

Functional and Comparative Genomics

Comparative and Functional Genomics

Session 4 Phylogenomics

Phylogenomics

Species tree reconstruction.

Genome-wide phylogenetic analysis (phylome).

Gene tree vs species tree.

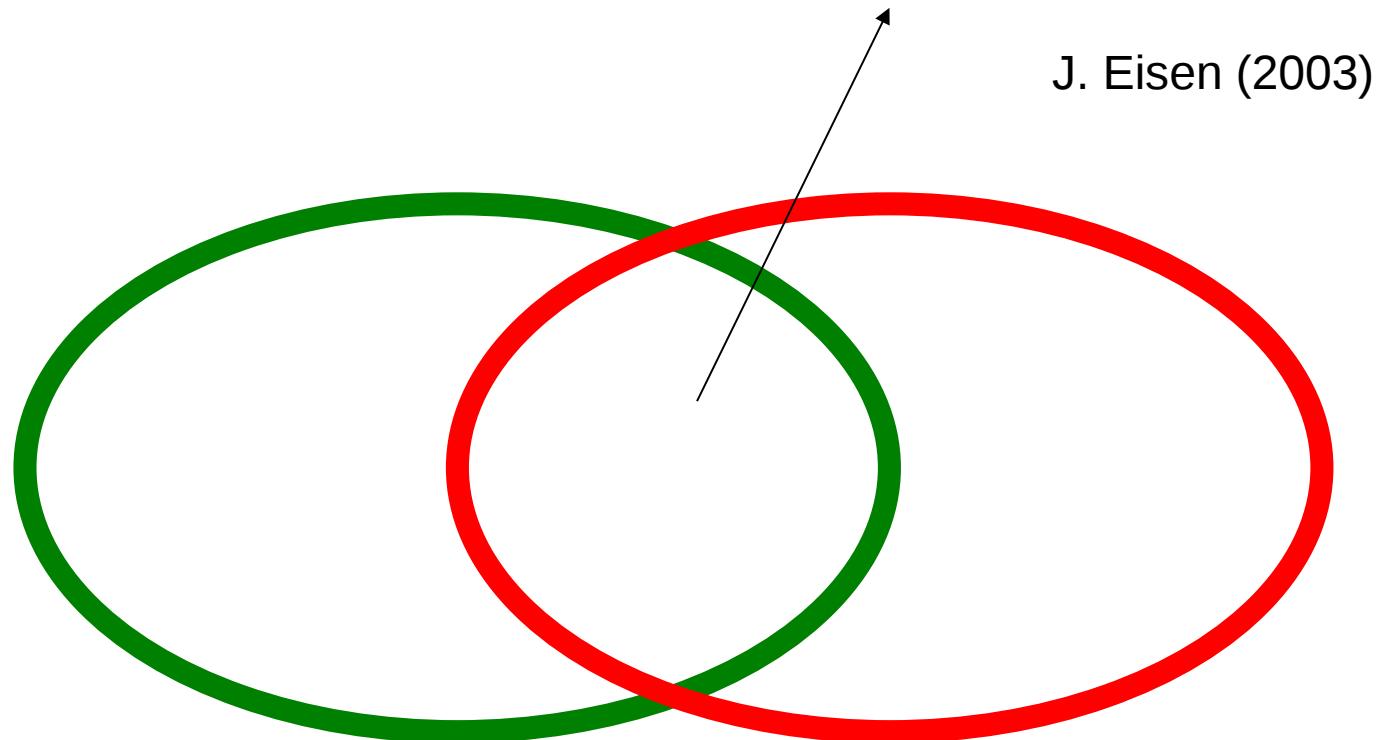
Non-vertical processes of evolution

horizontal gene transfer.

Hybridization

Whole genome duplication.

Phylogenomics is the intersection between **Genomics** and **Evolution**.



That is, looking at genomes from an evolutionary perspective, often using **phylogenetics**

protein ----> biochemical pathway ---> proteome

one species

few related species

a taxonomic kingdom

all life domains

phylogenetics

more trees

more sequences in the tree

phylogenomics

Phylogenomics:

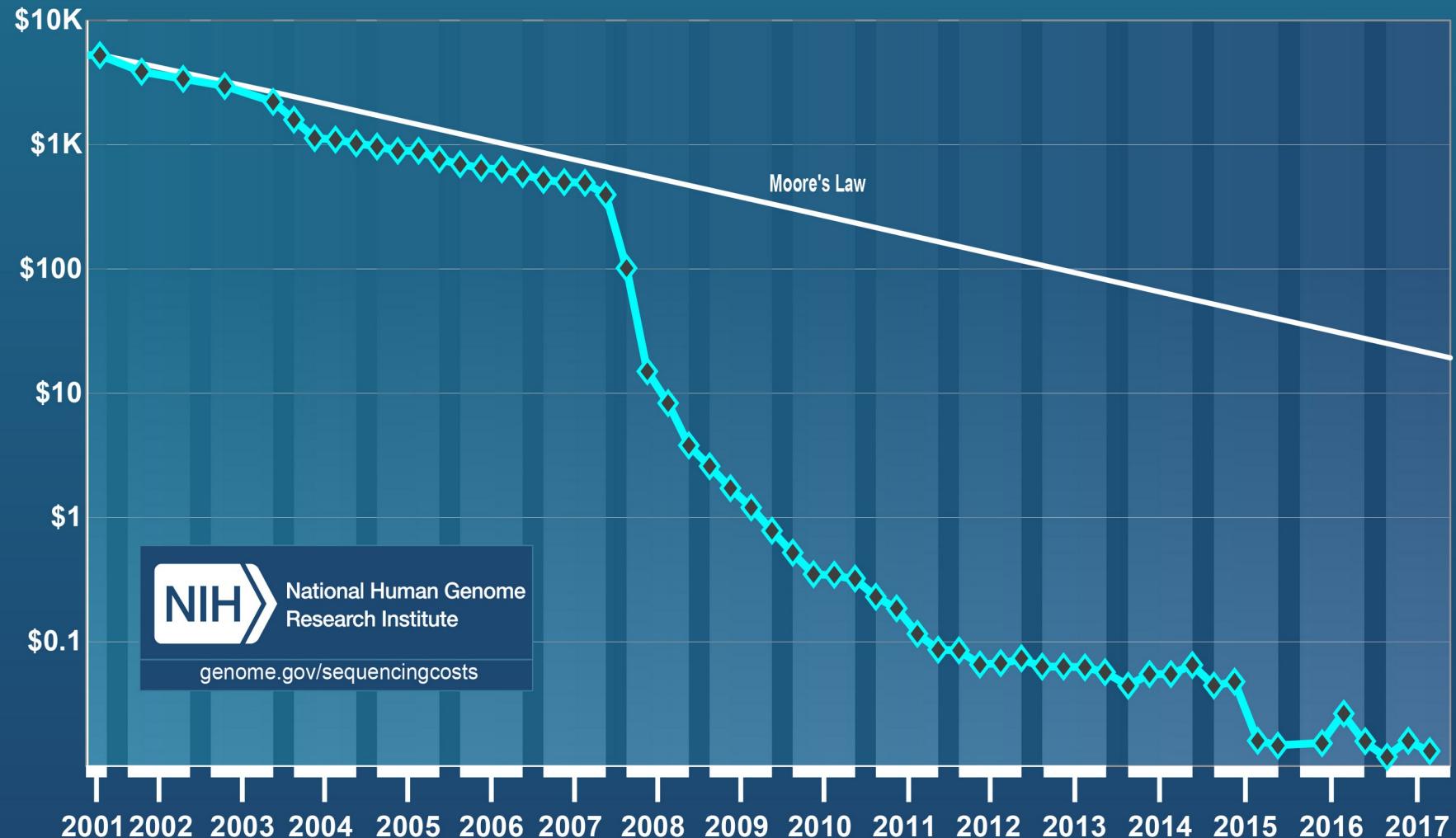
Opportunities:

- necessary to provide an evolutionary framework to the deluge of data generated
- useful to obtain biological knowledge from sequence data
- The more data, the more powerful

Challenges:

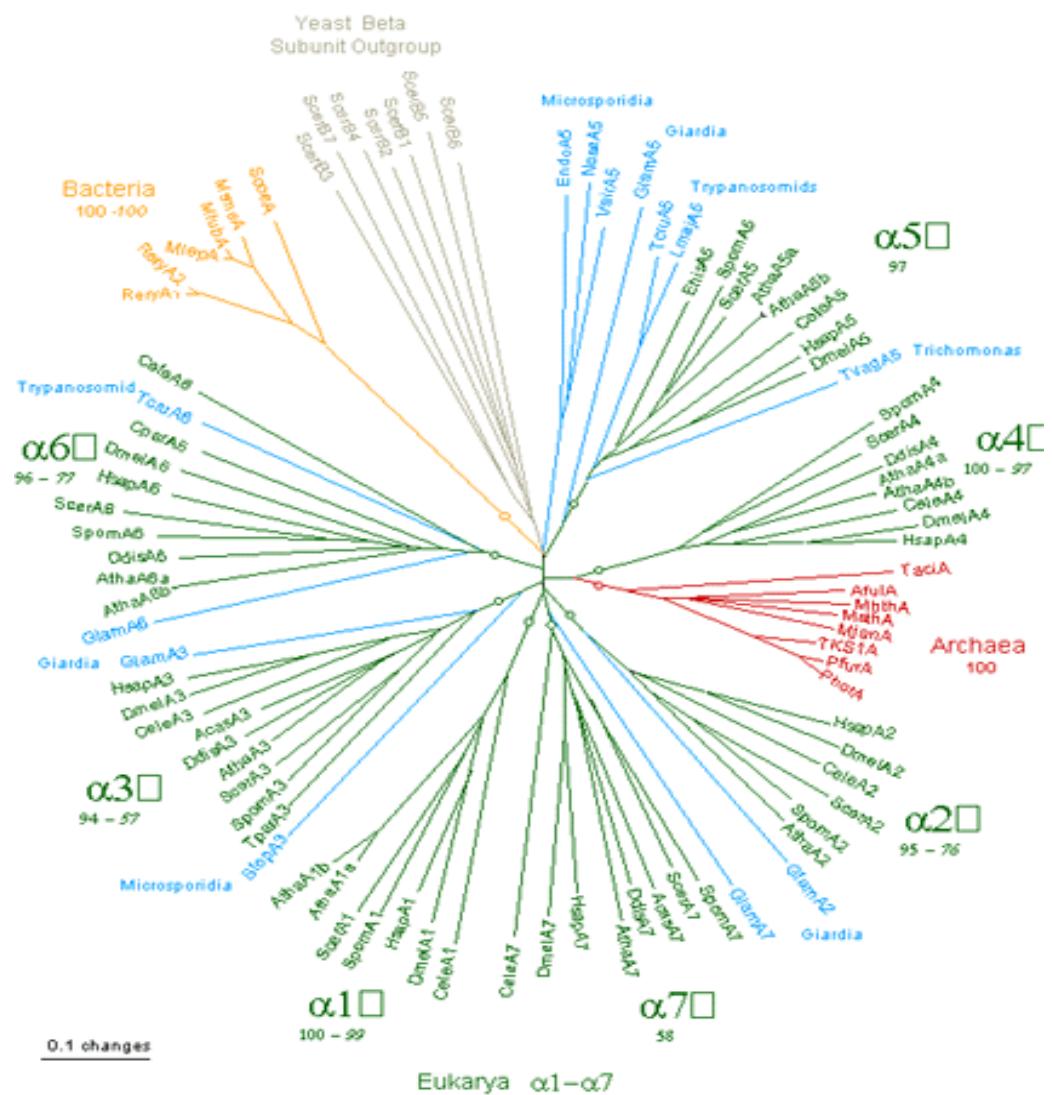
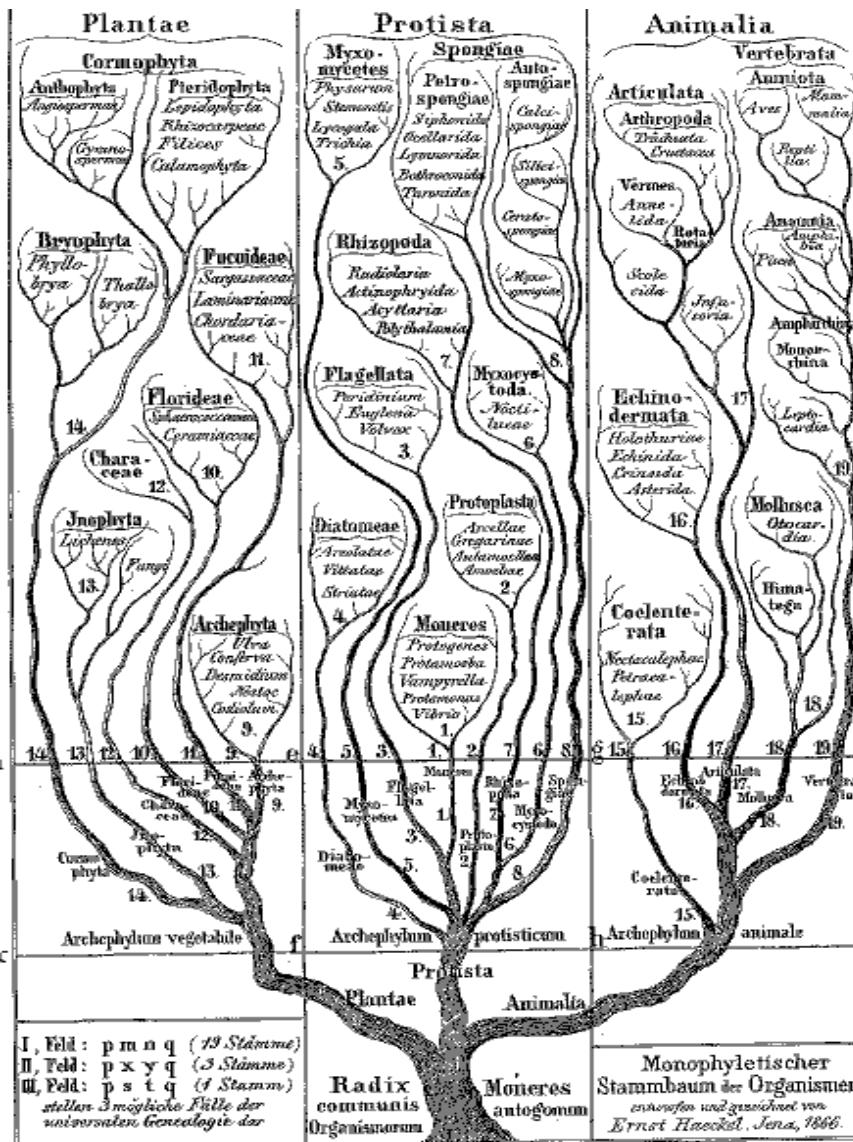
- Computationally demanding
- Need for automation
- Need for proper scalability, (e.g alignments, from a certain point, the more data, the more noisy)

