin* | *K* = 0.28

OG_1142 | *Or43* | *K* = 3.42

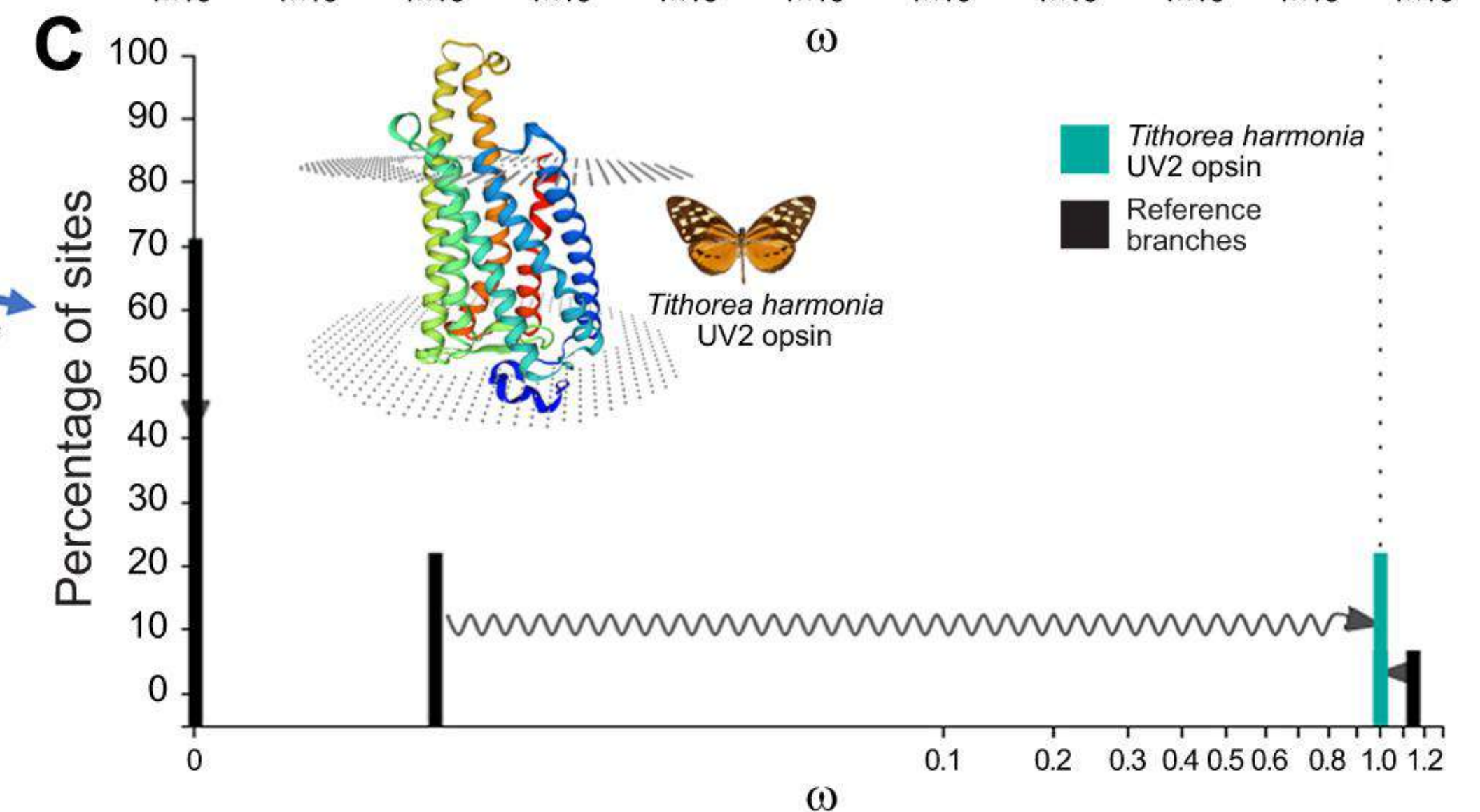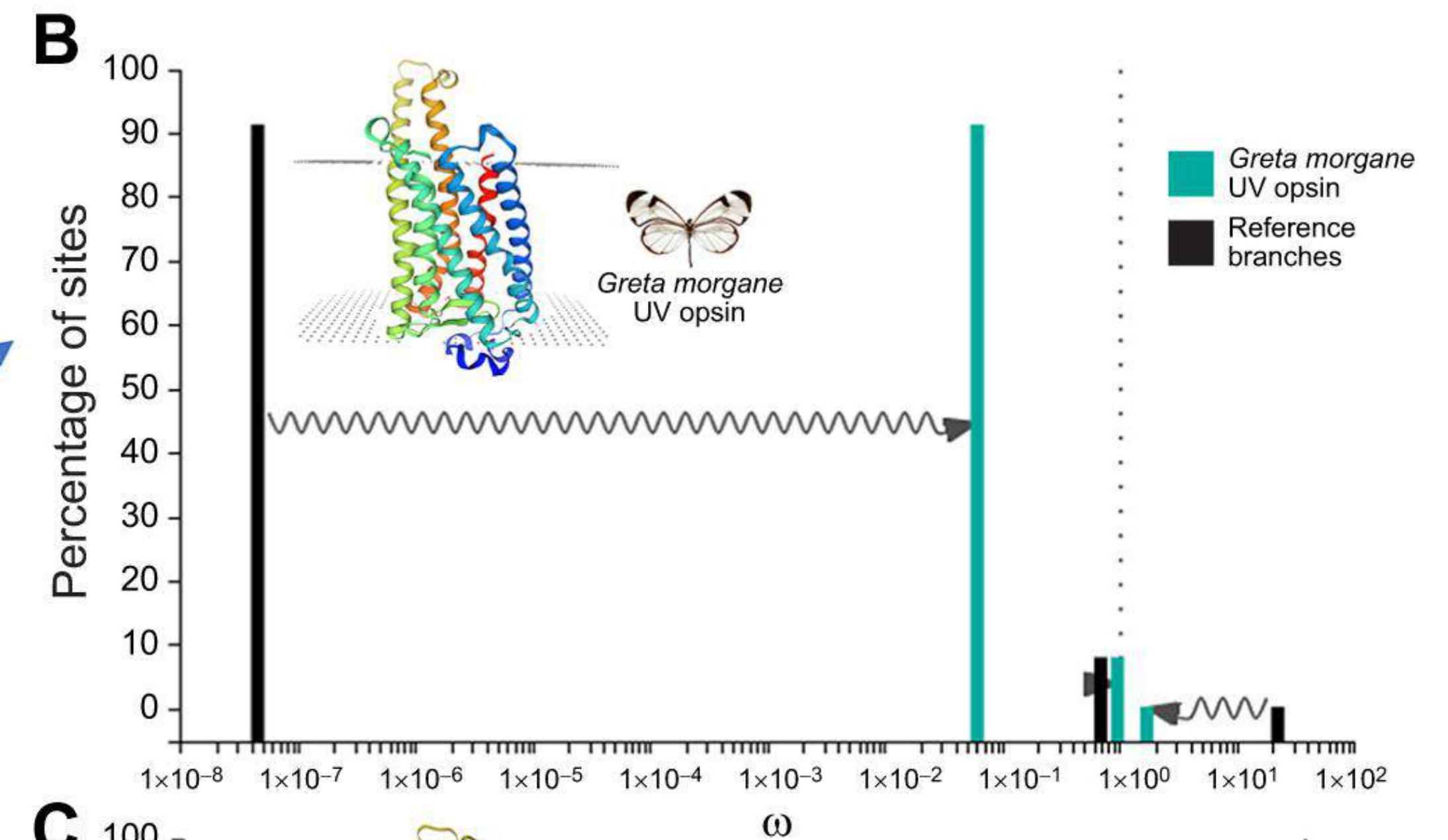Adjusted p values (log10)

