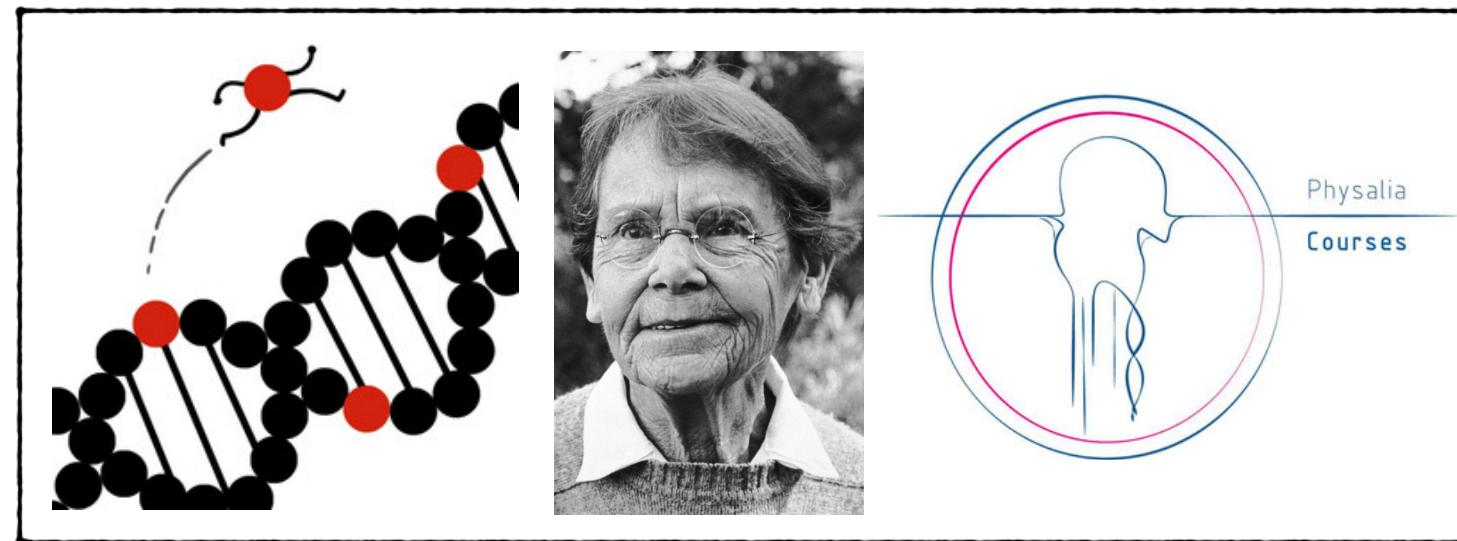


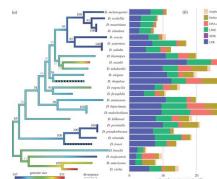
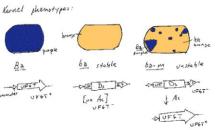
BIOINFORMATIC ANALYSIS OF TRANSPOSSABLE ELEMENTS

3rd-7th November 2025

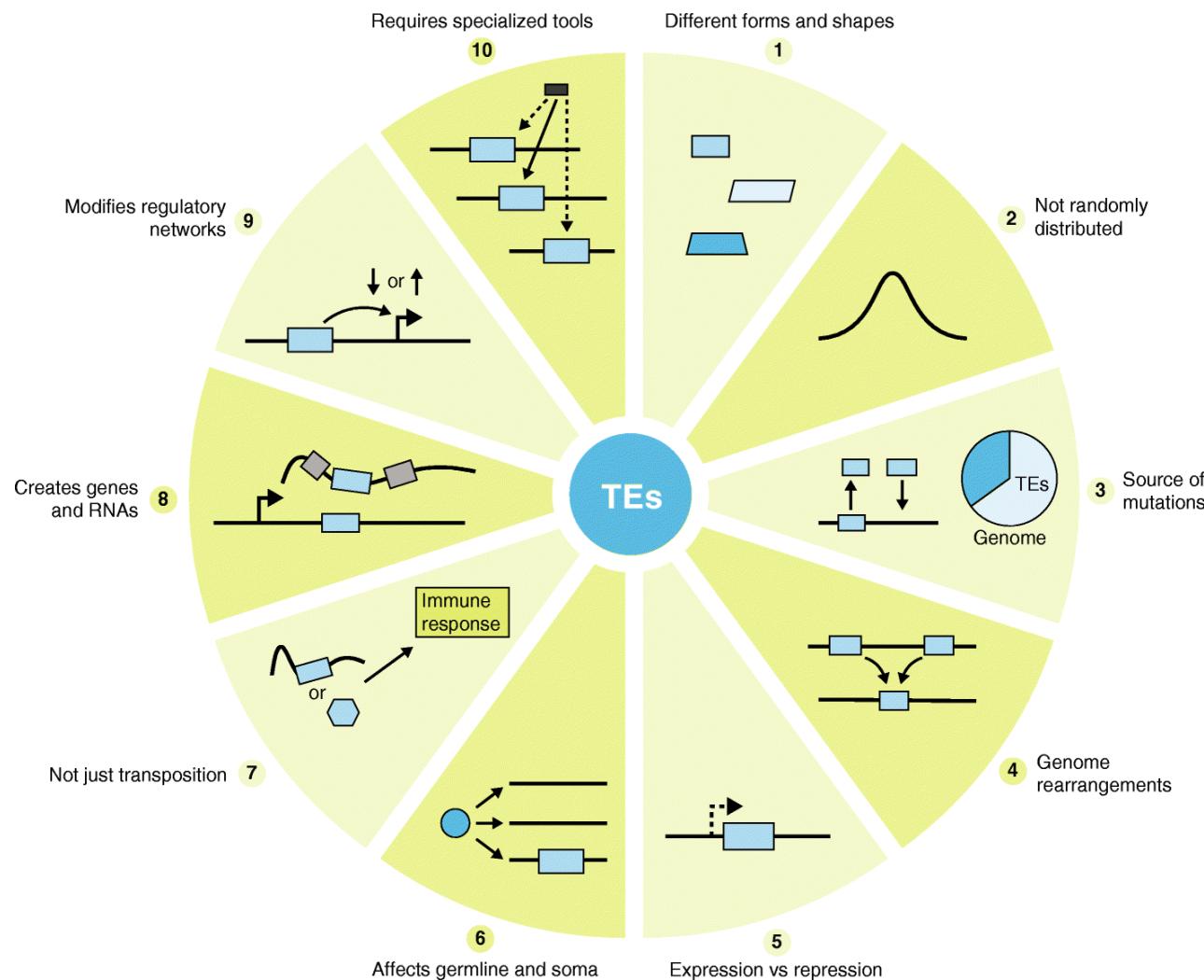
Lecture 4 TEs and genome evolution



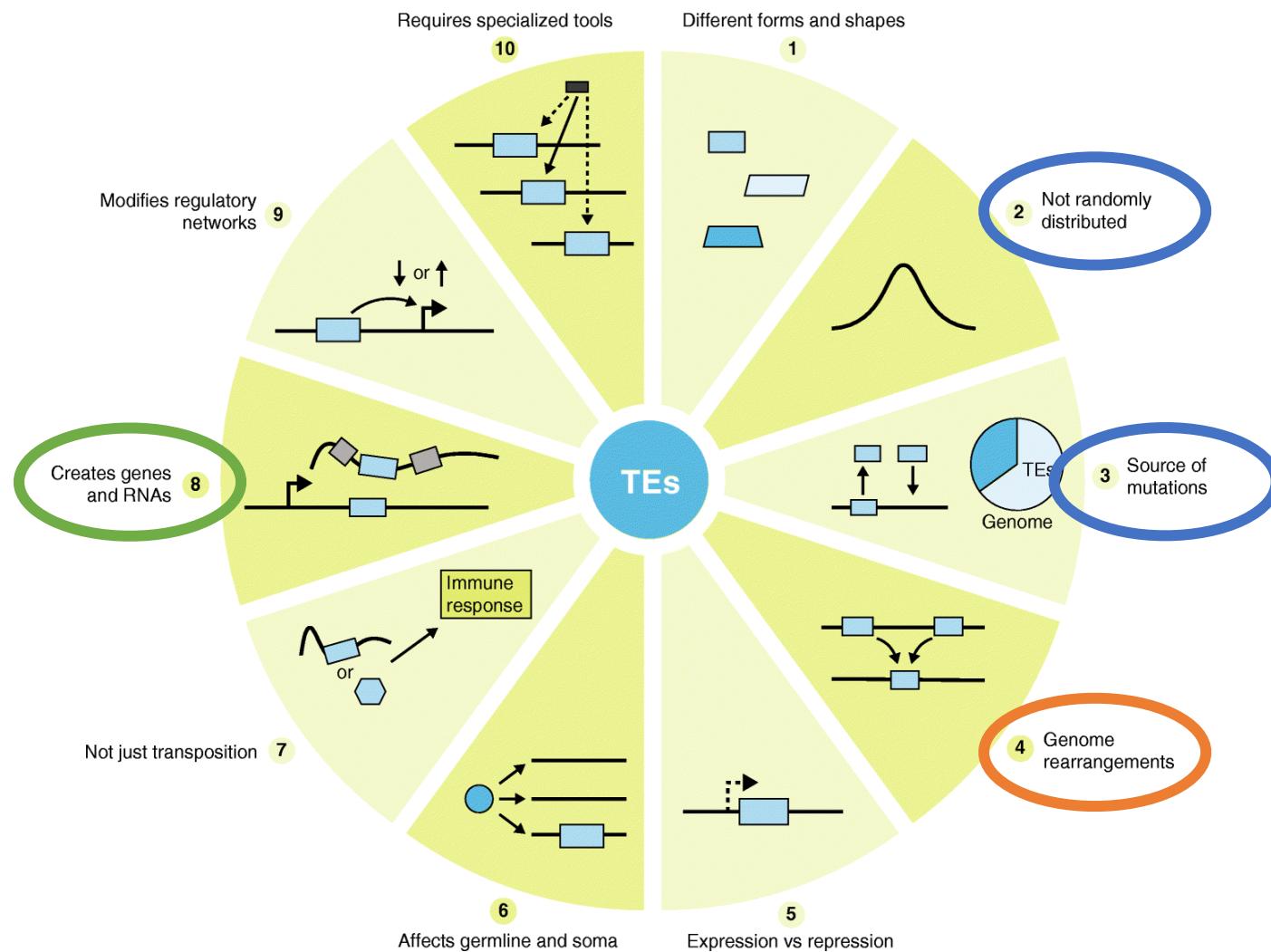
Valentina Peona



«Ten things you should know about TEs»



«Ten things you should know about TEs»