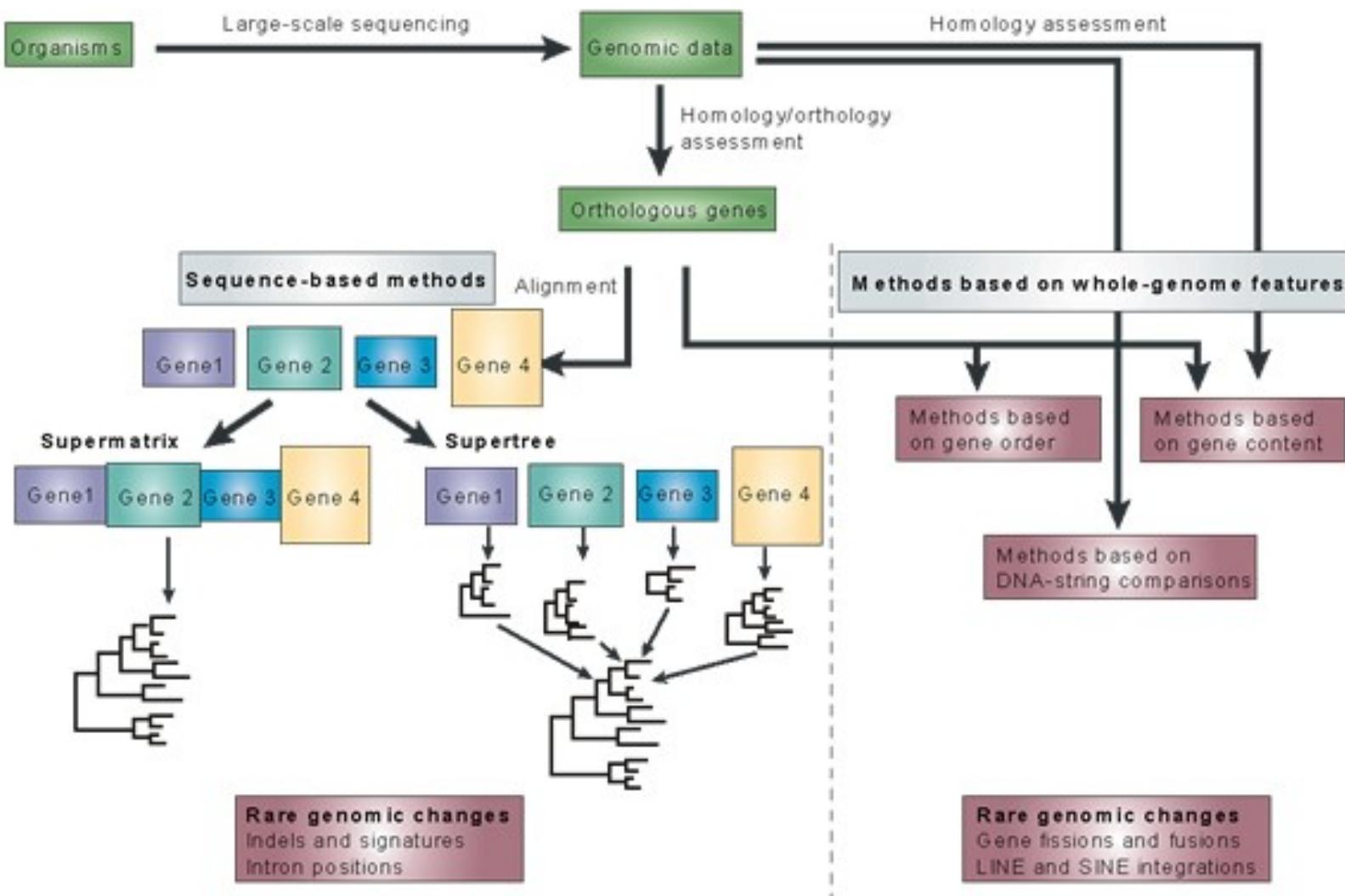
Cost per Raw Megabase of DNA Sequence



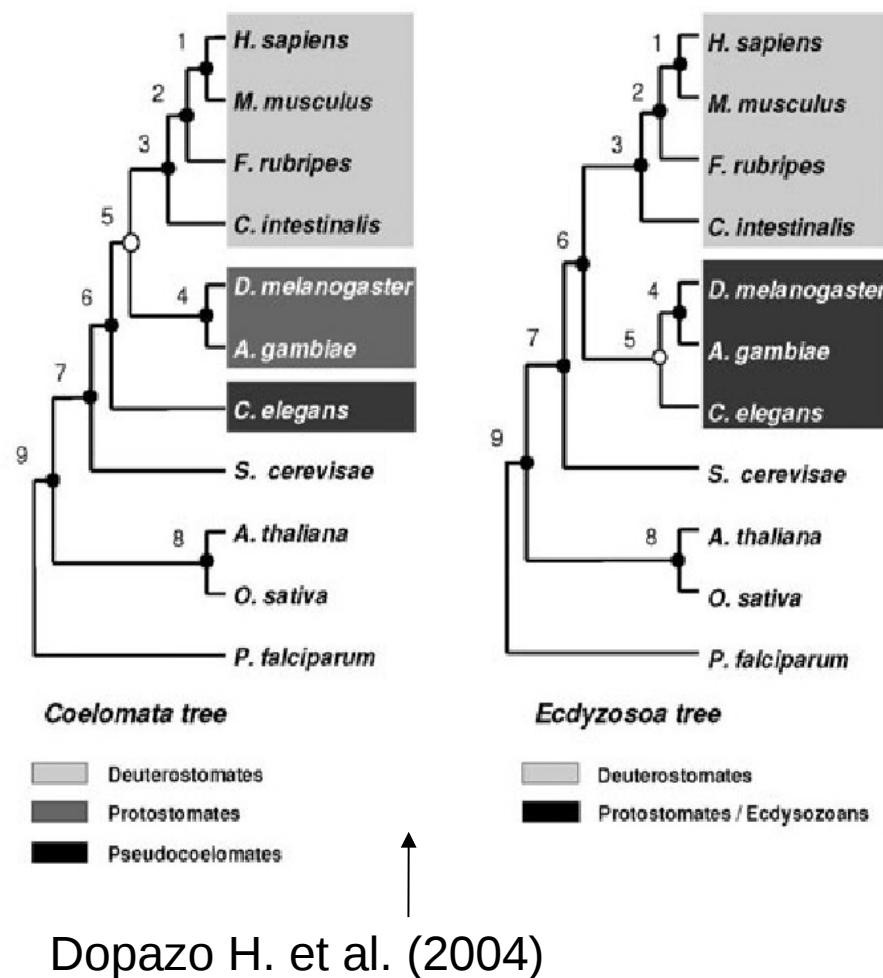
<https://www.genome.gov/sequencingcostsdata/>

Phylogenomics and species tree reconstruction

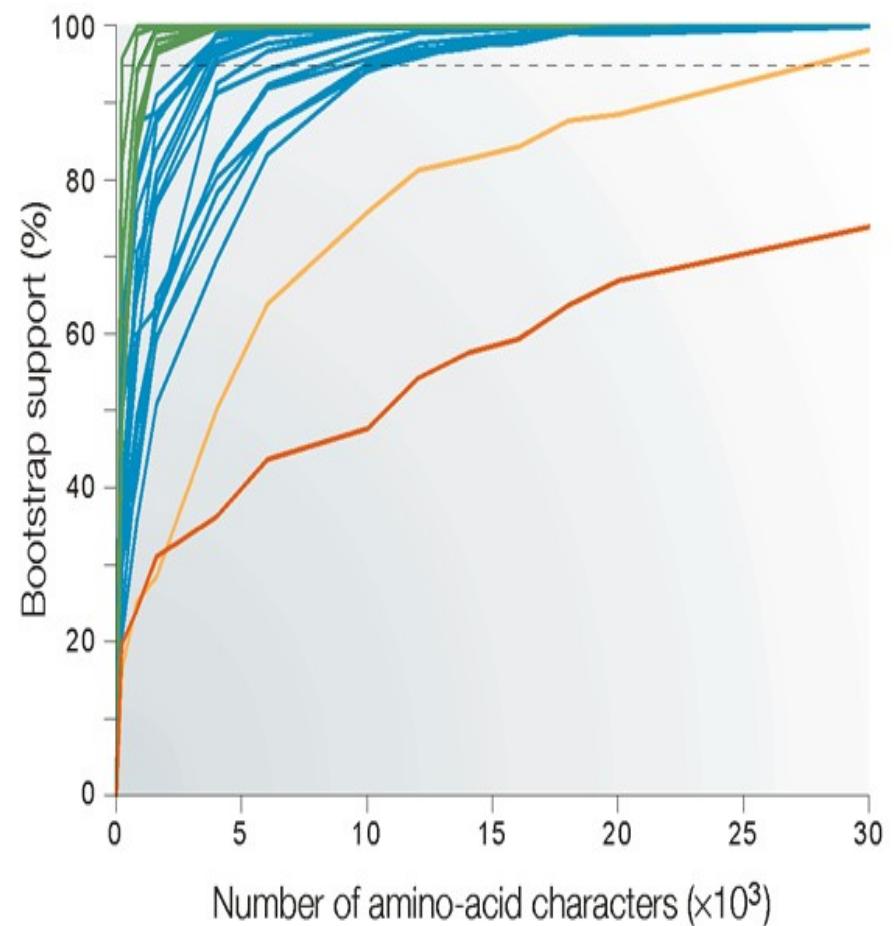




concatenated alignments allow using more sites (information) to resolve a phylogeny

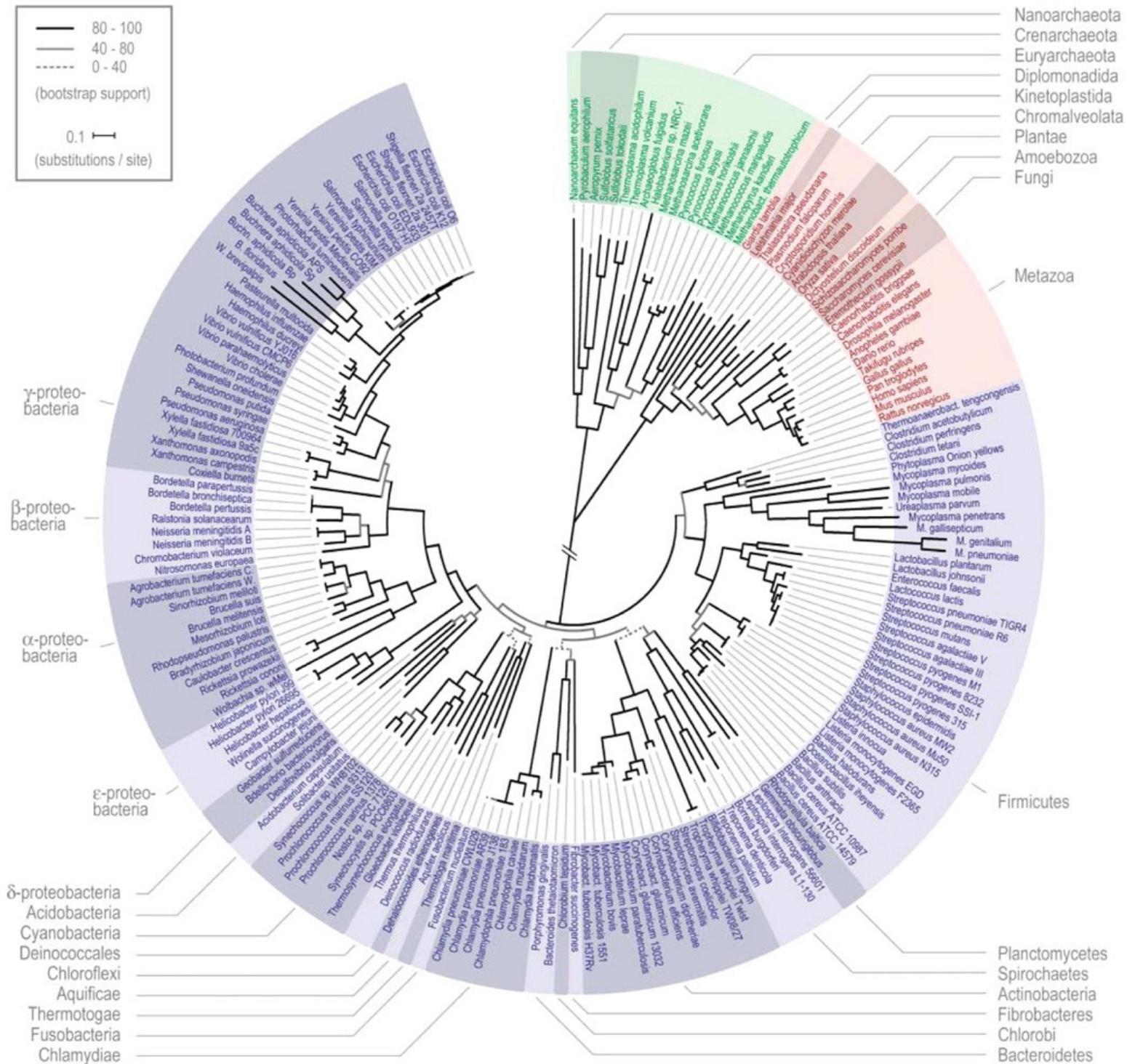


Delsuc et. al (2005)

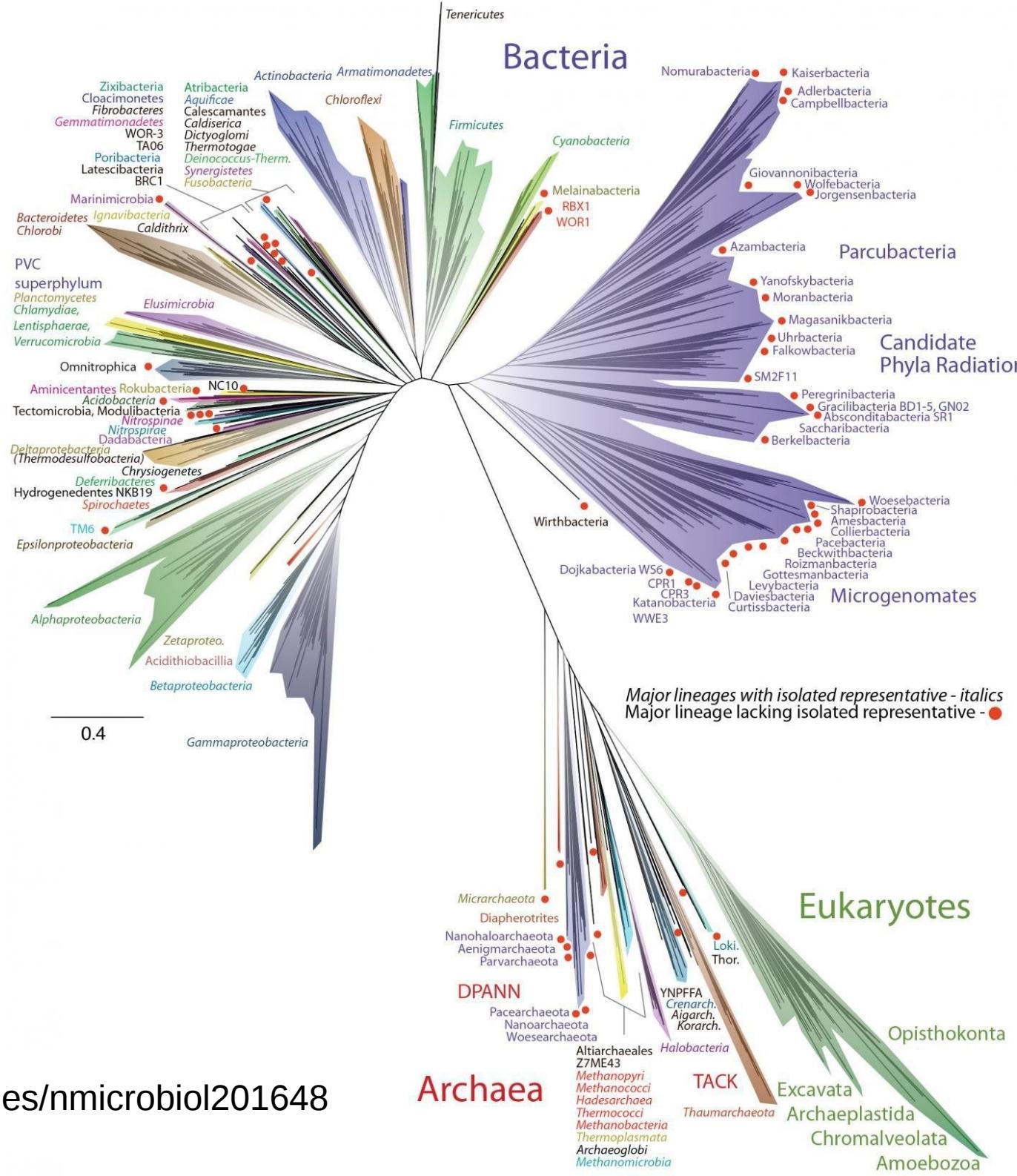


Cicarelli (2005)

