

# Genomics of eukaryotic microorganisms

Isheng Jason Tsai

[2019 version]



# Lecture objective

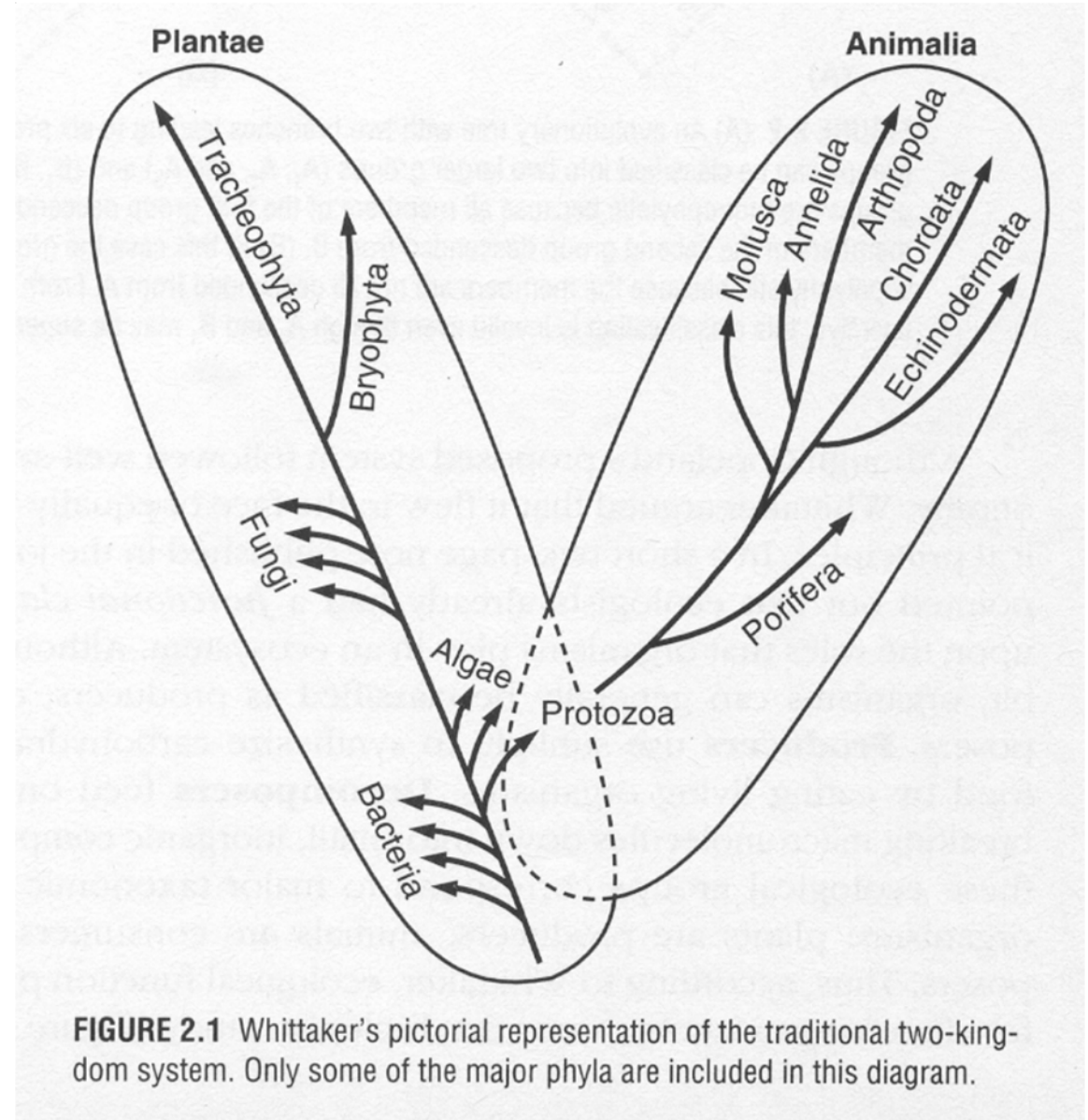
- Classification of eukaryotic microbes – a very general description
- Comparative genomics
  - Inferring orthology
  - Phylogenomics
- Case studies

# Classification of eukaryotic microorganisms – a history

Robert Whittaker



A handwritten signature of Robert Whittaker.



**FIGURE 2.1** Whittaker's pictorial representation of the traditional two-kingdom system. Only some of the major phyla are included in this diagram.

# Classification of eukaryotic microorganisms – a history

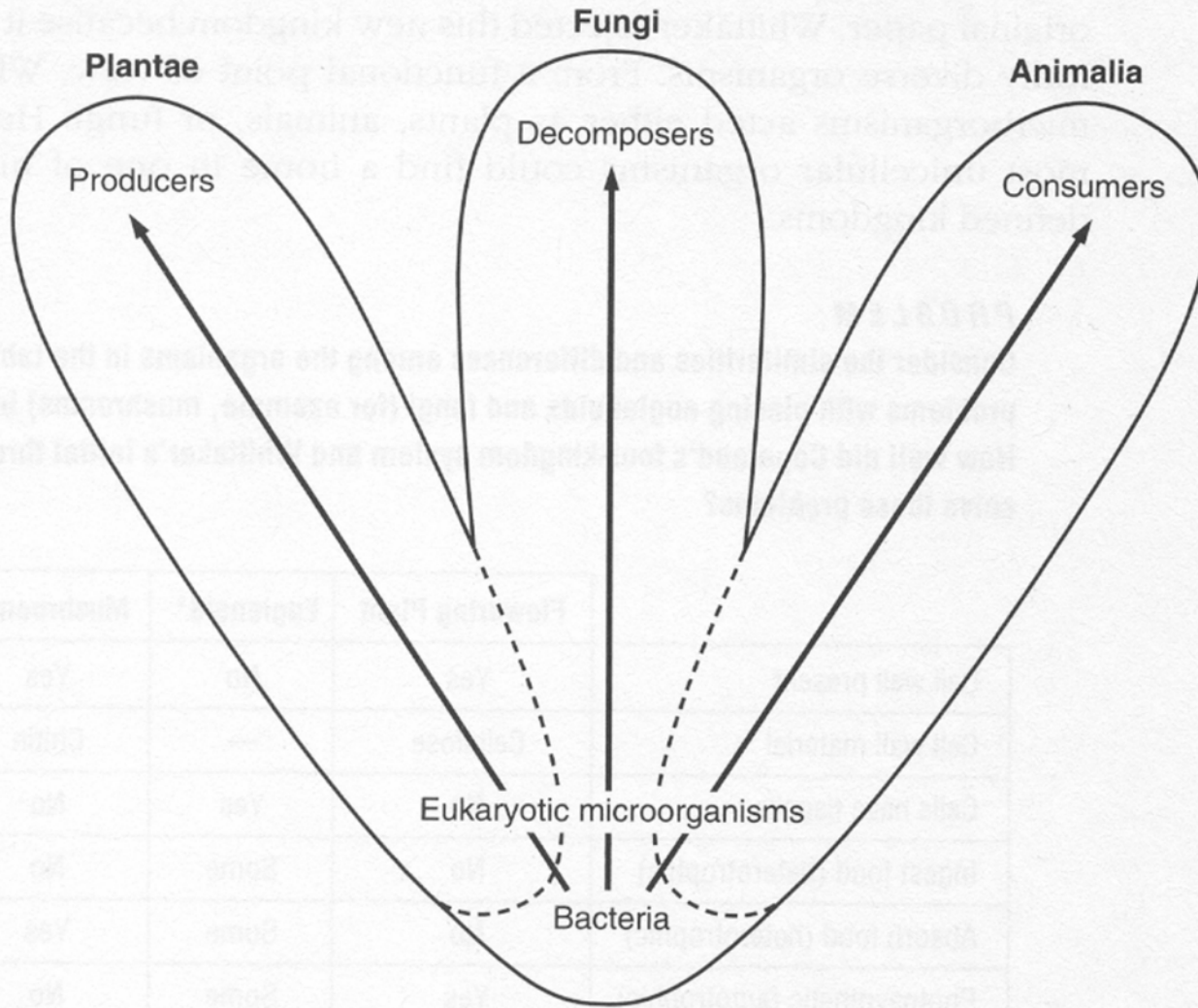
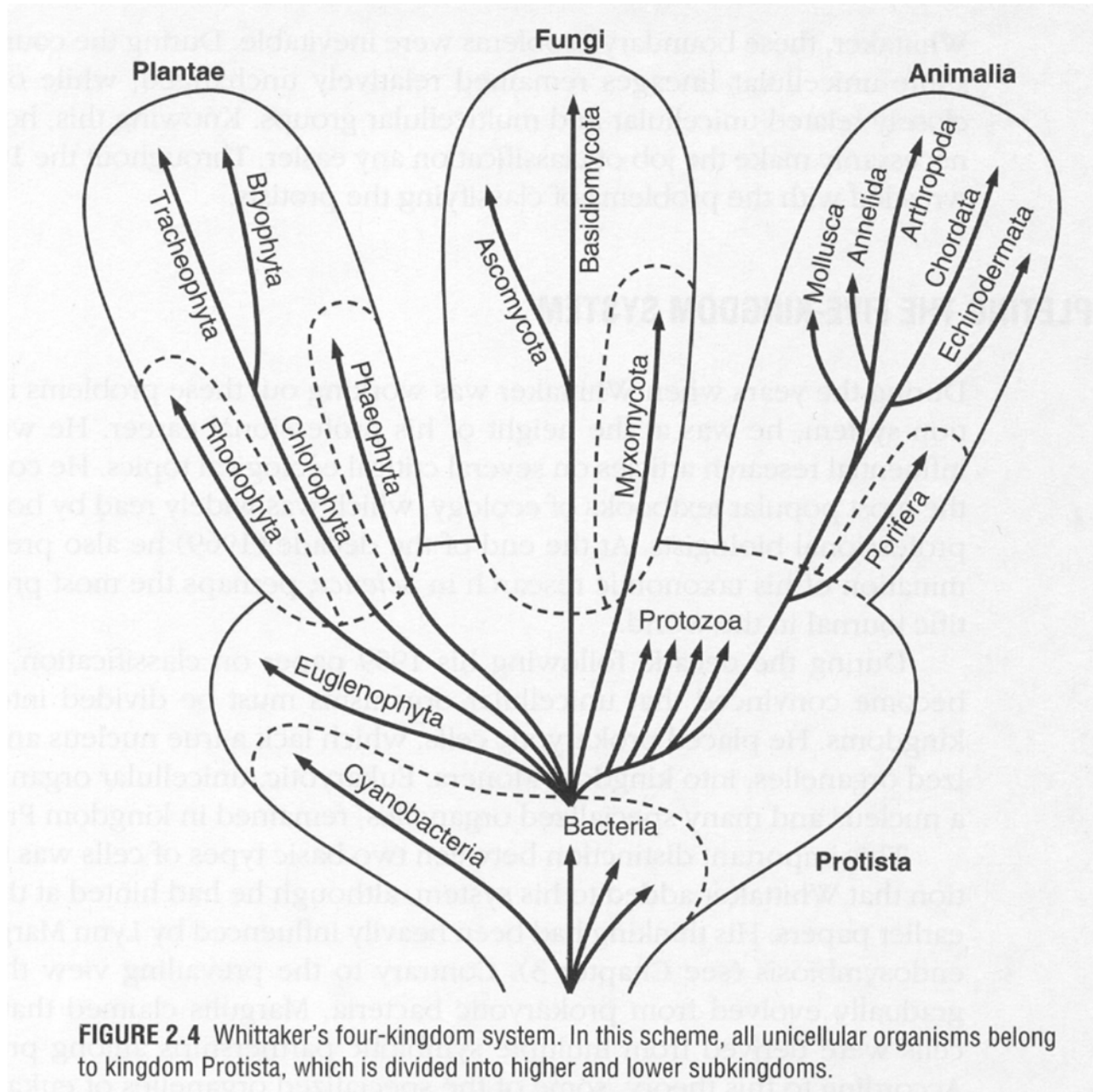


FIGURE 2.3 Whittaker's early three-kingdom system based upon ecological function.

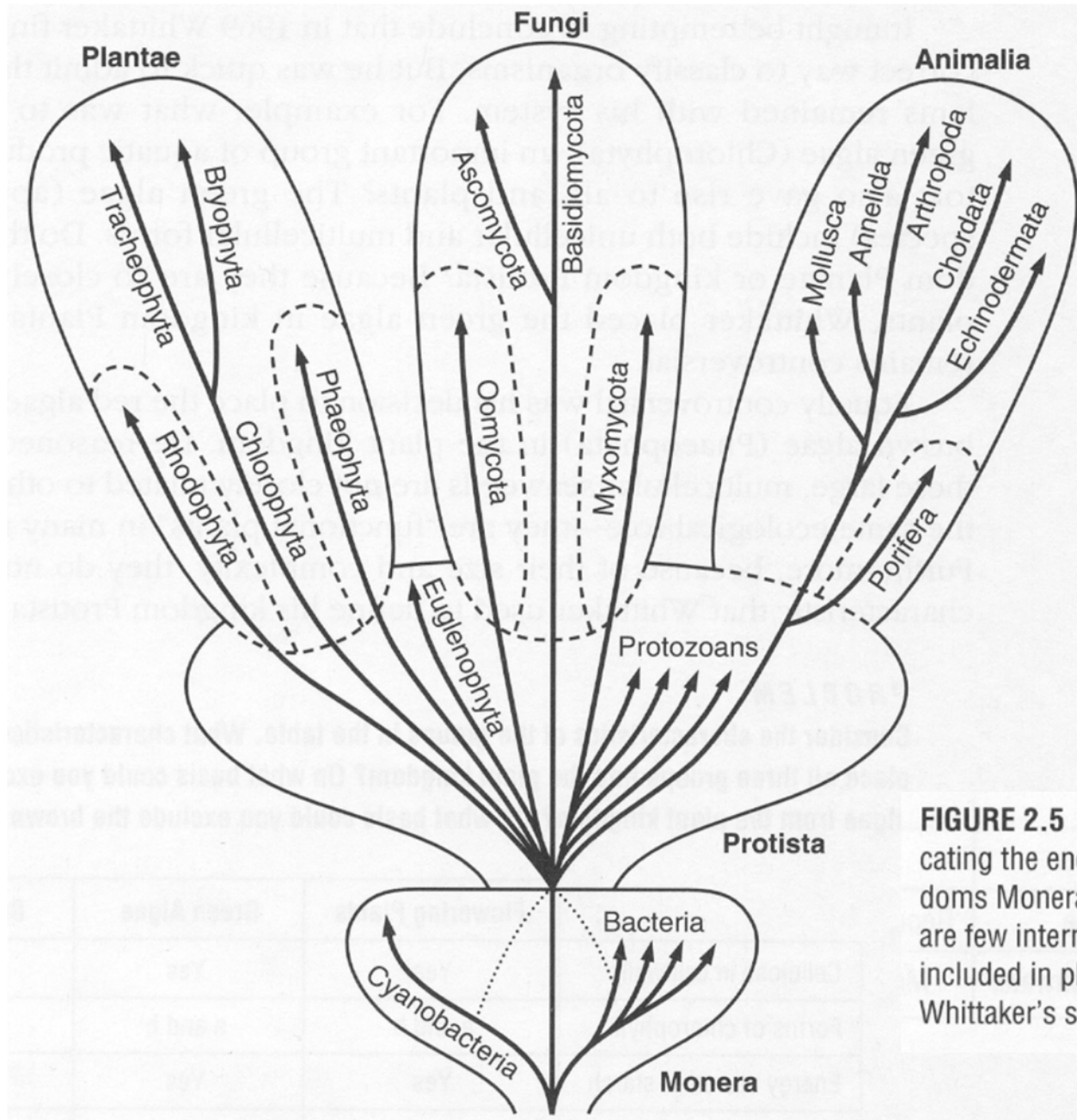
- Ecological classification should reflect three major branches on the evolutionary tree
- Justification was appealing but had serious problems
- Example – Need to place most bacteria in kingdom Fungi
- How about algae and protozoans?

# Classification of eukaryotic microorganisms – four kingdoms (1957)



**FIGURE 2.4** Whittaker's four-kingdom system. In this scheme, all unicellular organisms belong to kingdom Protista, which is divided into higher and lower subkingdoms.

# Classification of eukaryotic microorganisms – five kingdoms



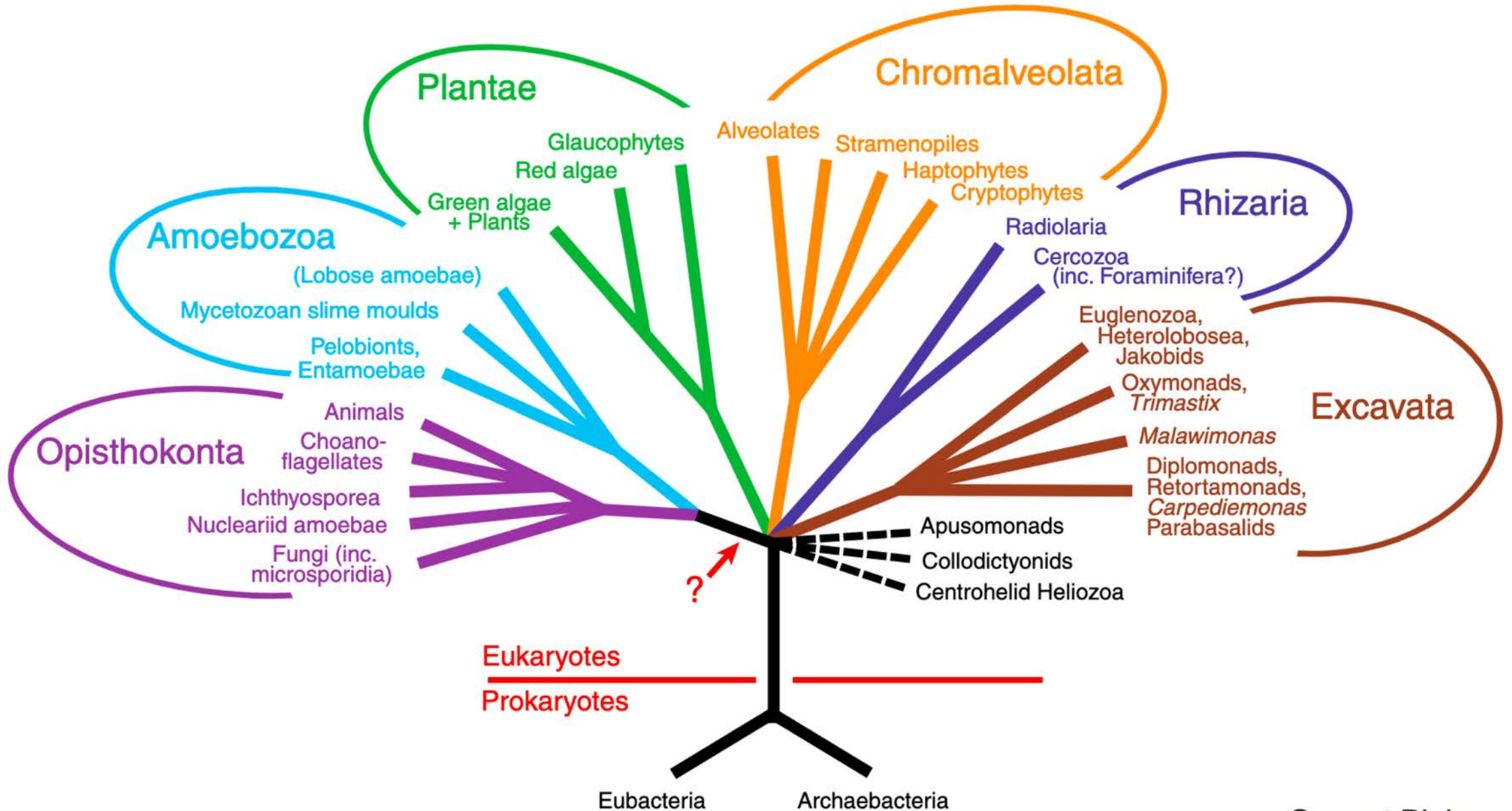
- Widely adopted by all biologists in 1970s
- Still problematic in certain groups possessing both unicellular and multicellular organisms (for example, green algae in plantae or protista?)

**FIGURE 2.5** Whittaker's five-kingdom system. Notice the dotted lines in kingdom Monera indicating the endosymbiotic origin of eukaryotic cells. Also notice that the boundary between kingdoms Monera and Protista is very narrow, because according to the endosymbiotic theory there are few intermediaries between prokaryotic and eukaryotic cells. Ambiguous "problem groups" included in plants, fungi, and animals make each of these multicellular kingdoms polyphyletic in Whittaker's scheme.

# Current standing of kingdoms

<b>Linnaeus</b> 1735 <sup>[29]</sup>	<b>Haeckel</b> 1866 <sup>[30]</sup>	<b>Chatton</b> 1925 <sup>[31][32]</sup>	<b>Copeland</b> 1938 <sup>[33][34]</sup>	<b>Whittaker</b> 1969 <sup>[35]</sup>	<b>Woese et al.</b> 1977 <sup>[36][37]</sup>	<b>Woese et al.</b> 1990 <sup>[38]</sup>	<b>Cavalier-Smith</b> 1993 <sup>[39][40][41]</sup>	<b>Cavalier-Smith</b> 1998 <sup>[42][43][44]</sup>	<b>Ruggiero et al.</b> 2015 <sup>[45]</sup>
2 kingdoms	3 kingdoms	2 empires	4 kingdoms	5 kingdoms	6 kingdoms	3 domains	8 kingdoms	6 kingdoms	7 kingdoms
	Protista	Prokaryota	Monera	Monera	Eubacteria	Bacteria	Eubacteria	Bacteria	Bacteria
					Archaeobacteria	Archaea	Archaeobacteria		Archaea
<i>(not treated)</i>		Eukaryota	Protista	Protista	Protista	Eucarya	Archezoa	Protozoa	Protozoa
							Protozoa		
Vegetabilia	Plantae		Plantae	Plantae	Plantae		Chromista	Chromista	Chromista
				Fungi	Fungi		Plantae	Plantae	Plantae
							Fungi	Fungi	Fungi
Animalia	Animalia		Animalia	Animalia	Animalia		Animalia	Animalia	Animalia

# The real 'kingdoms' of eukaryotes





# The Revised Classification of Eukaryotes (2012)

	Super-groups	Examples
Eukaryota	Amorphea	Amoebozoa Tubulinea Mycetozoa
		Opisthokonta Fungi Choanomonada Metazoa
		Apusomonada
		Breviata
		Excavata Metamonada Malawimonas Discoba
	Diaphoretickes	Cryptophyceae
		Centrohelida
		Telonemia
	Sar	Haptophyta
		Cercozoa Foraminifera "Radiolaria"
Alveolata Stramenopiles		
Archaeplastida	Glaucophyta	
	Rhodophyceae	
	Chloroplastida	
Incertae sedis Eukaryota	Incertae sedis, and table 3	

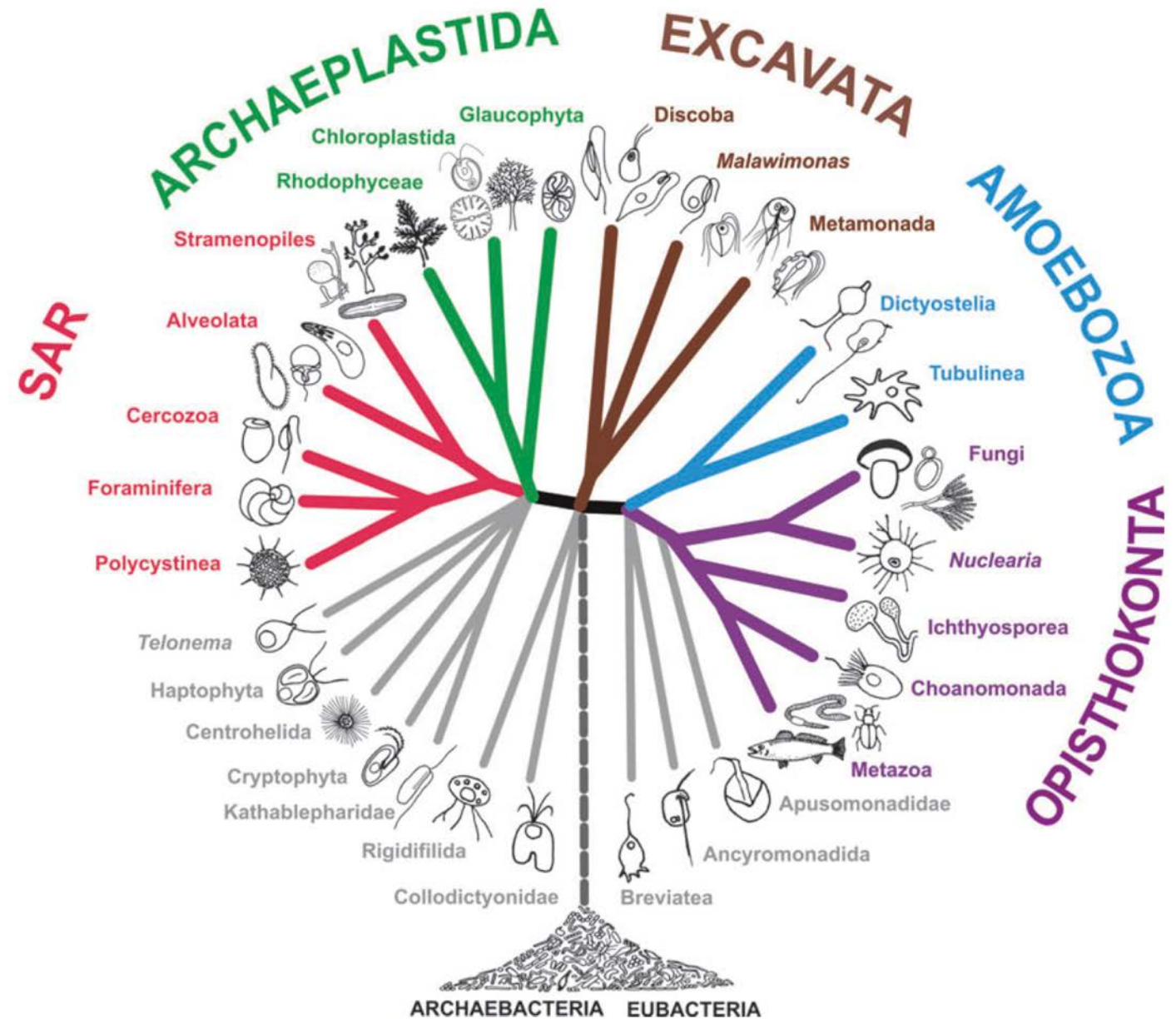


Fig. 1. A view of eukaryote phylogeny reflecting the classification presented herein.

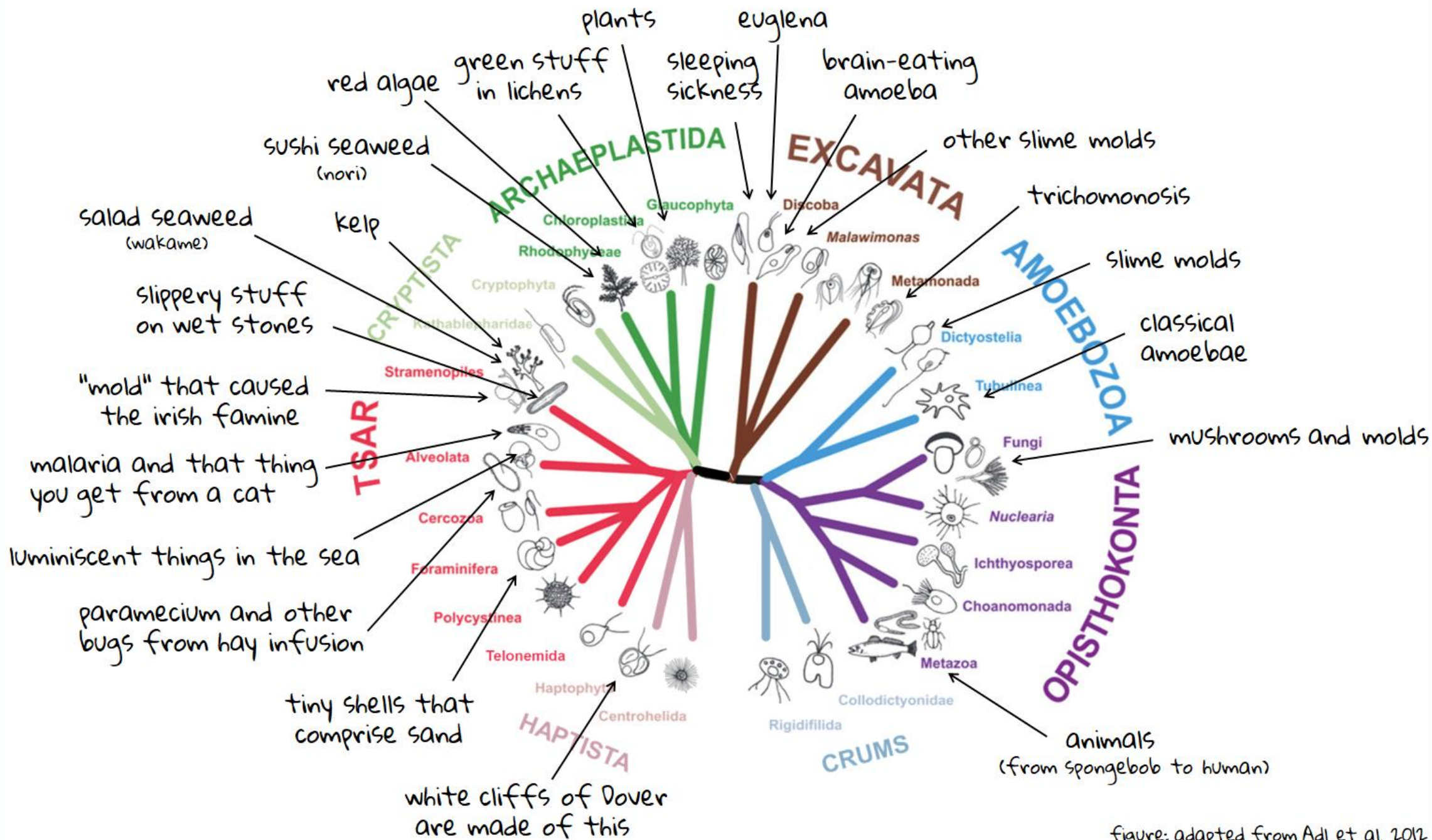
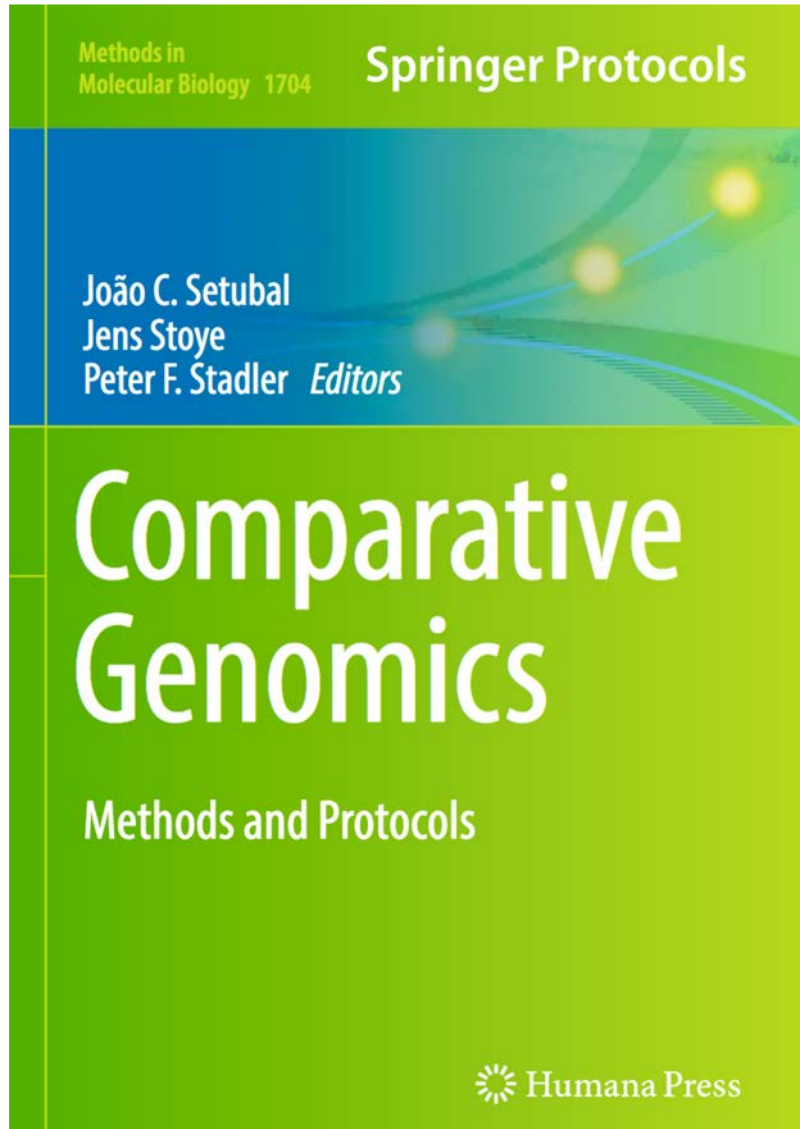


figure: adapted from Adl et al. 2012  
text: @euglenaria

How to establish all these relationships?

And what can we reveal from these relationships?

# Recommended book and paper

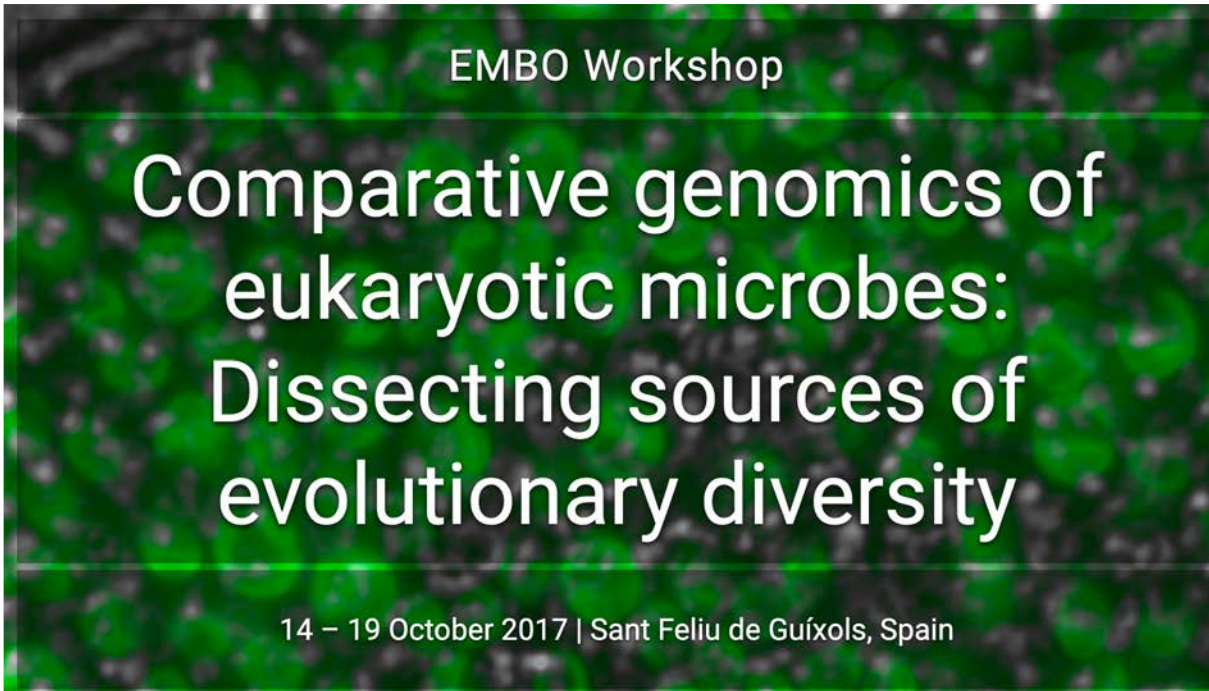


## Orthology: definitions, inference, and impact on species phylogeny inference

Rosa Fernández<sup>1</sup>, Toni Gabaldón<sup>1,2,3,\*</sup>, Christophe Dessimoz<sup>4,5,6,7,8,\*</sup>

**Abstract:** Orthology is a central concept in evolutionary and comparative genomics, used to relate corresponding genes in different species. In particular, orthologs are needed to infer species trees. In this chapter, we introduce the fundamental concepts of orthology relationships and orthologous groups, including some non-trivial (and thus commonly misunderstood) implications. Next, we review some of the main methods and resources used to identify orthologs. The final part of the chapter discusses the impact of orthology methods on species phylogeny inference, drawing lessons from several recent comparative studies.

# Recommended references



Unicellular eukaryotes comprise the overwhelming majority of eukaryotic cellular and genomic diversity, pervading all branches of the eukaryotic tree of life. **Recent sequencing efforts have significantly increased the number of unicellular eukaryotes for which genomic/cellular/proteomic data are available.**

**Metagenomics, transcriptomics, single-cell genomics, and other approaches are being applied to unravel the ecology, physiology, diversity and evolution of microbial eukaryotes and are shedding light on fundamental questions such as the origin of the eukaryotic cell, endosymbiosis, the origin of multicellularity and the evolution of major cellular systems in eukaryotes.**

Although there are conferences devoted to genomics of prokaryotes and that of plants and animals, this EMBO Workshop will be the only forum bringing together this diverse community with a range of expertise and which concentrates on microbial eukaryotes.

There is a need, particularly in organisms without large communities, to have a forum where approaches, new methodologies and datasets can be shared to the advancement of this field.

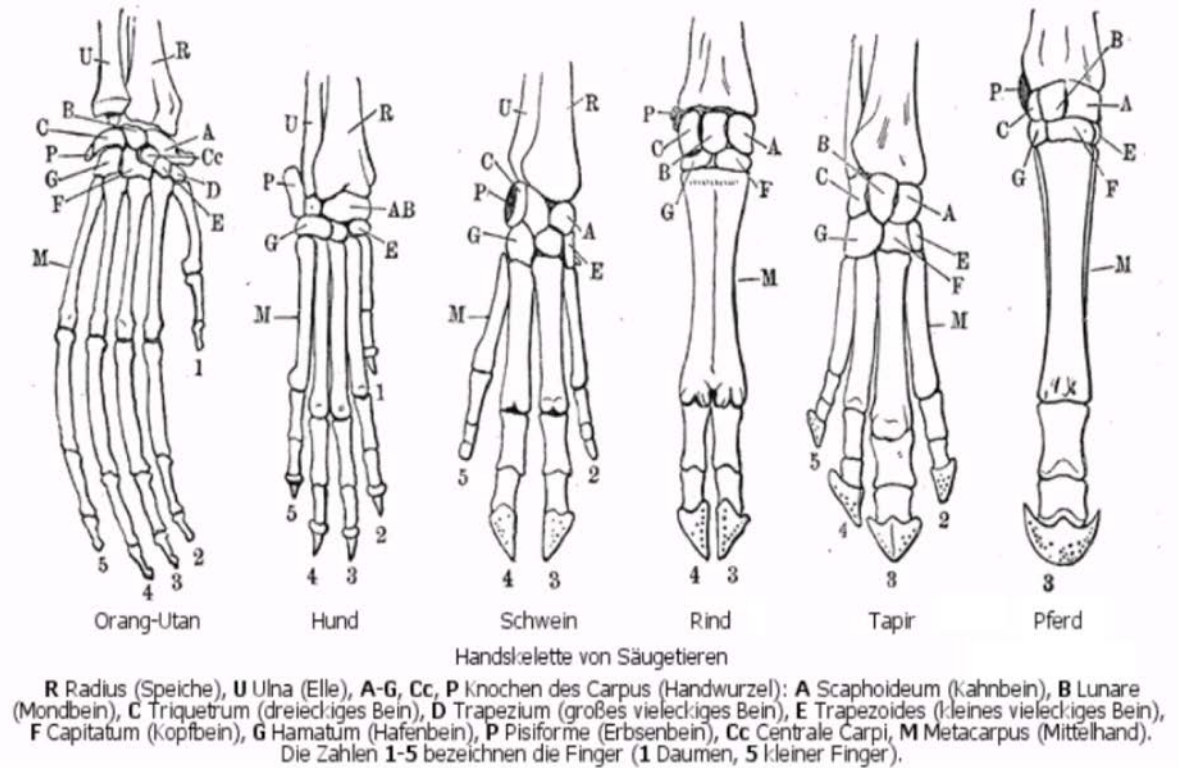
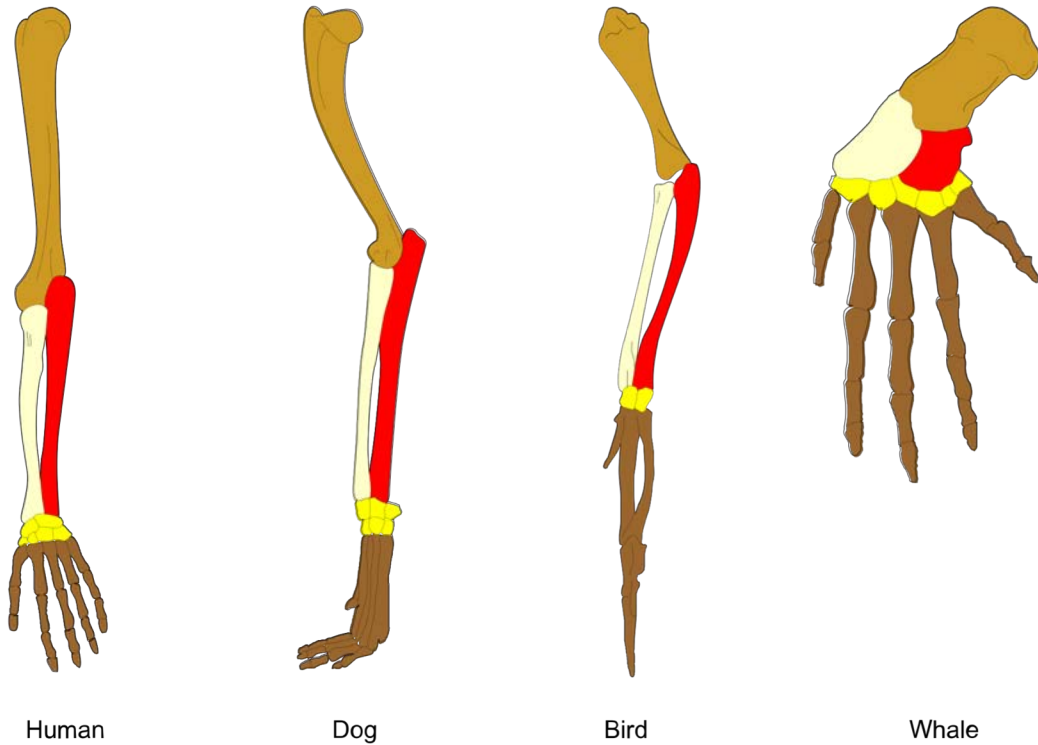
# Homology

# Termed before Darwin's time!



**Sir Richard Owen** [KCB](#) [FRS](#) (20 July 1804 – 18 December 1892) was an English [biologist](#), [comparative anatomist](#) and [paleontologist](#).

# Homology



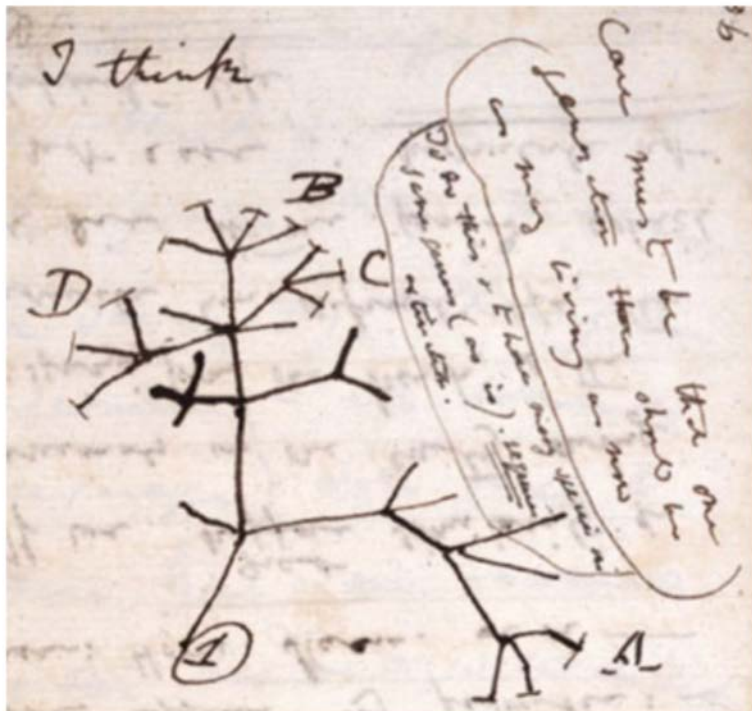
**“the same organ in different animals under every variety of form and function” – Richard Owen**

Owen 1843, p.379

[https://en.wikipedia.org/wiki/Homology\\_\(biology\)](https://en.wikipedia.org/wiki/Homology_(biology))



# Darwin later reformulated homology as a result of “descent with modification”



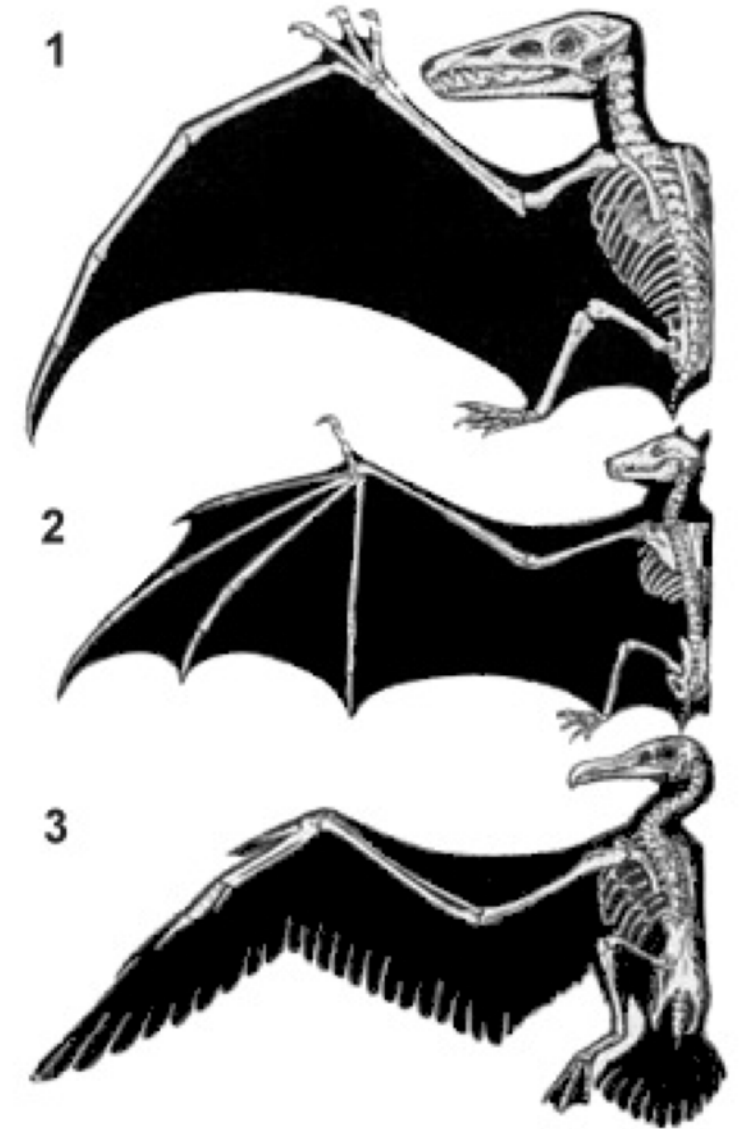
CHAPTER VI.  
DIFFICULTIES ON THEORY.  
Difficulties on the theory of descent with modification—Transitions—Absence or rarity of transitional varieties—Transitions in habits of life—Diversified habits in the same species—Species with habits widely different from those of their allies—Organs of extreme perfection—Means of transition—Cases of difficulty—Natura non facit saltum—Organs of small importance—Organs not in all cases absolutely perfect—The law of Unity of Type and of the Conditions of Existence embraced by the theory of Natural Selection, . . . . . 154

CHAPTER XIII.  
MUTUAL AFFINITIES OF ORGANIC BEINGS: MORPHOLOGY: EMBRYOLOGY: RUDIMENTARY ORGANS.  
CLASSIFICATION, groups subordinate to groups—Natural system—Rules and difficulties in classification, explained on the theory of descent with modification—Classi-

# Homology

The wings of pterosaur (1), bats(2) and birds (3) are **analogous** as wings, but **homologous** as forelimbs.

**Homologs** (any features: genes, trait, morphology) share **ancestry**



# Homology

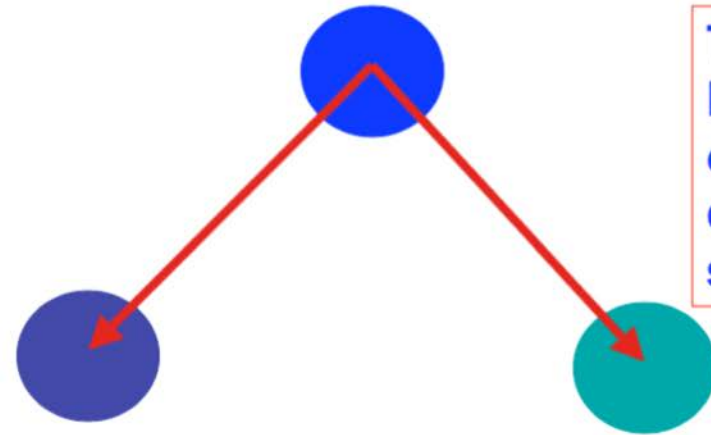
Question: How do we establish homology at sequence level?

Search for similarity , collinearity, conservation of morphological characters

## Search for similarity

One of the most frequent activity in Bioinformatics

Common ancestor



Two genes are homologs if and only if they derive from the same ancestor

Gene1

Gene2

Homology is almost uniquely inferred by sequence similarity

# Beware ; why?

~~Significant homology~~

55% married?  
45% grandmom?

~~Weak homology~~

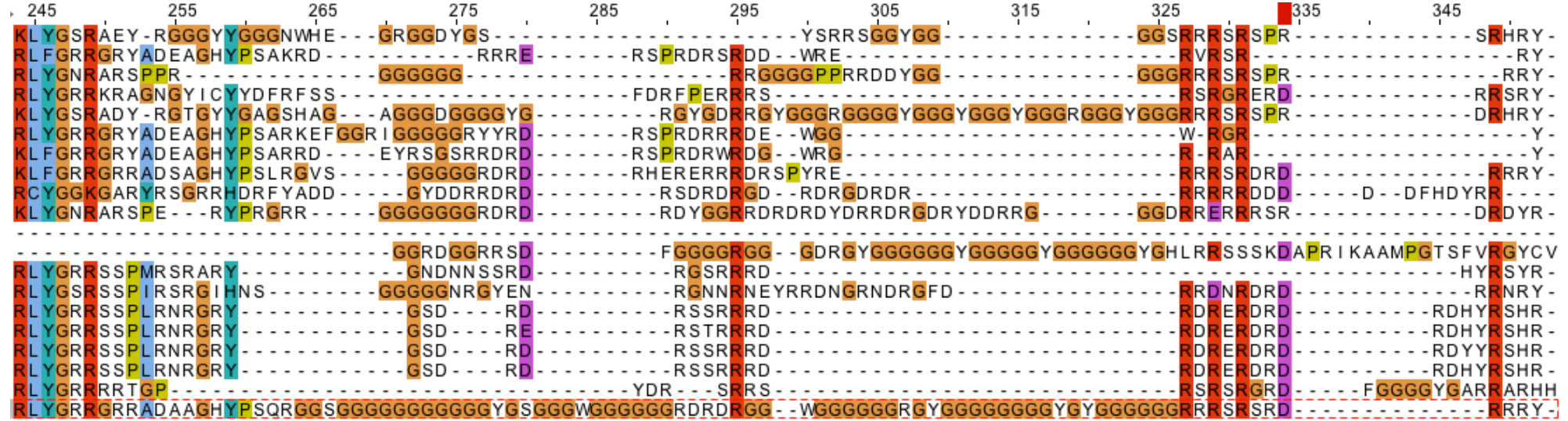
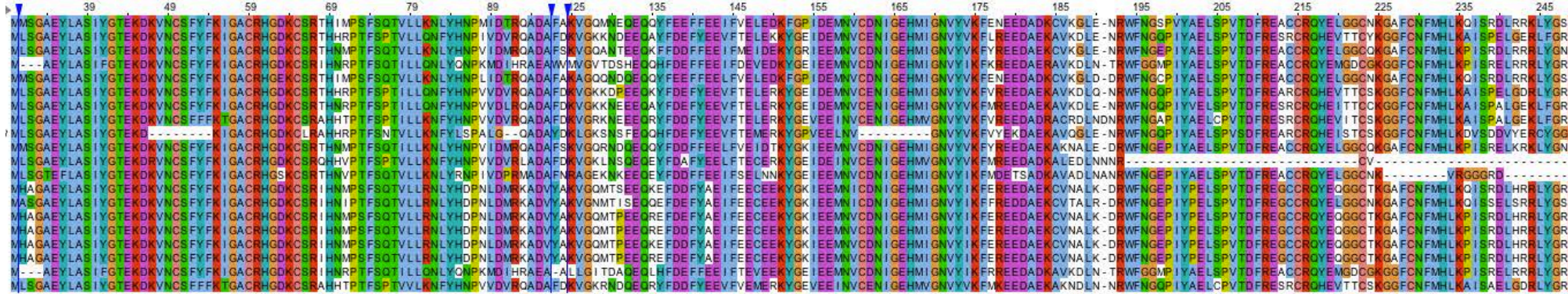
**If you think about the meaning of homology,  
then it really makes no sense**

Significant similarity

Weak similarity

# Extension of homology to sequences

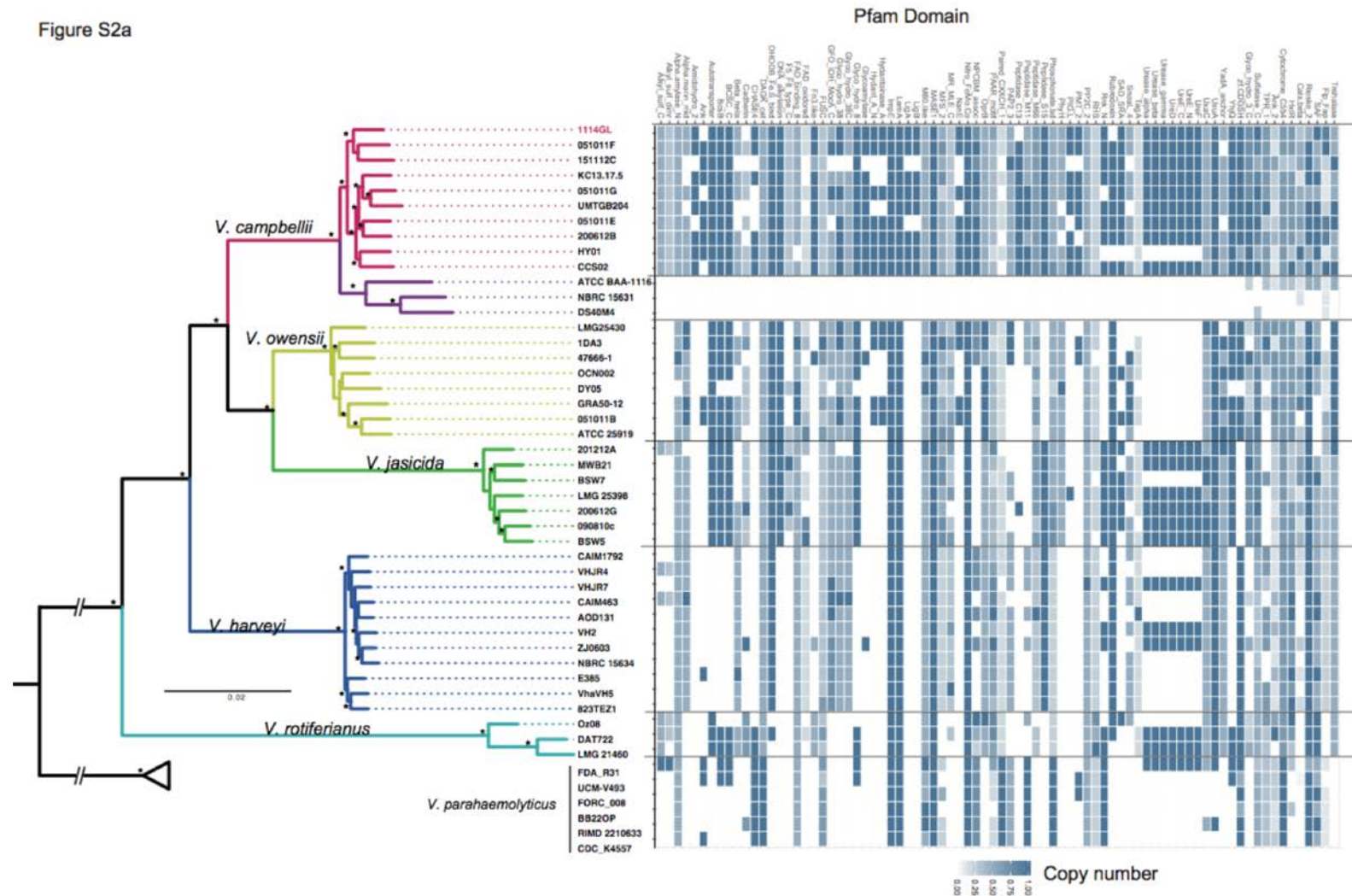
*Two sequences are homologous if they share the same a common ancestor*



# Extension of homology to genomes / species

Similarity of individual sequences at different levels (sequence similarity ; domain combinations)

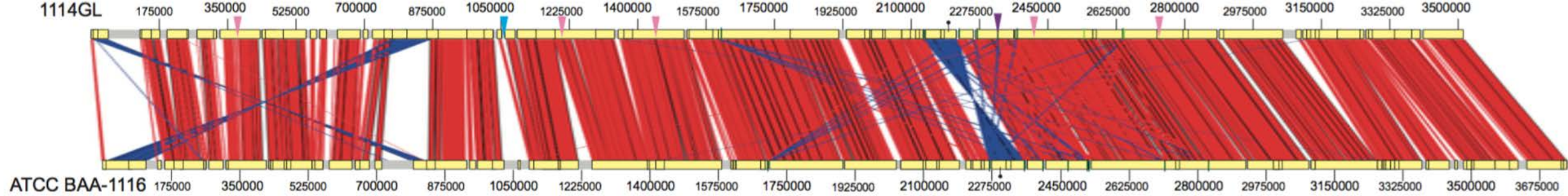
Figure S2a



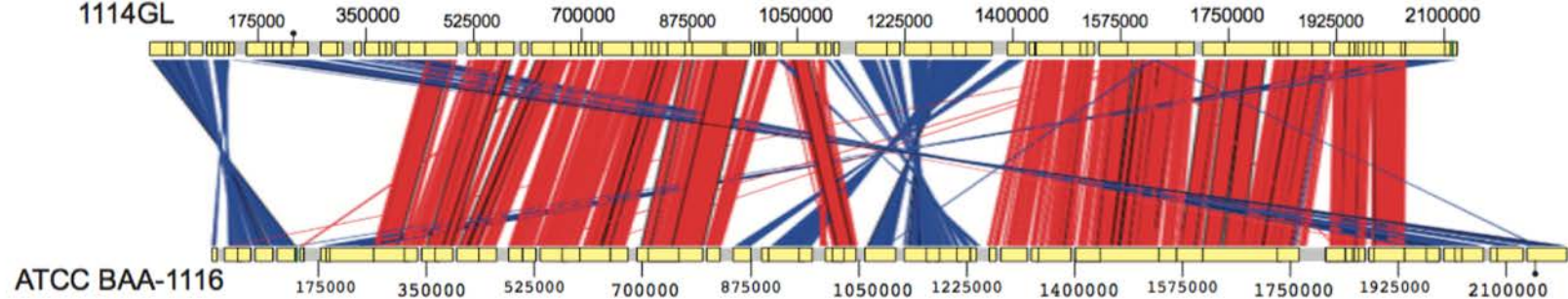
# Extension of homology to genomes / species

Similarity of individual features (ordering and rearrangement)

(a) Chromosome I



(b) Chromosome II



- Gap
- Inter-scaffold gap
- The insertion including two genes with Big\_2 domains
- Ori
- rRNA operon
- Partial rRNA operon



# HOMOLOGY, GENES, AND EVOLUTIONARY INNOVATION



GÜNTER P. WAGNER

Günter Wagner has thought long and hard about homology in relation to character identity, and in his new book he goes into great detail about why we should use **character identity as the basis for the homology of morphological characters**. For readers of *Systematic Biology*, the book is also a reminder that every **morphological character used in a phylogenetic analysis is a hypothesis of homology, and that great care is needed when deciding whether morphological characters in different organisms are likely to be homologs**.