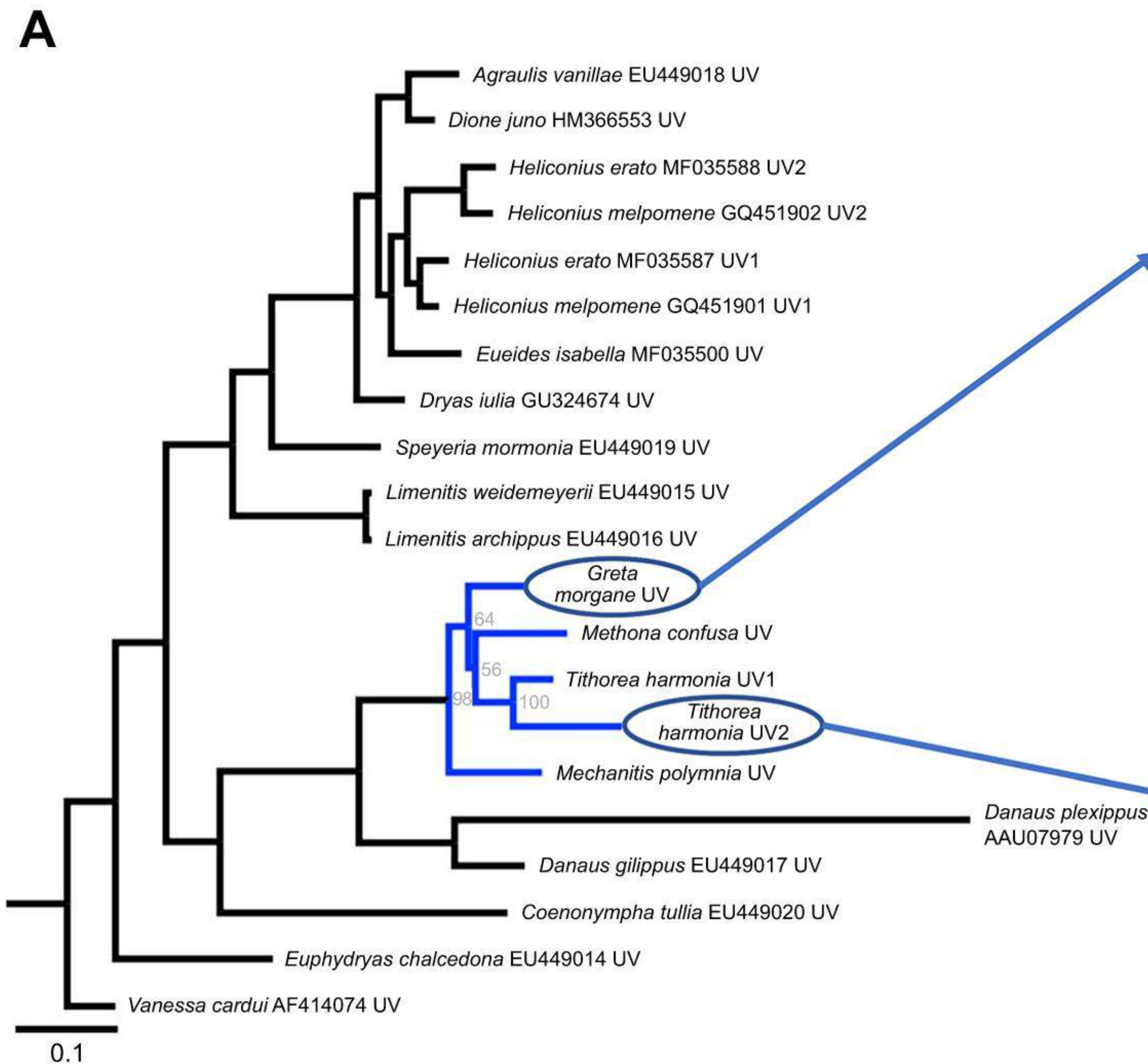*K* in branches of gregarious species

# » Loss / Gain of function »

# » Loss of function »
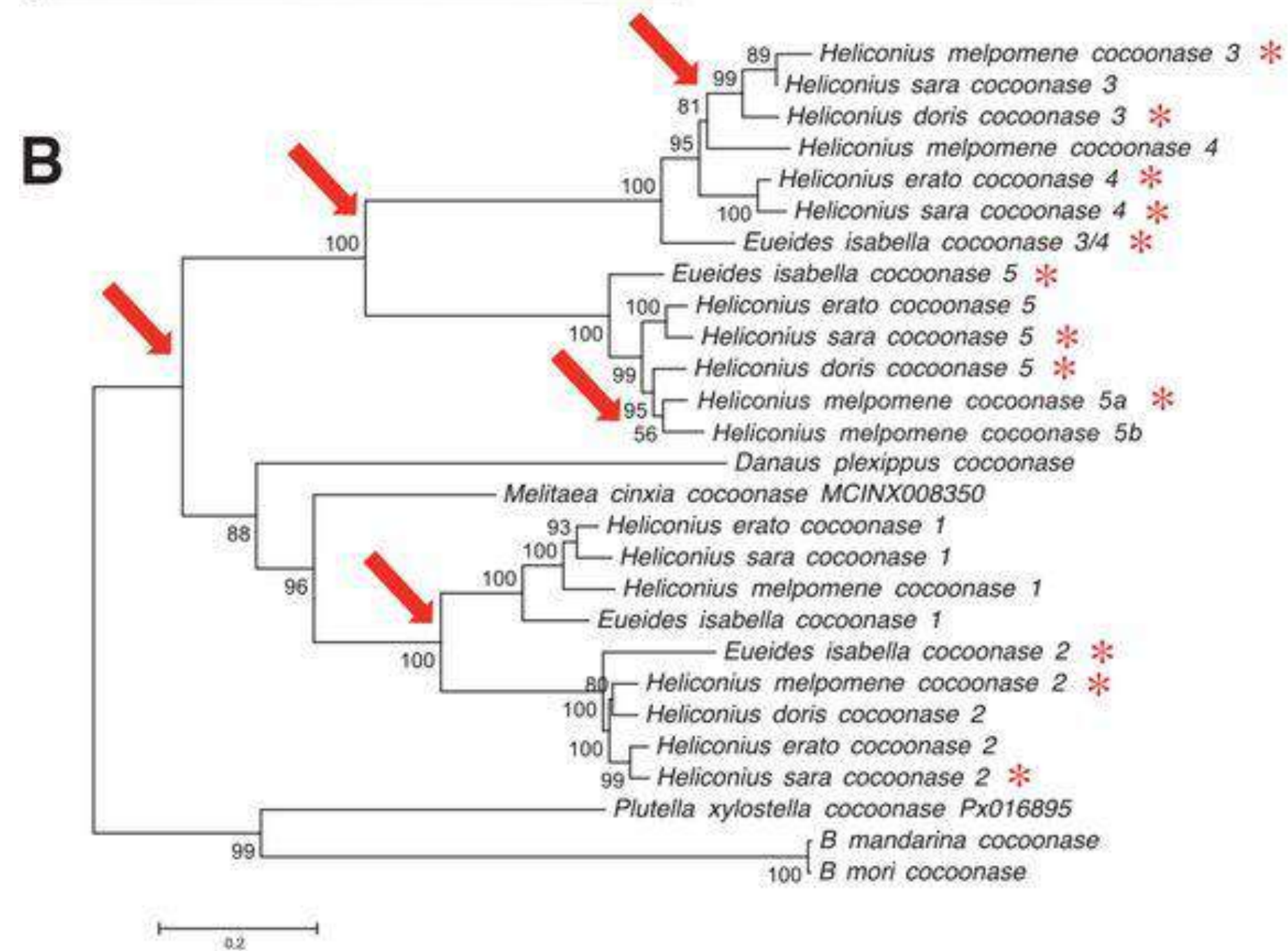
# » Loss of function »

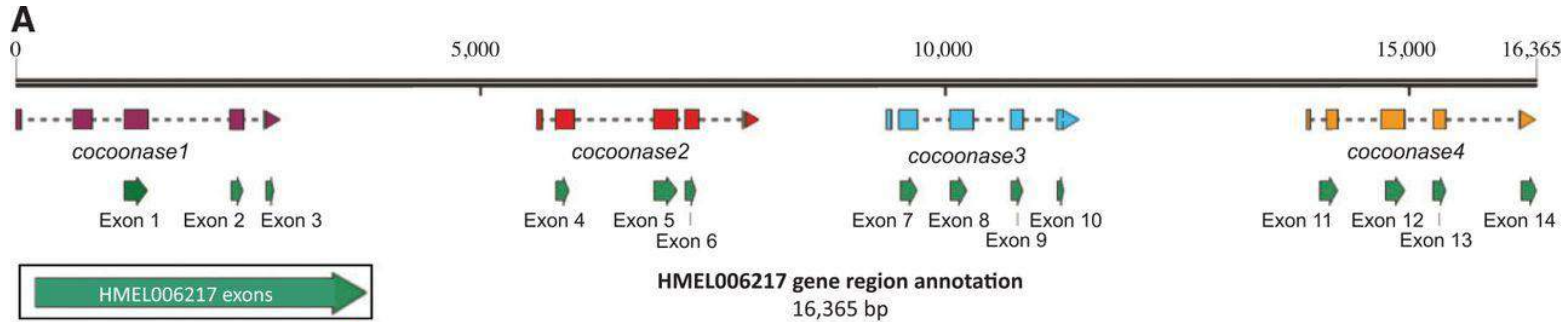# » (*Putative*) Gain of function »

**Cocoonase:** a protease secreted during the emergence of silk moths

# » (*Putative*) Gain of function »



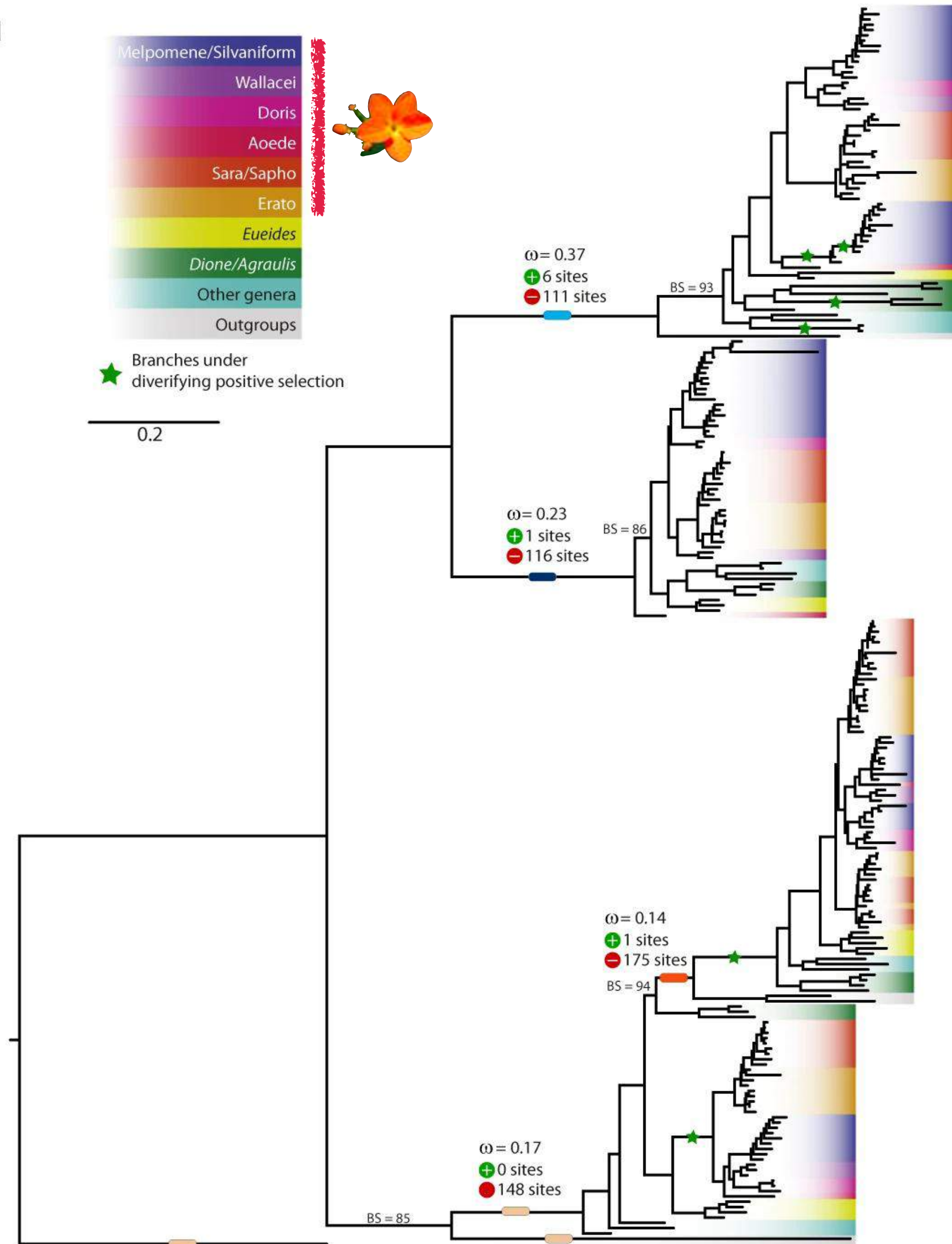Photo credit: @mena_sebas

# » (*Putative*) Gain of function »

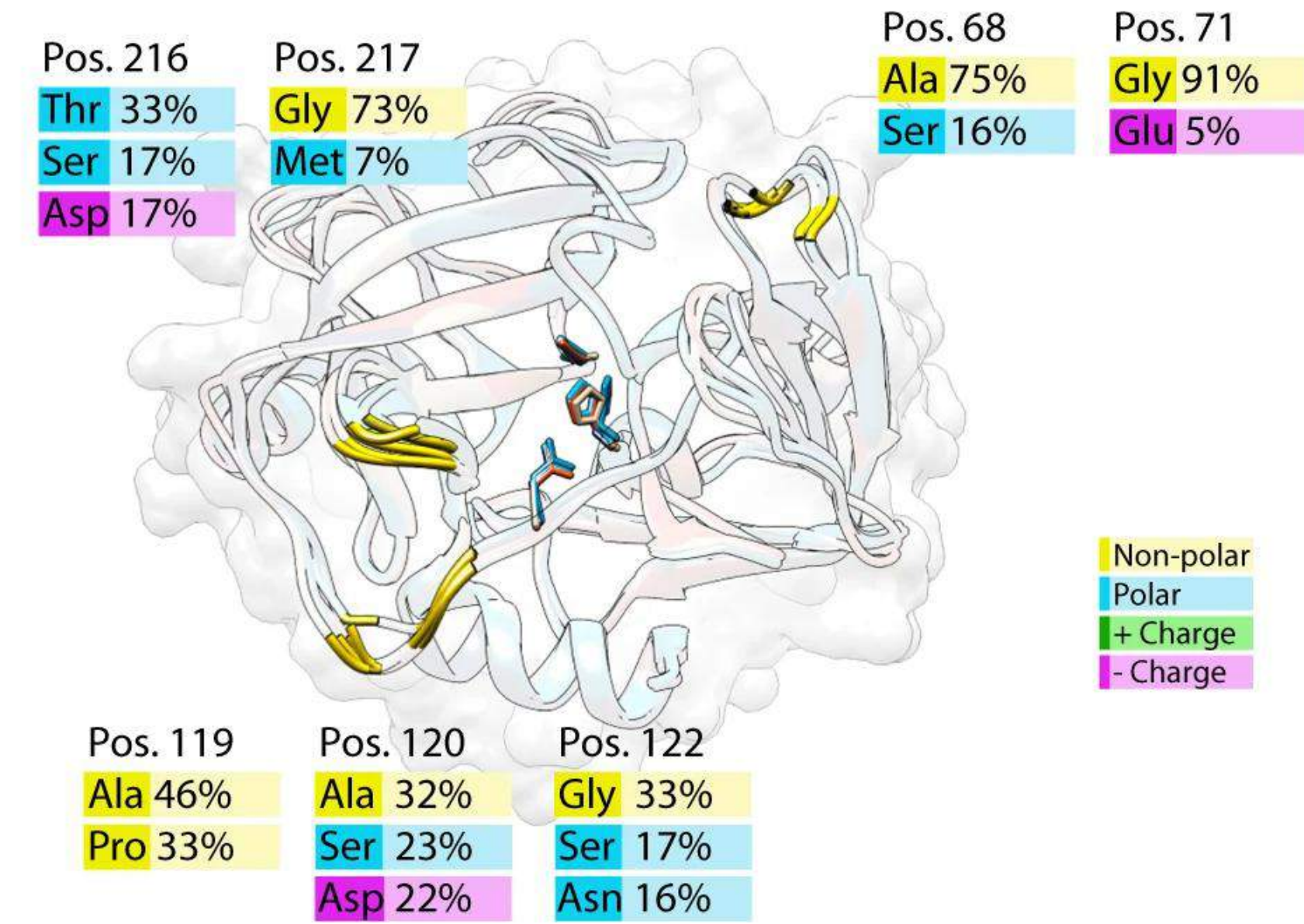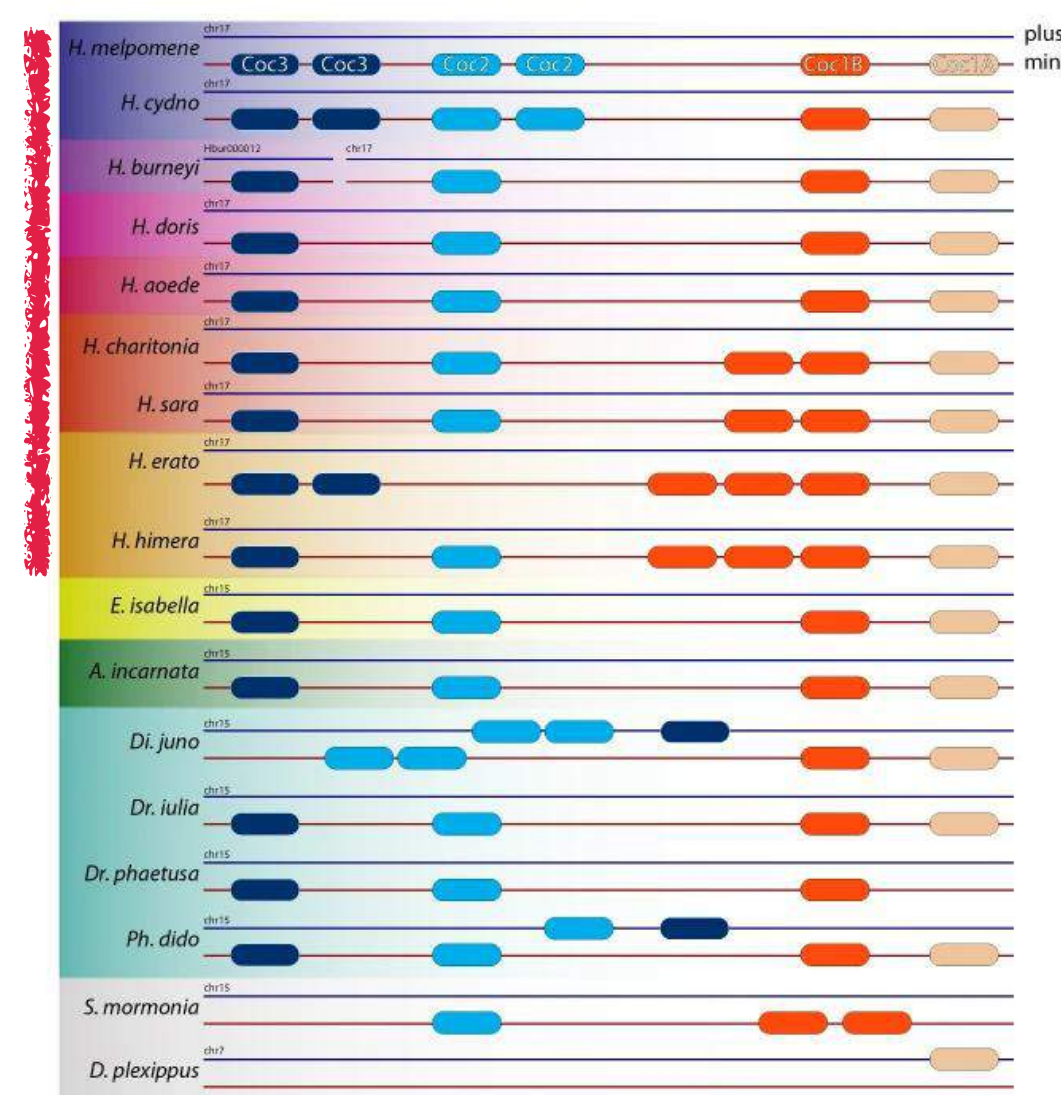# » (*Putative*) Gain of function »