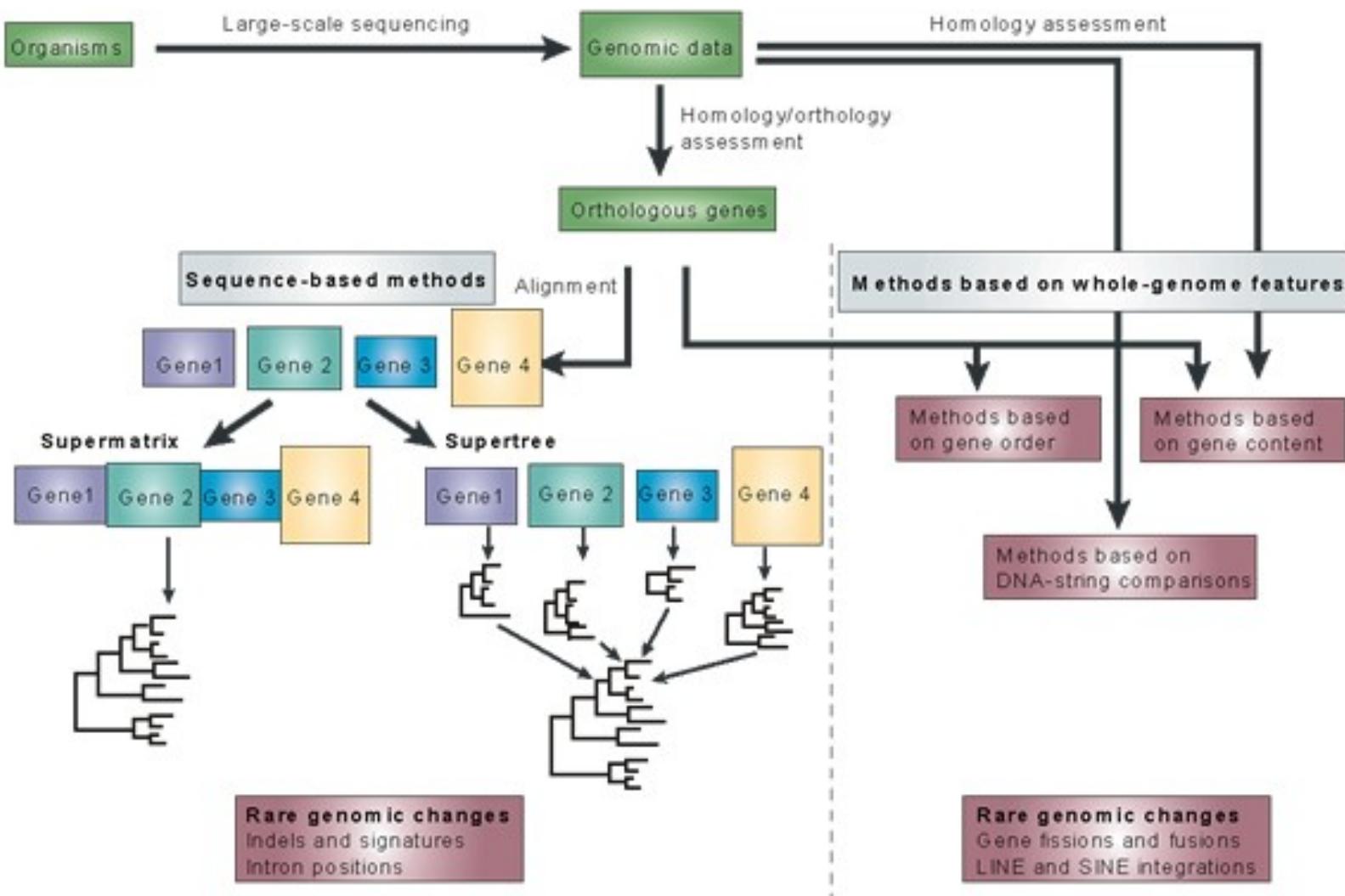
195 genomes
31 genes

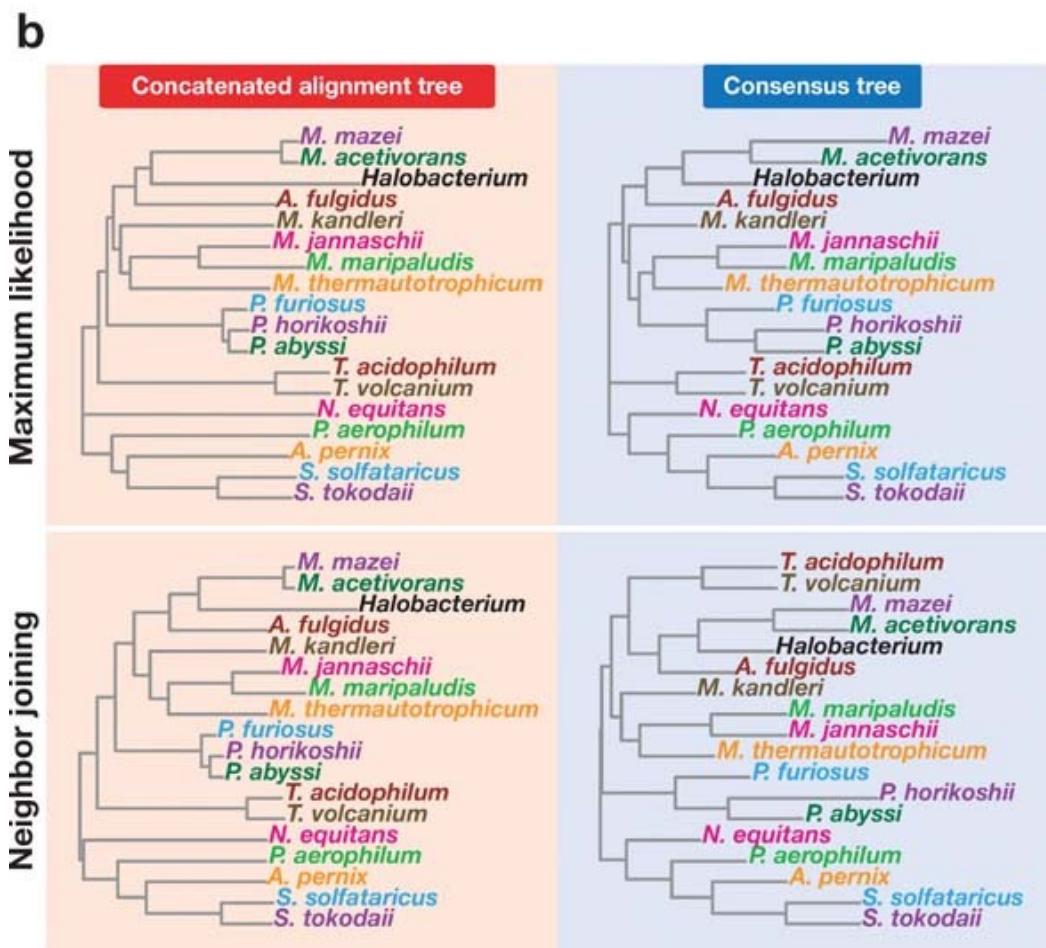
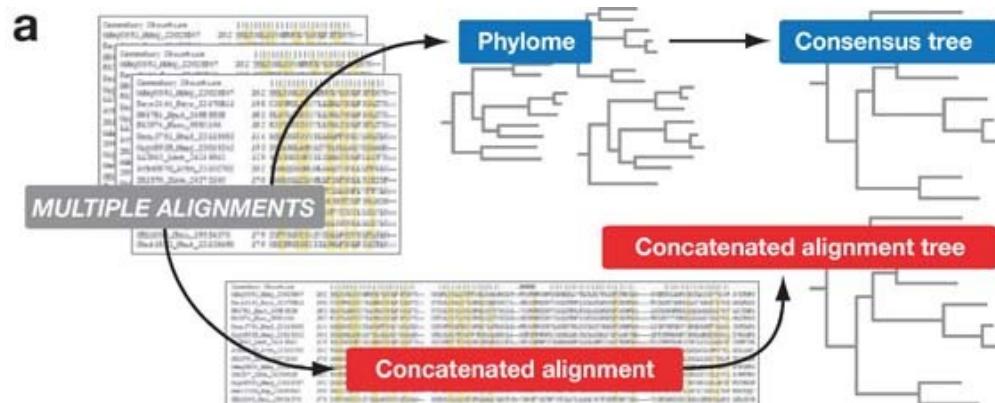


30,437 genomes
16 ribosomal genes



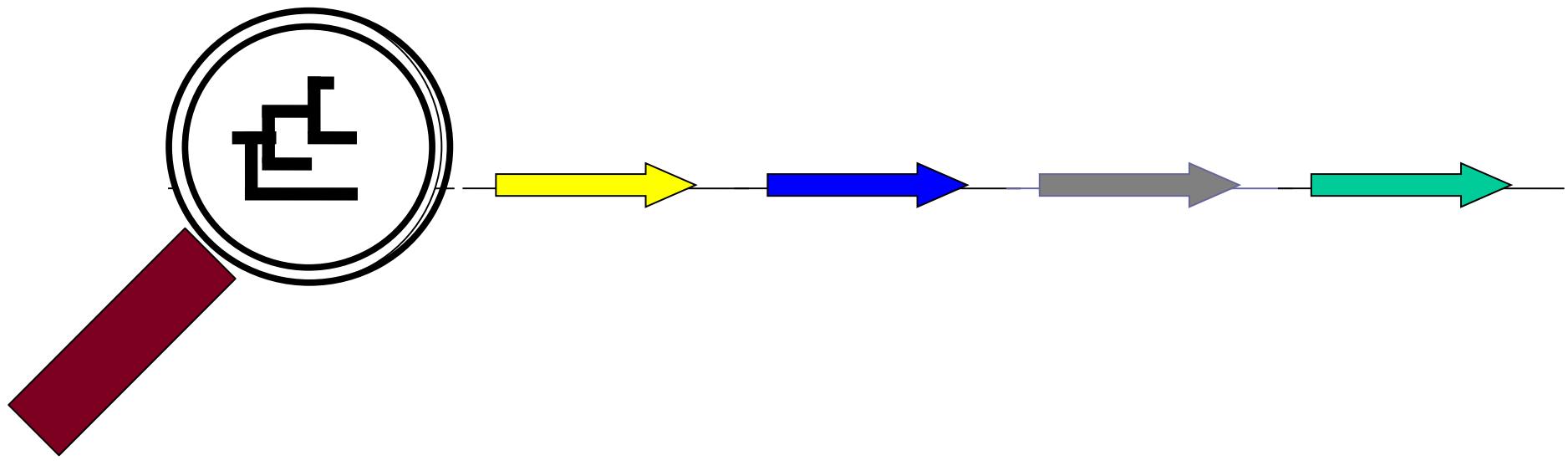
Hug. et. al. (2016)
<https://www.nature.com/articles/nmicrobiol201648>

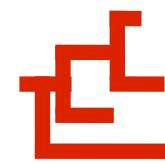
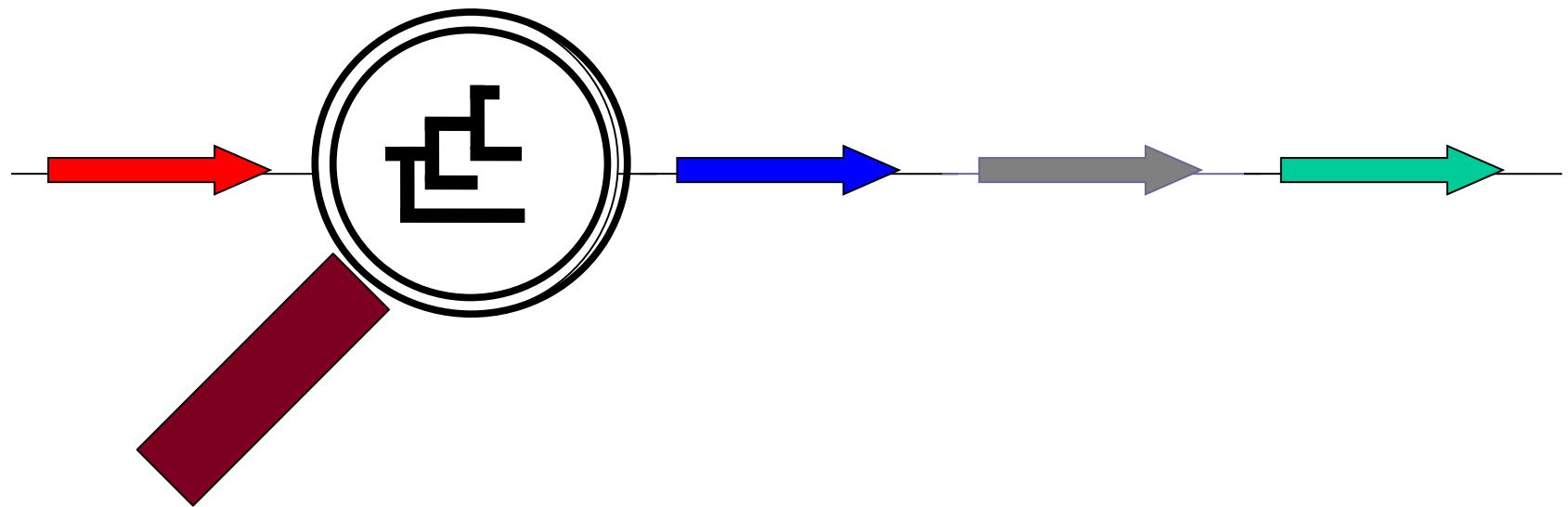


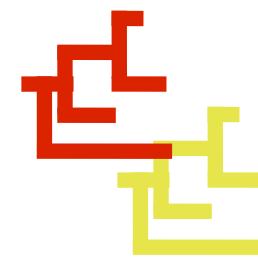
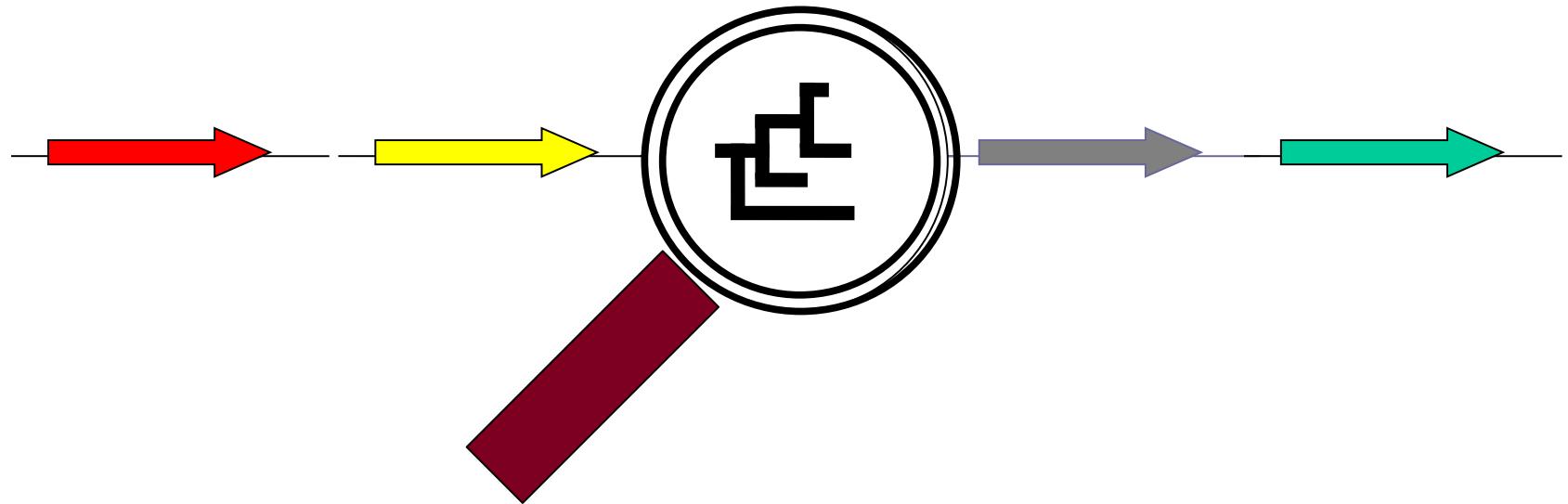


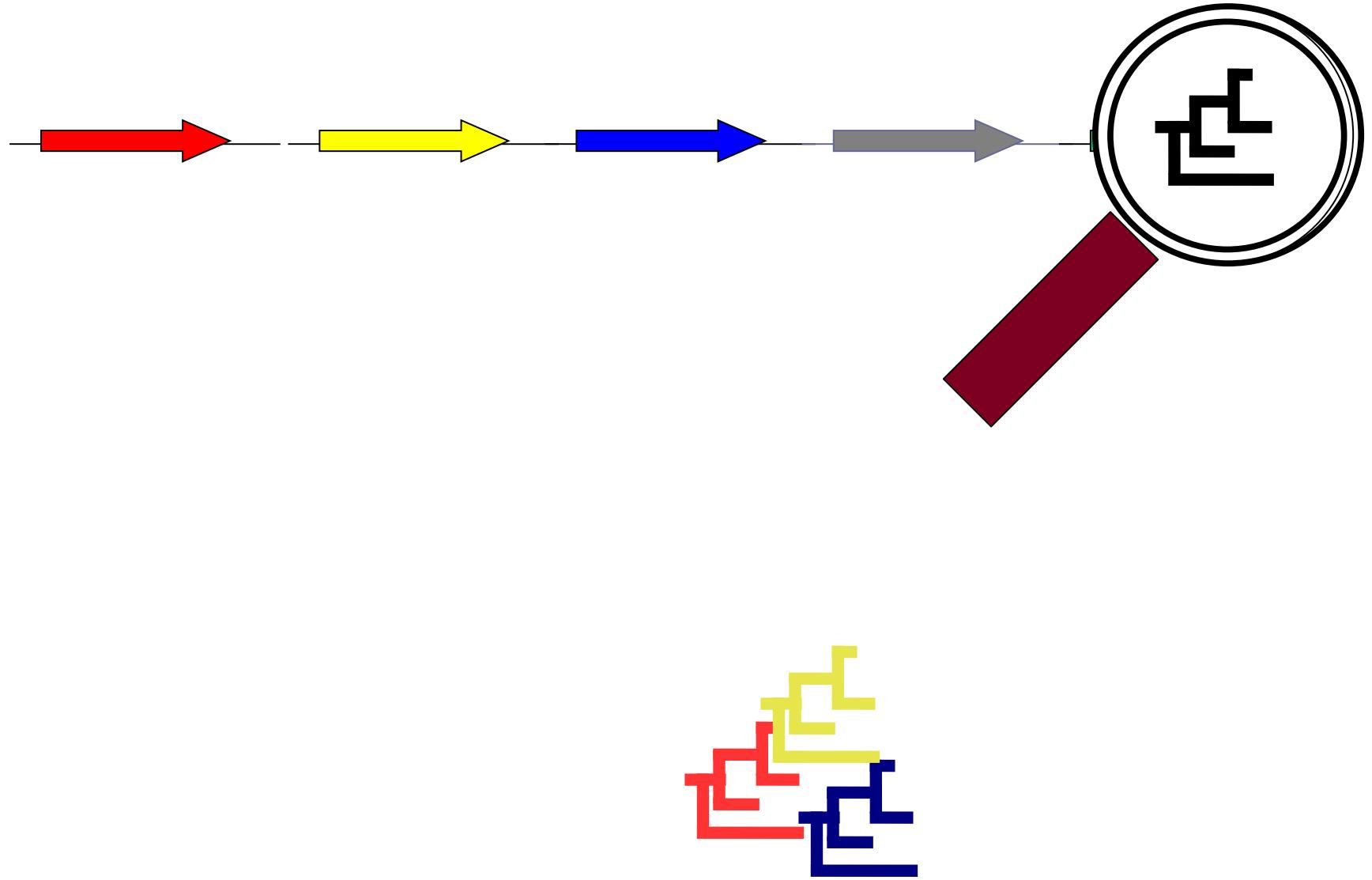
Genome-wide phylogenetic analysis (phylome).









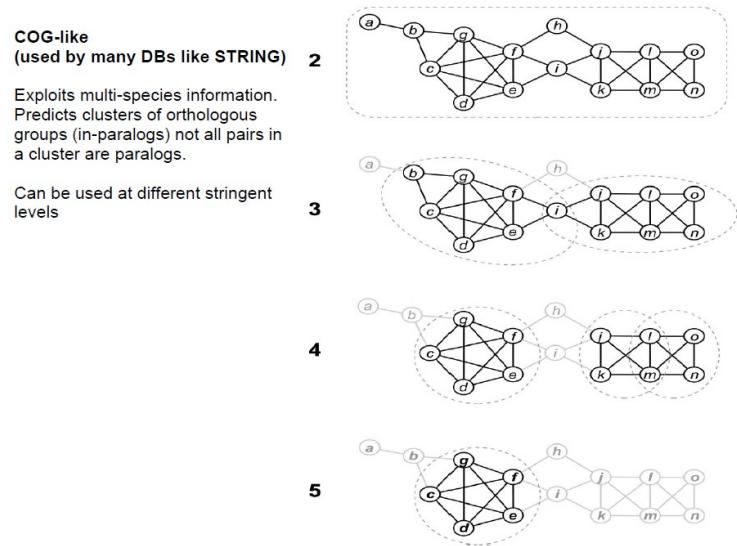


Phylome:

Complete collection of evolutionary histories of all genes encoded in a given genome

Phylome reconstruction

A) Family-based approach: first build families, then reconstruct one tree per family
(Ensembl, EggNOG..)