...He also writes that “This book, although ostensibly about homology, is really a book on evolutionary developmental biology” (p. 3). Wagner argues that “the origin of novel characters and novel body plans is one of the most important but least researched questions in evolutionary biology” (p. 3)....

# Why comparative genomics? – A summary

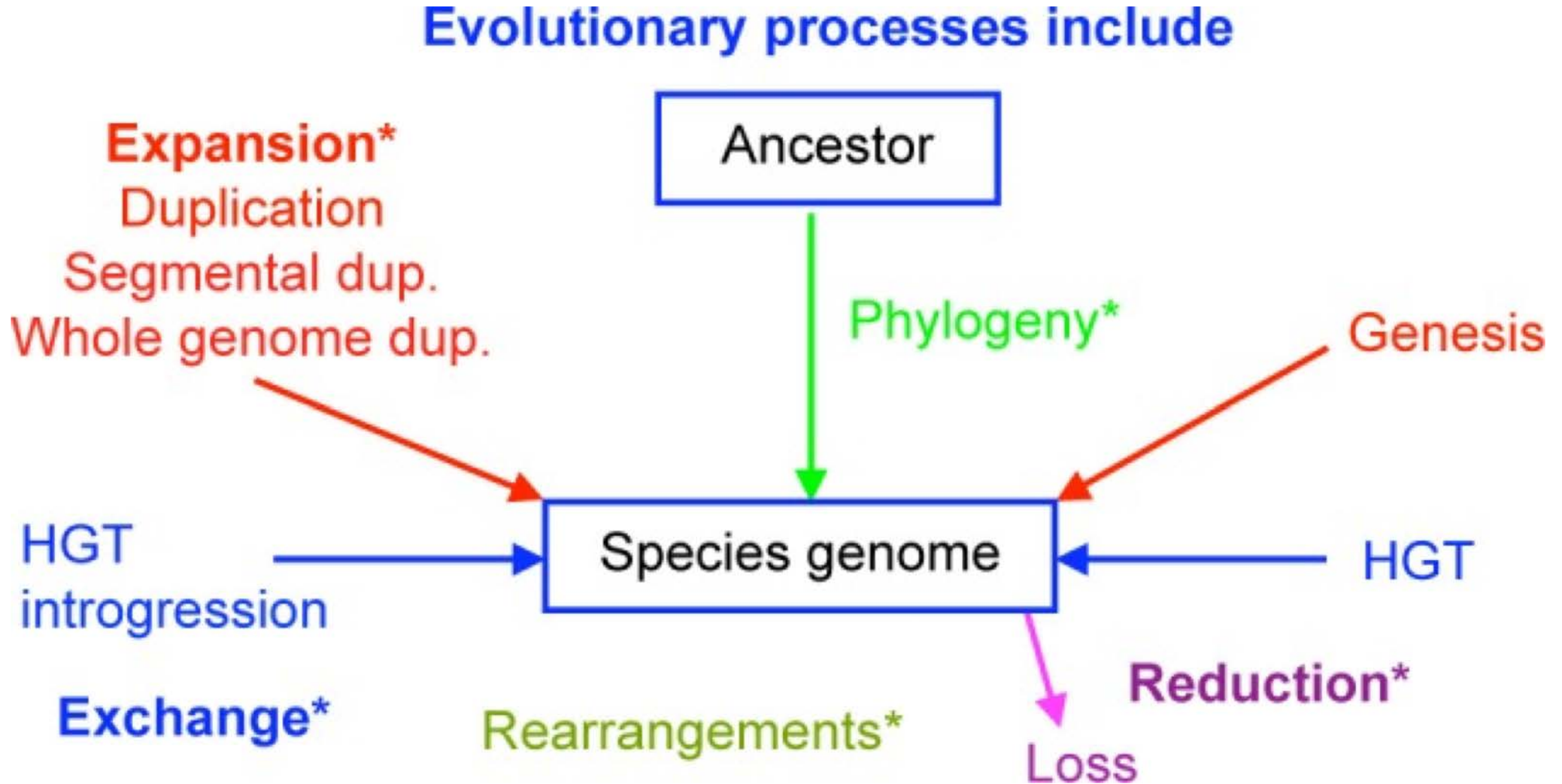
Compare multiple genomes now a norm

Similarity and differences between genomes

Use genomes to study evolution of these species:

- At various resolution (whole genome, chromosomes, regions, genes, base pairs)
- Identify the genomic basis of key phenotypes

# Evolution process of a genome

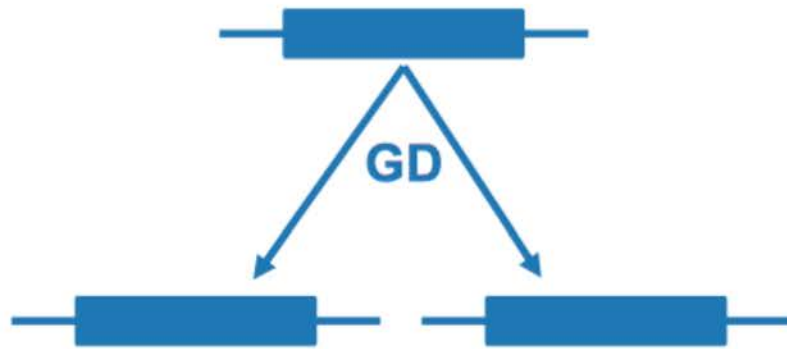


# Sources of gene innovation

(Intuitive as genome gain genes of new functions)

## Gene duplication (GD)

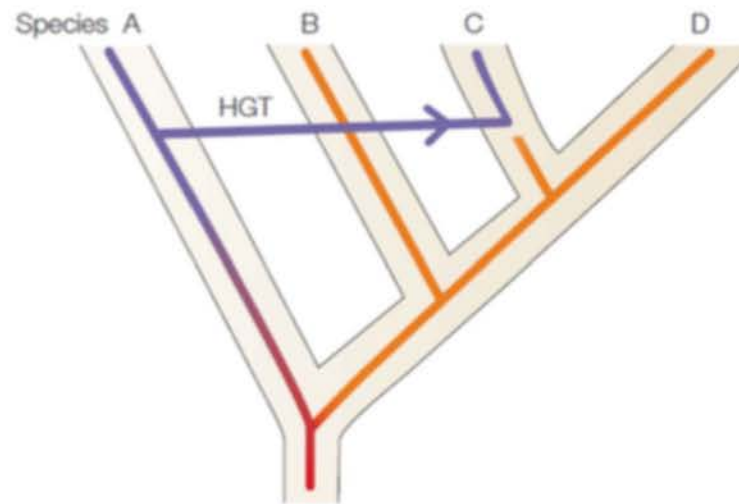
Any duplication of a region of DNA that contains a gene



- ❖ Plant organic material decay
- ❖ Starch catabolism
- ❖ Degradation of host tissues
- ❖ Toxin production

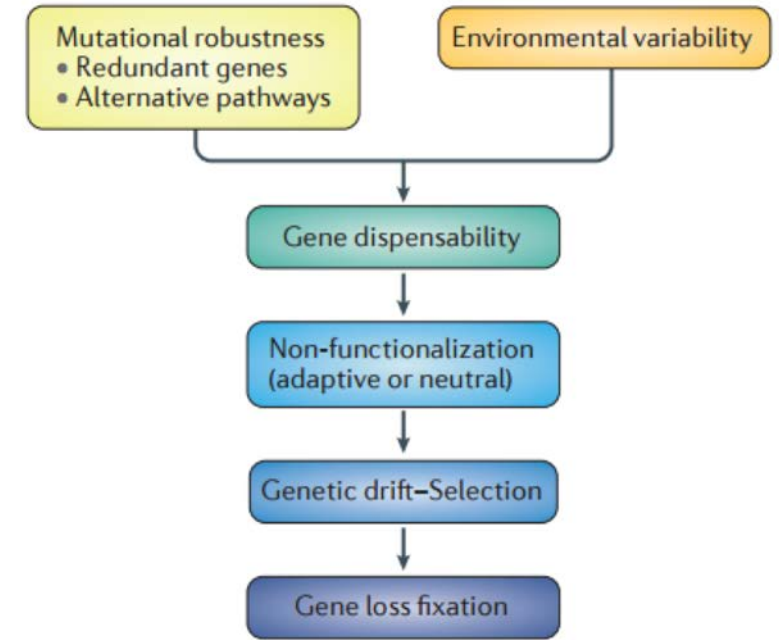
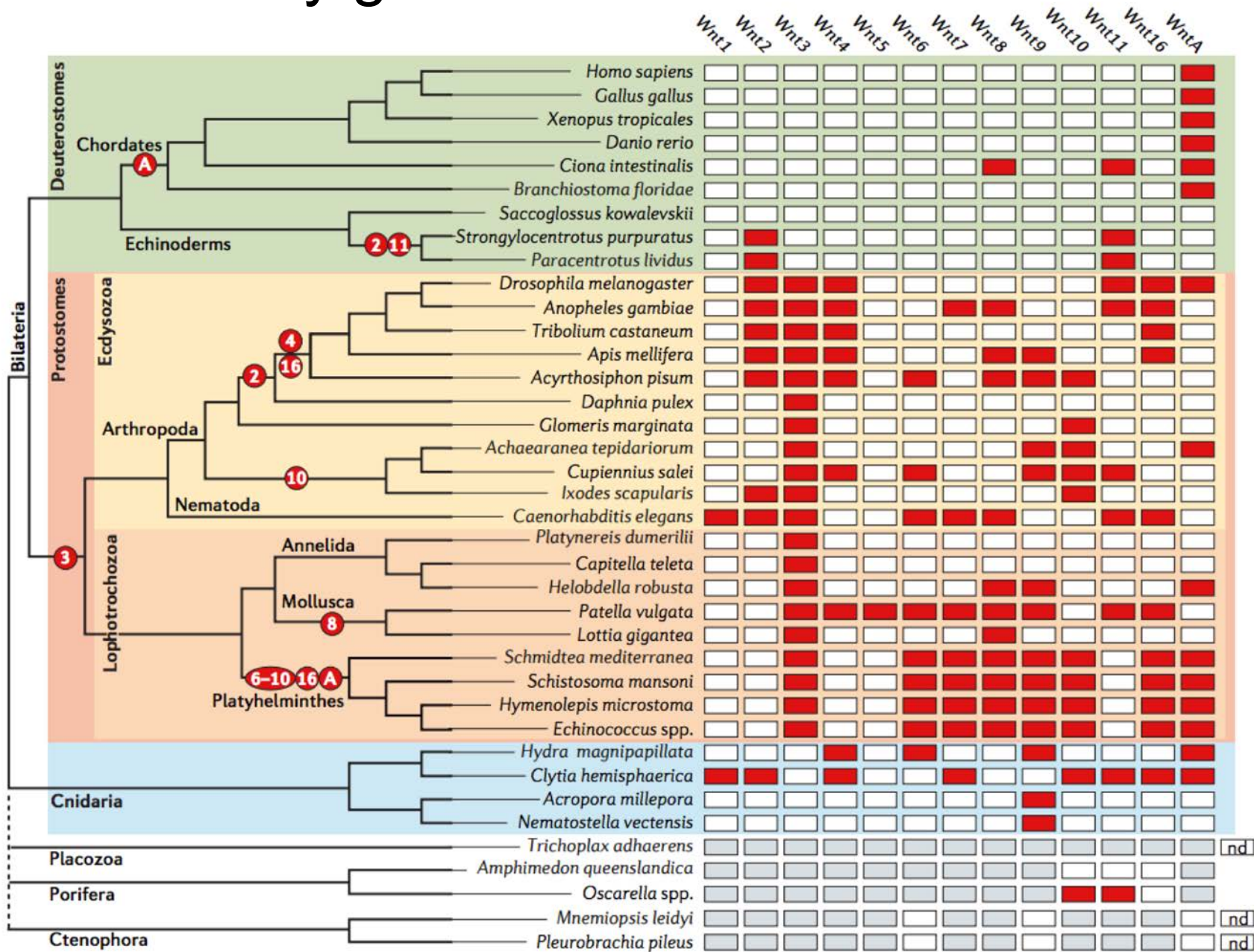
## Horizontal gene transfer (HGT)

Exchange of genes between organisms other than through reproduction

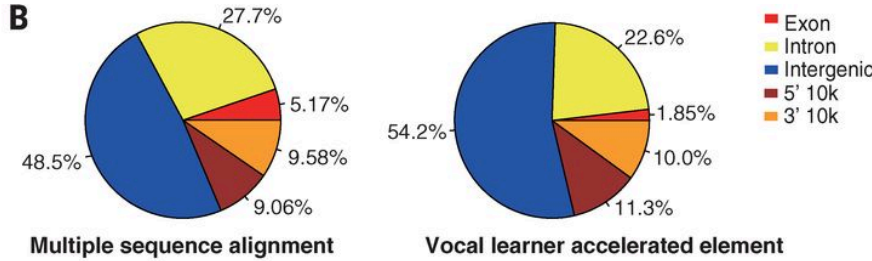
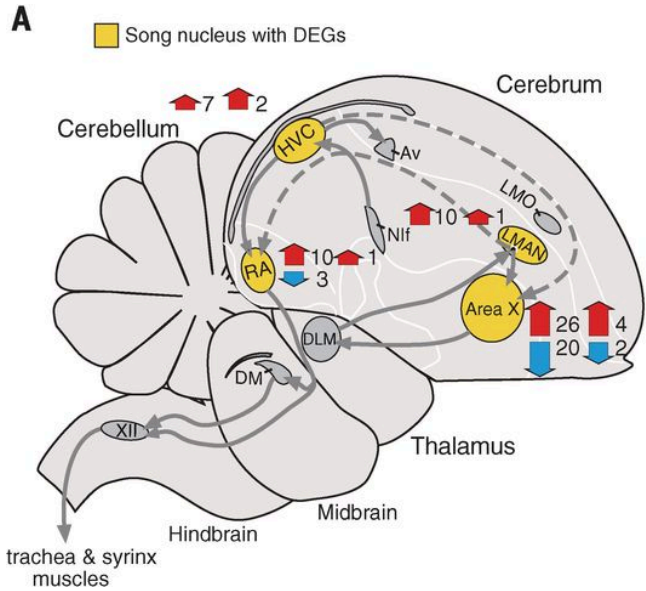
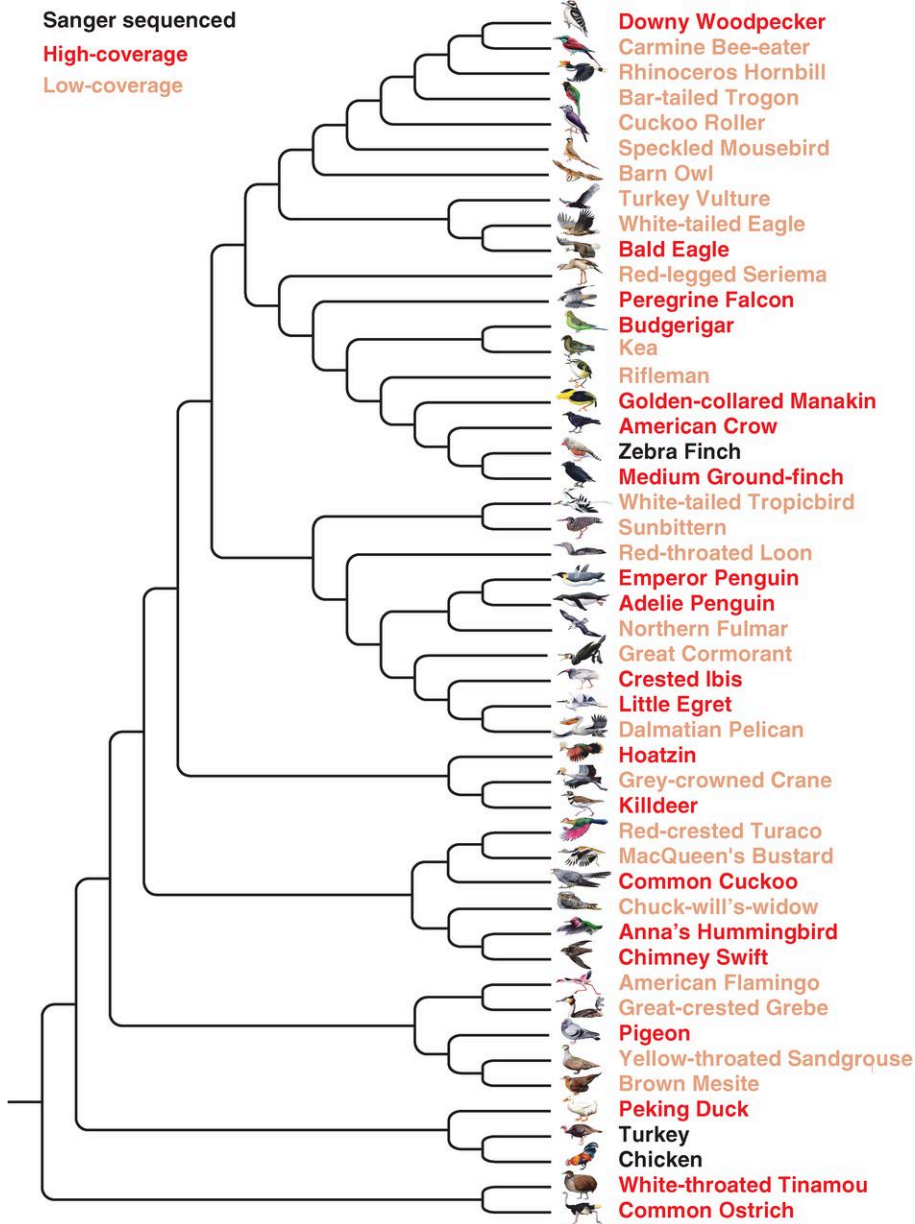


- ❖ Xenobiotic catabolism
- ❖ Toxin production
- ❖ Degradation of plant cell walls
- ❖ Wine fermentation

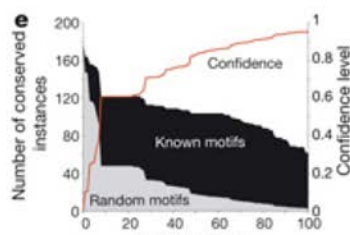
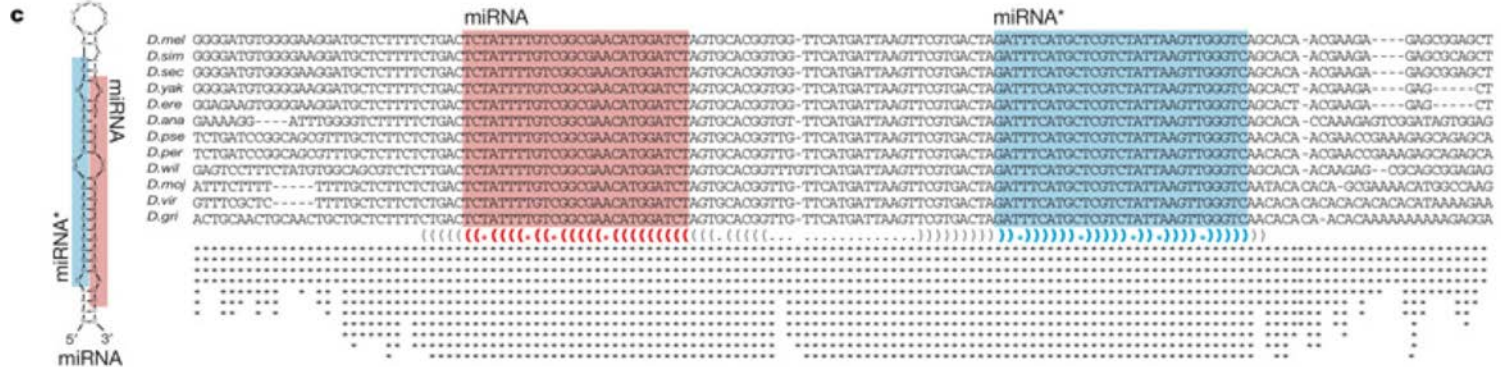
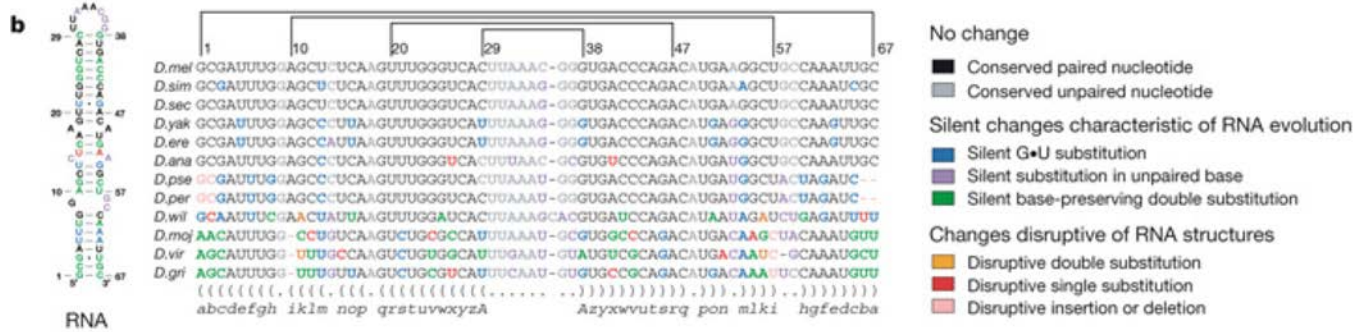
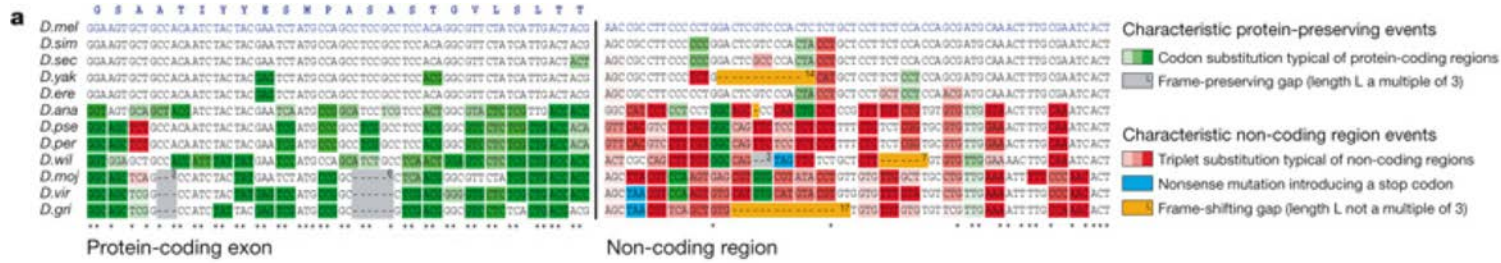
# Evolution by gene loss



# Reveal the evolutionary relationships among species

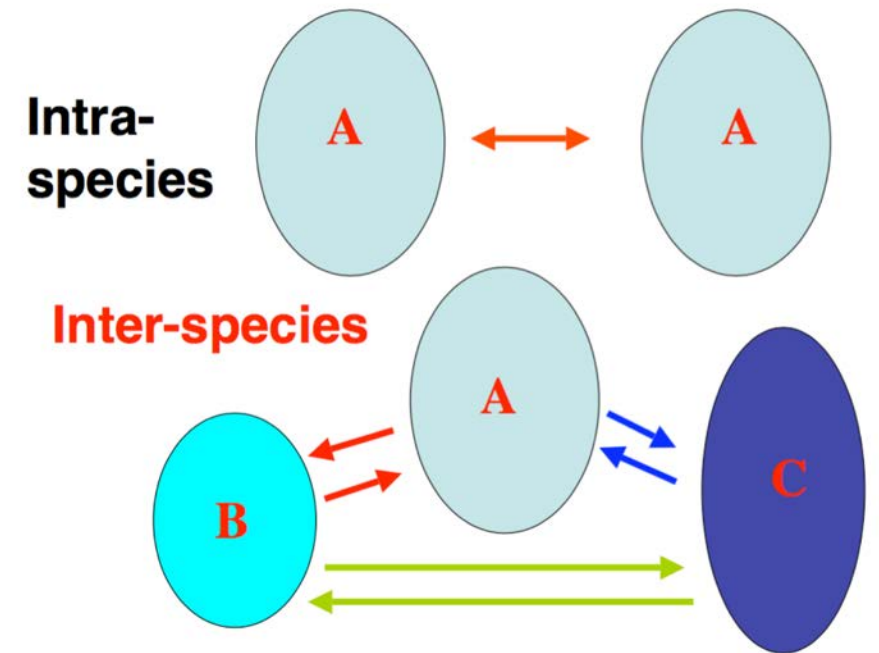


# Link evolutionary processes with function



# Comparing genomes

- Alignment of homologous regions
  - **Inter-genomic**: aligning genomic sequences from **different** species
  - **Intra-genomic** aligning genomic sequences from the **same** species
- Different levels of **resolution**
  - Comparative mapping (markers)
  - Synteny (~ gene content)
  - Colinearity (gene content + order conservation)
  - DNA-based alignments (base-to-base mapping)





Orthology

# Refining *how* homologous genes are related

## DISTINGUISHING HOMOLOGOUS FROM ANALOGOUS PROTEINS (1970)

WALTER M. FITCH



1929 - 2011

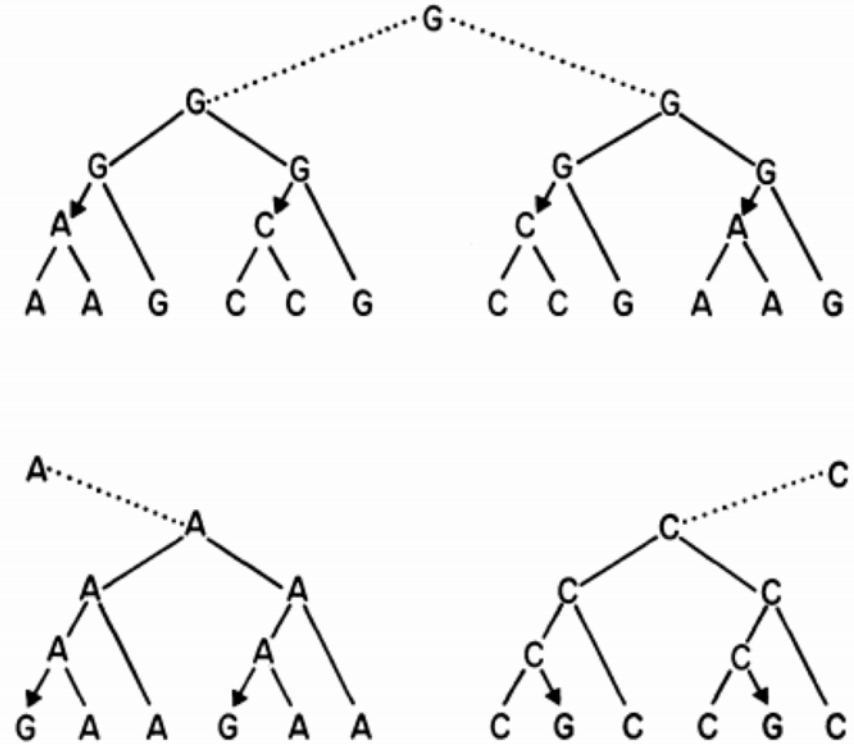


FIG. 1.—Distinguishing convergent from divergent types of nucleotide replacement patterns. Given are two groups of species (related within each group as shown by the solid lines) together with the nucleotide present at a specific position of the gene for each member species as shown at the branch tips. Given also the requirement that the ancestral nucleotide must permit the descendant nucleotides to be obtained in the minimum number of replacements, the ancestral nucleotide of the upper two groups must be set as G, with the required replacements indicated by the arrows. Were one to postulate a common ancestor for the two groups, no new mutations would need to be assumed; hence, this kind of pattern is called the divergent types. The lower two groups are identical except for rearranging the nucleotides at the branch tips, but now, in order to account for descendants in only four nucleotide replacements, the ancestral nucleotide of the lower two groups must be A and C. To postulate a common ancestor for these two groups would require, unlike the upper pair, an additional mutation. This situation shows different ancestral characters apparently converging toward the same descendant character, and hence is called the convergent type. One can calculate the frequency with which one might expect each type to be found in examining a large number of such nucleotide positions and compare that value to what is in fact found for a particular set of proteins. An abnormally large number of either type is evidence favoring that type of relation between the two groups examined.

# From homology to orthology

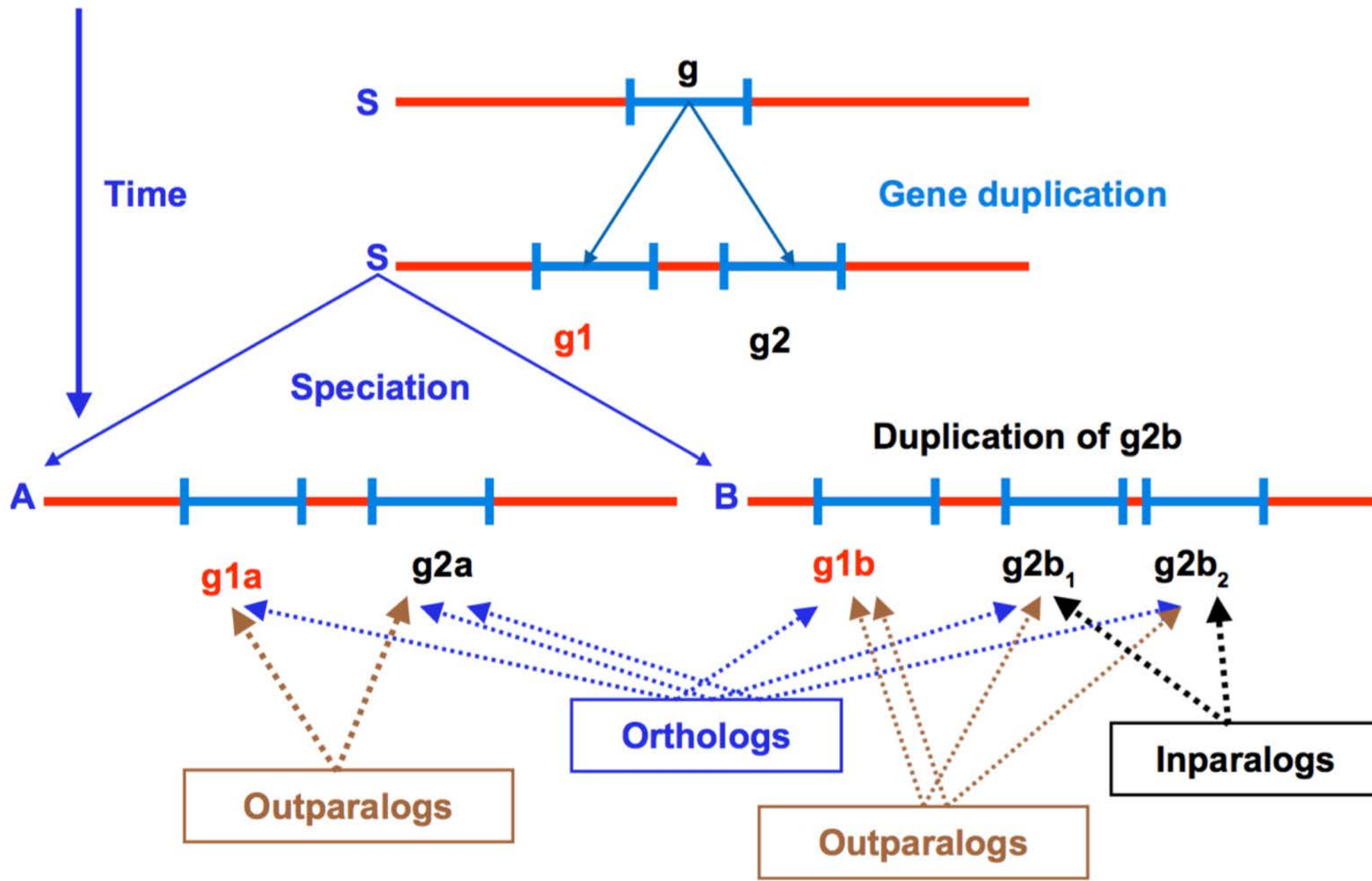
Homologues are sequences derived from a common ancestor...

- What are then orthologues? and paralogues?

Original definition of orthology and paralogy by Walter Fitch (1970, Systematic Zoology 19:99-113):

*"Where the homology is **the result of gene duplication** so that both copies have descended side by side during the history of an organism, (for example, alpha and beta hemoglobin) the genes should be called **paralogous** (para = in parallel).*

*Where the homology is **the result of speciation** so that the history of the gene reflects the history of the species (for example alpha hemoglobin in man and mouse) the genes should be called **orthologous** (ortho = exact)."*



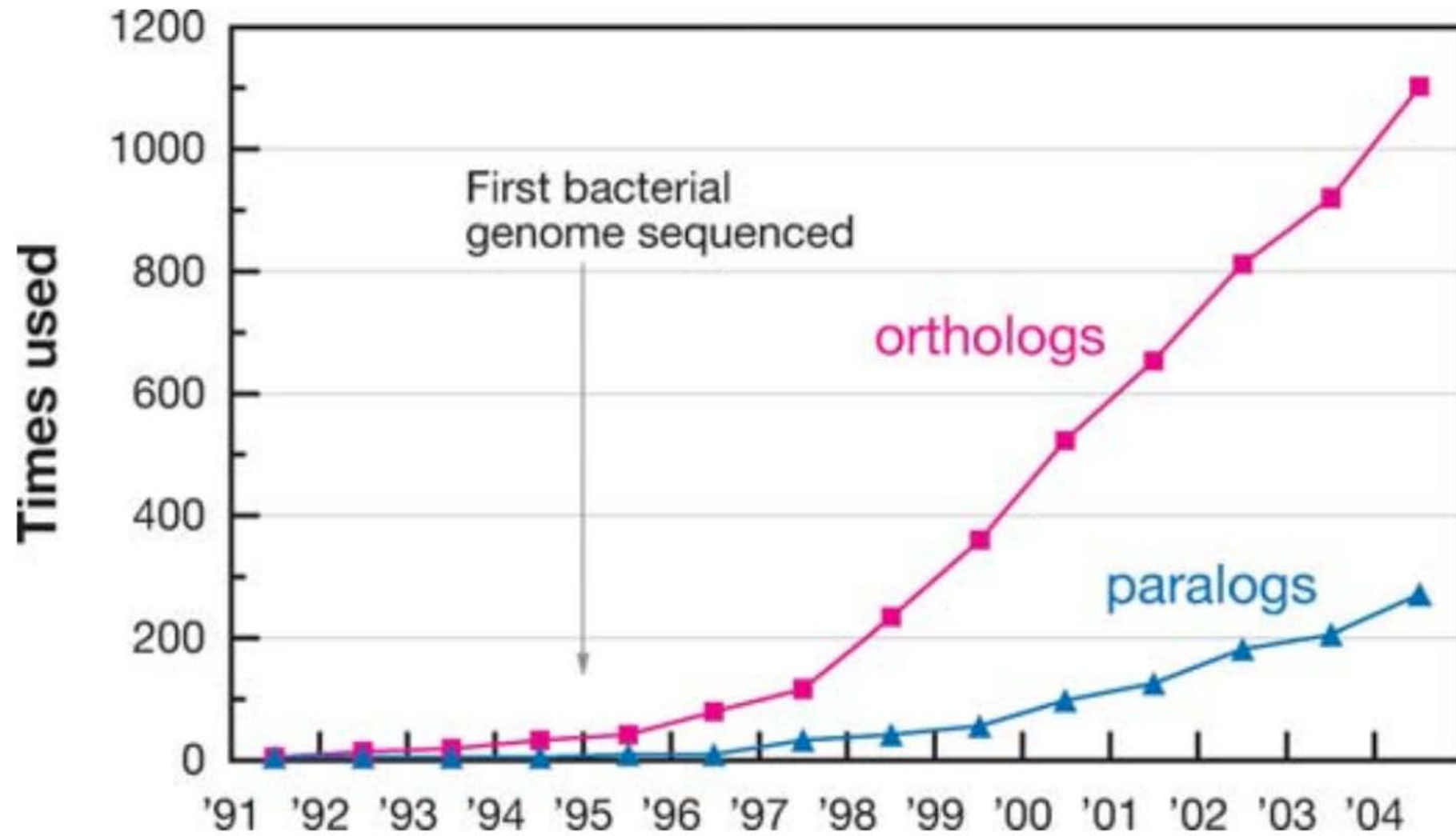
# Why is orthology important?

**Orthologs detection is of fundamental importance in:**

- Reconstruction of the evolution of species and their genomes (Phylogenomics);**
- Evolutionary studies of biological systems;**
- Annotation of newly sequenced organisms;**
- Functional genomics (transfer of functional annotation predicted on “orthology-function conjecture”);**
- Gene organization in a given species.**

**Accurate determination of evolutionary relationships between orthologous gene families is of utmost importance for such goals.**

# Usage of “ortholog” and “paralog”



# Corollary

- Orthology definition is purely on evolutionary terms (not functional, not synteny...)
- There is no limit on the number of orthologs or paralogs that a given gene can have (when more than one ortholog exist, there is nothing such as "*the true ortholog*")
- Many-to-Many orthology relationships do exist (co-orthology)
- No limit on how ancient/recent is the ancestral relationship of orthologs and paralogs
- Orthology is non-transitive (as opposed to homology)

# More precise definitions

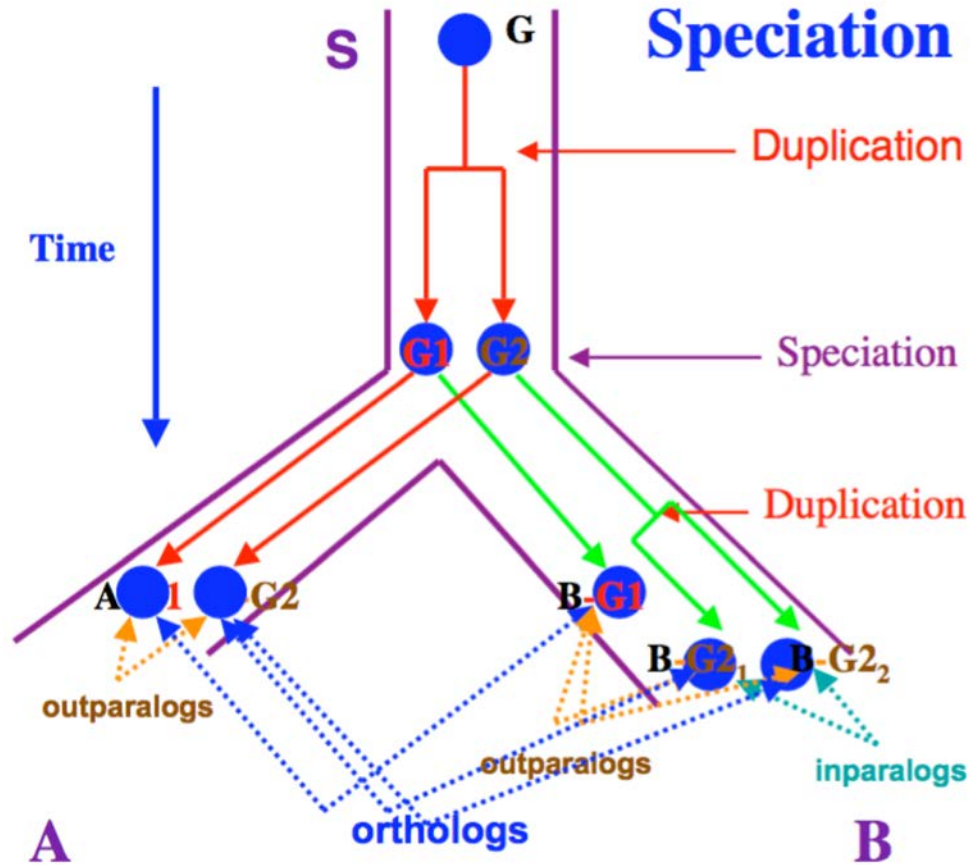


Table 1 Homology: terms and definitions

Homologs		Genes sharing a common origin
Orthologs	Genes originating from a single ancestral gene in the last common ancestor of the compared genomes.	
Pseudoorthologs	Genes that actually are paralogs but appear to be orthologous due to differential, lineage-specific gene loss.	
Xenologs	Homologous genes acquired via XGD by one or both of the compared species but appearing to be orthologous in pairwise genome comparisons.	
Co-orthologs	Two or more genes in one lineage that are, collectively, orthologous to one or more genes in another lineage due to a lineage-specific duplication(s). Members of a co-orthologous gene set are inparalogs relative to the respective speciation event.	
Paralogs		Genes related by duplication
Inparalogs (symparalogs)	Paralogous genes resulting from a lineage-specific duplication(s) subsequent to a given speciation event (defined only relative to a speciation event, no absolute meaning).	
Outparalogs (alloparalogs)	Paralogous genes resulting from a duplication(s) preceding a given speciation event (defined only relative to a speciation event, no absolute meaning).	
Pseudoparalogs	Homologous genes that come out as paralogs in a single-genome analysis but actually ended up in the given genome as a result of a combination of vertical inheritance and HGT.	



# Importance of assigning correct orthology

**Important implications for phylogeny:** only sets of orthologous genes are expected to reflect the underlying species evolution (although there are many exceptions)

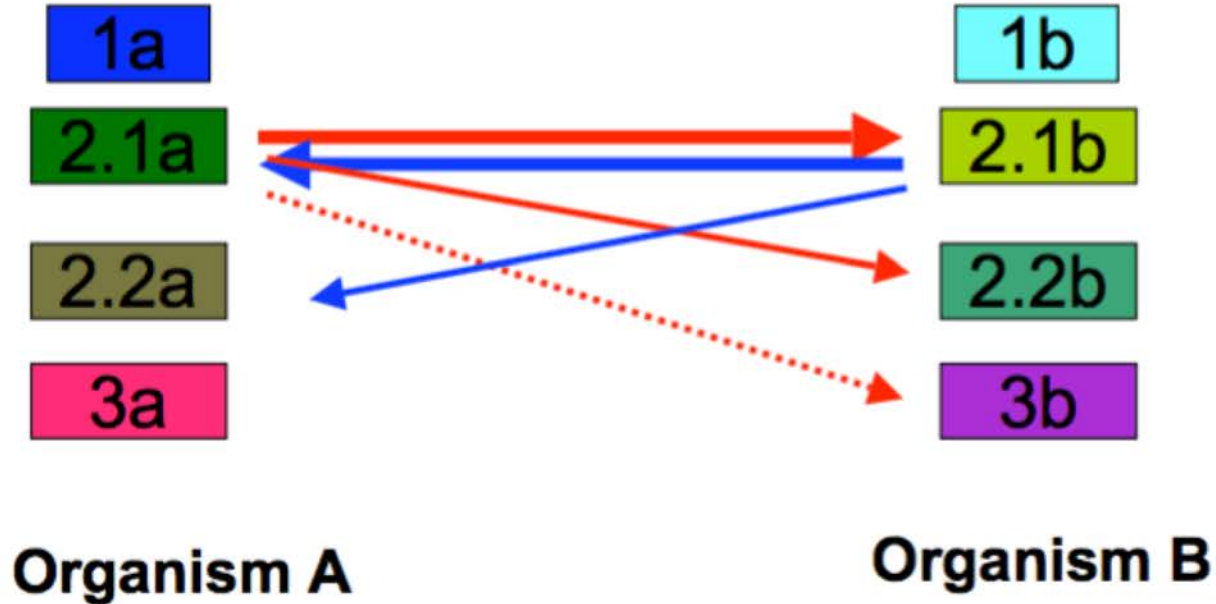
The most exact way of **comparing two (or more) genomes** in terms of their gene content. Necessary to uncover how genomes evolve.

Implications for **functional inference:** orthologs, as compared to paralogs, are more likely to share the same function

# Ortholog inference methods

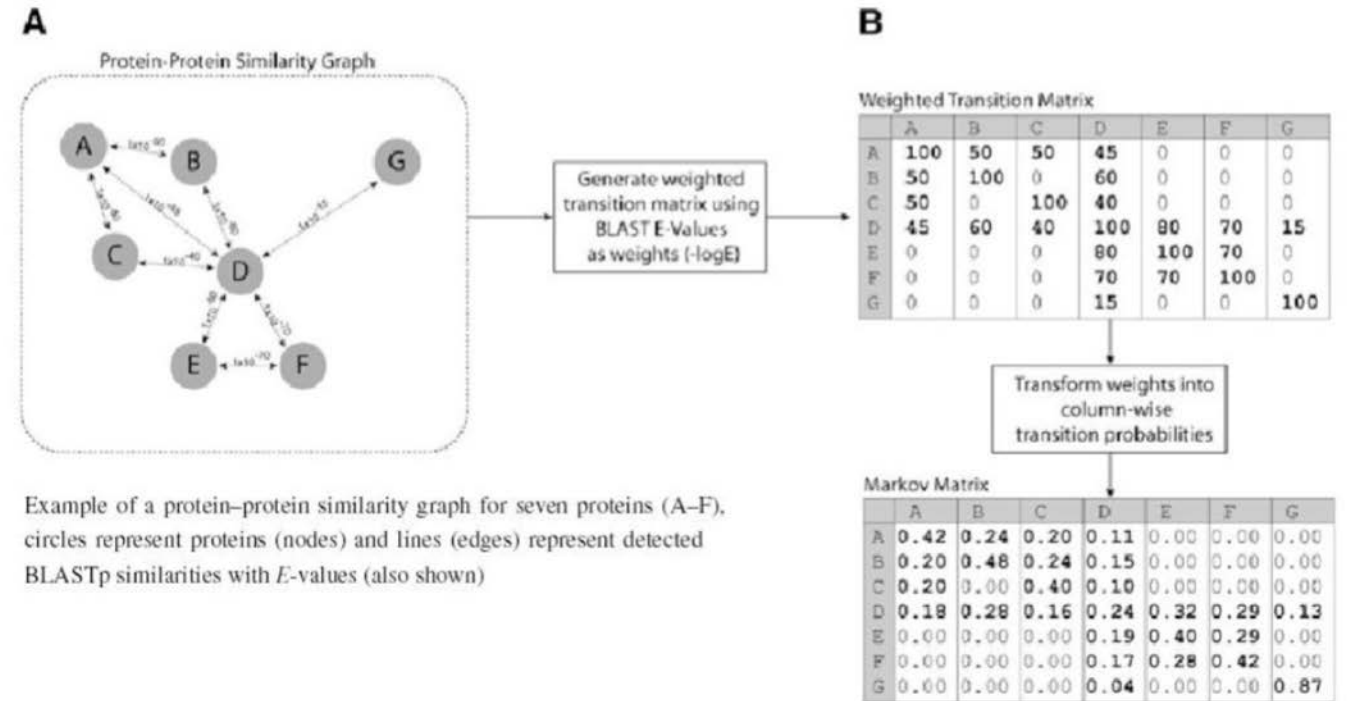
# How to detect orthologous genes?

- The most intuitive way: **Best Reciprocal Hit (RBH)**

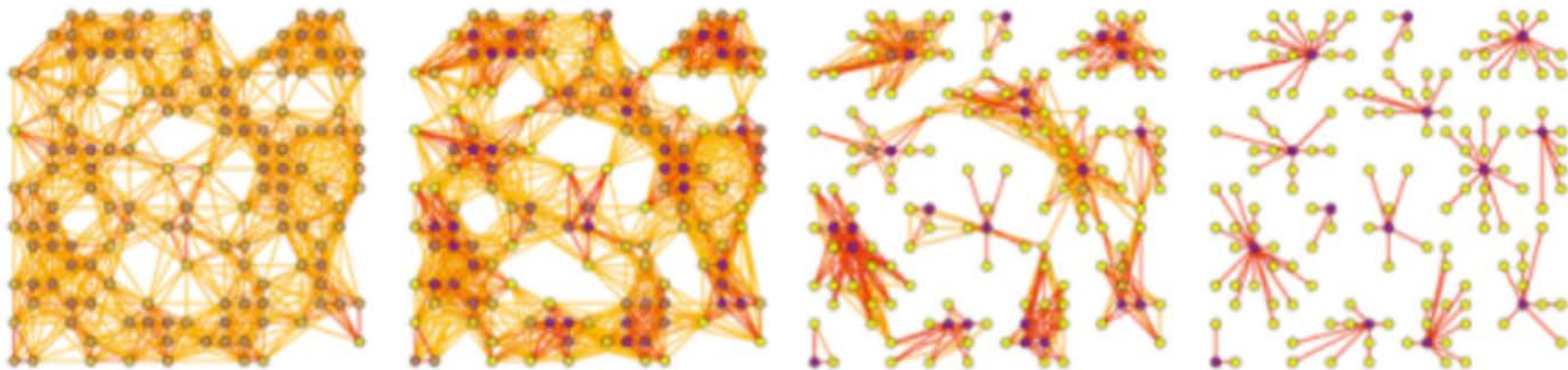


# Sequence by clustering

## mcl: The Markov Cluster Algorithm <http://micans.org/mcl/> (Stijn Van Dongen)



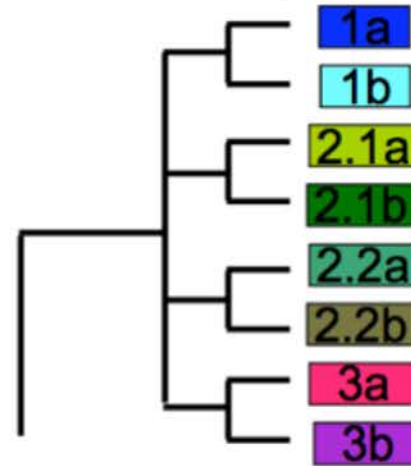
Produce clusters (gene families) using different inflation parameter



Weighted transition matrix and associated column stochastic Markov matrix for the seven proteins shown in (A).

# How to detect orthologous genes?

- more rigorous: make a phylogenetic tree of the gene family



- more rigorous: look at synteny conservation



--> In fact inferring orthology is much more complicated particularly when considering more than 2 genomes!

# Tree reconciliation

Detection of speciation and duplication events using a species tree and gene family tree

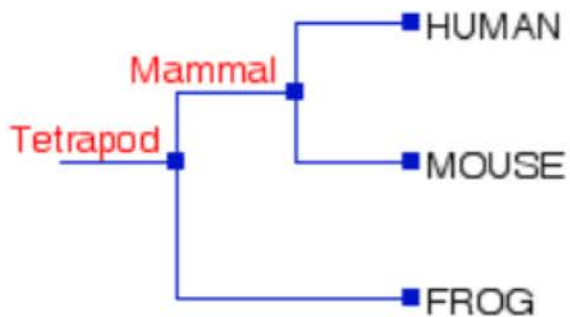
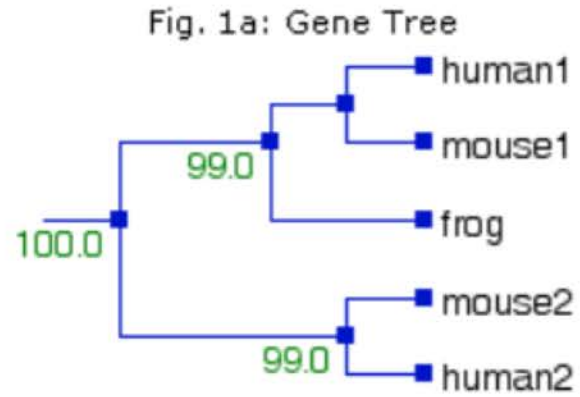


Fig. 1b: Species Tree

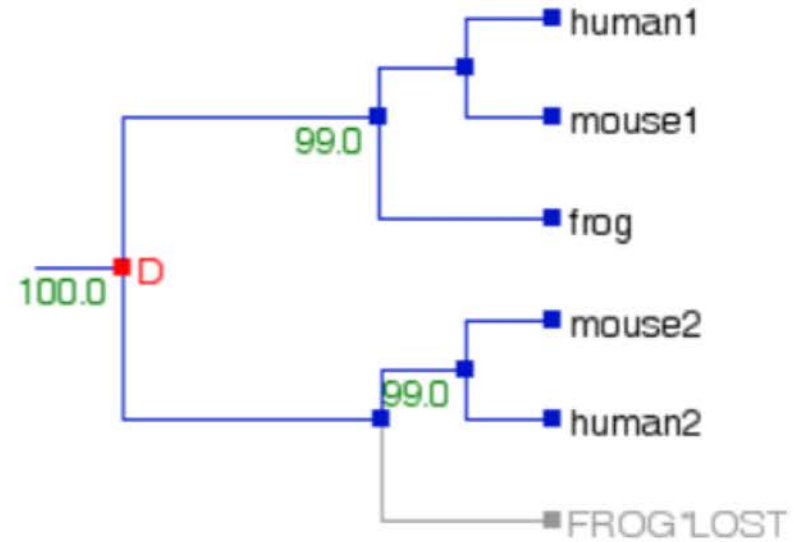
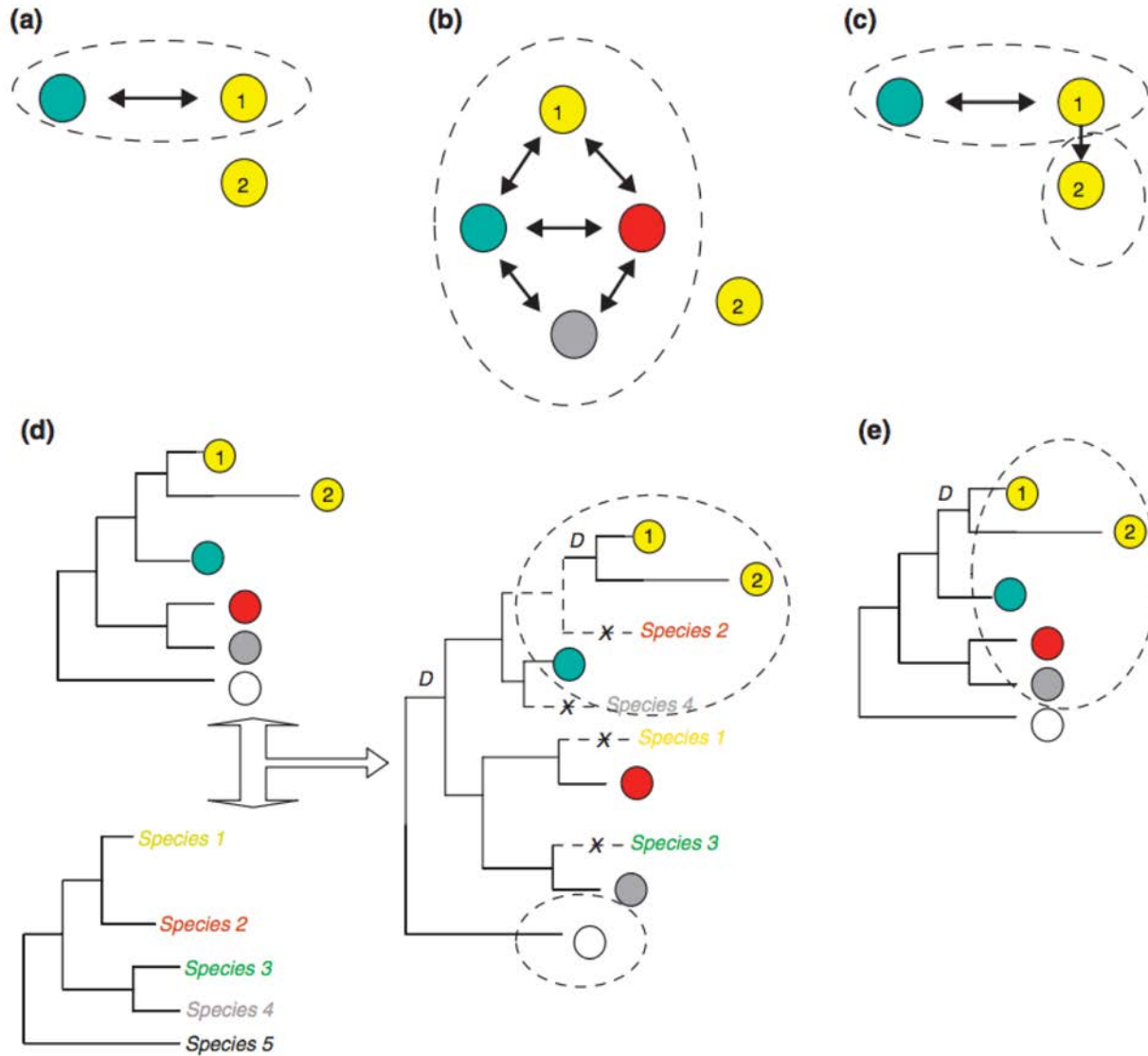


Fig. 1c: Reconciled Gene Tree

# Orthology prediction methods



- a) Best bidirectional hits
- b) COG, MCL-clustering approach
- c) InParanoid
- d) Tree reconciliation
- e) Species-overlap (PhylomeDB)

# Methods

## Similarity

Rely on genome comparisons and clustering of highly similar genes to identify orthologous groups **(suitable for large genome datasets)**

## Phylogeny

use candidate gene families determined by similarity and then rely on the reconciliation of the phylogeny of these genes with their corresponding species phylogeny to determine the subset of orthologs

**(Good and more interpretable for small set of genomes)**

## Others

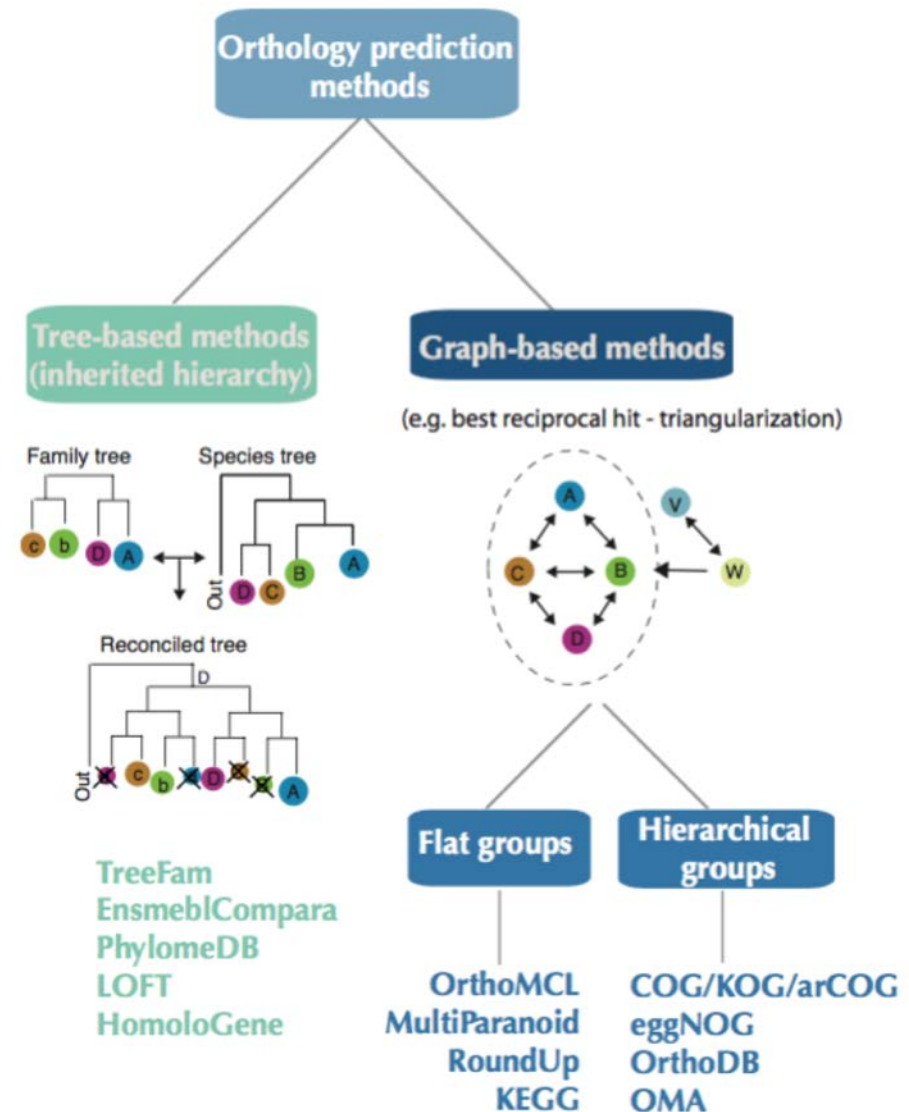
Combination of (1) and (2)