# » Single cell RNA data integration »



https://genome.cshlp.org/content/33/1/96

https://elifesciences.org/articles/66747

# » Questions ? »

# » Methods »

# » Phylogeny on Homologous loci (Tree-based) »



```
X1  ATGCTTAGGTCAGCTAGATAGTGGCTAATACCTAGCAGTTGAGTAA
Y1  ATGCTTAGGTCAGCATAGTGGCTAATACCTAGCAGTTGAGTAA
X2  ATGCTTAGGTCAGCTAGATAGTGGCTAATACCTAGCAGTTGAGTAA
Y2  ATGCTTAGGTCAGCATAGTGGCTAGTACCTAGCAGTTGAGTAA
Ya3 ATGCTTAGGTCAGCTAGATAGTGGCTAATACCTAGCAGTTGAGTAA
Yb3 ATGCTTAGGTCAGCATAGTGGCTAATACCTAGCAGTTGAGTAA
```

```
X1  ATGCTTAGGTCAGCTAGATAGTGGCTAA---TACCTAGCAGTTGAGTAA
Y1  ATGCTTAGGTCAGC---ATAGTGGCTAA---TACCTAGCAGTTGAGTAA
X2  ATGCTTAGGTCAGCTAGATAGTGGCTAA---TACCTAGCAGTTGAGTAA
Y2  ATGCTTAGGTCAGC---ATAGTGGCTAG---TACCTAGCAGTTGAGTAA
Ya3 ATGCTTAGGTCAGCTAGATAGTGGCTAA---TACCTAGCAGTTGAGTAA
Yb3 ATGCTTAGGTCAGC---ATAGTGGCTAA---TACCTAGCAGTTGAGTAA
```

**Substitution model**

# » Homology based-methods »

**BLAST** (**B**asic **L**ocal **A**lignment **S**earch **T**ool) is a widely used algorithm for comparing biological sequences.

# » Homology based-methods »

**BLAST** identifies regions of <u>*local similarity*</u> between sequences by breaking the query and database sequences into smaller segments (called words) and then aligning these segments.

# It does not identify ORTHOLOGY *per se!!!*

It is a **sequence similarity search tool**, and while it can help identify genes with similar sequences, it does not directly determine **orthology**

# » Homology based-methods »

You can implement a **reciprocal Best Hits (RBH)** to identify two genes which are each other's best hits across species, but this is an approximation and not proof.

It can fail in cases like gene duplication, loss, or incomplete genomes.

# » Hidden Markov Models (HMMs) »

It's a *probabilistic model* used to describe sequences (DNA, RNA, or proteins) based on their underlying statistical properties.

# » Hidden Markov Models (HMMs) »

# BUSCO



**Dataset: signature of BUSCO genes**

| AA consensus sequences | Block profiles | Profile HMMs | Score + length cutoffs |

AUGUSTUS-retraining

A) Genome Assembly → *tblastn* → *AUGUSTUS* → *hmmsearch* → BUSCO Classifier
\ contigs / \ candidate regions / \ candidate proteins /

B) Annotated Gene Set - candidate proteins - *hmmsearch* → BUSCO Classifier

C) Transcriptome → tblastn → *hmmsearch* → BUSCO Classifier
\ transcripts / \ longest ORFs as candidate proteins /

Report and outputs

BUSCO: Assessing Genome Assembly and Annotation Completeness | *https://link.springer.com/protocol/10.1007/978-1-4939-9173-0_14*

# » **Methods | Network based** (with some integration of the phylogeny) »



**Pairwise generated metric**

**E-value:** the number of alignments you would expect to find by chance in a database of a given size.
**Bitscore:** a normalized score that represents the quality of the alignment. It's independent of the database size.

# » Methods | Network based (with some integration of the phylogeny) »

# » Methods | Whole-genome alignment based »

## Box 3 | Genuine and false implications of orthology and paralogy relationships

- Orthologues form a clade (that is, they are monophyletic) in an accurate phylogenetic tree. This is a necessary corollary of the orthology definition (BOX 1).

- Orthology does not imply a one-to-one relationship between genes from different organisms. Lineage-specific gene duplications often lead to one-to-many and many-to-many co-orthology relationships (BOX 1).

- The molecular clock is not implicit in the definition of orthology: orthologues in different lineages may evolve at different (in principle, arbitrarily different) rates (BOX 1).

- Conservation of sequence, structure or genomic context is not implicit in the definition of orthology.

- Given the above, orthology does not necessarily imply that orthologous genes (even in the absence of lineage-specific duplications) are the most similar sequences or structures in compared genomes.

- The converse is not necessarily true either: genes that are most similar to each other in compared genomes (often denoted bidirectional best hits (BBHs)) might not be orthologous. The BBHs may represent cryptic paralogy after differential loss of ancestral paralogues in compared lineages or xenologues, whereby one of the genes in a BBH pair was acquired by horizontal gene transfer.

- Orthology does not necessarily imply conservation of gene function.

- The converse is not necessarily true either: genes with equivalent functions are not necessarily orthologous.

- All of the above caveats notwithstanding, the generalized orthology conjecture predicts that, as a genome-wide statistical trend, orthologues are the most similar genes in different species, in terms of sequence, structure and function.

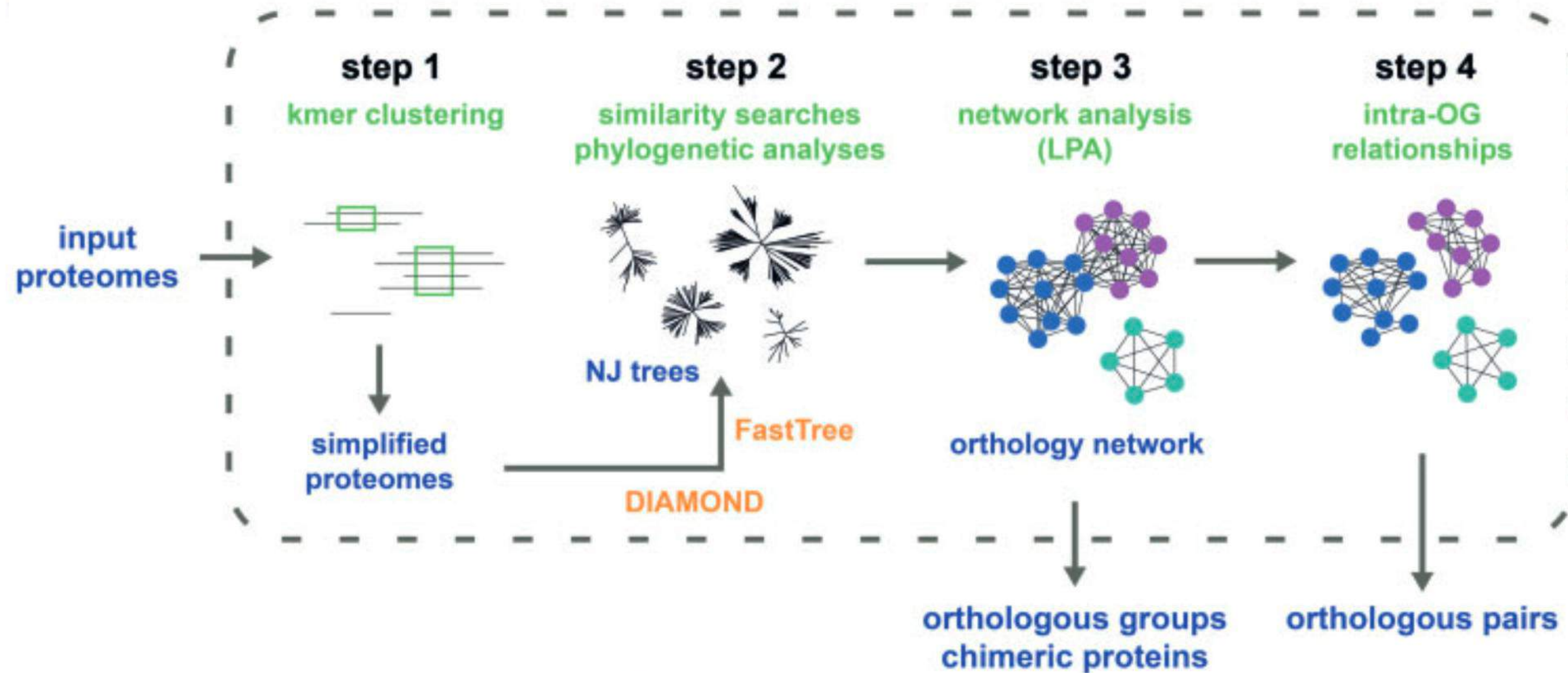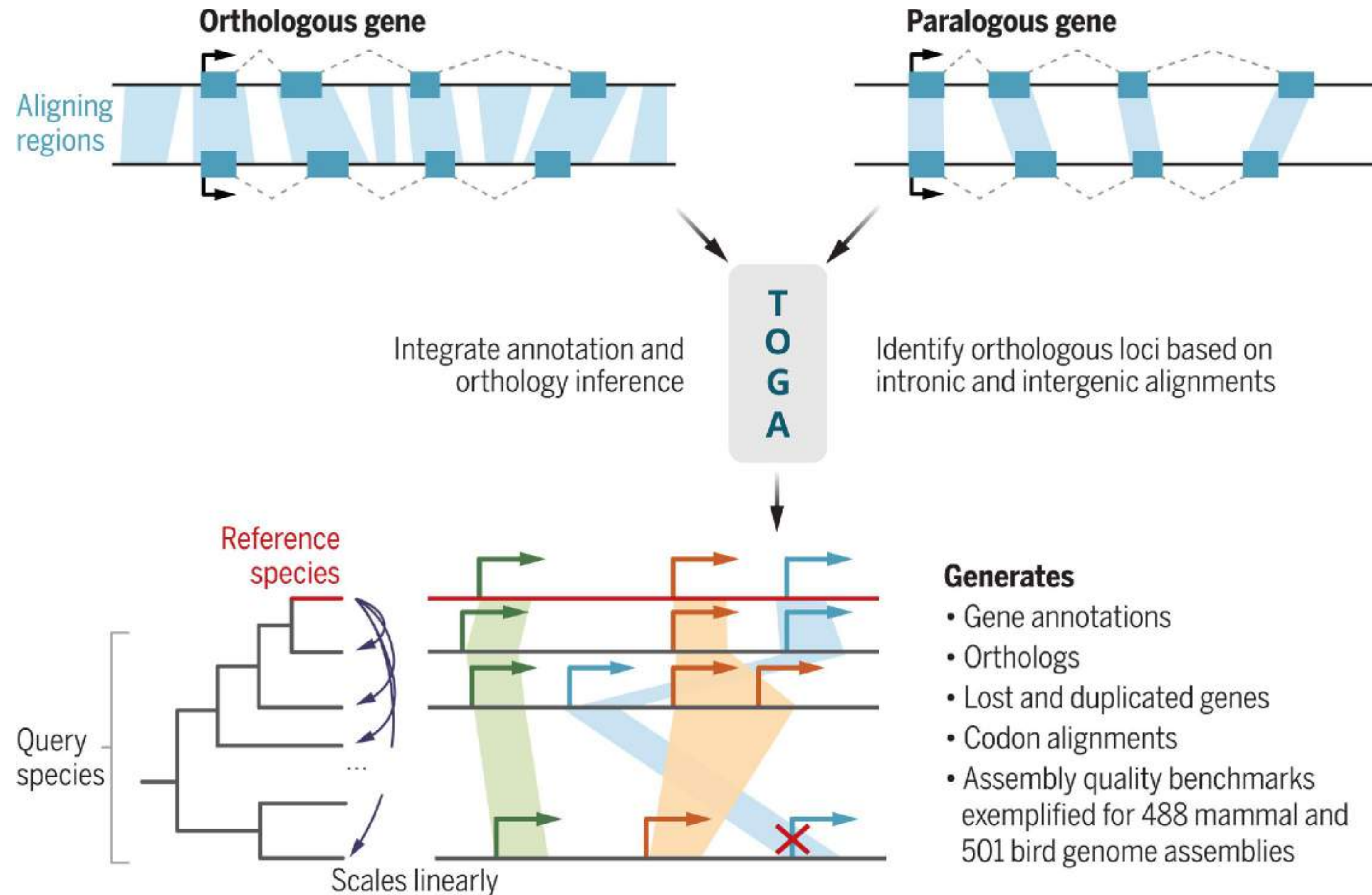- Paralogy applies to genes not only within species (as often assumed) but also between species; in cases of differential gene loss and complex evolutionary scenarios, distinguishing orthology and paralogy may be non-trivial (BOX 1).

- Paralogy does not necessarily imply functional divergence (as is often assumed): for instance, paralogy may contribute to protein dosage modulation.