Drivers of genome evolution

1

Structure: from size to 3D organization

2

Ecology: waves of transposition and differential accumulation

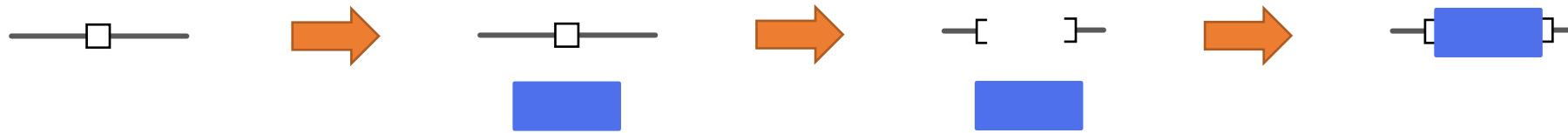
3

Gene evolution: TEs as innovators and constraint

Structural changes

Mechanisms

1. Non-homologous end joining DNA repair - NHEJ



2. Non-allelic homologous recombination - NAHR



3. Copy paste transposition

4. Cut and paste transposition

Structural changes

Mechanisms

1. Non allelic homologous recombination - NAHR

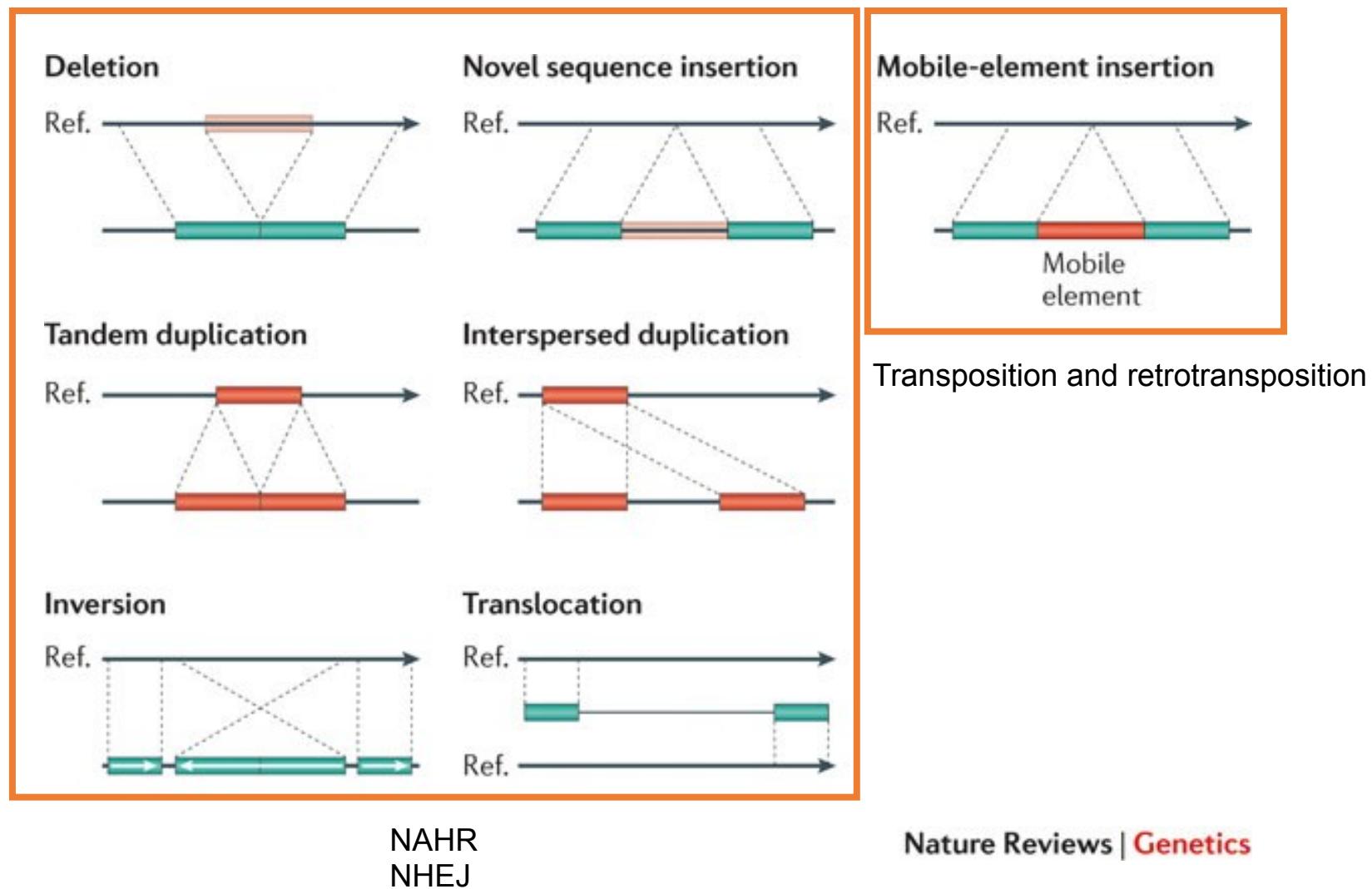


2. Copy paste transposition

3. Cut and paste transposition

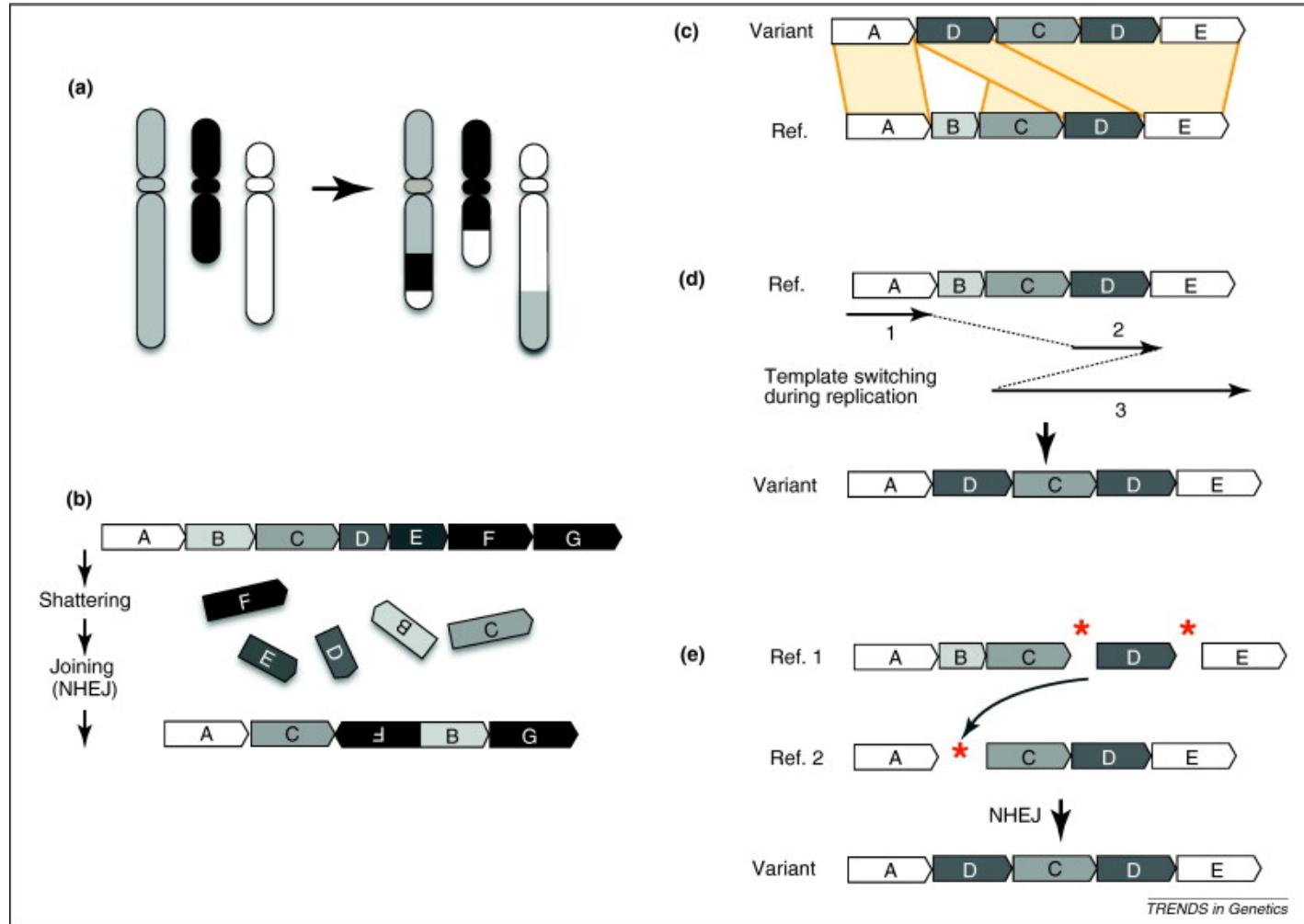
Genome rearrangements

Types



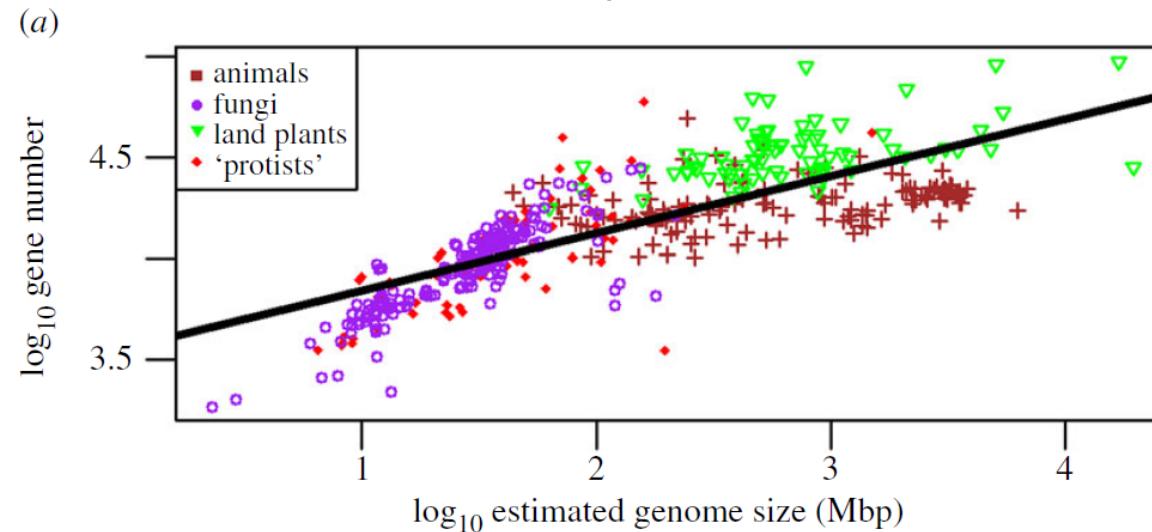
Nature Reviews | Genetics

Complex rearrangements



Genome size evolution

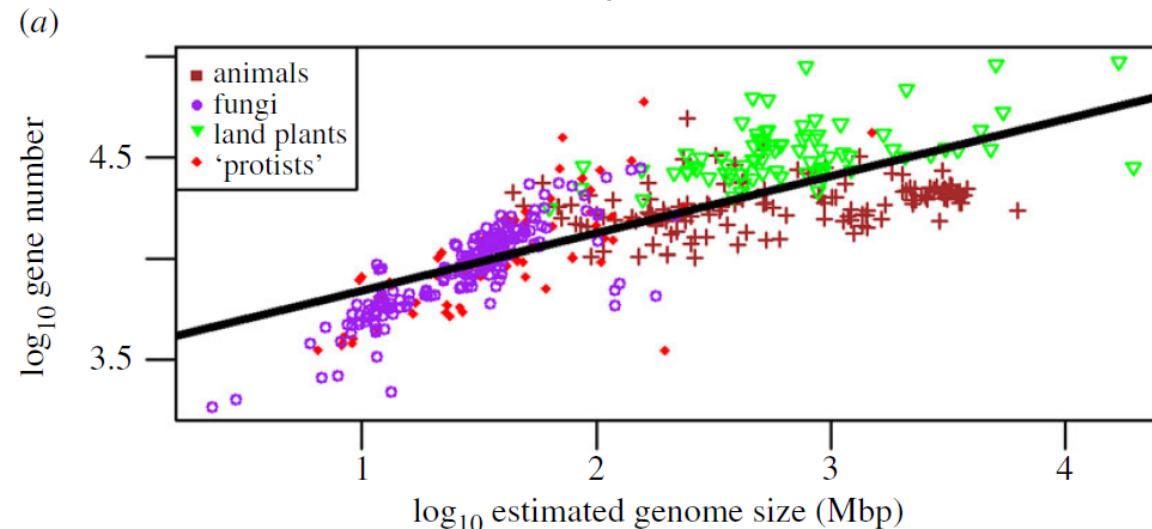
C-value paradox



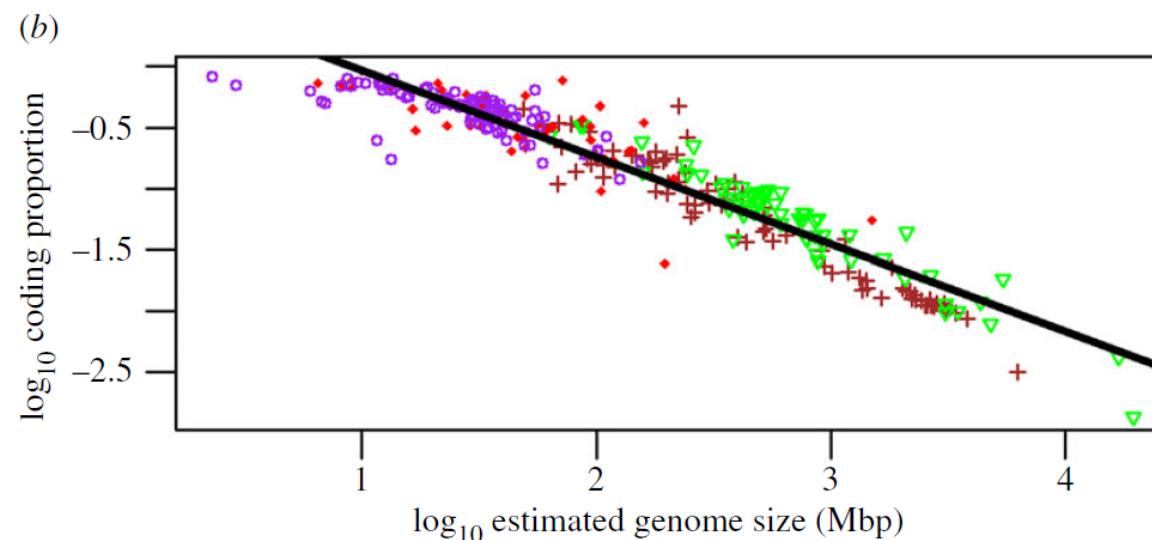
Positive relationship with number of genes

Genome size evolution

C-value paradox



Positive relationship with number of genes



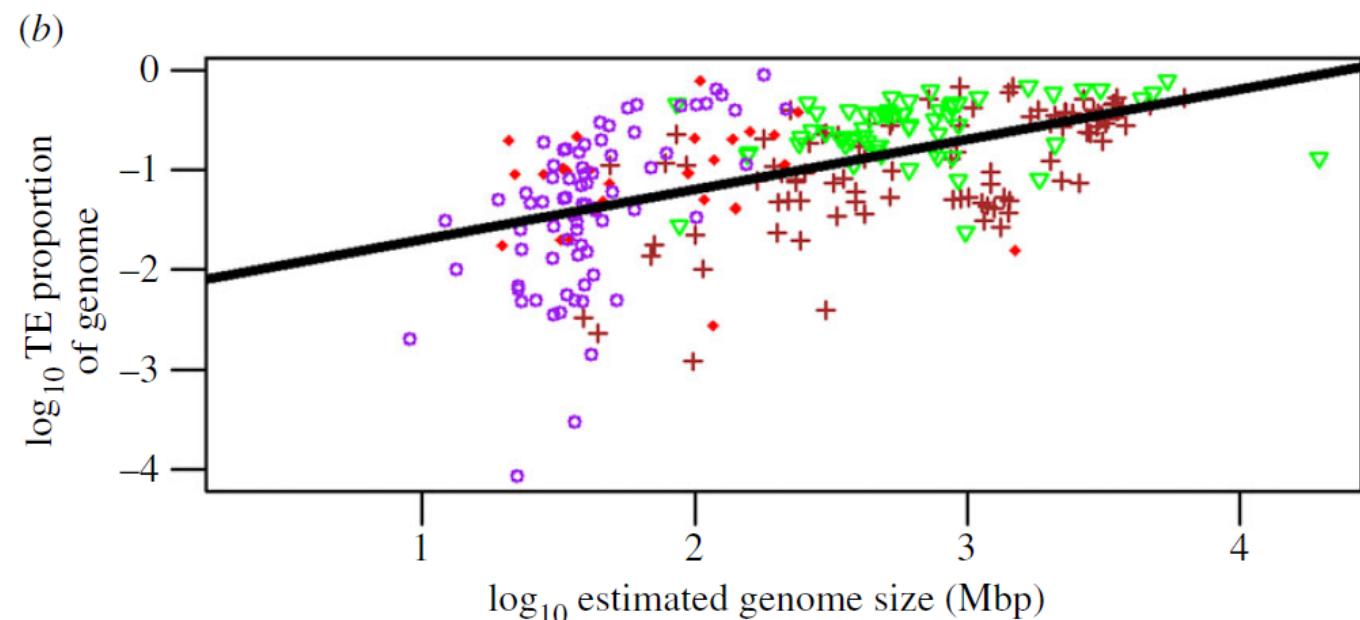
Negative relationship with proportion

Bigger genome = more non-functional (for the host) DNA

Genome size evolution

| Species | C-value (kb) |
|--------------------------------|--------------|
| <i>Navicula pelliculosa</i> | 35,000 |
| <i>Drosophila melanogaster</i> | 180,000 |
| <i>Gallus domesticus</i> | 1,200,000 |
| <i>Lampreta planrei</i> | 1,900,000 |
| <i>Xenopus laevis</i> | 3,100,000 |
| <i>Homo sapiens</i> | 3,400,000 |
| <i>Nicotiana tabaccum</i> | 3,800,000 |
| <i>Paramecium caudaatum</i> | 8,600,000 |
| <i>Allium cepa</i> | 18,000,000 |
| <i>Lilium formosanum</i> | 36,000,000 |
| <i>Amphiuma means</i> | 84,000,000 |
| <i>Protopterus aethiopicus</i> | 140,000,000 |
| <i>Amboeba proteus</i> | 290,000,000 |
| <i>Amoeba dubia</i> | 670,000,000 |

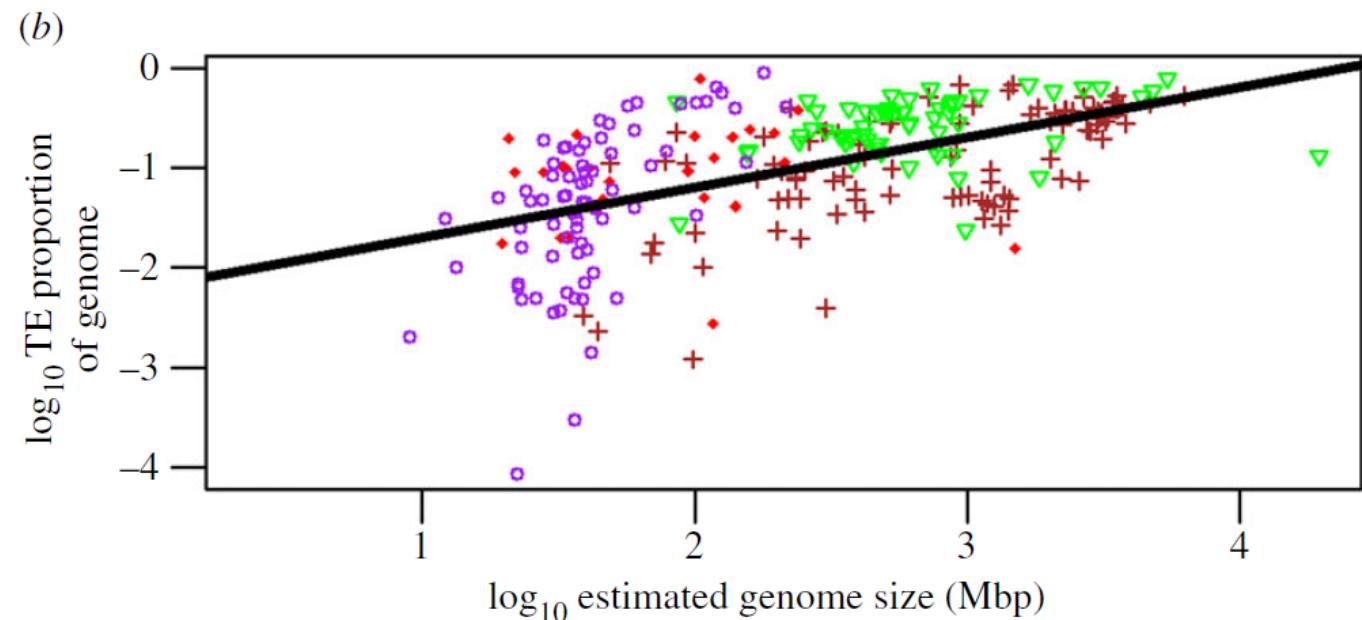
C-value paradox



Genome size evolution

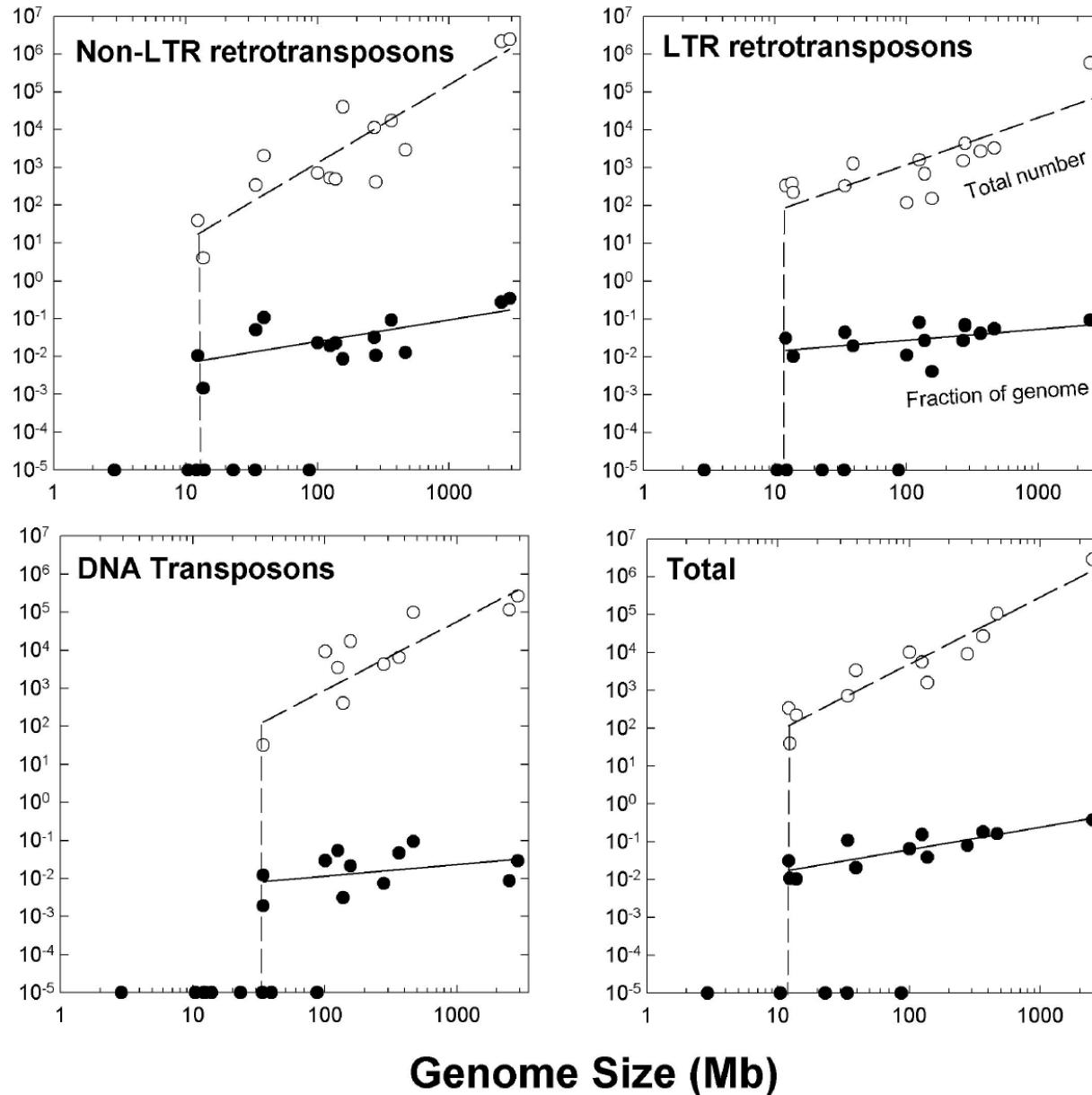
| Species | C-value (kb) |
|--------------------------------|--------------|
| <i>Navicula pelliculosa</i> | 35,000 |
| <i>Drosophila melanogaster</i> | 180,000 |
| <i>Gallus domesticus</i> | 1,200,000 |
| <i>Lampreta planrei</i> | 1,900,000 |
| <i>Xenopus laevis</i> | 3,100,000 |
| <i>Homo sapiens</i> | 3,400,000 |
| <i>Nicotiana tabaccum</i> | 3,800,000 |
| <i>Paramecium caudaatum</i> | 8,600,000 |
| <i>Allium cepa</i> | 18,000,000 |
| <i>Lilium formosanum</i> | 36,000,000 |
| <i>Amphiuma means</i> | 84,000,000 |
| <i>Protopterus aethiopicus</i> | 140,000,000 |
| <i>Amboeba proteus</i> | 290,000,000 |
| <i>Amoeba dubia</i> | 670,000,000 |

C-value paradox



- Just because we find a lot of TEs does not mean that the genome «needs» them, maybe genomes are just stuck with TEs and can deal with them
- Ultimate vs proximate cause for TE presence
- Important distinction for when we design hypotheses

Nothing in evolution makes sense...



Why do we have more/less TEs?

"Consistent with theoretical expectations, all three classes of mobile elements appear to have a threshold genome size below which mobile elements are unable to establish, an intermediate range in which only a fraction of species harbor them, and an upper threshold (~100 Mb) above which all species are infected (1) (Fig. 4). By extrapolation from the mutation rate cited above, the critical effective population size above which a unicellular eukaryote population appears to be immune to retrotransposon proliferation is $\sim 7 \times 10^7$, whereas for DNA-based transposons it is $\sim 2 \times 10^7$."

- Need space for jumping and for accumulating
- Genomes that are too small may not tolerate well the presence of TEs
- Effective population size, recombination rate, ...

... except in the light of population genetics

The frailty of adaptive hypotheses for the origins of organismal complexity

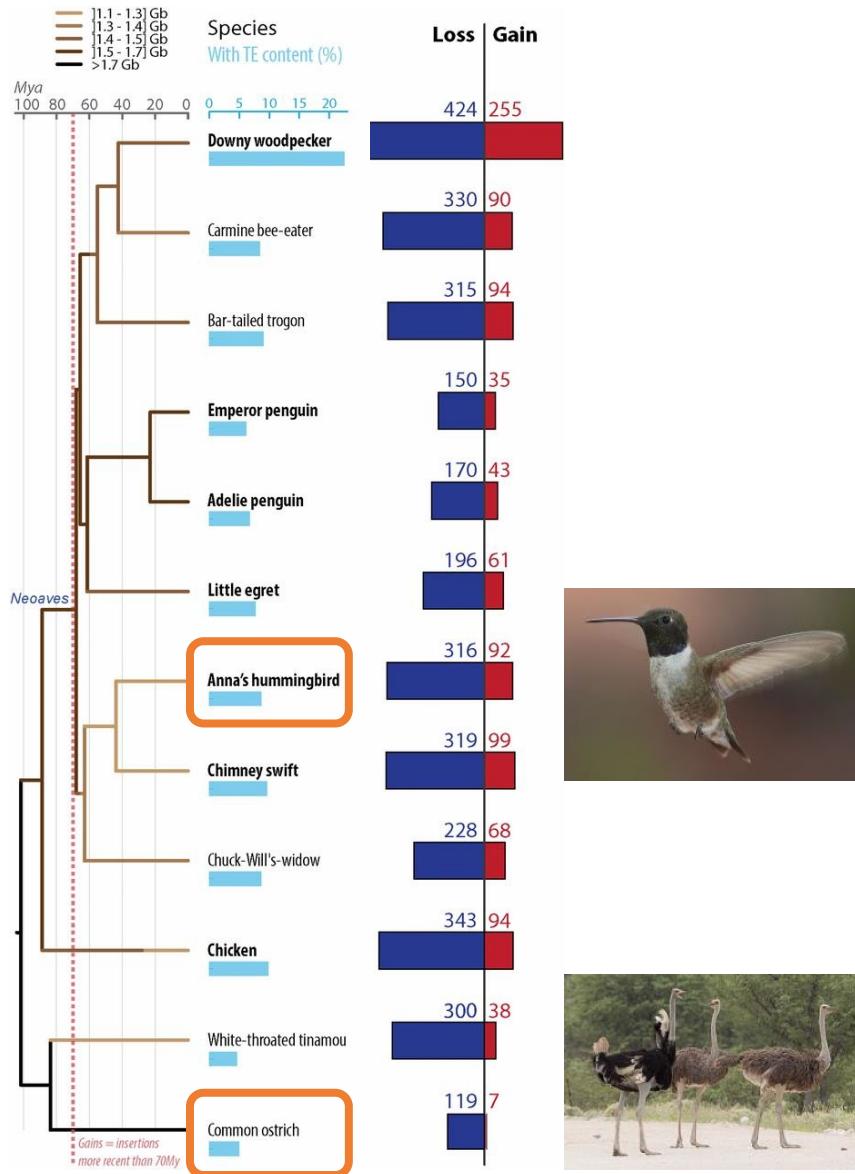
Michael Lynch*

Department of Biology, Indiana University, Bloomington, IN 47405

The vast majority of biologists engaged in evolutionary studies interpret virtually every aspect of biodiversity in adaptive terms. This narrow view of evolution has become untenable in light of recent observations from genomic sequencing and population-genetic theory. Numerous aspects of genomic architecture, gene structure, and developmental pathways are difficult to explain without invoking the nonadaptive forces of genetic drift and mutation. In addition, emergent biological features such as complexity, modularity, and evolvability, all of which are current targets of considerable speculation, may be nothing more than indirect by-products of processes operating at lower levels of organization. These issues are examined in the context of the view that the origins of many aspects of biological diversity, from gene-structural embellishments to novelties at the phenotypic level, have roots in nonadaptive processes, with the population-genetic environment imposing strong directionality on the paths that are open to evolutionary exploitation.