Some uses synteny



# Tools

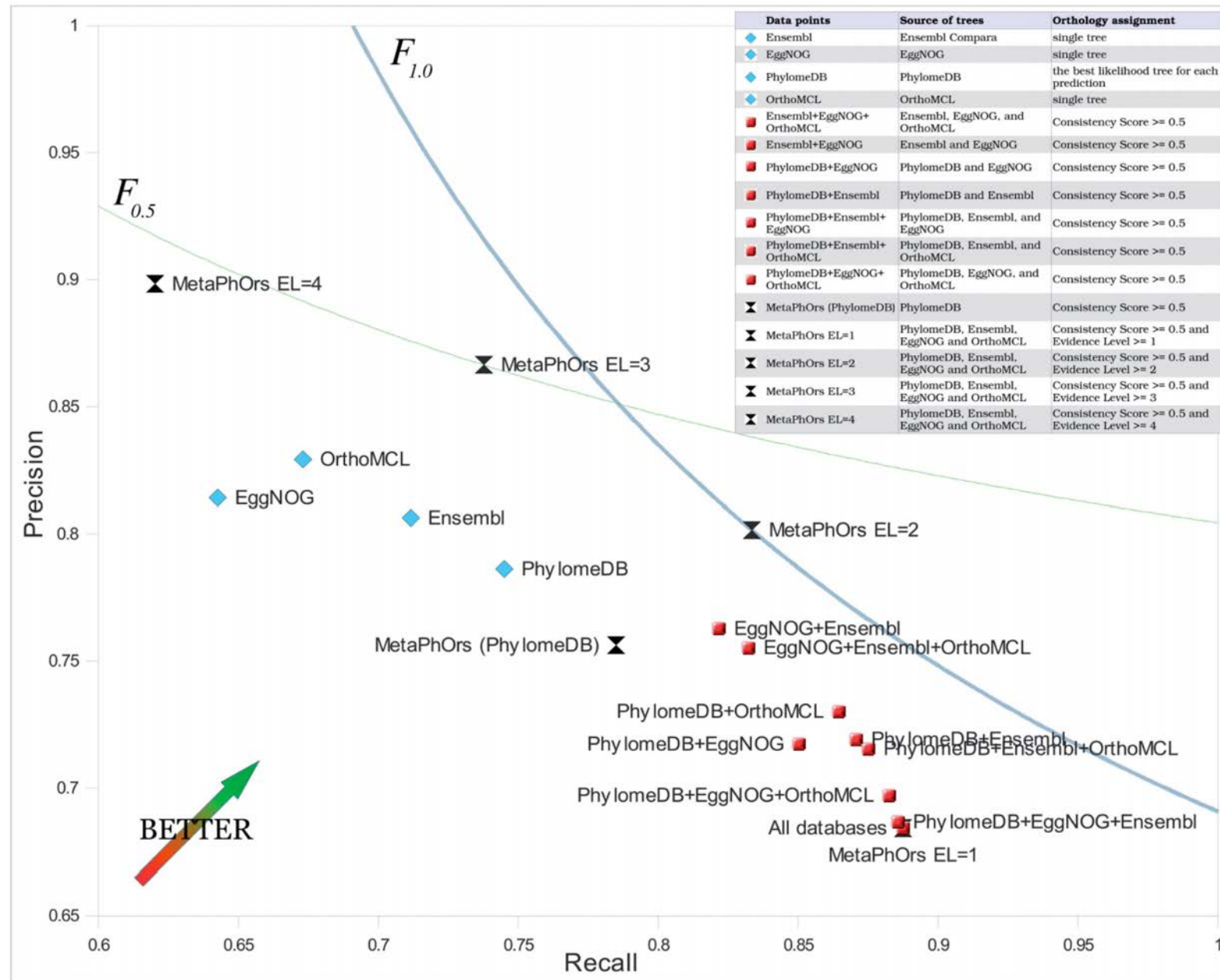
Method	Type	Comments	Reference
BUSCO	Graph	Based on precomputed “universal single-copy” genes (defined for a number of standard clades), and thus inherently limited to these. Originally developed to assess genome completeness.	(Waterhouse et al., 2017)
COG/KOG	Graph	One of the first methods, still widely used for prokaryotic data. Includes a manual curation step.	(Tatusov et al., 2003)
EggNOG	Hybrid	Originally developed as extension of COG/KOG. Recent versions also include tree-based refinements.	(Huerta-Cepas et al., 2016b)
ETE 3.0	Tree	General purpose tree analysis and visualisation package for Python, with species overlap function.	(Huerta-Cepas et al., 2016a)
Forester	Tree	General purpose tree analysis and visualisation software, including reconciliation function.	(Zmasek and Eddy, 2001)
GIGA	Tree	Gene/species tree reconciliation algorithm used in the PANTHER database. Also includes a heuristic for lateral gene transfer detection.	(Thomas, 2010)
GSR	Tree	Probabilistic gene/species tree reconciliation method	(Akerborg et al., 2009)
HaMSTR	Graph	The method uses a reference species to define one Hidden Markov Model per orthologous group, followed by reciprocal best hit within a family	(Ebersberger et al., 2009)
Hieranoid	Graph	Successor of Inparanoid to infer hierarchical orthologous groups from multiple species	(Kaduk et al., 2017)
Inparanoid	Graph	Infers orthologous groups independently for each pair of species.	(Sonhammer and Östlund, 2015)
MetaPhOrs	Hybrid	Meta-method integrating predictions from multiple sources.	(Pryszcz et al., 2011)
Notung	Tree	Gene/species tree reconciliation software, with optional support for lateral gene transfer inference.	(Chen et al., 2000)
OMA	Graph	Infers both types of groups reviewed in this chapter: strict groups (suitable as markers for species tree inference) and hierarchical orthologous groups.	(Altenhoff et al., 2018a)
OrthoDB	Graph	Infers hierarchical orthologous groups. Used to infer the single-copy universal gene models of BUSCO.	(Zdobnov et al., 2017)
OrthoFinder	Graph	Infers hierarchical orthologous group with respect to the deepest speciation level only (the last common	(Emms and Kelly, 2015)



Tekaia (2016)

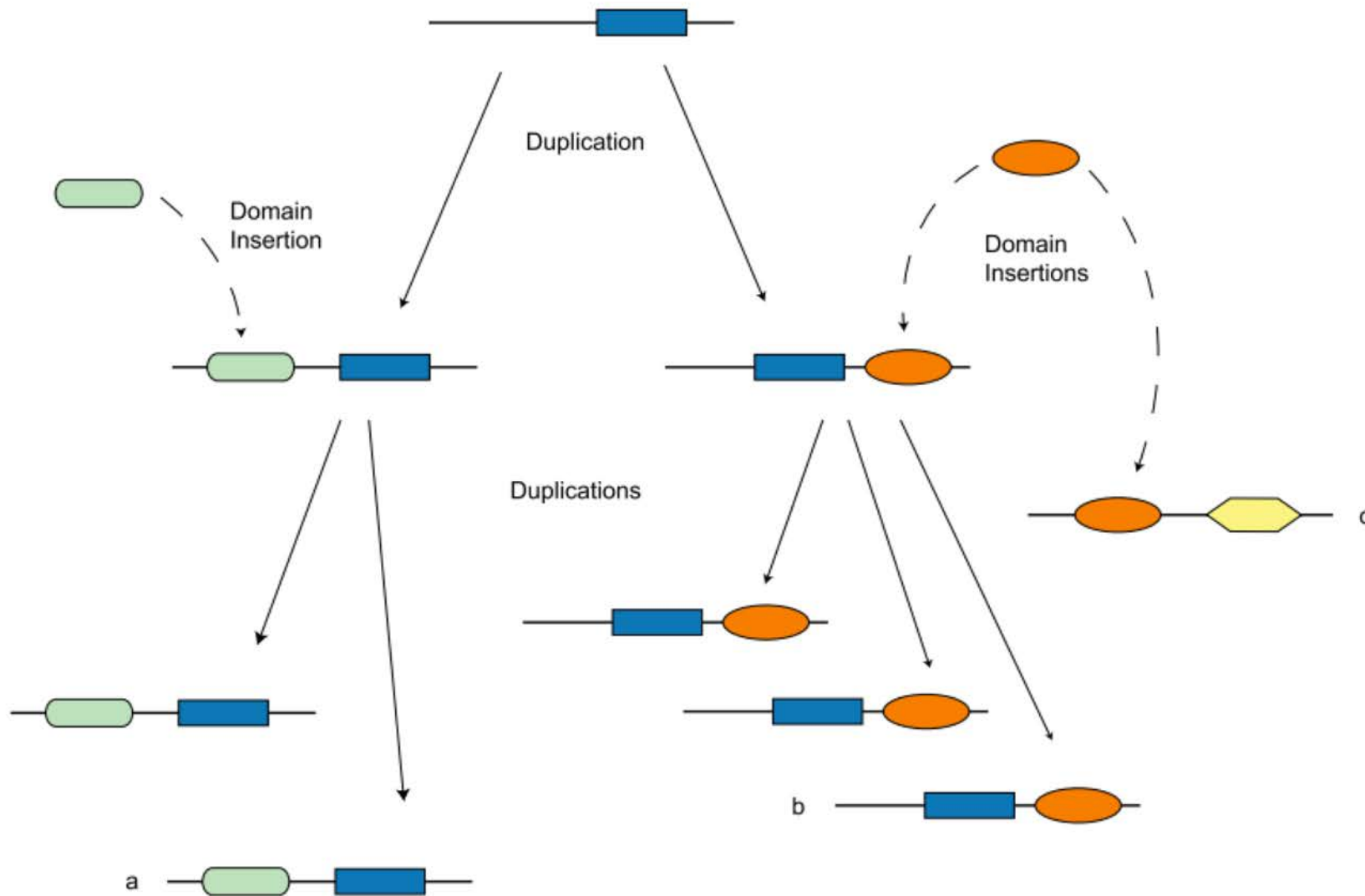
Fernández *et al* (2019)

# Every tool kind of disagrees...



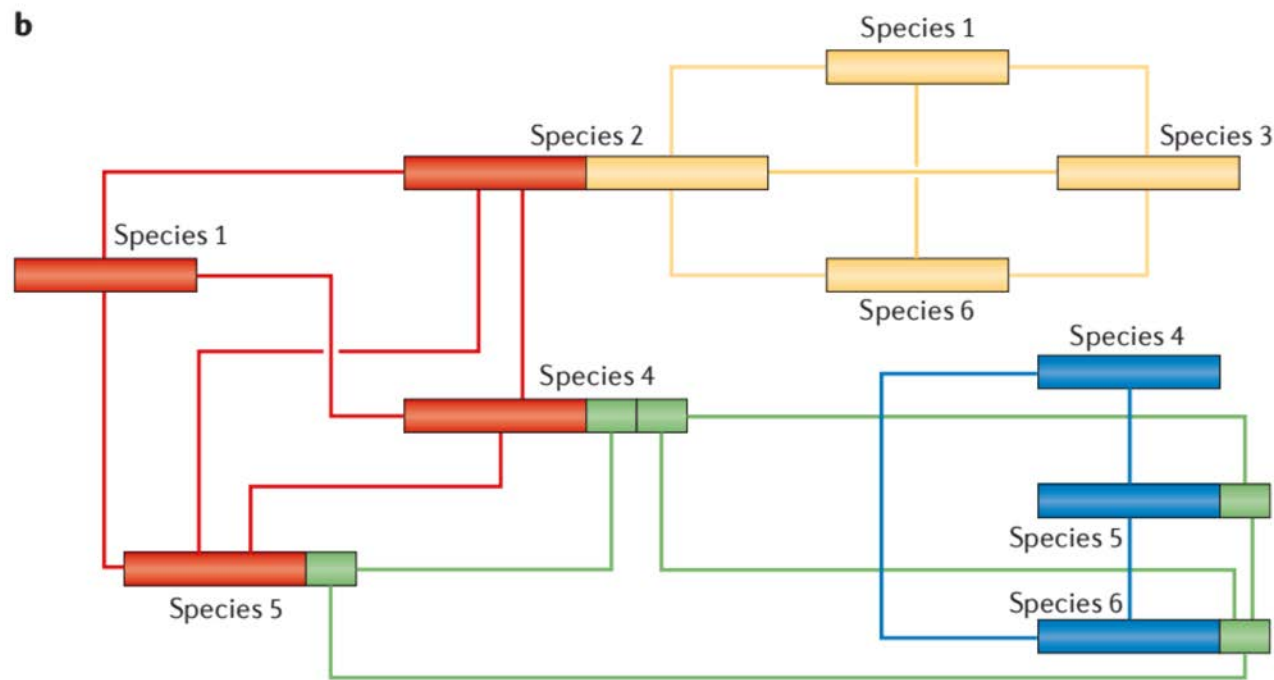
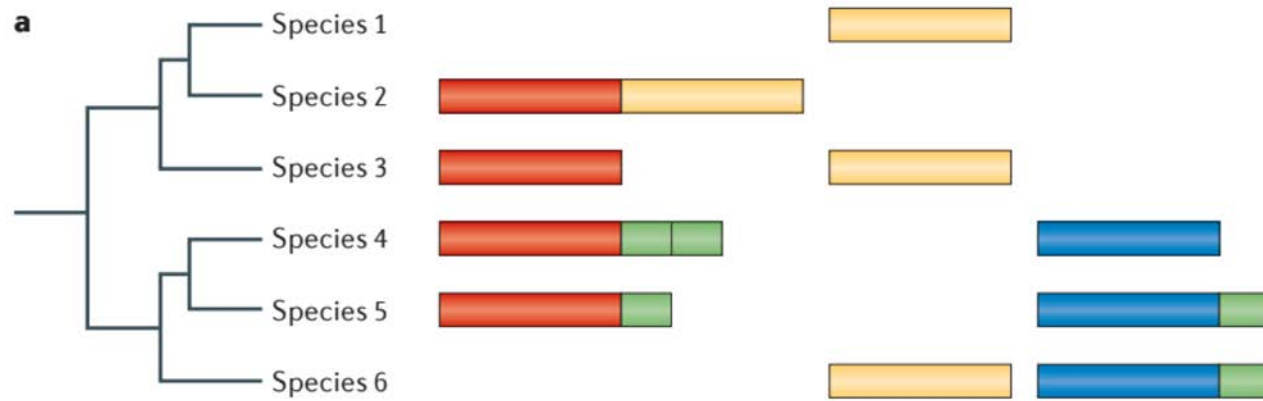
# Caveats

# Evolution of multi-domain proteins

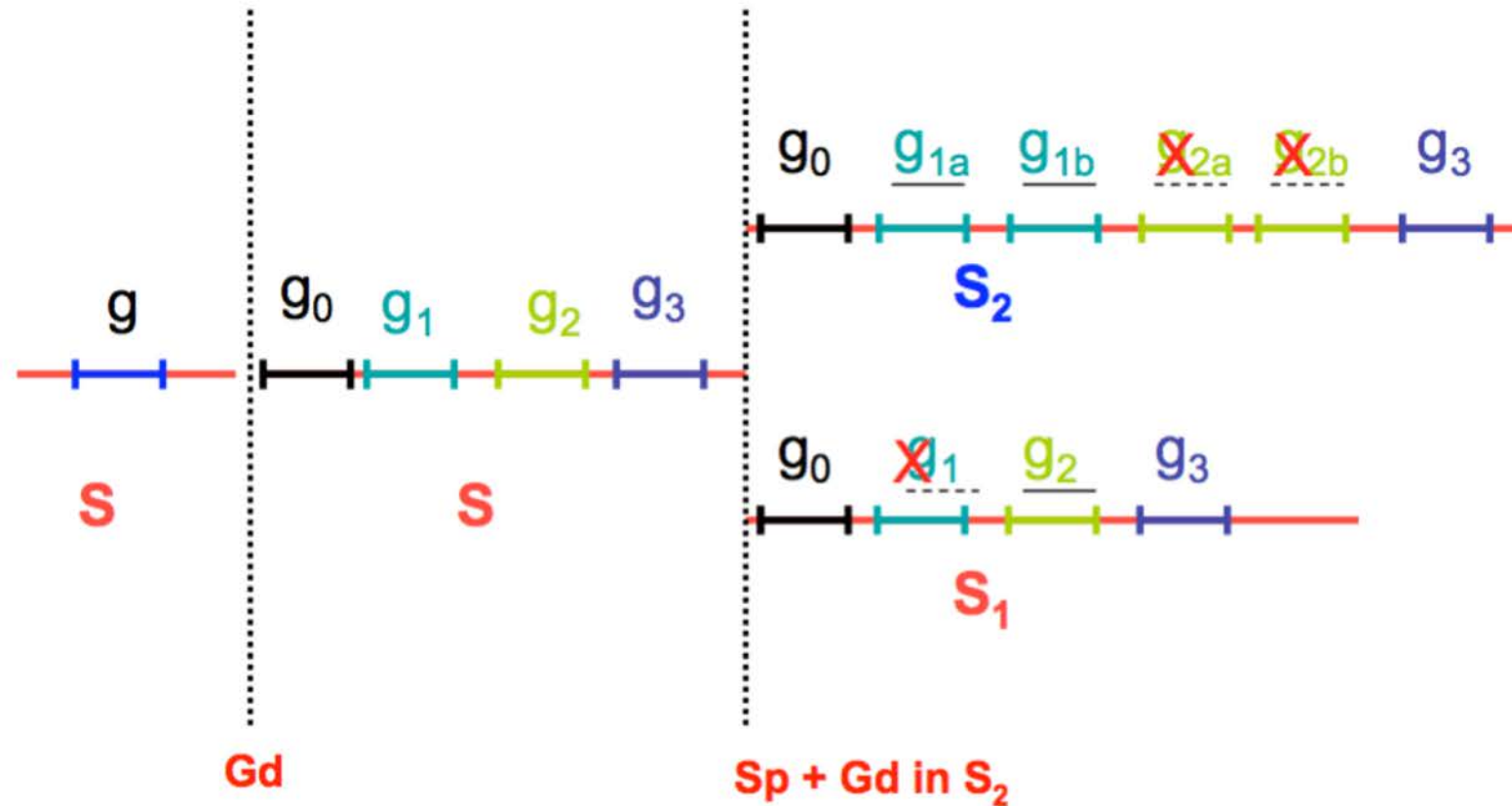


**Figure 1. The evolution of a hypothetical multidomain family by gene duplication and domain insertion.** Genes in the *a* and *b* subfamilies share a common ancestor but do not have identical domain composition. Gene *c* shares a homologous domain with genes in the *b* subfamily, but there is no gene that is ancestral to both *b* and *c*.  
doi:10.1371/journal.pcbi.1000063.g001

# Problem of clustering to assign gene families when comes to different domain combinations



# Detection can go wrong: Example of an orthology misleading situation



We assume that gene  $g_1$  (in S<sub>1</sub>) and genes  $g_{2a}$  and  $g_{2b}$  (in S<sub>2</sub>) are lost, similarity and phylogenetic methods for orthology detection will assign erroneously orthology to  $g_2$ ,  $g_{1a}$  and  $g_{1b}$ . Indeed these are not orthologous, because  $g_2$ ,  $g_{1a}$  and  $g_{1b}$  do not result from the same ancestral gene after the speciation event.

In this case solely the environment conservation, will help in detecting the gene duplication and loss event, and hypothesise their non-orthology.

# Summary point

## SUMMARY POINTS

1. Orthologs and paralogs are two types of homologous genes that evolved, respectively, by vertical descent from a single ancestral gene and by duplication.
2. Distinguishing between orthologs and paralogs is crucial for successful functional annotation of genomes and for reconstruction of genome evolution.
3. A finer classification of orthologs and paralogs has been developed to reflect the interplay between duplication and speciation events, and effects of gene loss and horizontal gene transfer on the observed homologous relationship.
4. Methods for identification of sets of orthologous and paralogous genes involve phylogenetic analysis and various procedures for sequence similarity-based clustering.
5. Analysis of clusters of orthologous and paralogous genes is instrumental in genome annotation and in delineation of trends in genome evolution.
6. Rearrangements of gene structure confound orthologous and paralogous relationships.
7. The gene-centered concepts of orthology and paralogy can be generalized downward, to the level of strings of nucleotides and even single base pairs, and upward, to multigene arrays.

# Phylogenomics

Phylogenomics aims at inferring detailed information about the evolutionary histories of organisms by using whole genomes rather than just a single gene or a few genes. The term was coined by Jonathan Eisen in the context of prediction of gene function

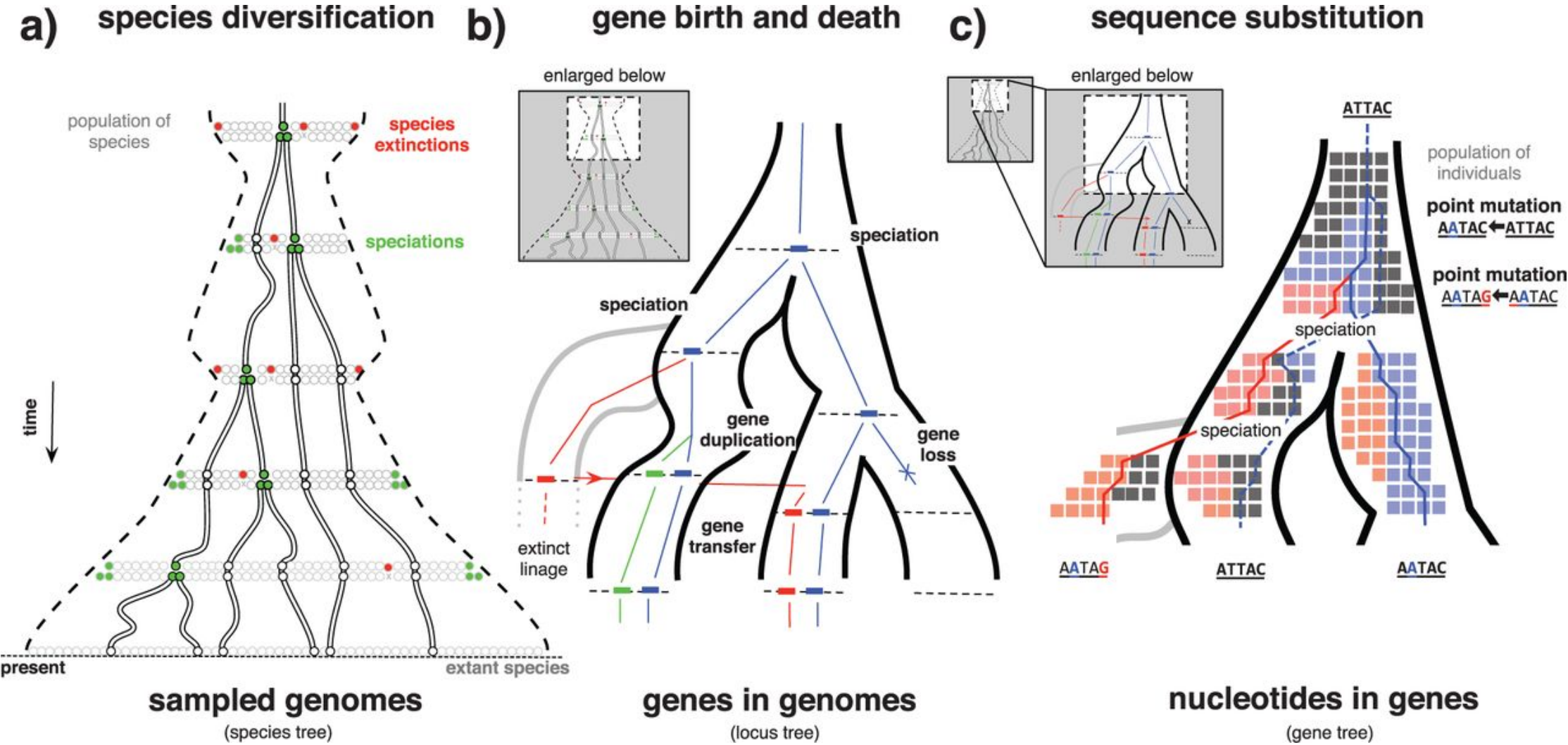
**It would be difficult or impossible to understand the evolutionary history of an organism, even having available its whole genome sequence, in isolation.** So it is always the case the phylogenomics is practiced for sets of genomes.



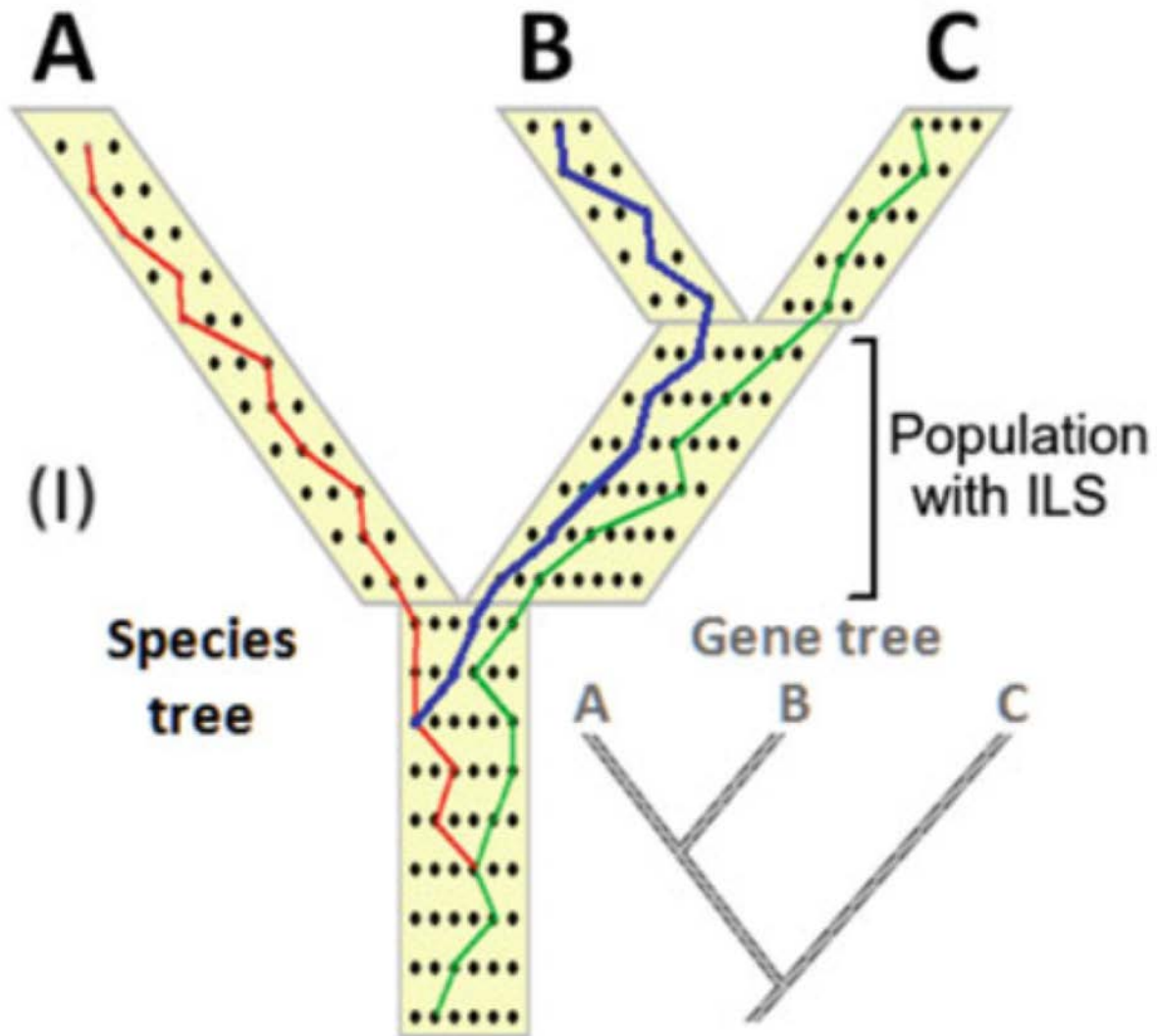
During the last 50 years, phylogeny has become more and more based on molecular data, increasingly **favoring homologous sequences over morphological characters**. This approach has been extremely fruitful, **producing constant improvement in the accuracy and resolution of phylogenetic reconstruction together with our understanding of evolutionary processes at the molecular level**.

However, we have known all along that we are barking up the wrong trees: with increasing sophistication in the models of sequence evolution, **we have been reconstructing trees describing the history of fragments of genomic sequence, which we will liberally call “gene” in this review, but never the history of species. Gene trees are not species trees** (Maddison 1997).

Each level of the hierarchy contributes to generating phylogenetic signal that can lead to differences between reconstructed gene trees.



# Processes that may induce gene trees that are different than the actual species tree

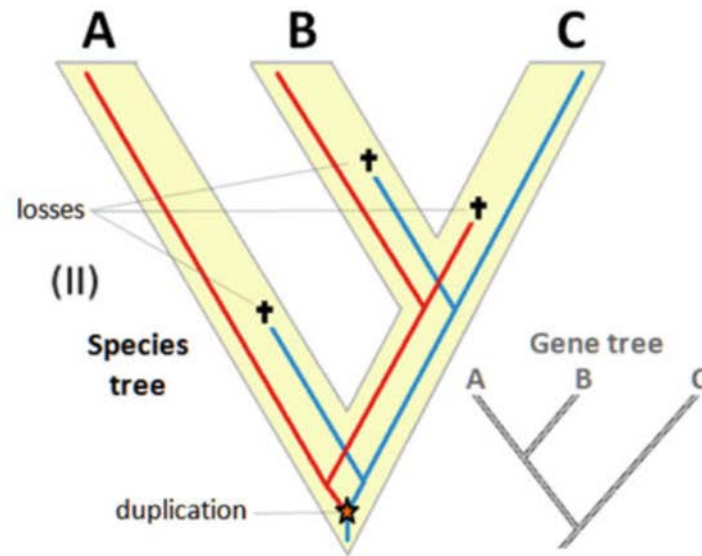


## i) Incomplete lineage sorting

When a species splits in two, allelic lineages sort into the two descendant species, and this lineage sorting varies along the genome.

If speciation events are close in time, the lineage sorting process may be incomplete at the second speciation event and lead to gene genealogies that do not match the species phylogeny

# Processes that may induce gene trees that are different than the actual species tree

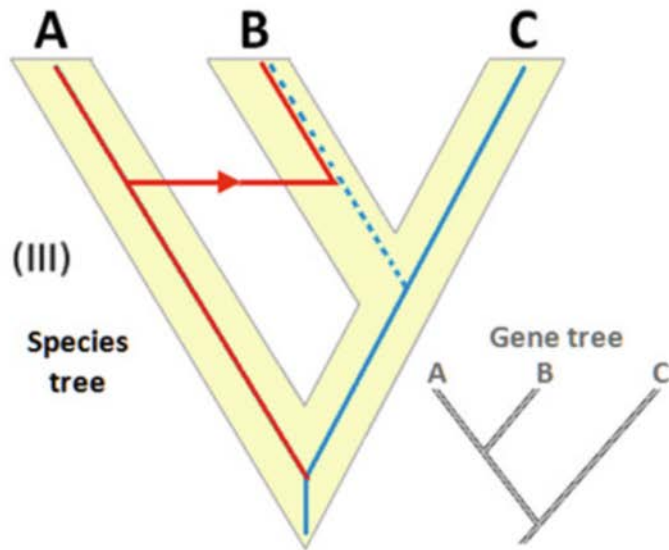


## (II) Duplication and Loss

a locus may generate a duplicate somewhere in the genome, and then both may be inherited or just a single copy is maintained in each lineage.

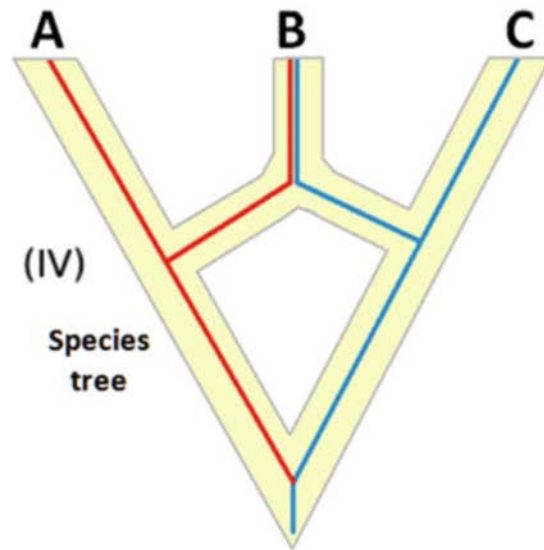
## (III) Horizontal Gene Transfer

(HGT): a donor DNA segment (from taxon A) is transmitted and incorporated into the host's genome (taxon B)

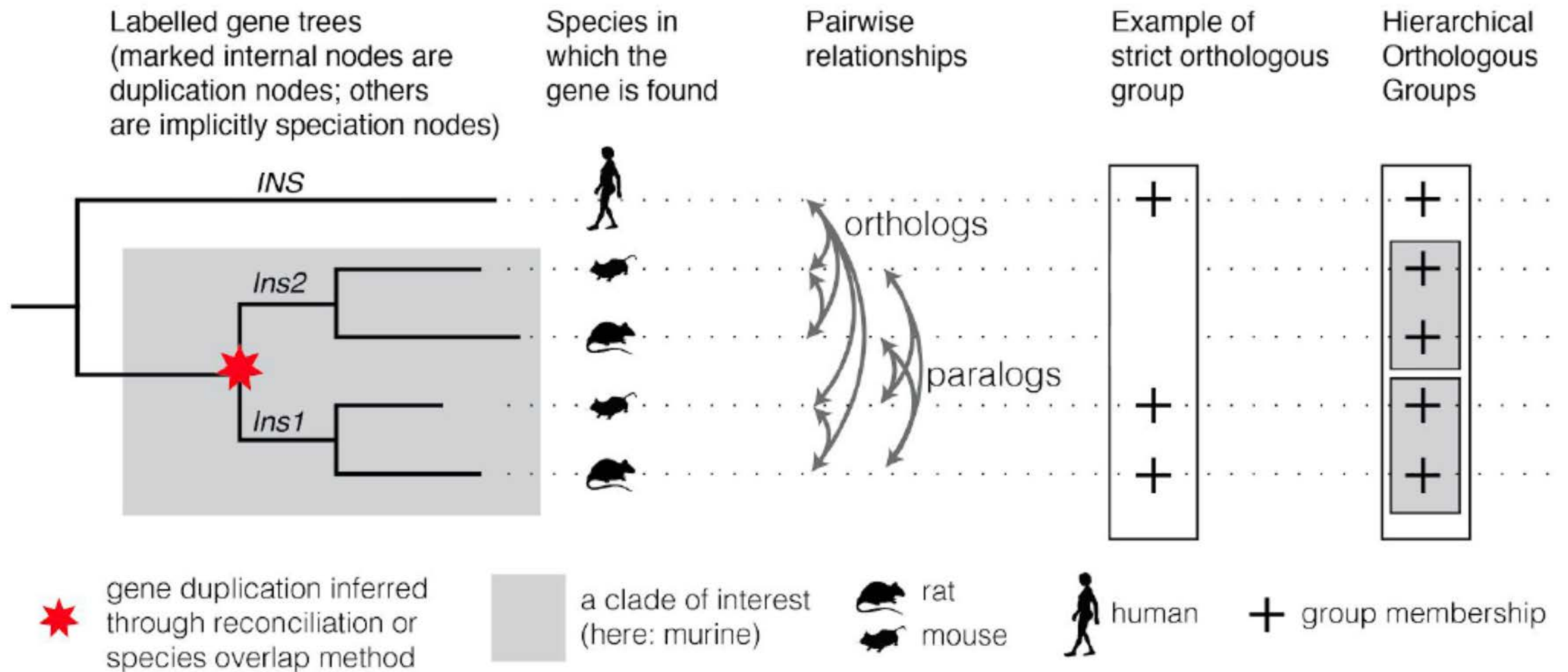


## (IV) Hybridization/Introgression

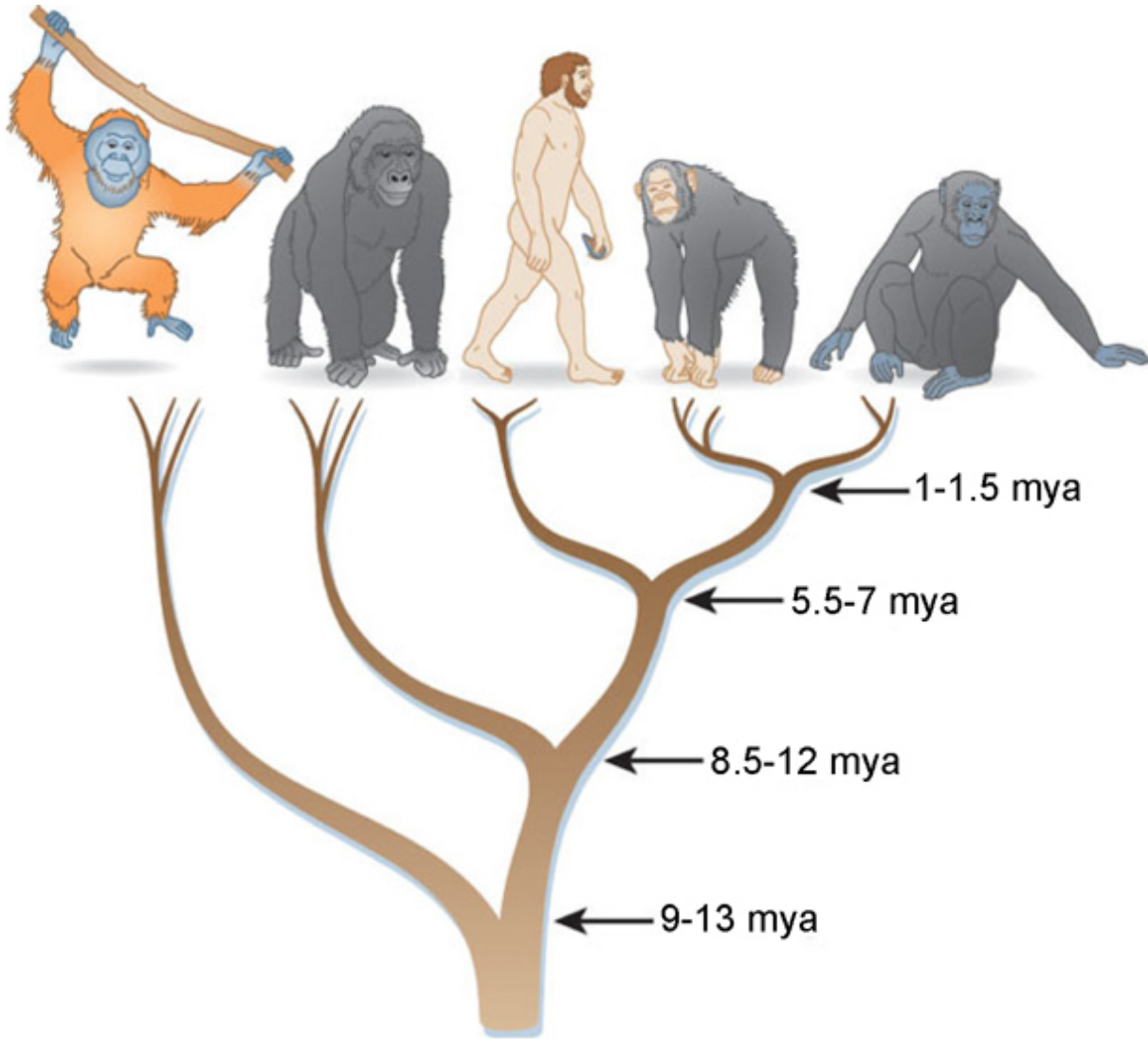
in extreme cases of lateral transfer, or upon mixing of related species, different regions of the genome will bear two distinct evolutionary histories;



# Problem of obtaining the 'true' orthologs for phylogenomics



# Why is Studying (Ape) Speciation Important? (Example)

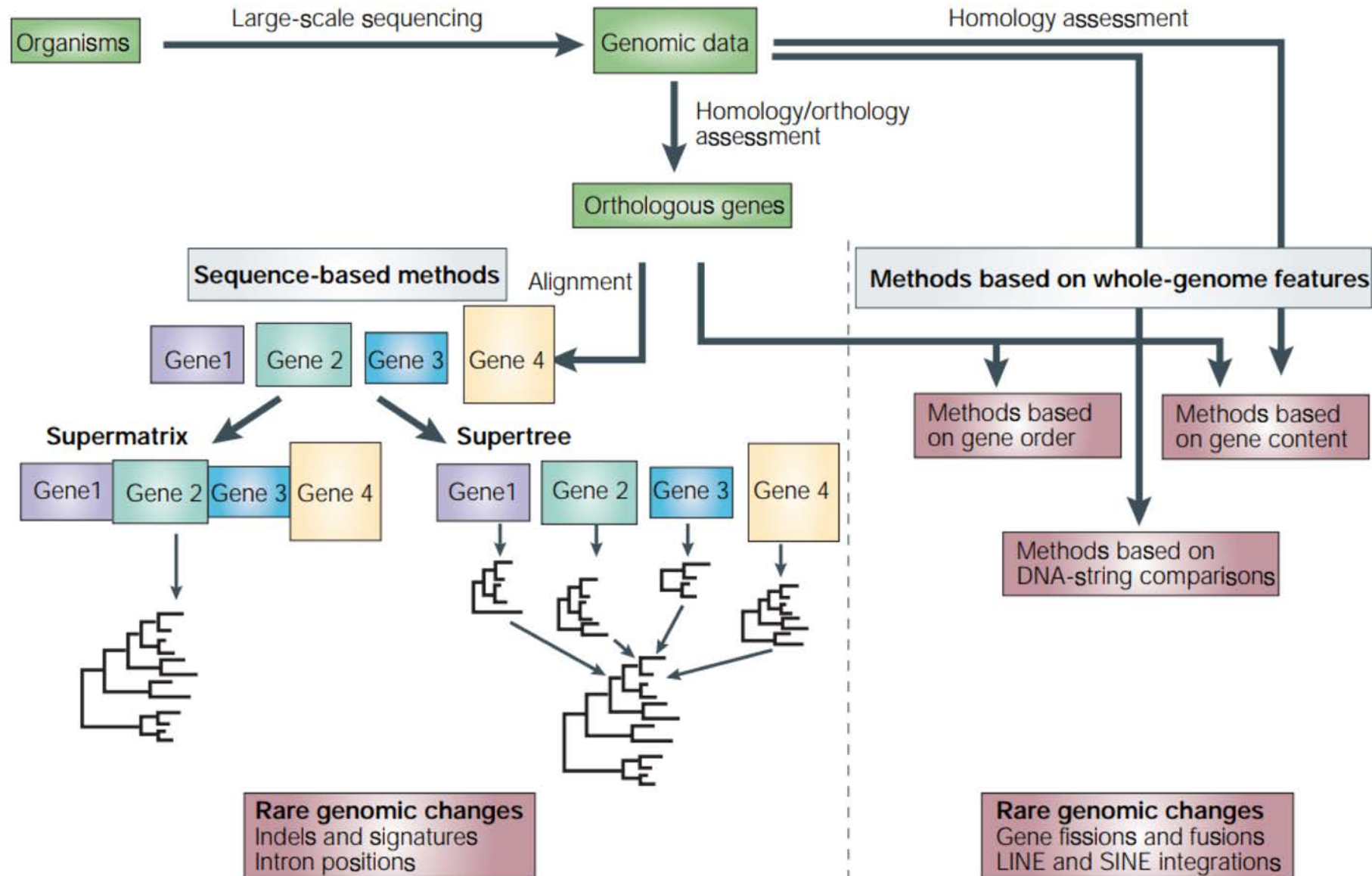


These studies also led to rich discussions about the suite of **factors that may have contributed to promoting speciation in the last common ancestor of humans and African apes**, as well as the **factors that might have contributed to creating the amazing diversity of Hominins that co-existed with each other during the Pliocene and Pleistocene** (Foley 2002).

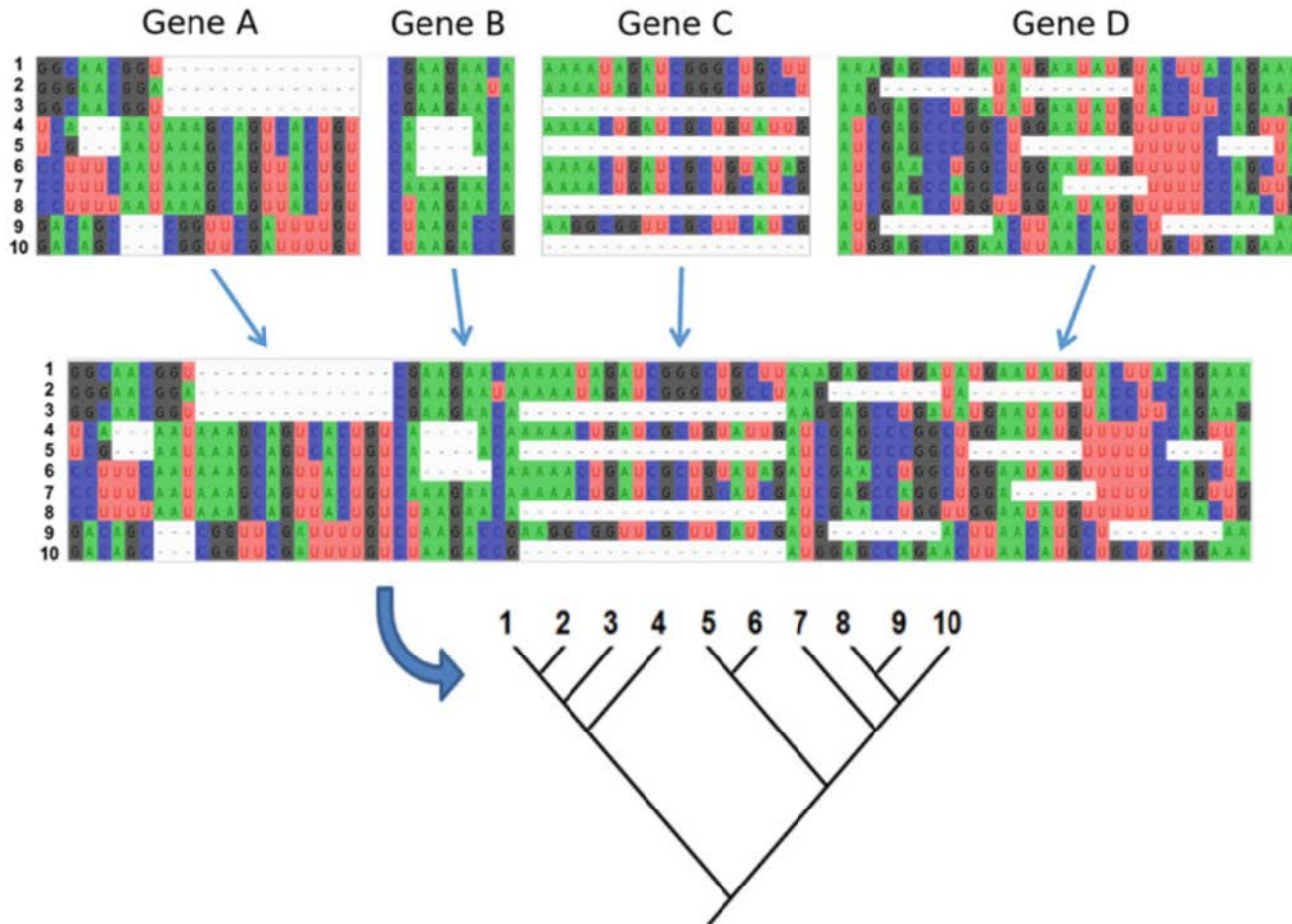
**For many years, there was considerable debate about which of the African apes is our closest relative....** The general consensus that emerged is that we share a more recent relationship with chimpanzees (*Pan troglodytes*) and bonobos (*Pan paniscus*) than we do with gorillas (*Gorilla gorilla*) (Ruvolo 1997, Chen & Li 2001).

Current estimates indicate that up to 30% of the sequence of the human genome is more closely related to Gorilla than to Chimpanzee due to this process (Scally et al. 2012).

# Probably the most common (easy) way to construct alignment of concatenated gene shared across all species



# Probably the most common (easy) way to construct alignment of concatenated gene shared across all species

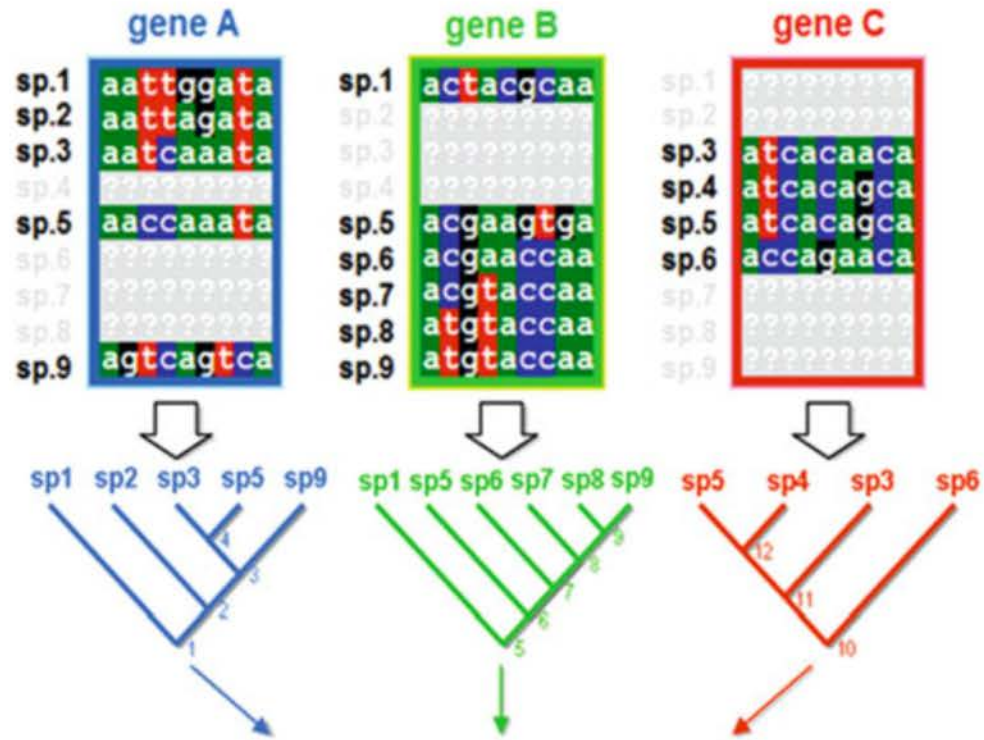


## Important drawbacks:

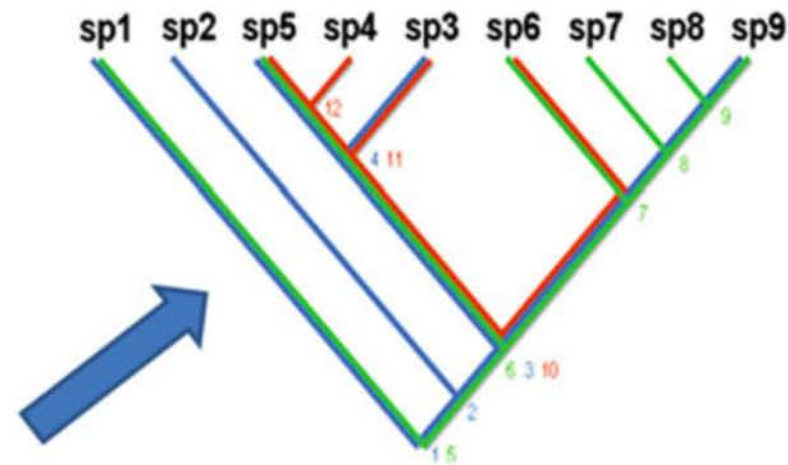
- (1) it hinders variation among gene trees by assuming implicitly that all of them conform to a single species tree;
- (2) if sampling was heterogeneous across species there may be too much missing data, which can affect topological reconstruction; Or limited number of genes shared among all species
- (3) large data sampling effects inflate credibility in some clades;
- (4) spurious hidden support can lead to support for non-existent clades; and
- (5) in case of moderate to severe levels of ILS, supermatrix can become statistically inconsistent.



# From genes to supertrees

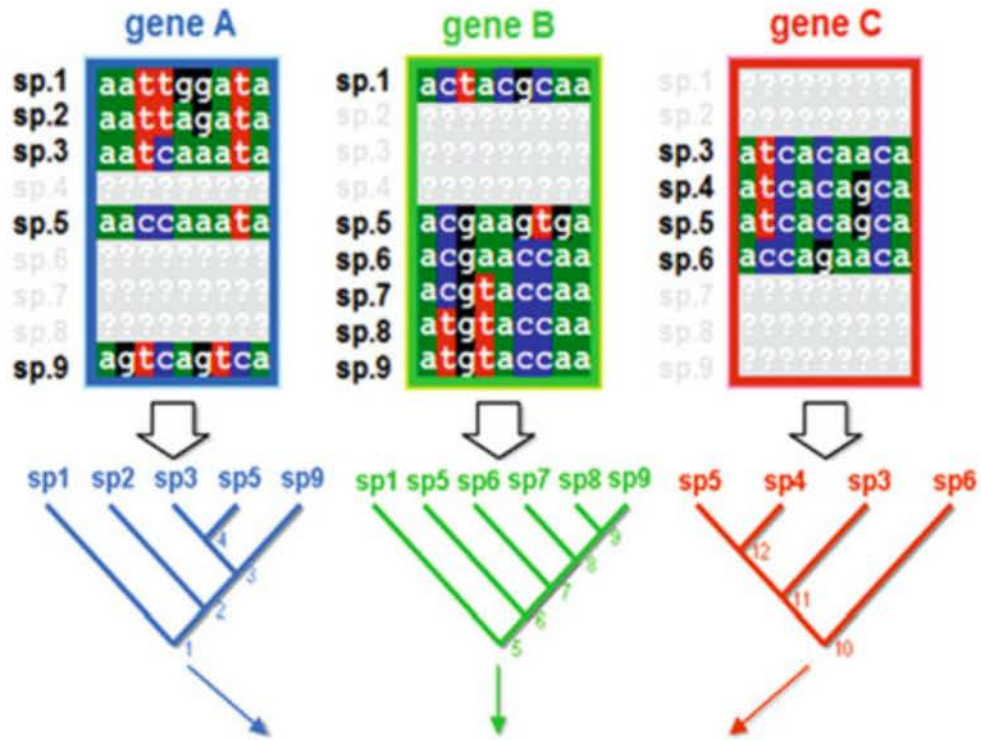


MRP	1	2	3	4	5	6	7	8	9	10	11	12
sp1	1	0	0	0	1	0	0	0	0	?	?	?
sp2	1	1	0	0	?	?	?	?	?	?	?	?
sp3	1	1	1	1	?	?	?	?	?	1	1	0
sp4	?	?	?	?	?	?	?	?	?	1	1	1
sp5	1	1	1	1	1	1	0	0	0	1	1	1
sp6	?	?	?	?	1	1	1	1	0	1	0	0
sp7	?	?	?	?	1	1	1	1	0	?	?	?
sp8	?	?	?	?	1	1	1	1	1	?	?	?
sp9	1	1	1	0	1	1	1	1	1	?	?	?

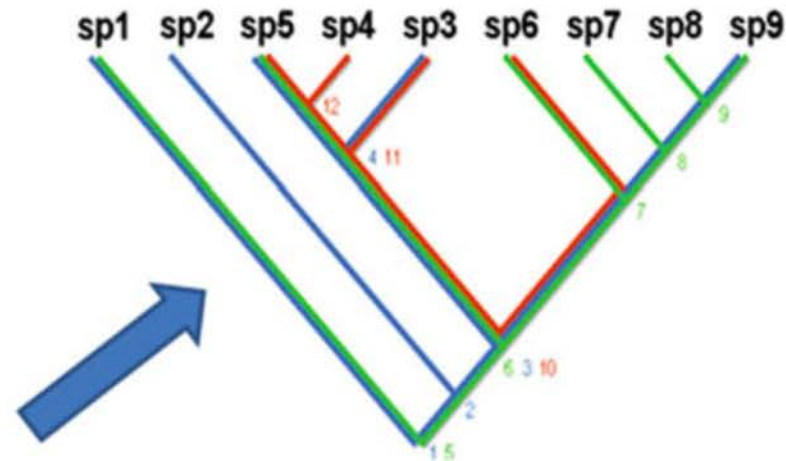


Instead of forcing all gene trees to comply to a single tree, **supertree methods infer the best topology for each gene (using the same phylogenetic method for each), and then a topological consensus is obtained.** Such methods are able to make consensus trees even if the number of leaves among gene trees differs but overlaps to some extent, for example when a gene has not been sequenced for some taxa

# Current methods



MRP	1	2	3	4	5	6	7	8	9	10	11	12
sp1	1	0	0	0	1	0	0	0	0	?	?	?
sp2	1	1	0	0	?	?	?	?	?	?	?	?
sp3	1	1	1	1	?	?	?	?	?	1	1	0
sp4	?	?	?	?	?	?	?	?	?	1	1	1
sp5	1	1	1	1	1	1	0	0	0	1	1	1
sp6	?	?	?	?	1	1	1	1	0	1	0	0
sp7	?	?	?	?	1	1	1	1	0	?	?	?
sp8	?	?	?	?	1	1	1	1	1	?	?	?
sp9	1	1	1	0	1	1	1	1	1	?	?	?