- Nevertheless, the generalized orthology conjecture implies that, as a general trend, paralogues are more functionally different than orthologues at the same level of sequence divergence.

# » Questions ? / Break ? »

# » Whole Genome alignment »

## » Conserved non-coding Elements (CNEs) »

## » CNEs enrichment (extra) »

» Whole Genome alignment »

» Conserved non-coding Elements (CNEs) »

» CNEEs enrichment (extra) »

# » Whole Genome alignment »

# » Whole Genome alignment »

> **Deeper Species-tree inference**
> **Comparative gene annotation**
> **Detection of selection / Conservation**

- **Multi-species map of genomic regions to a corresponding region in each other genome.**

- **Taking into account complex rearrangements and copy number changes.**

Common limitations are 'reference bias':
 - Constrains a multiple alignment to only regions present in reference genome.
 - Restricting the alignment to be 'single-copy', determining miss multiple-orthology relationships.

**Cactus (ProgressiveCactus) is a "reference *free*" whole genome aligner.**

# » Whole Genome alignment »



Armstron *et al* **2020** *Nature*

# » Some of the applications (annotation) »



Fiddes *et al* **2018** *Genome Research*

# » Some of the applications (orthology inference) »



Kirilenco *et al* **2023** *Science*

# » Phylogeny / Introgression »

# » Phylogeny / Introgression »

# » Phylogeny / Introgression »

# » Structural Rearrangements »

# » Structural Rearrangements »



Schrader *et al* **2021** *Nature Comm.*