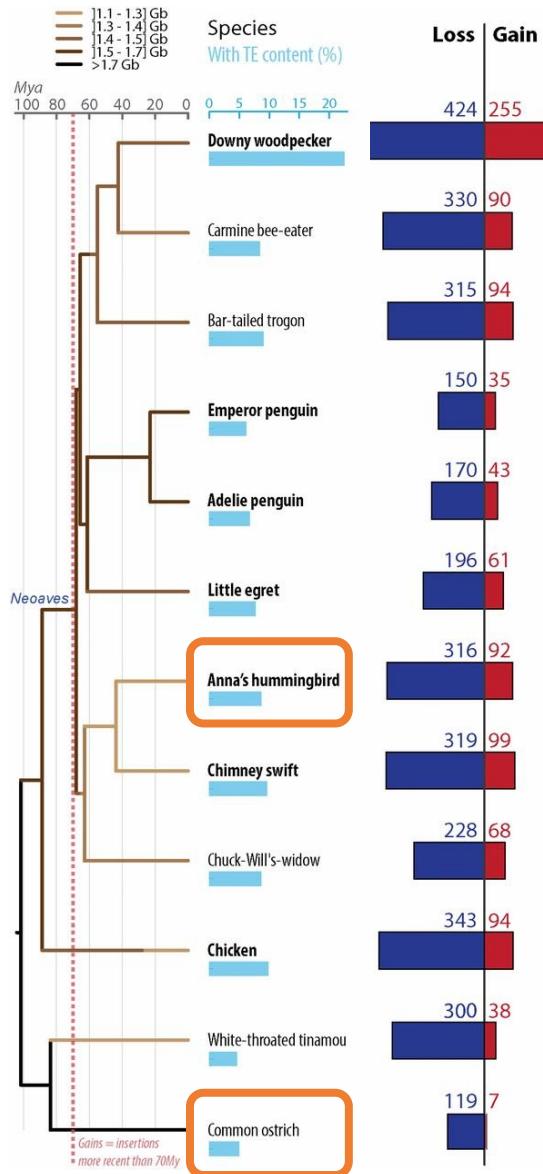
Genome size evolution

Accordion model



Genome size evolution

Accordion model

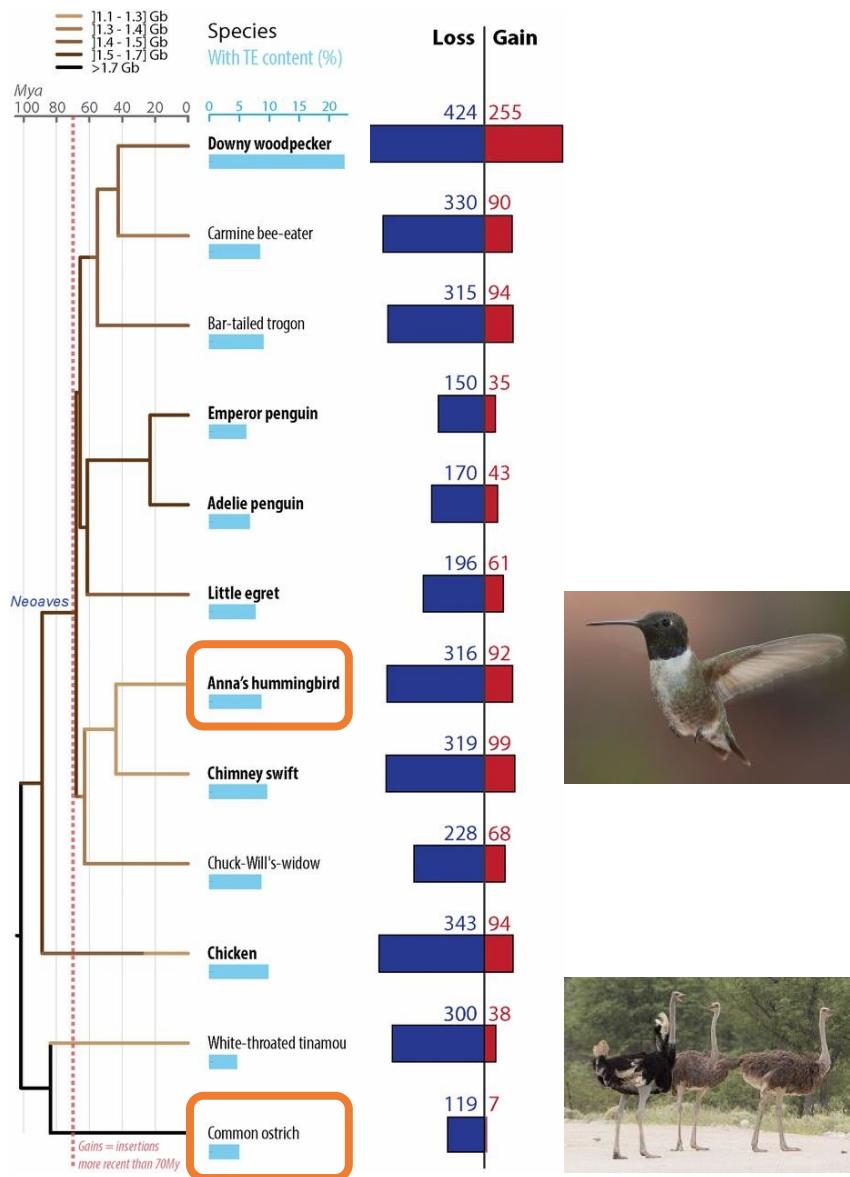


No constraint



Genome size evolution

Accordion model

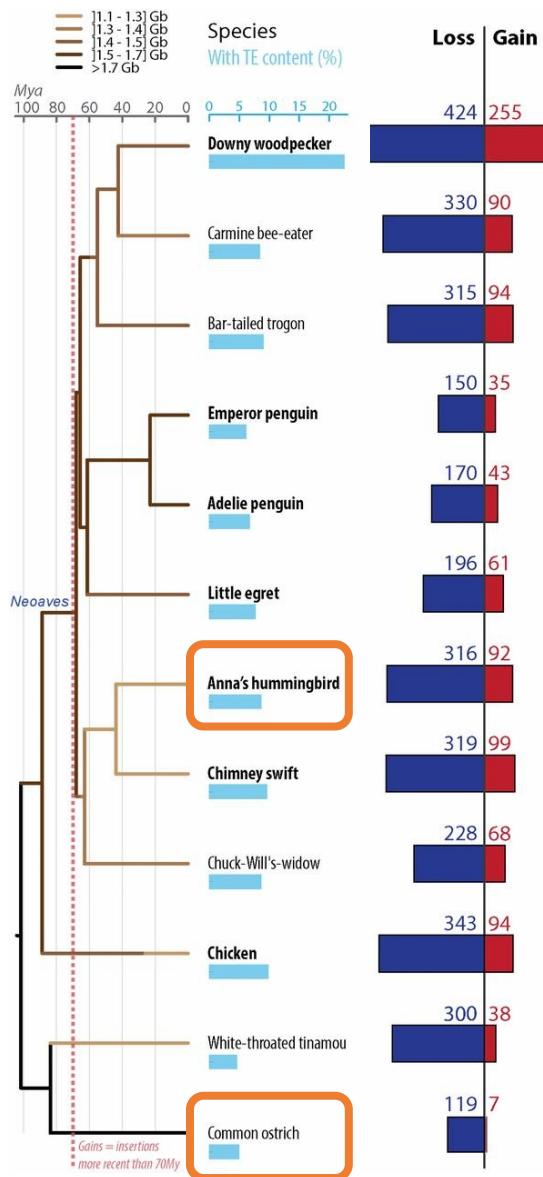


Consider not only host popgen,
but also TE popgen!

No constraint

Genome size evolution

Accordion model



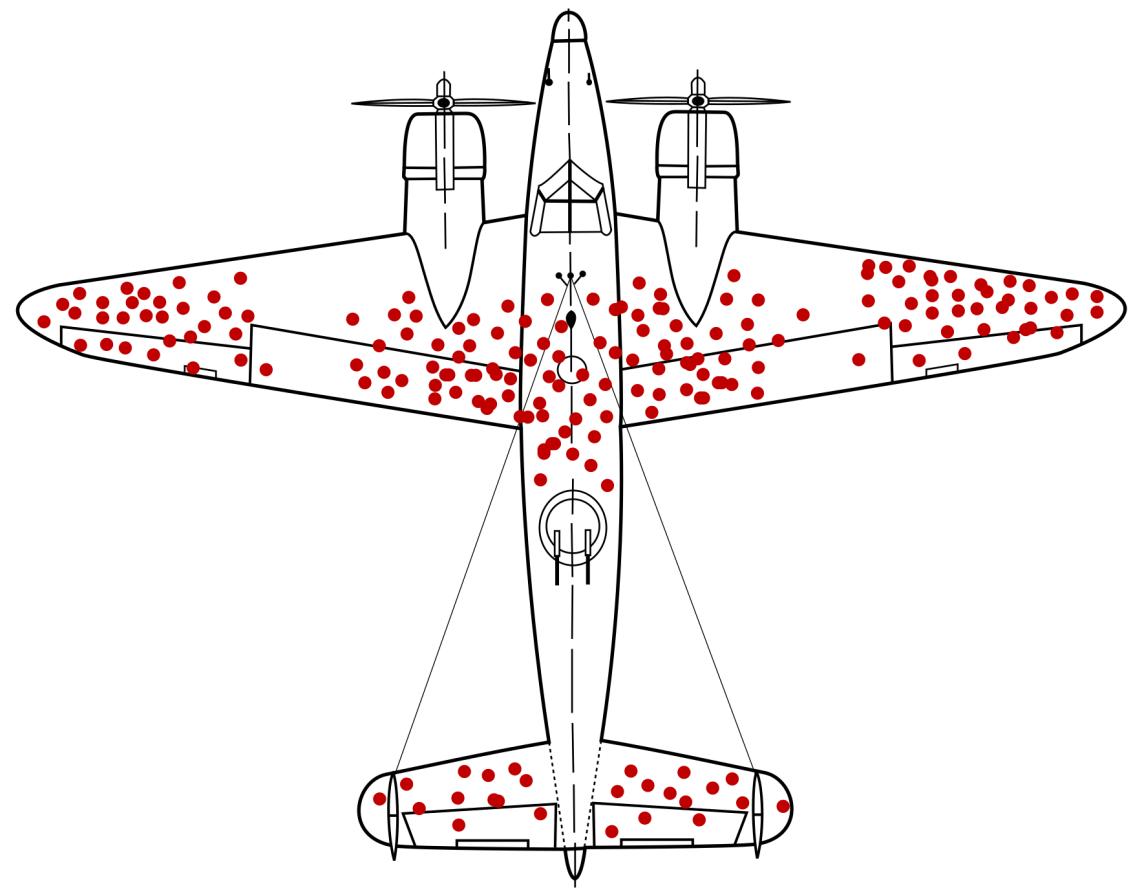
No constraint



Consider not only host popgen,
but also TE popgen!

Need to clean up their genomes
VS
genome structure reflects popgen of the host
and TEs

Survivorship bias



pattern != process

Survivorship bias

- More TEs because of more insertions or more fixations?

TE popgen vs host popgen

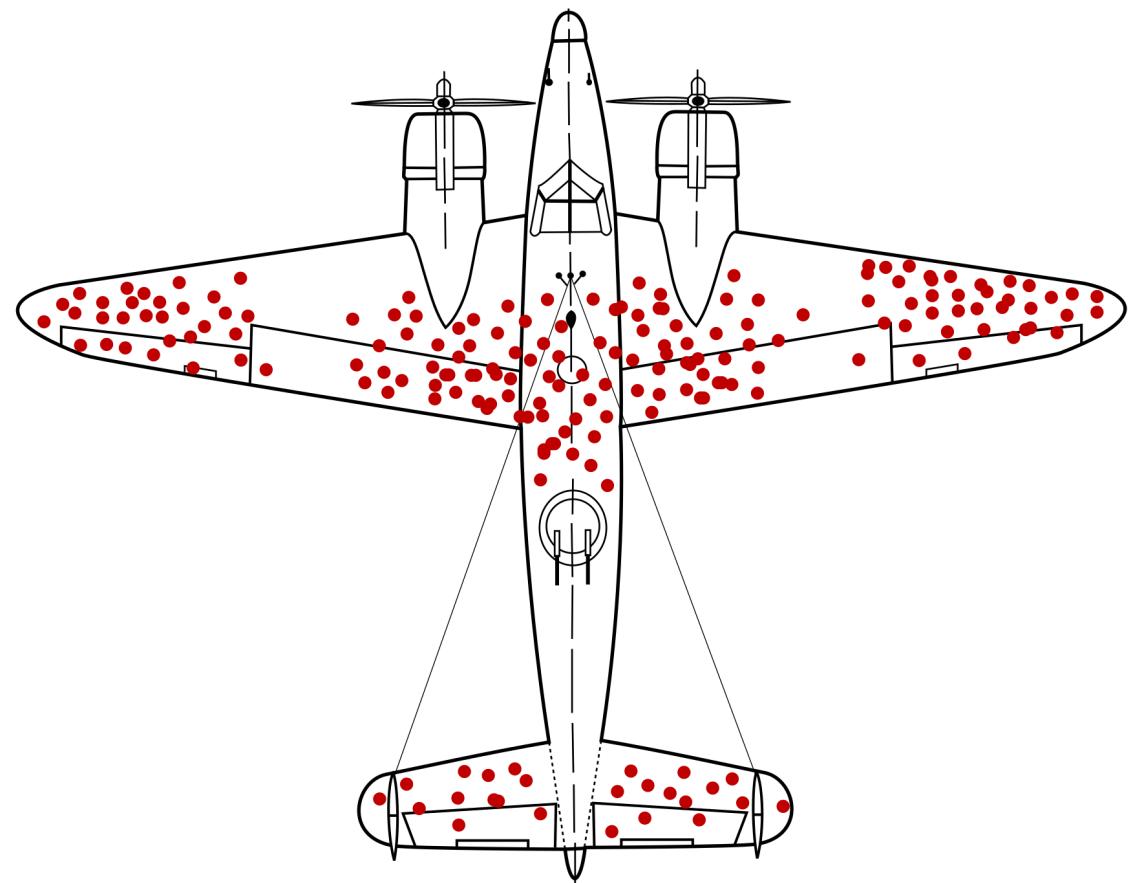
- Fewer TEs because of preferential insertion or negative selection?

you cannot see detrimental insertions

- Larger genome because of more TE accumulation or fewer deletions?

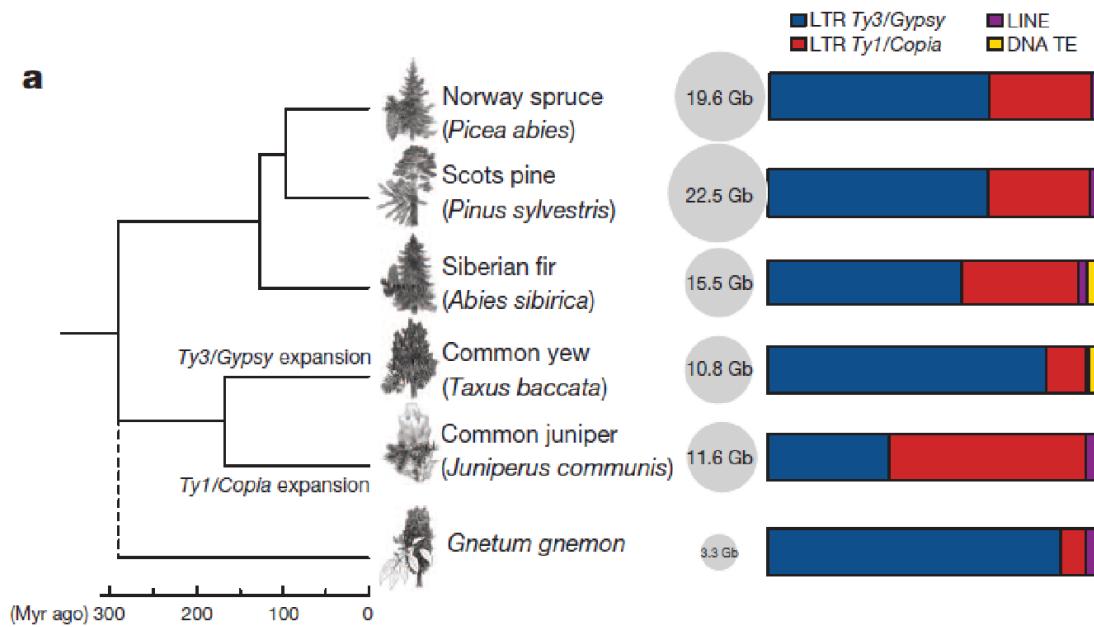
lower deletion rate

- Smaller genome because of lower TE accumulation or more deletions?

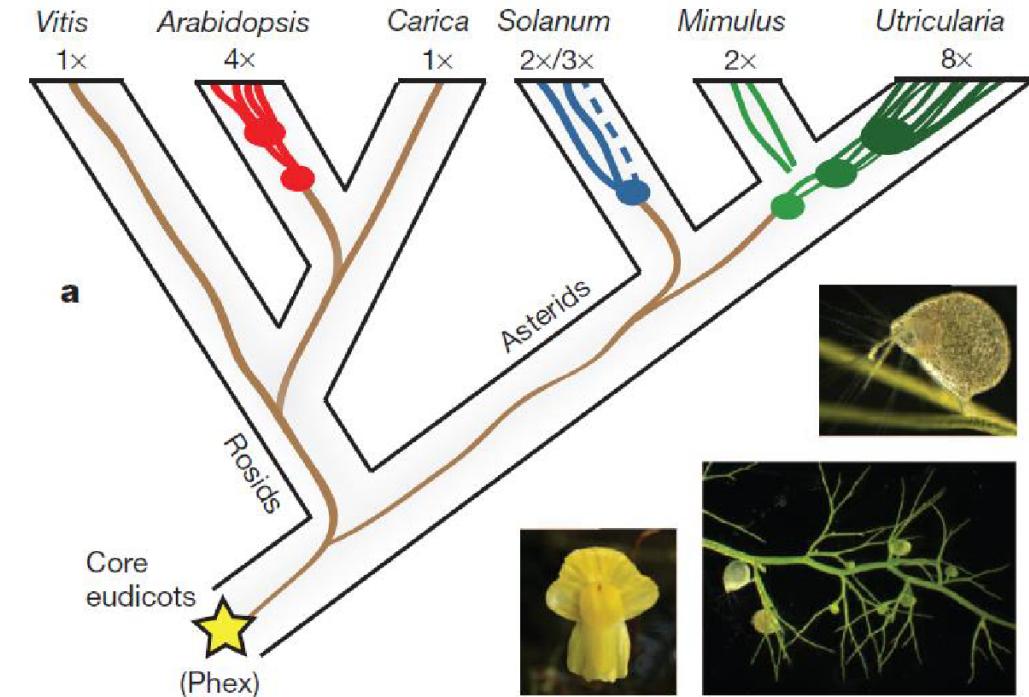


pattern != process

Genomic gigantism and dwarfism

a

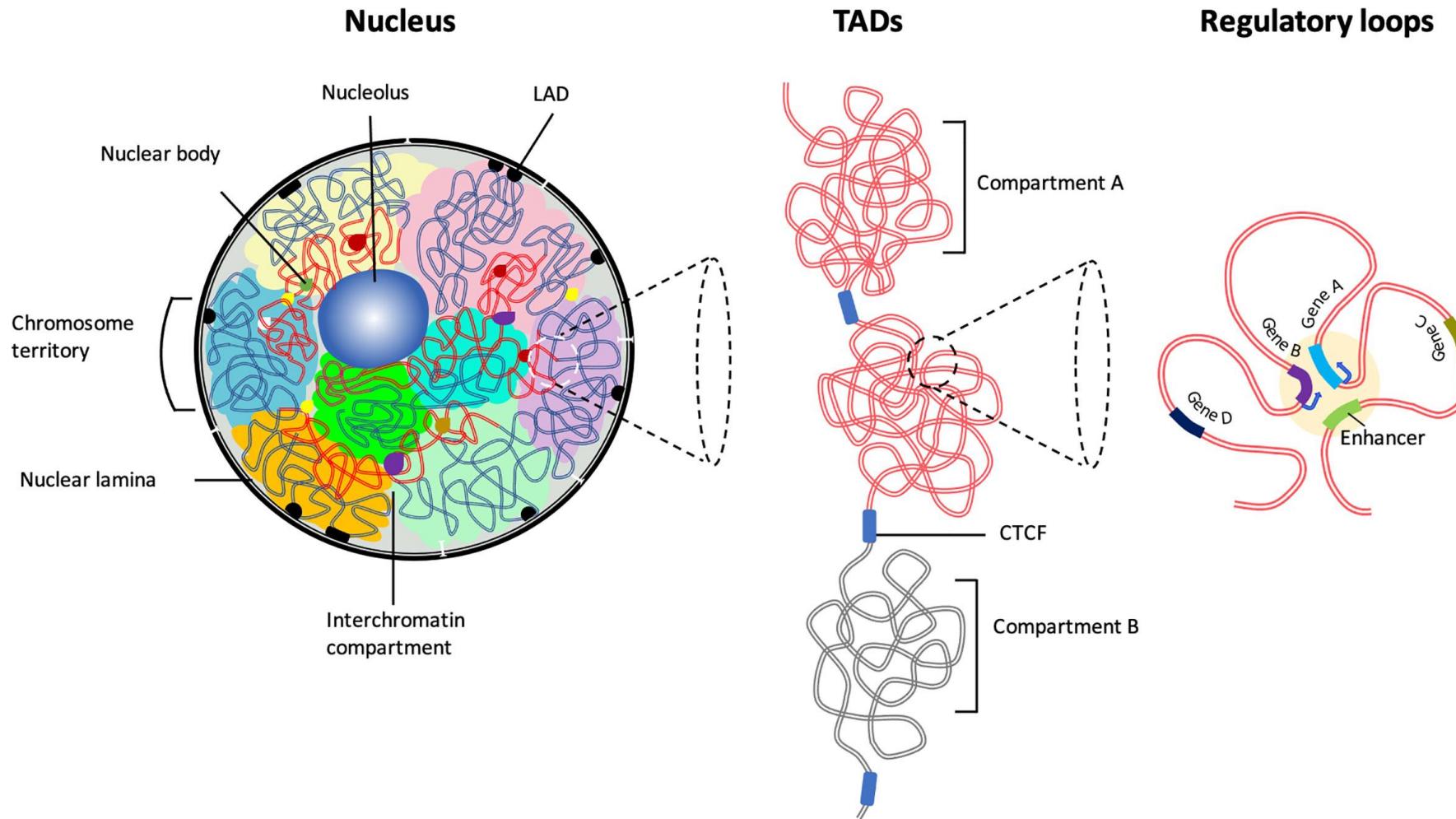
Tons of TE activity/accumulation or lack of deletions?