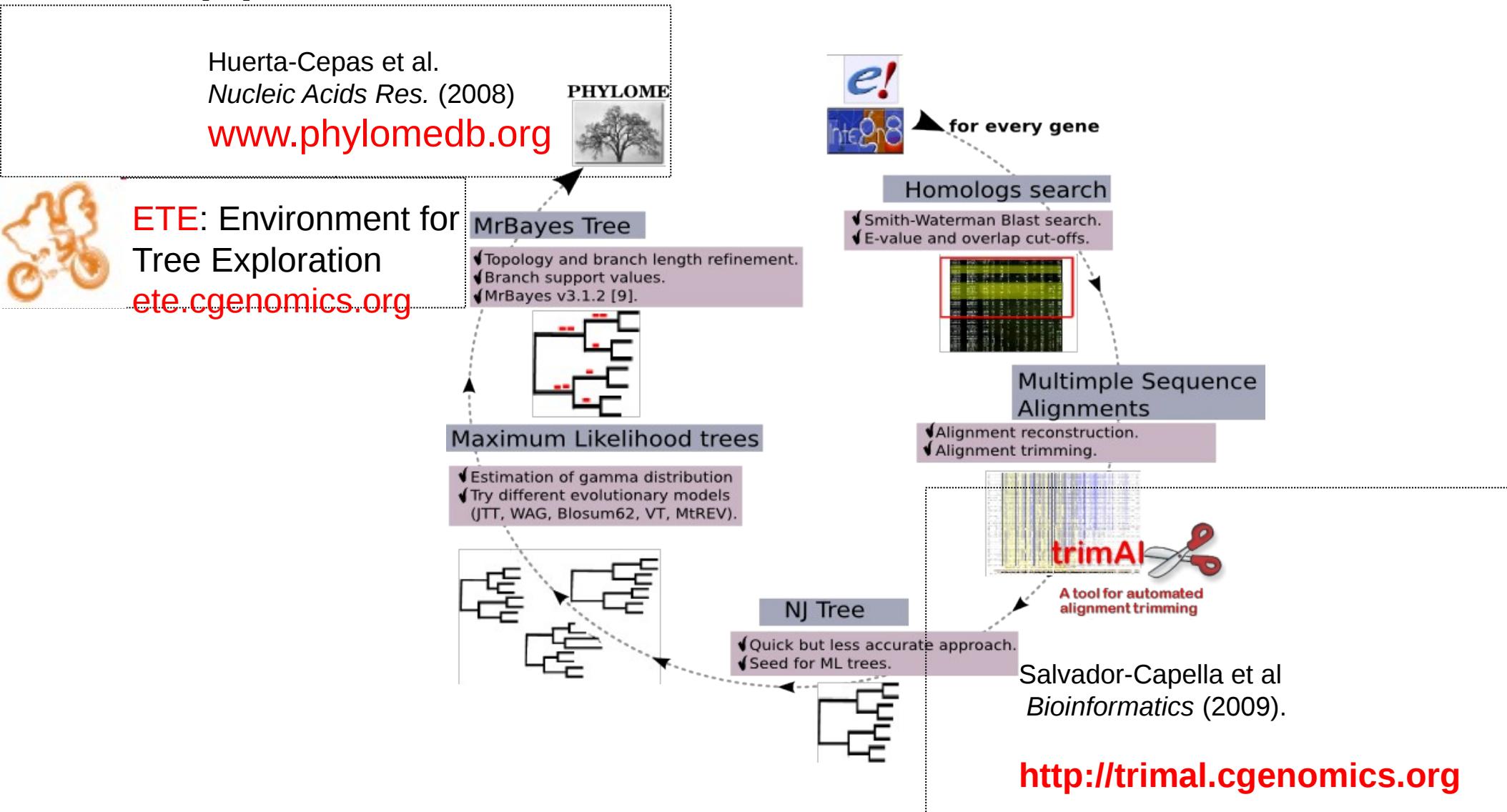


- 1- Build families
- 2- Alignment per family
- 3- Phylogenetic reconstruction per alignment

B) Gene-based approach: sequentially use each gene of interest as a seed to build a gene tree.
(PhylomeDB)

- 1- Search homologs
- 2- Alignment per gene + homologs
- 3- Phylogenetic reconstruction per alignment

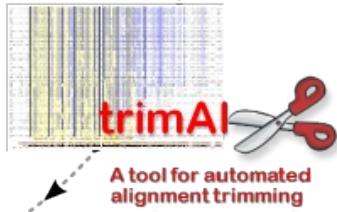
Our pipeline:



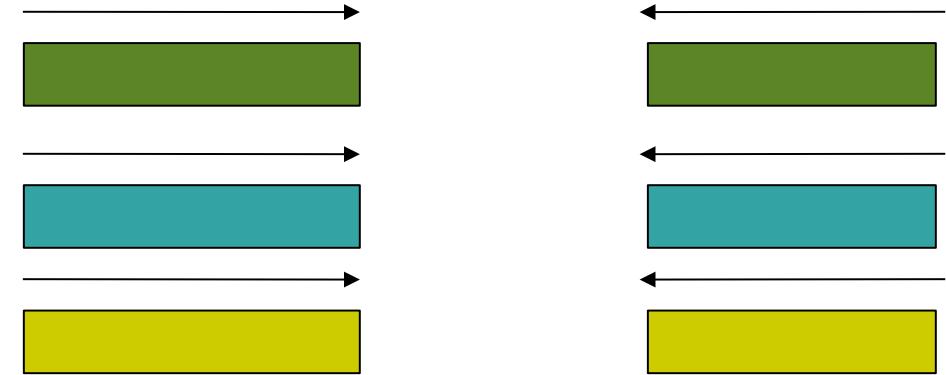
Pipeline described in Huerta-Cepas et al NAR (2011)

Multimle Sequence Alignments

- Alignment reconstruction.
- Alignment trimming.



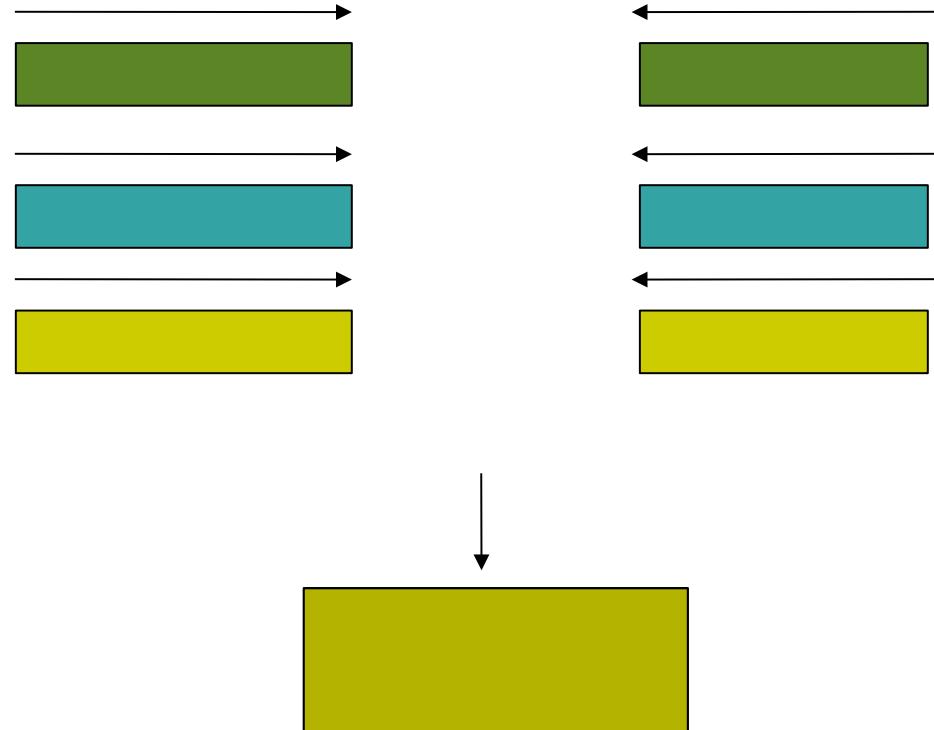
The set of homologous Sequences are aligned by 3 different aligners in forward and reverse modes (Head or Tails approach)



[Http://trimal.cgenomics.com](http://trimal.cgenomics.com)

The set of homologous
Sequences are aligned by 3 different aligners
in forward and reverse modes (Head or Tails
approach)

A consensus is built



| | |
|-----------------|---|
| sw_DSBA_PSES/1 | ---MRNL I I S A A L V A A S L F G M S A Q Q A E P I E S G K Q Y V - E L T S A V P V |
| sw_DSBA_SALTY/1 | ---M K K I W L A -- L A G M V L A F S A S A A Q I S D G K Q Y I - T L D K P - V |
| sw_DSBA_ENTAM/3 | A K W I N S I F K S V V L T A A L A L P F T A S - A F T E G T D Y M - V L E K P - |
| sw_DSBA_LEGPN/1 | -----L M P M T A L A T Q F I E G K D Y Q T V A S A Q - L S |

cons

| | |
|---|---|
| sw_DSBA_PSES/1 sw_DSBA_SALTY/1 sw_DSBA_ENTAM/3 sw_DSBA_LEGPN/1 | AVPGK-IEVI E LFWYGC P HCYAF E PTI---NPWVEKLPSDVNFVR --AGE-PQVLE FFS F S YCPHCYQFEEVLHVS D NVKKKLPEGTKMTK -IPDADKT L IKVFSYACPFCYKYDKAVT--GPVADKVADLVTFVP TNKD K TPL I TE FFS SYGCPWCYKIDAPLN-D-WATR M GKAHLER |
|---|---|

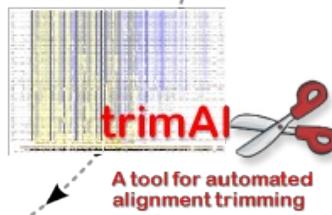
cons

https://www.semanticscience.org

Multimle Sequence Alignments

Alignment reconstruction.

Alignment trimming.



A tool for automated alignment trimming

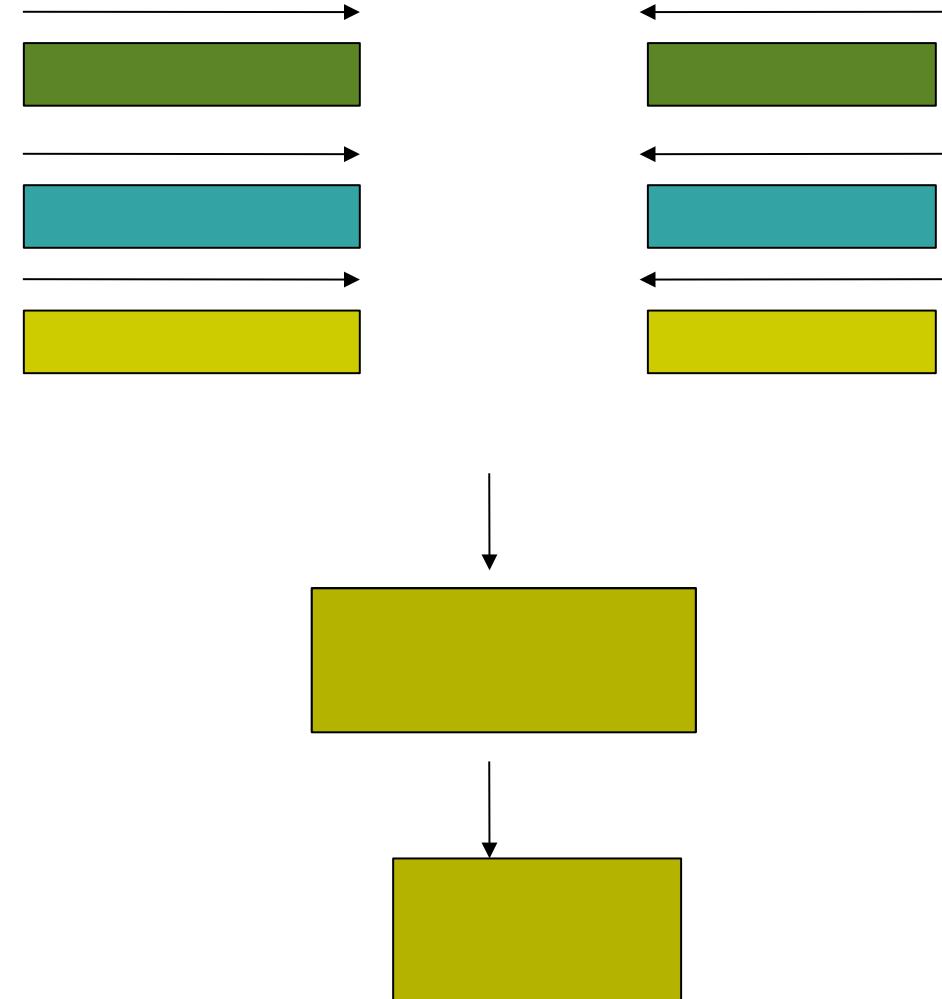
e approach.

The set of homologous Sequences are aligned by 3 different aligners in forward and reverse modes (Head or Tails approach)

A consensus is buit

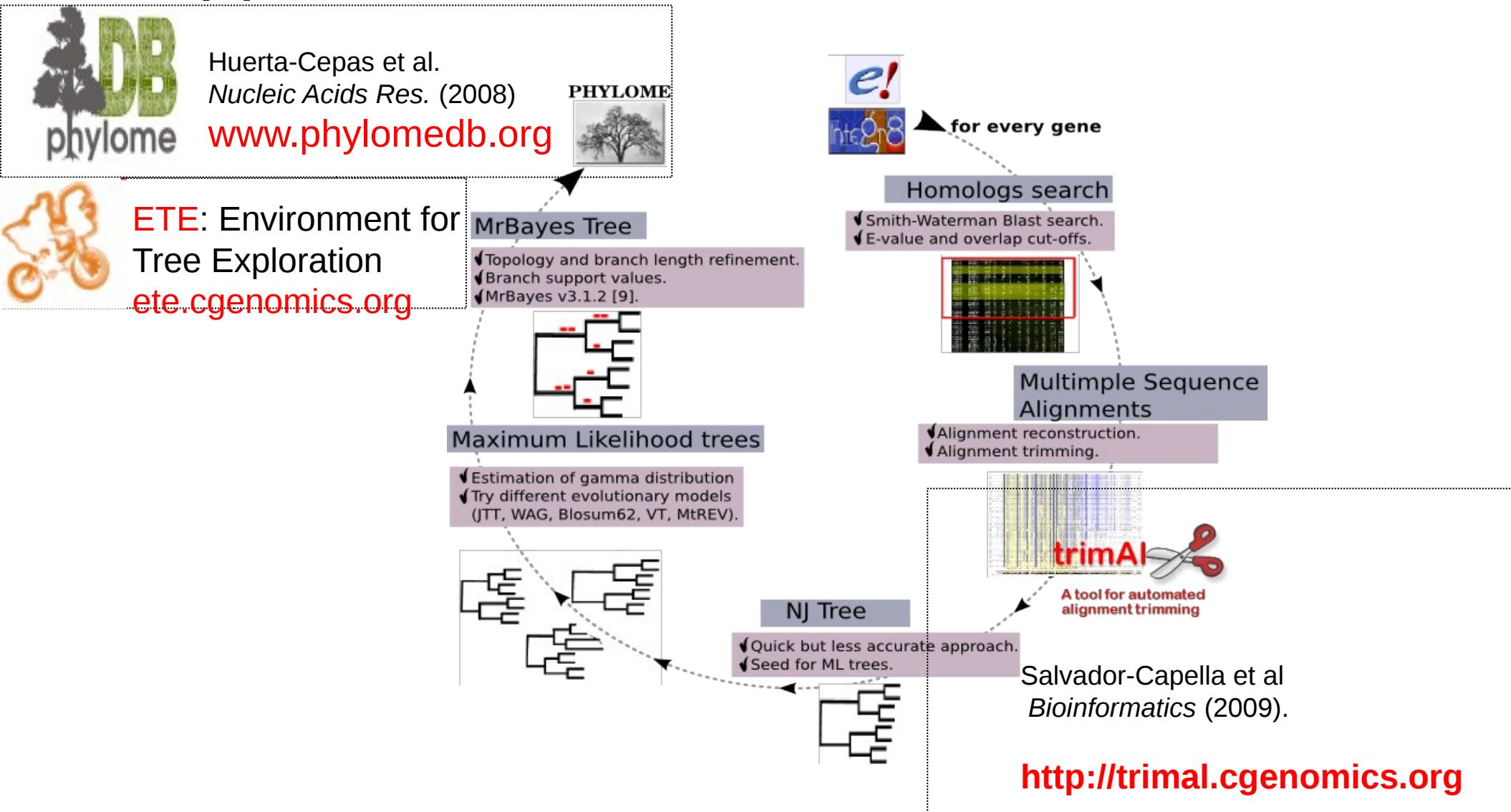
The consensus is trimmed (trimAl) based on:

- consistency across the 6 alignments
- gap content



[Http://trimal.cgenomics.com](http://trimal.cgenomics.com)

Our pipeline:



Pipeline described in Huerta-Cepas et al NAR (2011)



[Login]

Search in PhylomeDB

(i.e. ENSG00000139618, YBL058W,
TP53)

[RandomTree!](#)

[BLAST search](#)

Latest Phylomes

| | |
|-----------------------|------|
| Arxula adeninivorans | 2014 |
| Beta vulgaris | 2013 |
| Clogmia albipunctata | 2013 |
| Penicillium digitatum | 2012 |
| Schistosoma mansoni | 2012 |

[see all phylomes](#)

PhylomeDB uses



PhylomeDB cross linking



Latest story

New Zygomycete phylome: the human pathogen *Lichtheimia corymbifera*

Mon, 09/15/2014 - 21:09

PhylomeDB extends its repertoire of fungal phylomes with that of a genome of a poorly sample clade, that of the basal group zygomycetes. In this case the phylome (245) of the human pathogenic mucorales *Lichtheimia corymbifera* has served to reveal extensive past gene duplications in this group. *Lichtheimia* species are the second most important cause of mucormycosis in Europe. The sequencing of its genome and the comparison with other Zygomycete species, particularly of *Rhizopus delemar*, the main

Welcome to PhylomeDB 4!