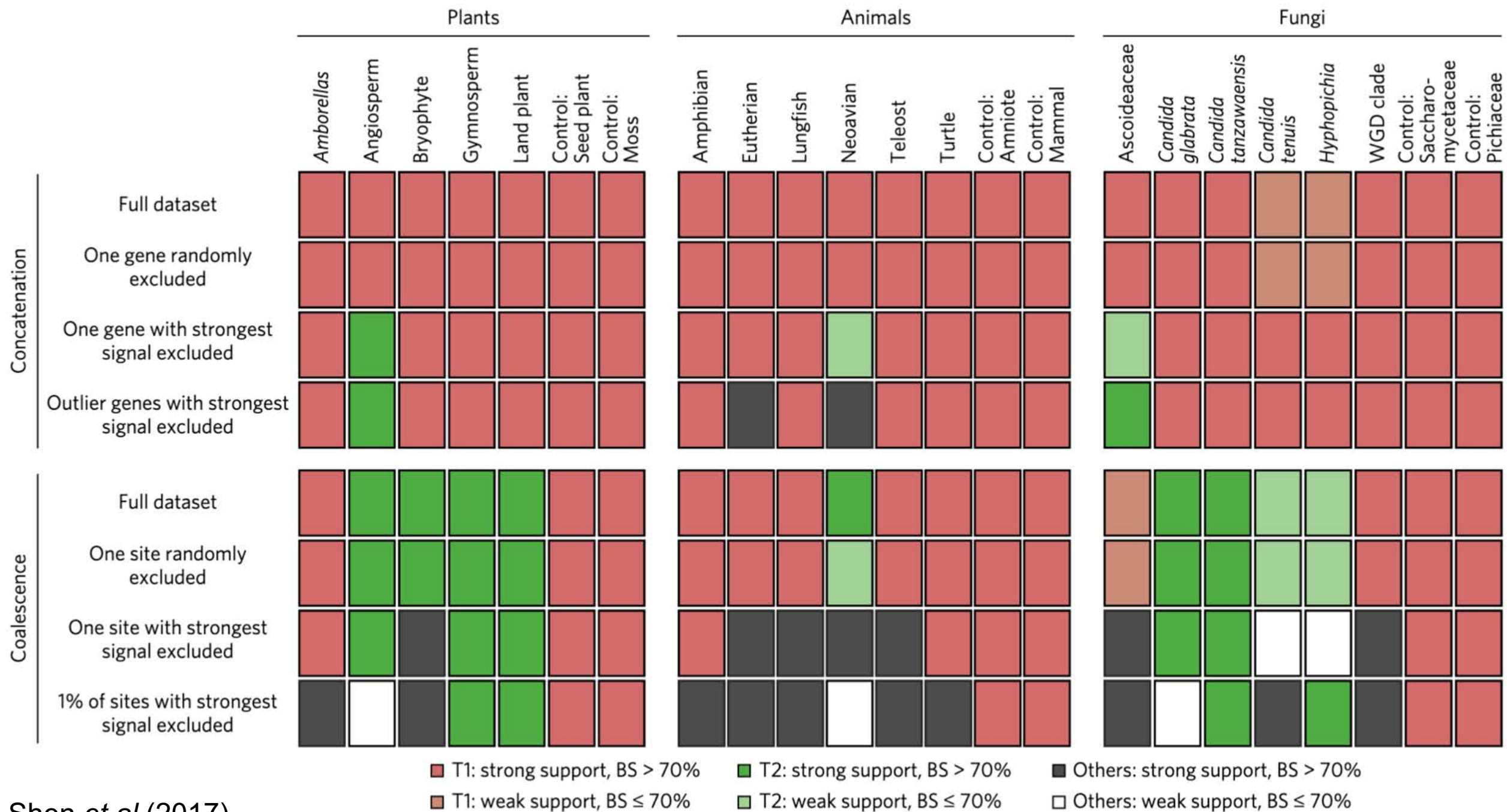


A step beyond supertrees is the use of methods that take into consideration specific evolutionary processes that may be responsible for differences in gene topologies, and then estimate the species tree which would most likely have generated such gene trees, under different scenarios

# Contentious relationships in phylogenomic studies can be driven by a handful of genes

Xing-Xing Shen<sup>1</sup>, Chris Todd Hittinger<sup>2</sup> and Antonis Rokas<sup>1\*</sup>

**Phylogenomic studies have resolved countless branches of the tree of life, but remain strongly contradictory on certain, contentious relationships. Here, we use a maximum likelihood framework to quantify the distribution of phylogenetic signal among genes and sites for 17 contentious branches and 6 well-established control branches in plant, animal and fungal phylogenomic data matrices. We find that resolution in some of these 17 branches rests on a single gene or a few sites, and that removal of a single gene in concatenation analyses or a single site from every gene in coalescence-based analyses diminishes support and can alter the inferred topology. These results suggest that tiny subsets of very large data matrices drive the resolution of specific internodes, providing a dissection of the distribution of support and observed incongruence in phylogenomic analyses. We submit that quantifying the distribution of phylogenetic signal in phylogenomic data is essential for evaluating whether branches, especially contentious ones, are truly resolved. Finally, we offer one detailed example of such an evaluation for the controversy regarding the earliest-branching metazoan phylum, for which examination of the distributions of gene-wise and site-wise phylogenetic signal across eight data matrices consistently supports ctenophores as the sister group to all other metazoans.**



Visualisation of gene content / families

# Extensive mosaic structure revealed by the complete genome sequence of uropathogenic *Escherichia coli*

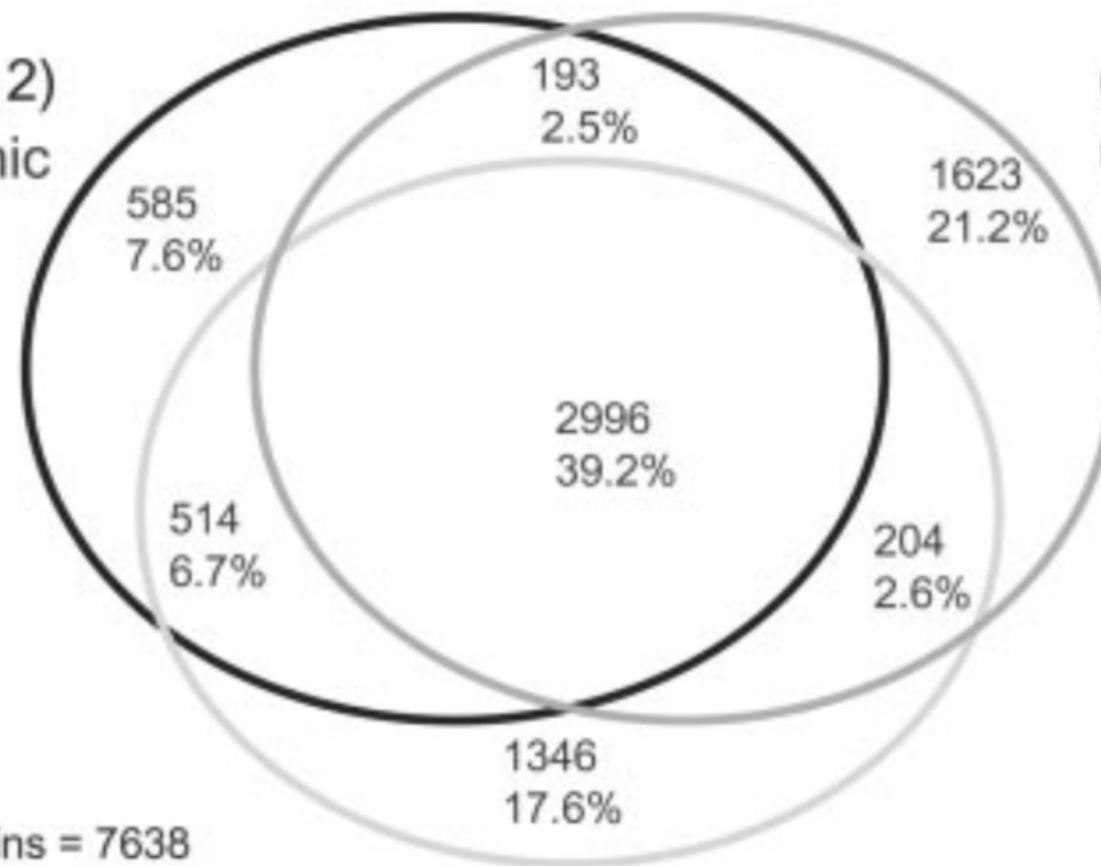
R. A. Welch\*, V. Burland<sup>††</sup>, G. Plunkett III<sup>†</sup>, P. Redford\*, P. Roesch\*, D. Rasko<sup>§</sup>, E. L. Buckles<sup>¶</sup>, S.-R. Liou<sup>¶||</sup>, A. Boutin<sup>†\*\*</sup>, J. Hackett<sup>†.††</sup>, D. Stroud<sup>†</sup>, G. F. Mayhew<sup>†</sup>, D. J. Rose<sup>†</sup>, S. Zhou<sup>†††</sup>, D. C. Schwartz<sup>†††</sup>, N. T. Perna<sup>§§</sup>, H. L. T. Mobley<sup>§</sup>, M. S. Donnenberg<sup>¶</sup>, and F. R. Blattner<sup>†</sup>

\*Department of Medical Microbiology,  
Sciences, University of Wisconsin  
Department of Medicine, University of Wisconsin

Edited by John J. Mekalanos, Harvard University

MG1655 (K-12)  
non-pathogenic

CFT073  
uropathogenic

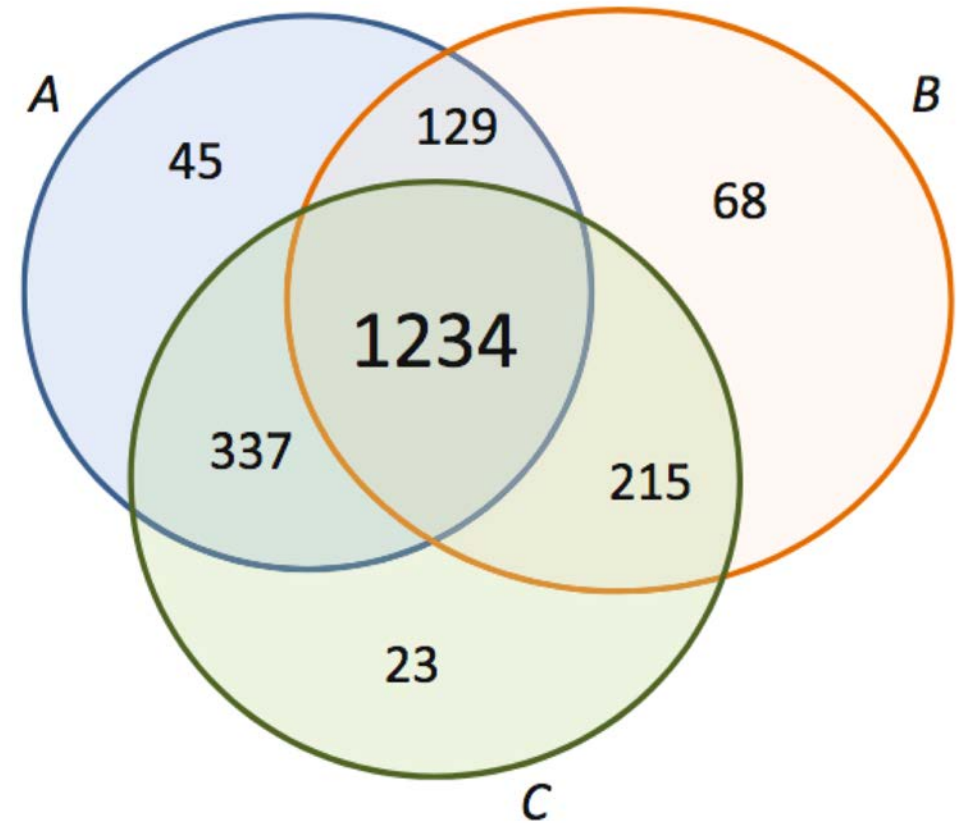


Total proteins = 7638  
2996 (39.2%) in all 3  
911 (11.9%) in 2 out of 3  
3554 (46.5%) in 1 out of 3

EDL933 (O157:H7)  
enterohaemorrhagic

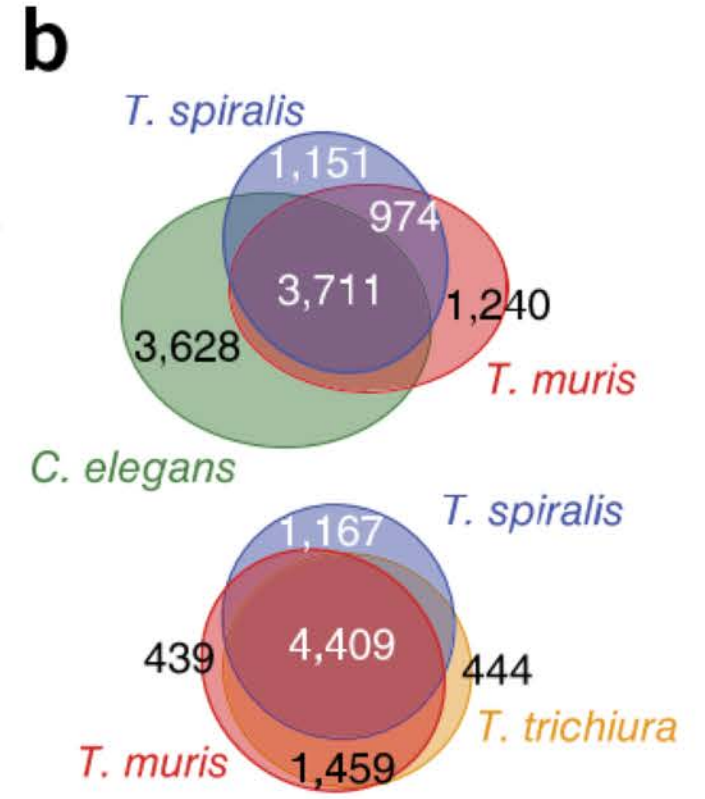
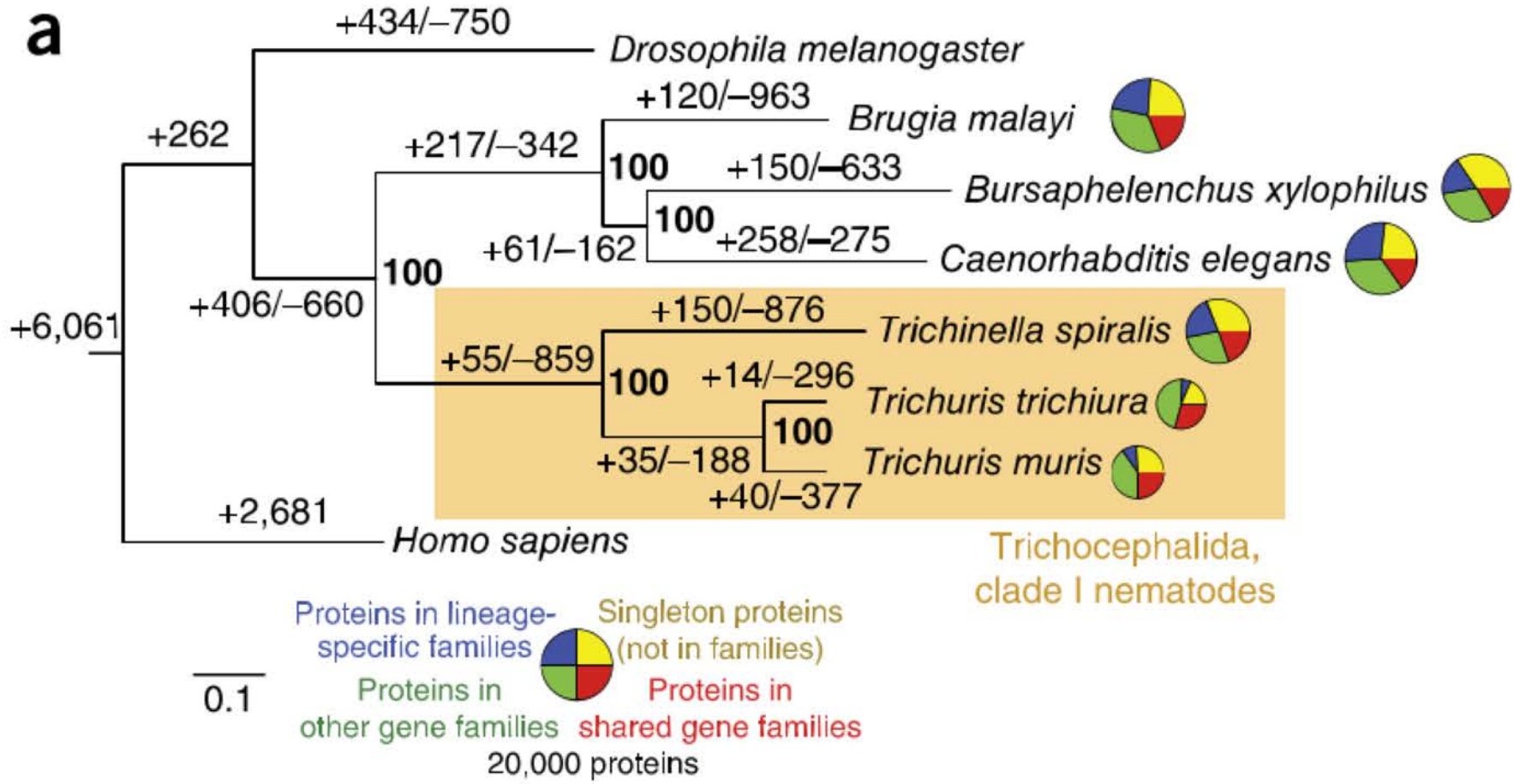
Illustration of a gene content Venn diagram for three hypothetical genomes A, B, and C

Gene	Genome						
	A	B	C	D	E	F	G
1	✓	✓			✓	✓	
2	✓		✓	✓	✓	✓	✓
3		✓		✓			
4		✓			✓		
5				✓			
6			✓		✓	✓	
7		✓		✓			✓



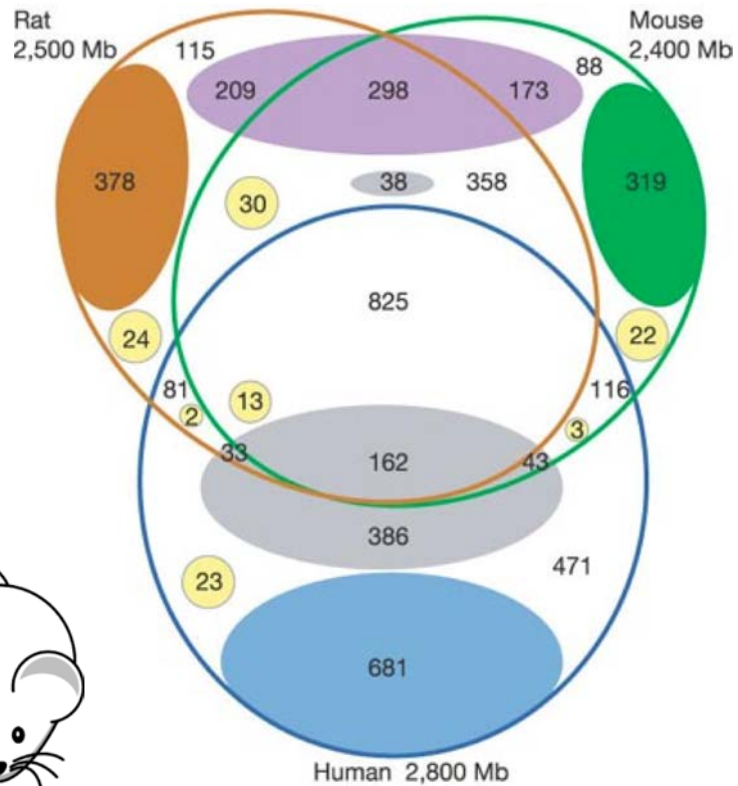
Schematic representation of a presence/absence gene matrix. Genomes are represented in columns, and gene families are represented in rows

# Phylogeny + Venn diagram to show expansion/loss





# Trend of venn diagram...



Genomic DNA	Rat	Mouse	Human
Repetitive DNA	Ancestral to human-mouse-rat	Rat-specific	Primate-specific
	Ancestral to mouse-rat	Mouse-specific	Simple

## A

### Dicots

- Arabidopsis thaliana*: 26304 / 24766
- Glycine max*: 36271 / 35969
- Populus trichocarpa*: 35516 / 33358
- Ricinus communis*: 30314 / 24039
- Theobroma cacao*: 28222 / 27154
- Vitis vinifera*: 24479 / 21795

### Basal

- Amborella trichopoda*: 24611 / 21191

### Early land plants

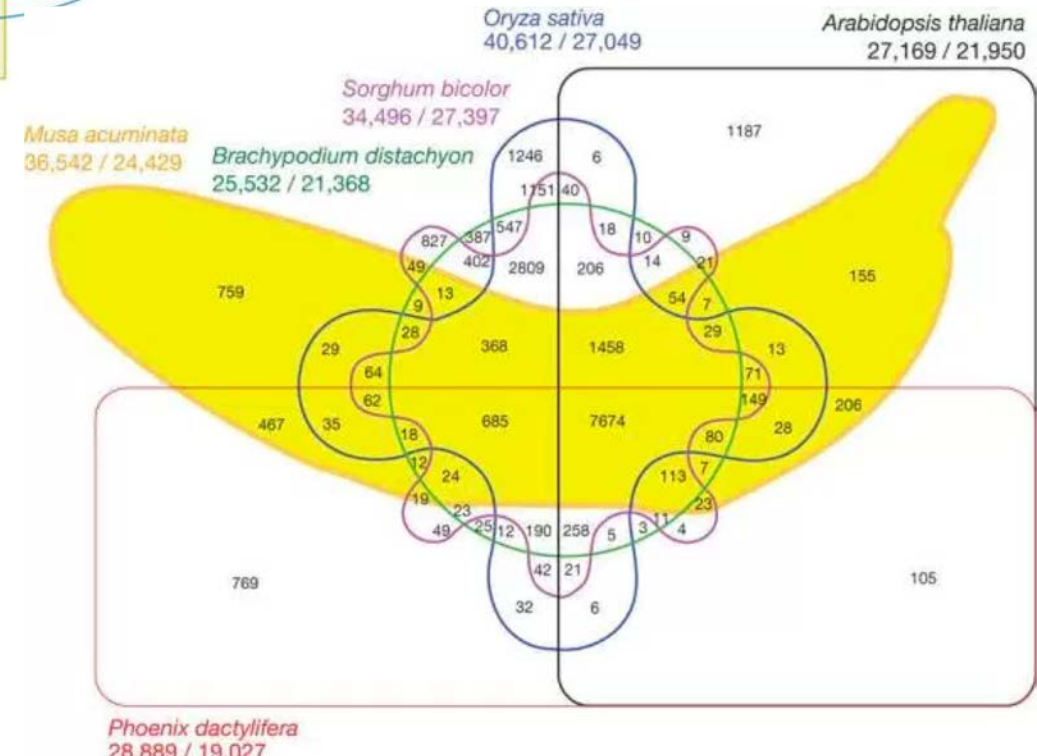
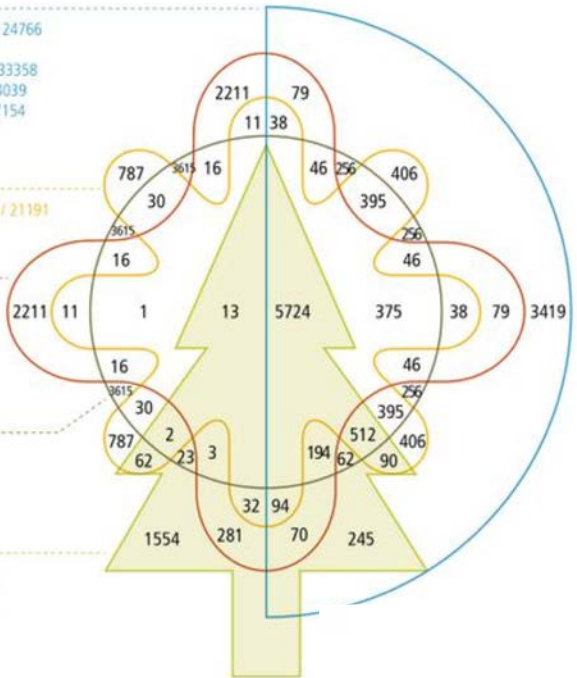
- Selaginella moellendorffii*: 16832 / 15909
- Physcomitrella patens*: 25938 / 19359

### Monocots

- Oryza sativa*: 39459 / 32660
- Zea mays*: 34586 / 30799

### Conifers

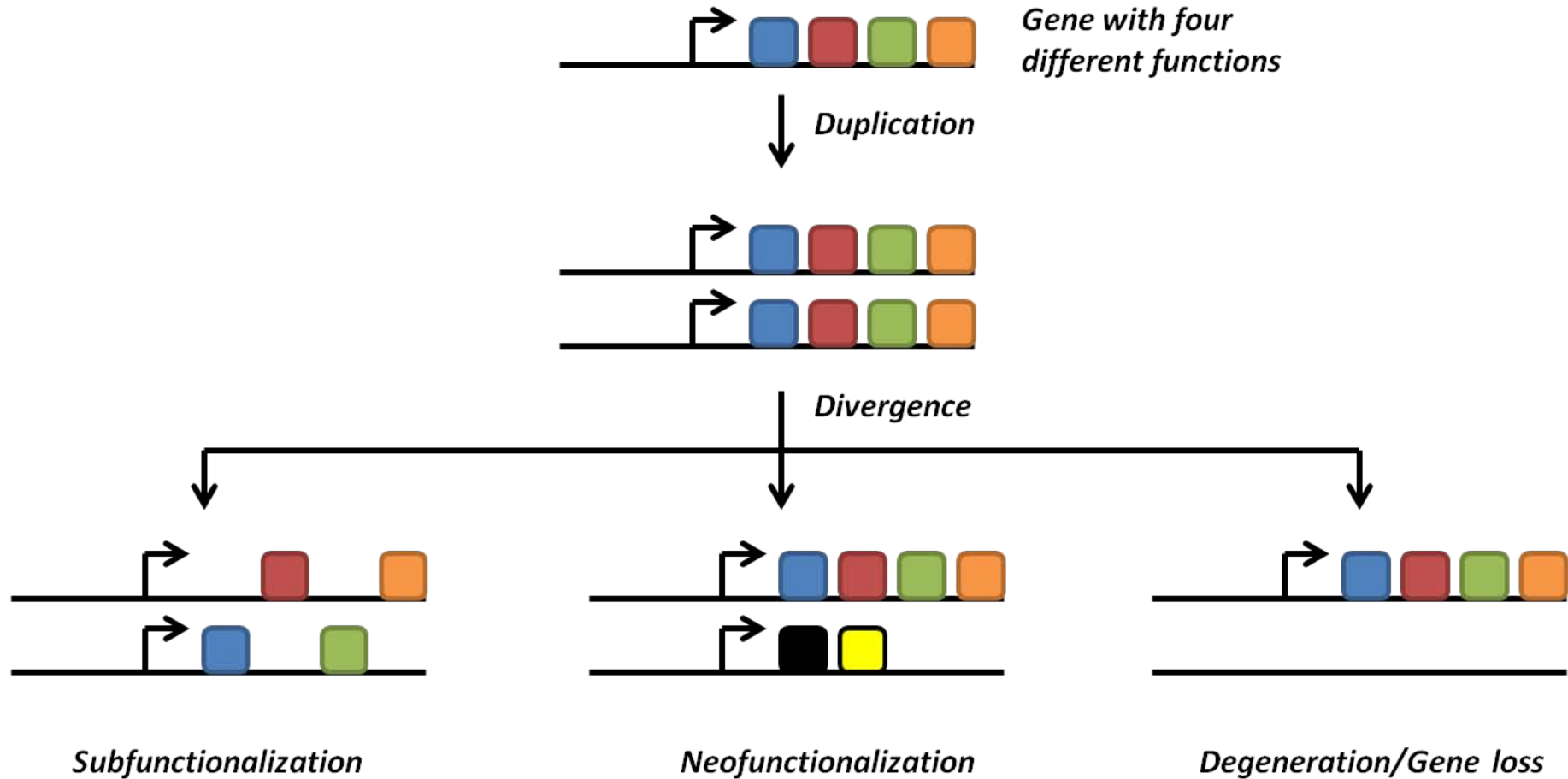
- Picea abies*: 20861 / 19934
- Picea sitchensis*: 8758 / 7780
- Pinus taeda*: 47207 / 46720



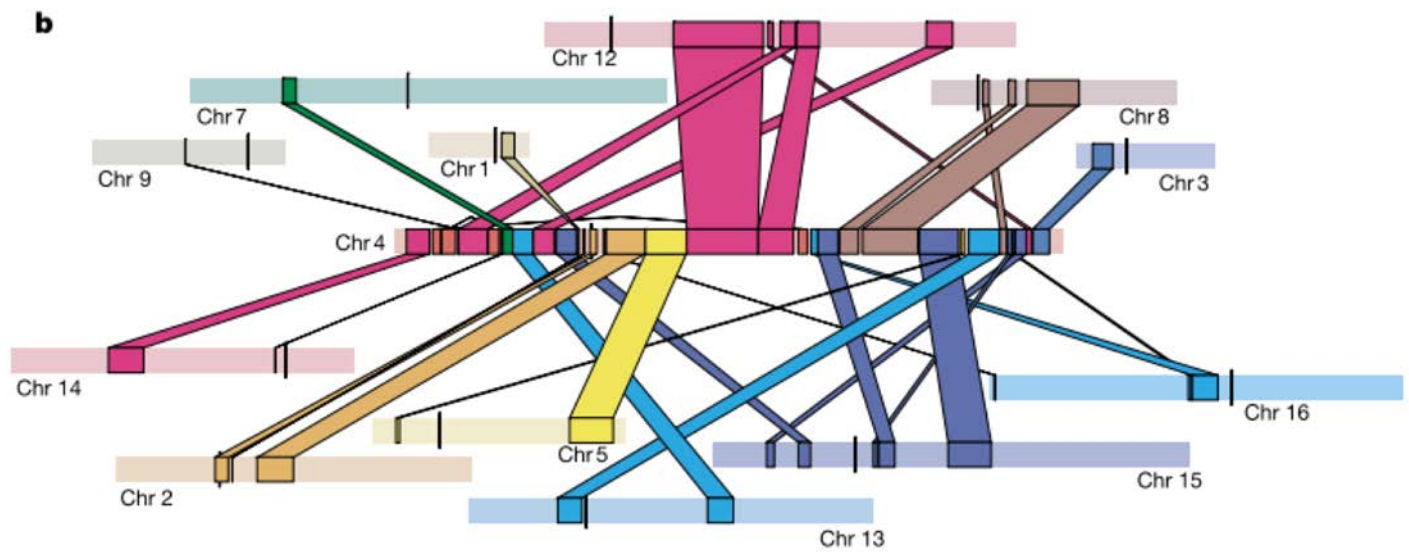
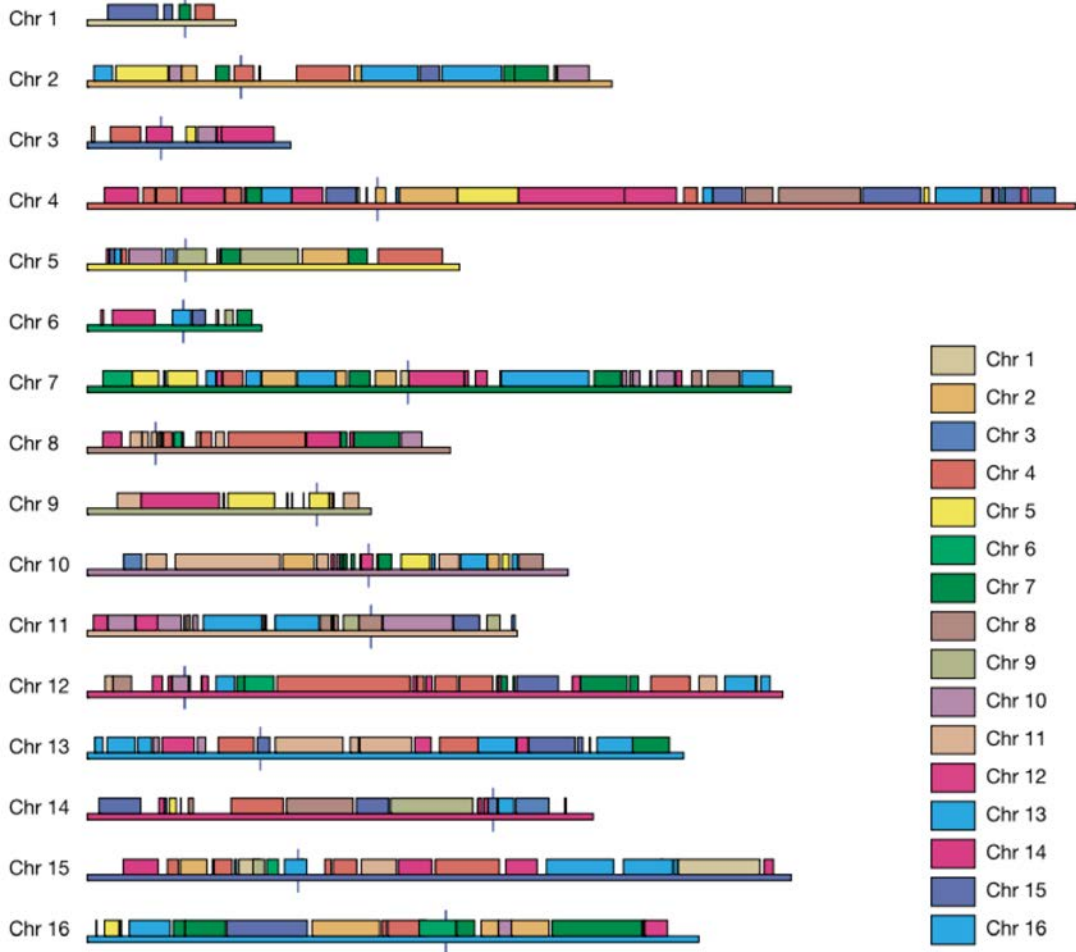
# Gene and genome duplication

# Why study gene duplication?

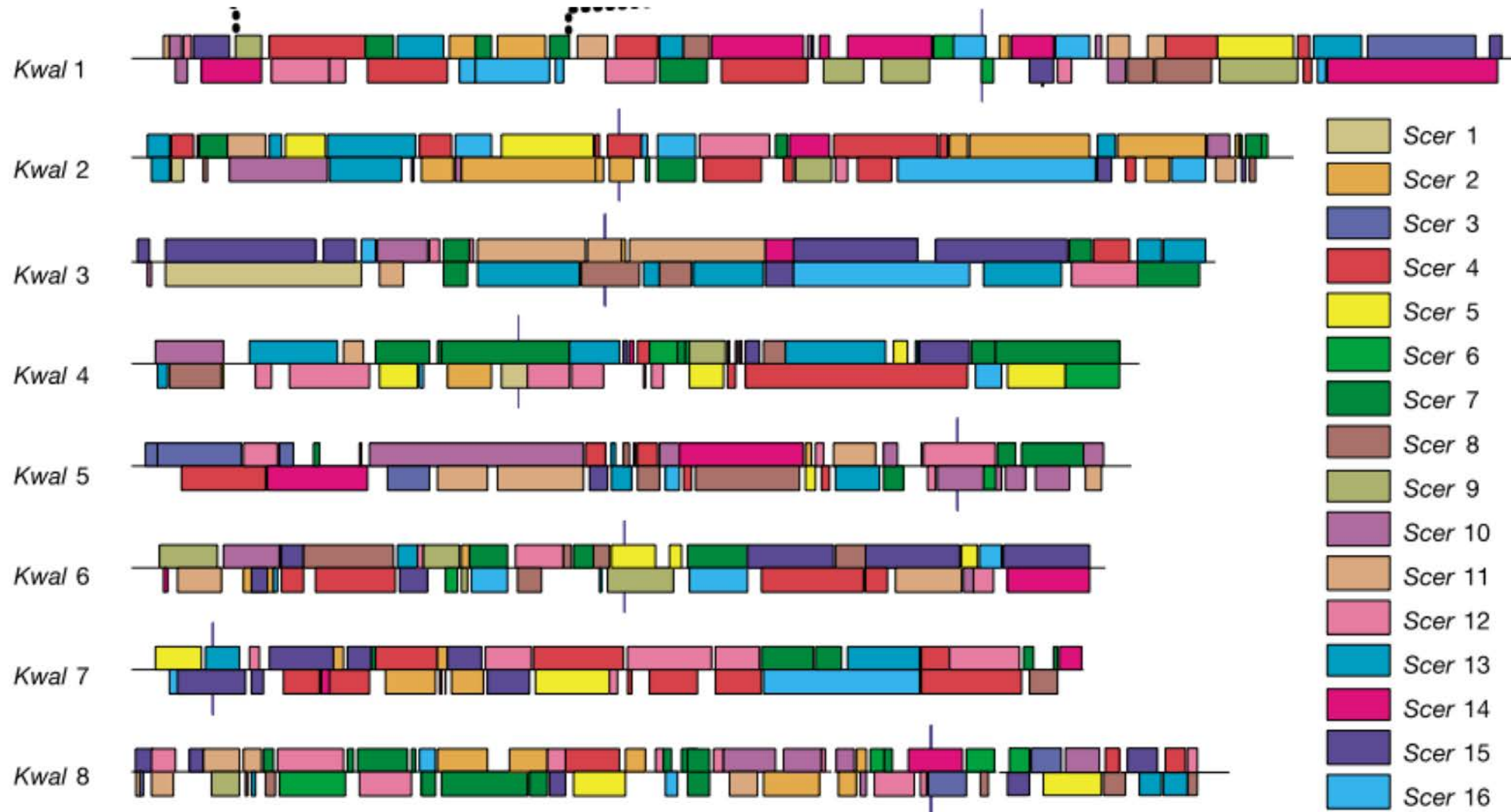
Gene duplications are traditionally considered as a major evolutionary source for protein new functions



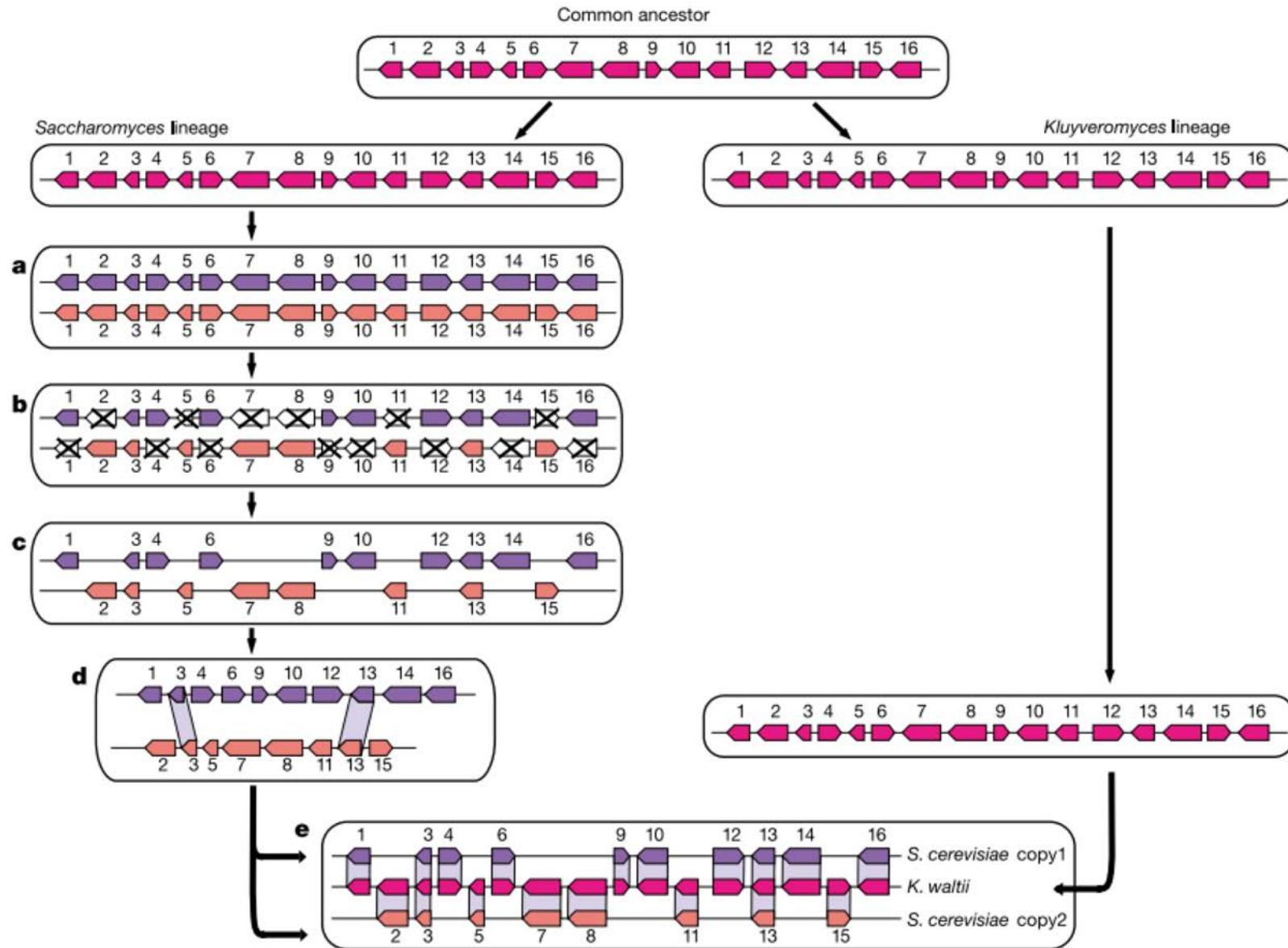
# Within species



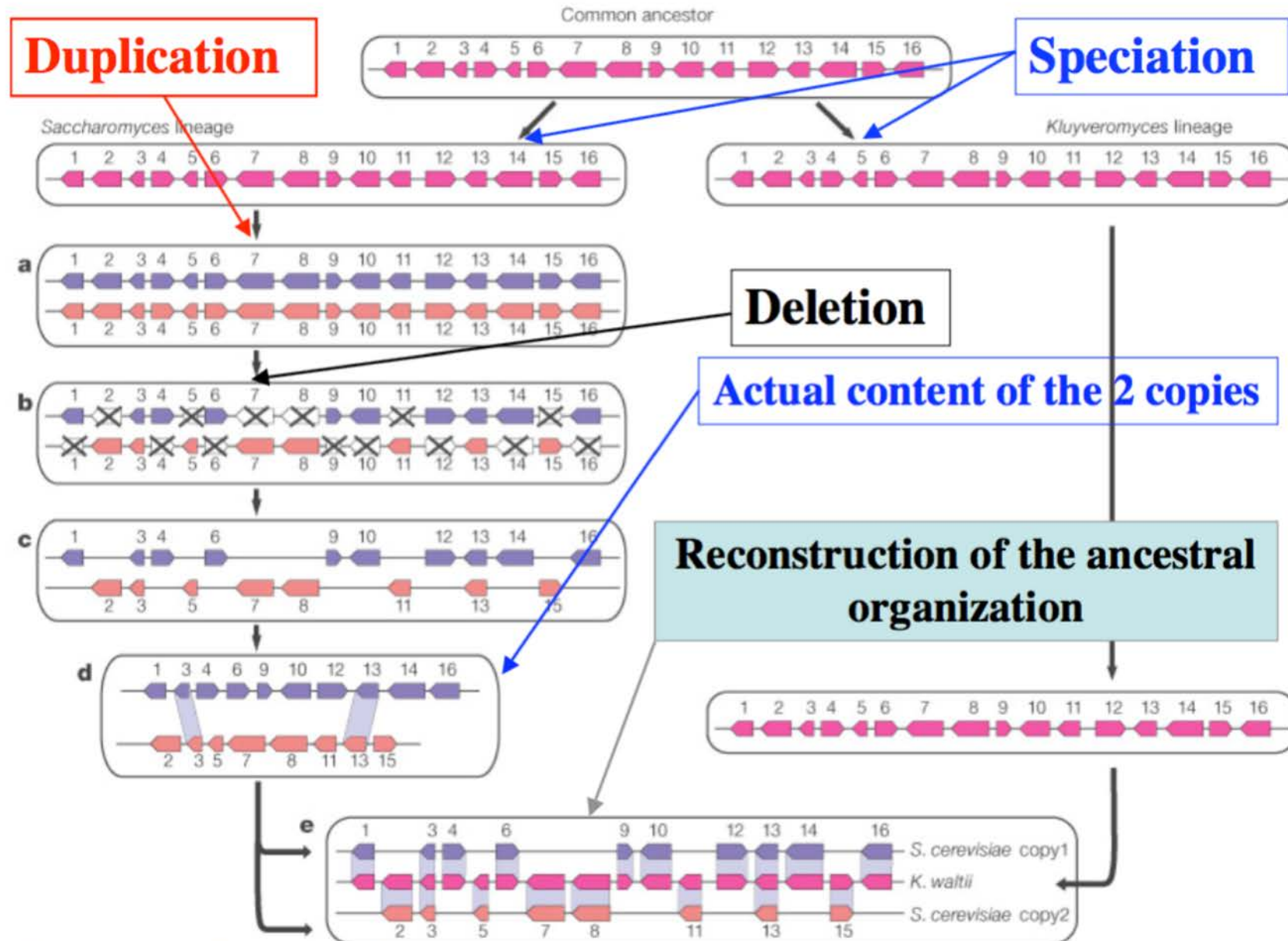
# Between species



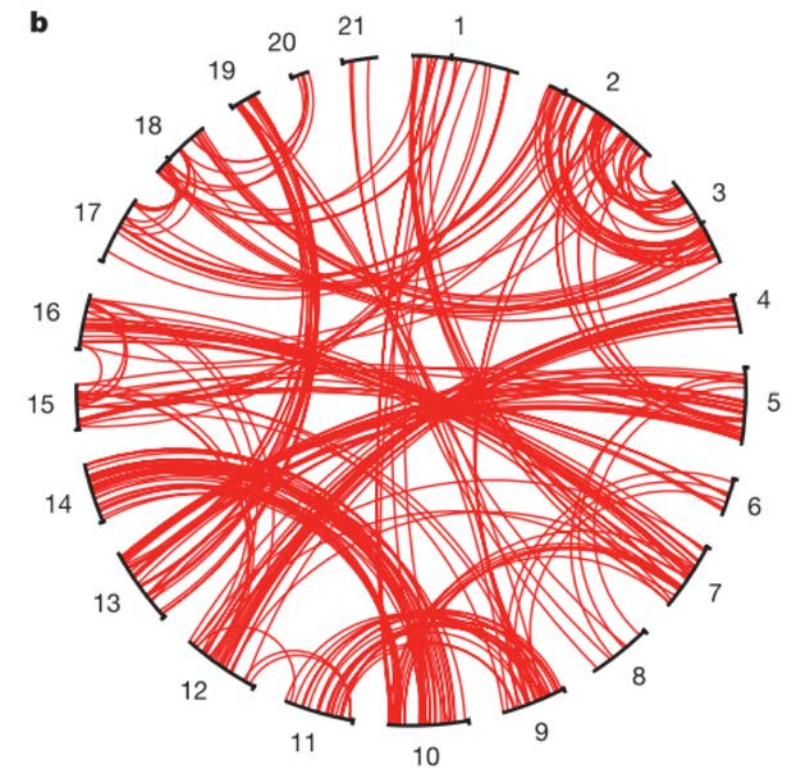
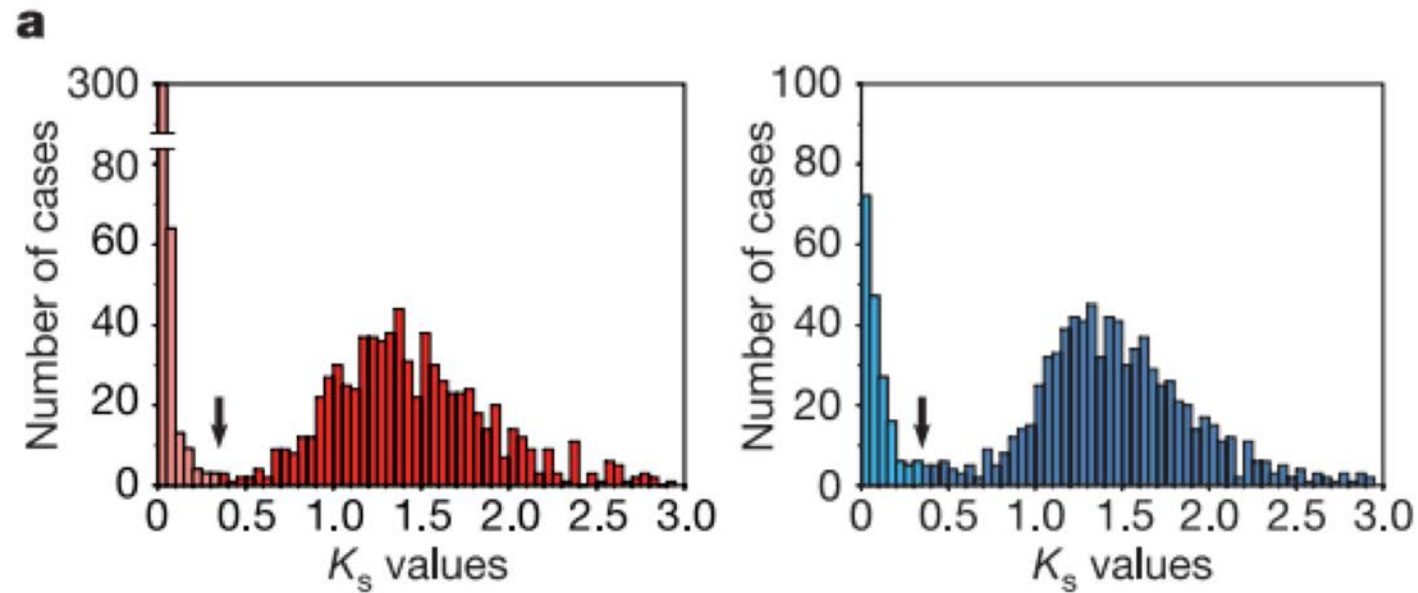
# Whole genome duplication model



# Determining ancestral conservation

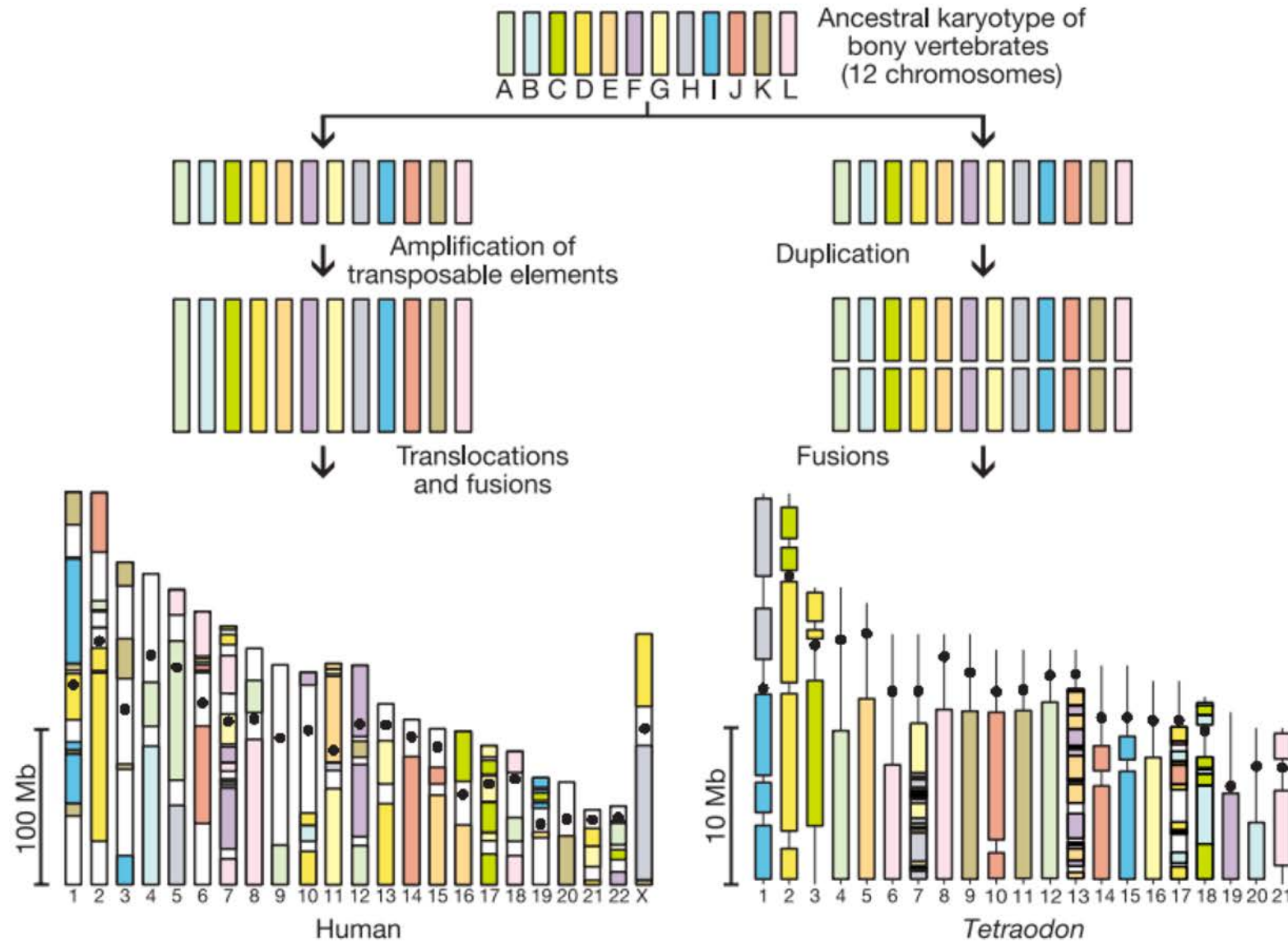


# Genome duplication in the teleost fish *Tetraodon nigroviridis* reveals the early vertebrate proto-karyotype



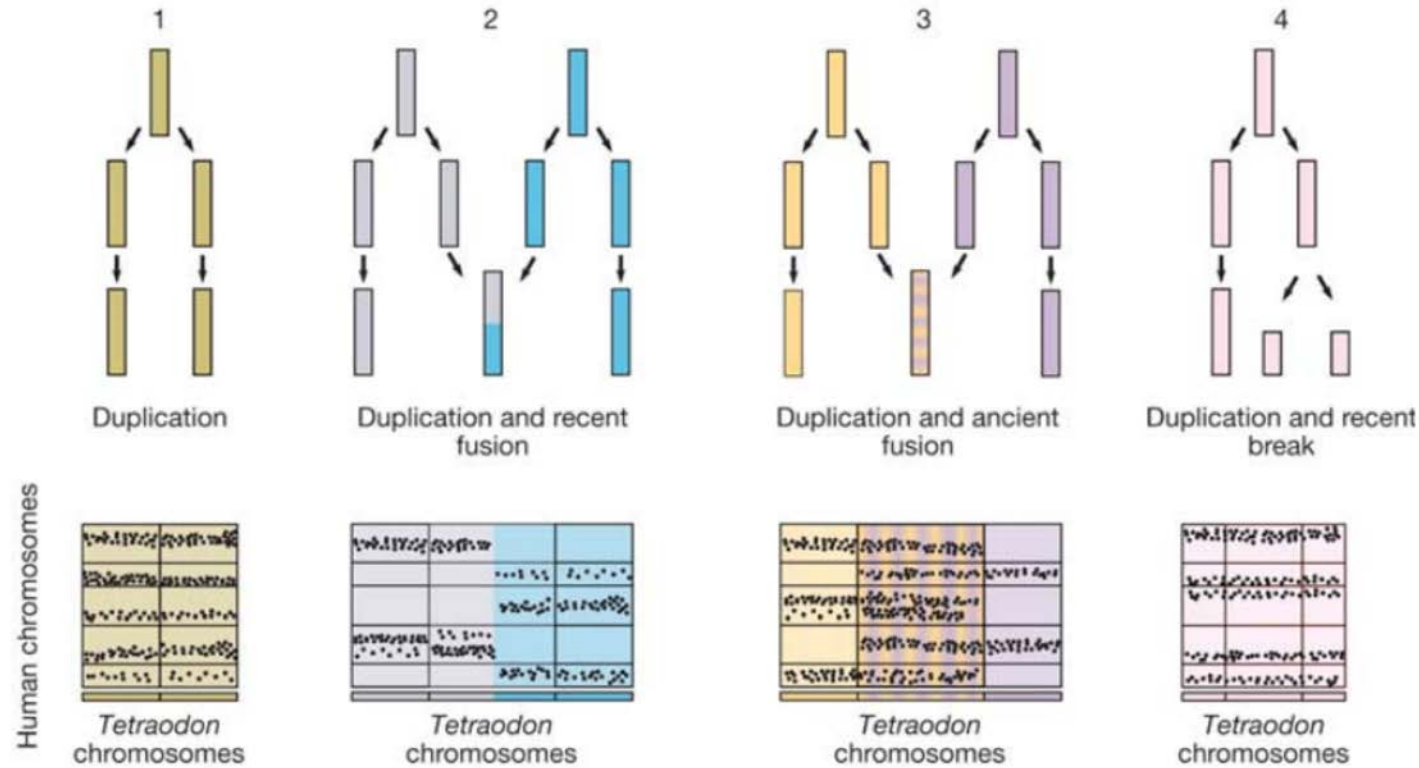


# Reconstructing ancient genome rearrangement

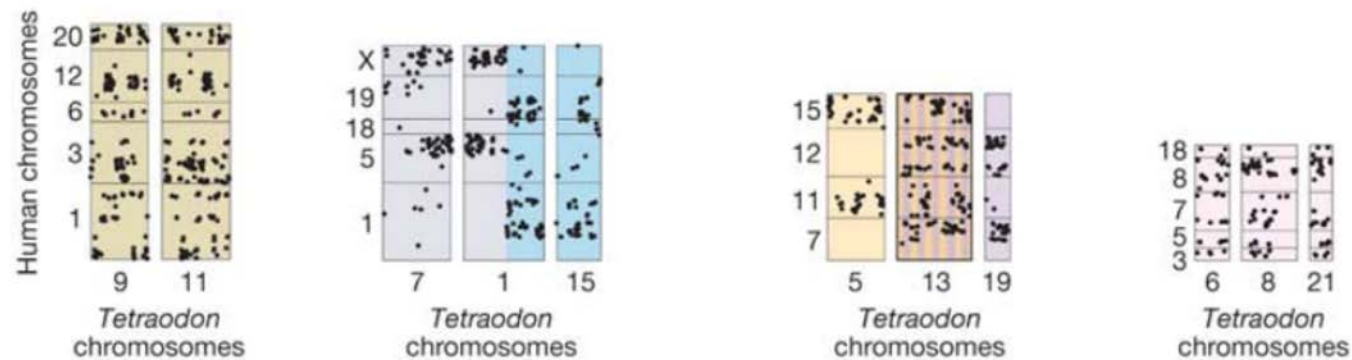


# Reconstructing ancient genome rearrangement

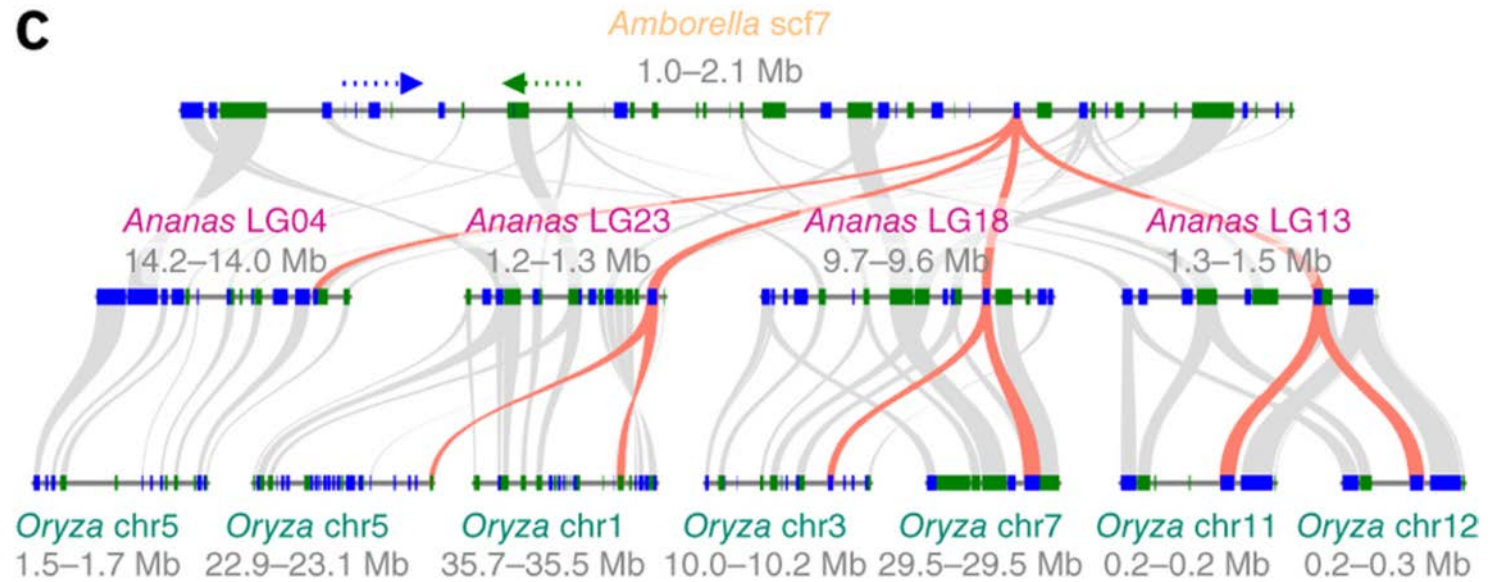
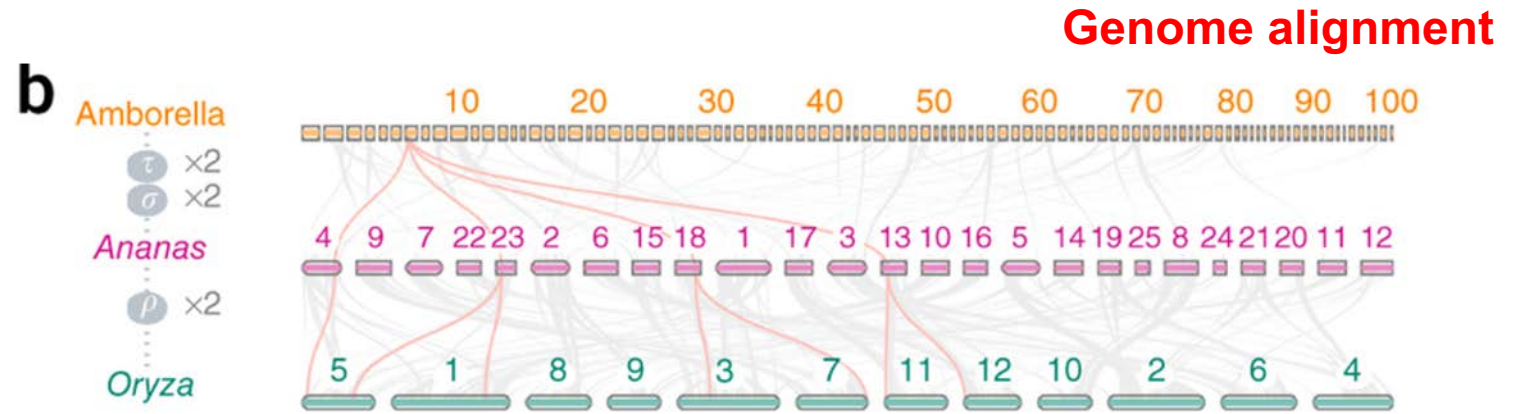
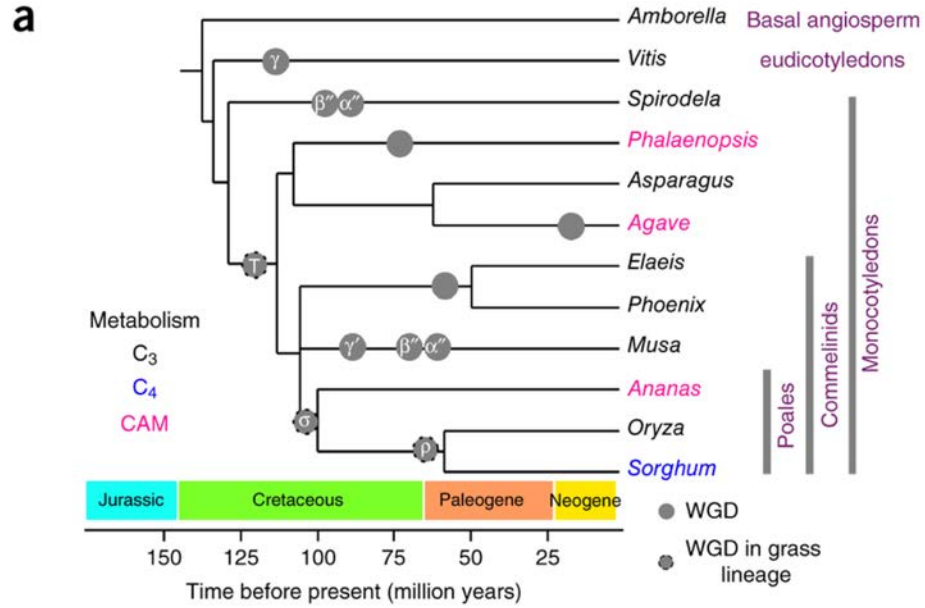
Model of chromosomal evolution