77 Mb genome despite 3 rounds of whole-genome duplications, only 3% repetitive DNA!

Lack of TE activity/accumulation or tons of deletions?

3D genome structure

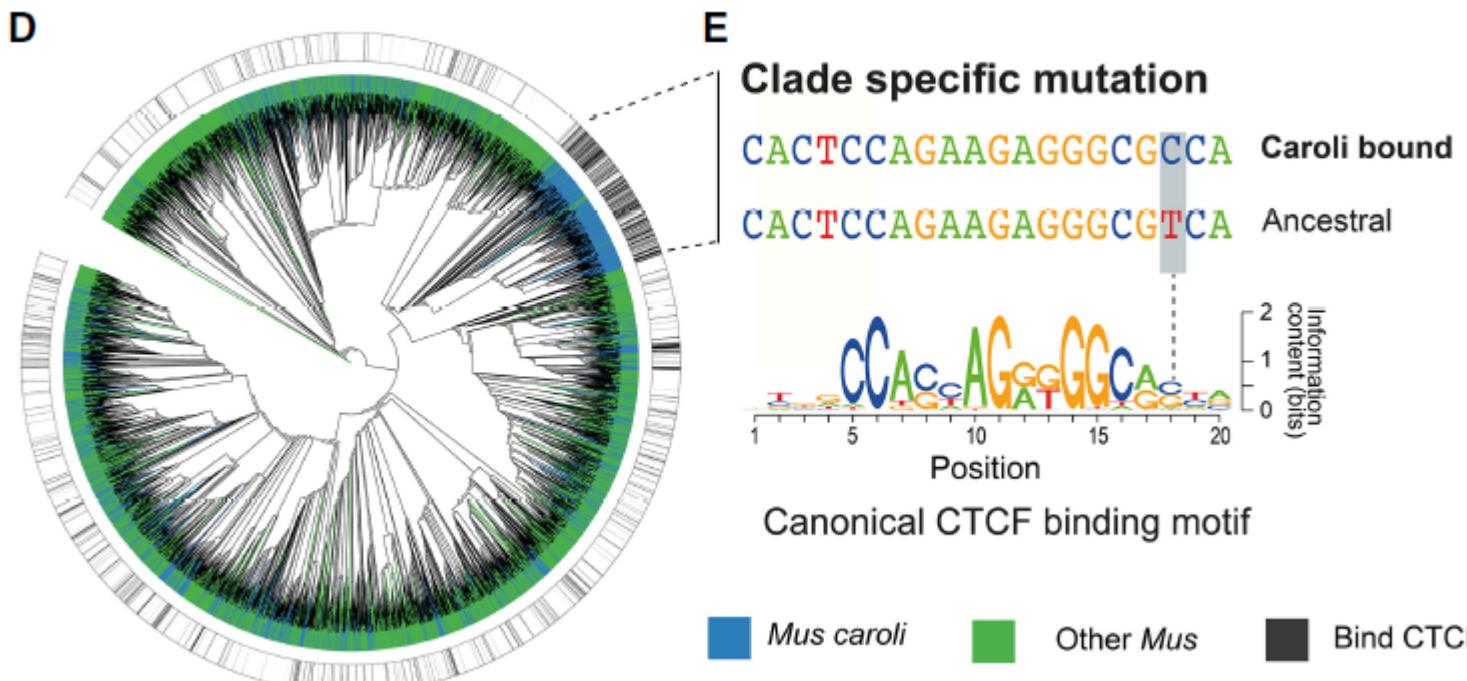


3D genome structure

Research

Repeat associated mechanisms of genome evolution and function revealed by the *Mus caroli* and *Mus pahari* genomes

David Thybert,^{1,2} Maša Roller,¹ Fábio C.P. Navarro,³ Ian Fiddes,⁴ Ian Streeter,¹ Christine Feig,⁵ David Martin-Galvez,¹ Mikhail Kolmogorov,⁶ Václav Janoušek,⁷



One SNP in the *Mus caroli* lineage turned SINE B2 into CTCF binding sites

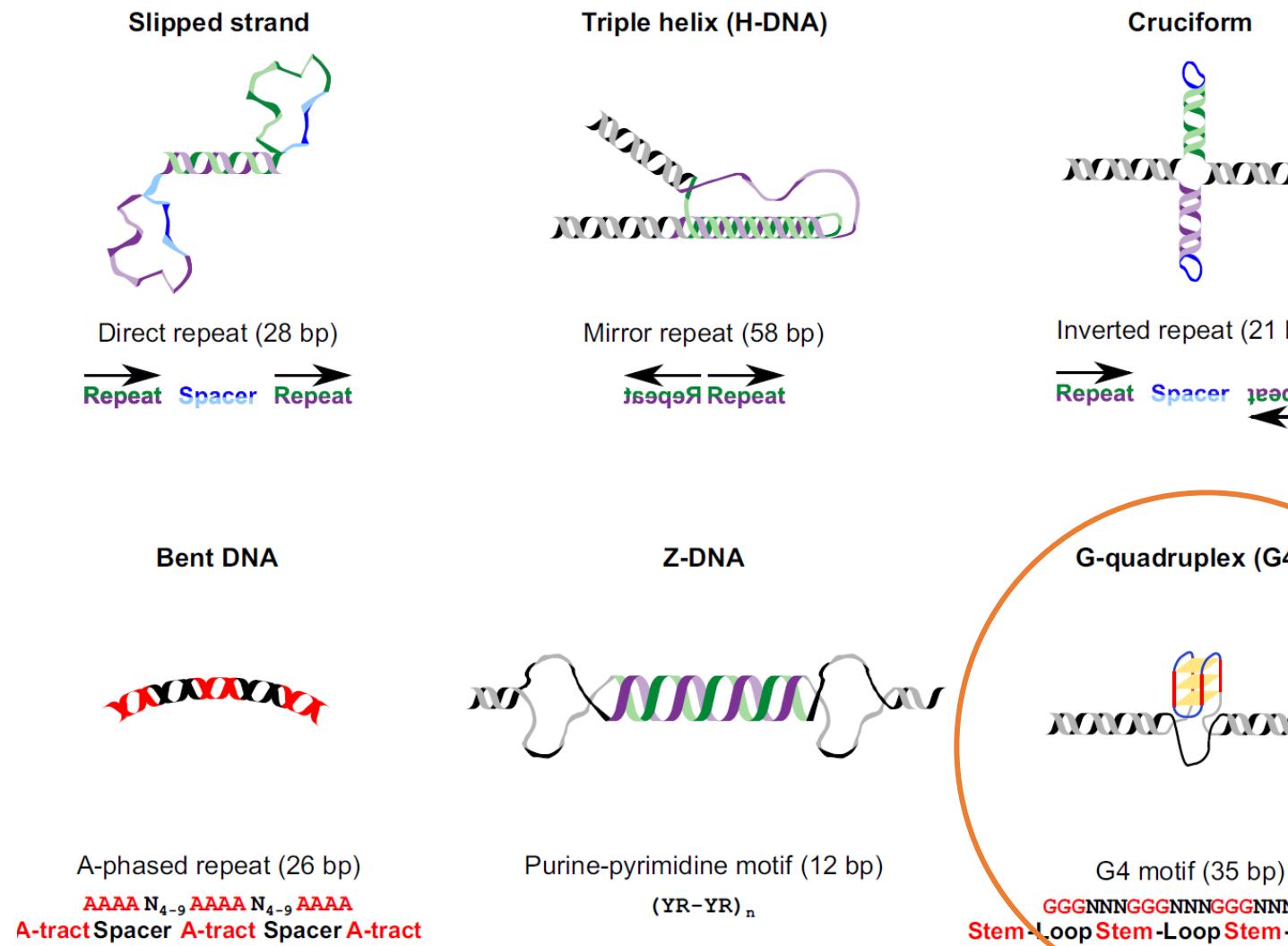
3D genome structure

MIR retrotransposon sequences provide insulators to the human genome

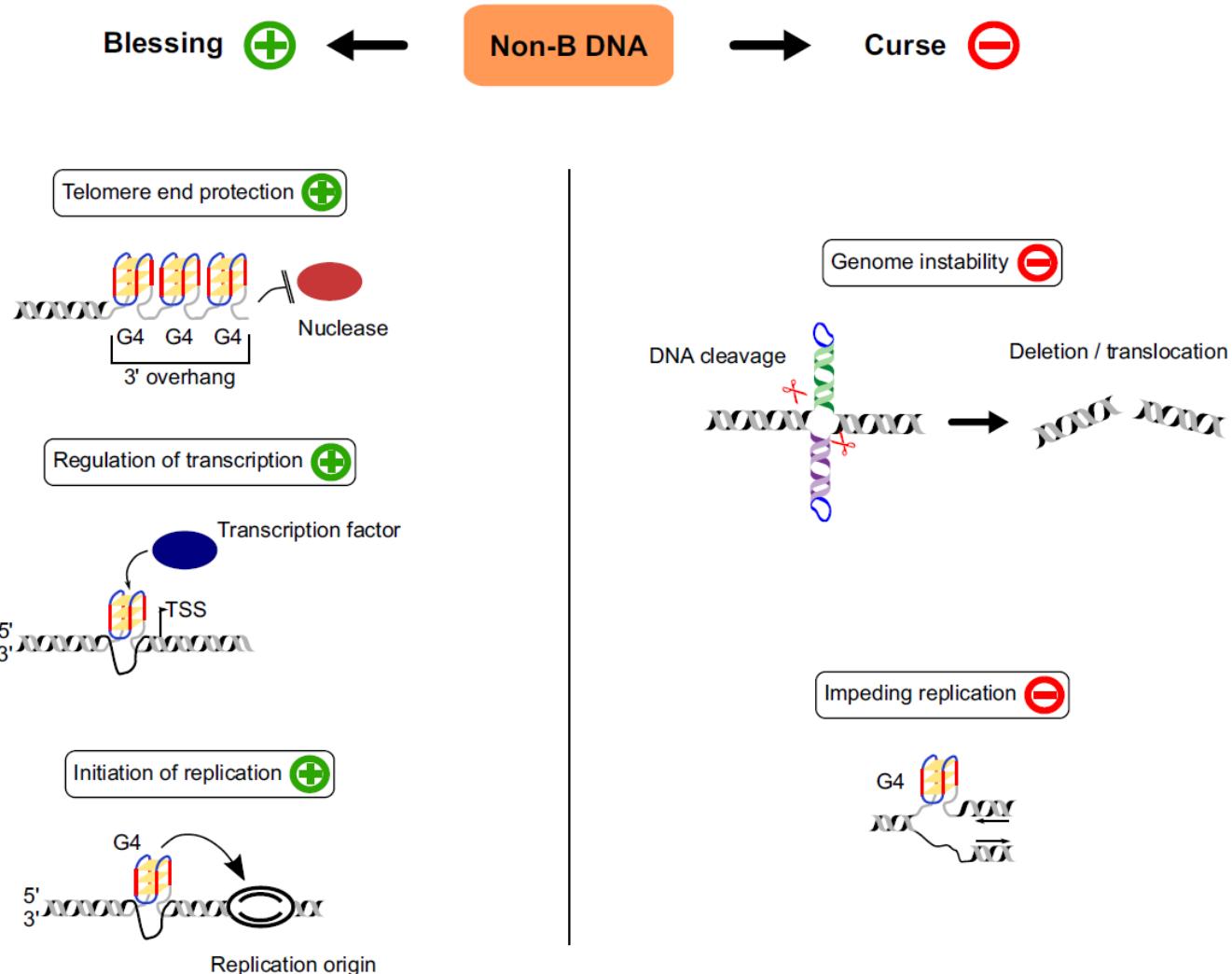
Jianrong Wang^a, Cristina Vicente-García^{b,c}, Davide Seruggia^{b,c}, Eduardo Moltó^{b,c}, Ana Fernandez-Miñán^d, Ana Neto^d, Elbert Lee^e, José Luis Gómez-Skarmeta^d, Lluís Montoliu^{b,c}, Victoria V. Lunyak^e, and I. King Jordan^{a,f,1}

MIR insulators appear to be CCCTC-binding factor (CTCF) independent and show a distinct local chromatin environment with marked peaks for RNA Pol III and a number of histone modifications, suggesting that MIR insulators recruit transcriptional complexes and chromatin modifying enzymes *in situ* to help establish chromatin and regulatory domains in the human genome. **The provisioning of insulators by MIRs across the human genome suggests a specific mechanism by which TE sequences can be used to modulate gene regulatory networks.**

Non-B DNA and TEs



Non-B DNA and TEs

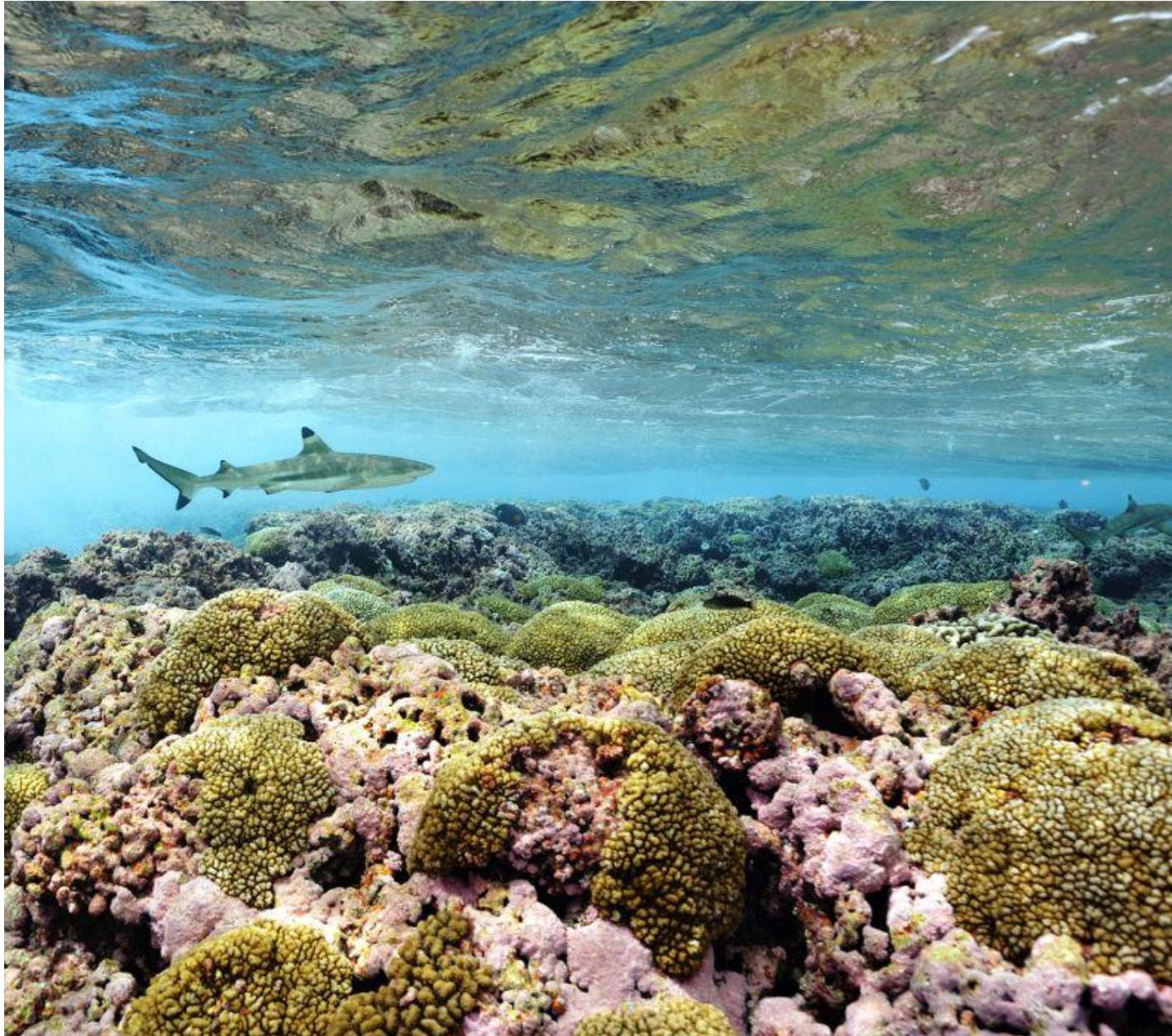


TEs were proposed to serve as vehicles of spreading non-B DNA across the genome

Non-B DNA located in TEs may play a role in:

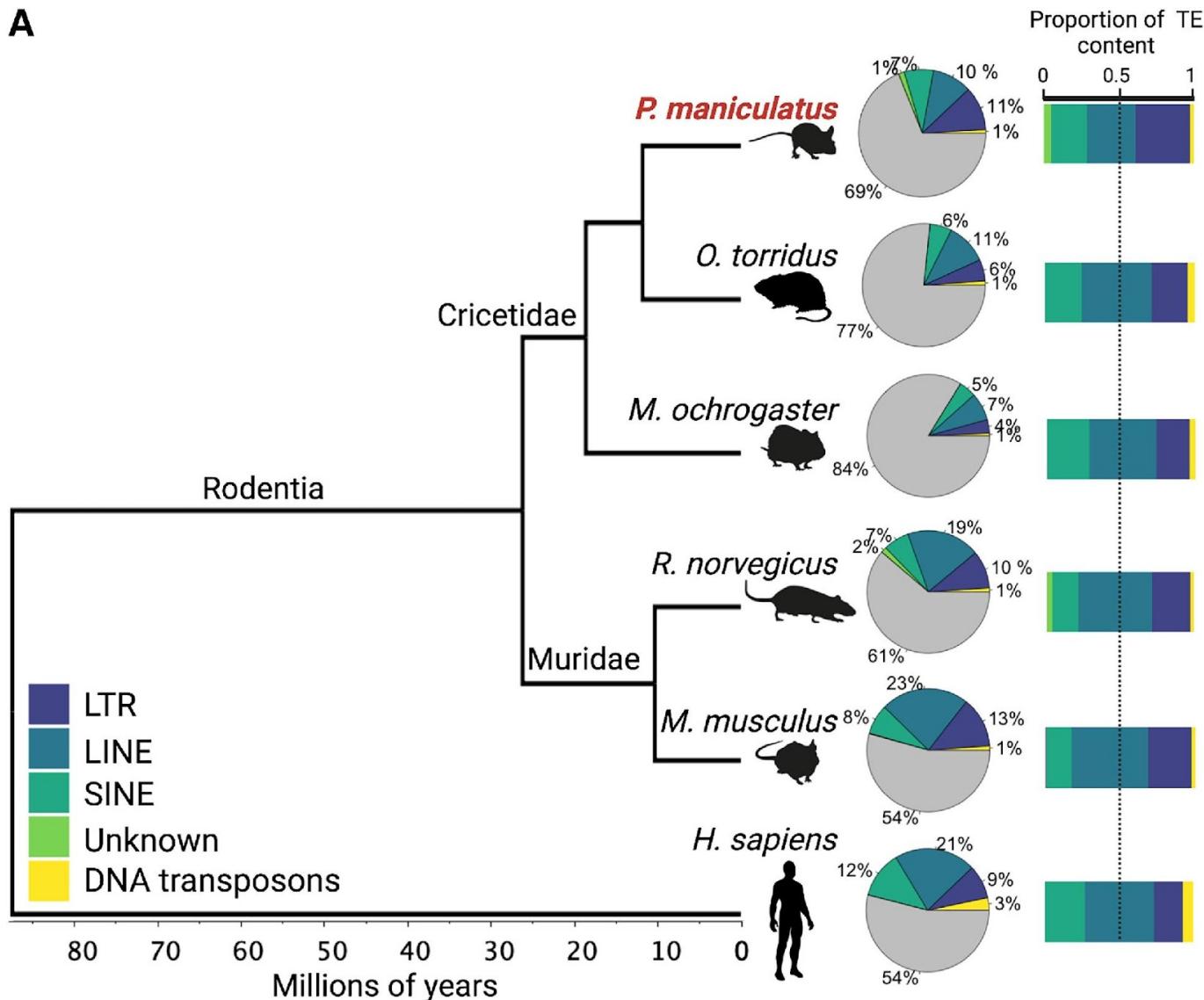
- regulating activity of neighboring genes by affecting the activity of promoters
- regulating epigenetic state of discrete genomic regions and facilitating genome silencing, including TE silencing, via spreading heterochromatin
- reshuffling of genomic DNA via recombination at TEs
- G4 motifs within hominid-specific SVA retrotransposons are enriched in cancer genome breakpoints