PhylomeDB is a public database for complete **catalogs of gene phylogenies** (phylomes). It allows users to interactively explore the evolutionary history of genes through the visualization of phylogenetic trees and multiple sequence alignments. Moreover, phylomeDB provides genome-wide orthology and paralogy predictions which are based on the analysis of the phylogenetic trees. The automated pipeline used to reconstruct trees aims at providing a high-quality phylogenetic analysis of different genomes, including Maximum Likelihood tree inference, **alignment trimming** and evolutionary model testing.

PhylomeDB includes also a public download section with the complete set of trees, alignments and orthology predictions, as well as a **web API** that facilitates cross linking trees from external sources. Finally, phylomeDB provides an advanced tree visualization interface based on the **ETE toolkit**, which integrates tree topologies, taxonomic information, domain mapping and alignment visualization in a single and interactive tree image.

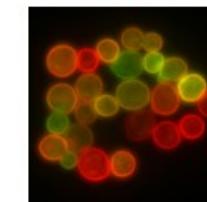
[What's new in phylomeDB 4?](#)

Popular Phylome Collections

Human



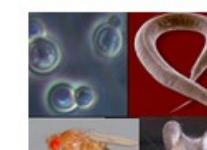
Fungi



Plants



Model Species



Latest News

New Zygomycete phylome: the human pathogen *Lichtheimia corymbifera*

Mon, 09/15/2014 - 21:09

New hemiascomycete phylome: *Blastobotrys (Axula) adeninivorans*, a yeast of biotechnological interest.

Mon, 05/19/2014 - 11:01

Help us to improve phylomeDB: complete our survey.

Thu, 02/20/2014 - 16:22

[show all](#)

PhylomeDB Twitter

Tweets

[Follow](#)

phylomedb @phylomedb

23 Oct

New birds, crocs, and fungal phylomes to come soon at phylomeDB. stay tuned!

[Expand](#)

phylomedb @phylomedb

15 Sep

New Zygomycete phylome: the human fungal pathogen *Lichtheimia corymbifera* phylomedb.org/?q=node/537

[Expand](#)

phylomedb @phylomedb

26 Aug

NOTICE: PhylomeDB will be down due to MAINTENANCE



[Login] Home

Collections All phylomes Downloads Help FAQ About

Search in PhylomeDB

(i.e. ENSG00000139618, YBL058W,

TP53)

Search

RandomTree!

BLAST search

Latest Phylomes

| | |
|------------------------------|------|
| Clogmia albipunctata | 2013 |
| Penicillium digitatum | 2012 |
| Schistosoma mansoni | 2012 |
| Cucumis melo | 2012 |

see all phylomes

PhylomeDB uses



TP53 tree in phylome 218

AS seed in Rat phylome

JTT (Ik:-18130.4)

-- in collateral trees --

Tree features

Search

Clear search

Image

Hard link

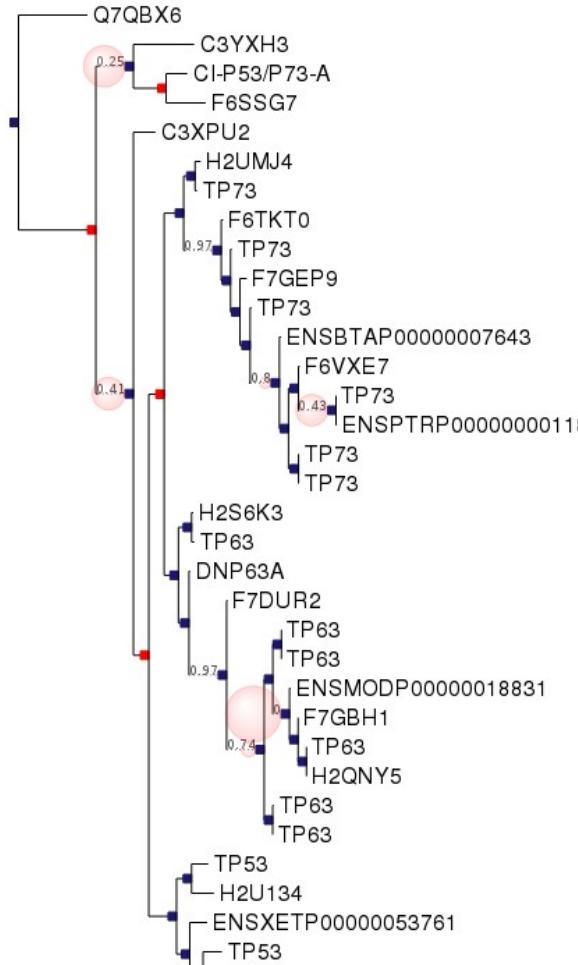
Download OrthoXML

See alignments

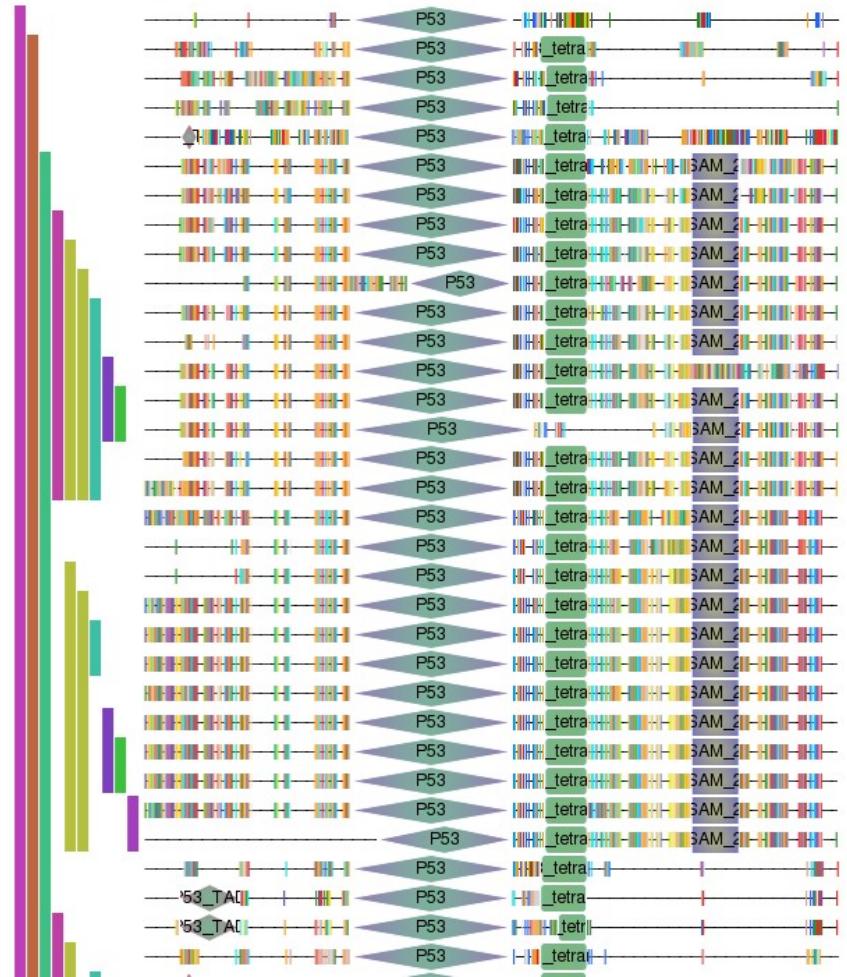
|

|

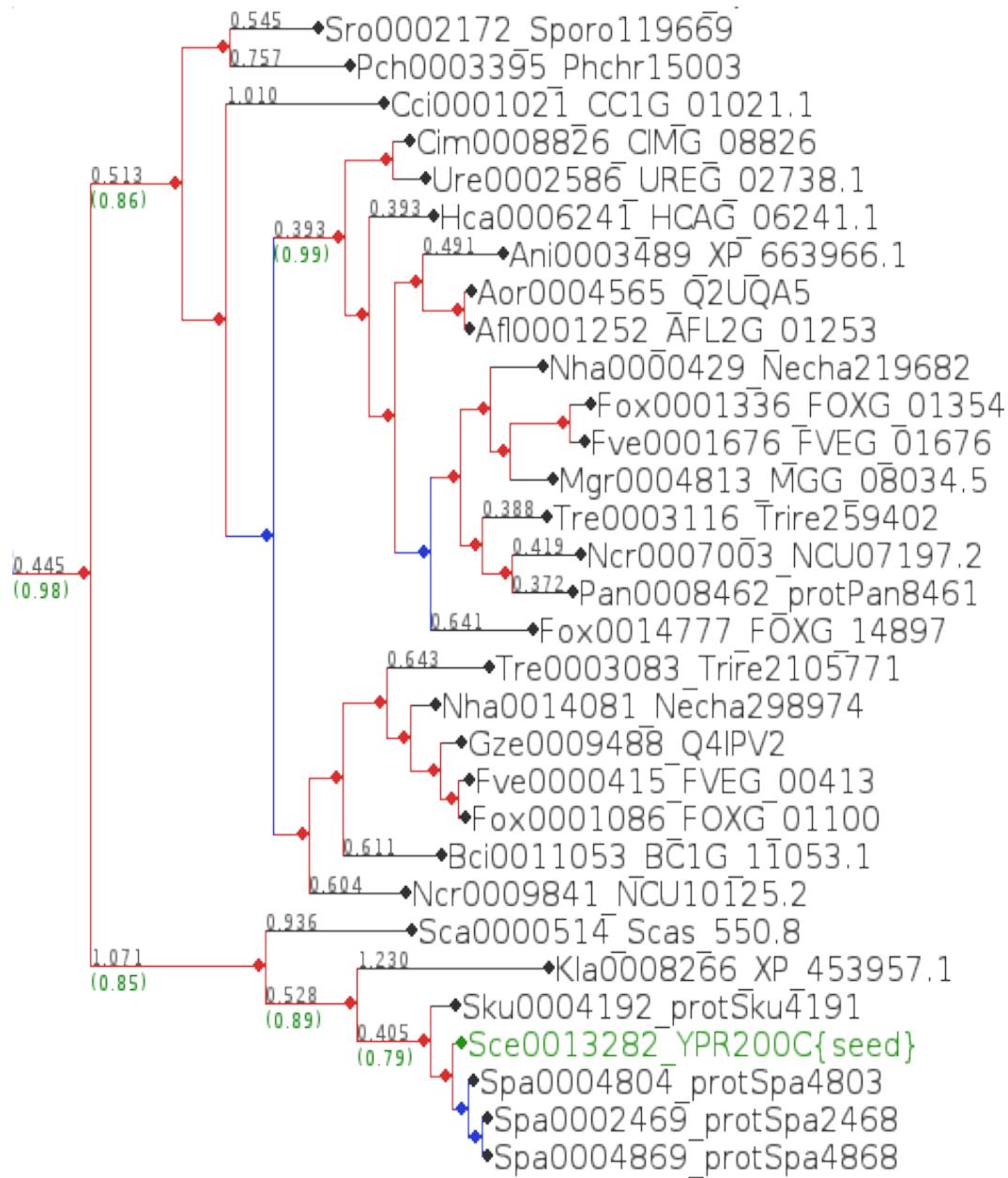
Download data.tar.gz



Anopheles gambiae
Branchiostoma floridae
Ciona intestinalis
Ciona intestinalis
Branchiostoma floridae
Takifugu rubripes
Danio rerio
Xenopus tropicalis
Gallus gallus
Monodelphis domestica
Canis familiaris
Bos taurus
Macaca mulatta
Homo sapiens
Pan troglodytes
Mus musculus
Rattus norvegicus
Takifugu rubripes
Danio rerio
Gallus gallus
Ornithorhynchus anatinus
Rattus norvegicus
Mus musculus
Monodelphis domestica
Macaca mulatta
Homo sapiens
Pan troglodytes
Canis familiaris
Bos taurus
Danio rerio
Takifugu rubripes
Xenopus tropicalis
Gallus gallus



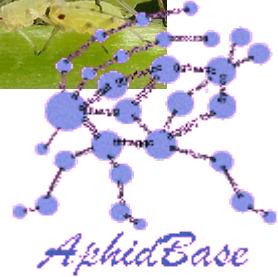
These phylomes can now be interrogated in many ways



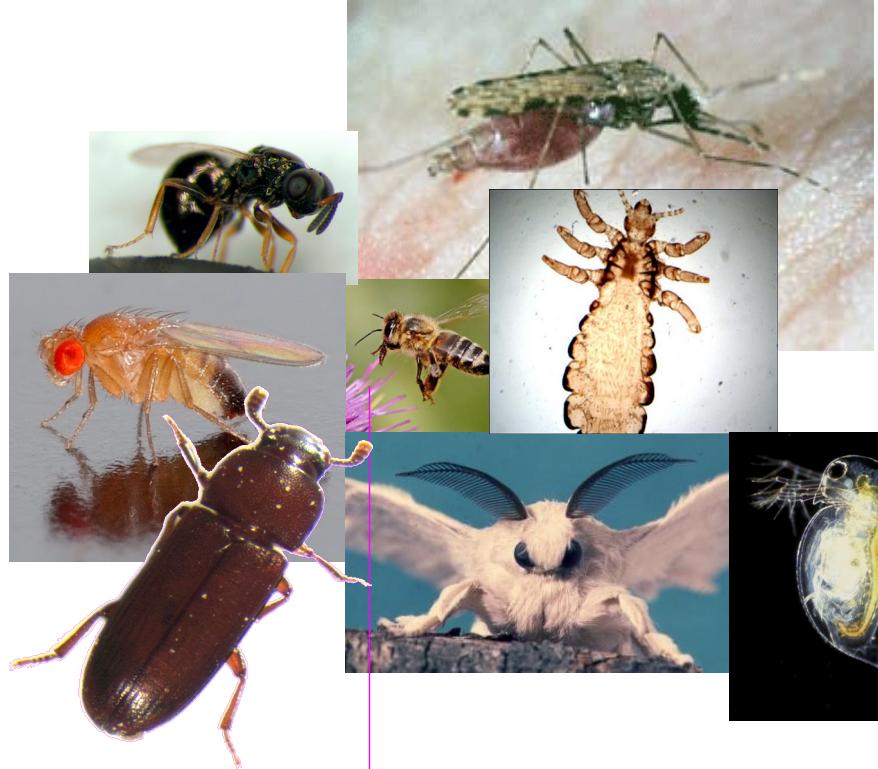
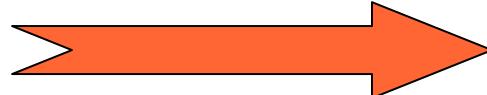
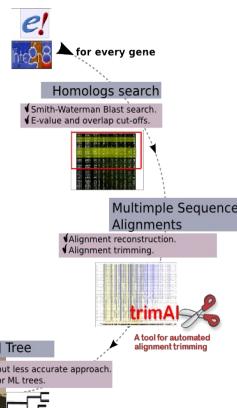
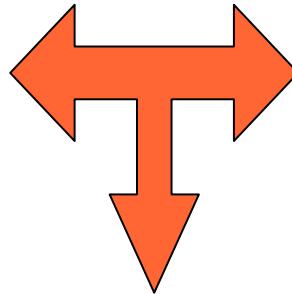
- Families that show a particular topology
- Detect and date duplication events
- Genes that have accelerated evolutionary rates at a particular lineage (positive/relaxed selection)
- Families expanded at particular lineages
- Footprints of horizontal gene transfer, lineage sorting, gene conversion and other evolutionary processes
- Search for co-evolving genes
- predict functional properties
- across-species prediction of orthology and paralogy

Large-scale phylogenetics to assist in the annotation of
newly sequenced genomes:
the *Acyrthosiphon pisum* genome .