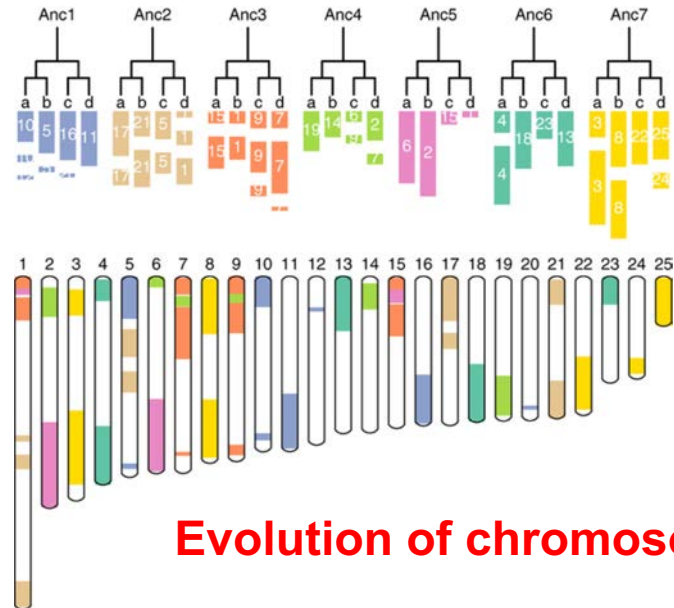
Observed distribution of orthologues between human and *Tetraodon*



# Pineapple genome



**Evolution of chromosomes**



**Colinearity of genes**

Case study: lost of gene families

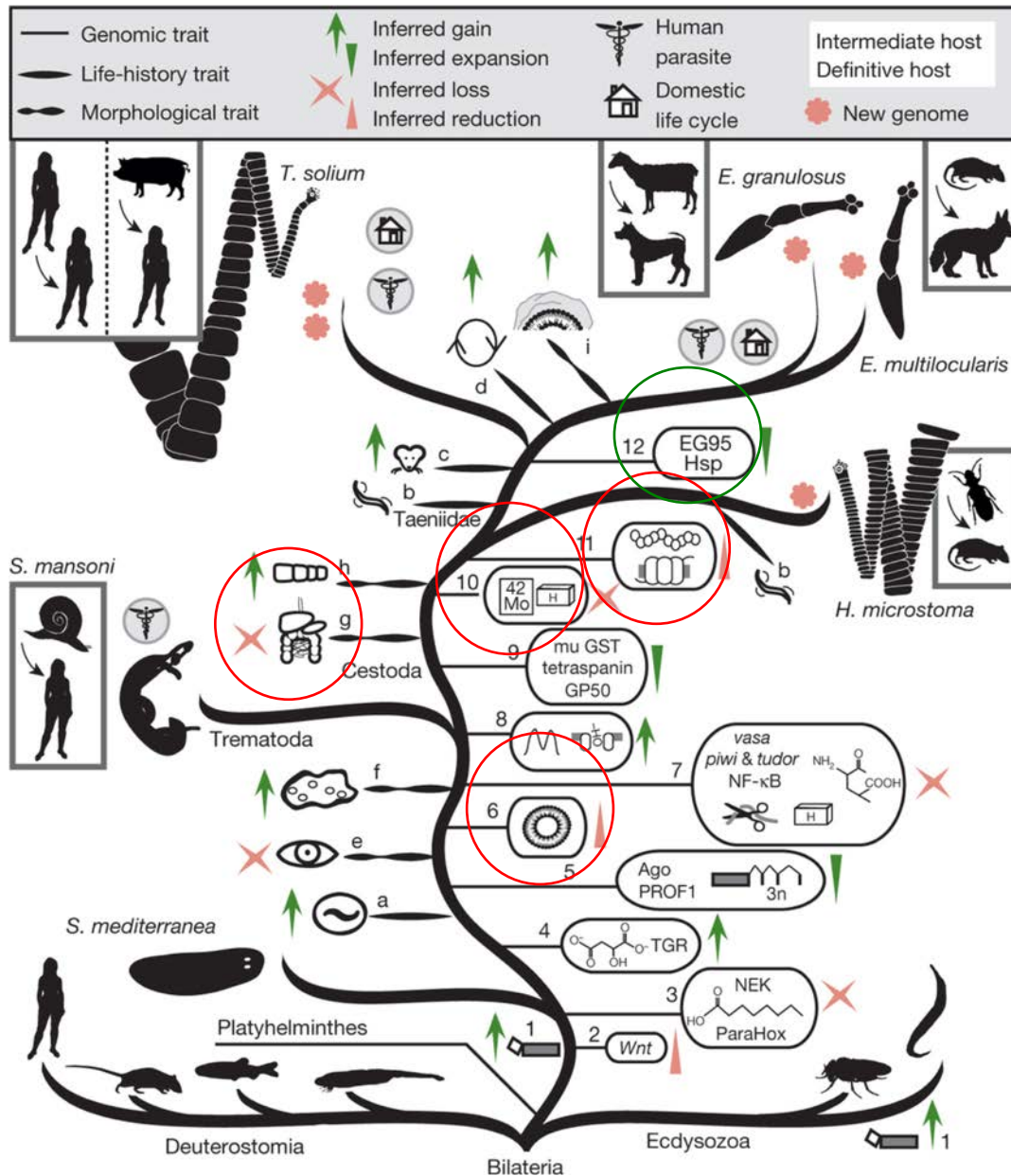
# Comparative genomics of tapeworms

tapeworms

Blood fluke

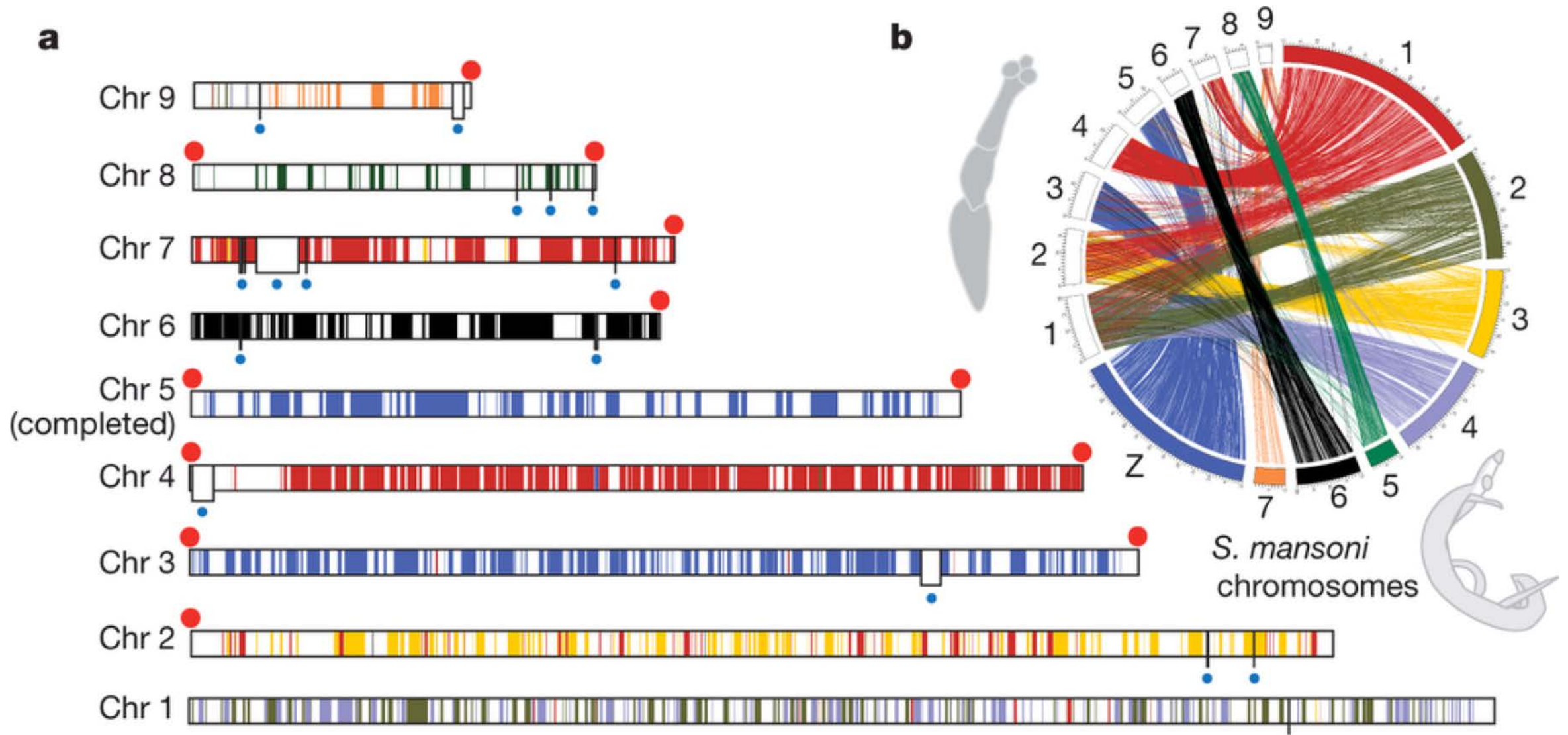
Free-living

Model

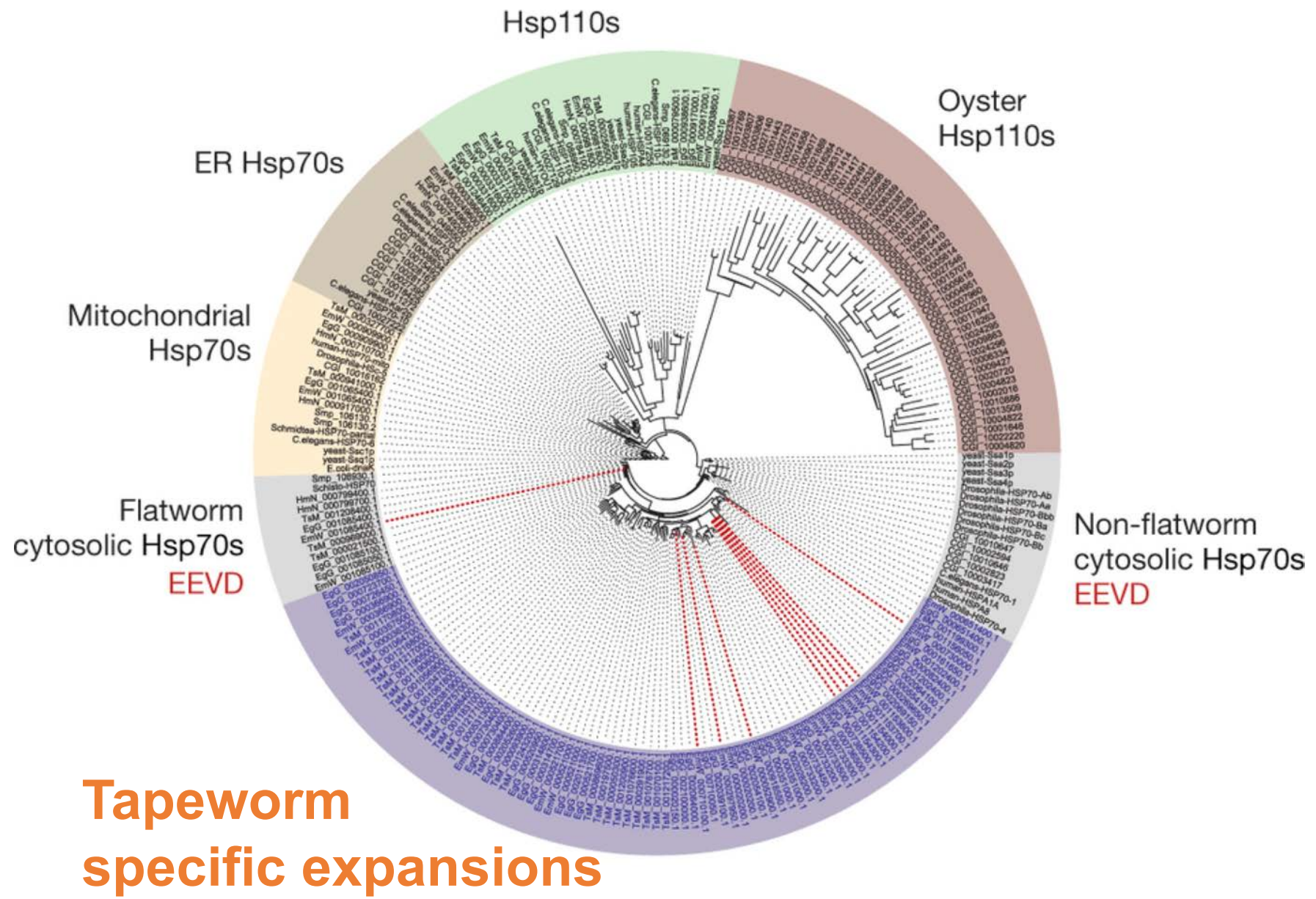


- A total of four tapeworm genomes were sequenced
- We compare with free-living and other parasite genomes
- ‘A route’ to complete parasitism

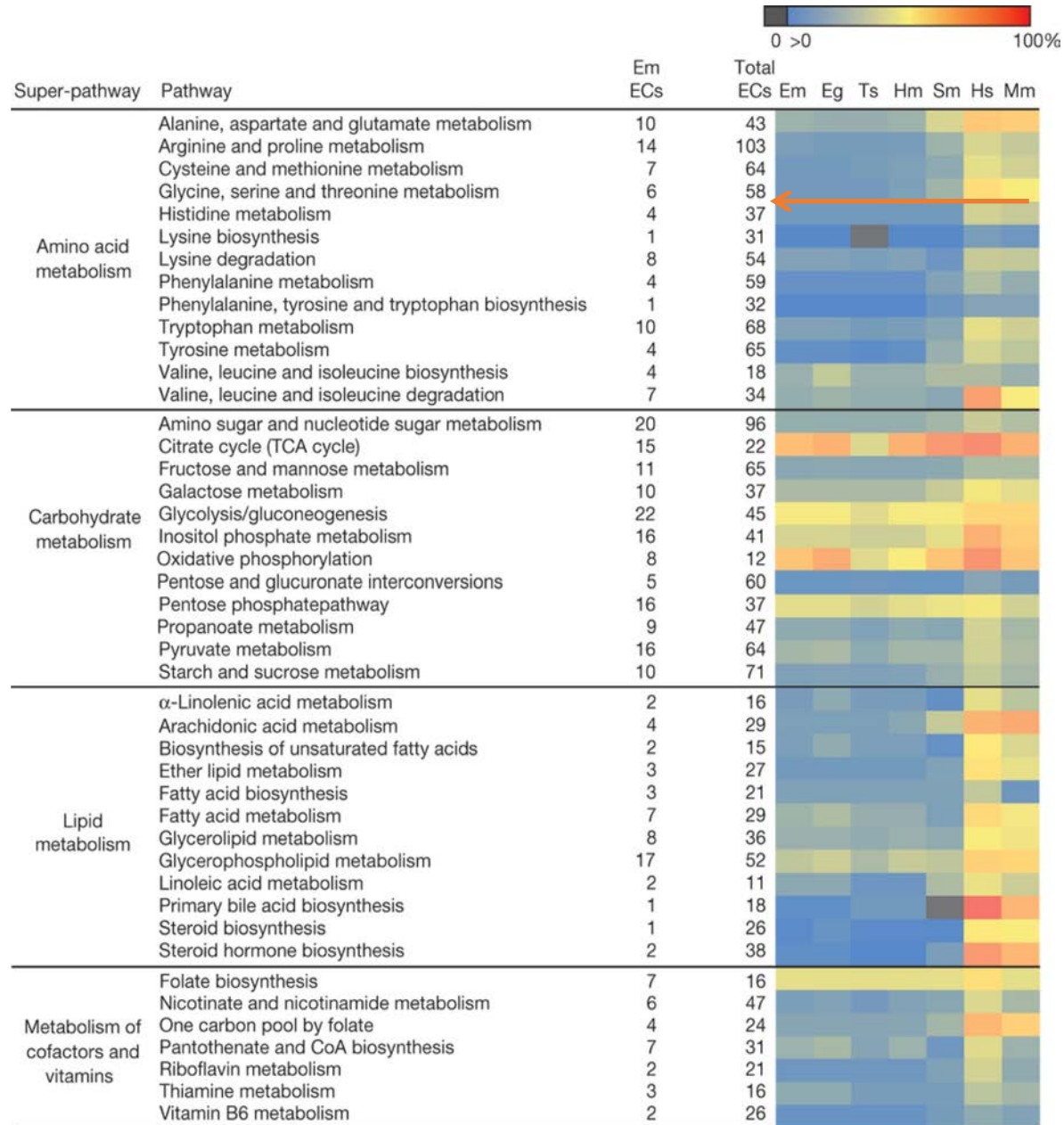
# Genome of *E. multilocularis*



# Heat shock protein expansion in tapeworms



# Reduced metabolism in tapeworms



Reduced metabolism



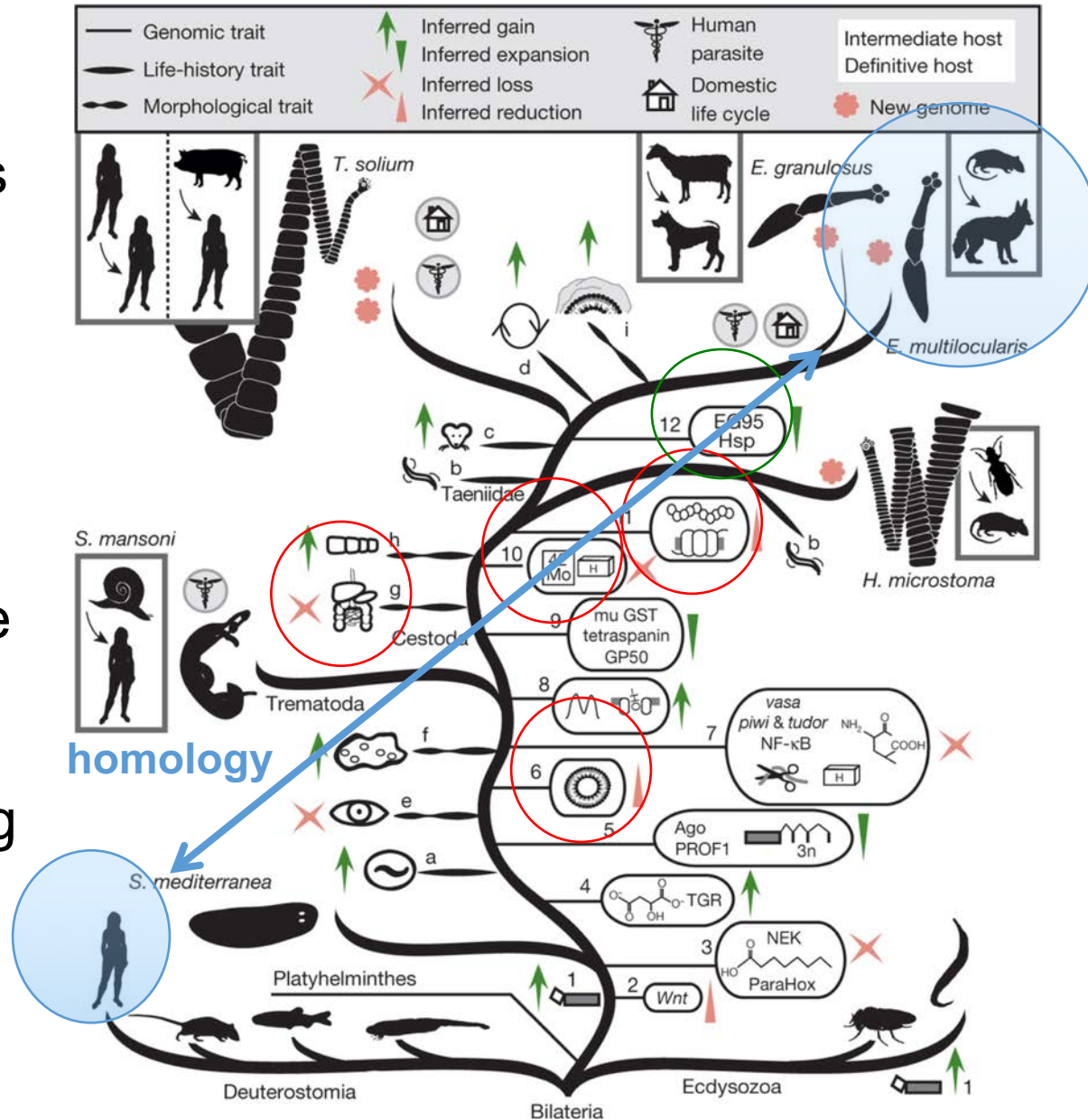
# Tapeworm's road to parasitism

tapeworms

Blood fluke

Free-living

Model

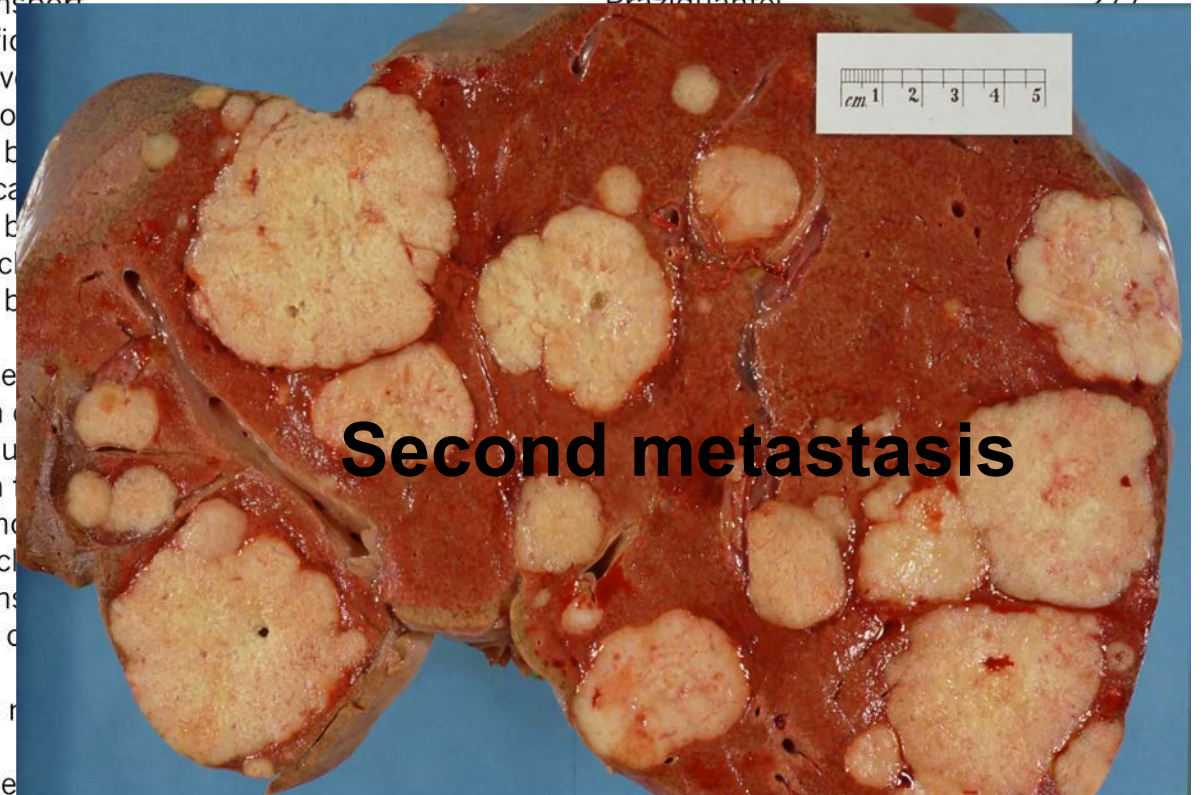


Predict candidate drugs

# Promising drug targets in tapeworms

**Table 1 | Top 20 promising targets in *E. multilocularis***

Target category	Target	Action	Expression	Drug	Rank
Current targets	Tubulin $\beta$ -chain	Cytoskeleton	M,A	Albendazole	406
	Voltage dependent calcium channel	Ion transport		Praziquantel	277



Elongation factor 2	Translation	M,A	Lorazepam)	
Cathepsin B	Protease	M	Experimental compounds	54
Dual-specificity mitogen activated protein	Signalling, activation of p38	M	Experimental compounds	55
Purine nucleoside phosphorylase	Purine metabolism	M,A	Didanosine	56
				63

<http://en.wikipedia.org/wiki/Metastasis>

<http://ocw.tufts.edu/data/>

Comparing genomes beyond gene (copy) numbers

# Extension of homology to genomes

**Gene family gains and losses** in previous lecture

Comparing genomes at **different resolution**

Synteny (gene content on the same chromosome )

Colinearity (gene content + order conservation)

DNA-based alignments (base-to-base mapping)

# Extension of homology to genomes: synteny

## Synteny Conservation and Chromosome Rearrangements During Mammalian Evolution

Jason Ehrlich,<sup>\*.1</sup> David Sankoff<sup>†</sup> and Joseph H. Nadeau<sup>\*,2</sup>

<sup>\*</sup>Jackson Laboratory, Bar Harbor, Maine 04609 and <sup>†</sup>Centre de recherches mathématiques,  
Université de Montréal, Montréal, Québec, H3C 3J7 Canada

Manuscript received December 13, 1996

Accepted for publication June 4, 1997

## *MAPS of LINKAGE and SYNTENY HOMOLOGIES between MOUSE and MAN*

JOSEPH H. NADEAU

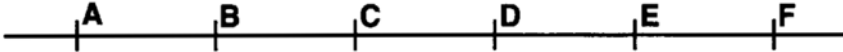
1989

*Synteny* refers to the occurrence of two or more genes on the same chromosome, whereas *conserved synteny* refers to two or more homologous genes that are syntenic in two or more species, regardless of gene order on each chromosome, i.e., synteny but not necessarily gene order is conserved (Figure 2; see also NADEAU 1989). *Conserved linkage* pertains to the conservation of both synteny and order of homologous genes between species (Figure 2; see also NADEAU 1989). A *disrupted synteny* refers to circumstances where a pair of genes are located on the same chromosome in one species but their homologues are located on different chromosomes in another species, i.e., the genes are syntenic in only one of the two species. Syntenic genes can be identified by examining published genetic maps and conserved segments can be identified by comparing

# Synteny

conservation of gene content

## A. Genetic map in reference species



Each unit is gene

### Conserved synteny and linkage

Gene arrangement:



Definition: Same gene order and similar genetic distances.

Count:  
 One conserved linkage involving genes A,B,C,E,F.  
 one conserved synteny involving genes A,B,C,E,F.

Possible cause:  
 No inter-chromosomal rearrangement.  
 No intra-chromosomal rearrangement.

### Conserved synteny, conserved linkage, disrupted linkage

Gene arrangement:

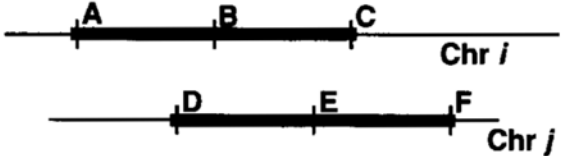


Count:  
 One conserved linkage involving genes B,C,D;  
 One conserved linkage involving genes E,F.  
 One disrupted linkage involving genes B,C,D vs E,F vs A.  
 One conserved synteny involving genes A,B,C,D,E,F.

Possible causes:  
 An intra-chromosomal rearrangement,  
 such as a paracentric inversion.

### Conserved synteny, disrupted synteny, conserved linkage, disrupted linkage

Gene arrangement:

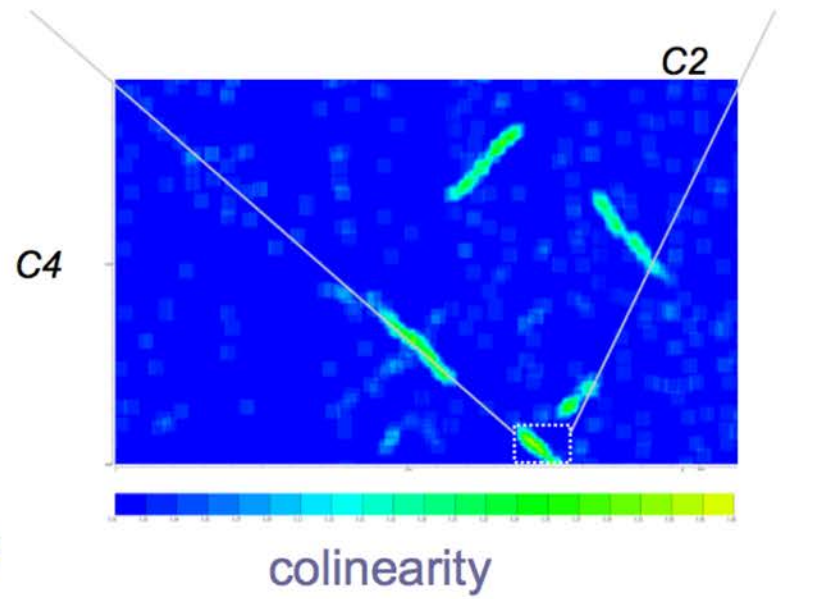
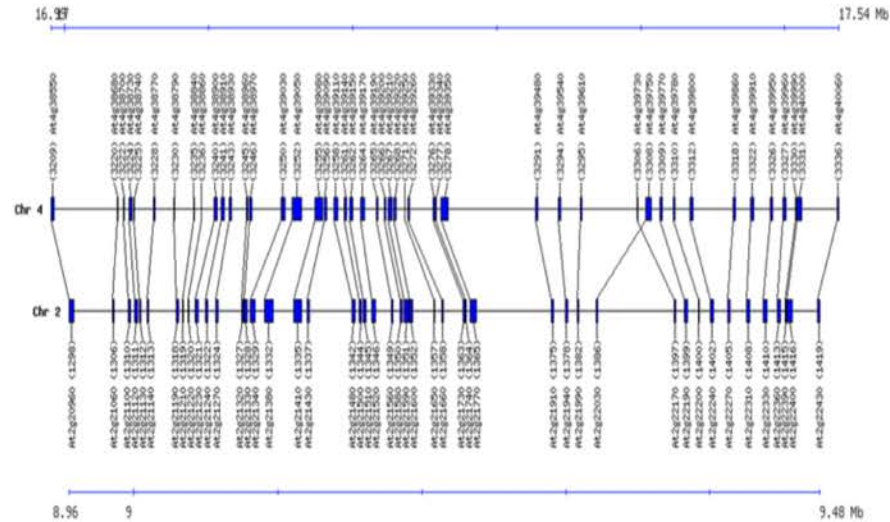


Count:  
 One conserved linkage involving genes A,B,C;  
 One conserved linkage involving genes D,E,F.  
 One disrupted linkage involving genes A,B,C vs D,E,F.  
 One conserved synteny involving genes A,B,C.  
 One conserved synteny involving genes D,E,F.  
 One disrupted synteny involving genes A,B,C vs D,E,F.

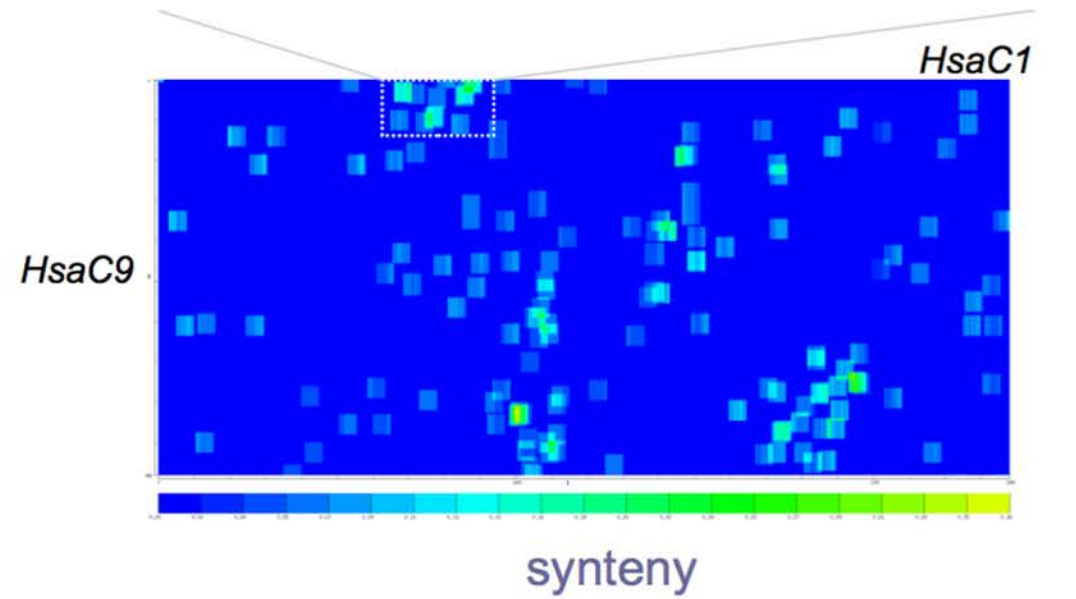
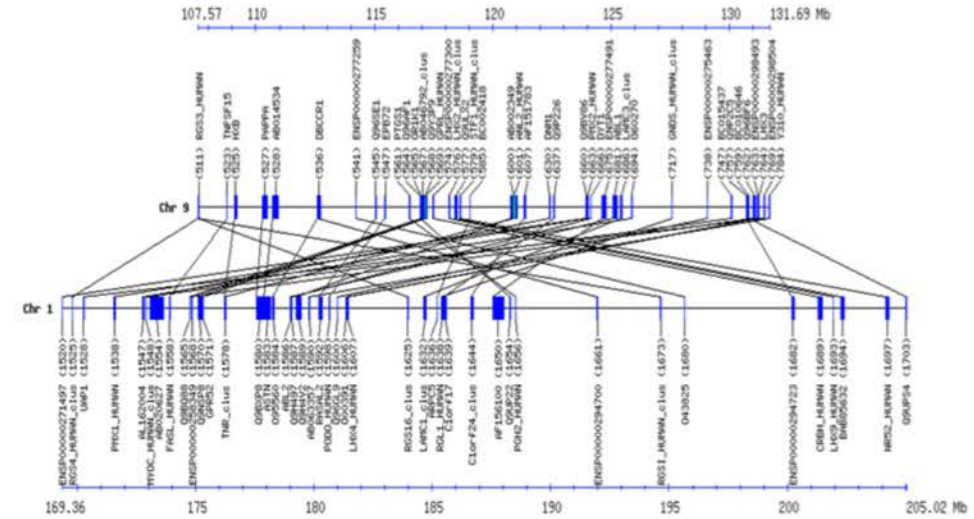
Possible causes:  
 An inter-chromosomal rearrangement,  
 such as a reciprocal translocation.

# Synteny and colinearity

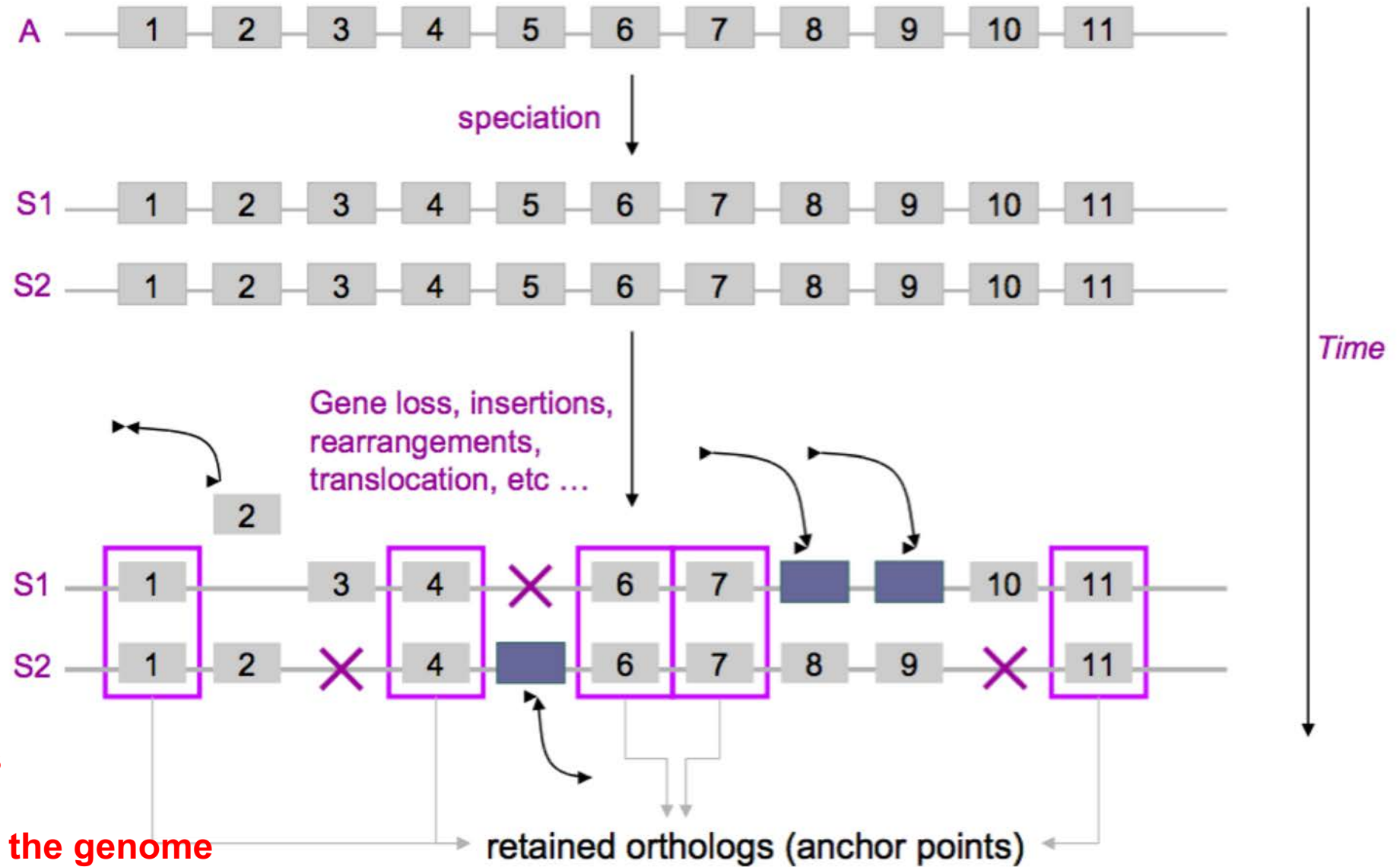
recent duplication



ancient duplication



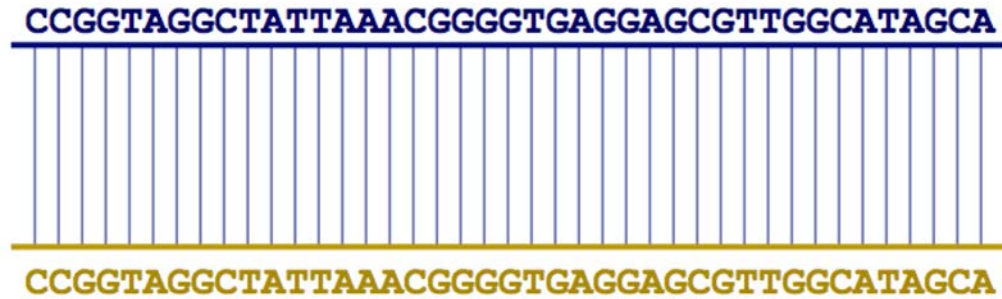
# Inferring gene collinearity



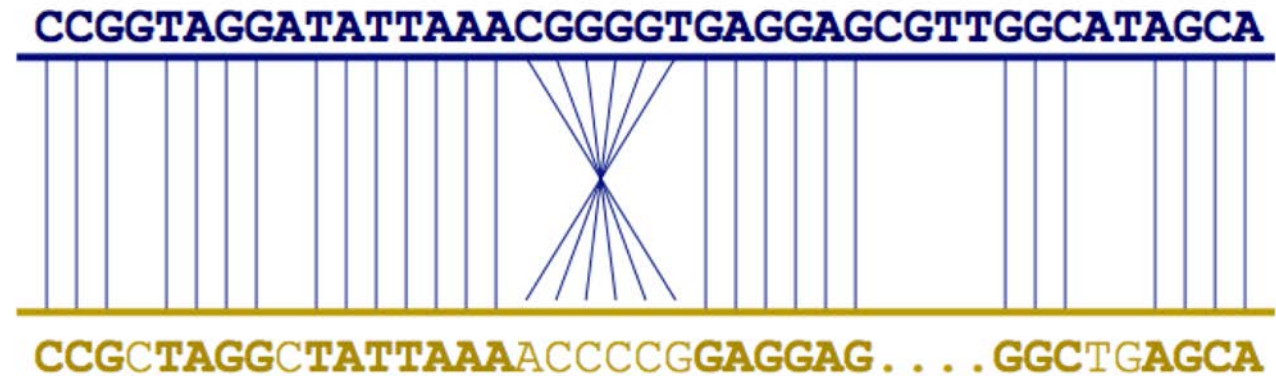


# Whole genome alignment

For two genomes, A and B, find a mapping from each position in A to its corresponding position in B

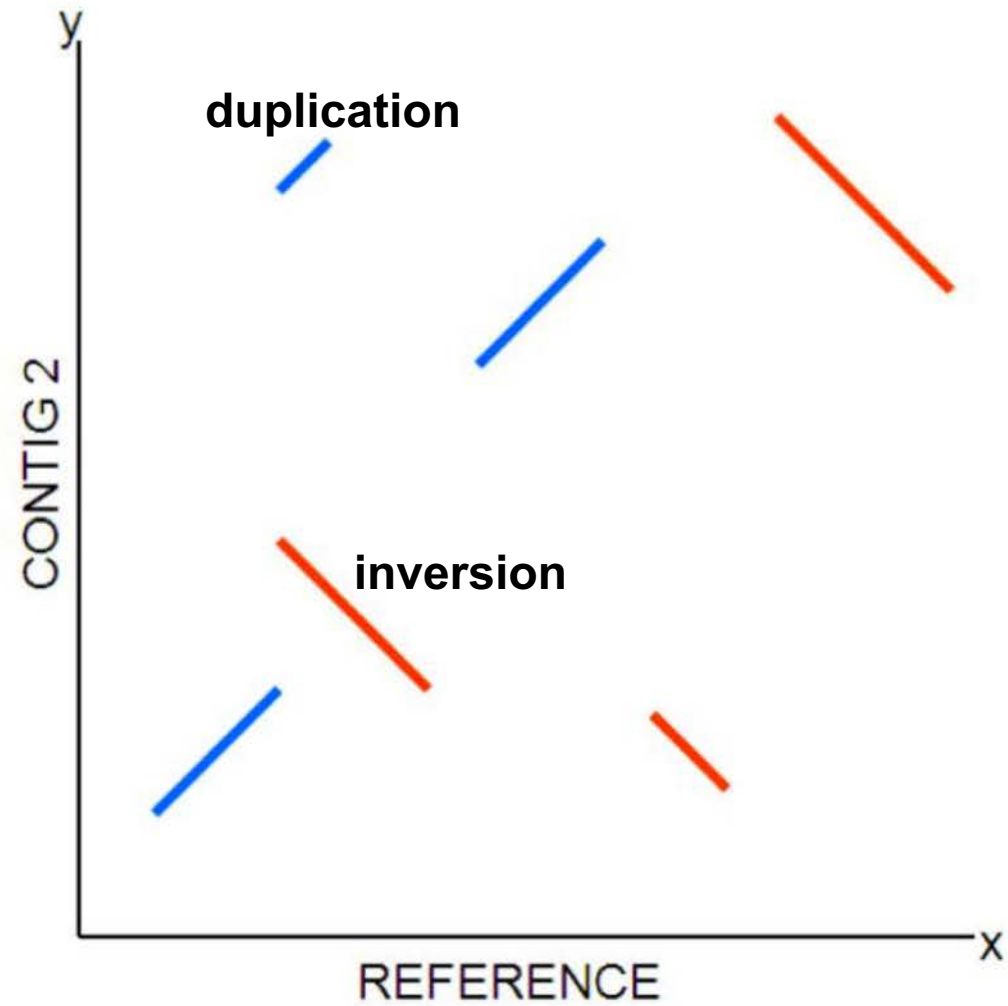
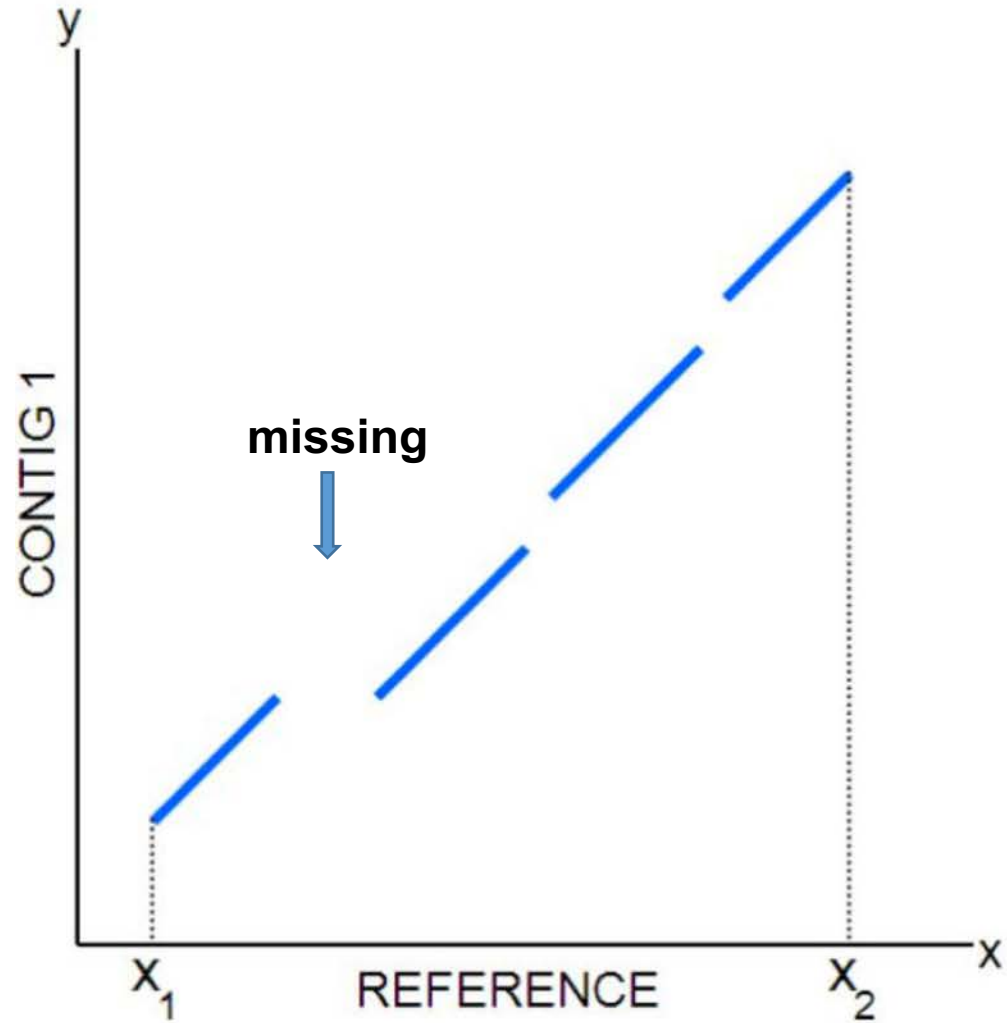


**In reality**, Genome A may have insertions, deletions, translocations, inversions, duplications or SNPs with respect to B (sometimes all of the above)



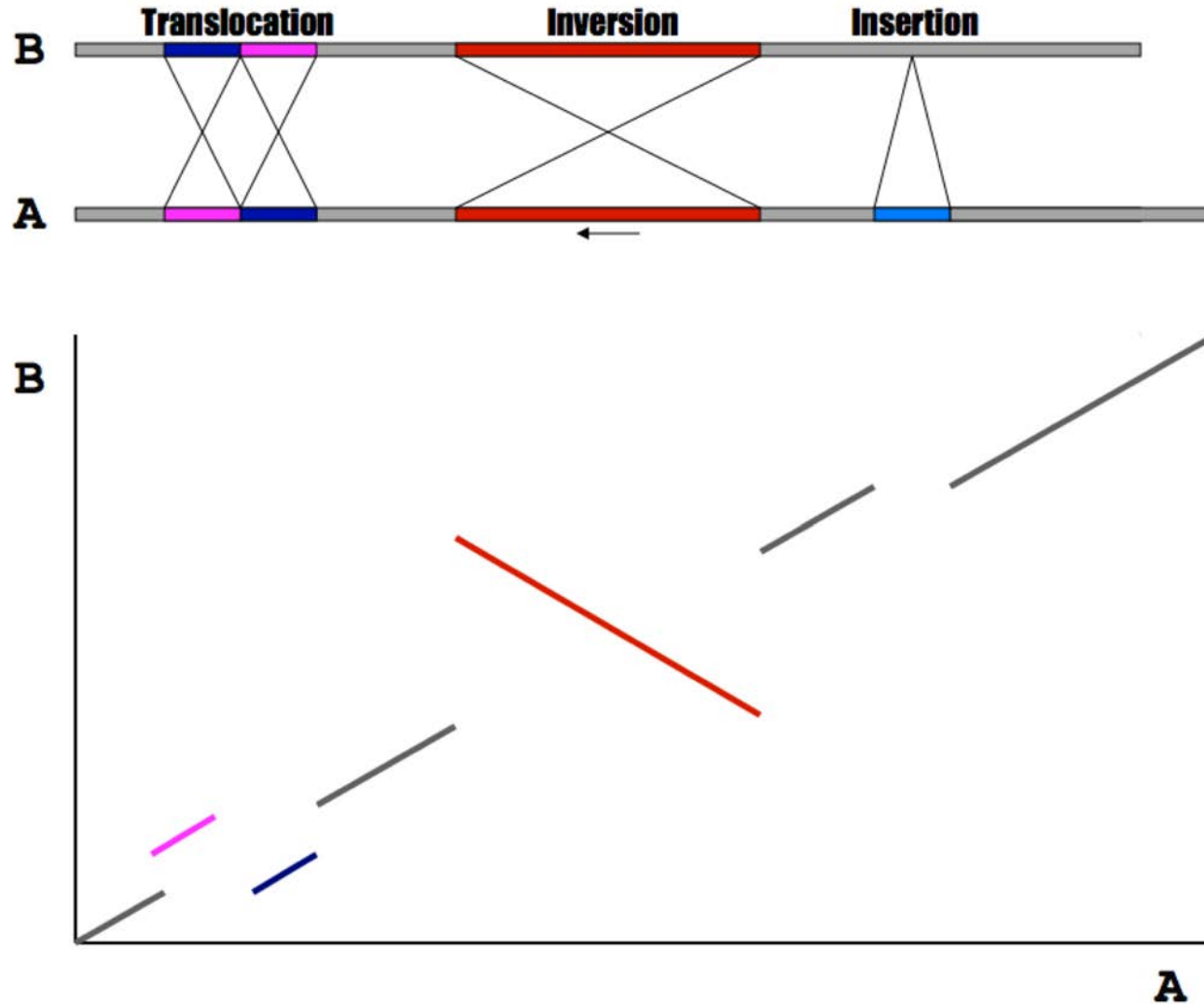
# Aligning genome at nucleotide / amino acid level

## Visualise through **dotplot**

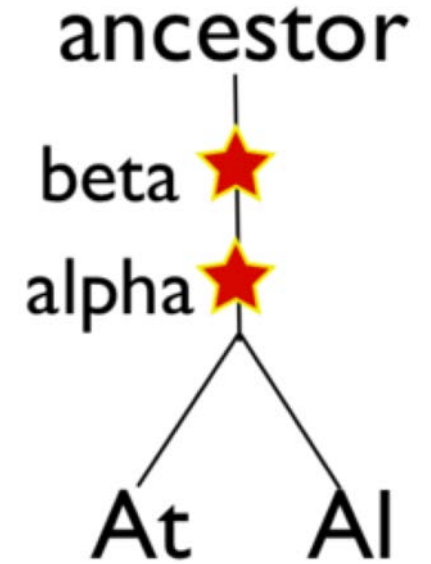
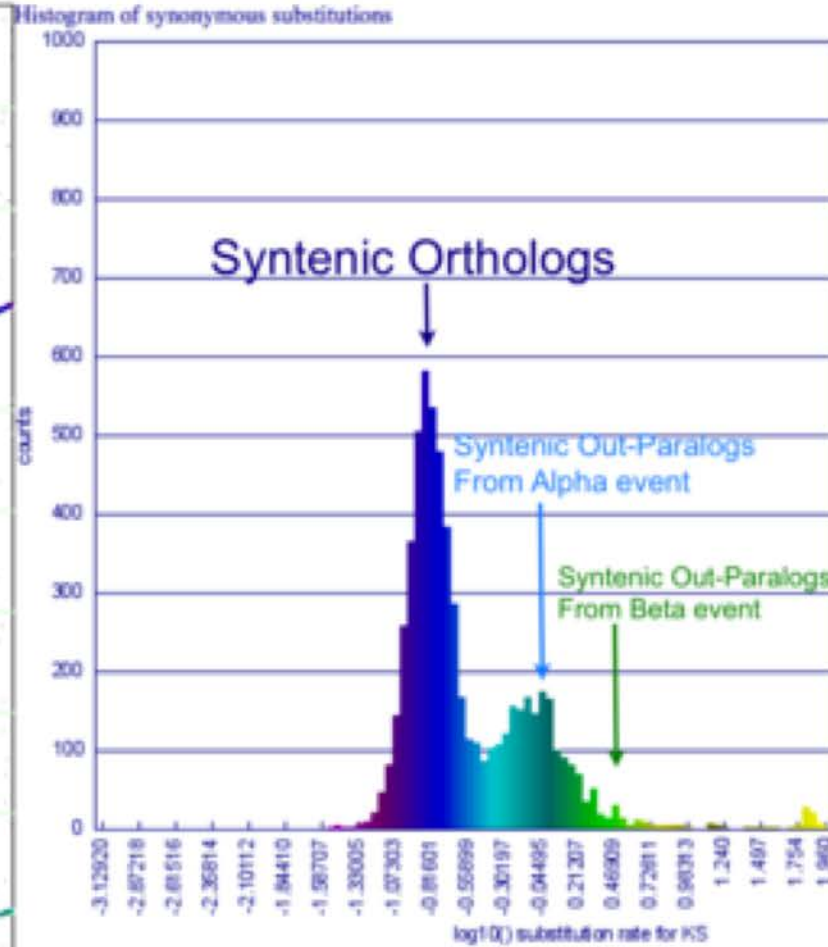
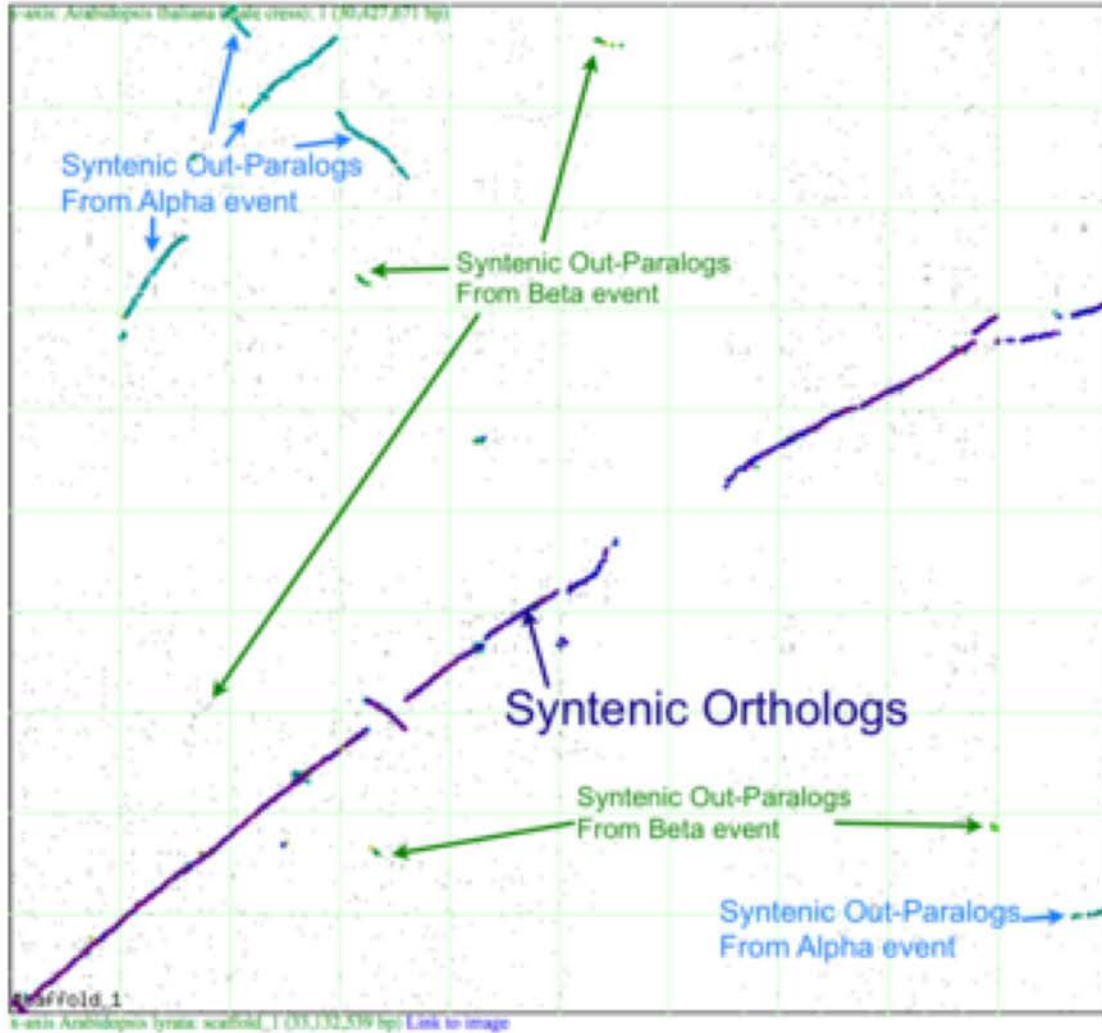