Increased Rate of Duplications

Increased Rate of Deletion

Expanded OGs

Contracted OGs

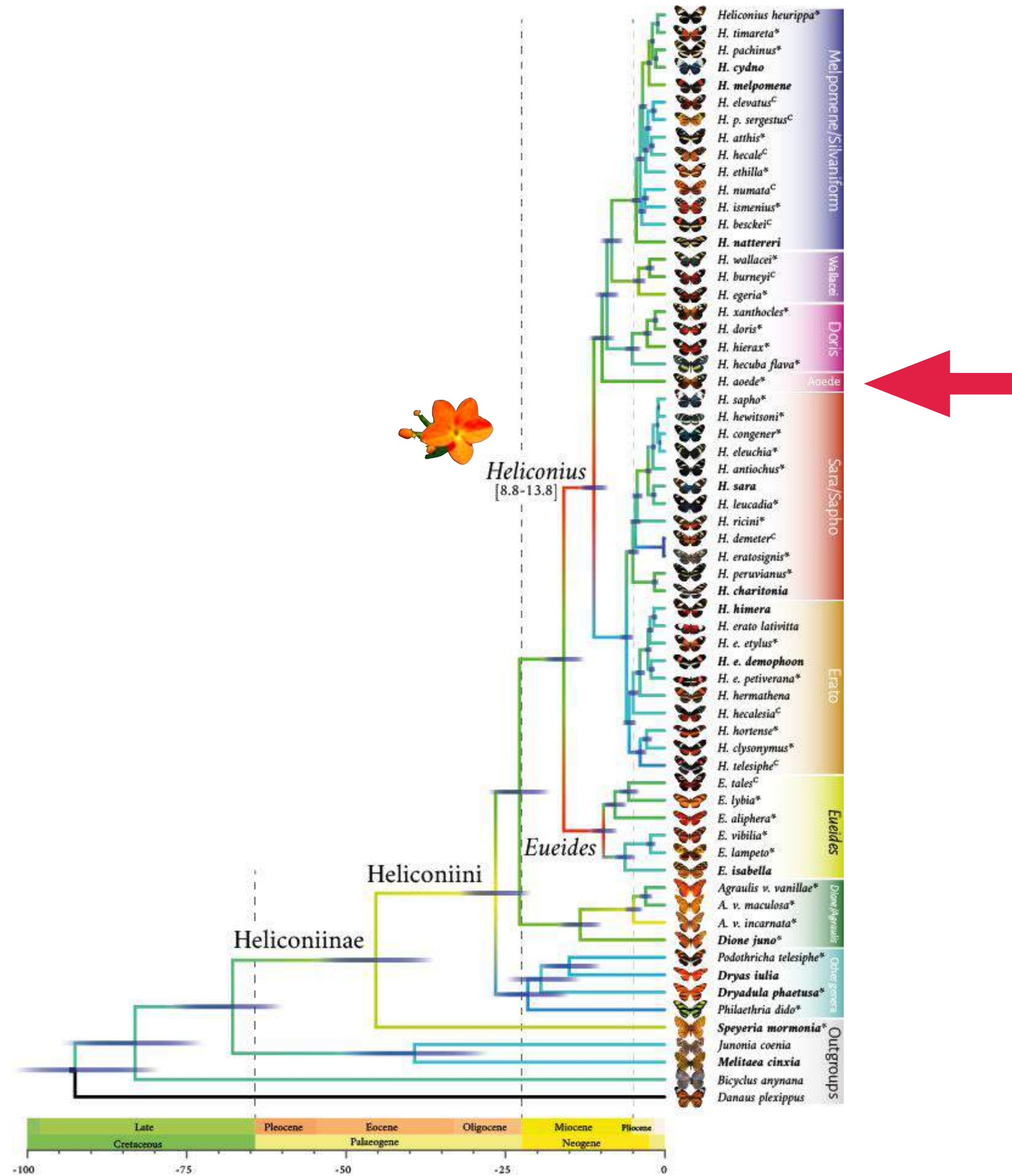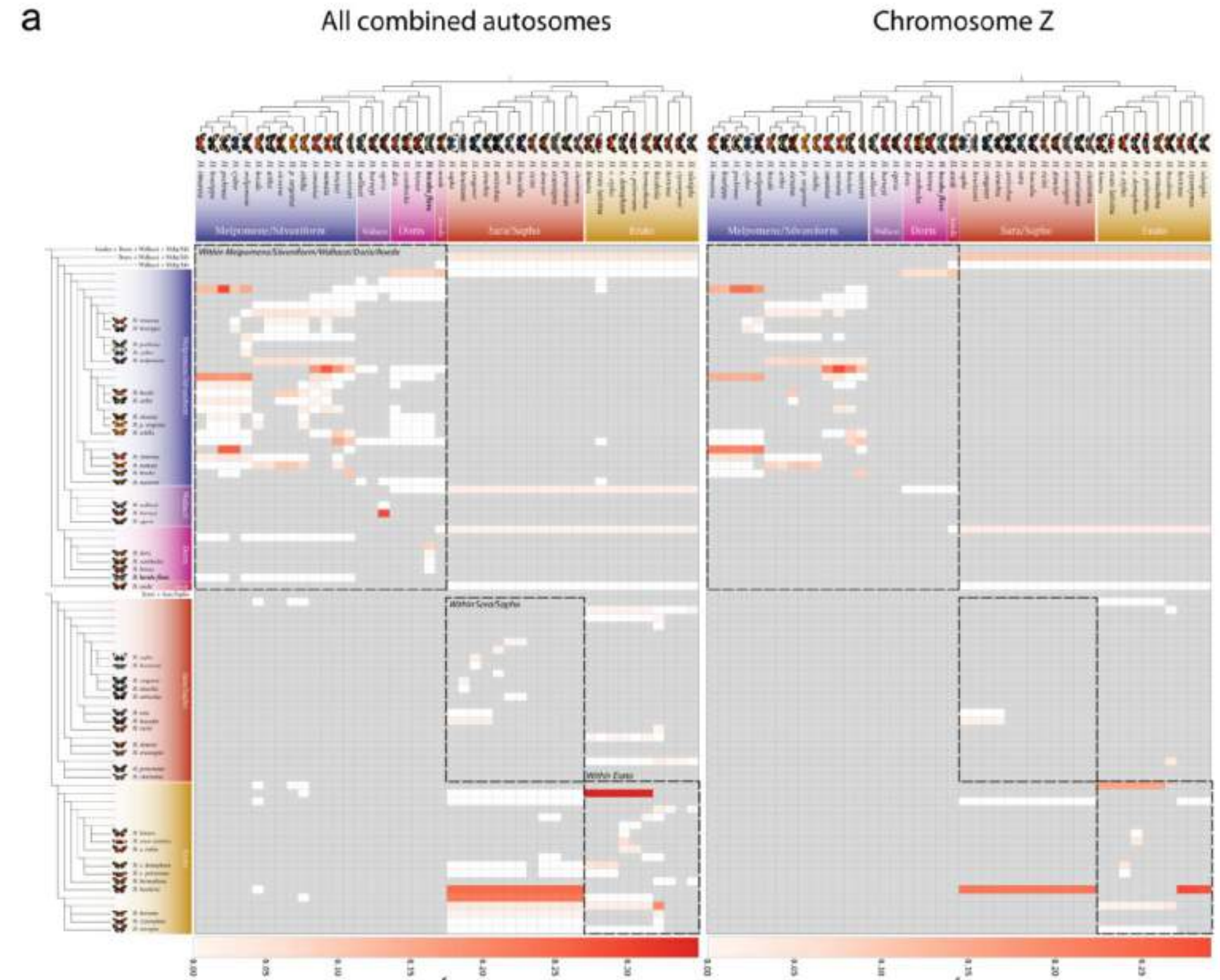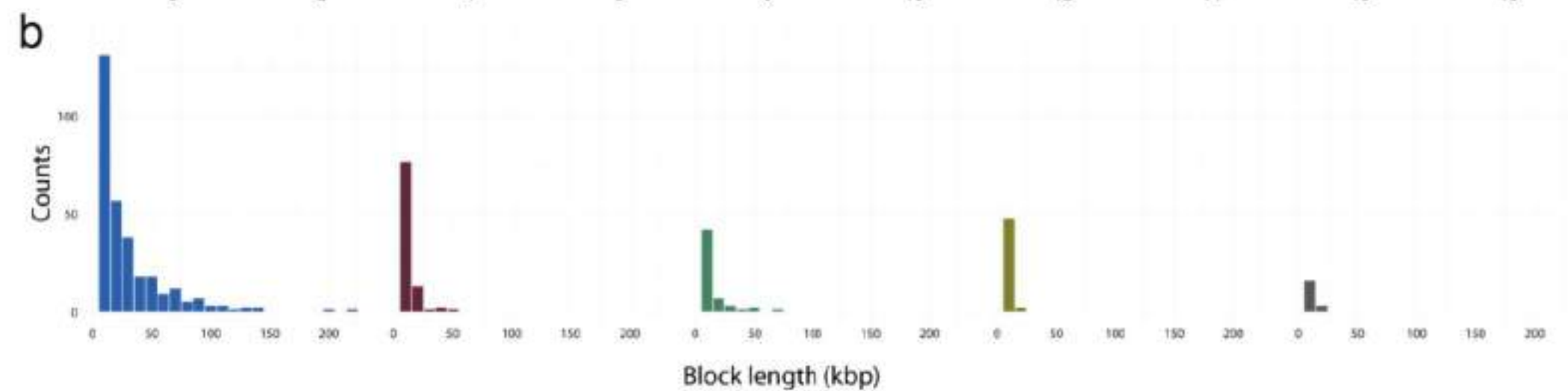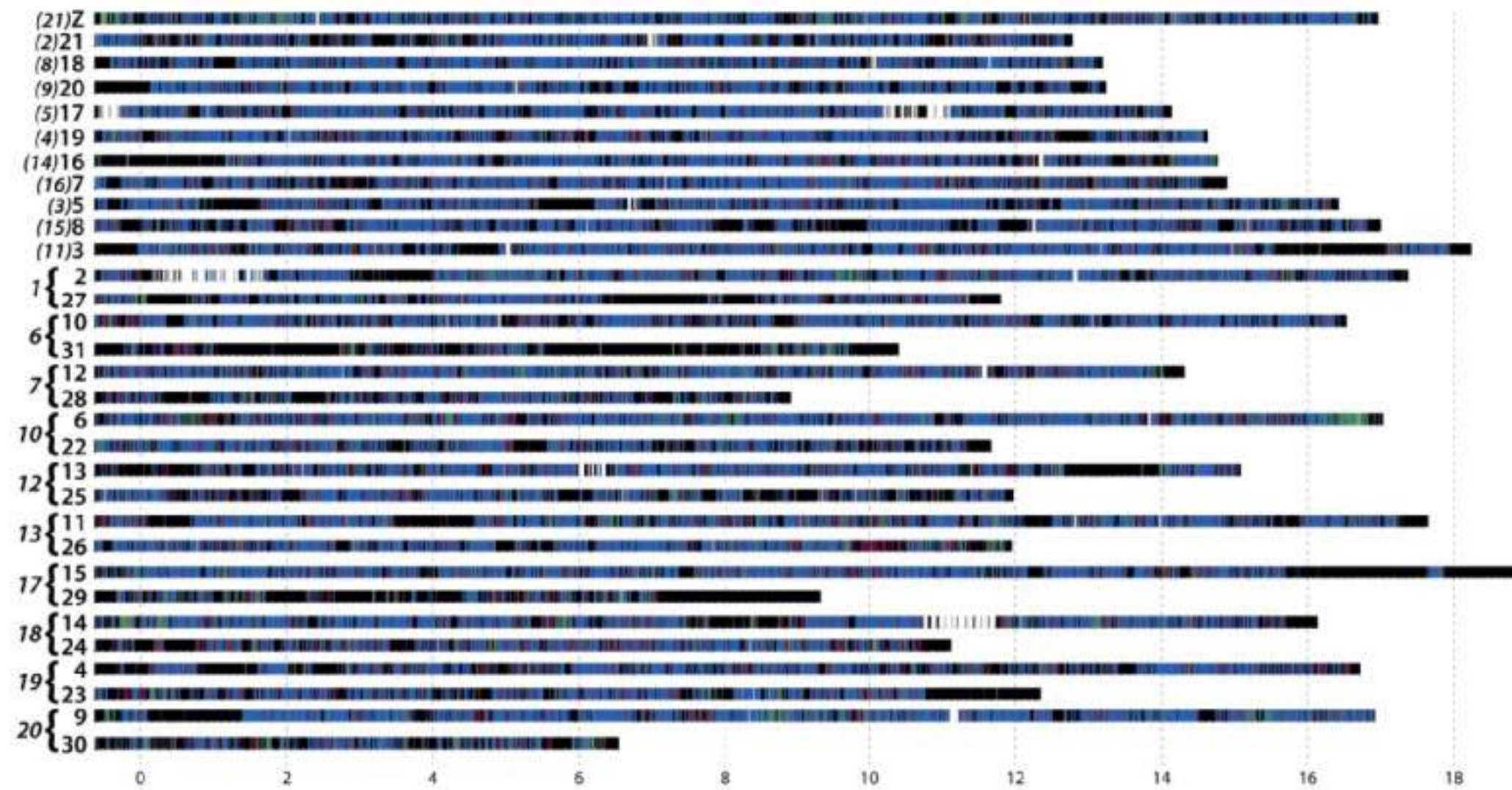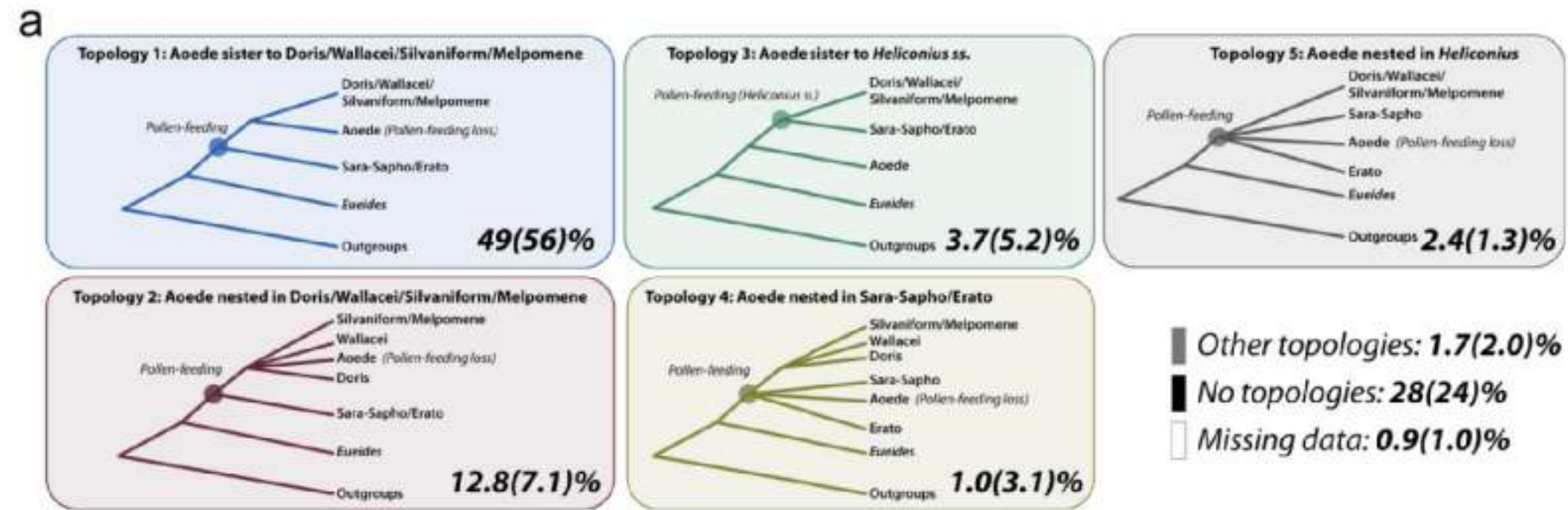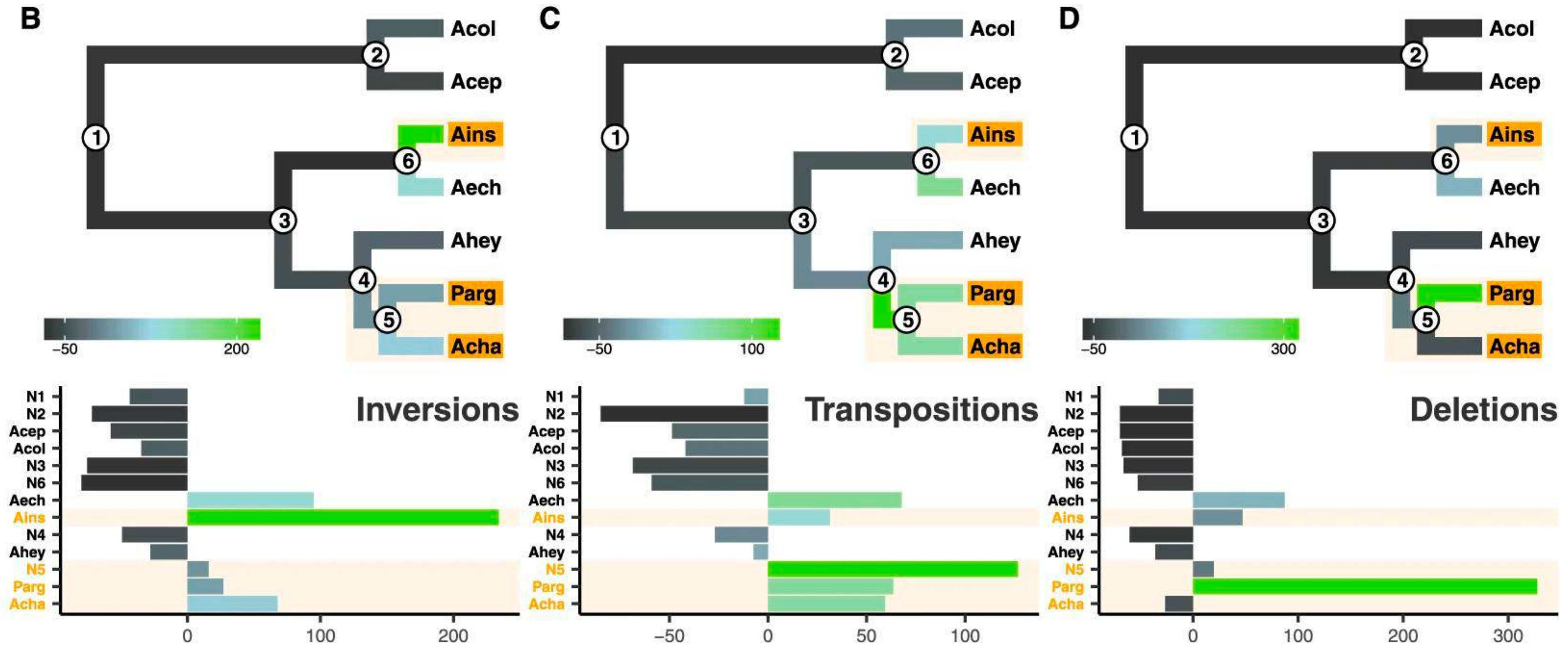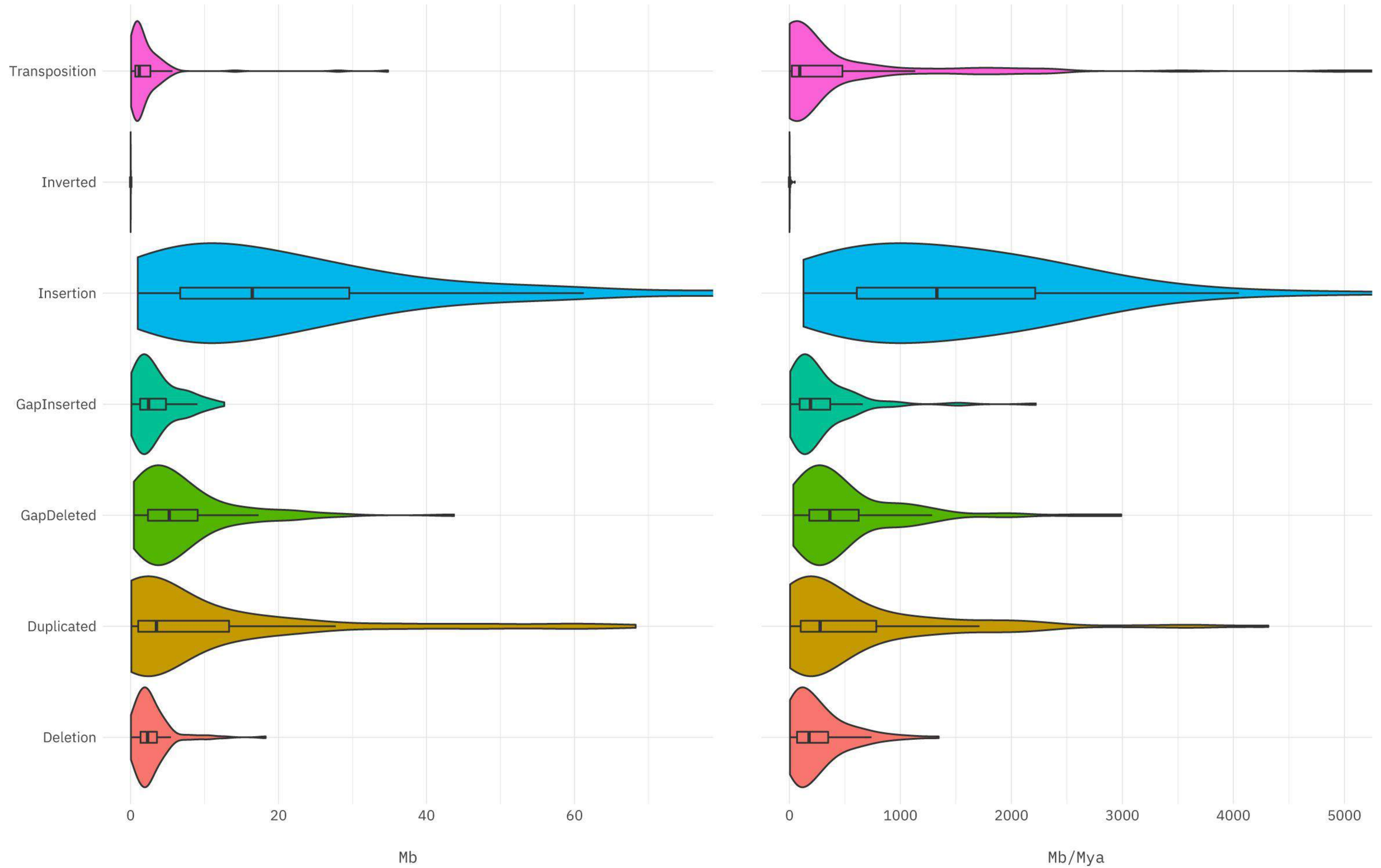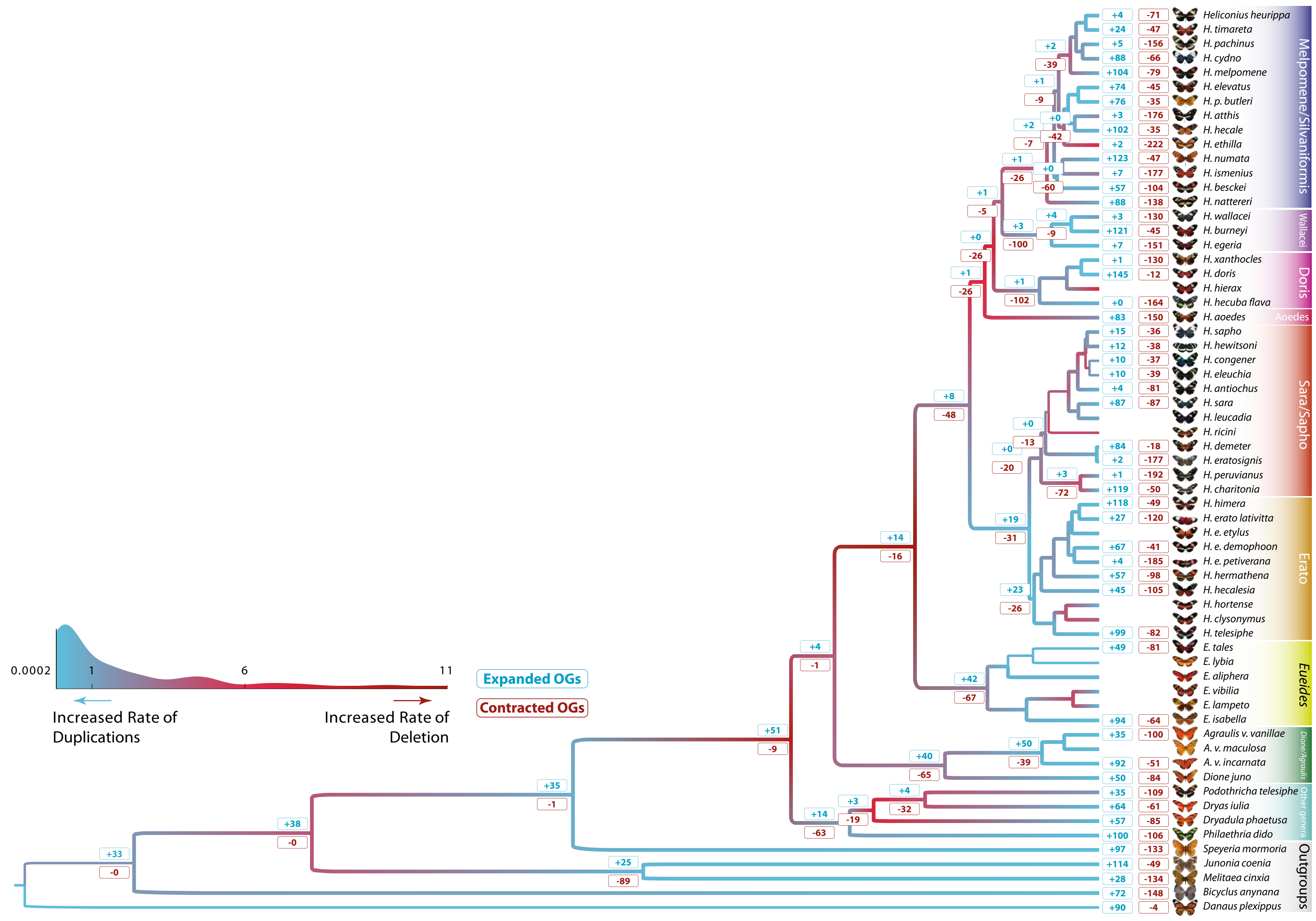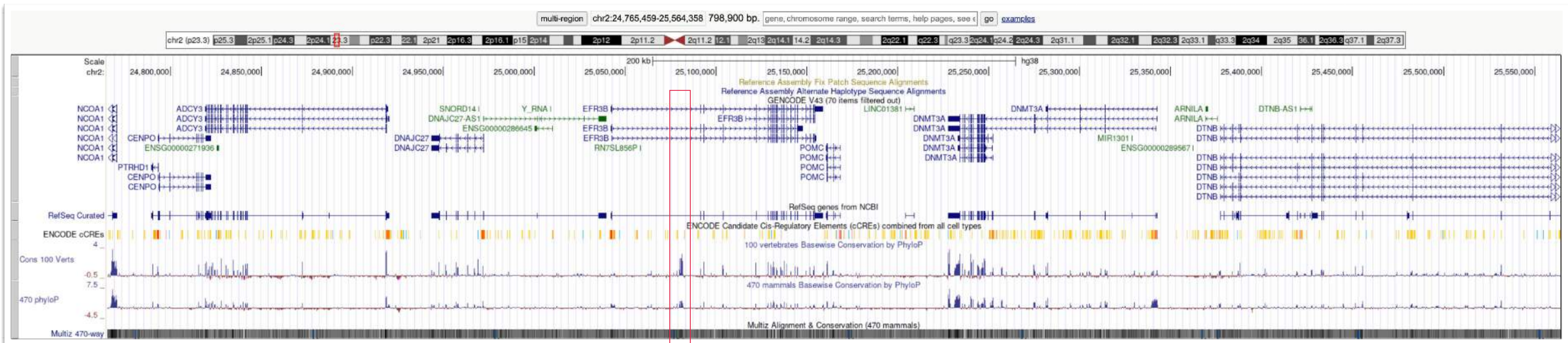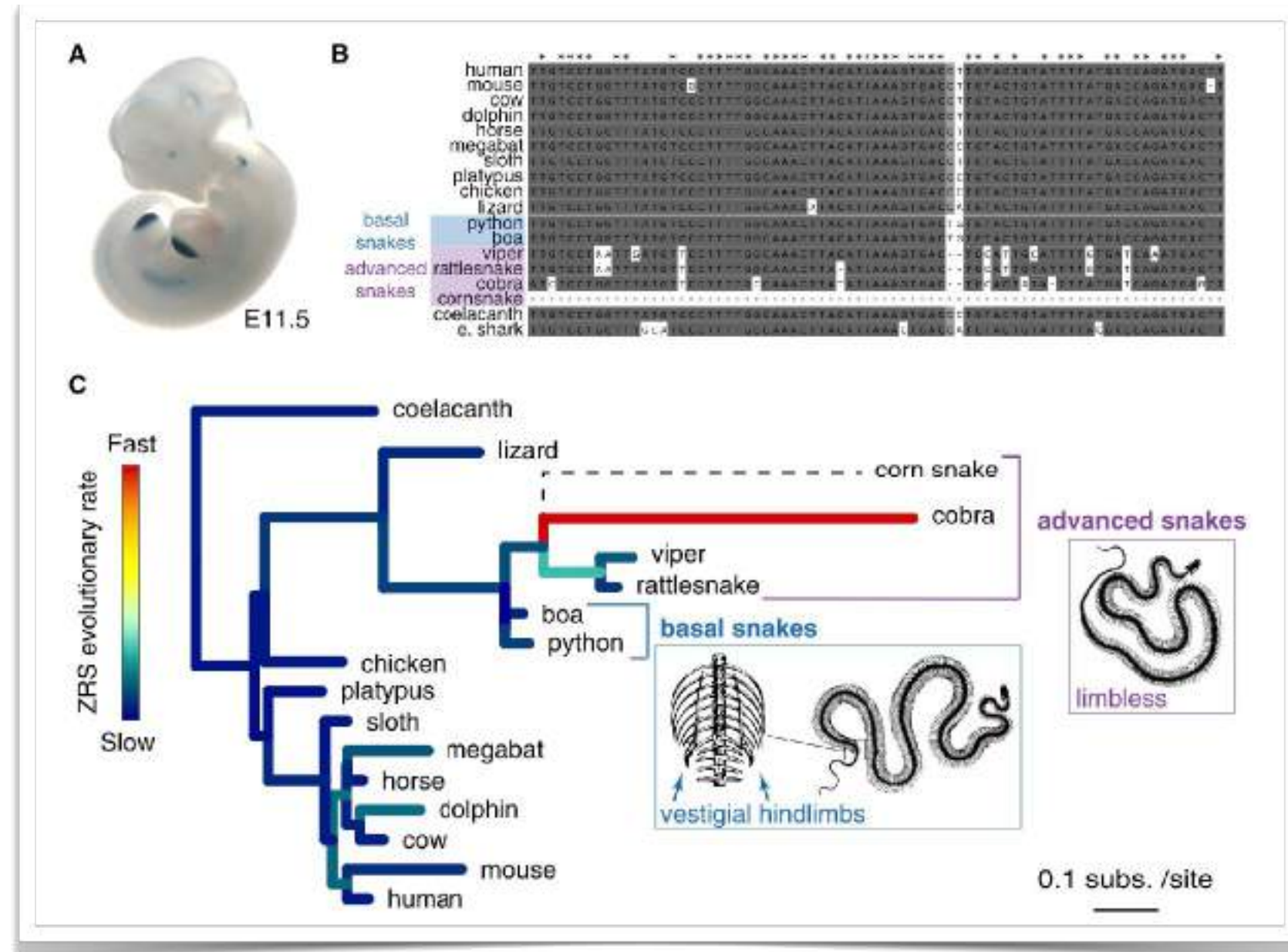| | | |
|---|---|---|
| +4 / −71 | *Heliconius heurippa* | Melpomene/Silvaniformis |
| +24 / −47 | *H. timareta* | |
| +2 / +5 / −156 | *H. pachinus* | |
| −39 / +88 / −66 | *H. cydno* | |
| +104 / −79 | *H. melpomene* | |
| +1 / +74 / −45 | *H. elevatus* | |
| −9 / +76 / −35 | *H. p. butleri* | |
| +0 / +3 / −176 | *H. atthis* | |
| +2 / +102 / −35 | *H. hecale* | |
| −7 / −42 / +2 / −222 | *H. ethilla* | |
| +1 / +123 / −47 | *H. numata* | |
| +0 / +7 / −177 | *H. ismenius* | |
| −26 / −60 / +57 / −104 | *H. besckei* | |
| +1 / +88 / −138 | *H. nattereri* | |
| −5 / +4 / +3 / −130 | *H. wallacei* | Wallacei |
| +3 / +121 / −45 | *H. burneyi* | |
| +0 / −9 / +7 / −151 | *H. egeria* | |
| −26 / +1 / −130 | *H. xanthocles* | Doris |
| +1 / +145 / −12 | *H. doris* | |
| −26 / +1 / −130 | *H. hierax* | |
| −102 / +0 / −164 | *H. hecuba flava* | |
| +83 / −150 | *H. aoedes* | Aoedes |
| +15 / −36 | *H. sapho* | |
| +12 / −38 | *H. hewitsoni* | |