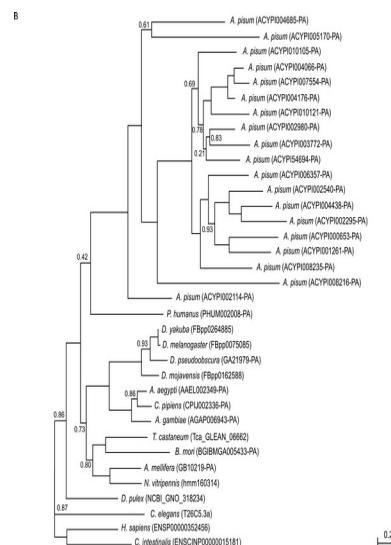




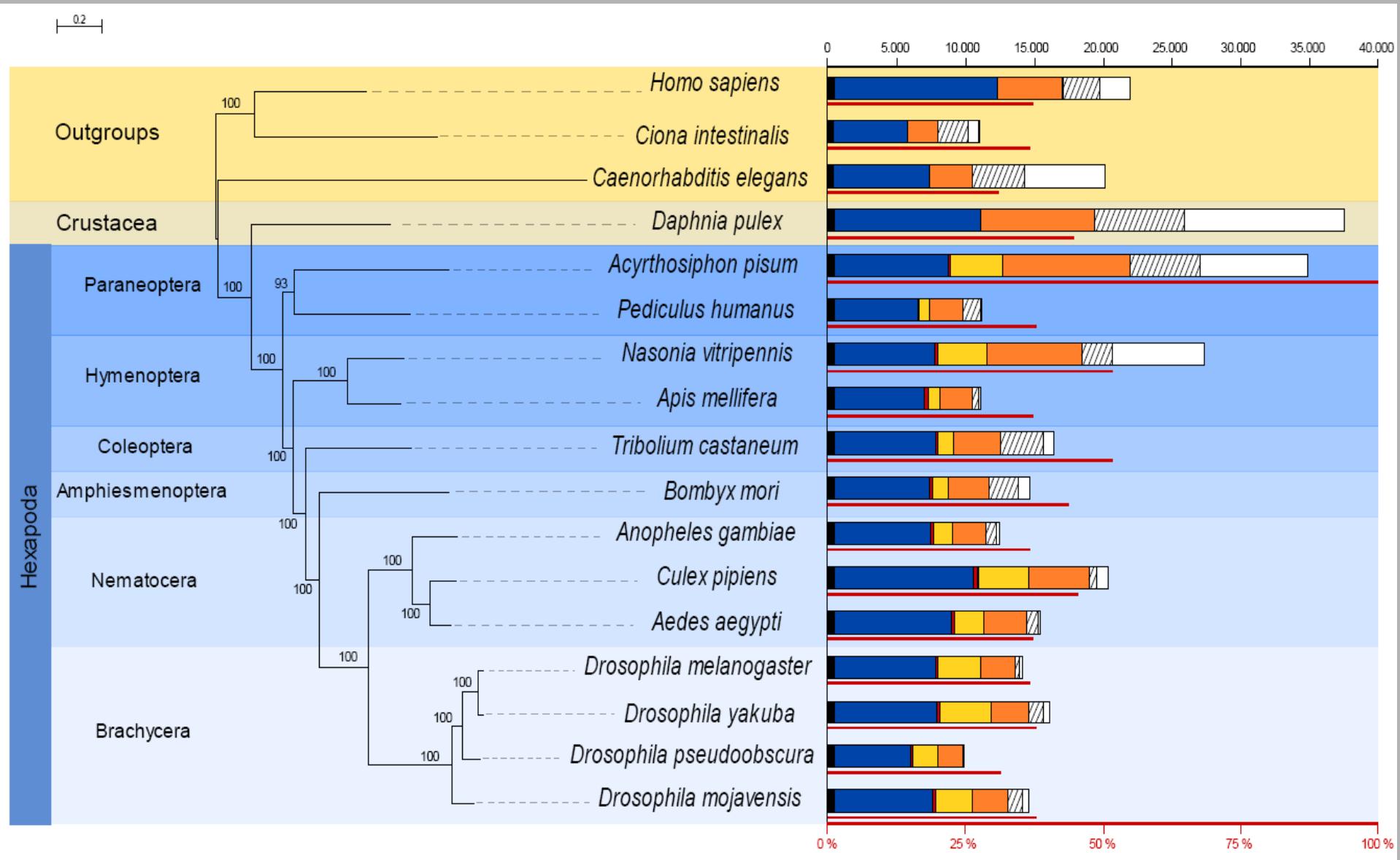
Acyr_1.0 (34,600 genes)



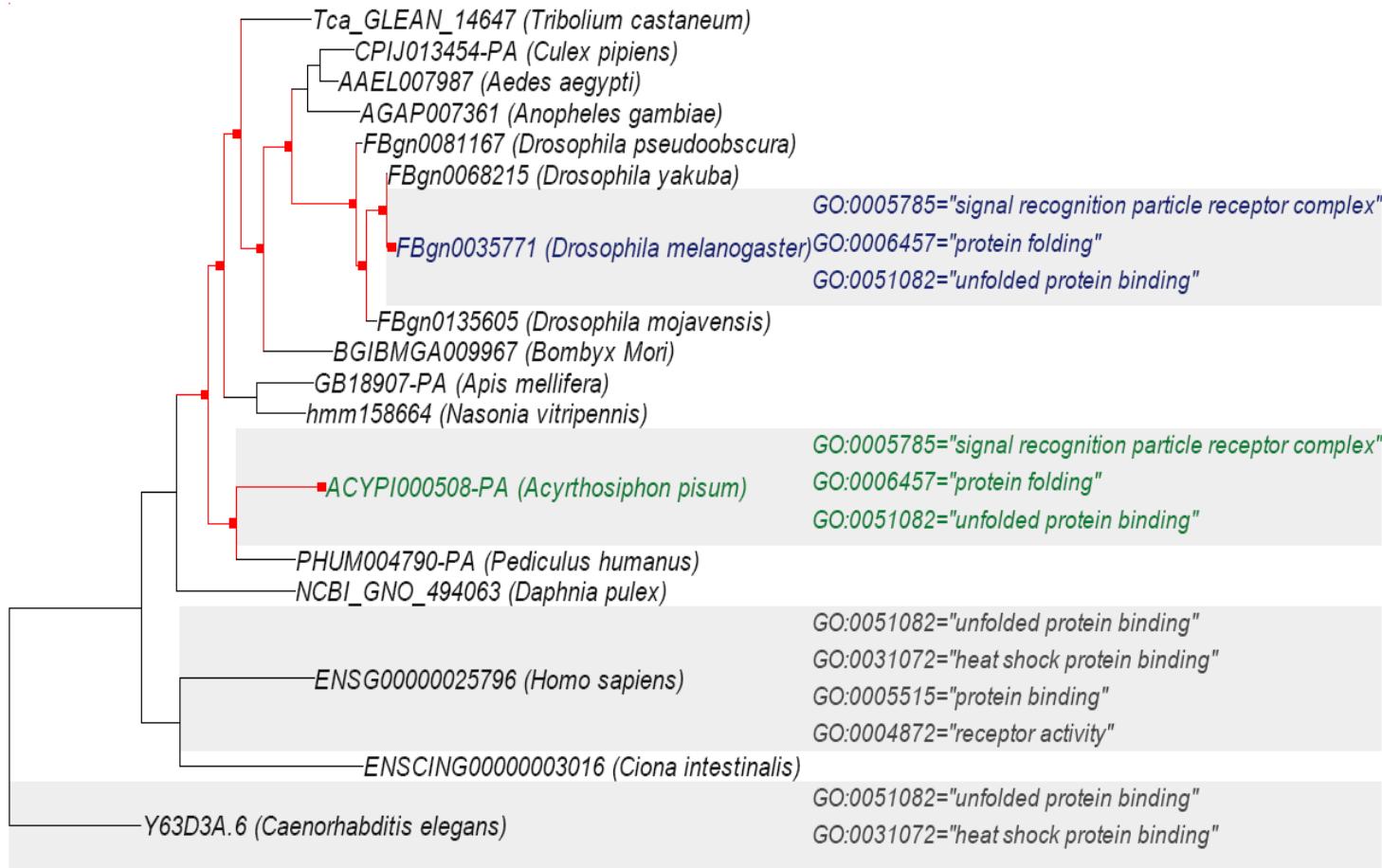
13 other sequenced arthropods and 3 out-groups



0.2



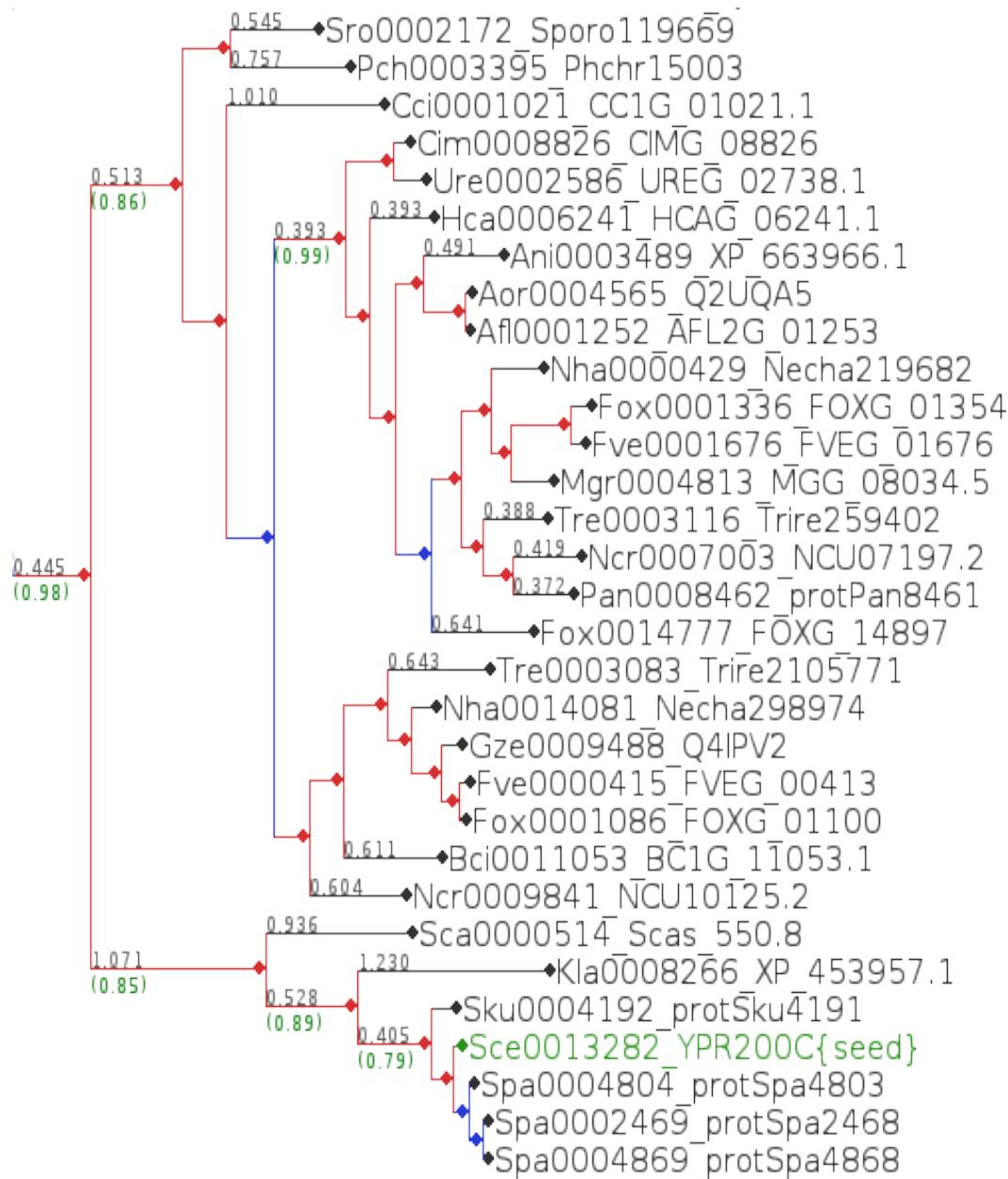
Tree based on the alignment concatenation of 197, widespread proteins



0.46

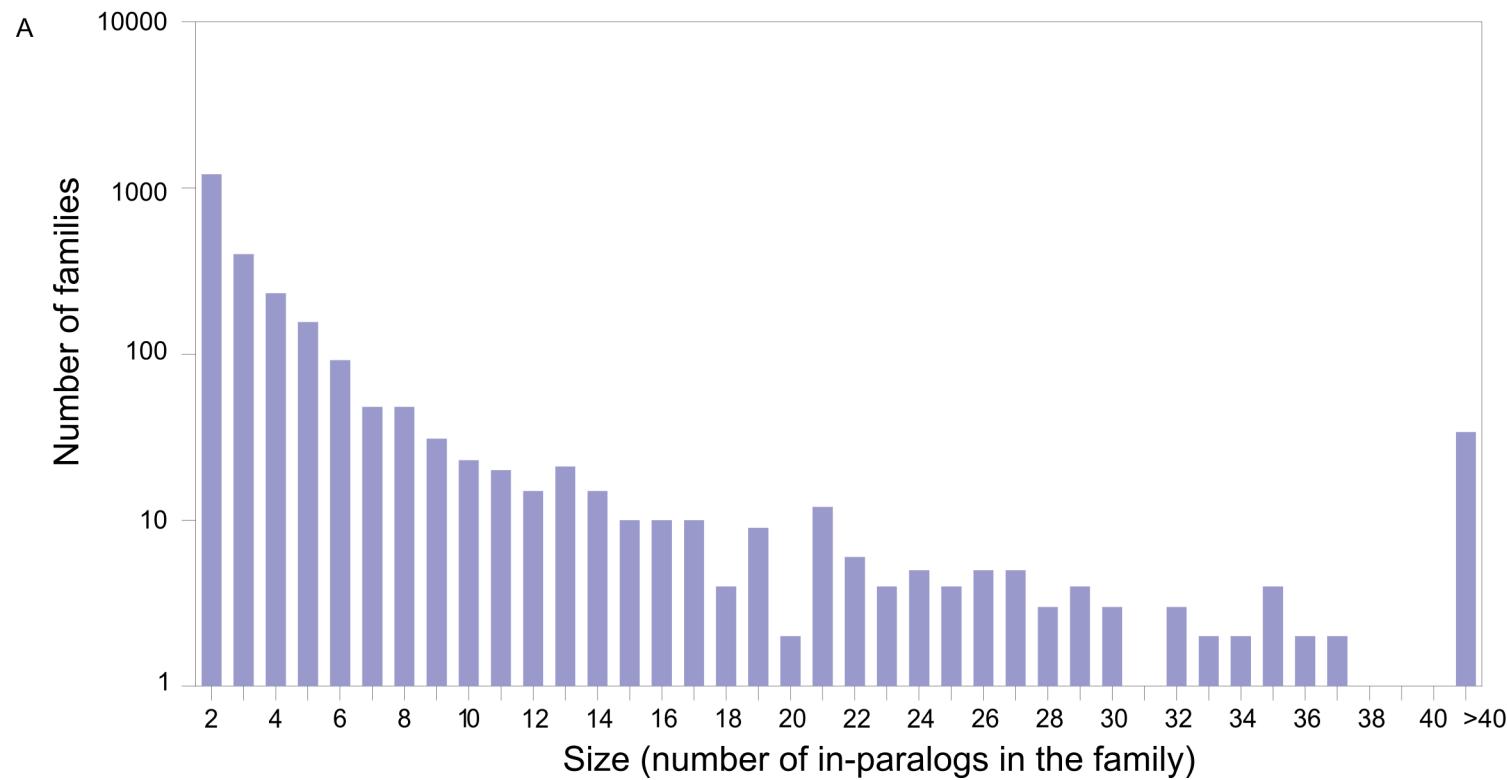
Orthologies with annotated *Drosophila melanogaster* genes:
4,059 (one-to-one), **2,282** (one-to-many, many-to-many or many-to-one)

These phylomes can now be interrogated in many ways



- Families that show a particular topology
- Detect and date duplication events
- Genes that have accelerated evolutionary rates at a particular lineage (positive/relaxed selection)
- **Families expanded at particular lineages**
- Footprints of horizontal gene transfer, lineage sorting, gene conversion and other evolutionary processes
- Search for co-evolving genes
- predict functional properties
- across-species prediction of orthology and paralogy

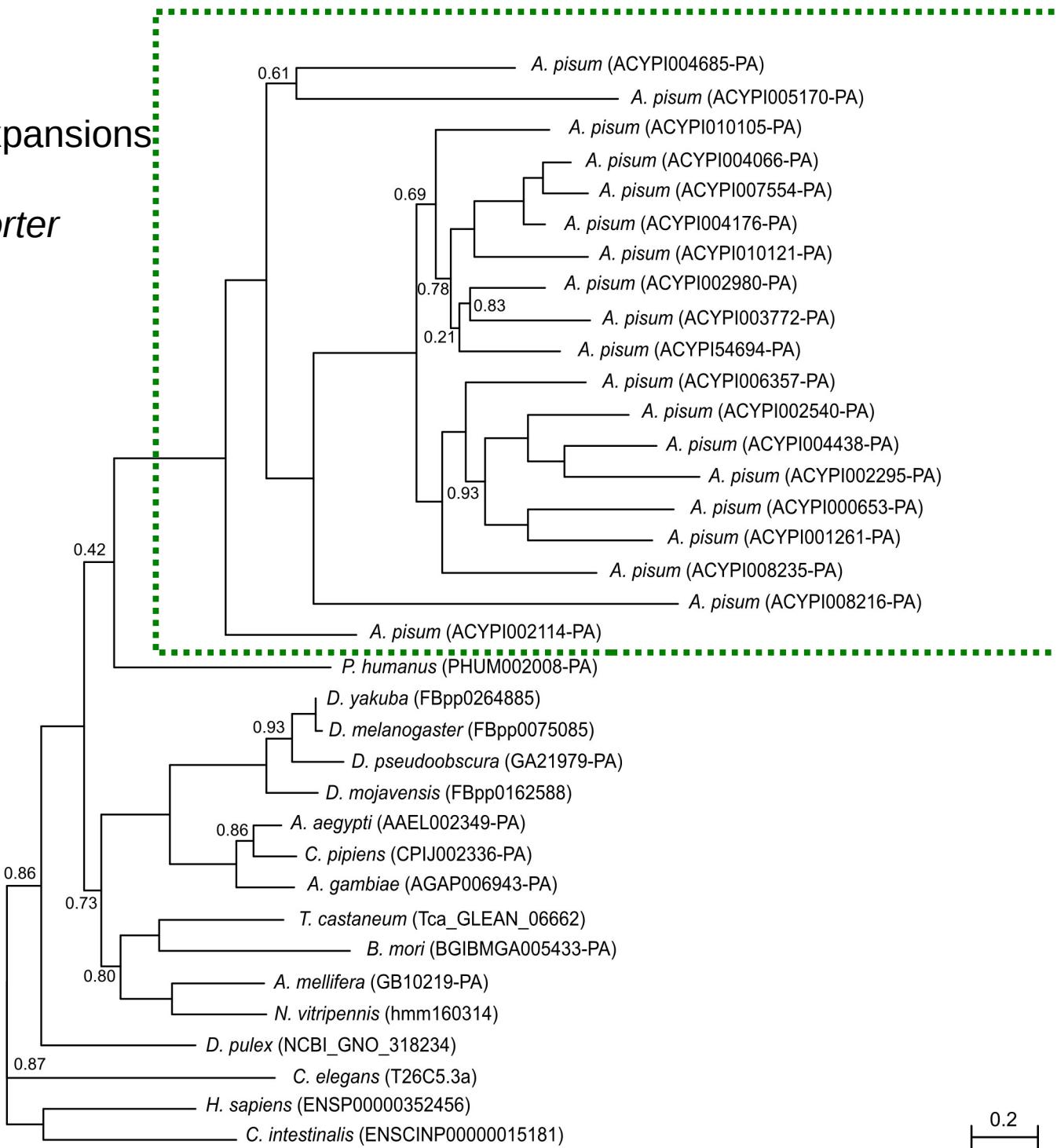
A wave of lineage-specific expansions in the pea aphid