# Aligning genome at nucleotide / amino acid level

## Visualise through **dotplot**

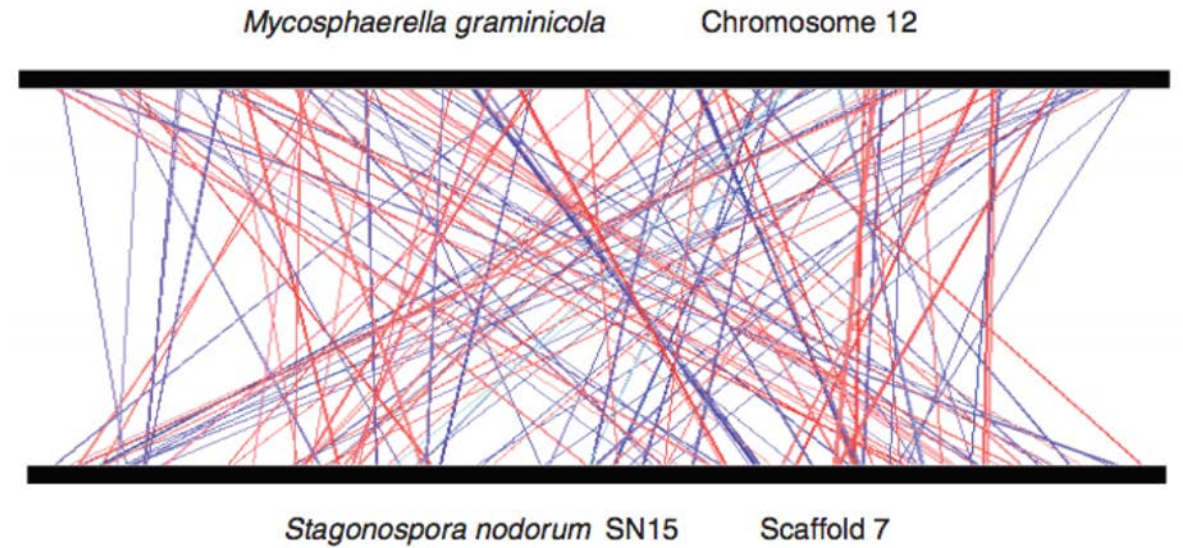
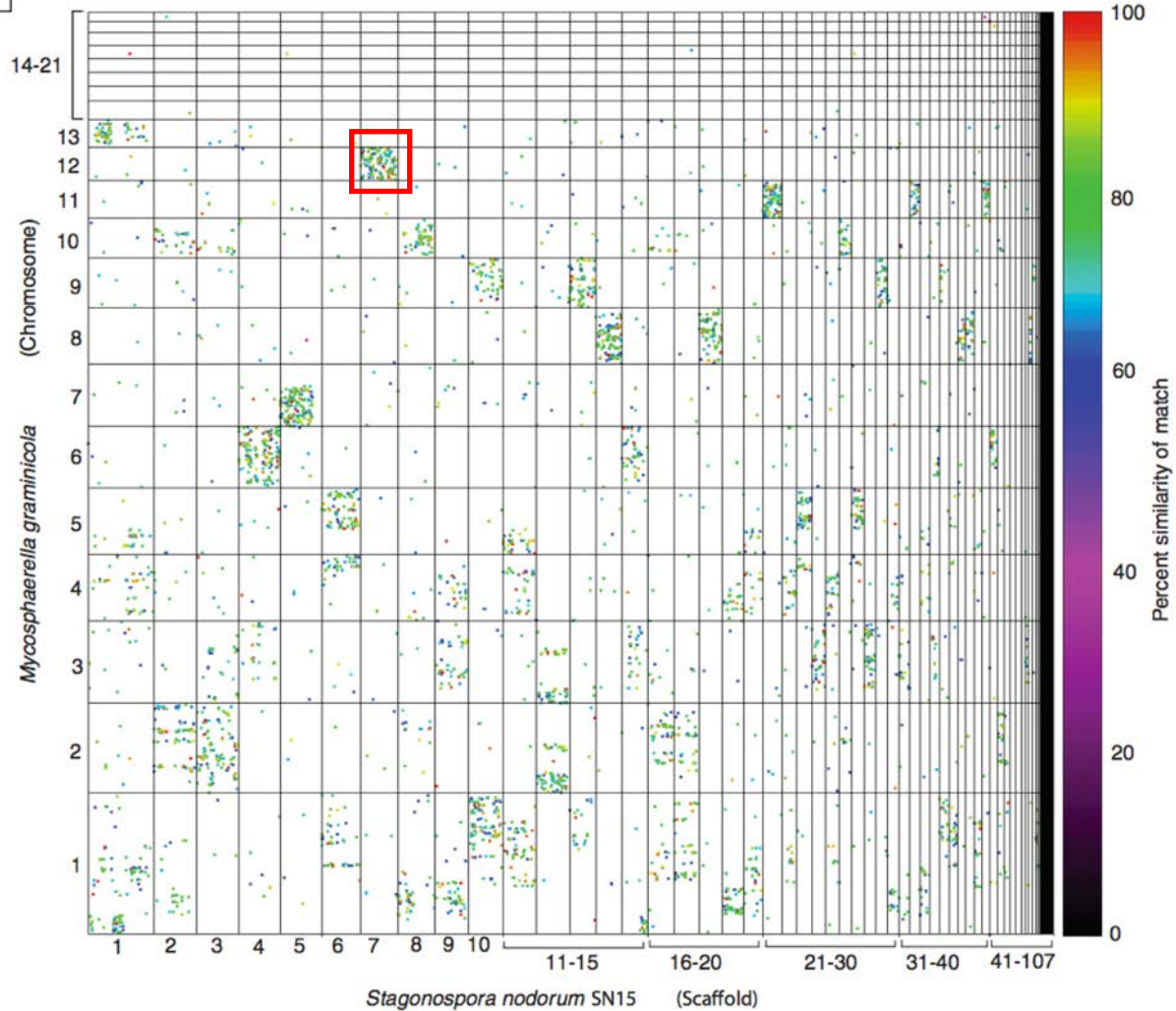


# Relationship between genome synteny, syntenic orthologs and duplications

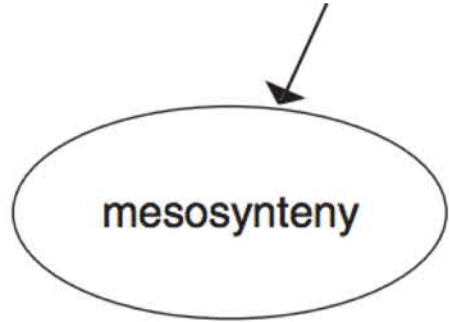


# Relationship between genome synteny, syntenic orthologs and duplications

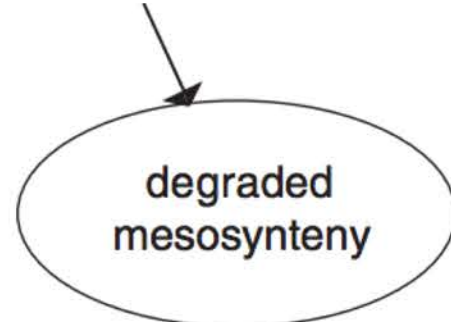
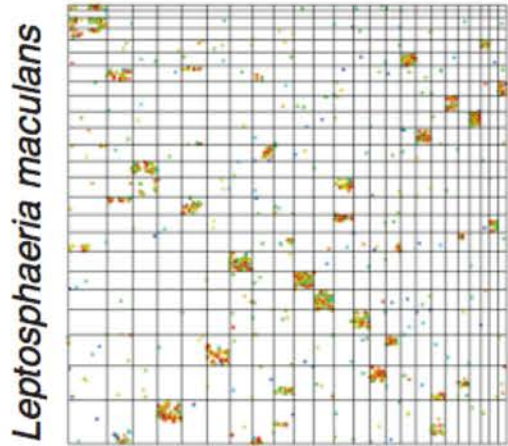
(a)



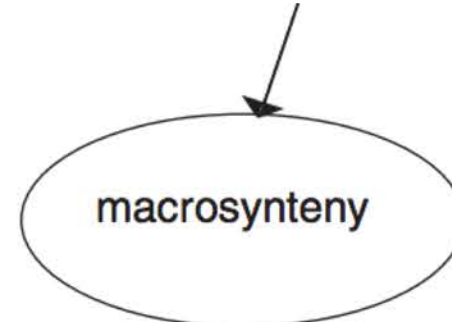
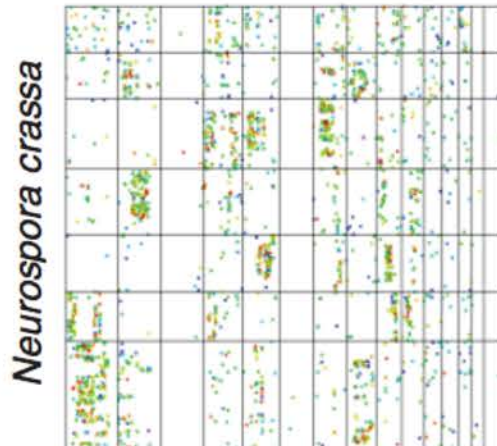
# Different kinds of genome synteny



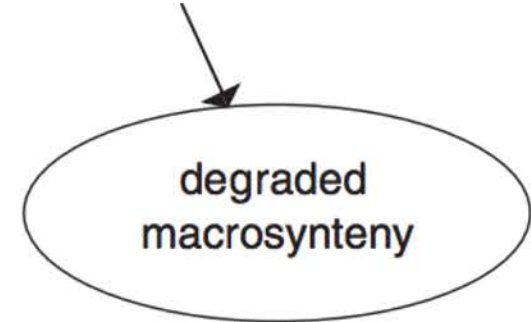
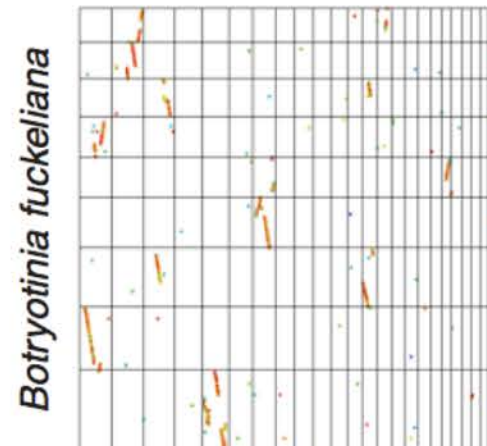
*Phaeosphaeria nodorum*



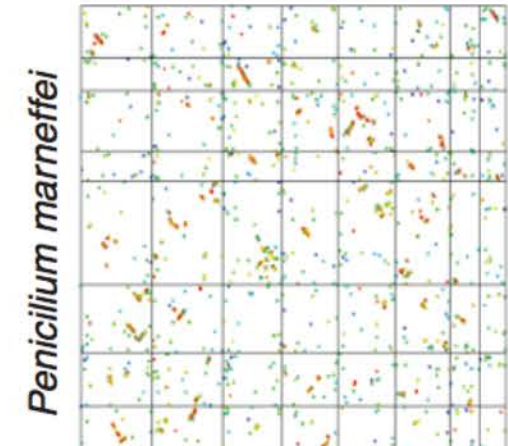
*Fusarium oxysporum*



*Sclerotinia sclerotiorum*



*Aspergillus fumigatus*

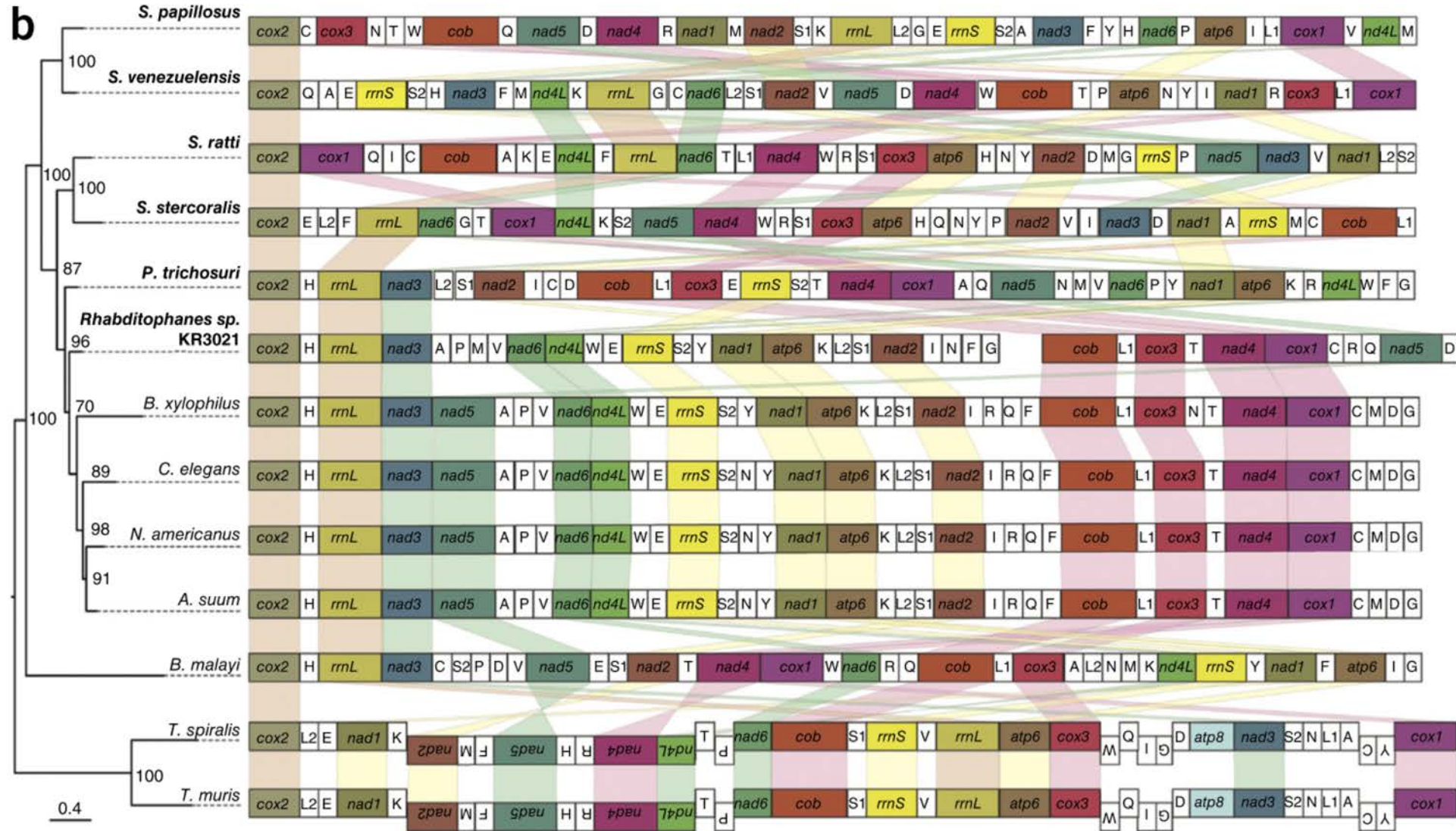


genes are conserved within homologous chromosomes, but with randomized orders and orientations

genes are conserved within homologous chromosomes, and with colinear gene regions

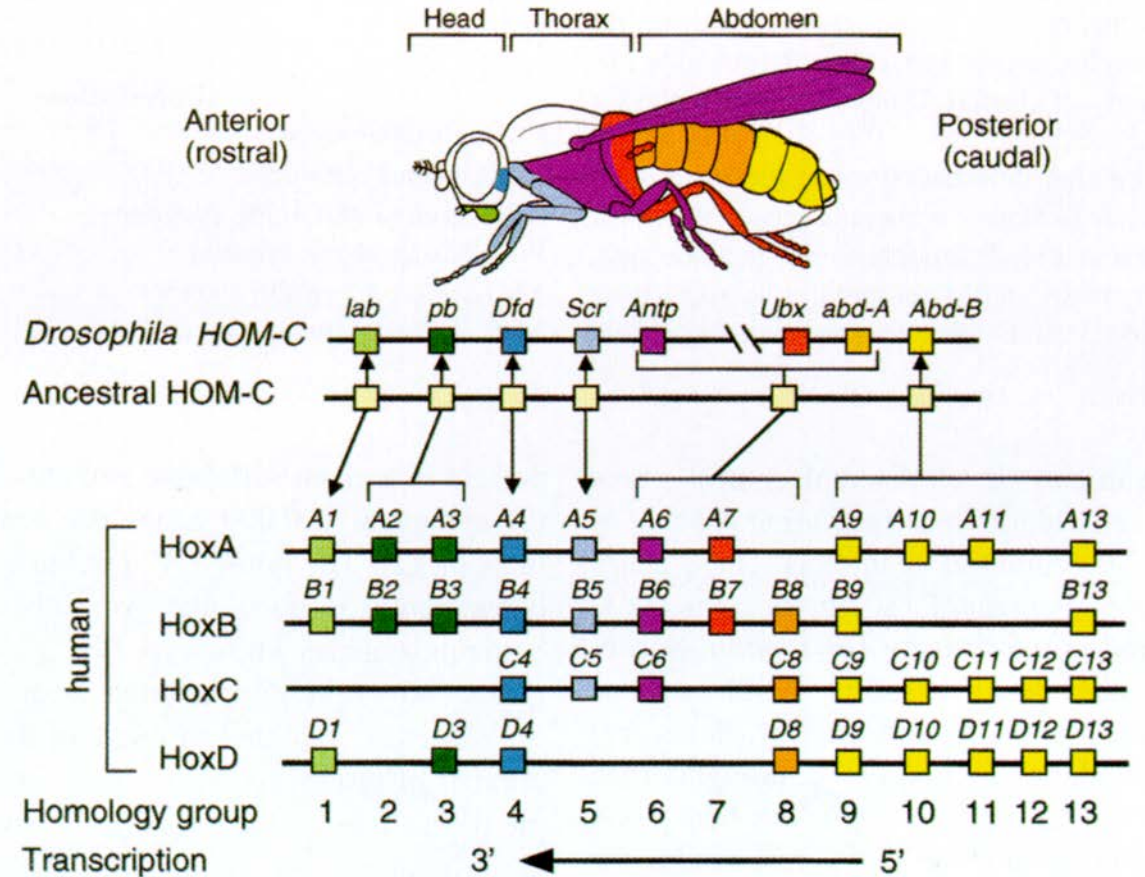
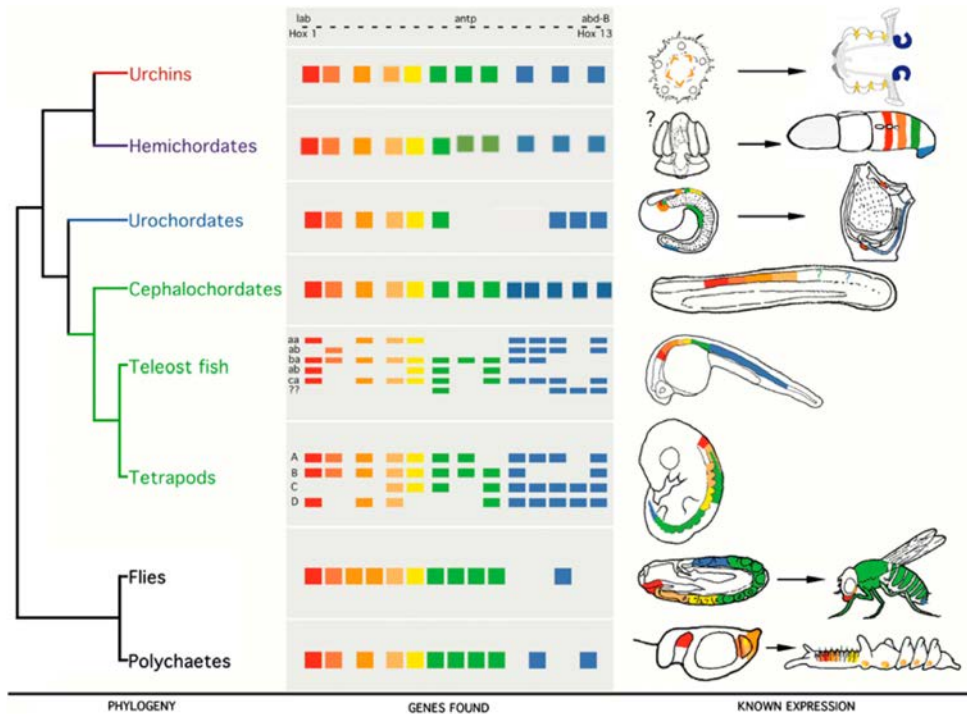
# Why are we interested in synteny and collinearity?

Establish relationship between species



# Why are we interested in synteny and collinearity?

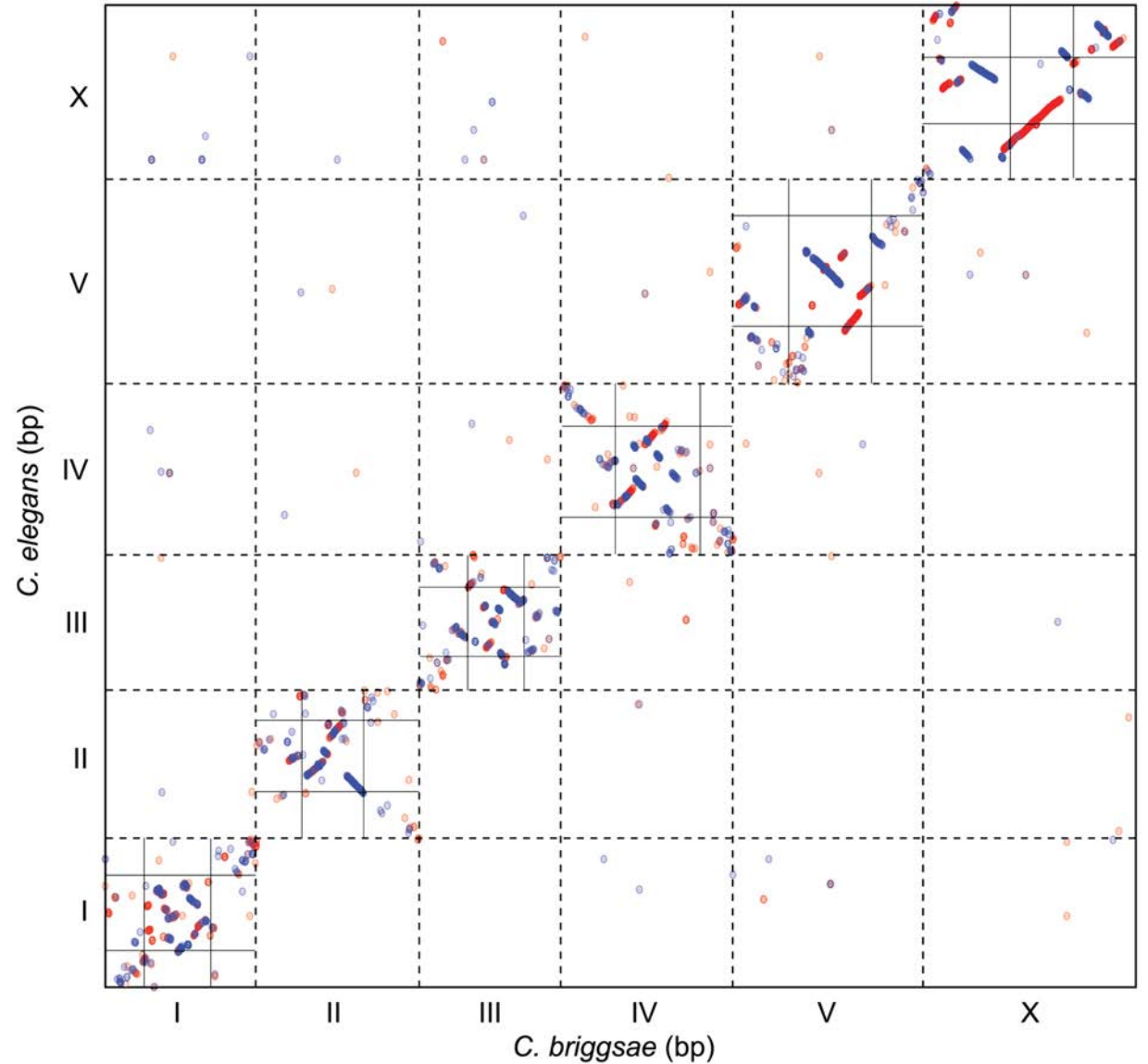
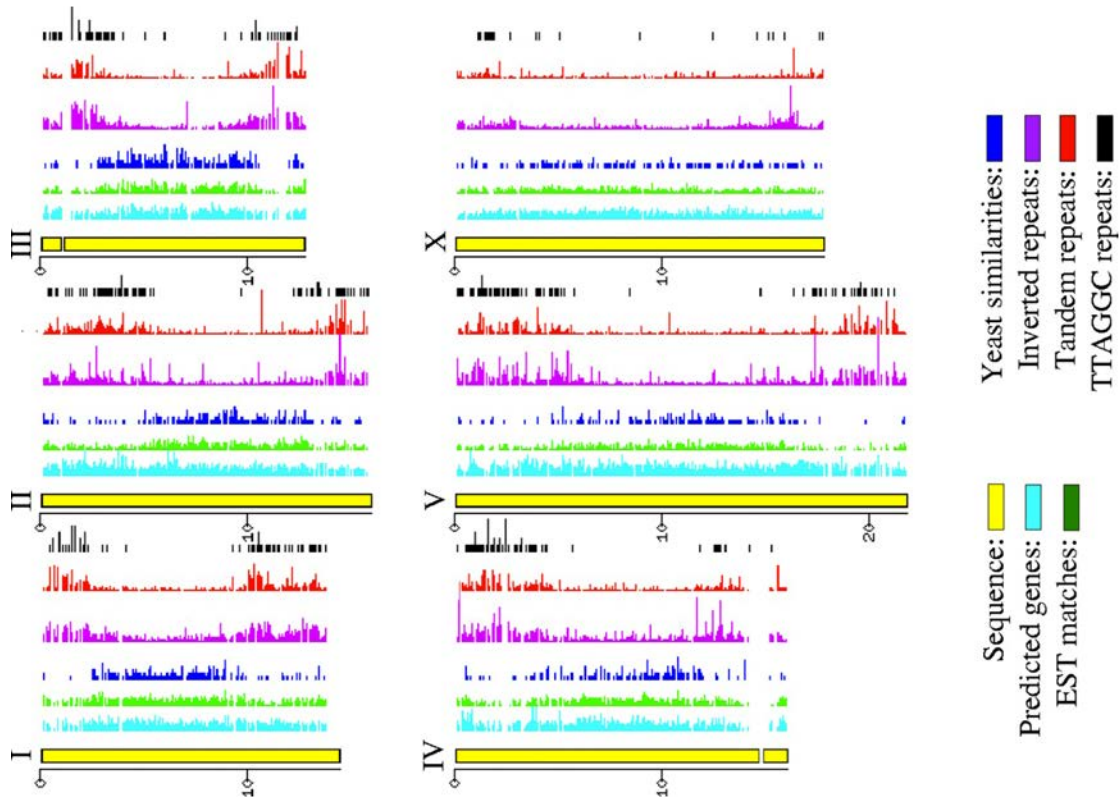
**Evolutionary conserved features (orthologs, synteny, collinearity)** are good indicators of functionally important genome regions





# Why are we interested in synteny and collinearity?

**Evolutionary conserved features (orthologs, synteny, collinearity) relate to genome biology**

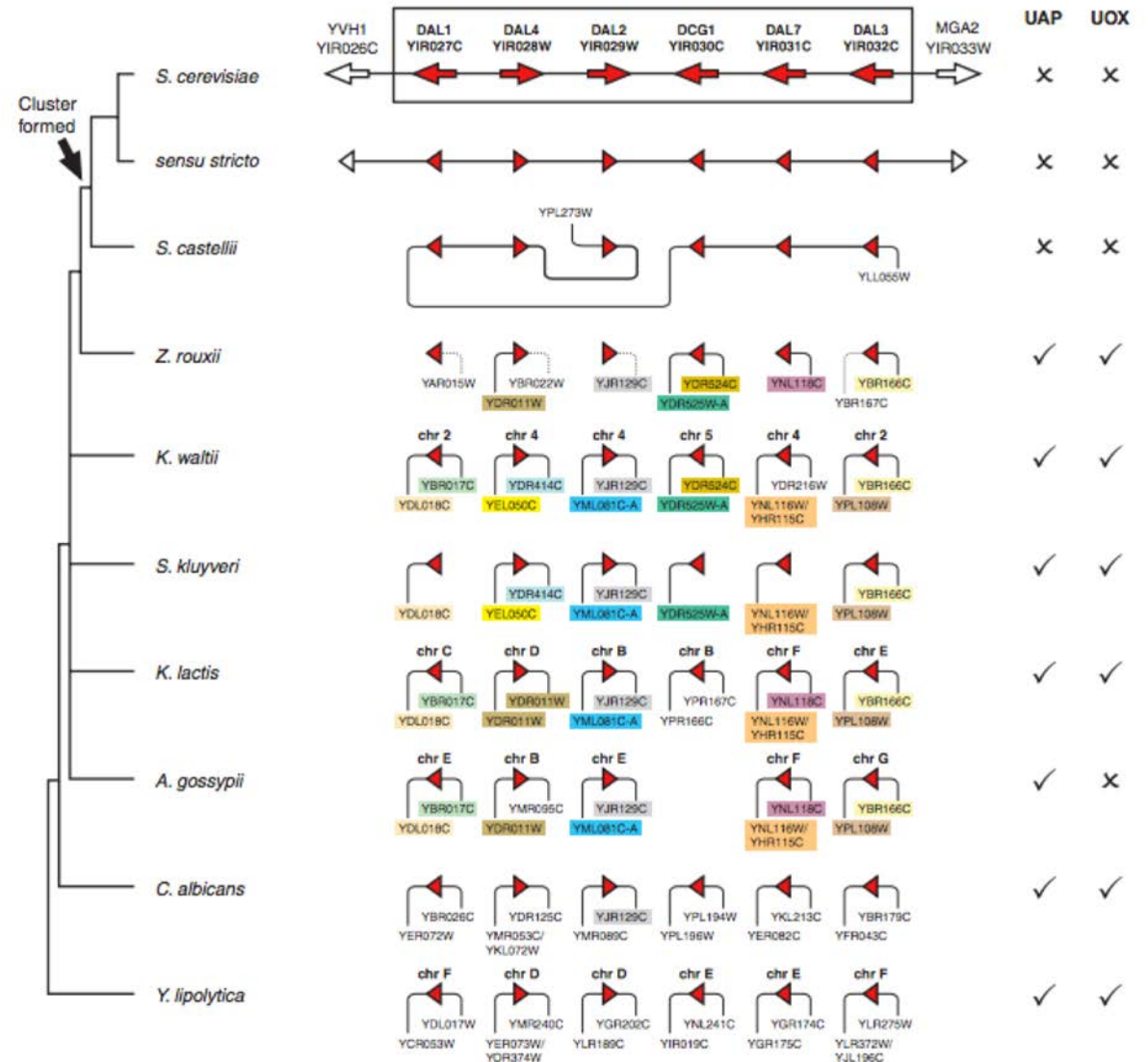


Stein *et al.*, PLOS Biology 2003

The *C. elegans* Sequencing Consortium Science 1998

# Why are we interested in synteny and collinearity?

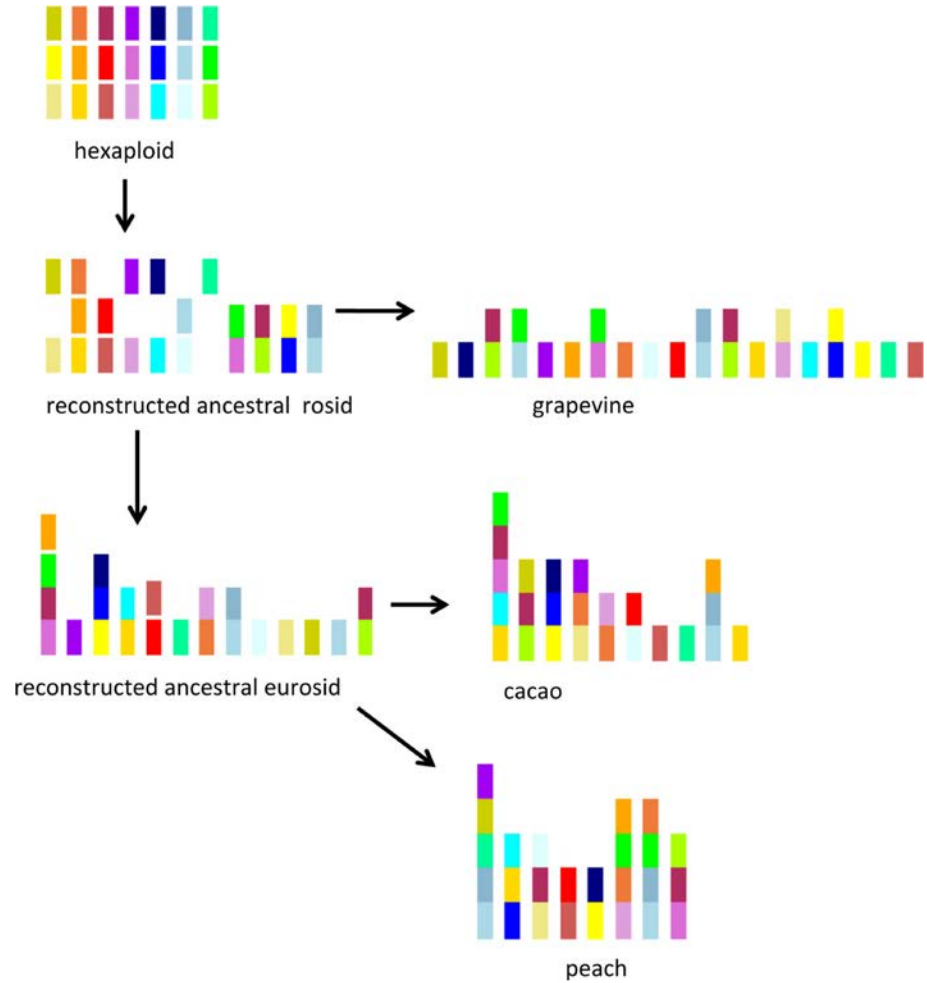
We can **reconstruct evolutionary histories of gene & gene families** and eventually lead to functioning of species



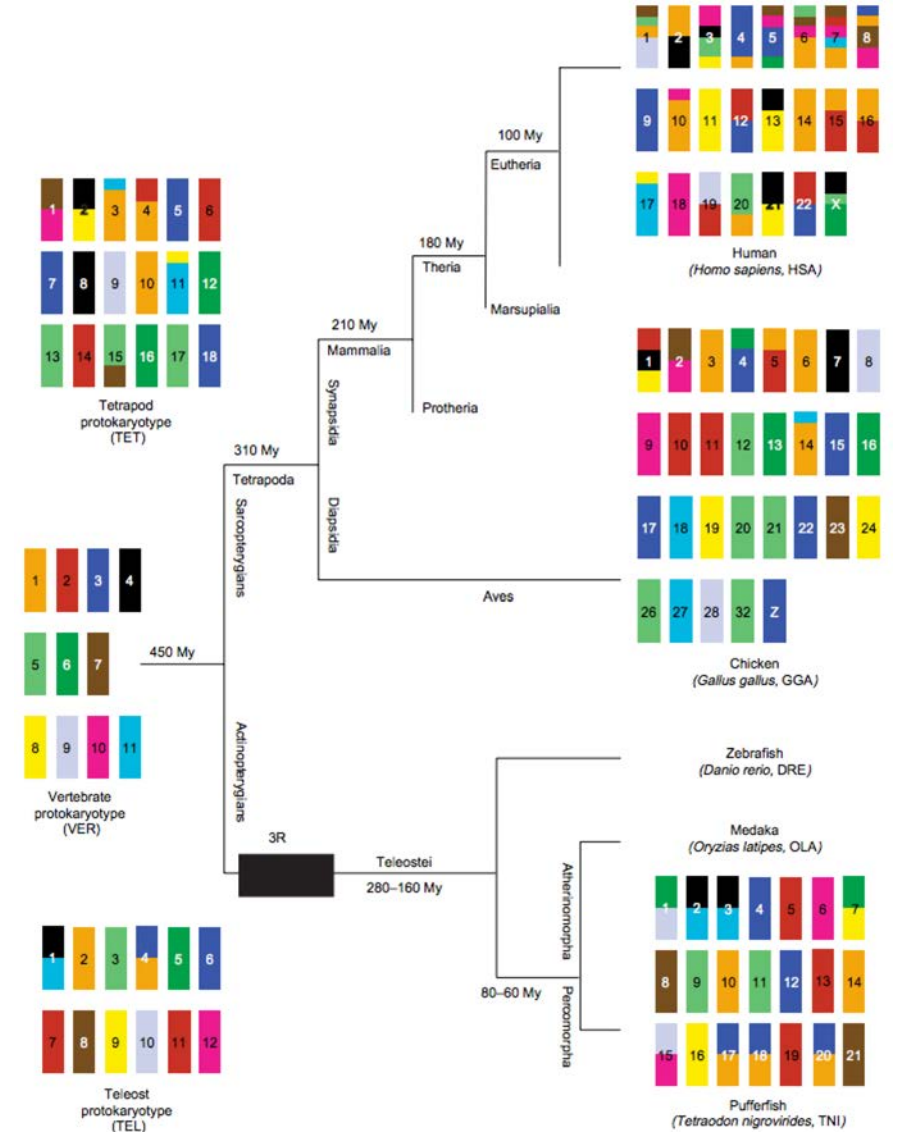
Birth of a metabolic gene cluster in yeast by adaptive gene relocation

# Why are we interested in synteny and collinearity?

We can **reconstruct ancient karyotypes** that eventually lead to better understanding of evolution of species



Zheng et al (2013)



Kohn et al (2006)

Case study:

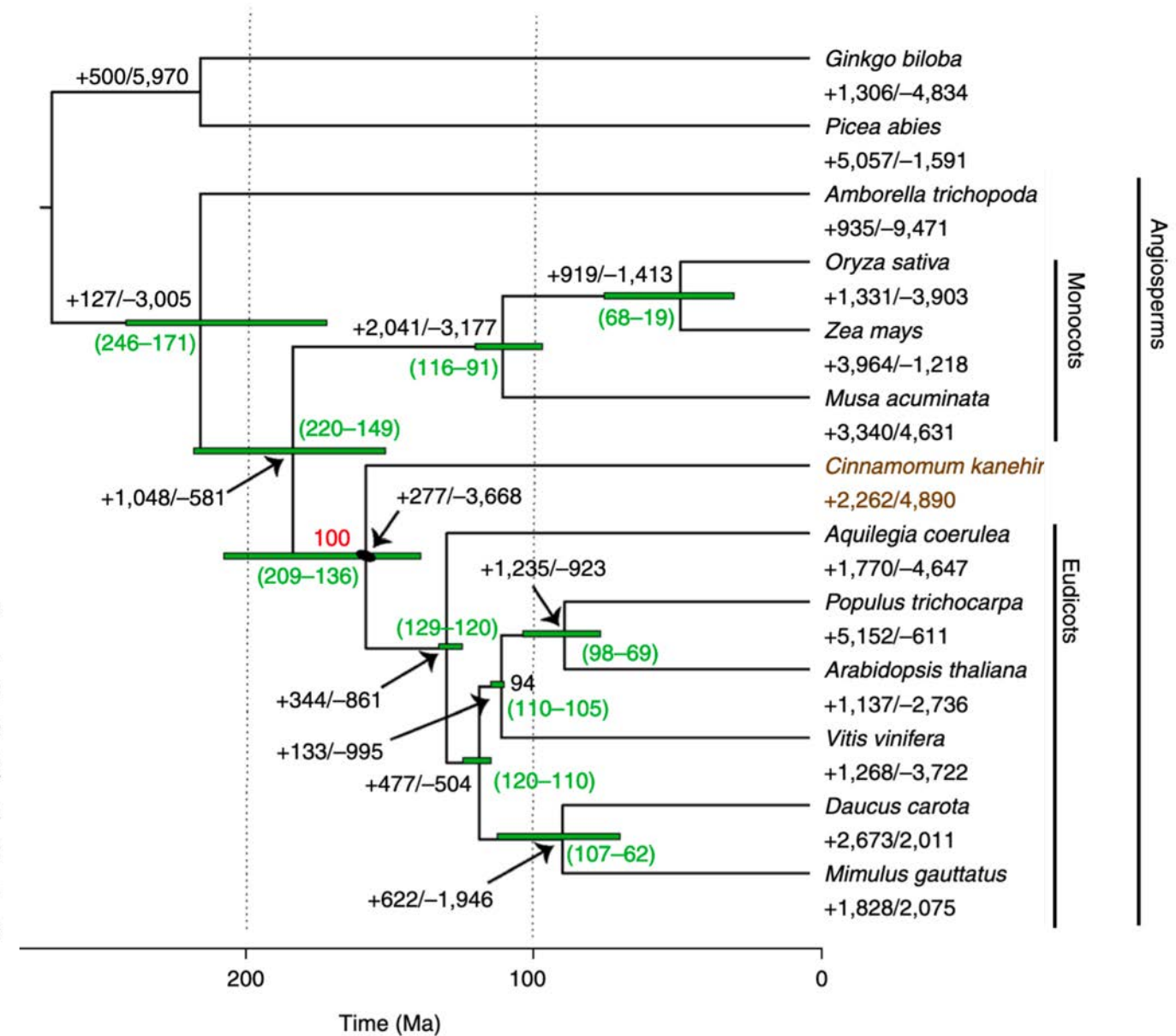
## Stout camphor tree genome fills gaps in understanding of flowering plant genome evolution

Shu-Miaw Chaw <sup>1,6\*</sup>, Yu-Ching Liu<sup>1</sup>, Yu-Wei Wu<sup>2</sup>, Han-Yu Wang<sup>1</sup>, Chan-Yi Ivy Lin<sup>1</sup>, Chung-Shien Wu<sup>1</sup>, Huei-Mien Ke<sup>1</sup>, Lo-Yu Chang<sup>1,3</sup>, Chih-Yao Hsu<sup>1</sup>, Hui-Ting Yang<sup>1</sup>, Edi Sudianto <sup>1</sup>, Min-Hung Hsu<sup>1,4</sup>, Kun-Pin Wu<sup>4</sup>, Ling-Ni Wang<sup>1</sup>, James H. Leebens-Mack<sup>5</sup> and Isheng J. Tsai <sup>1,6\*</sup>

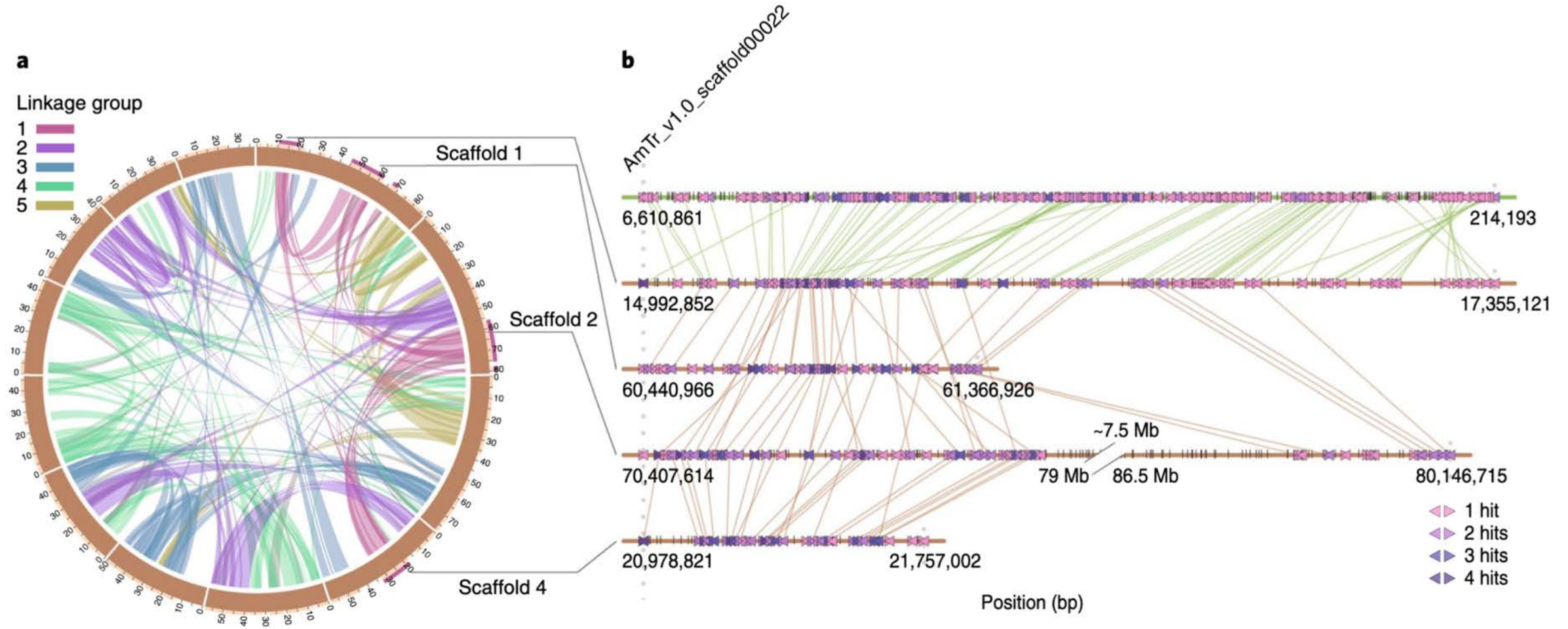
# Phylogenomic placement

**Phylogenomic placement of *C. kanehirae* sister to eudicots.** To resolve the long-standing debate over the phylogenetic placement of magnoliids relative to other major flowering plant lineages, we constructed a phylogenetic tree based on 211 strictly single-copy orthologue sets (that is, one and only one homologue in all species) identified through OrthoFinder<sup>21</sup> gene family circumscription of all gene models from the SCT and 12 other seed plant genomes (see Methods). A single species tree was recovered through maximum likelihood analysis<sup>27</sup> of a concatenated supermatrix of the single-copy gene alignments and coalescent-based analysis using the 211

gene trees<sup>28</sup> (Fig. 2 and Supplementary Fig. 11). SCT, representing the magnoliid lineage, was placed as sister to the eudicot clade (Fig. 2). This topology remained robust when we included a transcriptome data set of an additional 22 species of magnoliids order from the 1,000 plants initiative<sup>29</sup> (1KP), although lower bootstrap support was obtained (Supplementary Fig. 12). Using MCMCtree<sup>30</sup> with fossil calibrations, we calculated a 95% confidence interval for the time of divergence between magnoliids and eudicots to be 136.0–209.4 Ma (Fig. 2), which overlaps with two other recent estimates (114.8–164.1 Ma<sup>31</sup> and 118.9–149.9 Ma<sup>32</sup>).

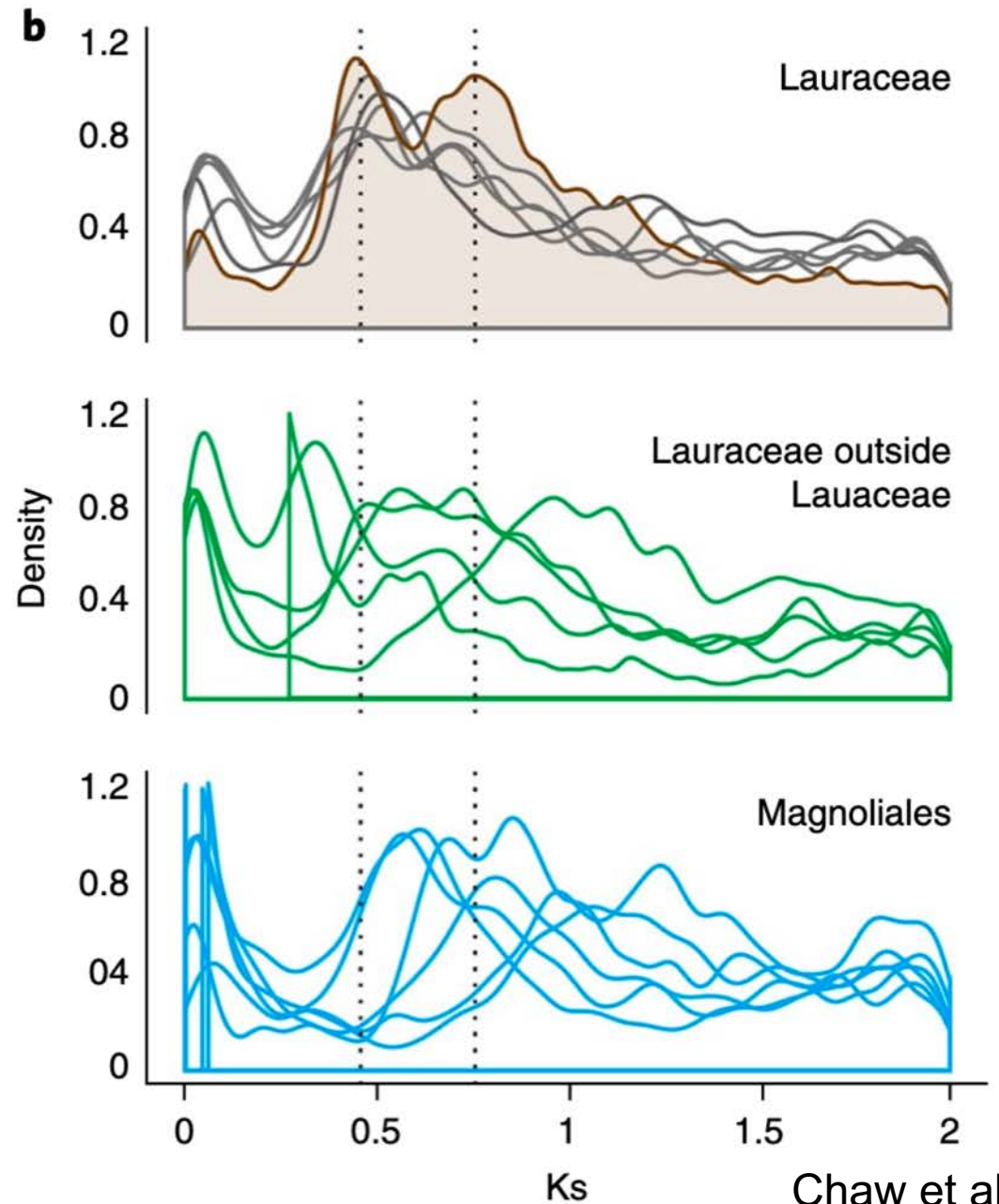
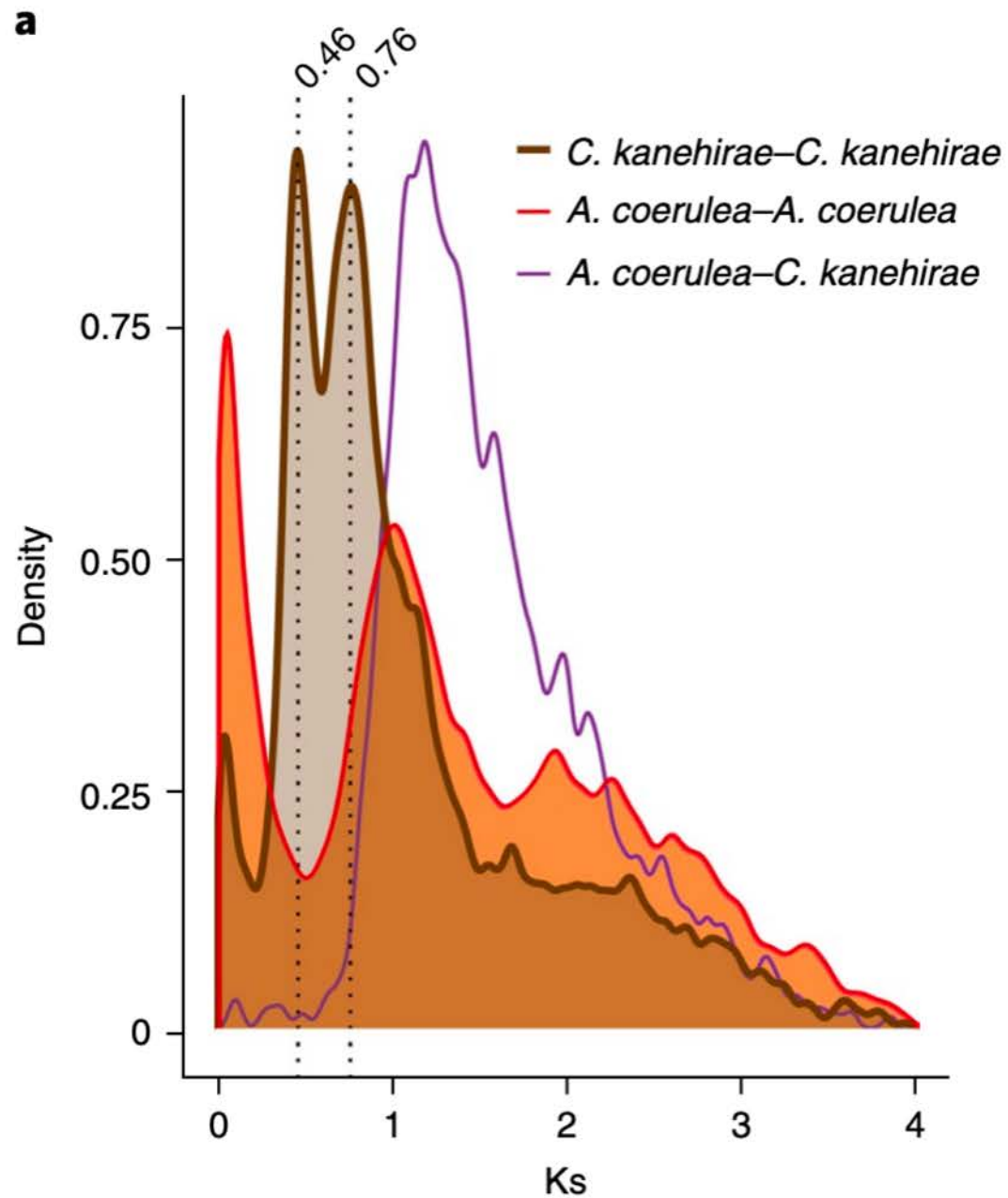


# Signature of whole genome duplication



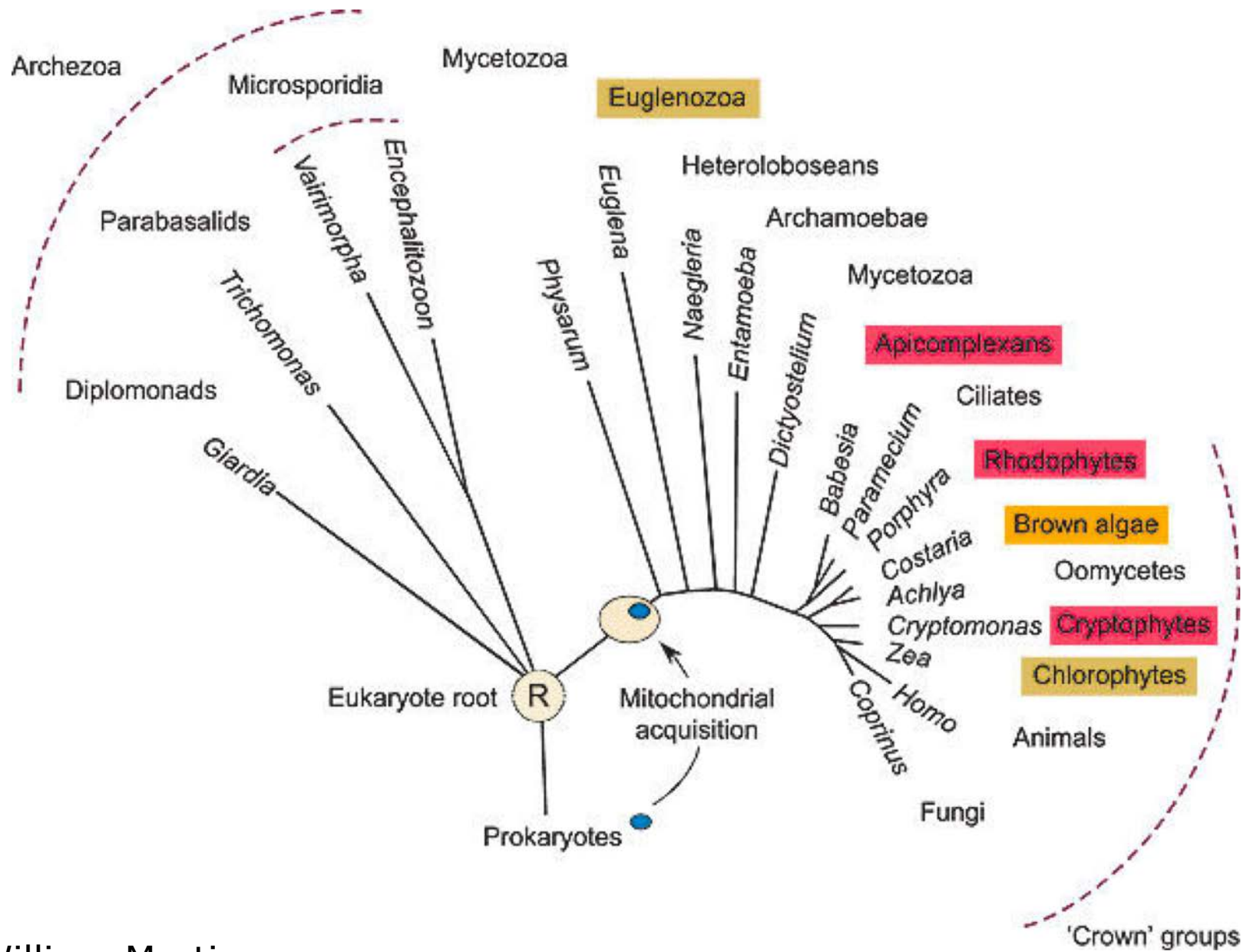
**Fig. 3 | Evolutionary analysis of the SCT genome.** **a**, Schematic representation of the intragenomic relationship among the 637 syntenic blocks in the SCT genome. Syntenic blocks (denoted by peach blocks) were assigned unambiguously into five linkage clusters representing ancient karyotypes and are colour coded. Purple blocks denote the syntenic block assigned in the first linkage group (see also Supplementary Fig. 13). **b**, Schematic representation of the first linkage group within the SCT genome and their corresponding relationship in *A. trichopoda*.

# Signature of whole genome duplication



How do we study origin of organelles?

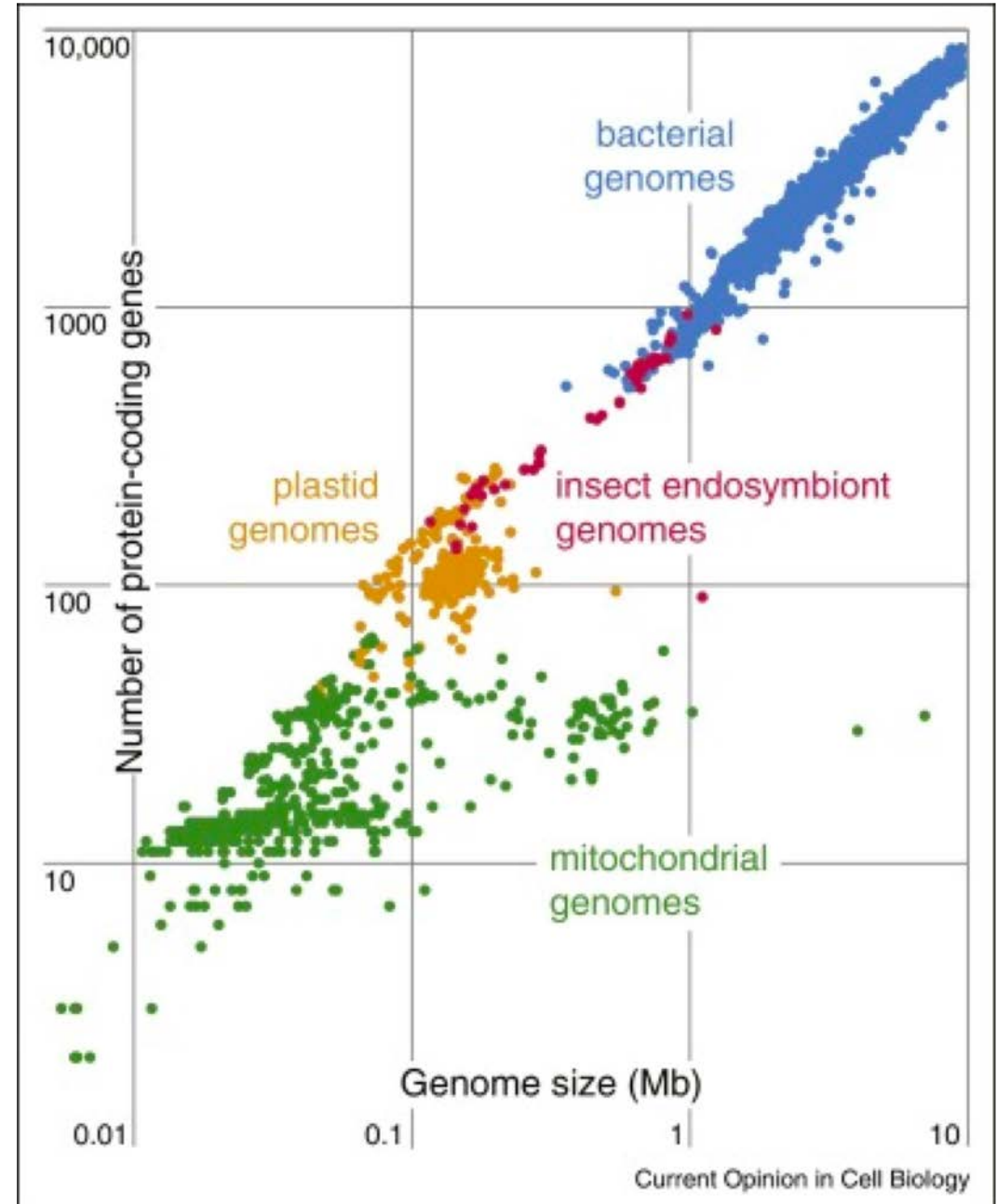




Martin Embley & William Martin  
*Nature* **440**, 623-630(30 March 2006)

Genomes from bacteria, insect endosymbionts, chloroplasts, and mitochondria form an unbroken continuum of size and coding density. The plot is truncated at 10 Mb and 10,000 genes.