| +10 / −37 | *H. congener* | |
| +10 / −39 | *H. eleuchia* | |
| +8 / −48 / +4 / −81 | *H. antiochus* | Sara/Sapho |
| +87 / −87 | *H. sara* | |
| +0 / *H. leucadia* | | |
| +0 / −13 / +84 / −18 | *H. ricini* | |
| +2 / −177 | *H. demeter* | |
| −20 / +3 / +1 / −192 | *H. eratosignis* | |
| −72 / +119 / −50 | *H. peruvianus* | |
| +118 / −49 | *H. charitonia* | |
| +14 / −16 / +19 / +27 / −120 | *H. himera* | |
| −31 / *H. erato lativitta* | | Erato |
| +67 / −41 | *H. e. etylus* | |
| +4 / −185 | *H. e. demophoon* | |
| +23 / +57 / −98 | *H. e. petiverana* | |
| −26 / +45 / −105 | *H. hermathena* | |
| *H. hecalesia* | | |
| +99 / −82 | *H. hortense* | |
| +4 / +49 / −81 | *H. clysonymus* | |
| −1 / *H. telesiphe* | | |
| +42 / *E. tales* | | *Eueides* |
| −67 / *E. lybia* | | |
| +94 / −64 | *E. aliphera* | |
| +51 / −9 / +50 / +35 / −100 | *E. vibilia* | |
| *E. lampeto* | | |
| +40 / +92 / −51 | *E. isabella* | |
| −39 / *Agraulis v. vanillae* | | Dione/Agraulis |
| −65 / +50 / −84 | *A. v. maculosa* | |
| +4 / +35 / −109 | *A. v. incarnata* | |
| +3 / +64 / −61 | *Dione juno* | |
| +14 / −32 / +57 / −85 | *Podothricha telesiphe* | Other genera |
| −19 / −63 / +100 / −106 | *Dryas iulia* | |
| +35 / −1 / +97 / −133 | *Dryadula phaetusa* | |
| +38 / −0 / +114 / −49 | *Philaethria dido* | |
| +33 / −0 / +25 / −89 | *Speyeria mormoria* | Outgroups |
| +28 / −134 | *Junonia coenia* | |
| +72 / −148 | *Melitaea cinxia* | |
| +90 / −4 | *Bicyclus anynana* | |
| | *Danaus plexippus* | |