Genomes as ecosystems

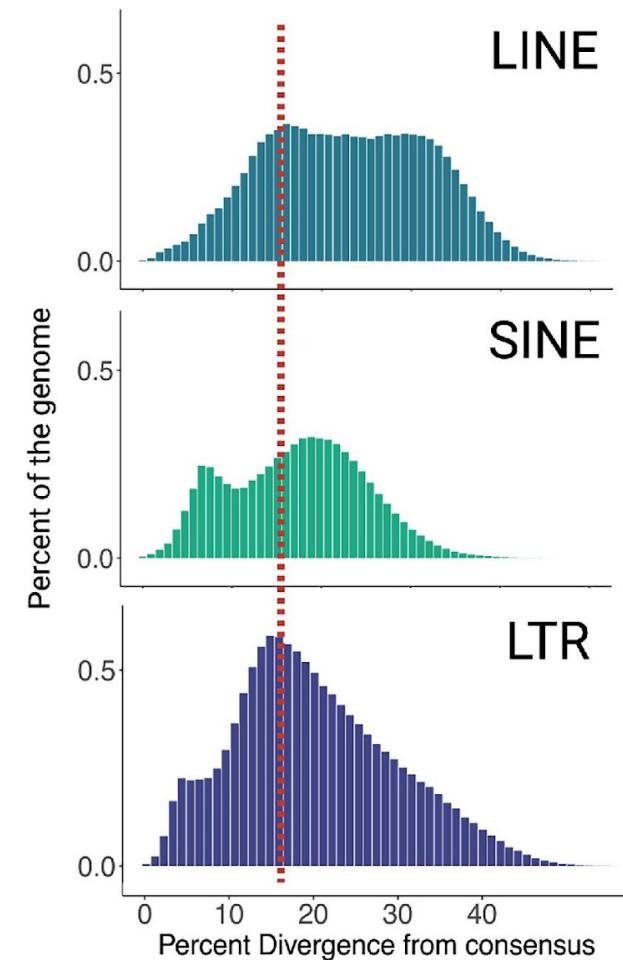


Genomes as ecosystems

A

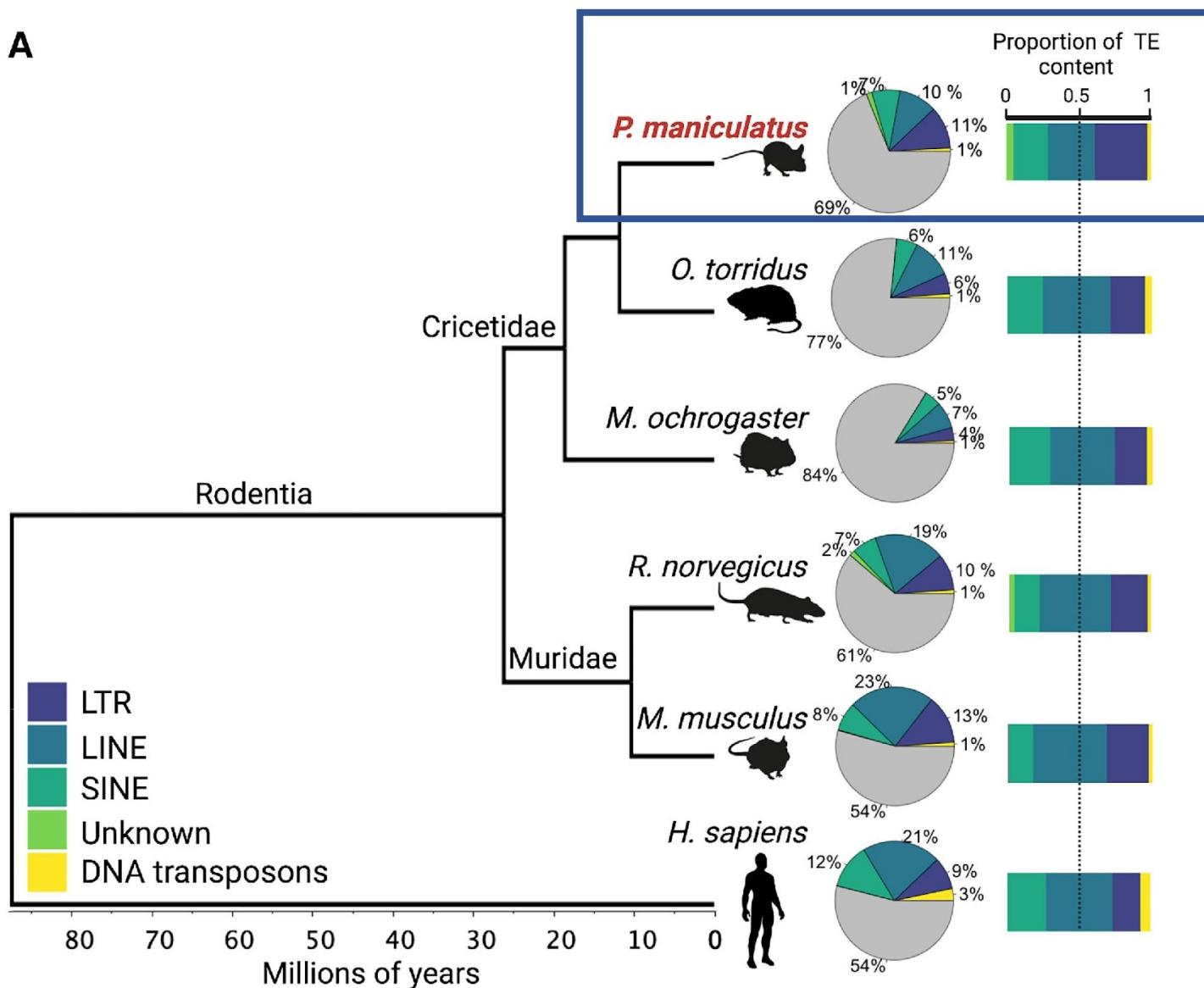


B

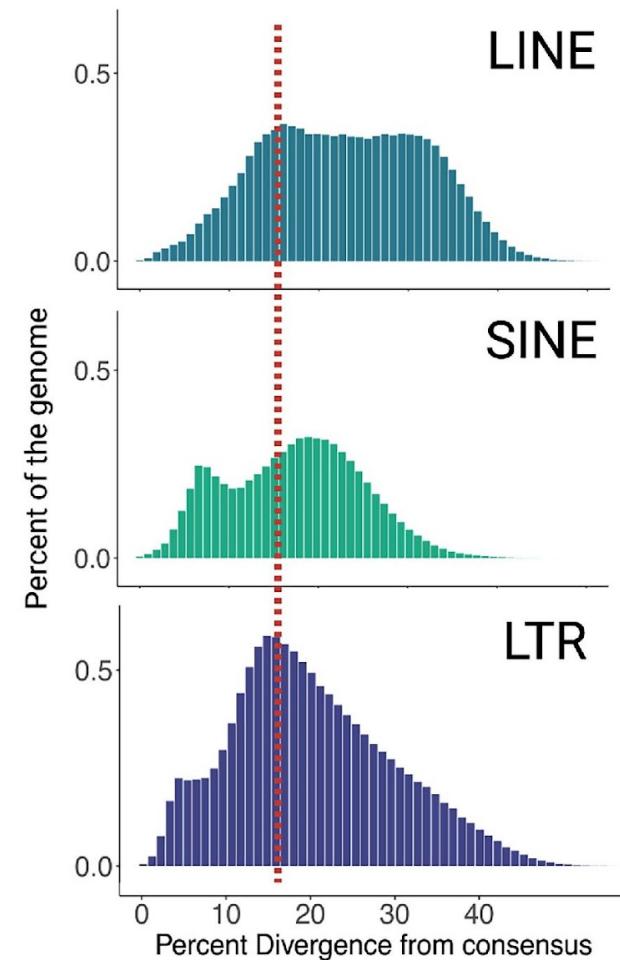


Genomes as ecosystems

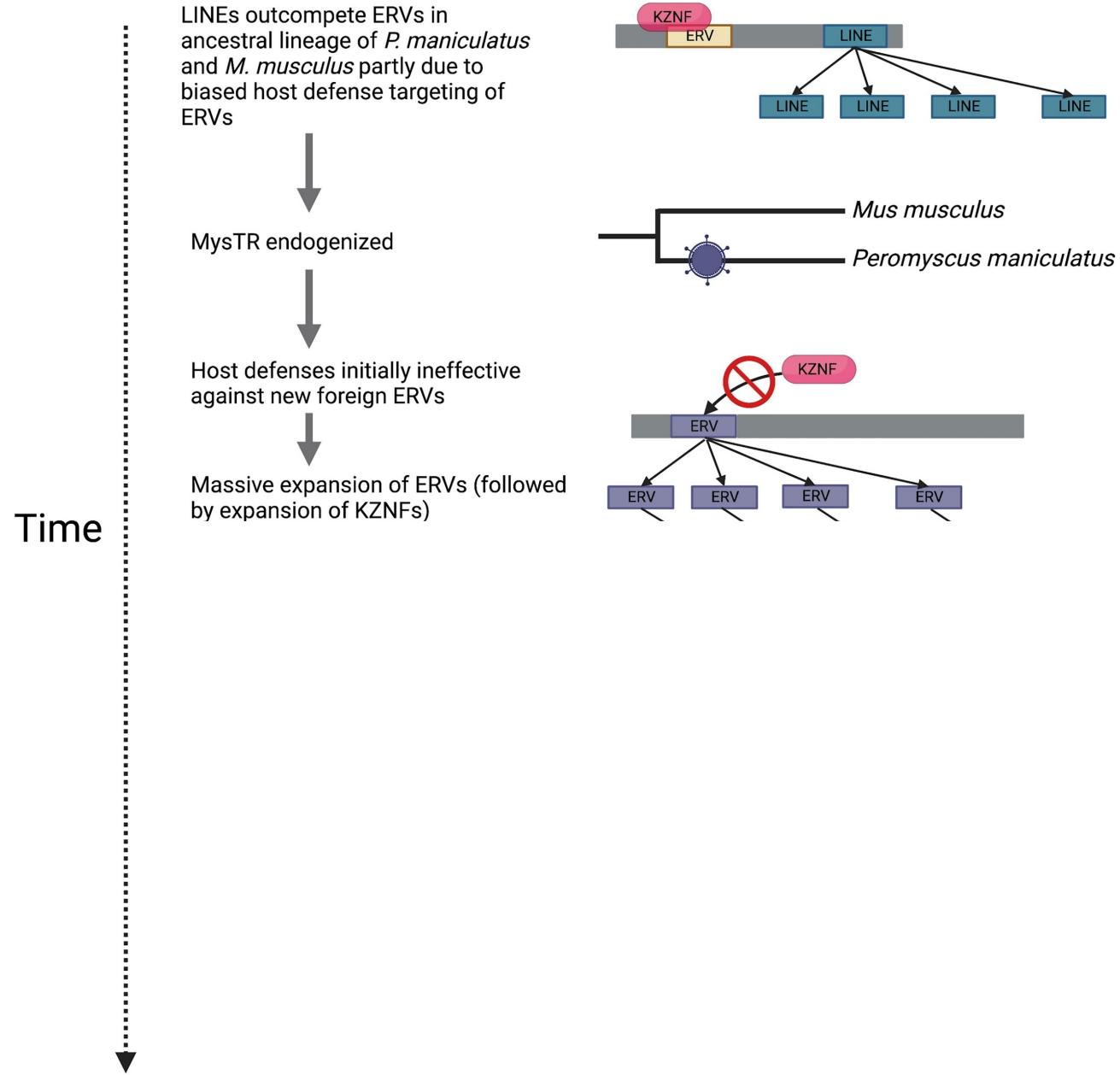
A



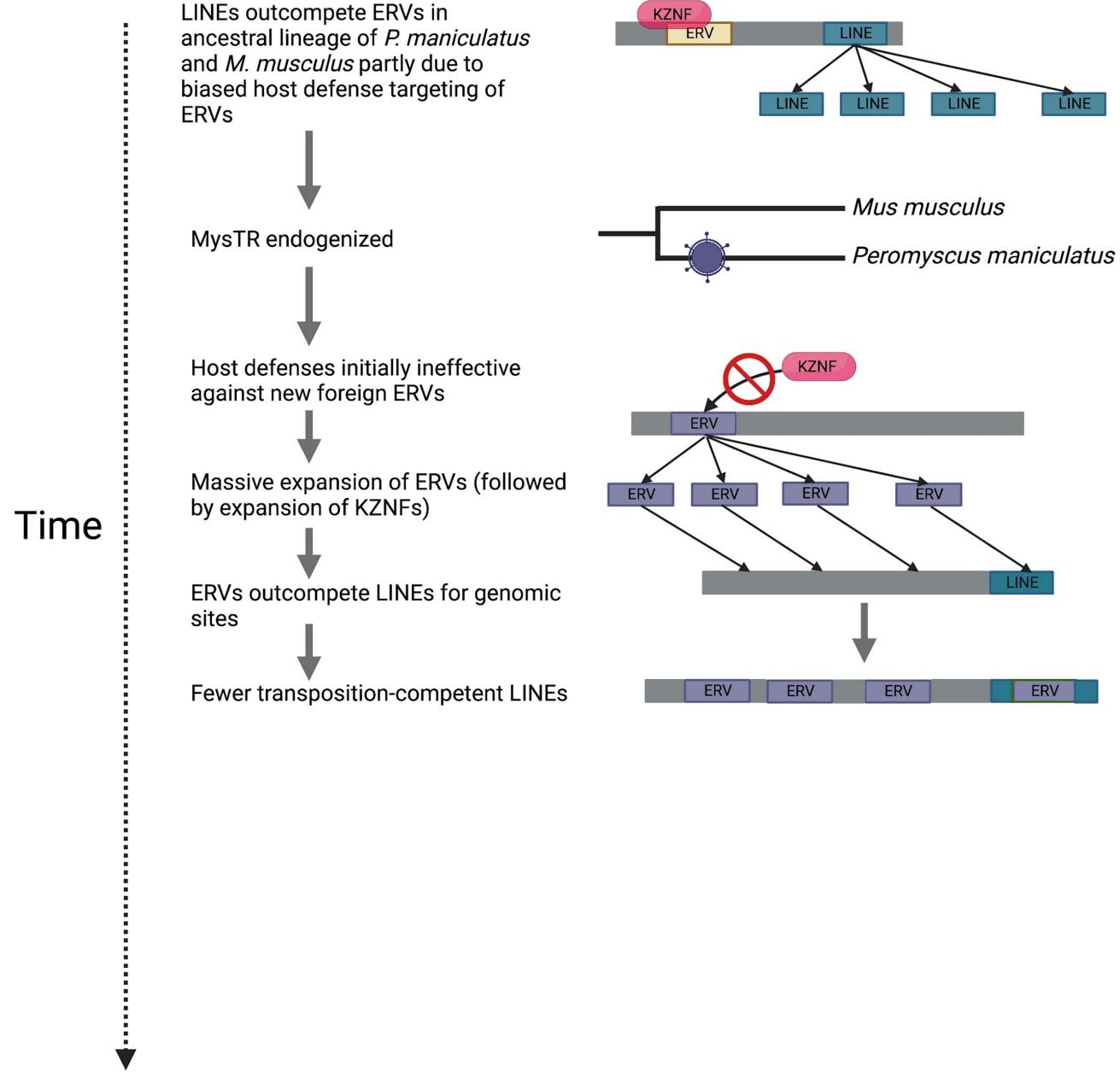
B



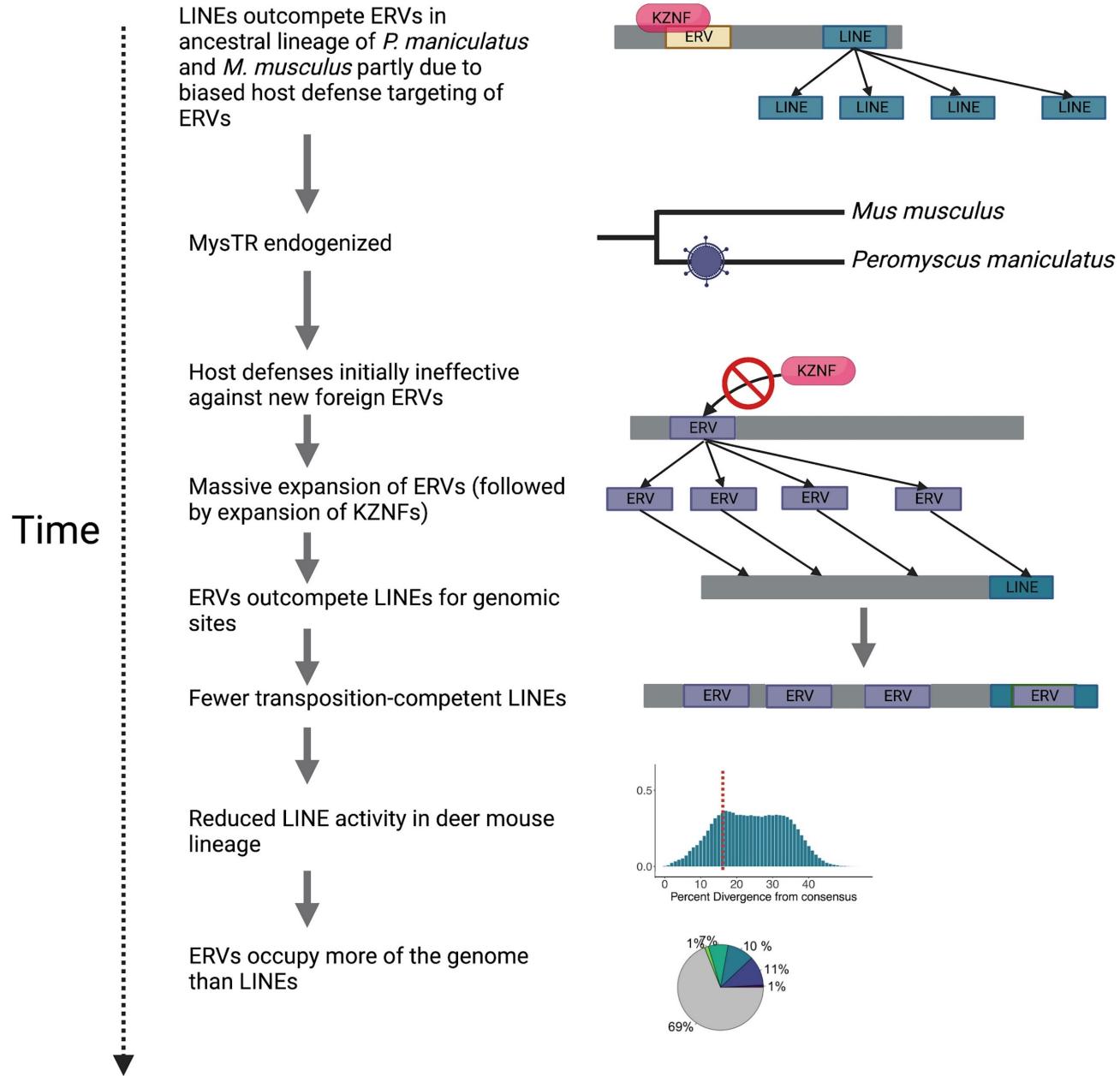
Genomes as ecosystems



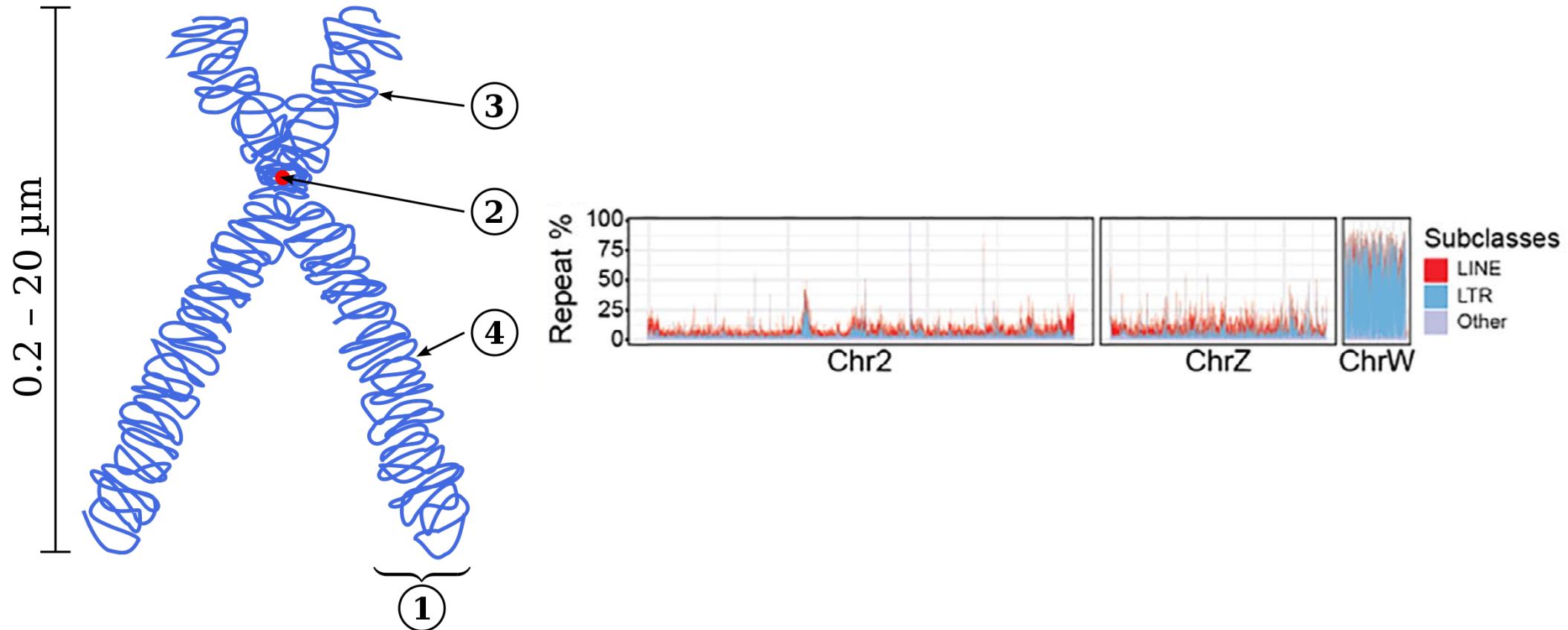
Genomes as ecosystems



Genomes as ecosystems



TEs are not uniformly distributed



Target specificity

R2

 | RESEARCH ARTICLE | STRUCTURAL BIOLOGY

Structure of the R2 non-LTR retrotransposon initiating target-primed reverse transcription

MAX E. WILKINSON , CHRIS J. FRANGIEH, RHIANNON K. MACRAE , AND FENG ZHANG  [Authors Info & Affiliations](#)

SCIENCE • 6 Apr 2023 • Vol 380, Issue 6642 • pp. 301-308 • DOI: 10.1126/science.adg7883

RESEARCH ARTICLE | GENETICS | 



The retrotransposon R2 maintains *Drosophila* ribosomal DNA repeats

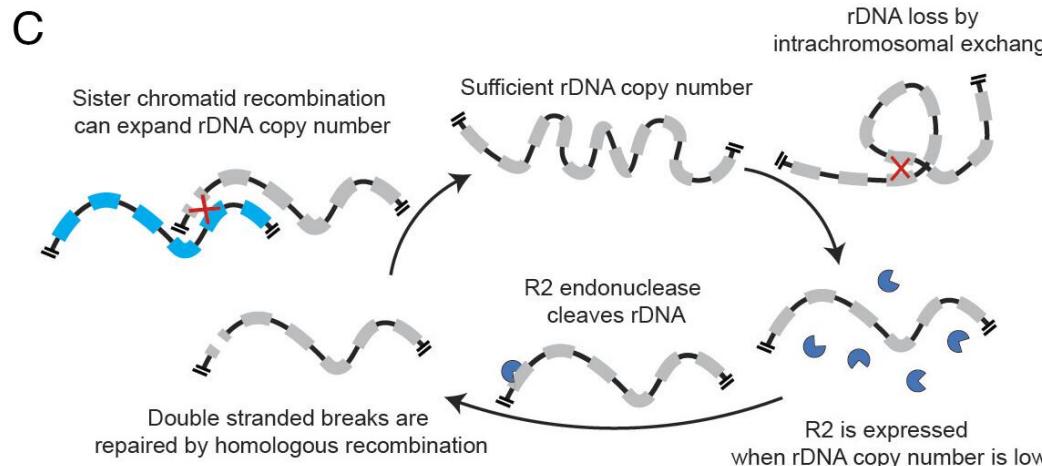
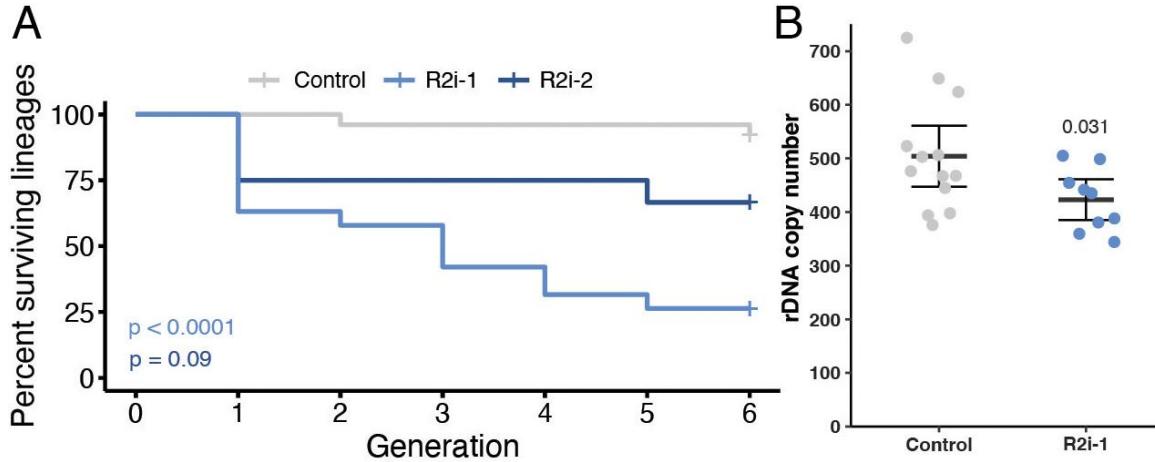
Jonathan O. Nelson , Alyssa Slicko , and Yukiko M. Yamashita   [Authors Info & Affiliations](#)

Edited by R. Scott Hawley, Stowers Institute for Medical Research, Kansas City, MO; received December 20, 2022; accepted May 3, 2023