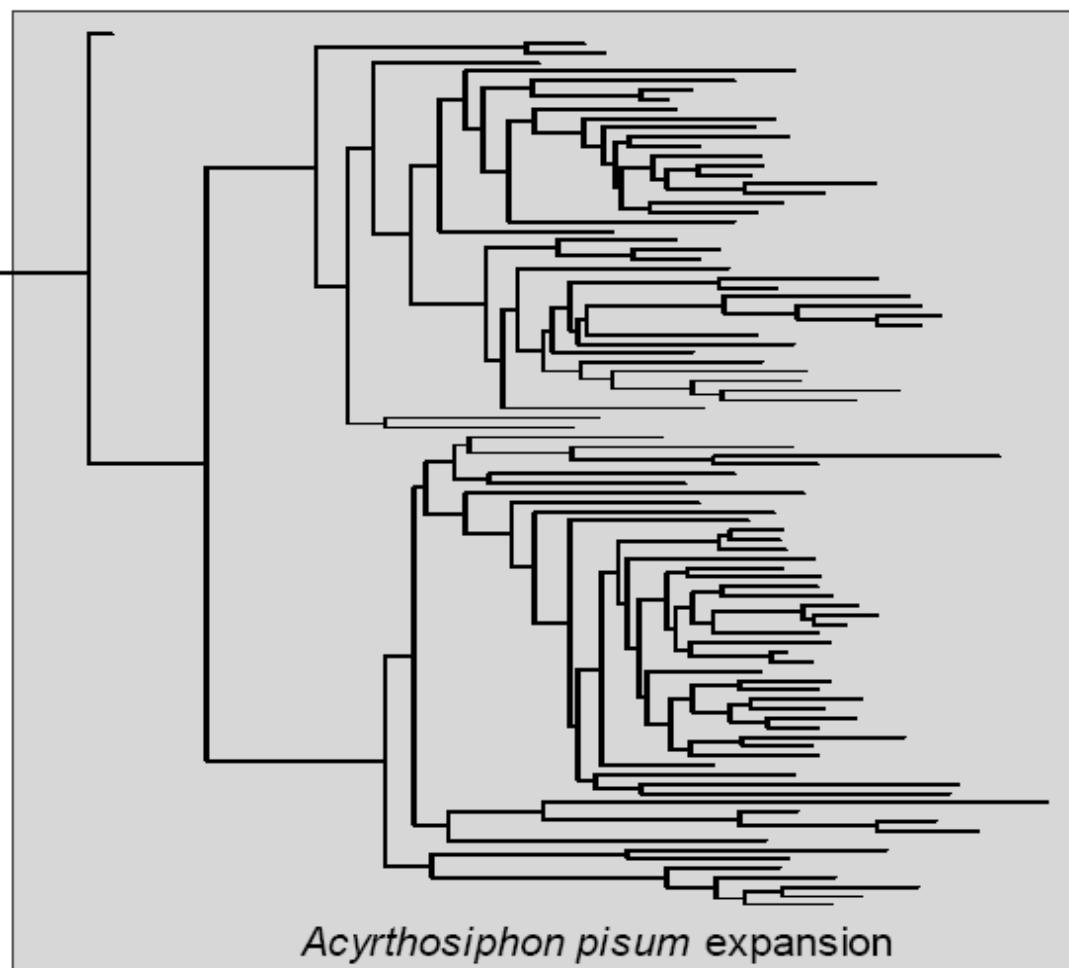
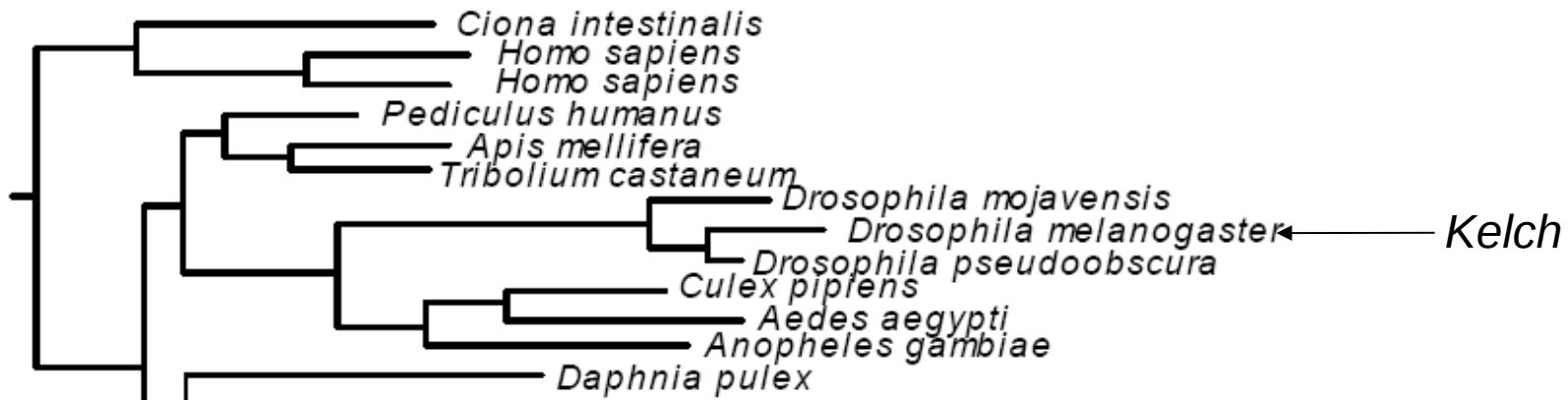


B

Lineage-specific expansions

Acetyl-CoA transporter





In *Drosophila*, *kelch* protein is involved in the organization and morphology of the ovarian ring channel.

A particularity of pea aphids is a complex life cycle with reproductive polyphenism and extensive differences in ovarian morphology between the different female morphs.

Is the *kelch* family expansion in aphids related to such diversity?

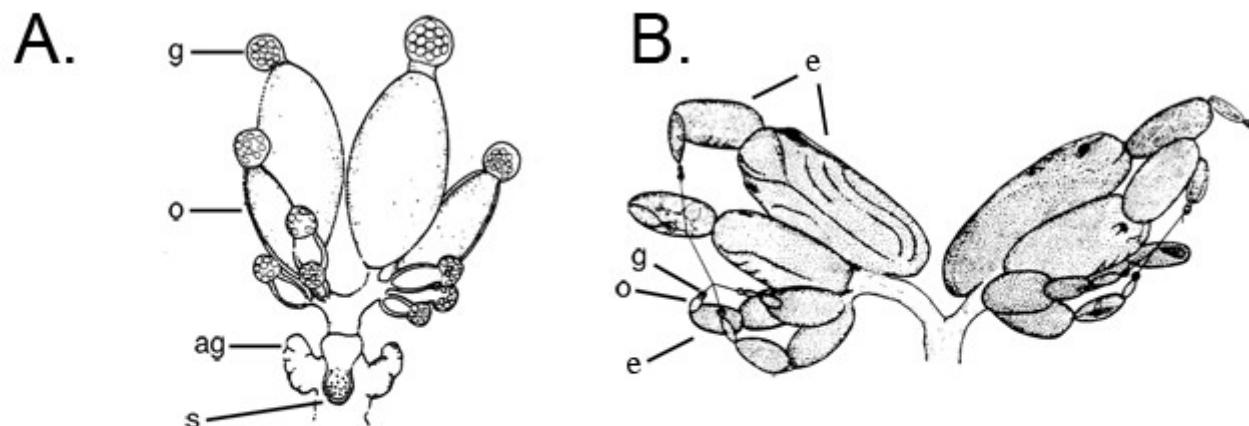
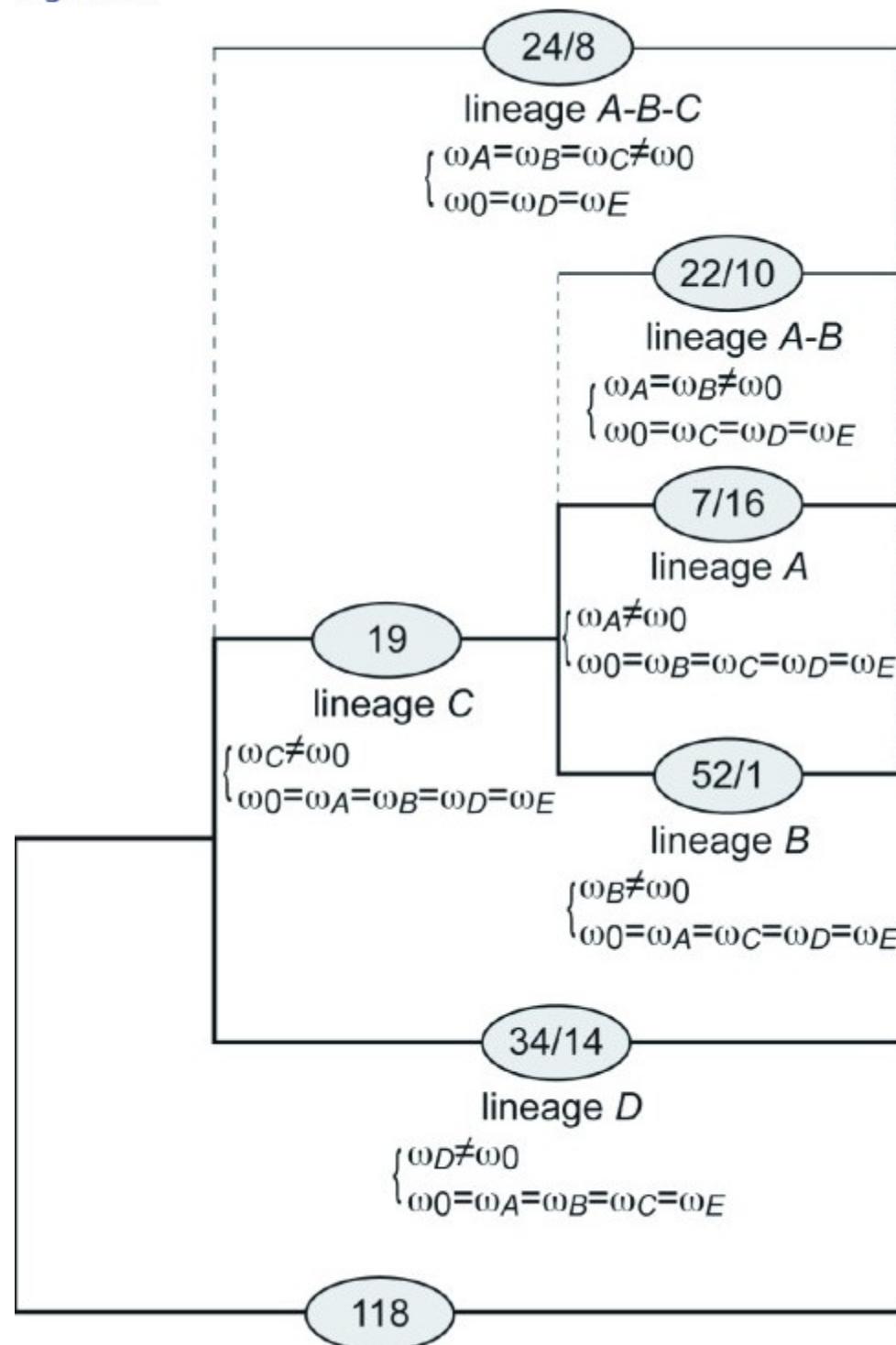


Figure 2. Viviparous and oviparous development. Oviparous (A) and viviparous (B) ovaries differ not only as to whether they possess embryos, accessory glands and spermathecae, but also in the relative size of germaria and oocytes. Abbreviations: g is germarium, o is oocyte, e is viviparous embryo, ag is accessory gland, s is spermatheca. Images are modified from Blackman, 1987.

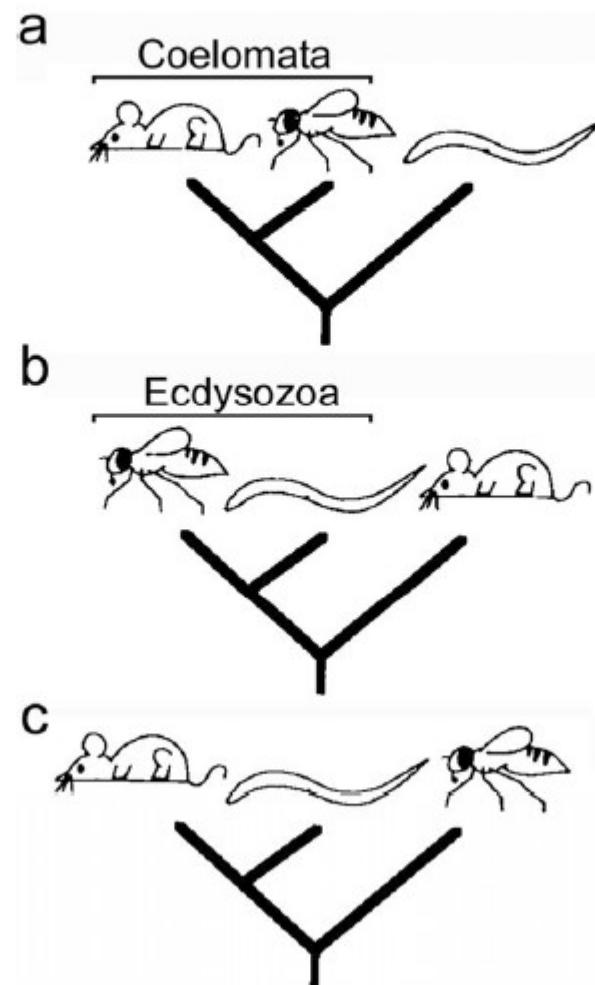
Figure 2.



Going beyond topologies

Gene tree vs species trees

Uncertainty in species trees and topological variability in gene trees



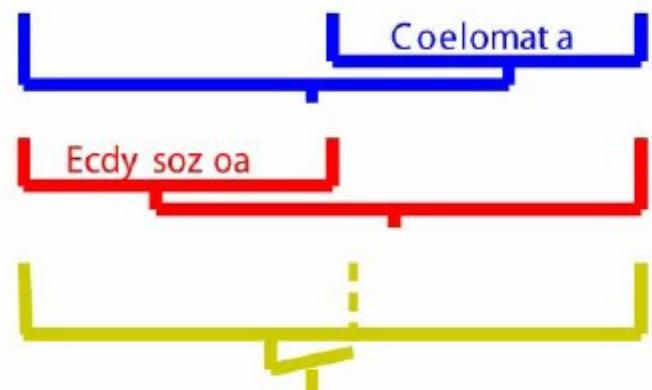
Nematodes



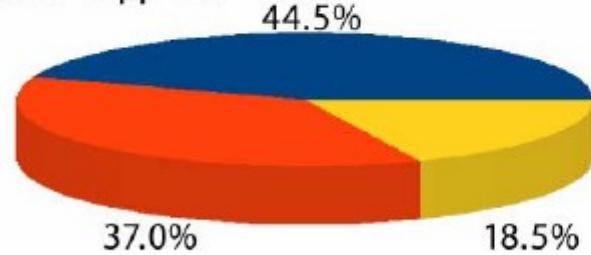
Arthropods



Chordates



Phylome supp ort:



What percentage of gene trees from the human phylome support each topology?