10Mya

# » Conserved non-coding Elements (CNEs) »

# » Conserved non-coding Elements (CNEs) »

> A class of *non-protein-coding* genomic sequences with elevated degree of conservation.

> CNEs are non-randomly distributed, clustering in the vicinity of genes with regulatory roles.

> Organised into functional ensembles (regulatory blocks), which coordinate the expression of shared target genes.

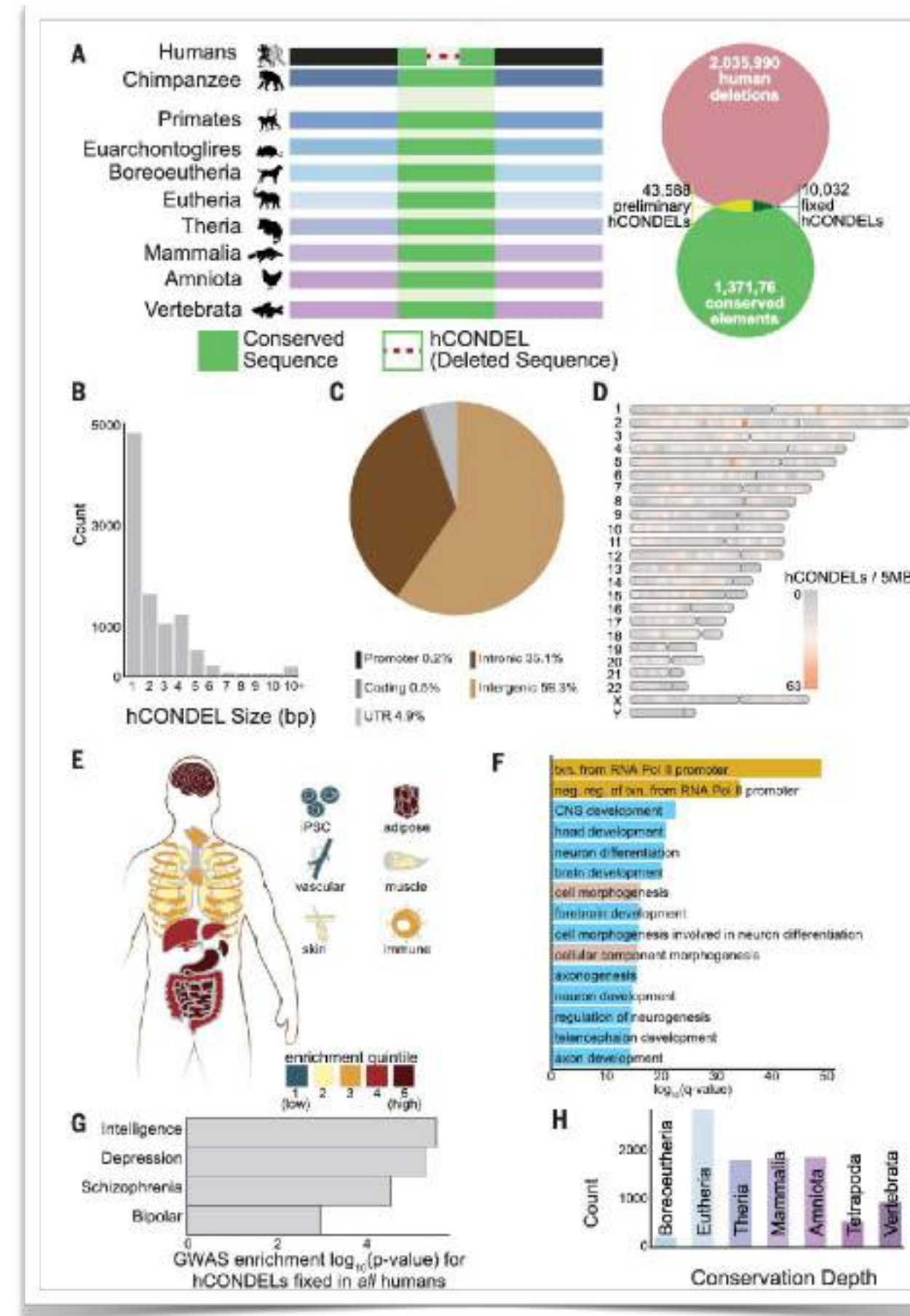> The disruption of these elements contribute to diseases linked with development, and cancer.
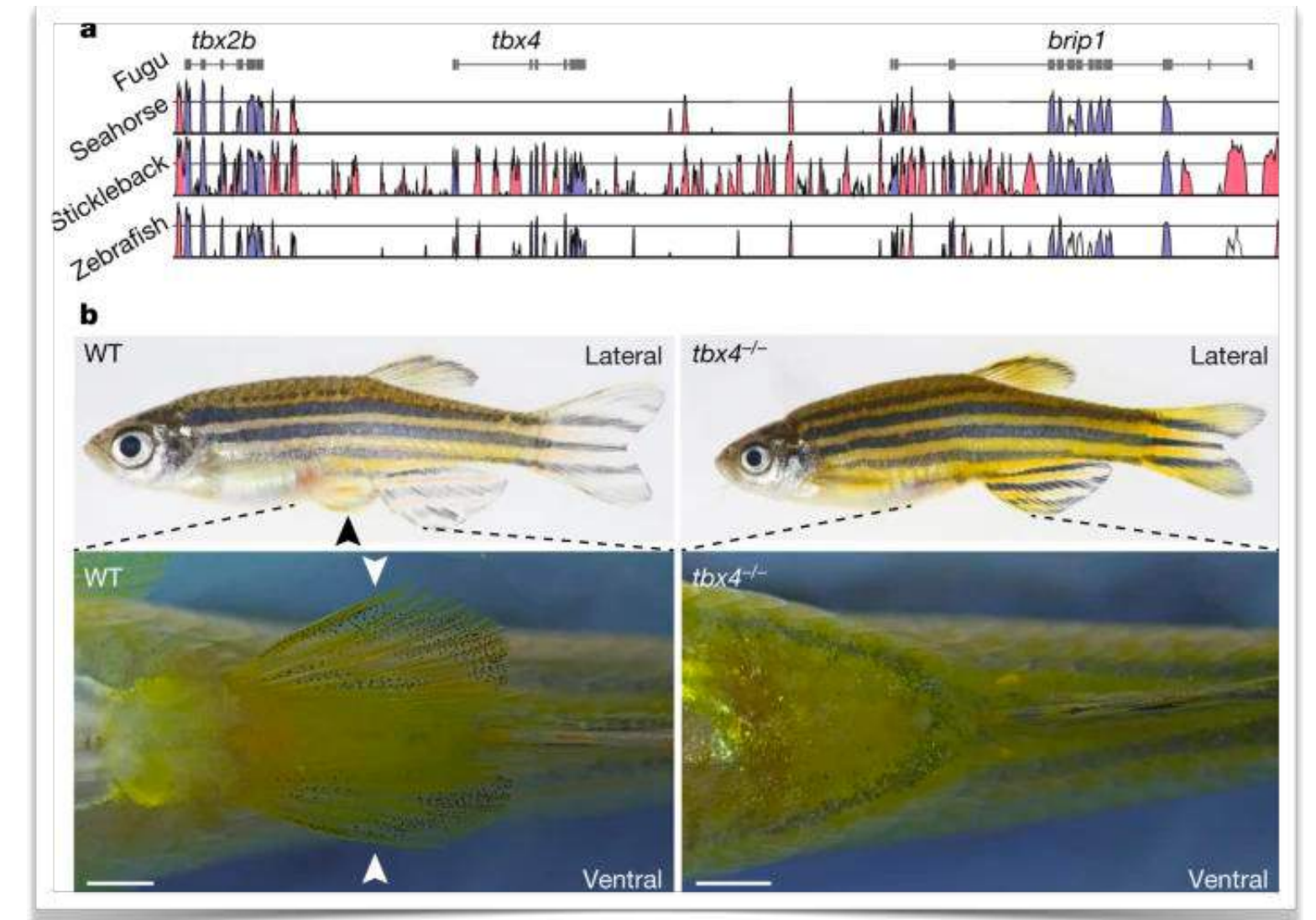


Candidate Cis-Regulatory Elements (cCREs)

# » Conserved non-coding Elements (CNEs) »



Kvon *et al* **2016** *Cell*



Xue *et al* **2023** *Science*
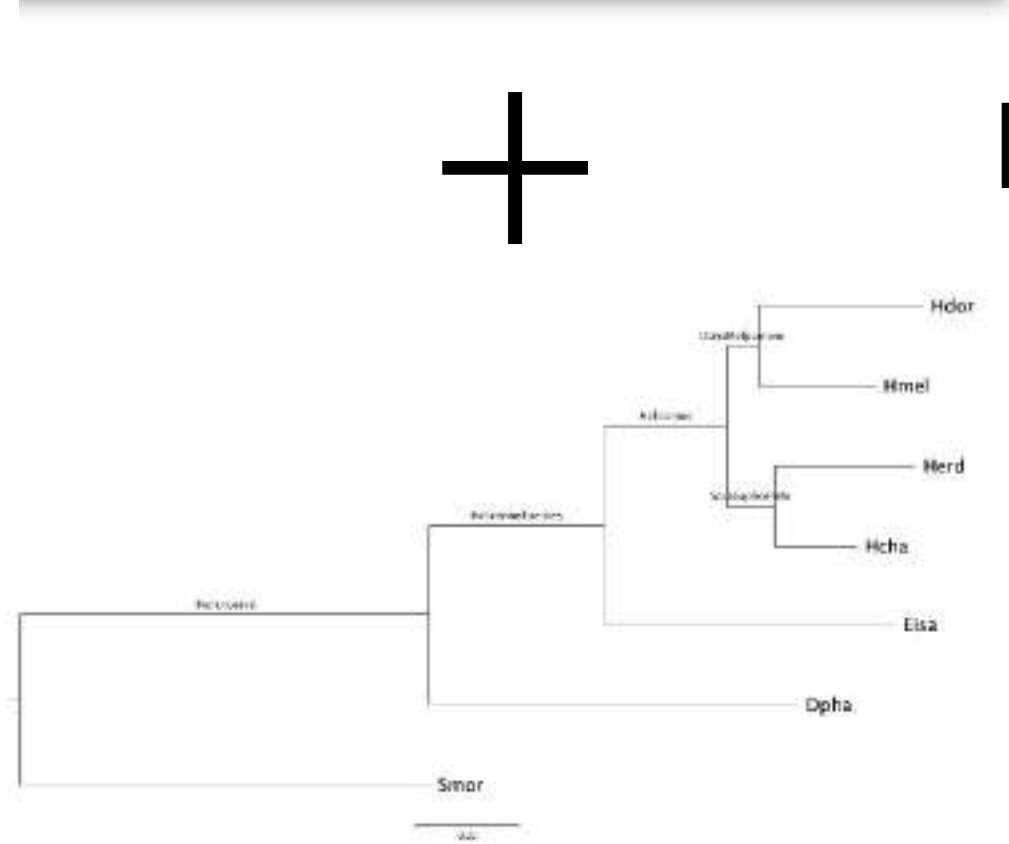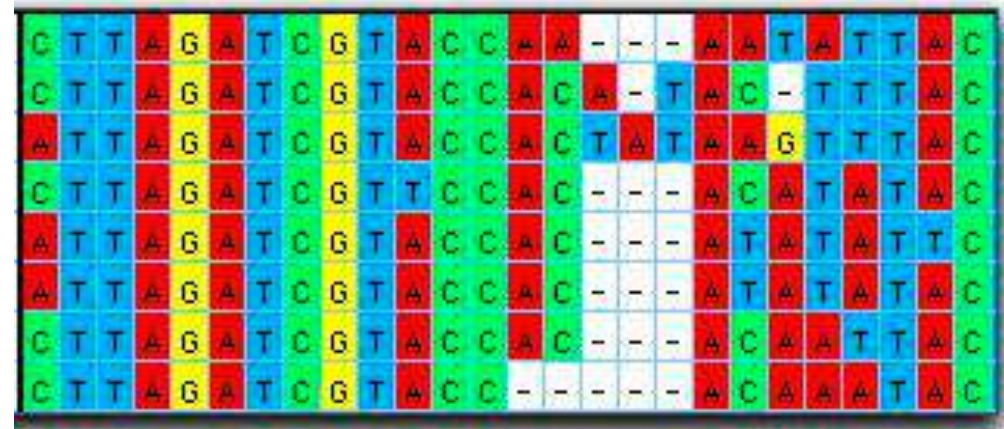


Lin *et al* **2016** *Nature*

# » How to identify these regions »

1. **A Model of DNA sequence evolution**

2. **Phylogenetic tree**

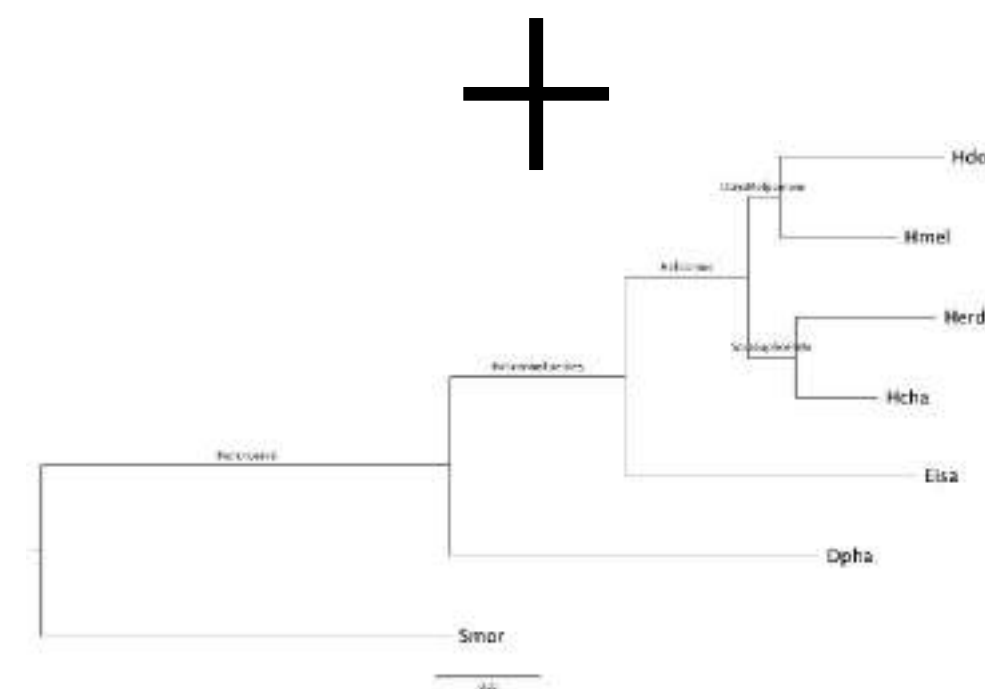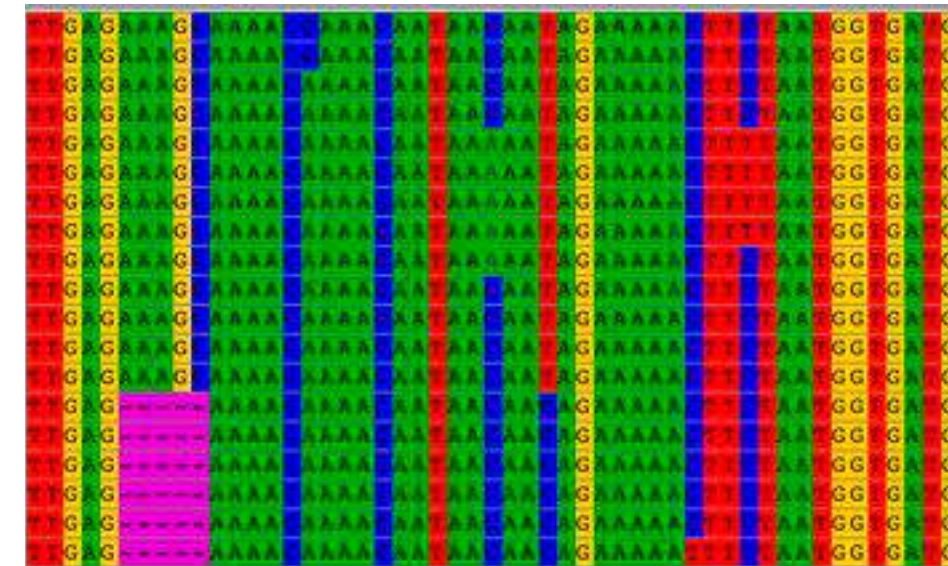# » How to identify these regions »