May 30, 2023 | 120 (23) e2221613120 | <https://doi.org/10.1073/pnas.2221613120>

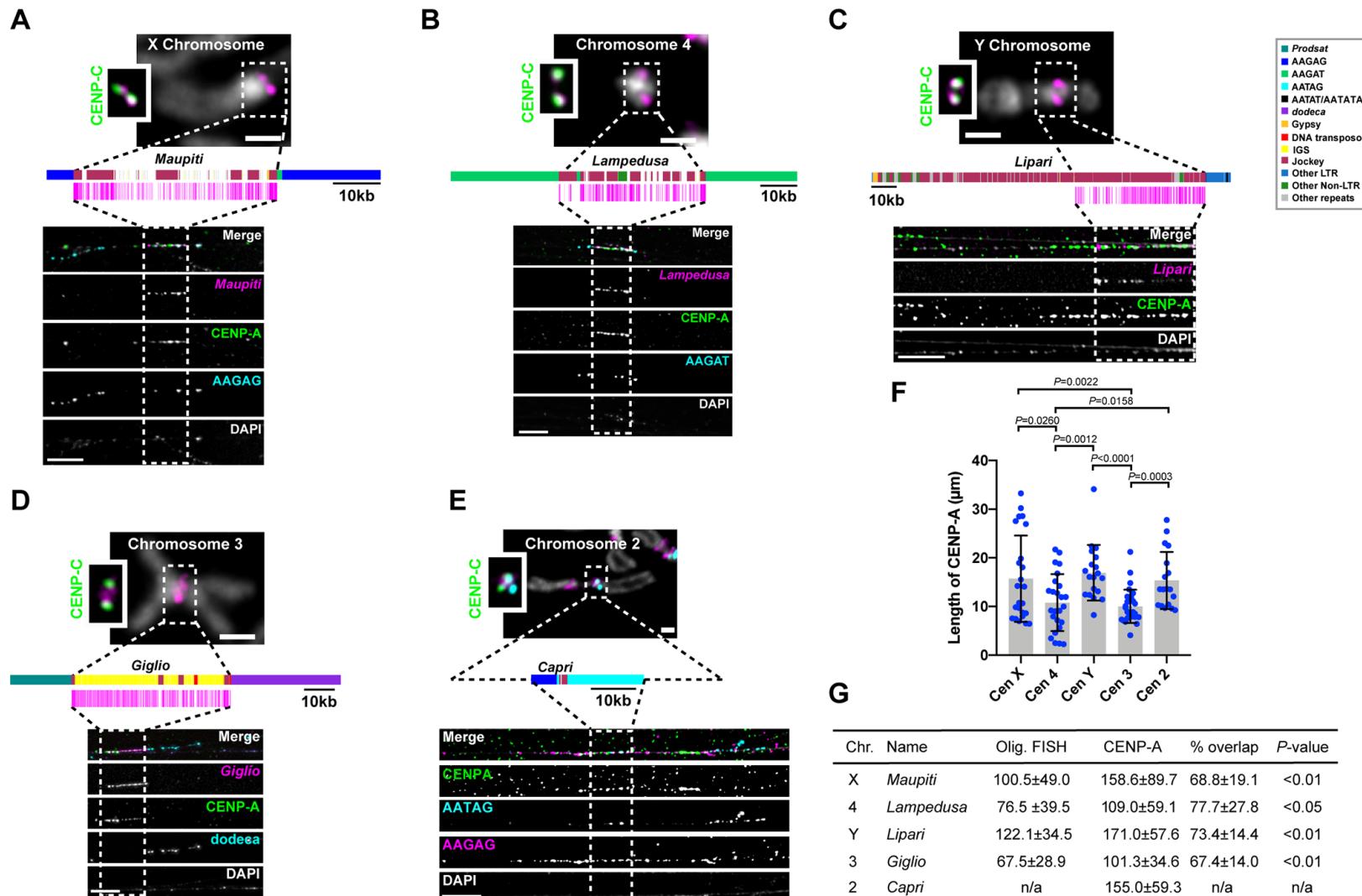
Target specificity

R2



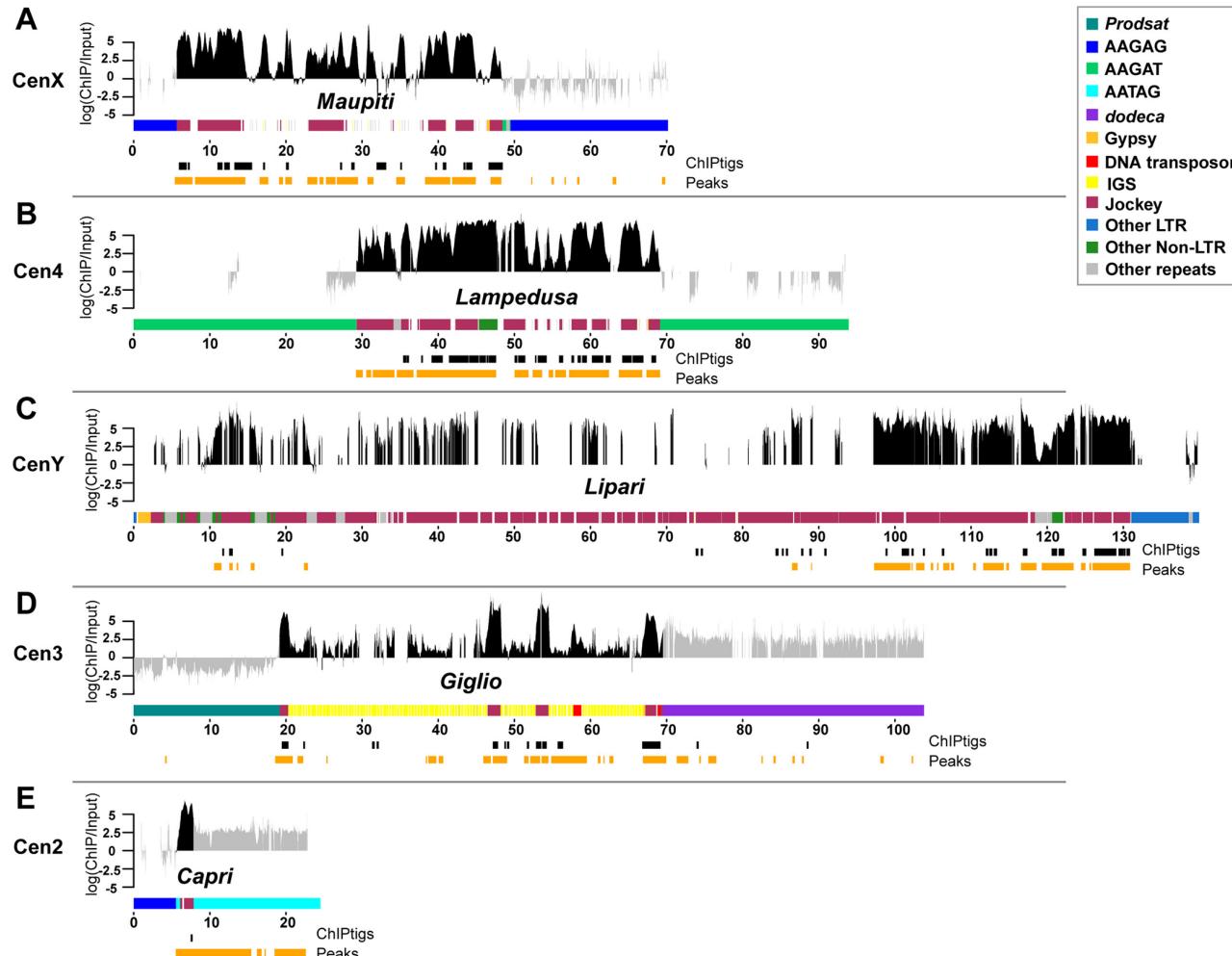
Non-recombining regions

Centromeres



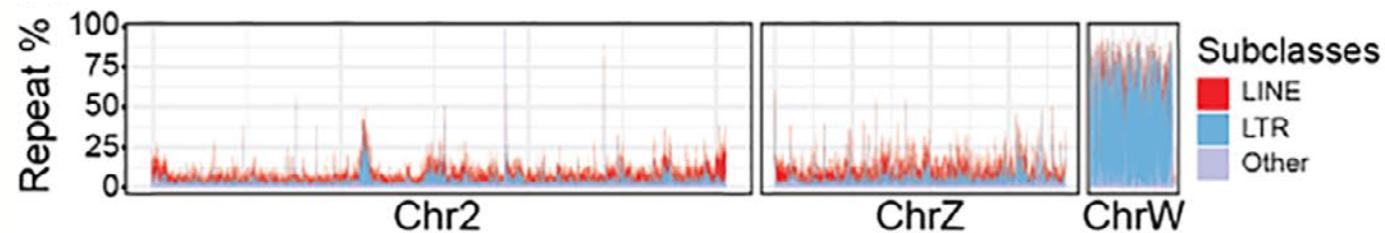
Non-recombining regions

Centromeres



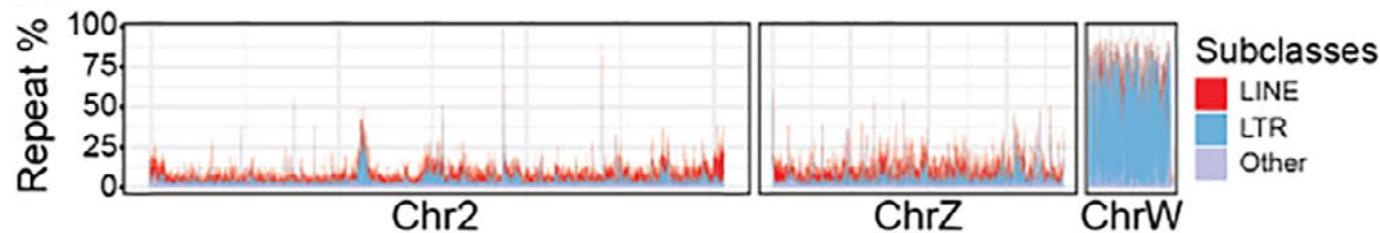
Non-recombining regions

Sex chromosomes

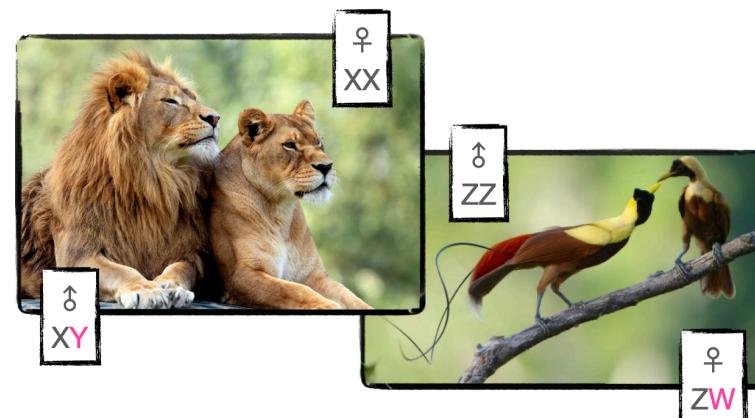


Non-recombining regions

Sex chromosomes

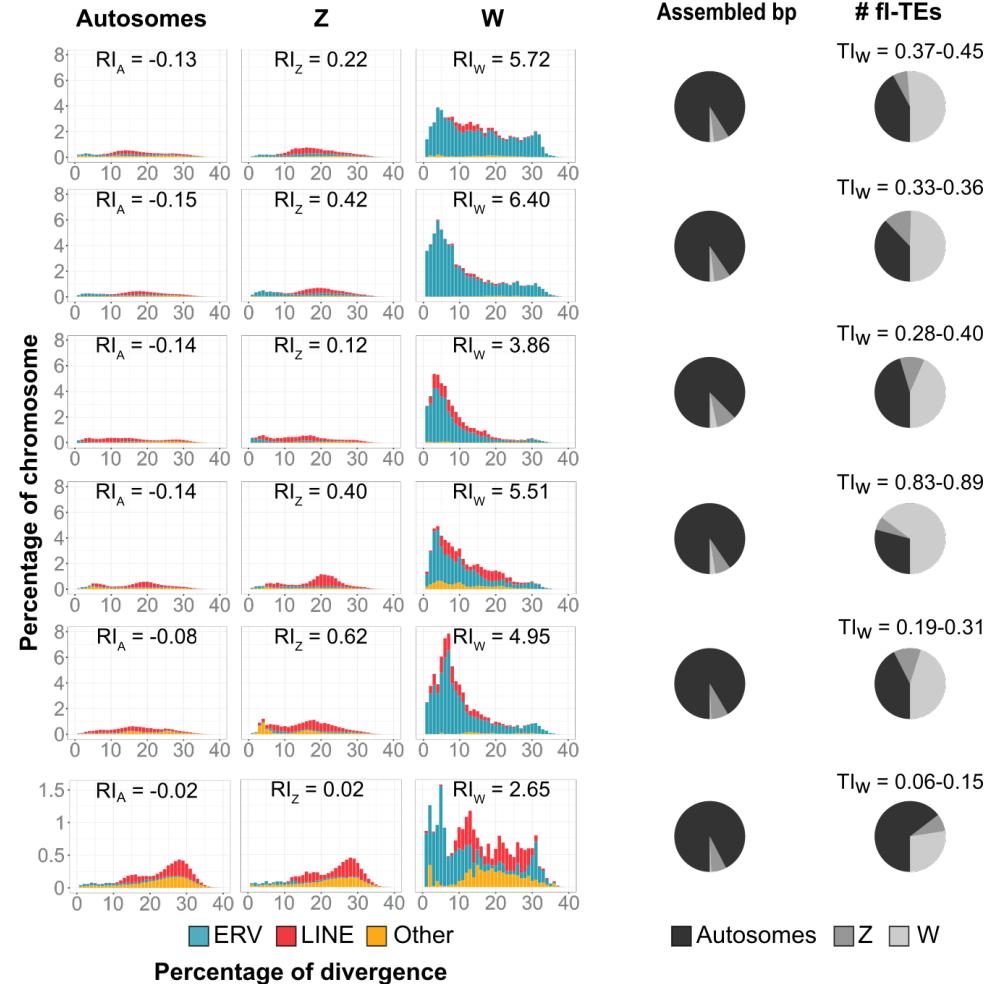
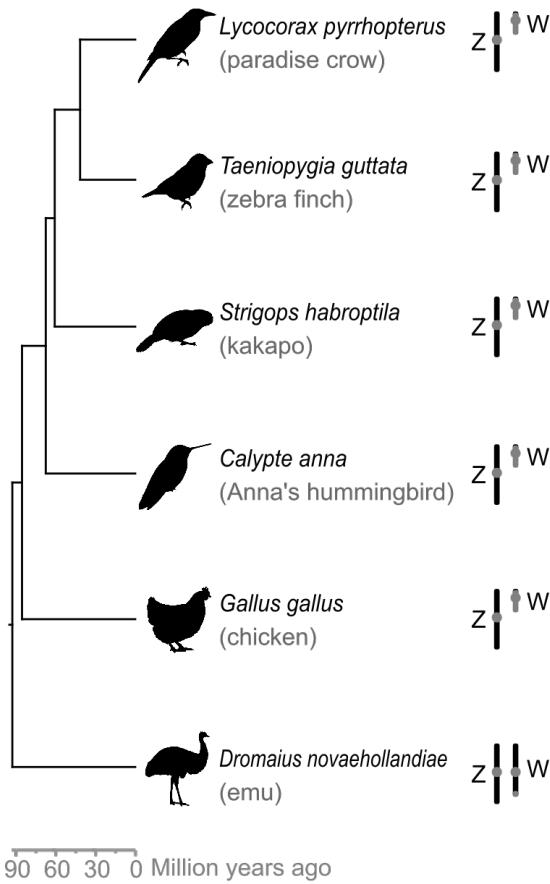


Genetic sex determination



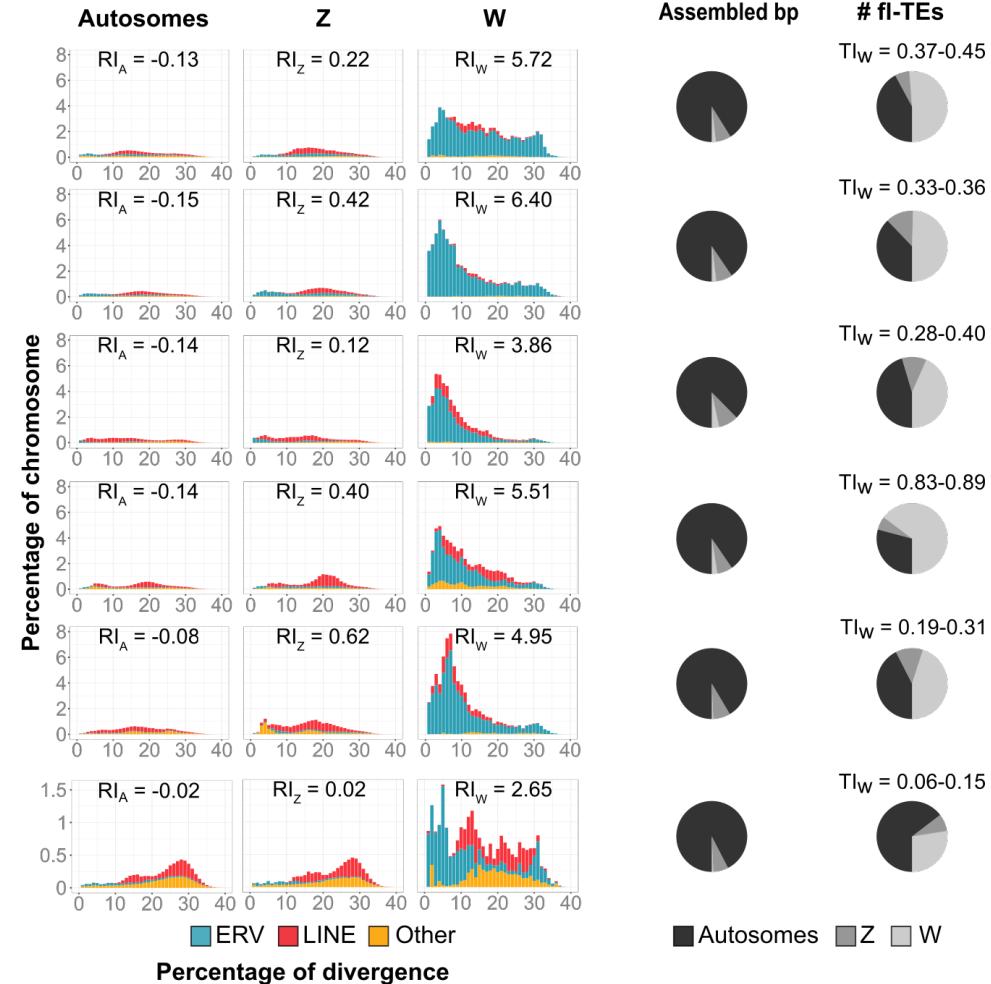
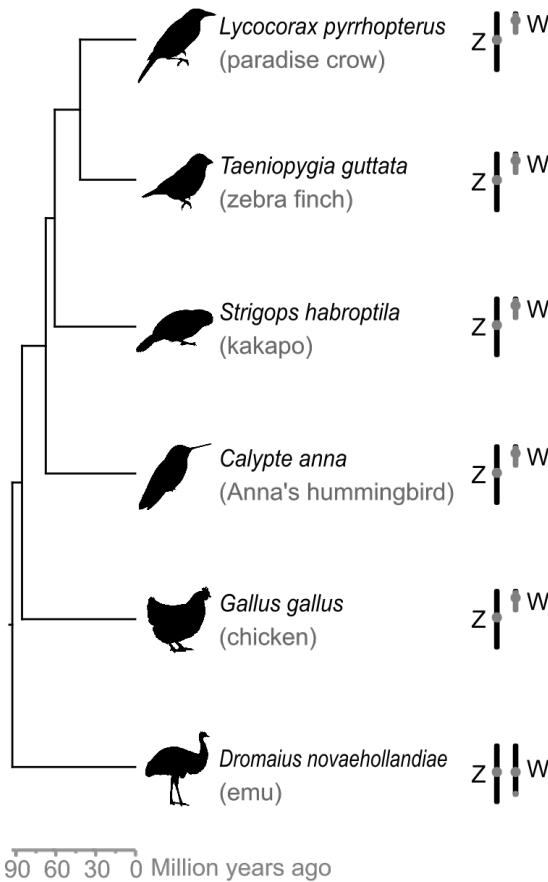
Non-recombining regions

Sex chromosomes



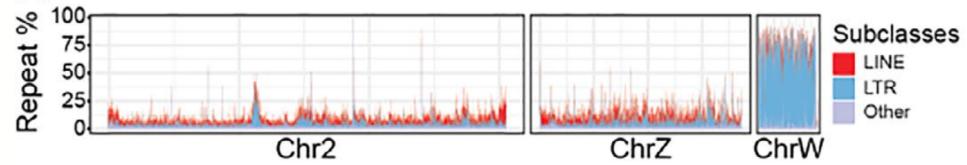
Non-recombining regions

Sex chromosomes



Half of all intact TEs is on W alone

Active transposable elements

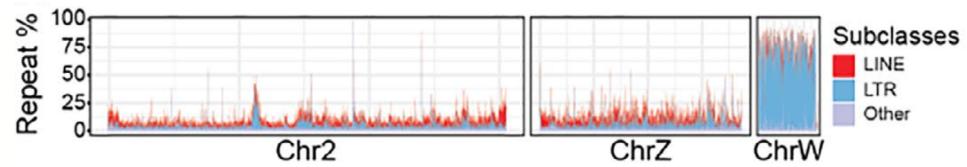


W



TE = transposable element

Active transposable elements

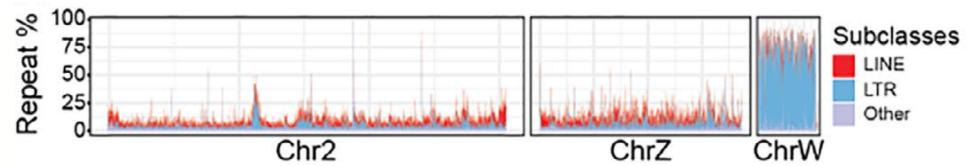


W



TE = transposable element

Active transposable elements



W

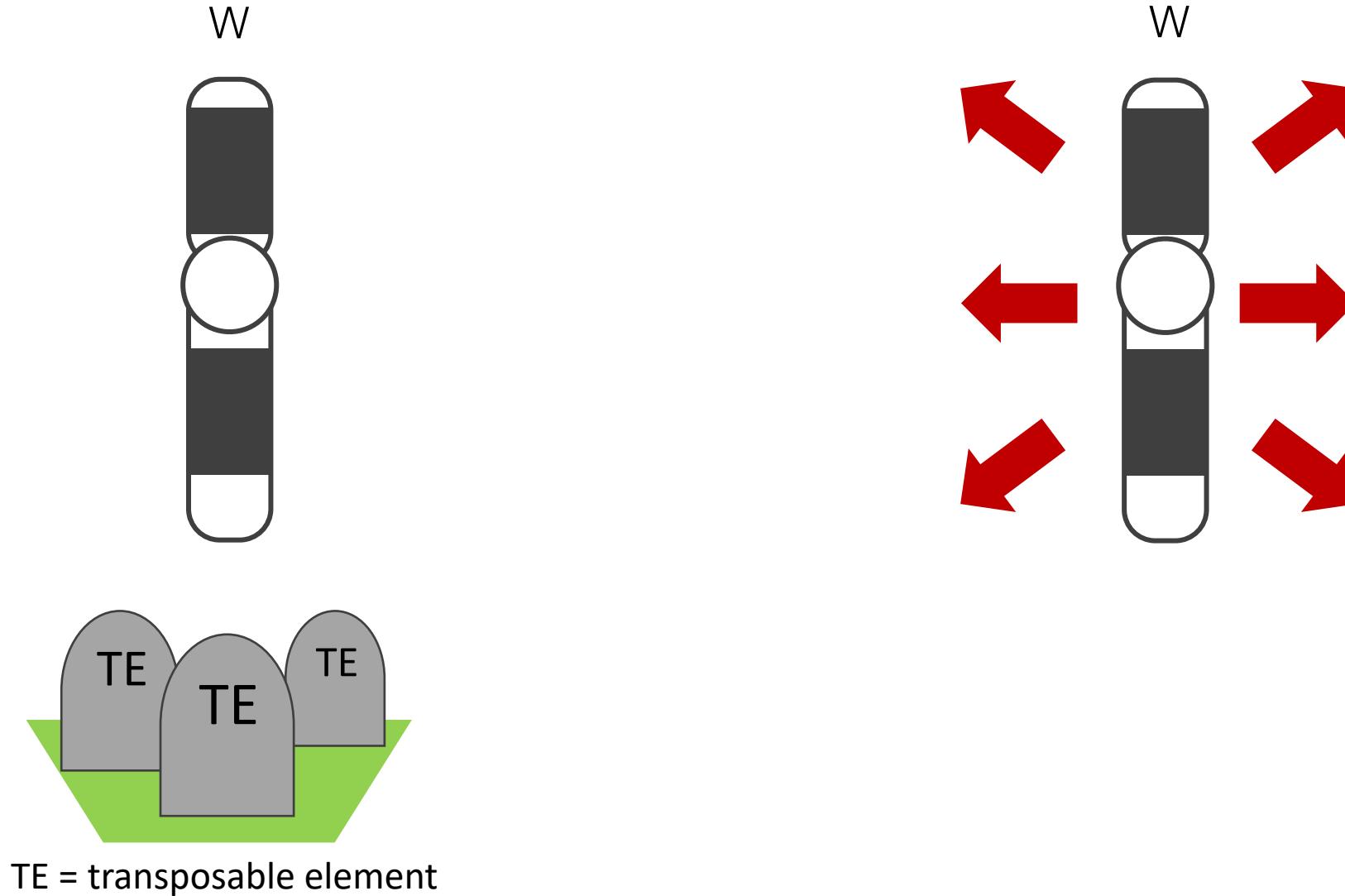
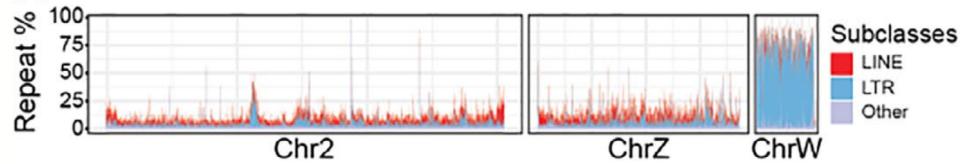


W

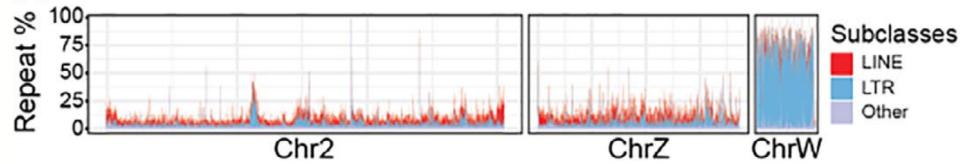


TE = transposable element

Active transposable elements



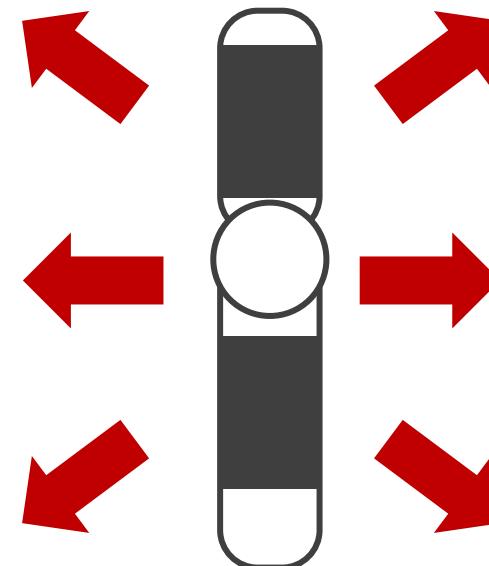
Active transposable elements



W



W



TE = transposable element



Find W-linked TEs

Find W-linked TEs

Identification of W-linked variants of
ITRs

Find W-linked TEs

| Identification of W-linked variants of ITRs | | |
|---|---------|-------------------------|
| | Genomic | s |
| Repeat library | | ACTAAGCCTAACTTG |
| Male | A | ACTAAGCCTAACTTG |
| | Z | ACTAAGCCTAACTTG |
| Female | A | ACTAAGCCTAACTTG |
| | Z | ACTAAGCCTAACTTG |
| | W | GCTAAGCCTAAACATG |

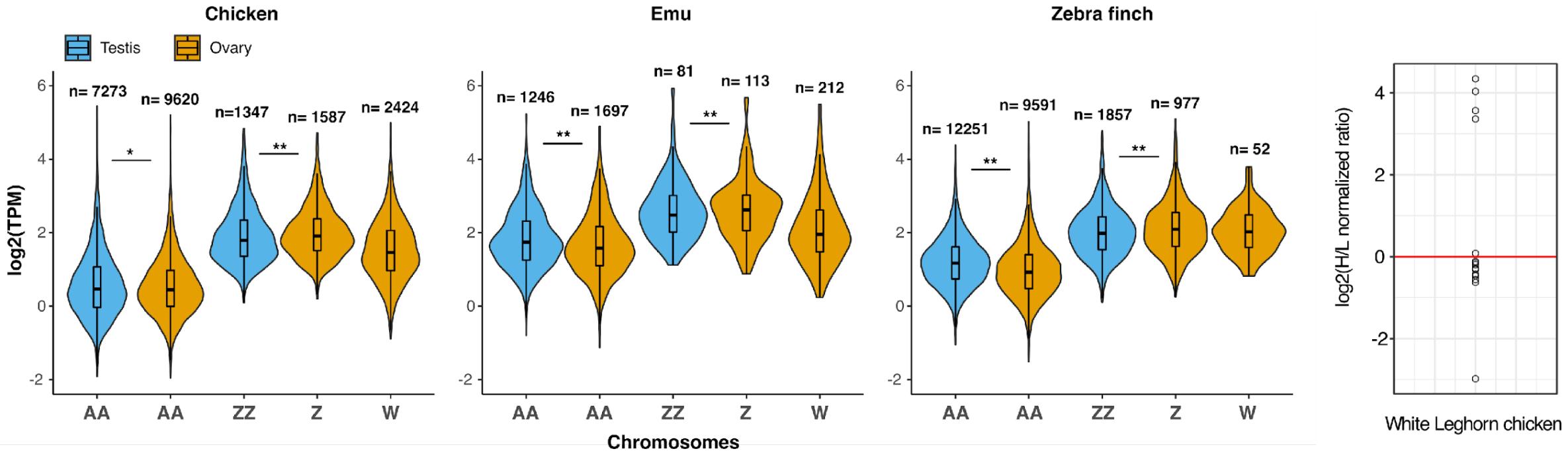
Find W-linked TEs

| | | Identification of W-linked variants of | | |
|----------------|---|--|------------------|------------------|
| | | Genomic | LTRs | Transcriptomics |
| | | s | | |
| Repeat library | | | ACTAAGCCTAACATTG | ACUAAGCCUAACUUG |
| Male | A | | ACTAAGCCTAACATTG | ACUAAGCCUAACUUG |
| | Z | | ACTAAGCCTAACATTG | ACUAAGCCUAACUUG |
| Female | A | | ACTAAGCCTAACATTG | ACUAAGCCUAACUUG |
| | Z | | ACTAAGCCTAACATTG | ACUAAGCCUAACUUG |
| | W | | GCTAAGCCTAACATG | GUUAAGCCUAACAUUG |