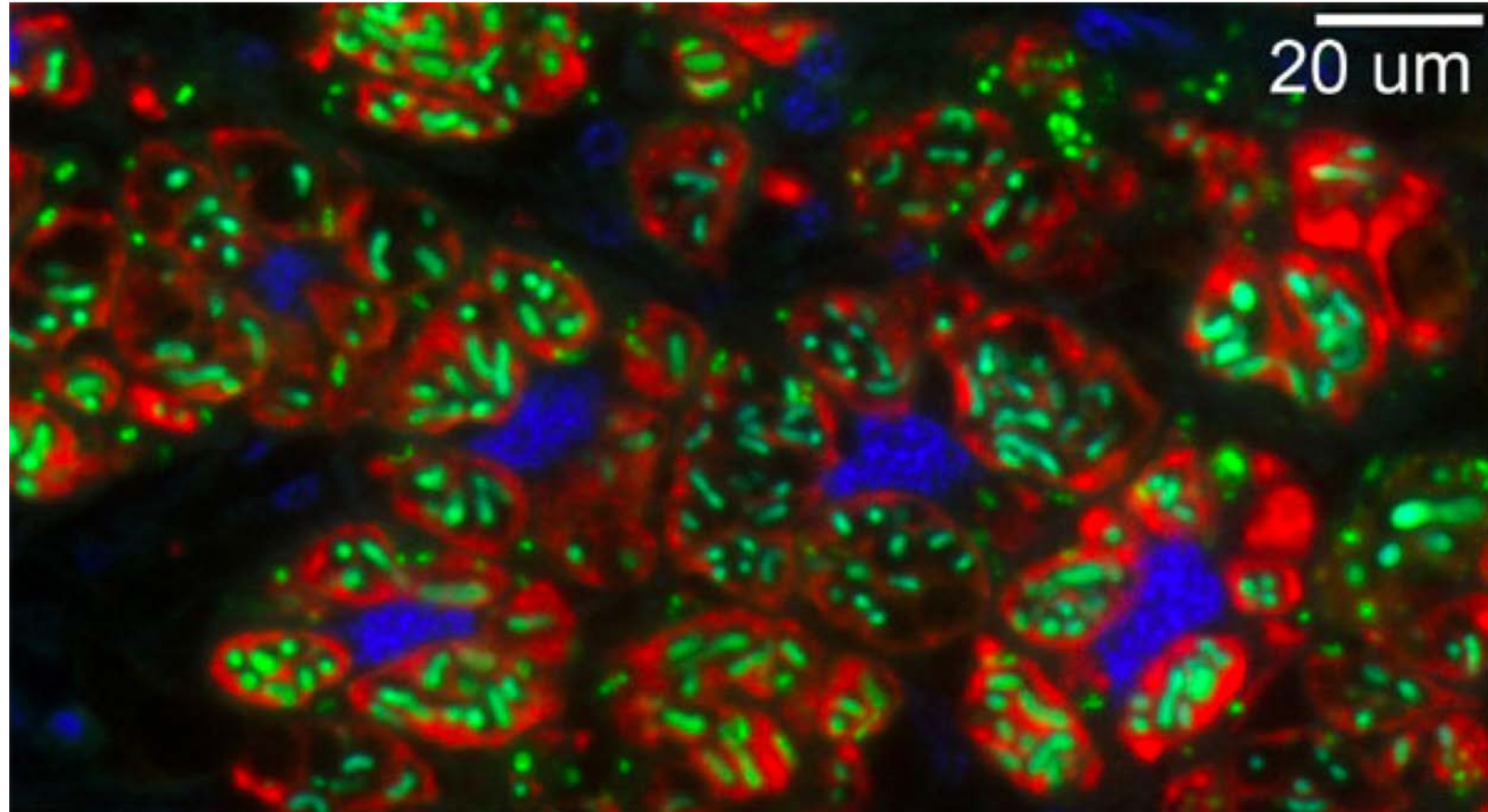
“Insect endosymbionts are missing (genomic) links between bacteria and organelles. It is now widely appreciated that all animals form symbioses with bacteria. Insects are especially interesting in this regard because they form many intracellular symbioses — that is, they allow bacteria to live inside their cells — that are not pathogenic from the host perspective”



# Case study: Mealybugs

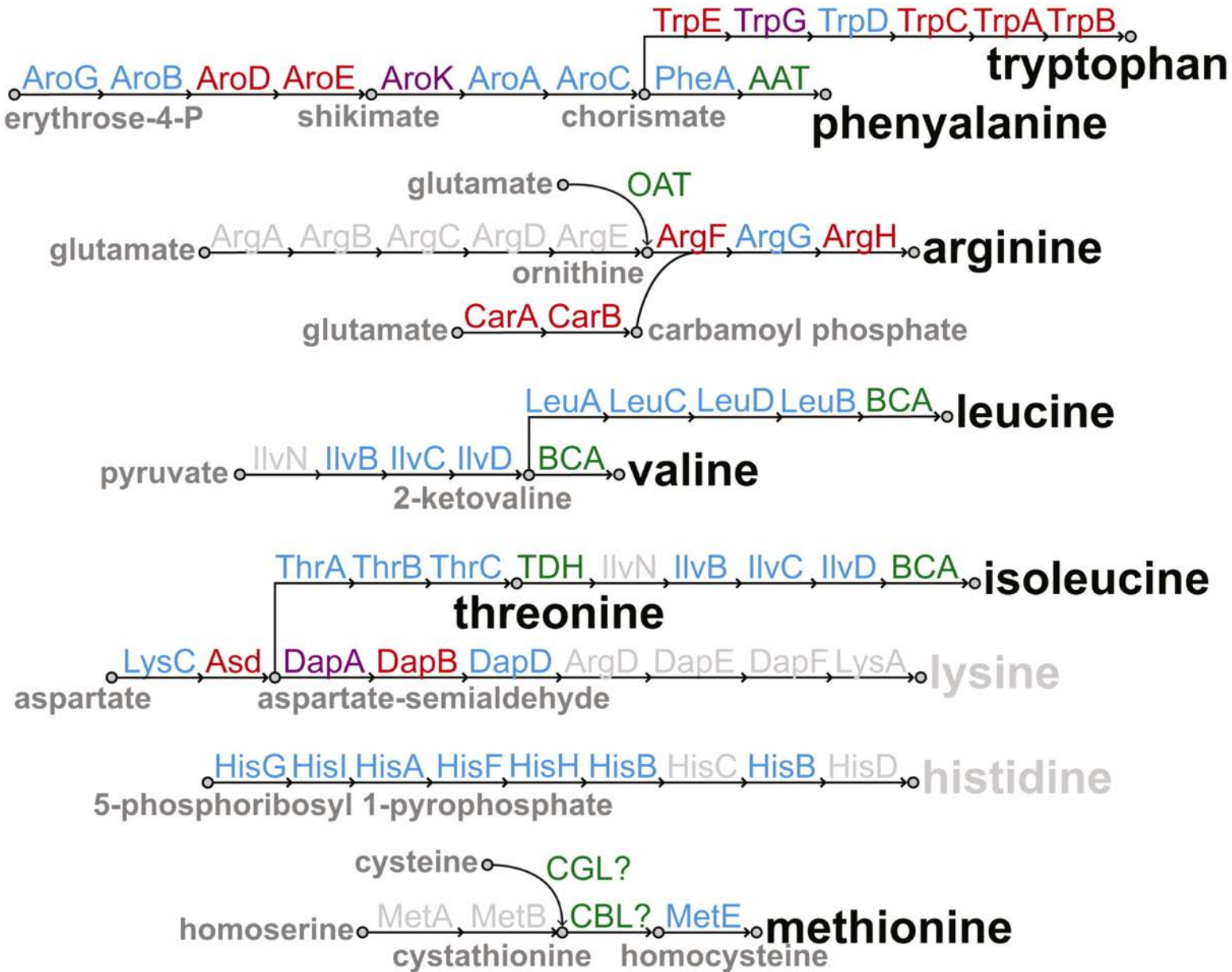


# Triple Symbiotic Relationship between Mealybugs, *Tremblaya princeps*, and *Moranella endobia*



Mealybug cells, showing *Tremblaya* (red), *Moranella* (green) and mealybug nuclei (blue).  
Credit: Ryuichi Koga, National Institute of Advanced Industrial Science and Technology, Japan

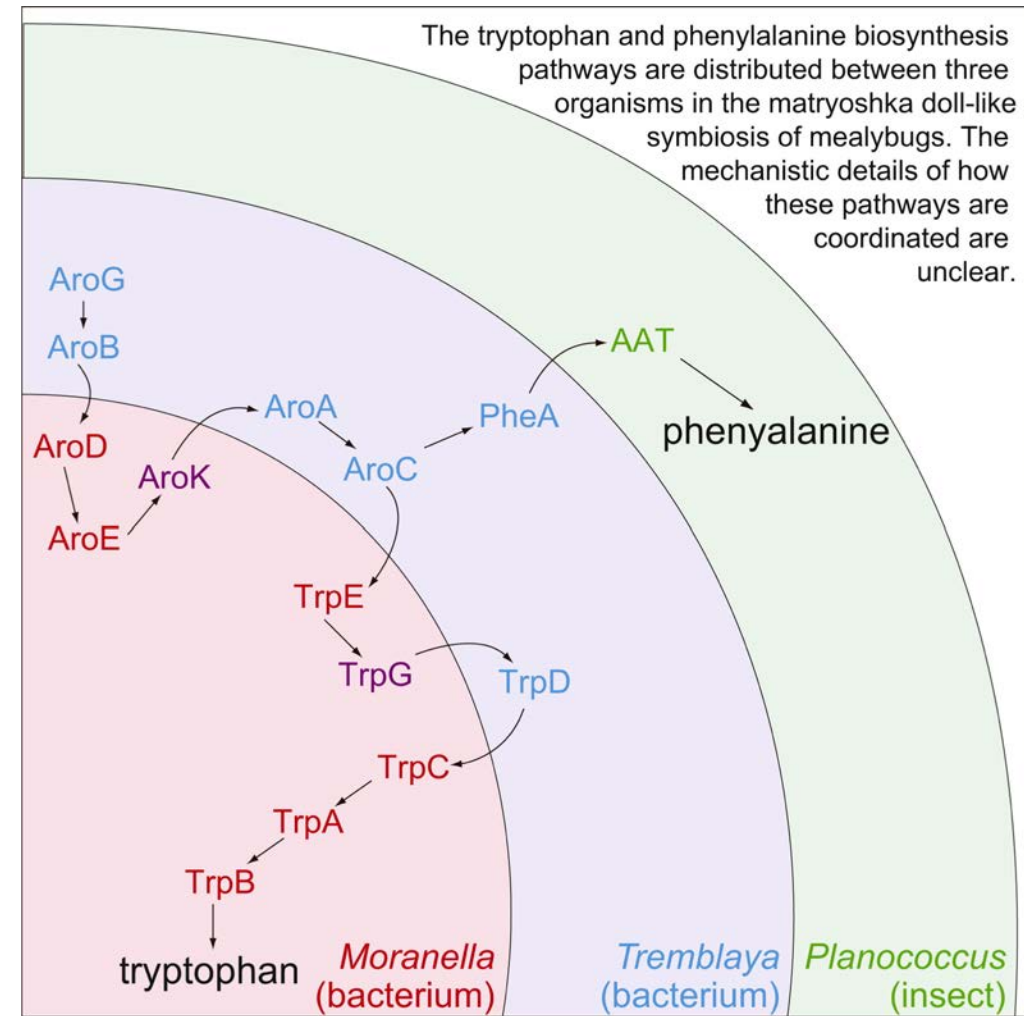
Predicted Essential Amino Acid Metabolic Contributions of the Mealybug-Tremblaya-Moranella Symbiosis



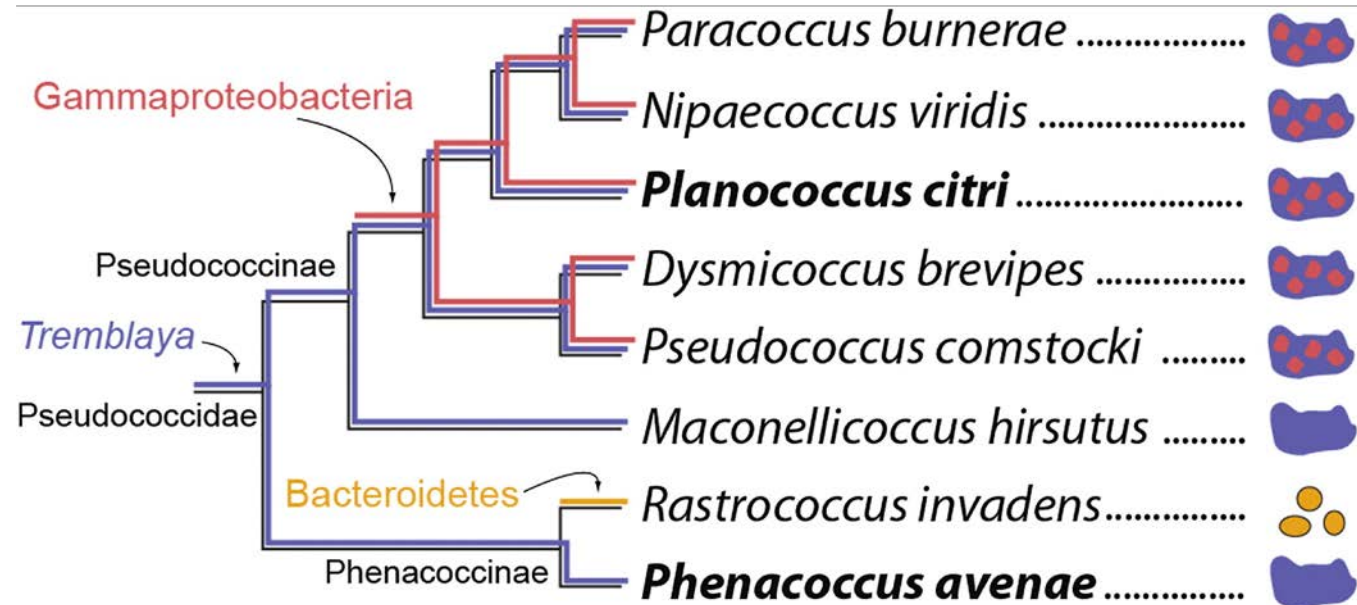
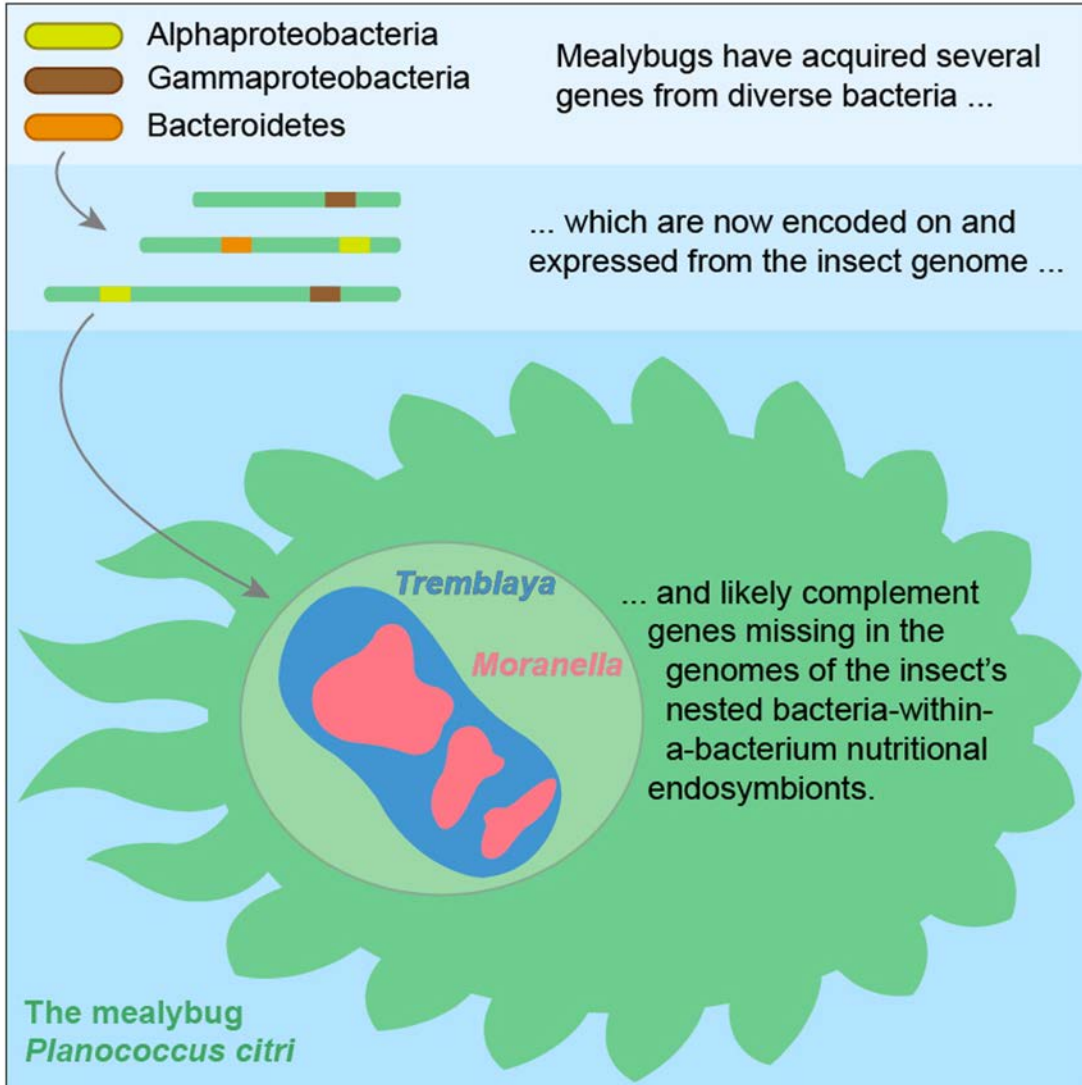
Gene homologs found in the Tremblaya genome are blue; the Moranella genome, red; both the Tremblaya and Moranella genomes, purple; neither the Tremblaya nor the Moranella genome, gray; activities not found in either bacterial genome but predicted to be encoded in the mealybug genome, green.

Genome degeneracy of a bacterial endosymbiont is driven by its own endosymbiont

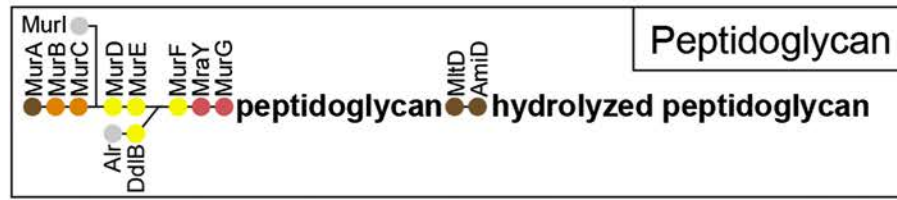
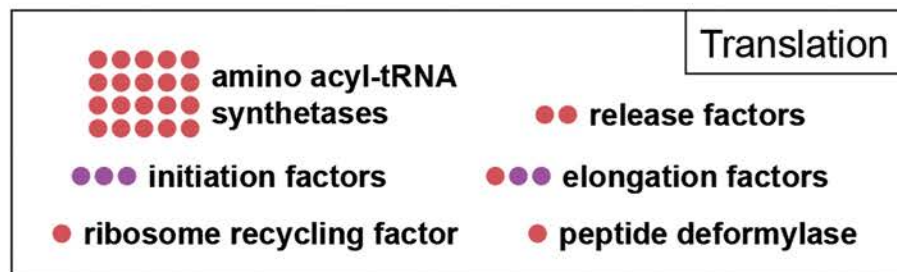
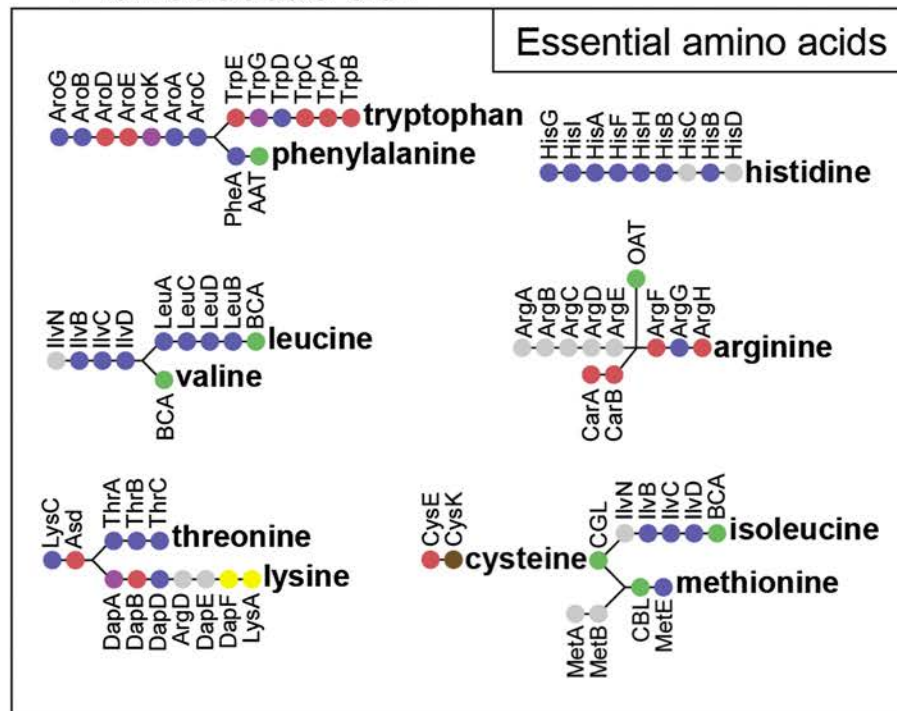
- HGT from diverse bacteria to the insect host genome support the three-way symbiosis
- Endosymbiont genomes can massively degrade without transfer of genes to the host



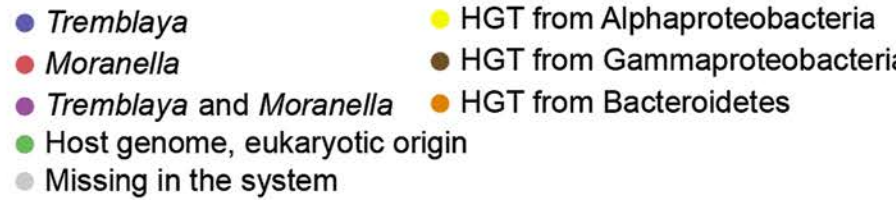
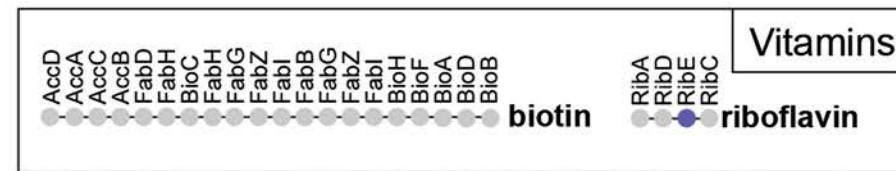
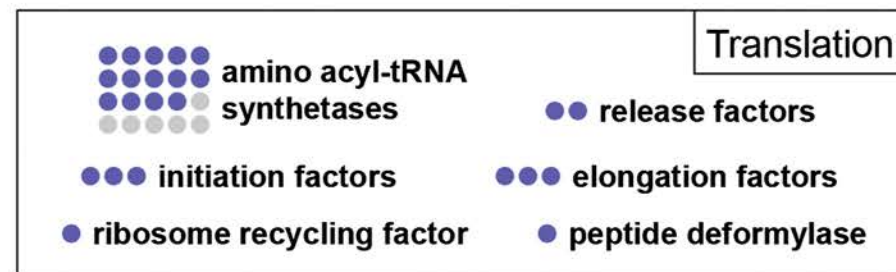
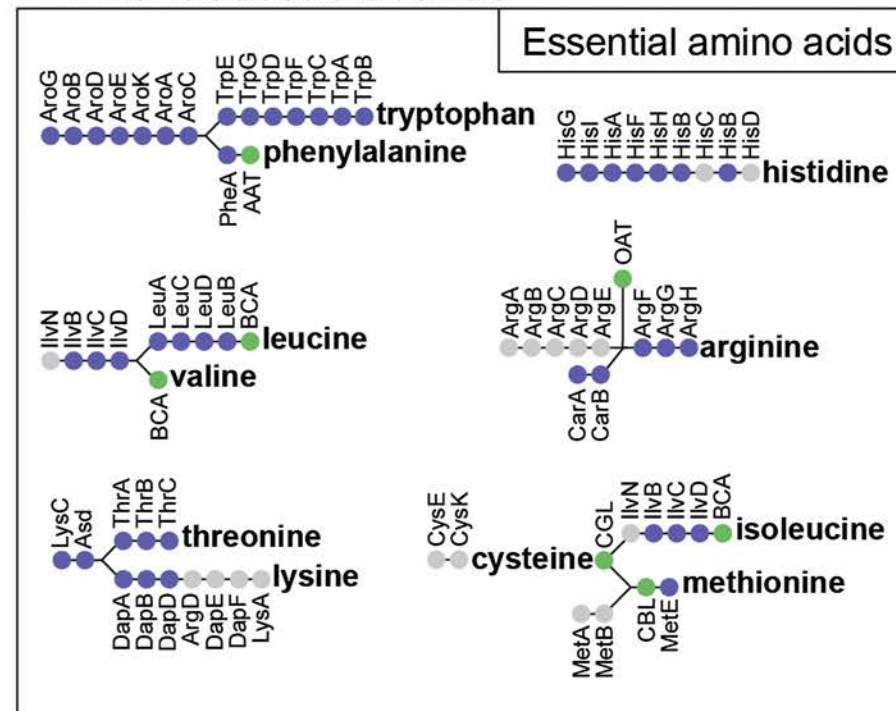
# Horizontal Gene Transfer from Diverse Bacteria to an Insect Genome Enables a Tripartite Nested Mealybug Symbiosis



### A *Planococcus citri*



### B *Phenacoccus avenae*





# Even more fascinating case

Cell

## Sympatric Speciation in a Bacterial Endosymbiont Results in Two Genomes with the Functionality of One

James T. Van Leuven,<sup>1</sup> Russell C. Meister,<sup>2</sup> Chris Simon,<sup>2</sup> and John P. McCutcheon<sup>1,3,\*</sup>

<sup>1</sup>Division of Biological Sciences, University of Montana, Missoula, MT 59812, USA

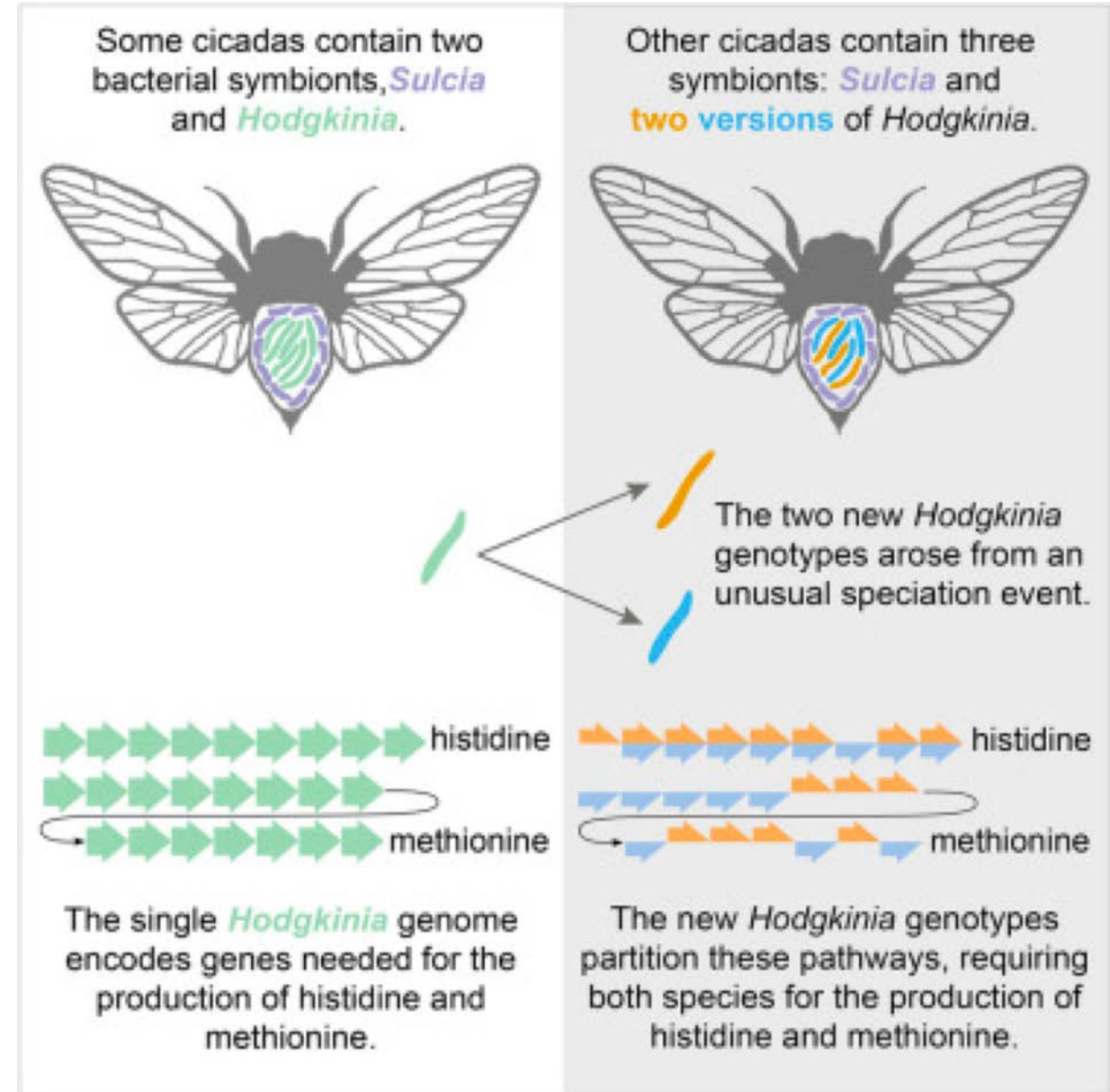
<sup>2</sup>Department of Ecology and Evolutionary Biology, University of Connecticut, Storrs, CT 06269, USA

<sup>3</sup>Canadian Institute for Advanced Research, CIFAR Program in Integrated Microbial Biodiversity, Toronto, ON M5G 1Z8, Canada

\*Correspondence: [john.mccutcheon@umontana.edu](mailto:john.mccutcheon@umontana.edu)

<http://dx.doi.org/10.1016/j.cell.2014.07.047>

<https://www.youtube.com/watch?v=XRI2JxTzJ-0&list=UUISV2Tk7x-wBBXP6-VCNbNw>



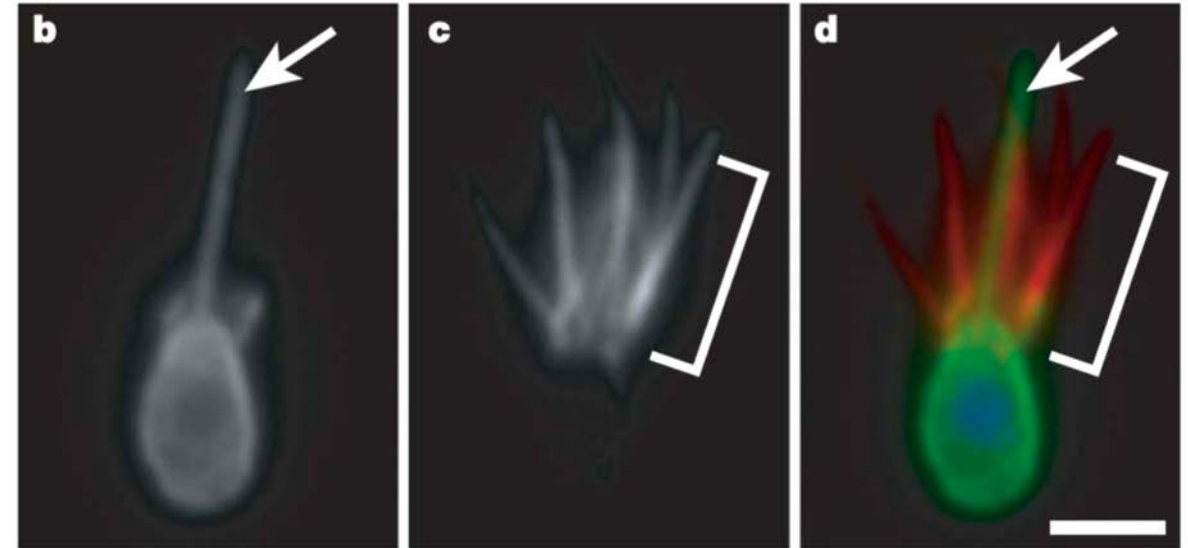
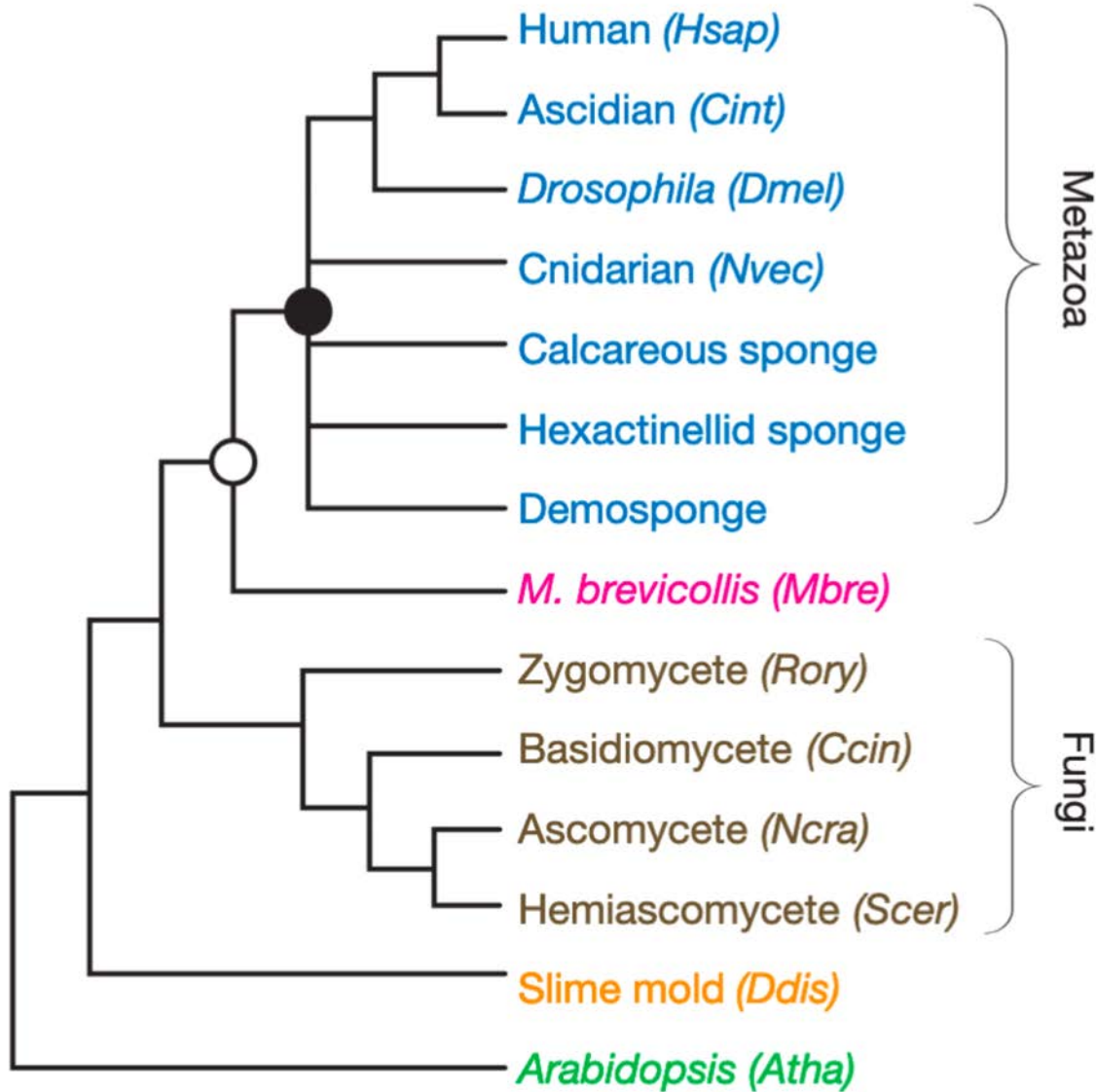
How do we study origin of animals (metazoans)?

# Question: what constitute the first animal?

- Unicellular -> Multicellular
  - What's needed for multicellularity?
    - Interactions between cells
    - Formation of aggregates? How?

# Choanoflagellate

**a**



*Ddis*) and *Arabidopsis* (*Arabidopsis thaliana*, *Atha*). **b–d**, Choanoflagellate cells bear a single apical flagellum (arrow, **b**) and an apical collar of actin-filled microvilli (bracket, **c**). **d**, An overlay of  $\beta$ -tubulin (green), polymerized actin (red) and DNA localization (blue) reveals the position of the flagellum within the collar of microvilli. Scale bar, 2  $\mu$ m.

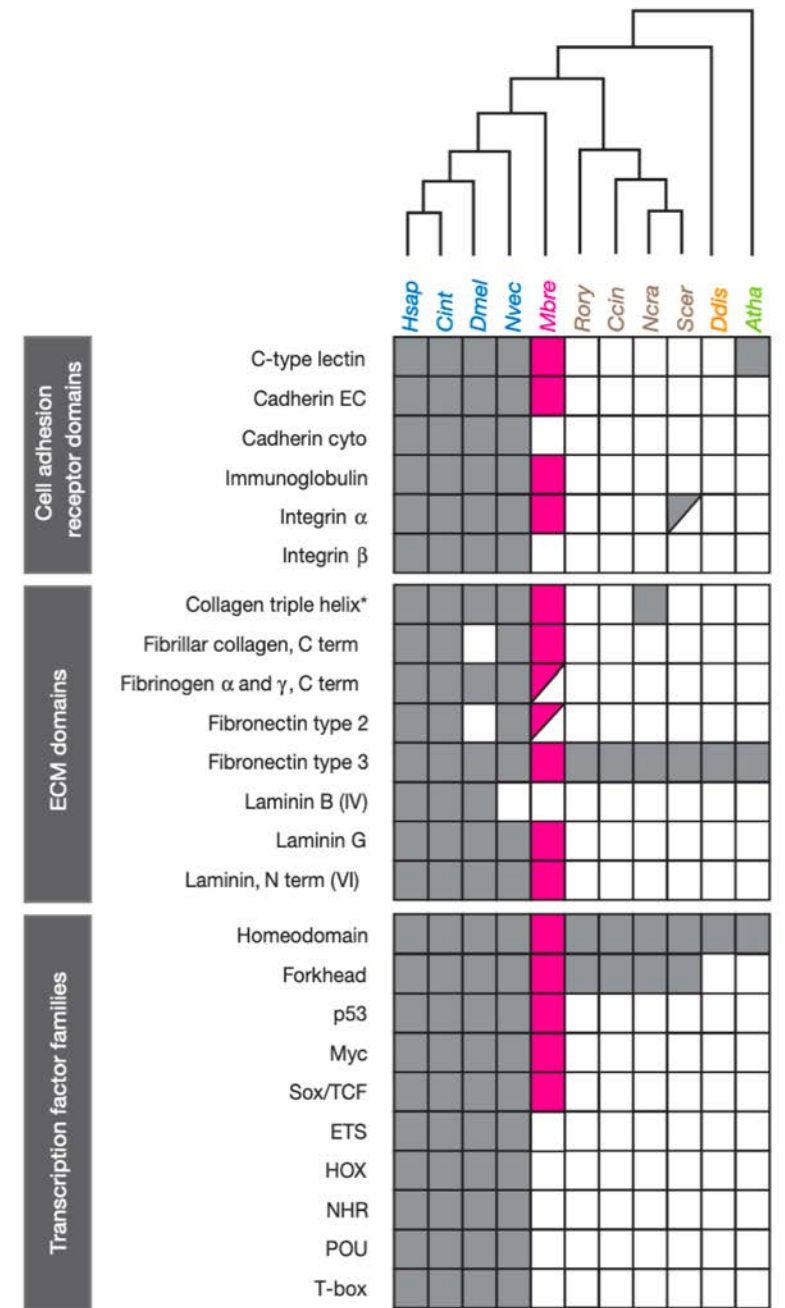
# Choanoflagellate – example findings

## An abundance of cell adhesion domains

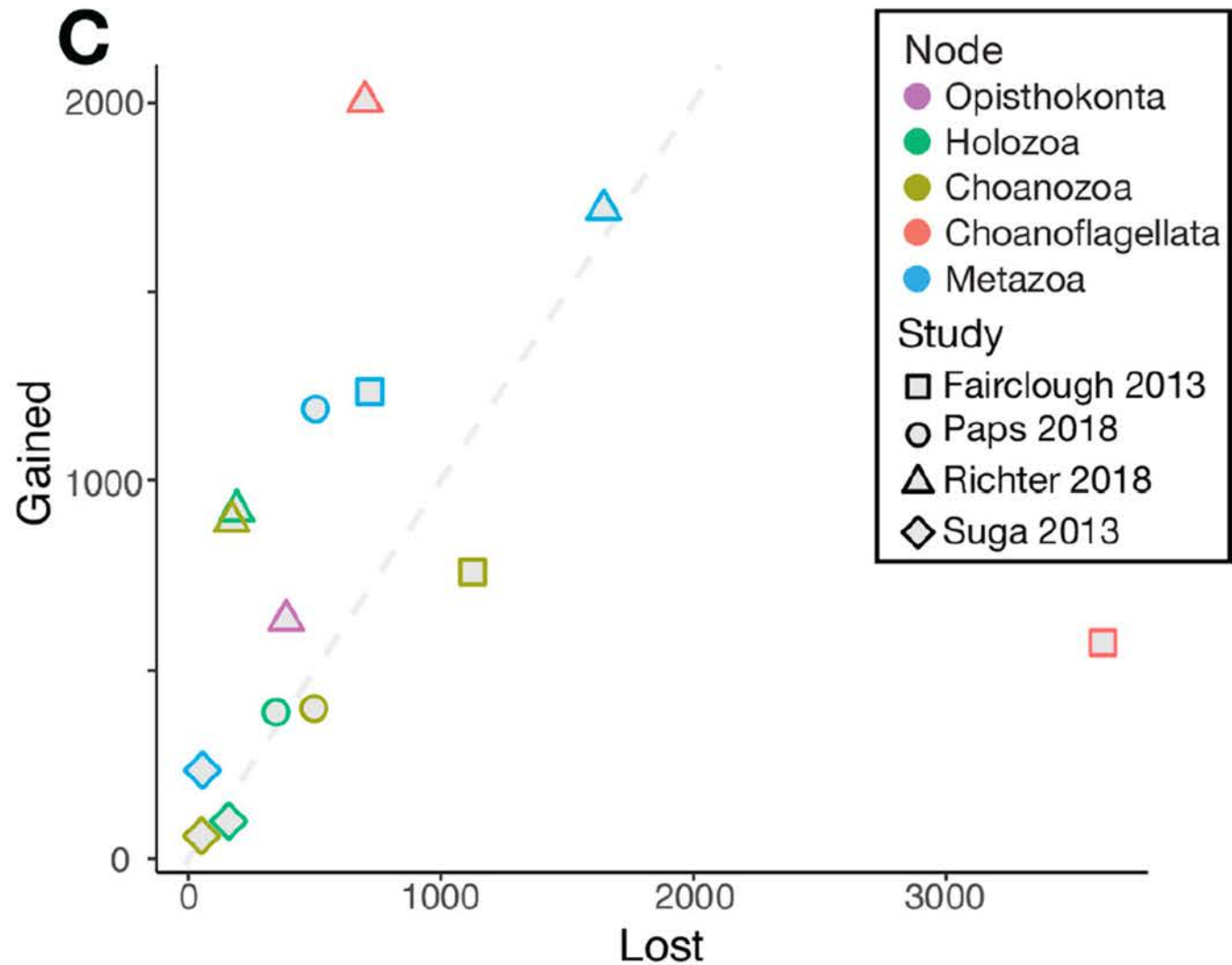
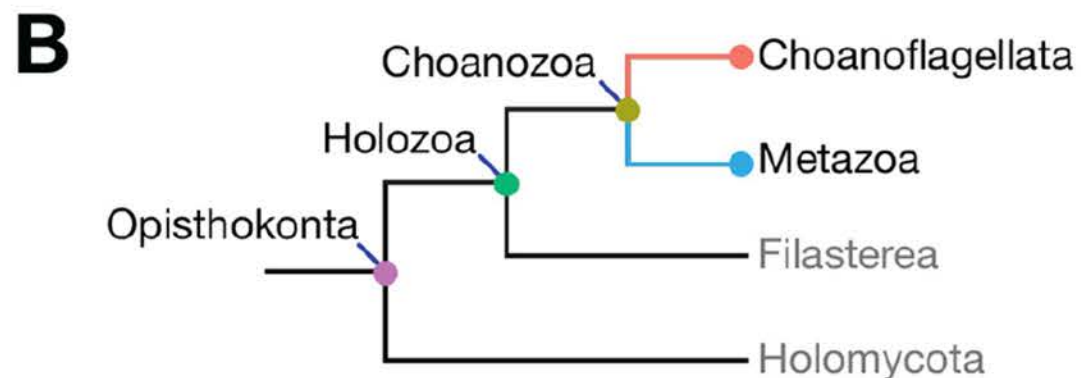
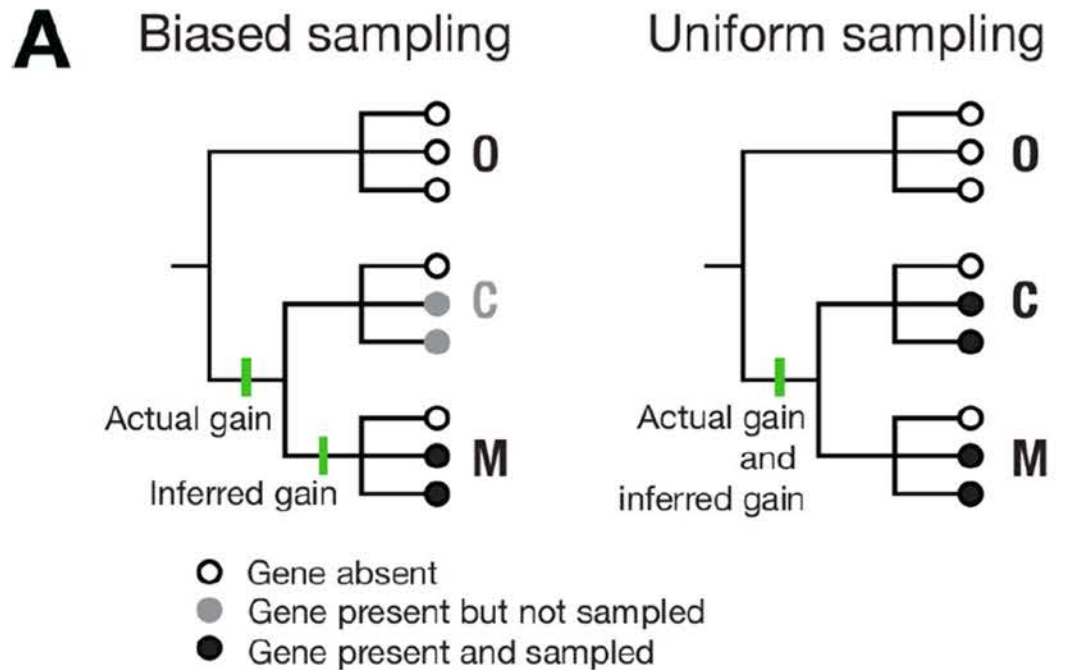
A critical step in the transition to multicellularity was the evolution of mechanisms for stable cell adhesion. *M. brevicollis* encodes a diverse array of cell adhesion and extracellular matrix (ECM) protein domains previously thought to be restricted to metazoans (Fig. 3).

The finding in *M. brevicollis* of cell adhesion domains that were previously known only in metazoans has two important implications. First, the common ancestor of metazoans and choanoflagellates possessed several of the critical structural components used for multicellularity in modern metazoans. Second, given the absence of evidence for stable cell adhesion in *M. brevicollis*, this also suggests that homologues of metazoan cell adhesion domains may act to mediate interactions between *M. brevicollis* and its extracellular environment.

The discovery of putatively secreted ECM proteins in a free-living choanoflagellate suggests that elements of the metazoan ECM evolved in contact with the external environment before being sequestered within an epithelium. Although some choanoflagellates secrete extracellular structures or adhere to form colonial assemblages<sup>19,33,34</sup>, *M. brevicollis* is not known to do so. Instead, these ECM protein homologues in *M. brevicollis* may mediate an analogous process such as substrate attachment.

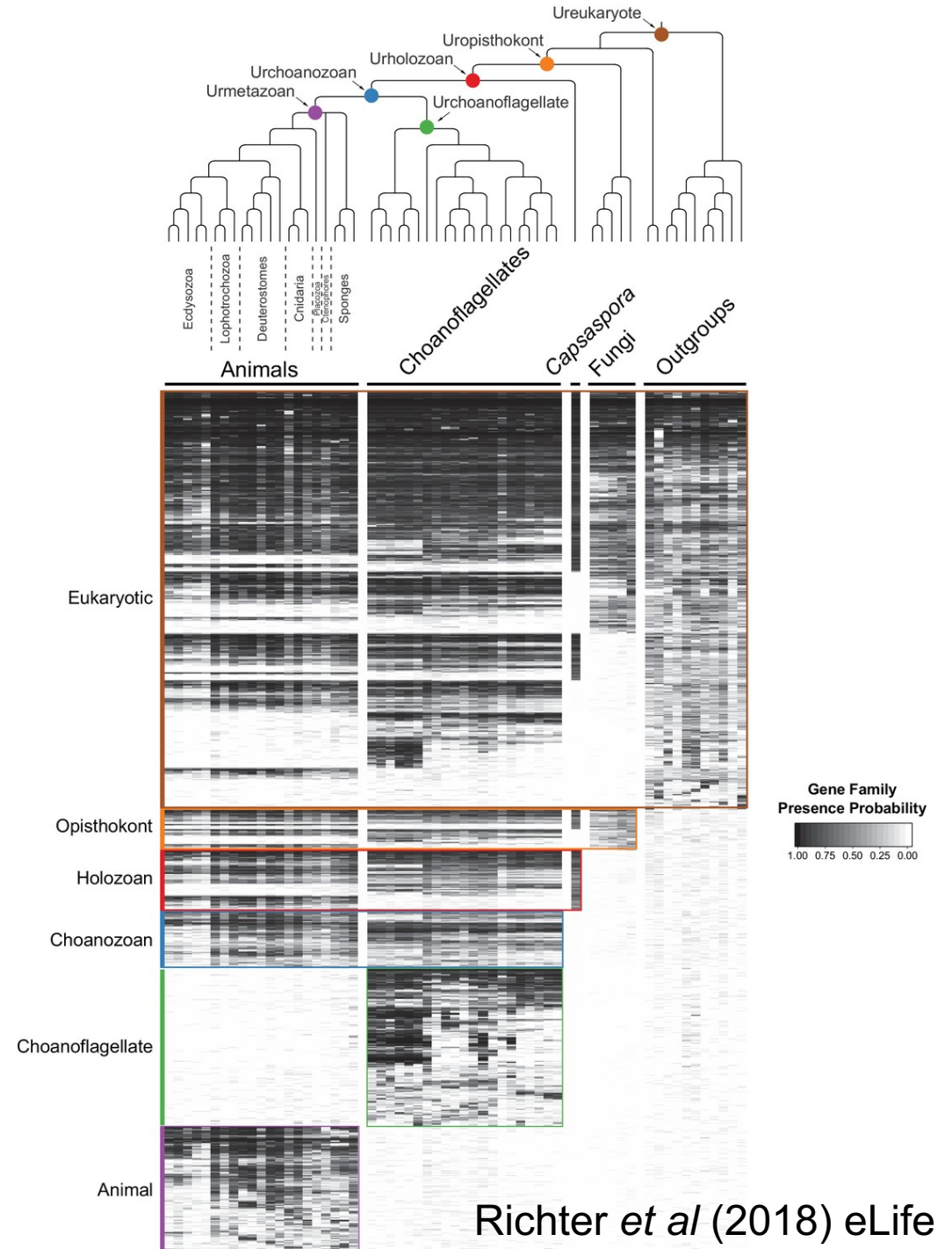


# Choanoflagellate – biased sampling can lead to different results

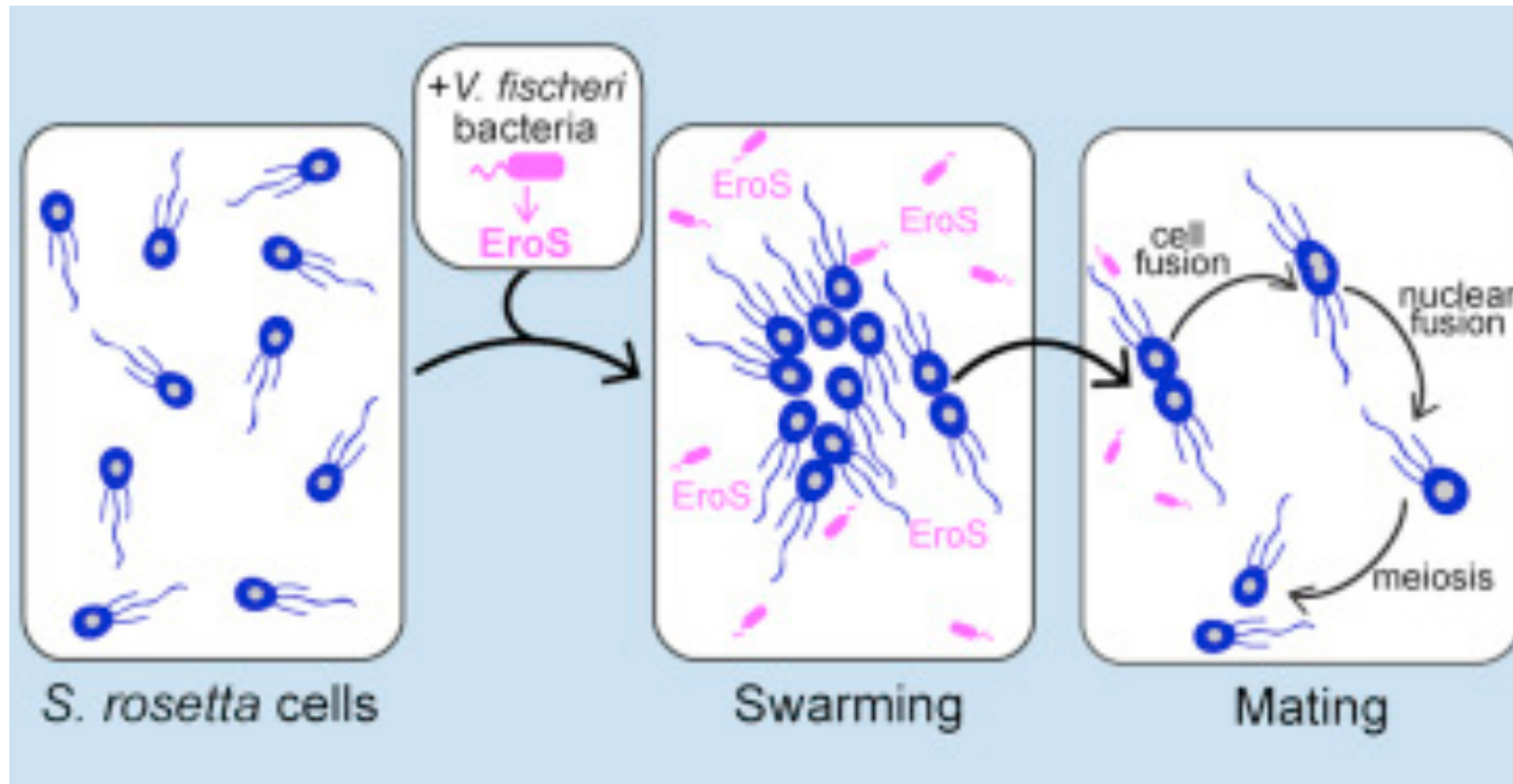


# Deeper sampling reveal Differential retention and loss of ancestral gene families in extant animals and choanoflagellates

“The patchwork ancestry of the Urmethazoan genome is illustrated by the fact that many gene families responsible for animal development, immunity and multicellular organization evolved through shuffling of protein domains that first originated in the choanozoan stem lineage together with ancient or animal-specific domains”



# Choanoflagellate – environmental cues to aggregate

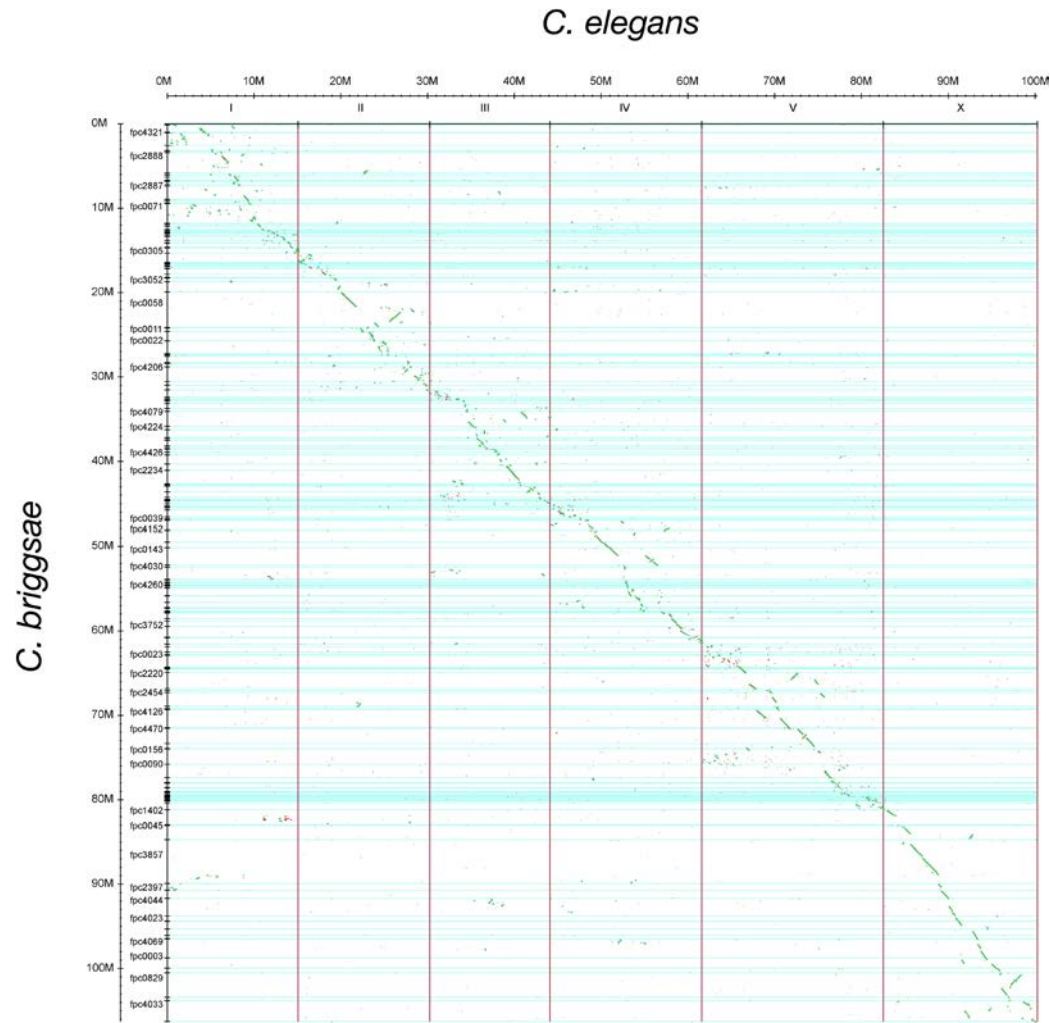


- The bacterium *Vibrio fischeri* induces mating in the choanoflagellate *S. rosetta*
- The “aphrodisiac” produced by *V. fischeri* is a chondroitinase that we name EroS
  - The enzymatic activity of EroS is required for this function
  - Chondroitin sulfate, the EroS substrate, evolved before the origin of animals