Neutral evolving regions
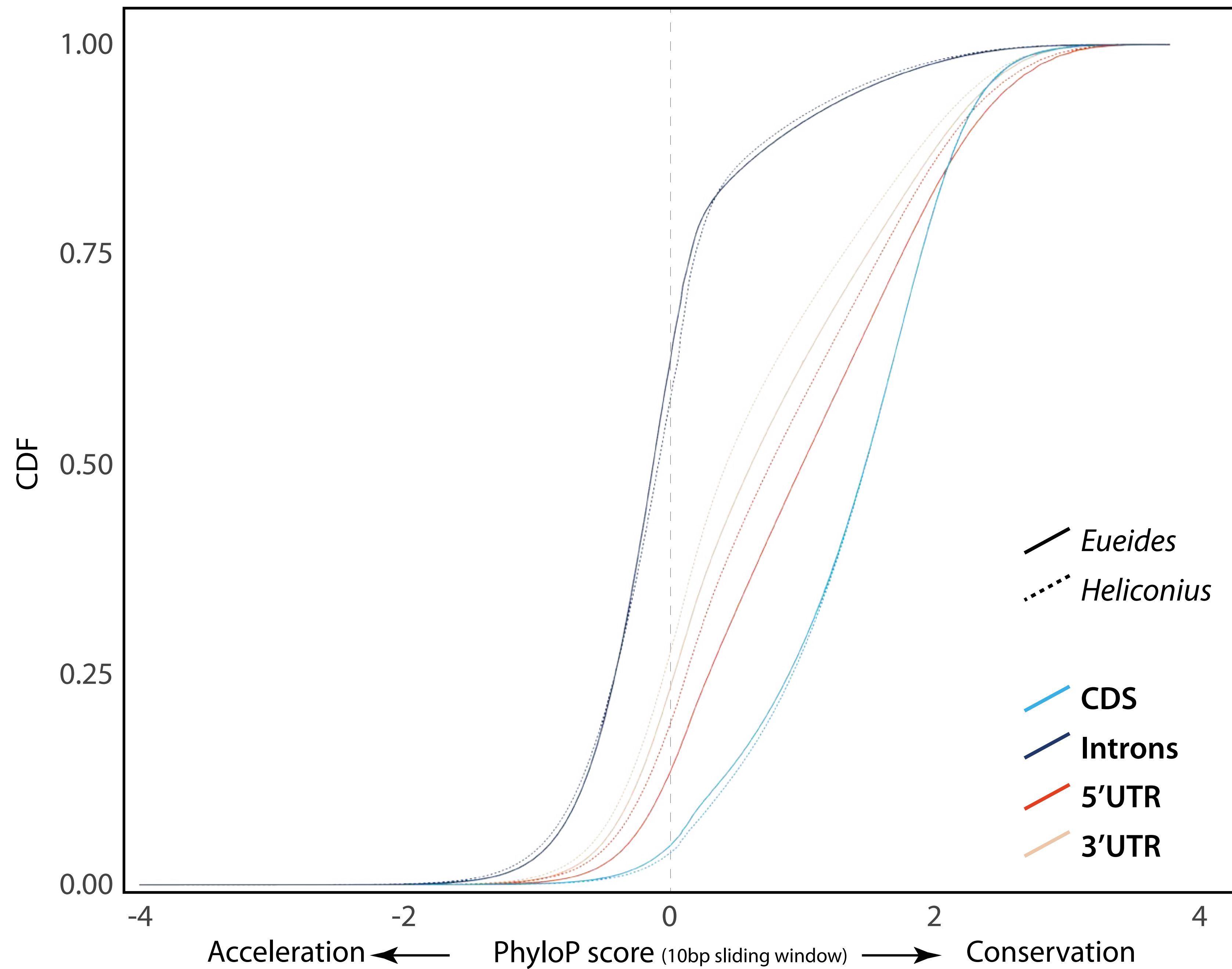


WGA



+

**Neutral model**

Other metrics
(*PhyloP*)

+

+

**neutral model**

**Conservation model**

**Non-conserved model**

CDF

*Eueides*
*Heliconius*

CDS
Introns
5'UTR
3'UTR

Acceleration ← PhyloP score (10bp sliding window) → Conservation

-4    -2    0    2    4

0.00    0.25    0.50    0.75    1.00

# » Let's start the exercise »

**Some tools you probably going to need:**

- Cactus
- halSummarizeMutations
- halAlignmentDepth
- wigToBigWig
- samtools
- .... and more ...

***Don't forget IGV!***

# » ATAC-seq (**A**ssay for **T**ransposase-**A**ccessible **C**hromatin using **seq**uencing) »



Assay for transposase accessible chromatin (ATAC-Seq)

Open DNA

Tn5 Transposome

Insert in regions of open chromatin

Fragmented and primed

DNA purification Amplification

DNA

» ATAC-seq (**A**ssay for **T**ransposase-**A**ccessible **C**hromatin using **seq**uencing) »

# » ATAC-seq (**A**ssay for **T**ransposase-**A**ccessible **C**hromatin using **seq**uencing) »



Assay for transposase accessible chromatin (ATAC-Seq)

Open DNA

Tn5 Transposome

Insert in regions of open chromatin

Fragmented and primed

DNA purification Amplification

DNA

## **Tagmentation**

The activity Tn5 (hyperactive) Transposase that <u>inserts</u> sequencing adapters into open regions of the genome and <u>cleaves</u>

# » ATAC-seq (**A**ssay for **T**ransposase-**A**ccessible **C**hromatin using **seq**uencing) »



Assay for transposase accessible chromatin (ATAC-Seq)

Open DNA

Tn5 Transposome

Insert in regions of open chromatin

Fragmented and primed

DNA purification Amplification

DNA

**Purification**

# » ATAC-seq (**A**ssay for **T**ransposase-**A**ccessible **C**hromatin using **seq**uencing) »



Assay for transposase accessible chromatin (ATAC-Seq)

Open DNA

Tn5 Transposome

Insert in regions of open chromatin

Fragmented and primed

DNA purification Amplification

DNA

**Short-Reads Sequencing**

# » ATAC-seq (**A**ssay for **T**ransposase-**A**ccessible **C**hromatin using **seq**uencing) »



Van Belleghem *et al* **2023** *Science*

# » Let's continue »

# » Notes »

**Some UNIX tools you probably going to need:**

- "**>**" Redirect

- "**|**" Pipe

- "`sed`" a stream editor

**For the extra task:**

- "`seq`" unix command to generate a sequence of numbers

- *For loop*

# » Test for Acceleration »

Heliconius

[8.8-13.8]

Eueides+Heliconius

Dione+Agraulis+
Eueides+Heliconius

Heliconiini        Eueides

Heliconiinae

## » *Mushroom body expansion (Brain)*



*Heliconius hecale*

*Eueides isabella*



Chromosome 20

8300kb          8325kb          8350kb          8375kb          8400kb

» *Mushroom body expansion (Brain)*
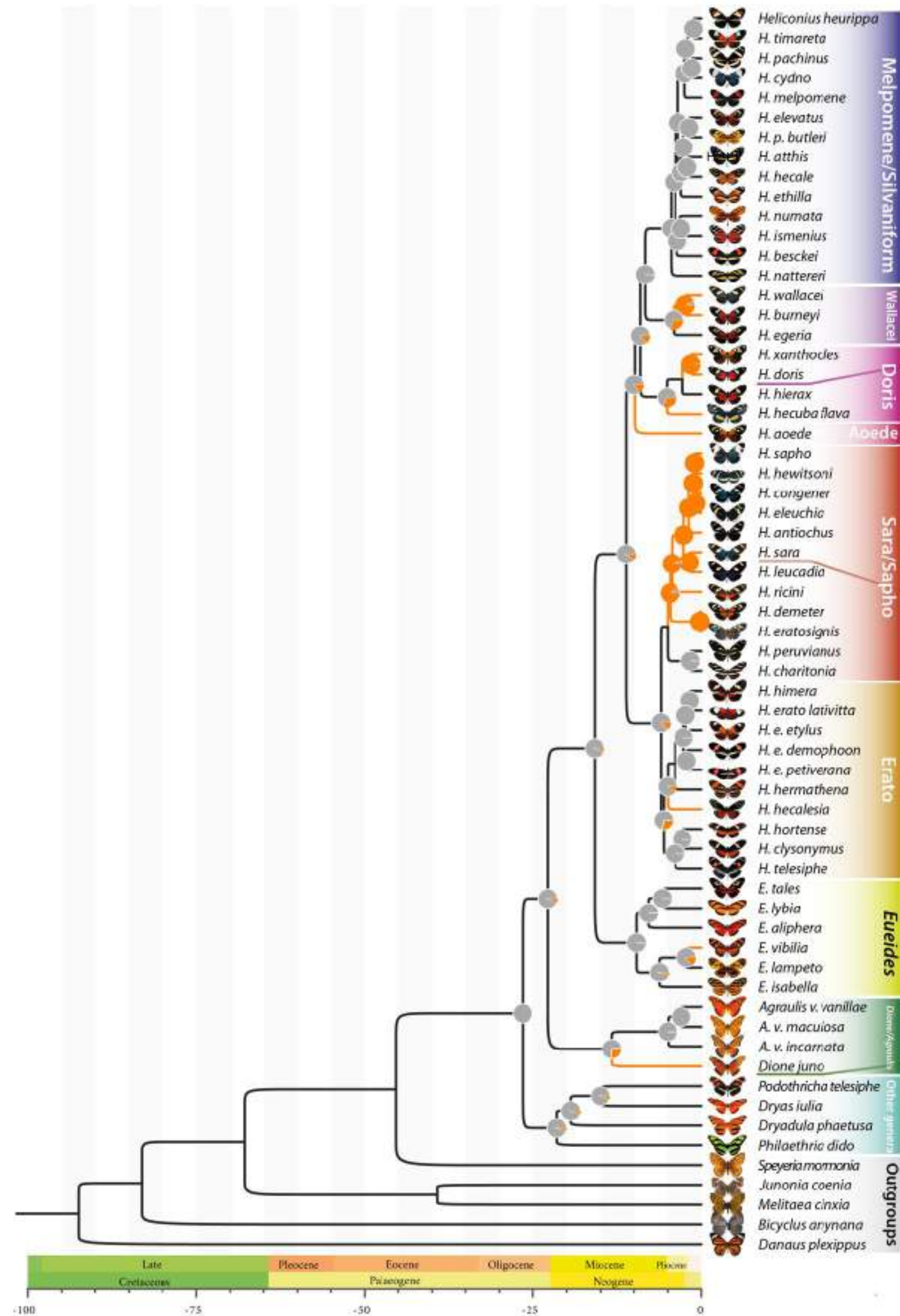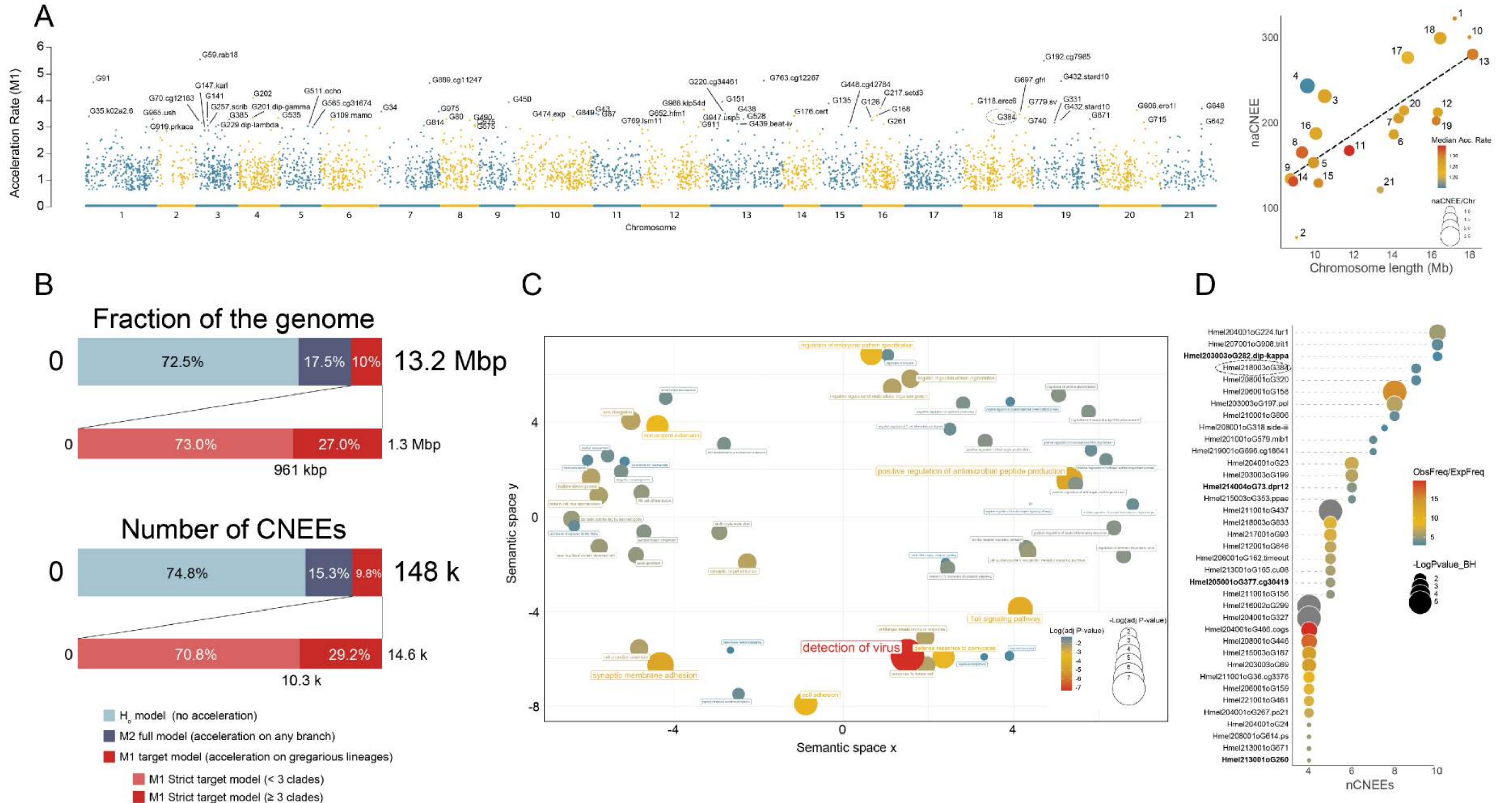
» *Mushroom body expansion (Brain)*

# » Test for Convergence »

# » Test for Convergence »