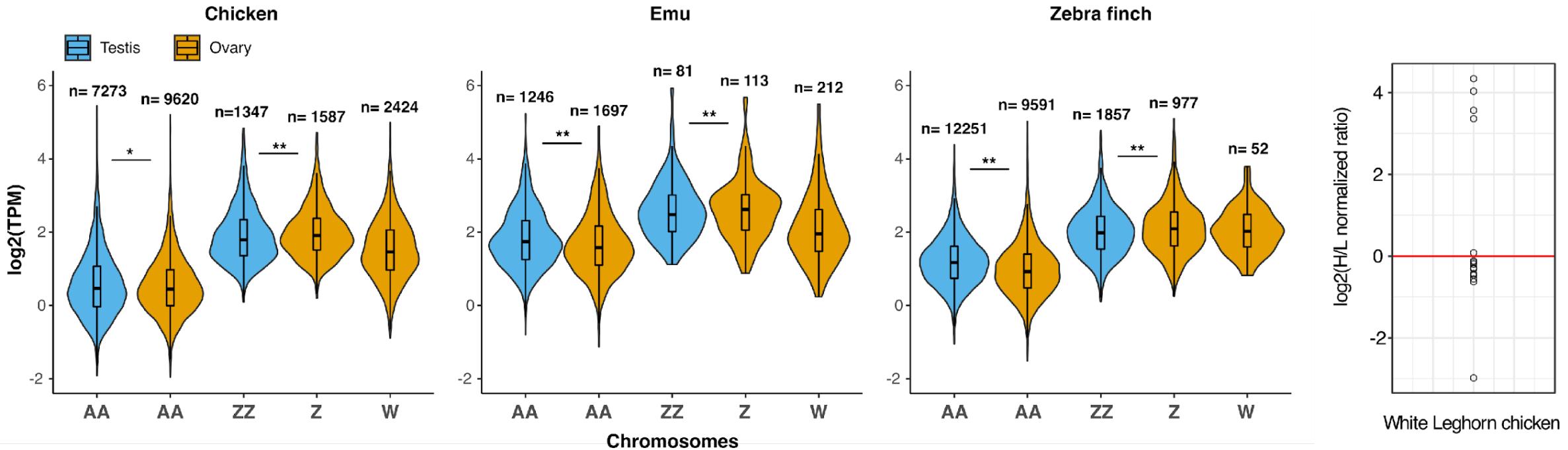
Find W-linked TEs

| | | Identification of W-linked variants of | | |
|-------------------|---|--|-------------------------|-----------------------|
| | | Genomic s | LTs | Transcriptomics |
| | | | | Proteomics |
| Repeat library | | ACTAAGCCTAACTTG | | |
| Male | A | ACTAAGCCTAACTTG | ACUAAGCCUAACUUG | FDANTKPNLDLQVL |
| | Z | ACTAAGCCTAACTTG | ACUAAGCCUAACUUG | FDANTKPNLDLQVL |
| Female | A | ACTAAGCCTAACTTG | ACUAAGCCUAACUUG | FDANTKPNLDLQVL |
| | Z | ACTAAGCCTAACTTG | ACUAAGCCUAACUUG | FDANTKPNLDLQVL |
| | W | GCTAAGCCTAAACATG | GUUAAGCCUAACAUUG | FDANAKPNMDLQVL |

W-linked TEs are more expressed

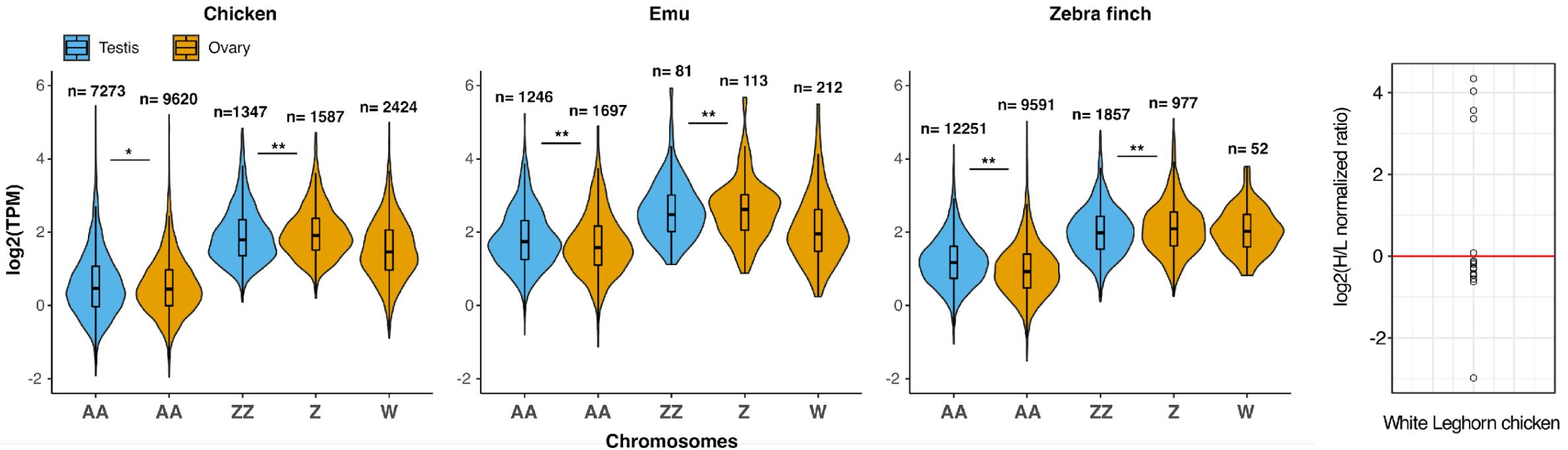


W-linked TEs are more expressed



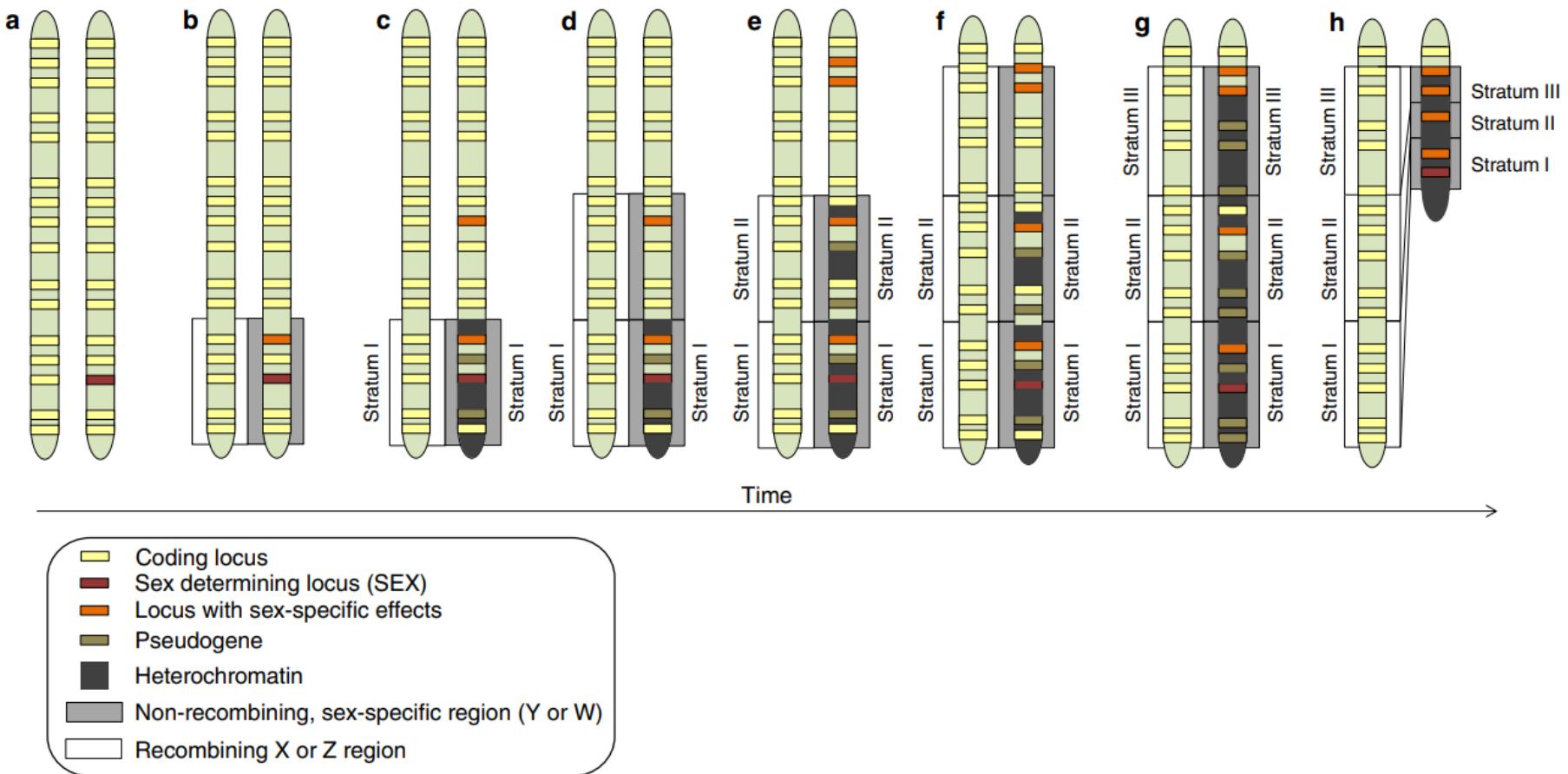
A reservoir (refugium) for active TEs that can have fitness effects on females

W-linked TEs are more expressed



- A reservoir (refugium) for active TEs that can have fitness effects on females
- According to short reads no full-length elements were present on the W

Sex chromosome evolution



Gene evolution

Innovators VS
constraint

Gene structure

Introns

Prp8, the pivotal protein of the spliceosomal catalytic center, evolved from a retroelement-encoded reverse transcriptase

MENSUR DLAKIĆ¹ and ARCADY MUSHEGIAN^{2,3}

¹Department of Microbiology, Montana State University, Bozeman, Montana 59717, USA

²Stowers Institute for Medical Research, Kansas City, Missouri 64110, USA

³Department of Microbiology, Kansas University Medical Center, Kansas City, Kansas 66160, USA

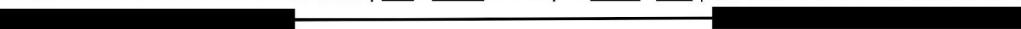
Gene structure

Introners

A *Alternaria alternata*

GTCCAGATGTTGCTCAG | **GTGAGT** // 52bp // **CTGTAG** | CTAACGGTATAAACATC


B *Symbiodinium microadriaticum*

TTCCGAGGTCTTGACAC**GTTG** | CTGGGT // 95bp // ACCTAG | **TTGGT**GCTGCACAAGGA


C *Chrysochromulina sp.*

CTCTTCCAGTCGGT**GCAG** | **GTACGG** // 58bp // GTTTCT | **GCAG**CGGCTGCCGCTTGGCTGC


D *Acanthoeca sp.*

TGGGCACCGCCGCGTTCT | **GTGAGCCGGGC**TTCCATC // 58bp // GATGGAAACGCCGGACAG | GTGTCGTTGGGATCGGG

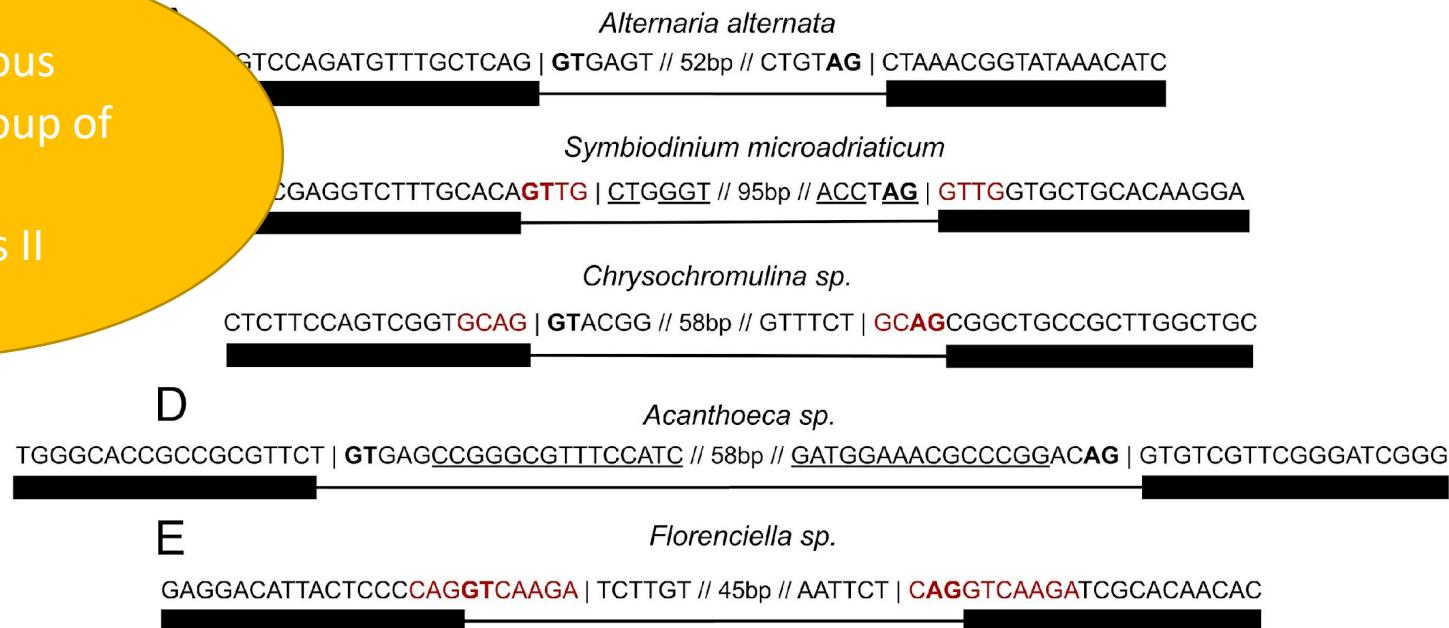

E *Florenciella sp.*

GAGGACATTACTCCCC**CAGGTCAAGA** | TCTTGT // 45bp // AATTCT | **CAGGTCAAGATCGCACAAACAC**

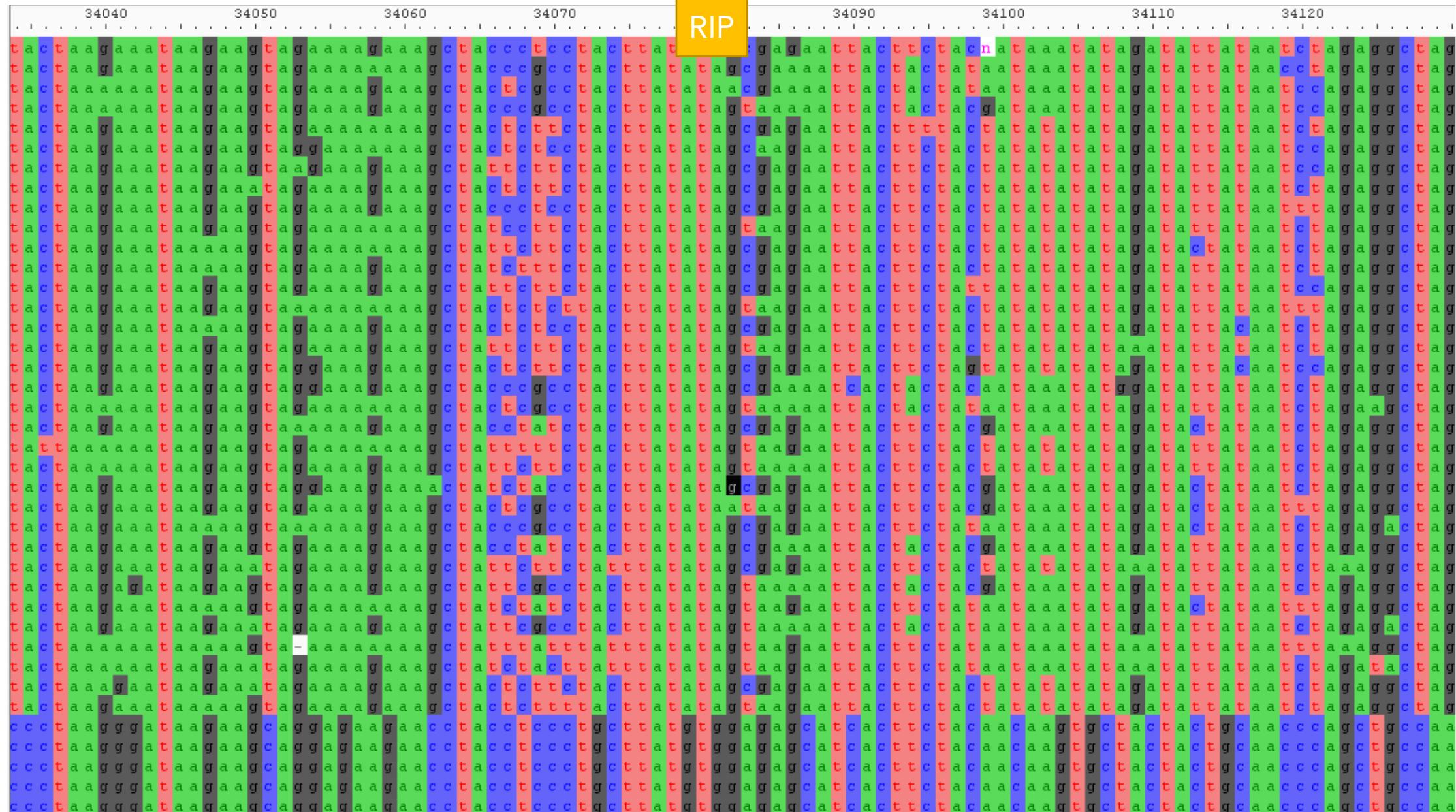

Gene structure

Introners

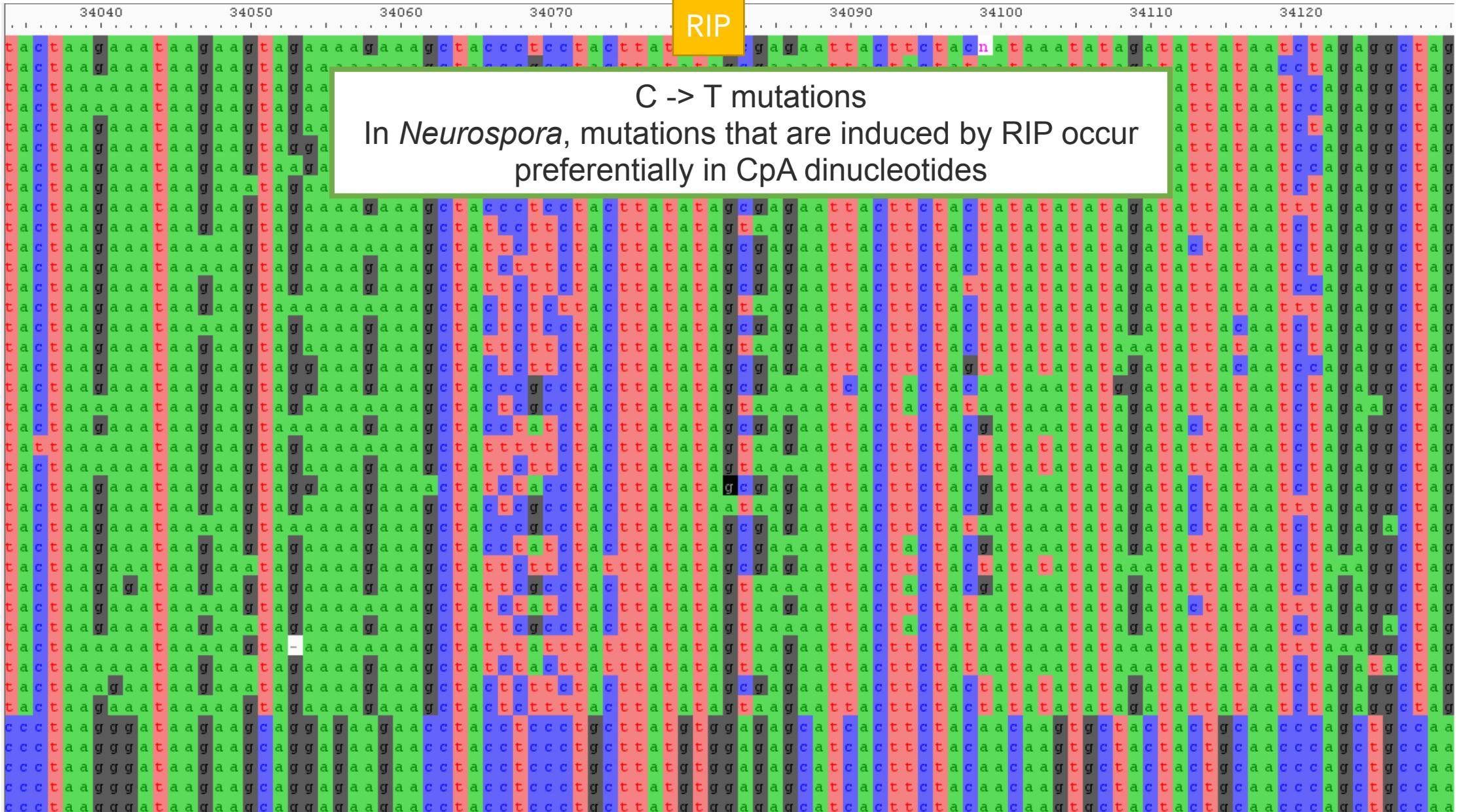
Heterogeneous
polyphyletic group of
TEs.
Mostly Class II



Defences against TEs block gene duplication



Defences against TEs block gene duplication



Defences against TEs block gene duplication