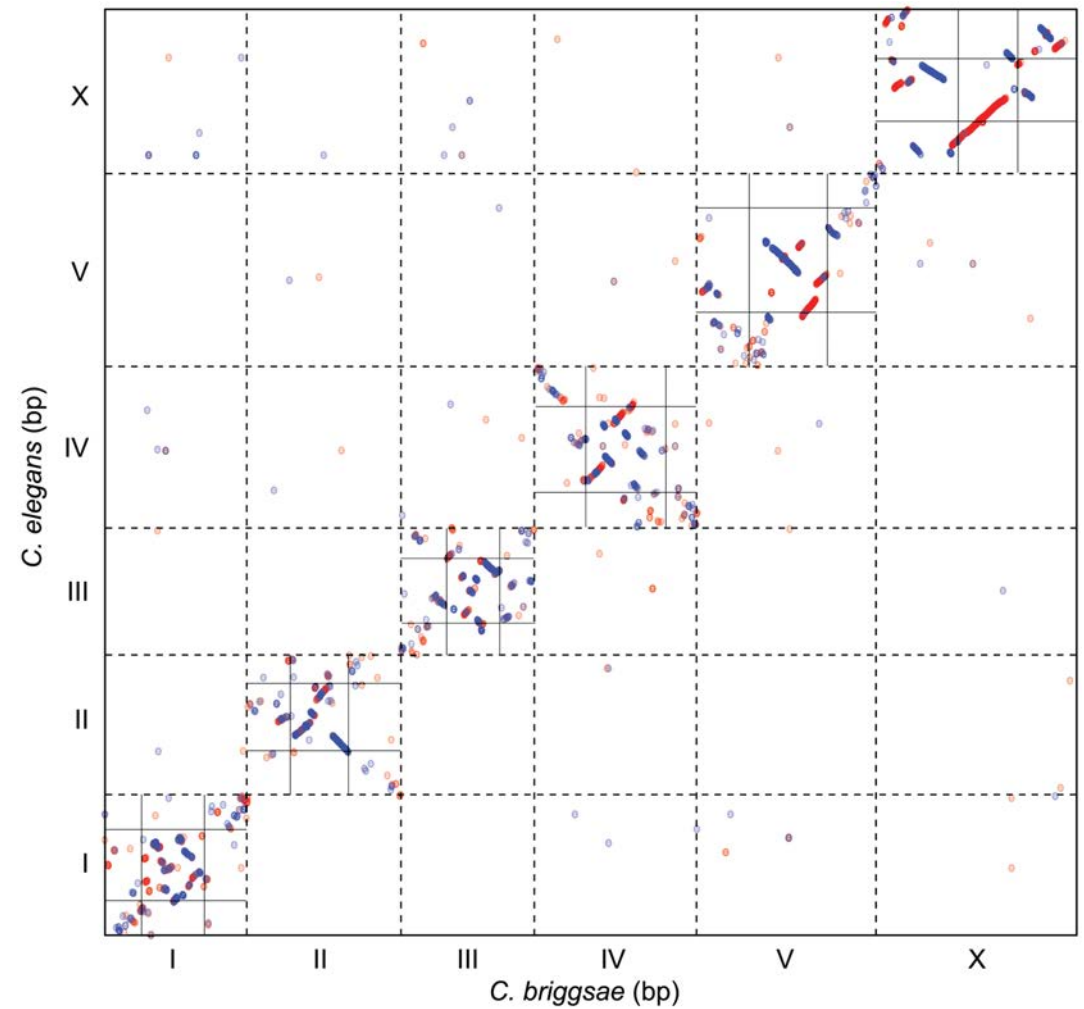


Some caveats

# Assembly quality likely to influence synteny observation



Stein *et al.*, PLOS Genetics (2003)



Ross *et al.*, PLOS Genetics (2011)

# Syteny based scaffolding: use with caution

Tang et al. *Genome Biology* (2015) 16:3  
DOI 10.1186/s13059-014-0573-1



METHOD

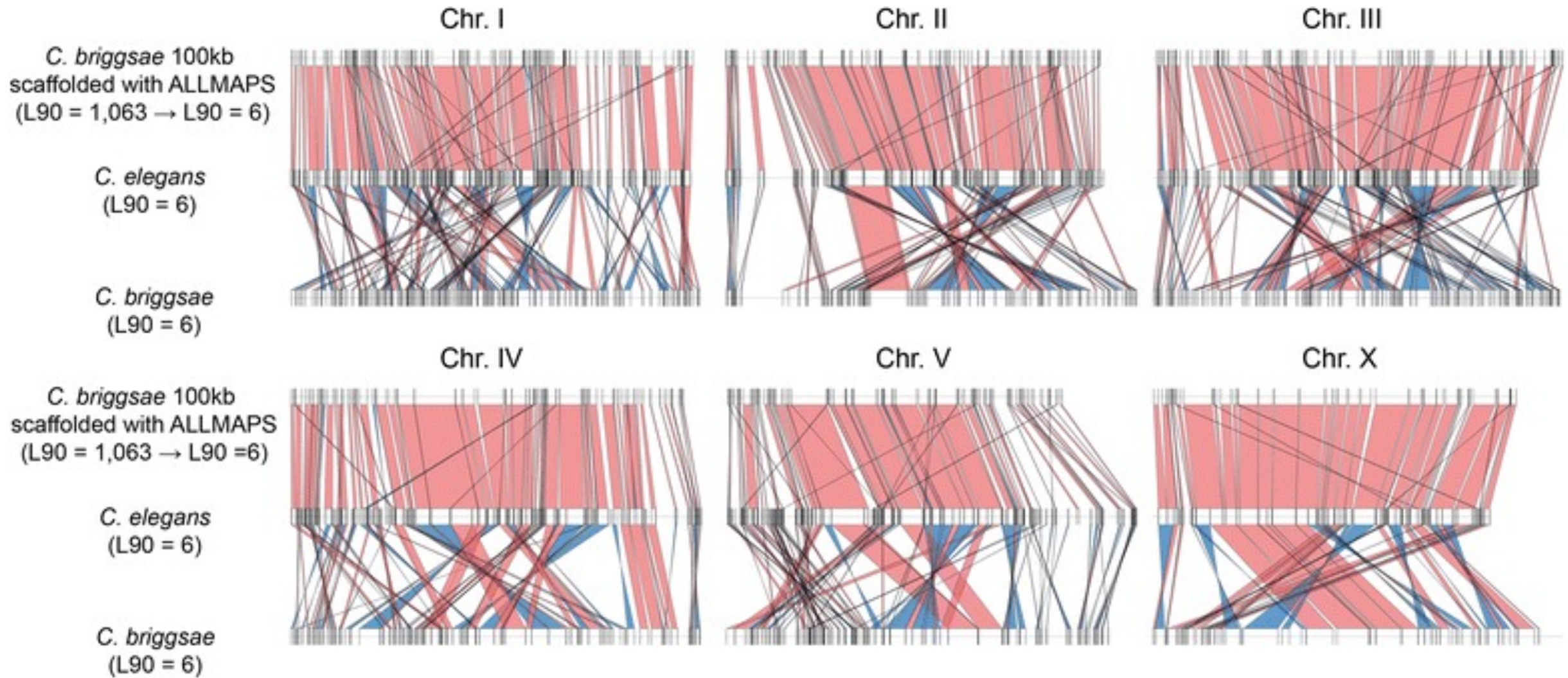
Open Access

## ALLMAPS: robust scaffold ordering based on multiple maps

Haibao Tang<sup>1,2,3\*</sup>, Xingtian Zhang<sup>4</sup>, Chenyong Miao<sup>1</sup>, Jisen Zhang<sup>1</sup>, Ray Ming<sup>1</sup>, James C Schnable<sup>3,5</sup>, Patrick S Schnable<sup>3,6</sup>, Eric Lyons<sup>2</sup> and Jianguo Lu<sup>7</sup>

for example, in ‘orphan’ species where there is little research investment in the past, **we can still create consensus chromosomal assemblies based on comparative maps against multiple, closely-related genomes as a collection of ‘references’ ... Correct?**

# Syteny based scaffolding: use with caution



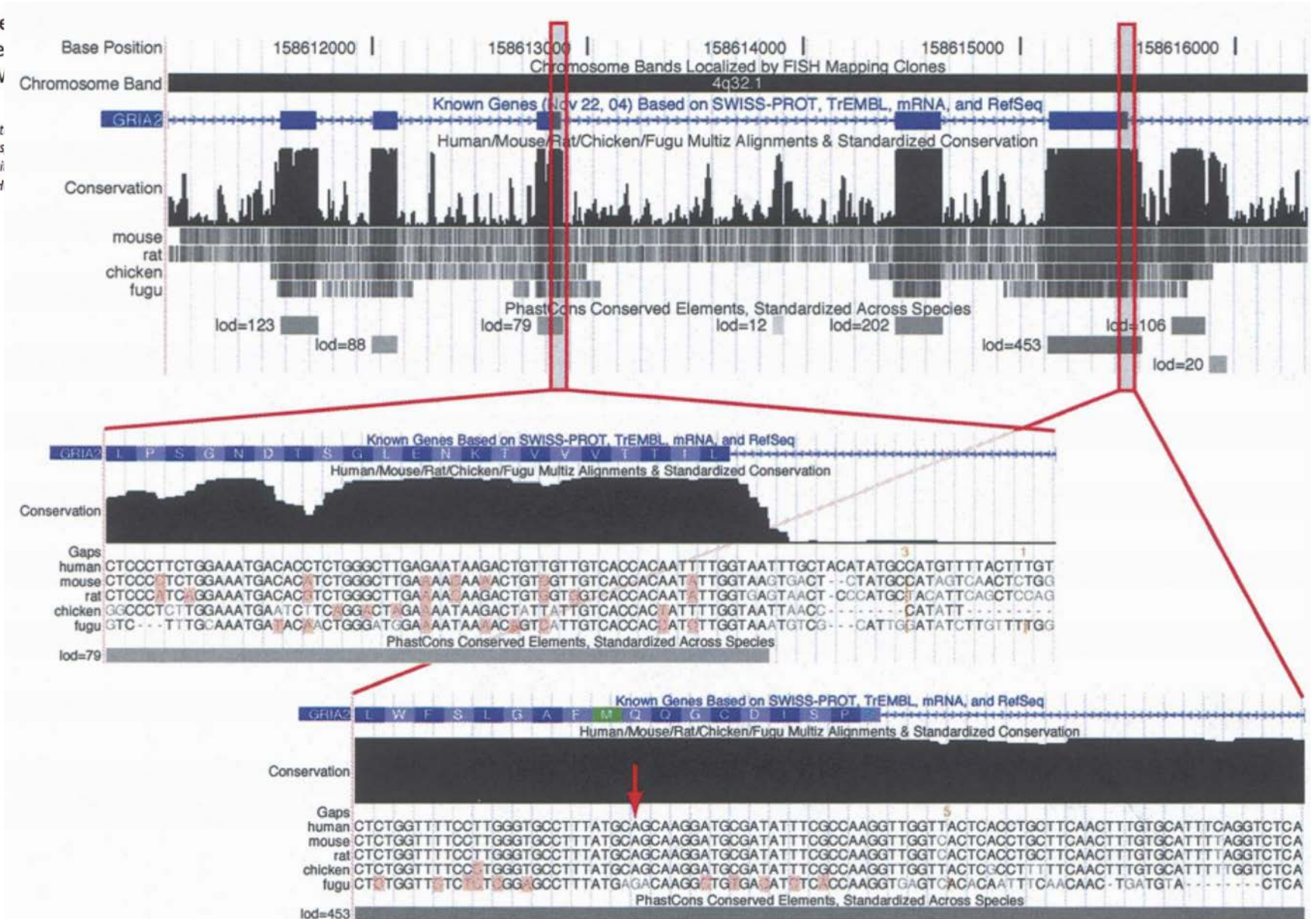
Homology beyond level of genes

# Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes

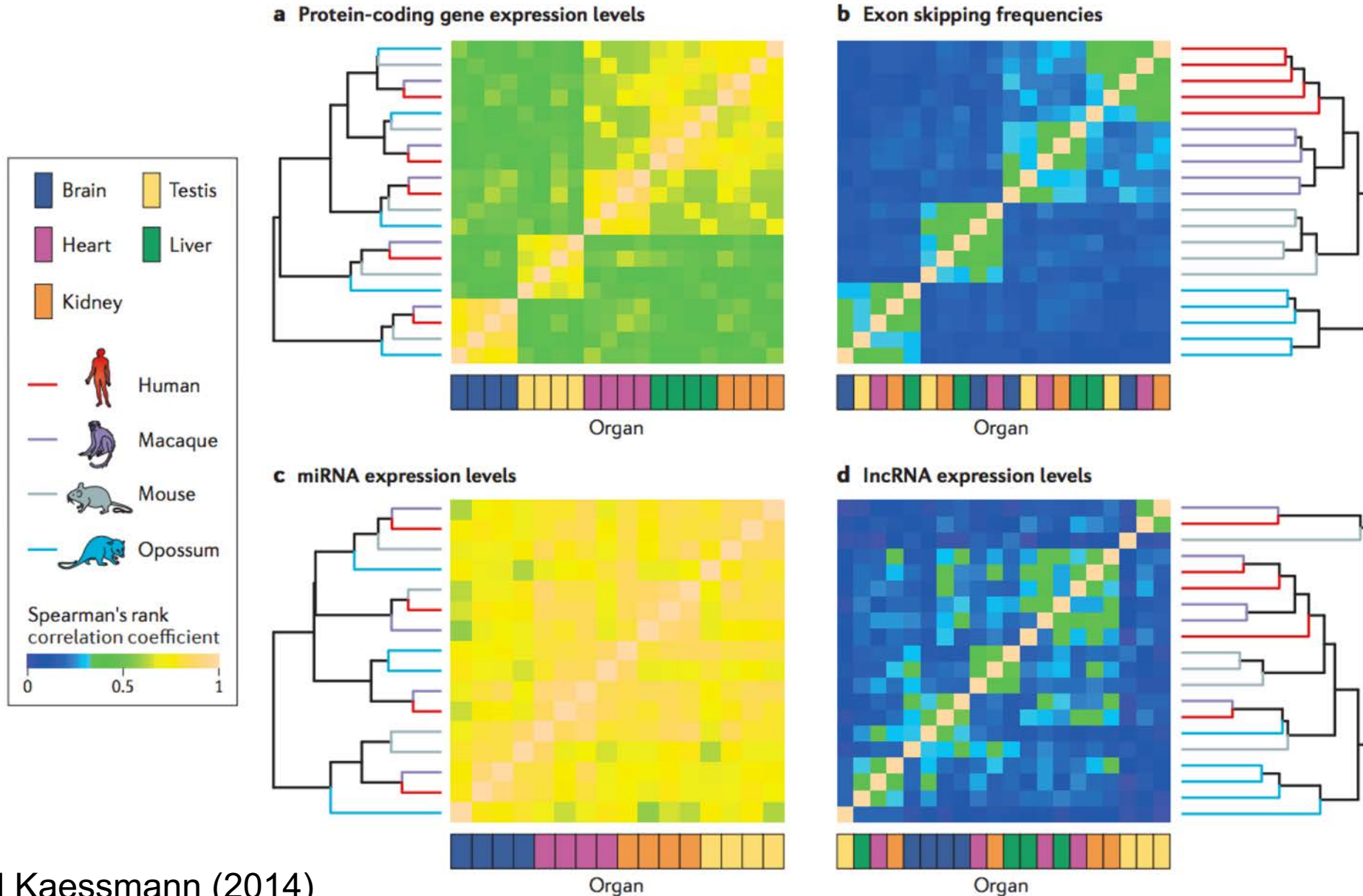
Adam Siepel,<sup>1,6</sup> Gill Bejerano,<sup>1</sup> Jakob S. Pedersen,<sup>1</sup> Angie Kate Rosenbloom,<sup>1</sup> Hiram Clawson,<sup>1</sup> John Spieth,<sup>4</sup> LaDe Stephen Richards,<sup>5</sup> George M. Weinstock,<sup>5</sup> Richard K. W. W. James Kent,<sup>1</sup> Webb Miller,<sup>3</sup> and David Haussler<sup>1,2</sup>

<sup>1</sup>Center for Biomolecular Science and Engineering, <sup>2</sup>Howard Hughes Medical Institut Cruz, California 95064, USA; <sup>3</sup>Center for Comparative Genomics and Bioinformatics Park, Pennsylvania 16802, USA; <sup>4</sup>Genome Sequencing Center, Washington Universi 63108, USA; <sup>5</sup>Human Genome Sequencing Center, Department of Molecular and H Houston, Texas 77030, USA

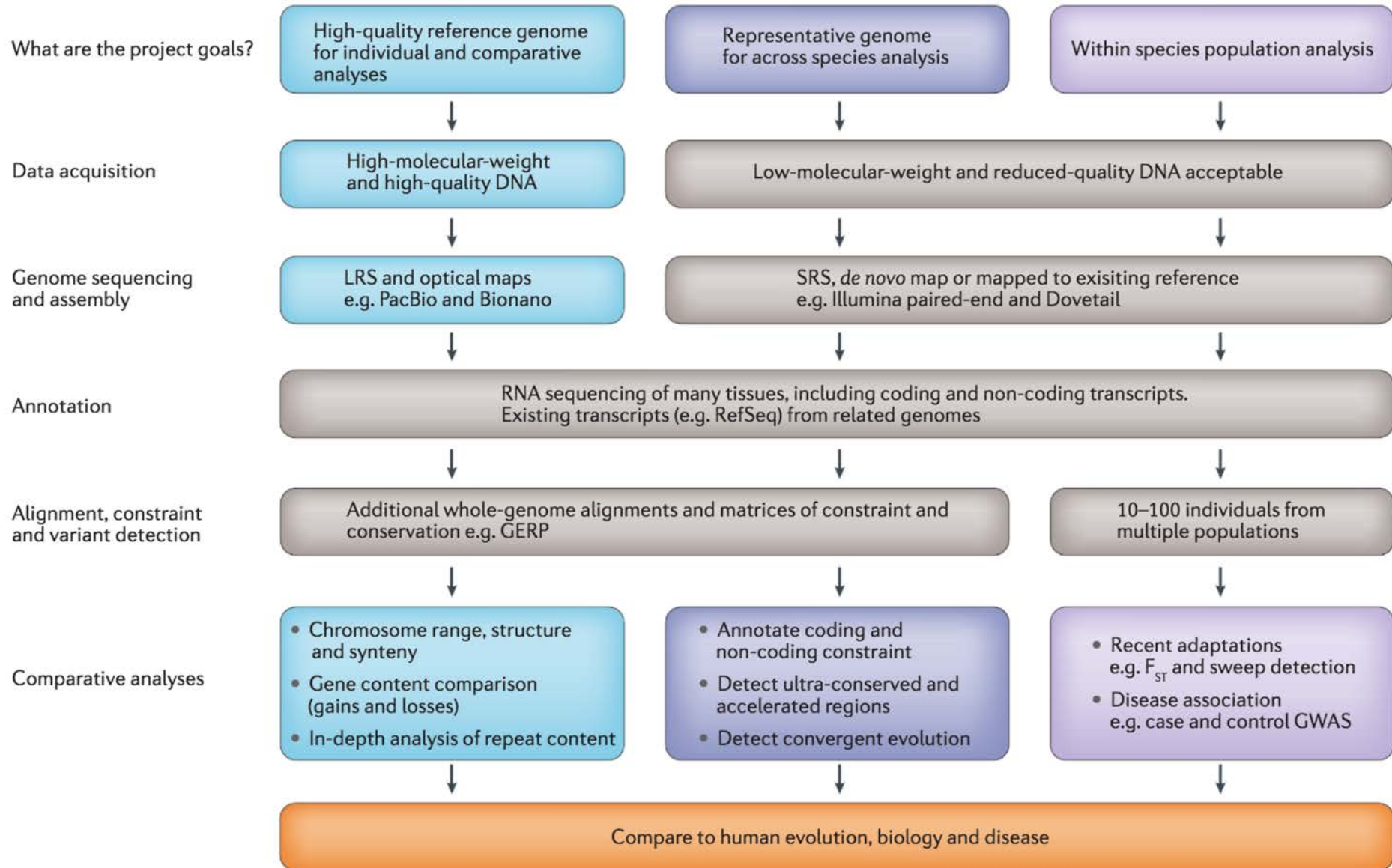
## PhastCons



# Global patterns of evolution for different aspects of the transcriptome



# Designing a sequencing project: 2017 version





## Evolution at chromosome level

- Autosomes vs. Sex chromosomes
- Euchromatin / Heterochromatin
- Chromosome arm versus center
- Large and small chromosomes
- Ploidy and polysomy

# Looking back in 2003

Group	Species	Common	Size (Mb)	Chromosome (1N)	Gene no.	Repeat %
Mammal	<i>Homo sapiens</i>	Human	2900	23	30,000	46
Mammal	<i>Mus musculus</i>	House mouse	2500	20	30,000	38
Fish	<i>Takifugu rubripes</i>	Tiger pufferfish	400	22 (?)	30,000	<10
Urochordate	<i>Ciona intestinalis</i>	Sea squirt	155	14	16,000	~10
Insect	<i>Anopheles gambiae</i>	Malaria mosquito	280	3	14,000	16
Insect	<i>Drosophila melanogaster</i>	Fruit fly	137	4	13,600	2
Nematode	<i>Caenorhabditis elegans</i>	Nematode worm	97	6	19,100	<1
Apicomplexa	<i>Plasmodium falciparum</i>	Human malaria parasite	23	14	5,300	<1
Apicomplexa	<i>Plasmodium yoelli</i>	Rodent malaria parasite	25	14	5,300	<1
Dictyosteliida	<i>Dictyostelium discoideum</i> *	Social amoeba	34	6	2,800	<1
Protozoan	<i>Leishmania major</i> *	Intracellular parasite	34	36	9,800	<1
Fungi	<i>Saccharomyces cerevisiae</i>	Brewer's yeast	12	16	5,700	2.4
Fungi	<i>Schizosaccharomyces pombe</i>	Fission yeast	13.8	3	4,900	0.35
Microsporidium	<i>Encephalitozoon cuniculi</i>	Intracellular parasite	2.5	11	2,000	<0.1
Angiosperm	<i>Arabidopsis thaliana</i>	Mustard weed	125	5	25,500	14
Angiosperm	<i>Oryza sativa</i>	Rice	400	12	32000–50000	?

**Chromosomal Rearrangements and Repeats: Cause or Consequence?**

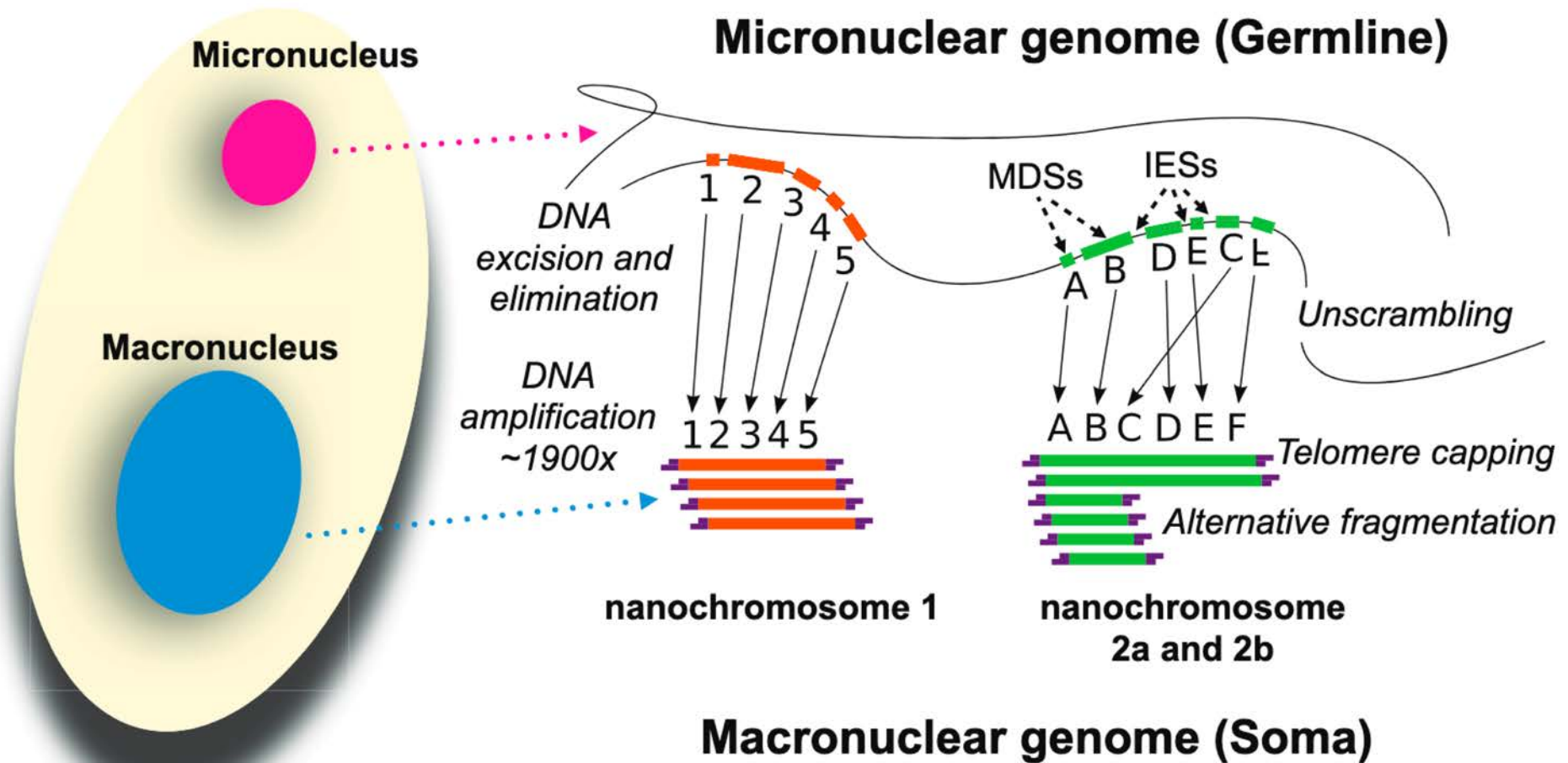
**Duplications: Engines of Gene and Genome Evolution?**

**Centromeric and Telomeric Regions— Sites of Rapid Genomic Change**

**Synteny: Fragile Versus Random Breakage Model?**

# The *Oxytricha trifallax* Macronuclear Genome: A Complex Eukaryotic Genome with 16,000 Tiny Chromosomes

Estienne C. Swart<sup>1</sup>,  
 Jaspreet S. Khurana<sup>2,3</sup>,  
 Robert S. Fulton<sup>2,3</sup>,  
 John C. Matese<sup>7</sup>, La  
 Thomas A. Jones<sup>6</sup>,  
 Elaine R. Mardis<sup>2,3</sup>,



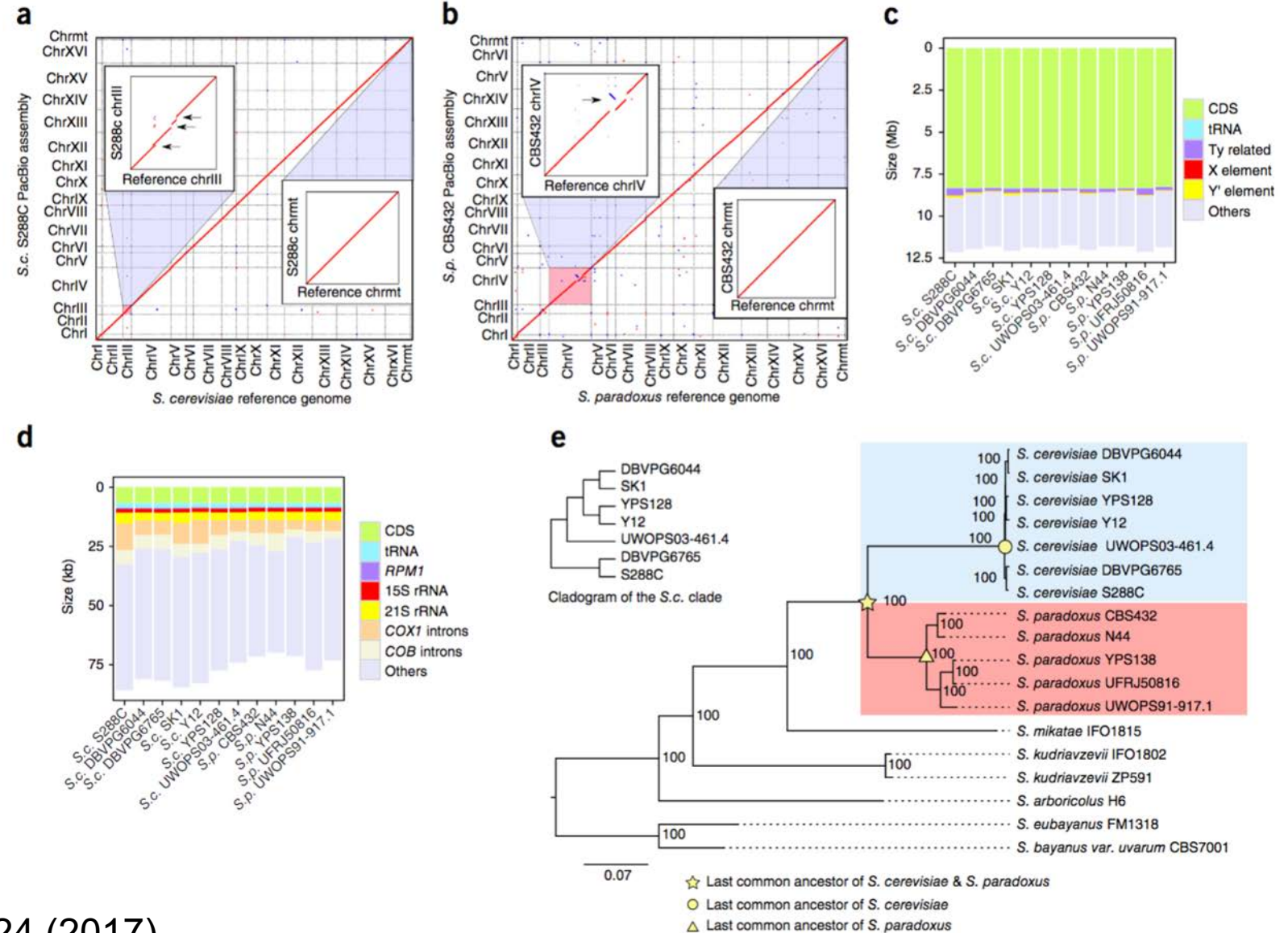
	Genome size	Genes	Chromosomes	Ploidy	Alternative fragmentation
<i>Oxytricha</i>	~50 Mb	~18,400	~15,600	Variable ~1,900 <sup>a</sup>	Yes
<i>Stylonychia</i>	~50 Mb <sup>b</sup>	~12,000 <sup>c</sup>	~10-15,000 <sup>b</sup>	Variable ~15,000 <sup>b</sup>	Yes
<i>Euplotes</i>	~50 Mb?	?	nano ?	Variable ~2,000 <sup>d</sup>	No?
<i>Nyctotherus</i>	~50 Mb <sup>e</sup>	?	nano ?	Variable <sup>e</sup> ?	?
<i>Tetrahymena</i>	105 Mb <sup>f</sup>	24,700 <sup>g</sup>	225 <sup>f</sup>	45 <sup>f</sup>	limited <sup>h,i</sup>
<i>Ichthyophthirius</i>	49 Mb <sup>j</sup>	8,100 <sup>j</sup>	71 <sup>j</sup>	~12,000 <sup>j</sup>	?
<i>Paramecium</i>	72 Mb <sup>k</sup>	40,000 <sup>k</sup>	~200 <sup>k</sup>	~800 <sup>l</sup>	limited <sup>l</sup>
<i>Perkinsus</i>	87 Mb	23,700	?	1	NA
<i>Plasmodium</i>	23 Mb <sup>m</sup>	5,300 <sup>m</sup>	14 <sup>m</sup>	1	NA

0.08

# Contrasting evolutionary genome dynamics between domesticated and wild yeasts

Jia-Xing Yue<sup>1</sup>, Jing Li<sup>1</sup>, Louise Aigrain<sup>2</sup>, Johan Hallin<sup>1</sup>, Karl Persson<sup>3</sup>, Karen Oliver<sup>2</sup>, Anders Bergström<sup>2</sup>, Paul Coupland<sup>2,5</sup>, Jonas Warringer<sup>3</sup>, Marco Cosentino Lagomarsino<sup>4</sup>, Gilles Fischer<sup>4</sup>, Richard Durbin<sup>2</sup> & Gianni Liti<sup>1</sup>

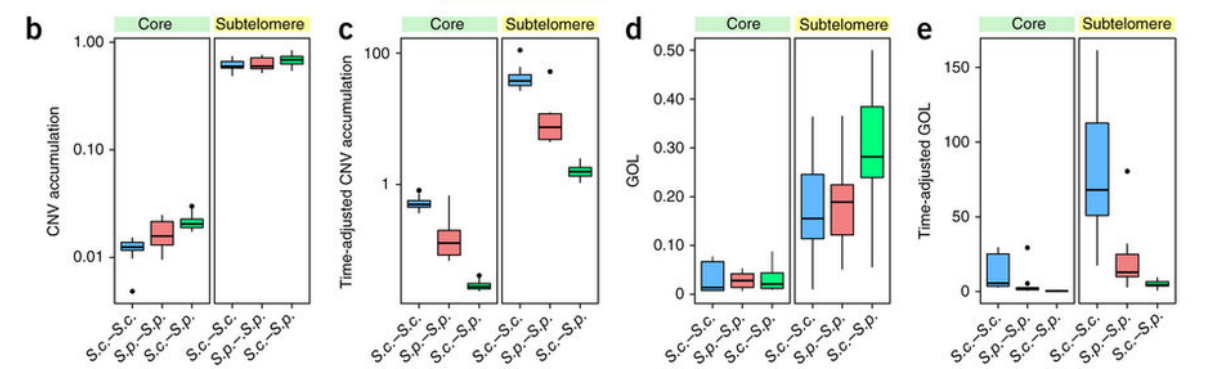
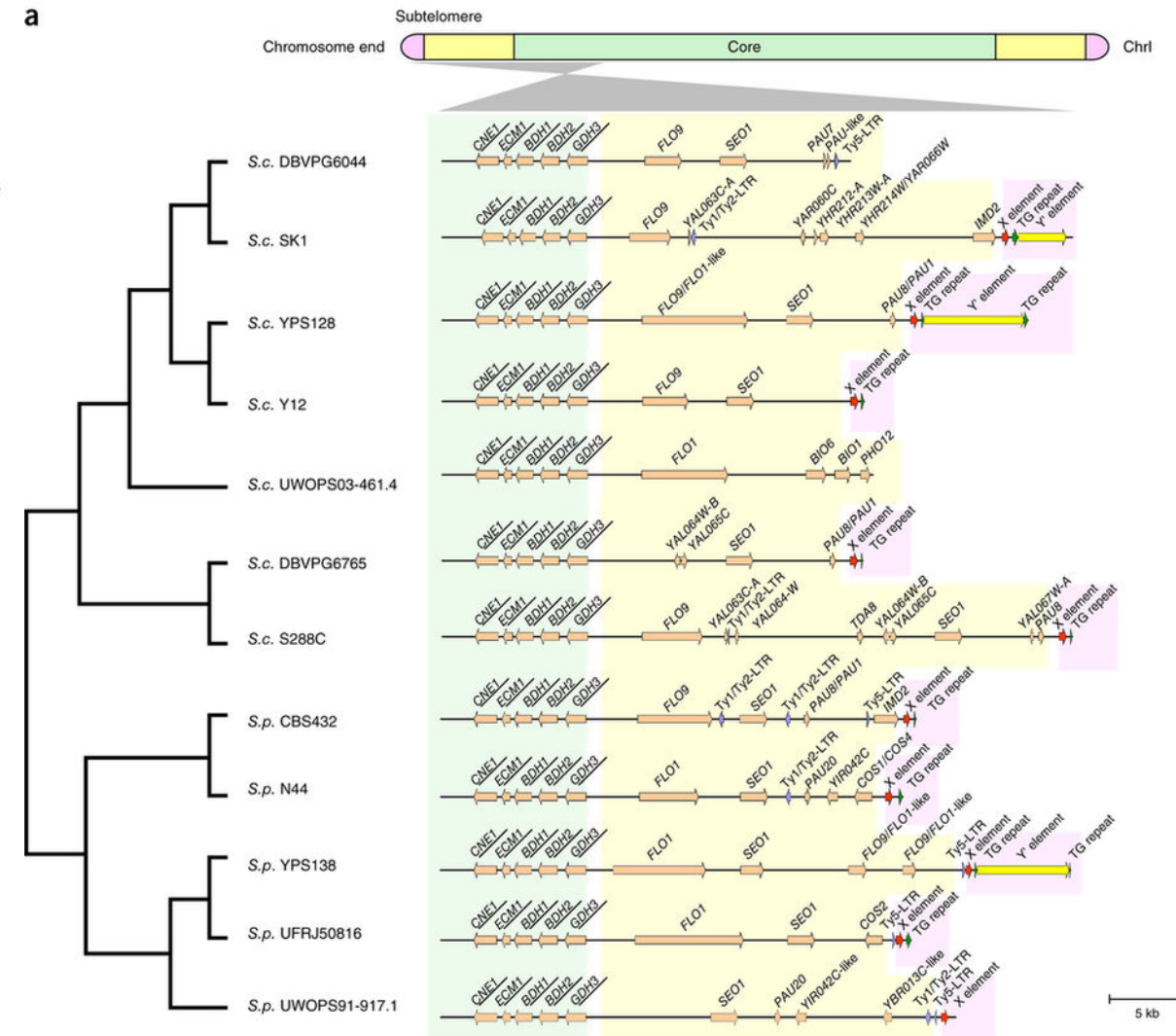
- long-read sequencing to generate **end-to-end genome assemblies** for **12 strains** representing major subpopulations of the partially domesticated yeast *Saccharomyces cerevisiae* and its wild relative *S. paradoxus*.



# Contrasting evolutionary genome dynamics between domesticated and wild yeasts

Jia-Xing Yue<sup>1</sup>, Jing Li<sup>1</sup>, Louise Aigrain<sup>2</sup>, Johan Hallin<sup>1</sup>, Karl Persson<sup>3</sup>, Karen Oliver<sup>2</sup>, Anders Bergström<sup>2</sup>, Paul Coupland<sup>2,5</sup>, Jonas Warringer<sup>3</sup>, Marco Cosentino Lagomarsino<sup>4</sup>, Gilles Fischer<sup>4</sup>, Richard Durbin<sup>2</sup> & Gianni Liti<sup>1</sup>

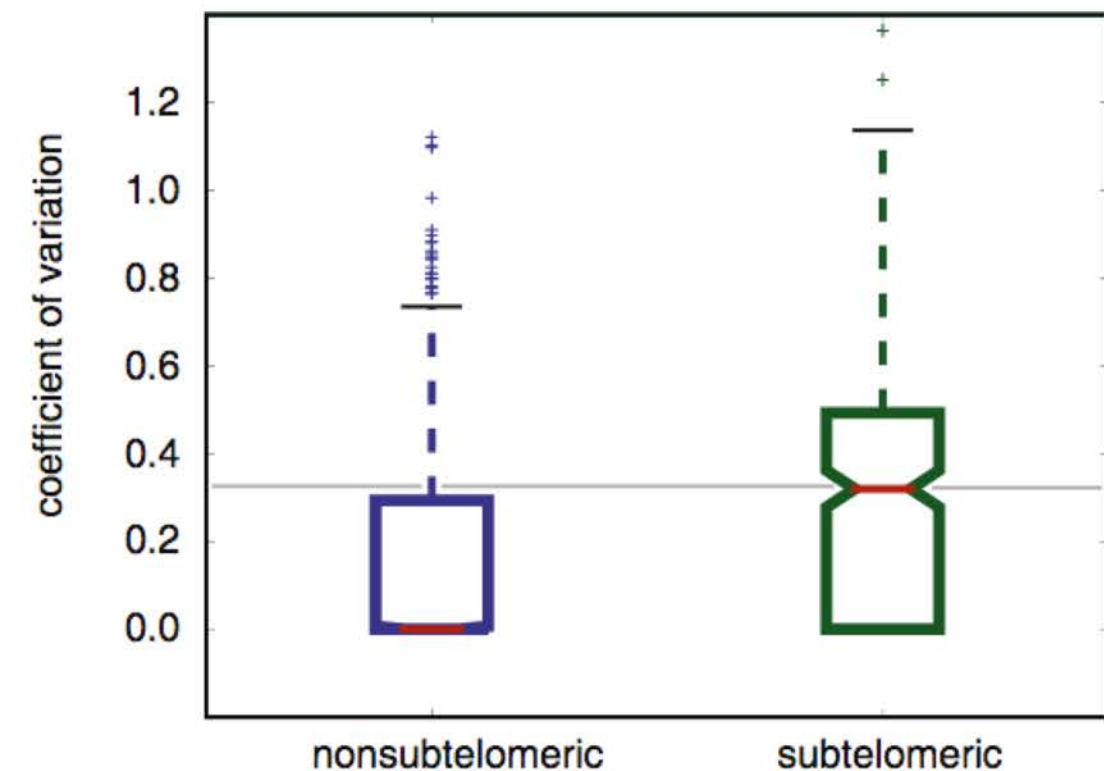
- enable precise definition of chromosomal boundaries between cores and subtelomeres
- *S. paradoxus* shows faster accumulation of balanced rearrangements (inversions, reciprocal translocations and transpositions), *S. cerevisiae* accumulates unbalanced rearrangements (novel insertions, deletions and duplications) more rapidly.
- Such striking contrasts between wild and domesticated yeasts are likely to reflect the influence of human activities on structural genome evolution.



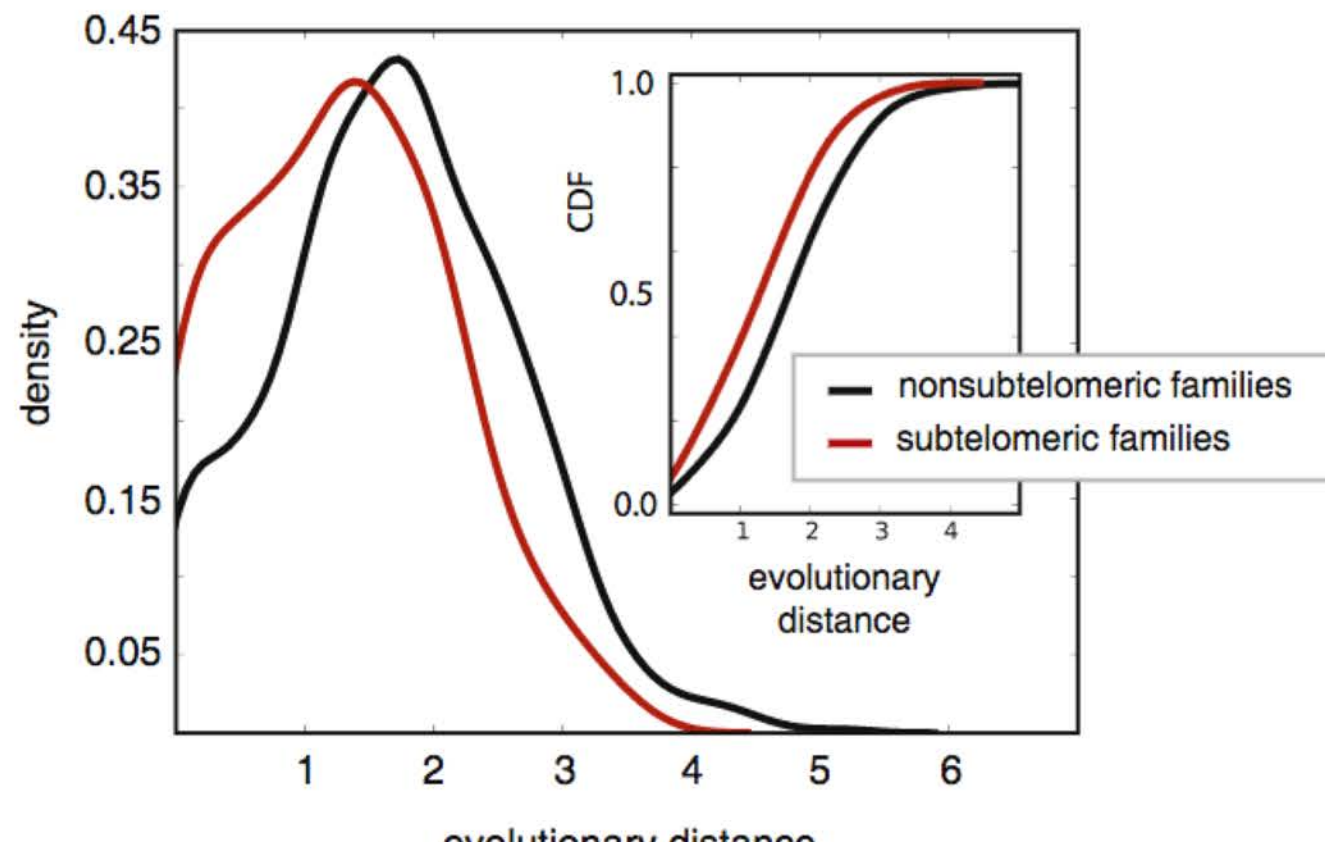
# Rapid Expansion and Functional Divergence of Subtelomeric Gene Families in Yeasts

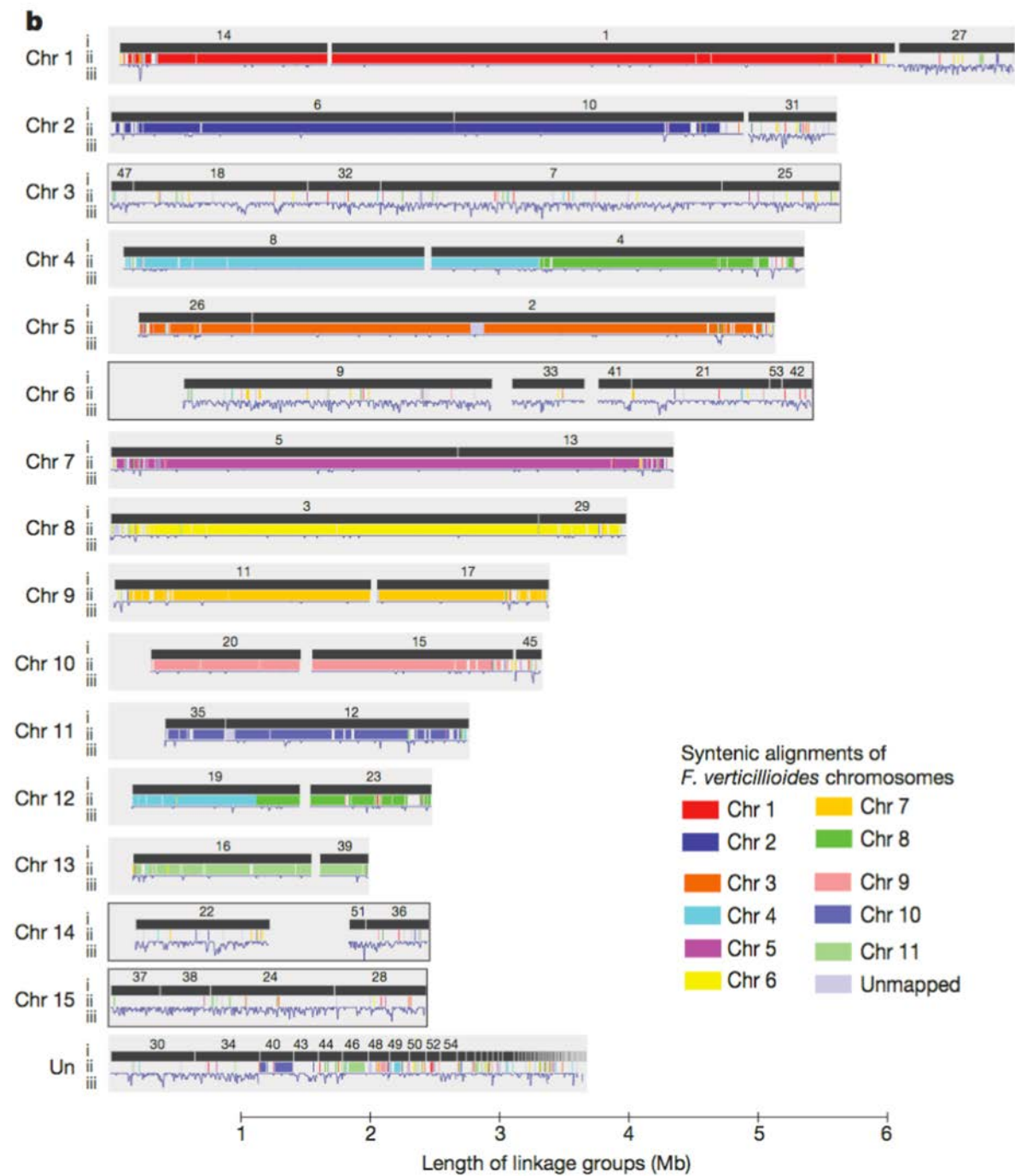
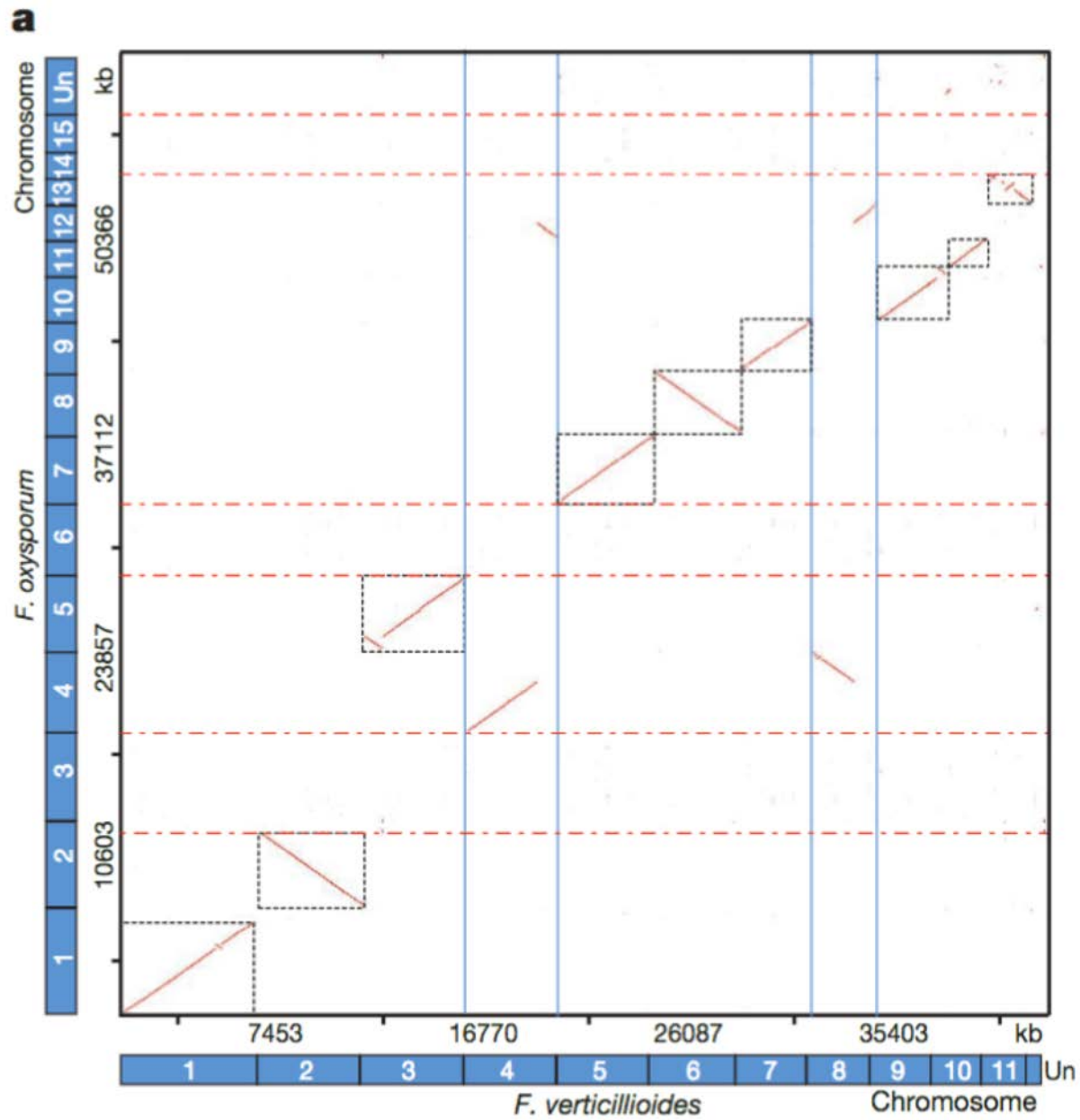
...our computational and experimental analyses show that the **extraordinary instability of eukaryotic subtelomeres supports rapid adaptation to novel niches by promoting gene recombination and duplication followed by functional divergence of the alleles**

**C** Subtelomeric Families Show More Copy Number Variation Between Species



**D** Subtelomeric Families Show More Recent Duplications

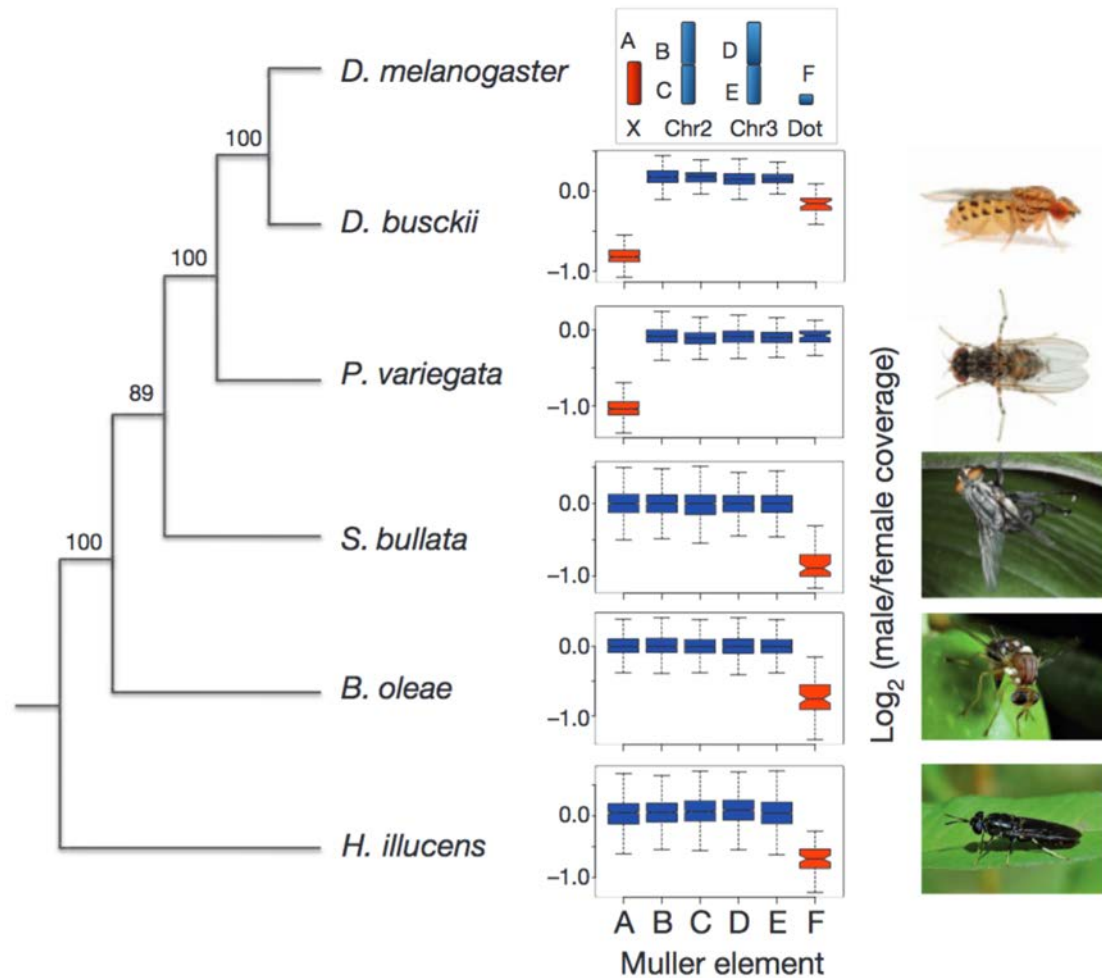






# Reversal of an ancient sex chromosome to an autosome in *Drosophila*

Beatriz Vicoso<sup>1</sup> & Doris Bachtrog<sup>1</sup>



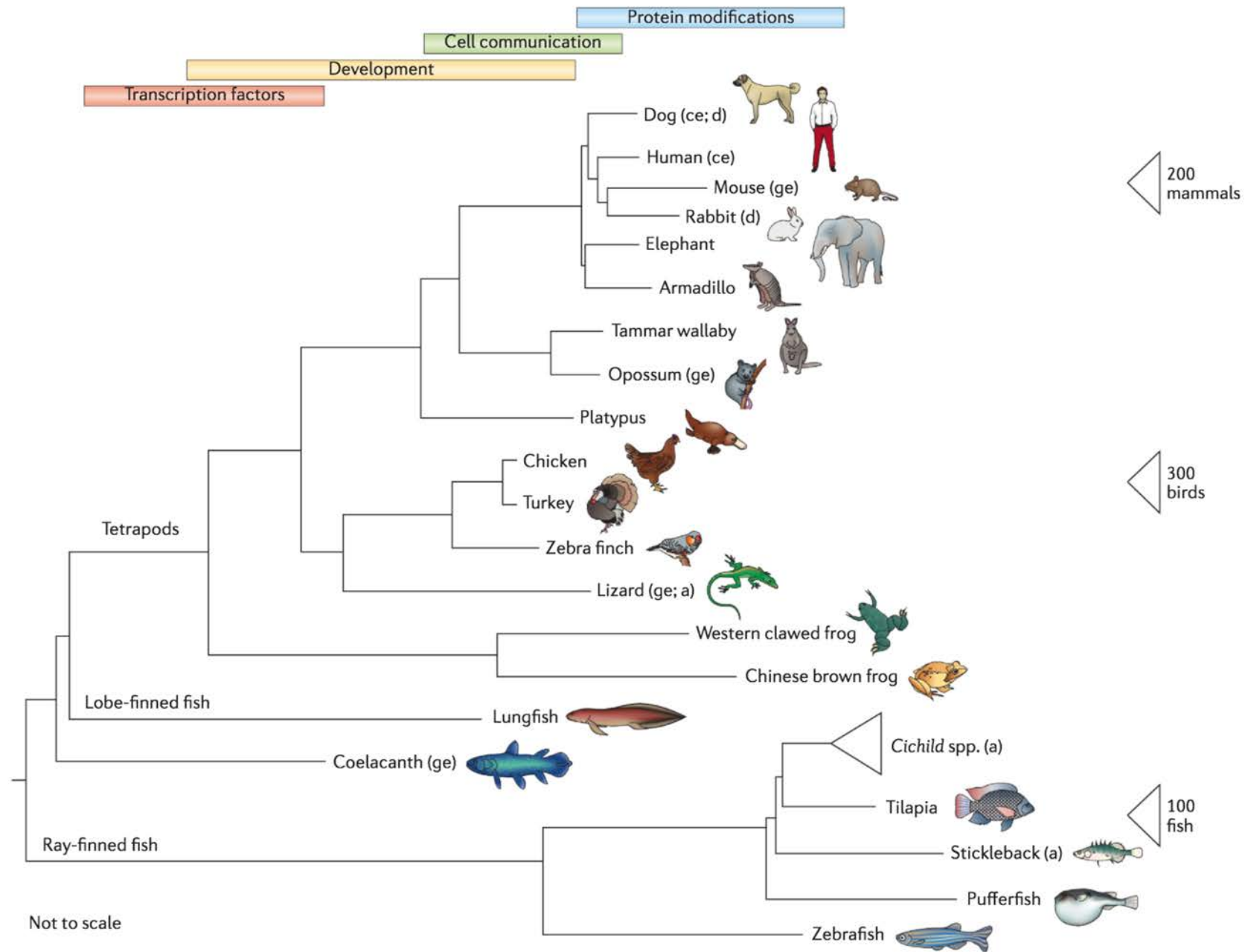
**Figure 1 | Sex chromosomes in higher Diptera revealed by genome analysis.** Evolutionary relationship inferred from 185 conserved protein-coding genes (93,134 amino acids) using PhyML (with bootstrap values indicated at the nodes), and male-to-female coverage ratio across chromosome elements (Muller elements A–F) in the Diptera species studied. X chromosomes (red) have only half the read coverage in males versus females. Boxes extend from the first to the third quartile and whiskers to the most extreme data point within 1.5 times the interquartile range.

# Good review – recent update

## Dissecting evolution and disease using comparative vertebrate genomics

*Jennifer R. S. Meadows<sup>1</sup> and Kerstin Lindblad-Toh<sup>1,2</sup>*

Abstract | With the generation of more than 100 sequenced vertebrate genomes in less than 25 years, the key question arises of how these resources can be used to inform new or ongoing projects. In the past, this diverse collection of sequences from human as well as model and non-model organisms has been used to annotate the human genome and to increase the understanding of human disease. In the future, comparative vertebrate genomics in conjunction with additional genomic resources will yield insights into the processes of genome function, evolution, speciation, selection and adaptation, as well as the quantification of species diversity. In this Review, we discuss how the genomics of non-human organisms can provide insights into vertebrate biology and how this can contribute to the understanding of human physiology and health.



# Why comparative genomics? – a summary

## How genome evolved; How genome functions

- At various resolution (whole genome, chromosomes, regions, genes, base pairs)
- Conservation, Duplication, Species specific genes
- Inferring Orthologs and paralogs
- Gene families (clusters) of paralogs, of orthologs
- Conserved or specialized domains in clusters of paralogs, orthologs
- Gene transfer, introgression between species
- Relate genotypes to phenotypes
- Identify the genomic basis of key phenotypes

# Genomics of Eukaryotic microorganisms – a summary

## How genome evolved; How genome functions

- We are only at the beginning phase of
- Untapped diversity and mechanisms waiting to be discovered
- Arguably much more fascinating (?) than animals

<https://www.nature.com/subjects/comparative-genomics>