Similar results for

Primates

Rodents

laurasatheria

Huerta-Cepas et. al. (2007)

Possible sources of gene tree vs species tree discordance

Analytical factors

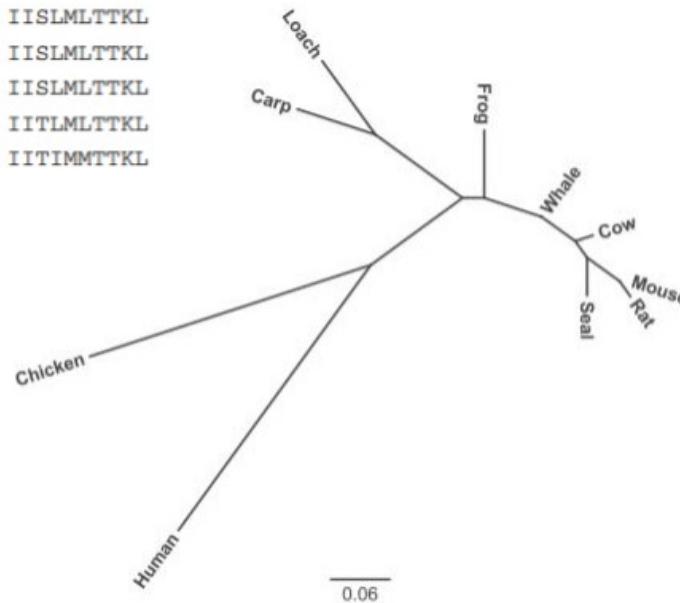
They lead to failure in accurately inferring a gene tree; these can be either due to **stochastic error** (e.g., insufficient sequence length or taxon samples) or due to **systematic error** (e.g., observed data far depart from model assumptions)

Biological factors

They lead to gene trees that are topologically distinct from each other and from the species tree. Known factors include **stochastic lineage sorting, hidden paralogy, horizontal gene transfer, recombination and natural selection**

10 50

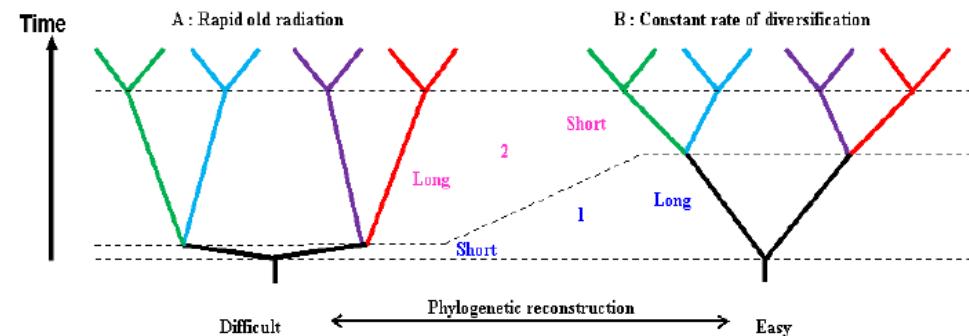
| | |
|---------|--|
| Cow | MAYPMQLGFQ DATSPIMEEL LHFHDHTLMI VFLISSLVLY IISLMLTTKL |
| Carp | MAHPTQLGFQ DAAMPVMEEL LHFHDHALMI VLLISTLVLY IITAMVSTKL |
| Chicken | MANHSQSQLGFQ DASSPIMEEL VEFHDHALMV ALAICSLVLY LLTLMLEKEL |
| Human | MAHAAQVGLQ DATSPIMEEL ITFHDHALMI IFLICFLVLY ALFLTLTTKL |
| Loach | MAHPTQLGFQ DAASPVMEEL LHFHDHALMI VFLISALVLY VIITTVSTKL |
| Mouse | MAYPFQLGLQ DATSPIMEEL MNFHDHTLMI VFLISSLVLY IISLMLTTKL |
| Rat | MAYPFQLGLQ DATSPIMEEL TNFHDHTLMI VFLISSLVLY IISLMLTTKL |
| Seal | MAYPLQMGLQ DATSPIMEEL LHFHDHTLMI VFLISSLVLY IISLMLTTKL |
| Whale | MAYPFQLGFQ DAASPIEEL LHFHDHTLMI VFLISSLVLY IITLMLTTKL |
| Frog | MAHPSQLGFQ DAASPIEEL LHFHDHTLMA VFLISTLVLY IITIMMMTKL |



Alignment can be seen as a sampling
It has stochastic variation, so **sampling “errors”**

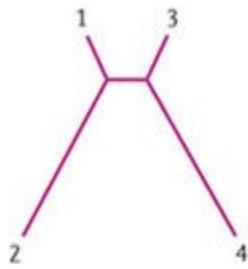
Can be expected, particularly at:

- Short sequences
- Short internodes (i.e. in fast radiations)

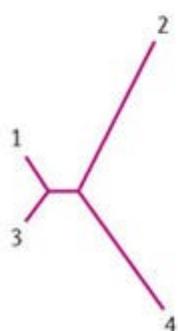


Systematic errors (i.e. long branch attraction artifact)

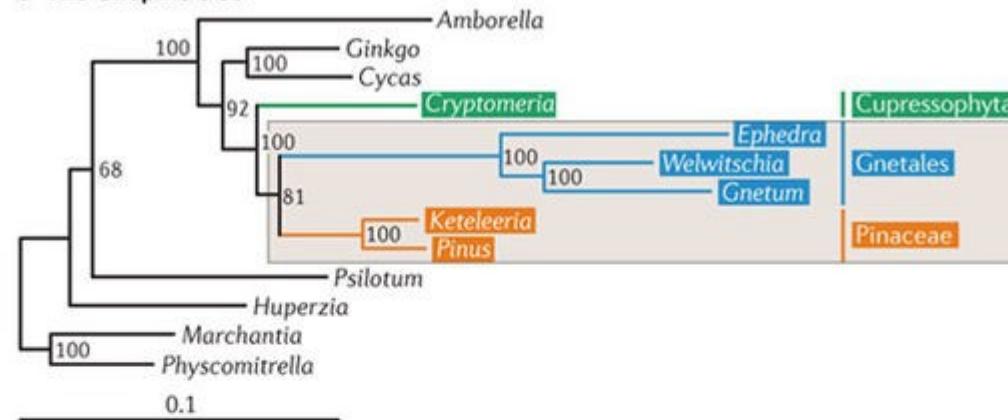
a Correct tree, T_1



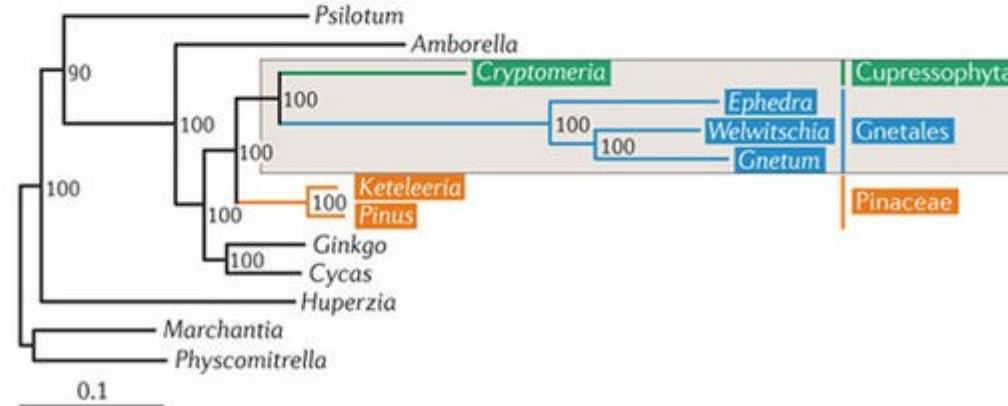
b Wrong tree, T_2



c The Gnepine tree



d The GneCup tree



Possible sources of gene tree vs species tree discordance

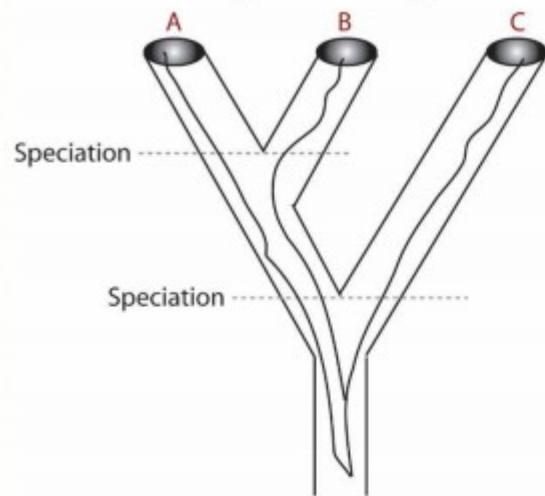
Analytical factors

They lead to failure in accurately inferring a gene tree; these can be either due to **stochastic error** (e.g., insufficient sequence length or taxon samples) or due to **systematic error** (e.g., observed data far depart from model assumptions)

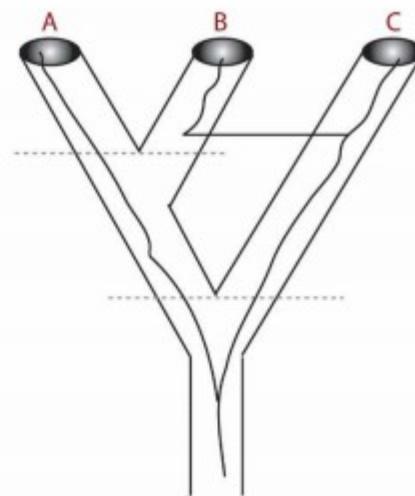
Biological factors

They lead to gene trees that are topologically distinct from each other and from the species tree. Known factors include **stochastic lineage sorting, hidden paralogy, horizontal gene transfer, recombination and natural selection**

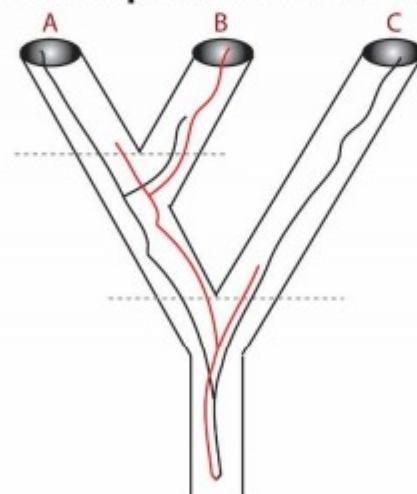
Lineage Sorting



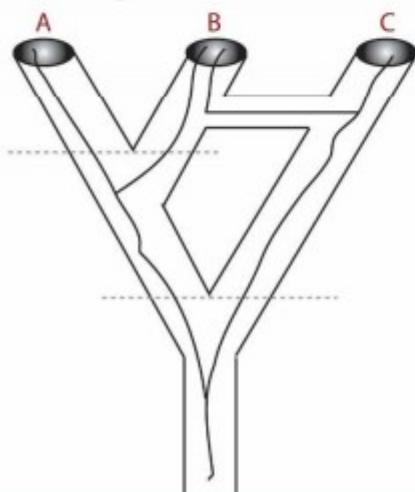
Horizontal Gene Transfer



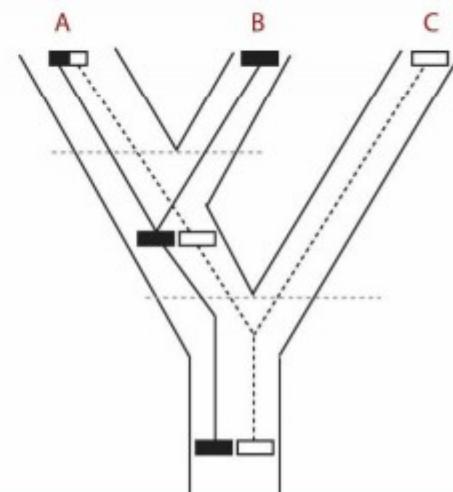
Gene Duplication and Loss



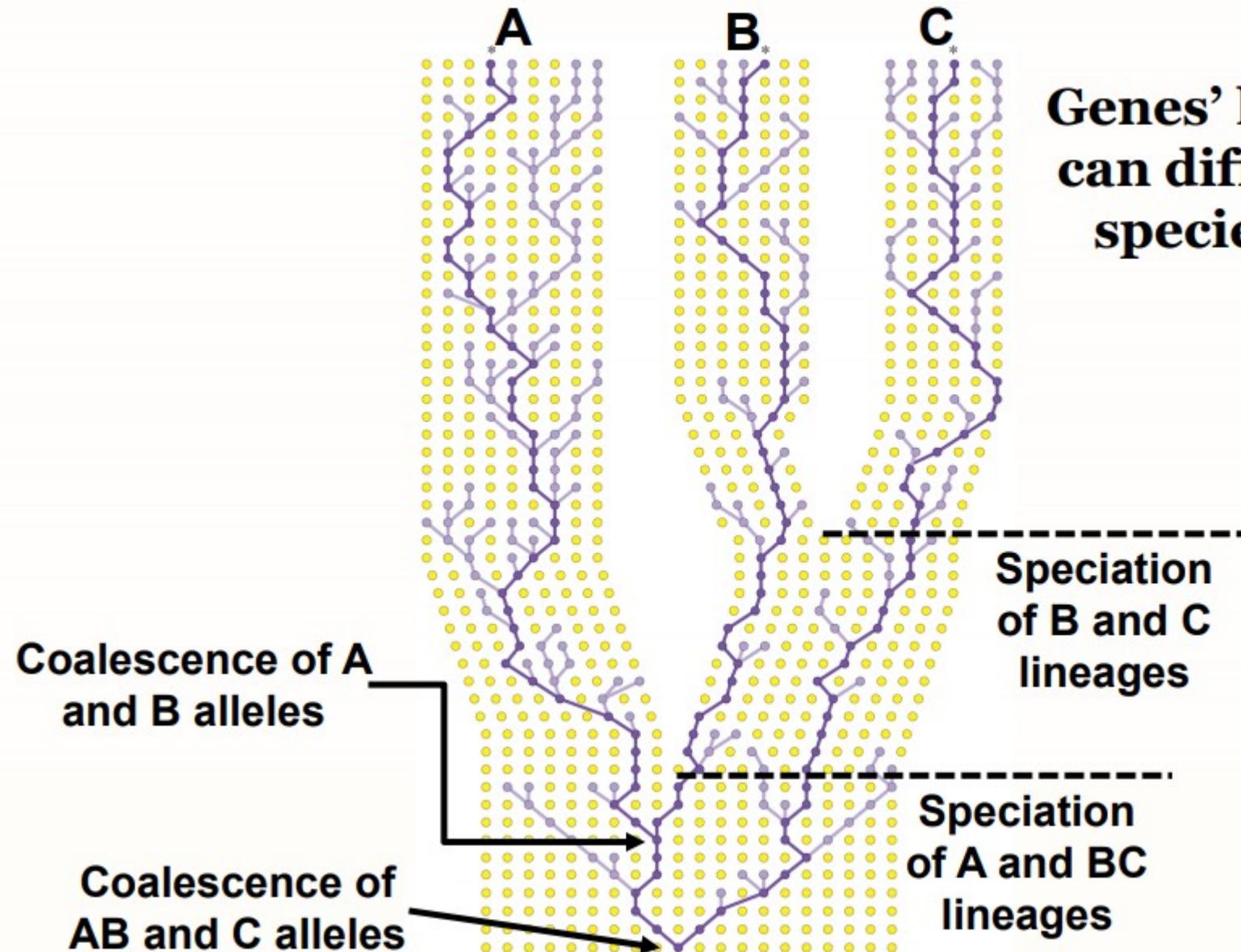
Hybridization



Recombination



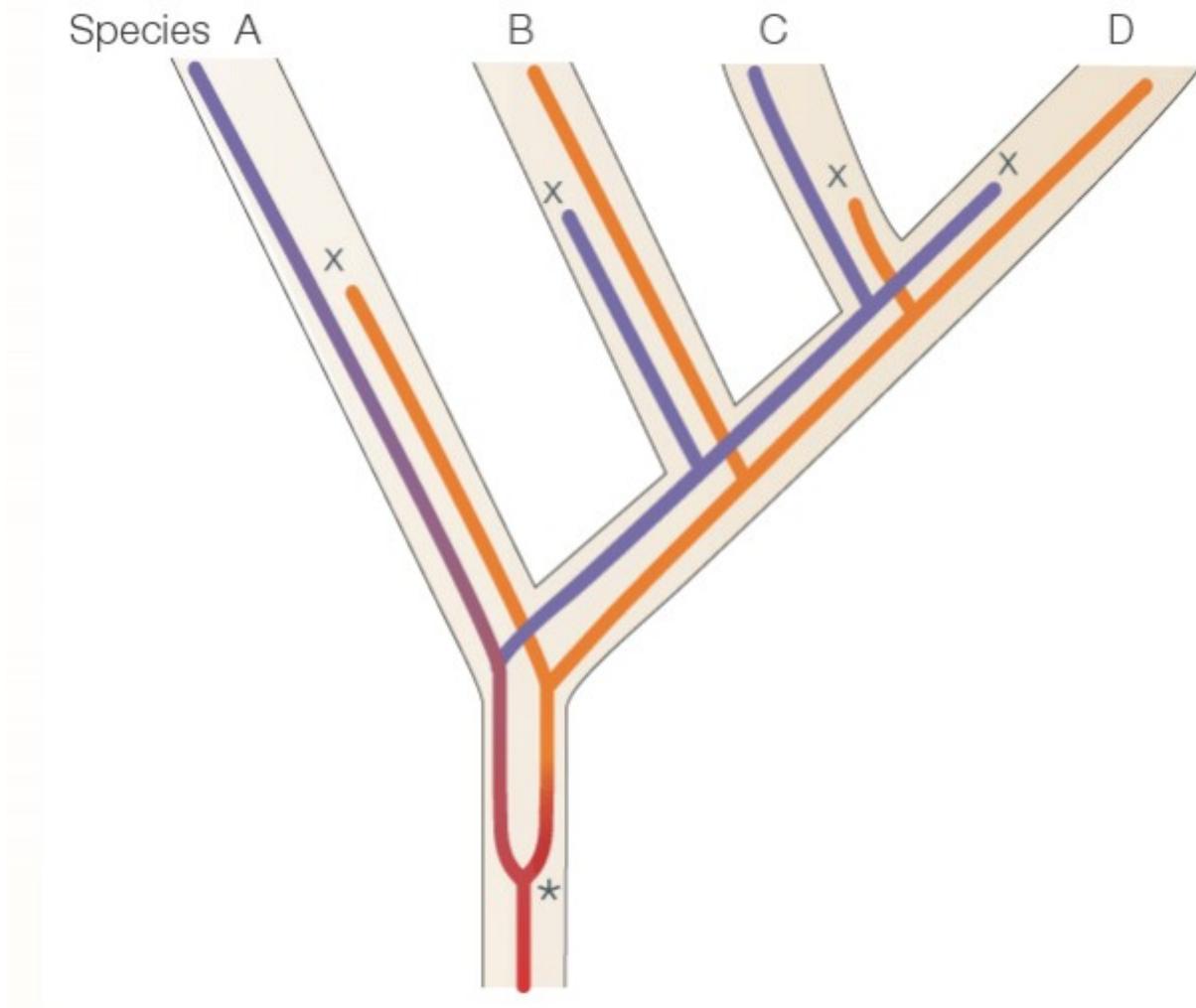
Incomplete lineage sorting



Genes' histories can differ from species ones

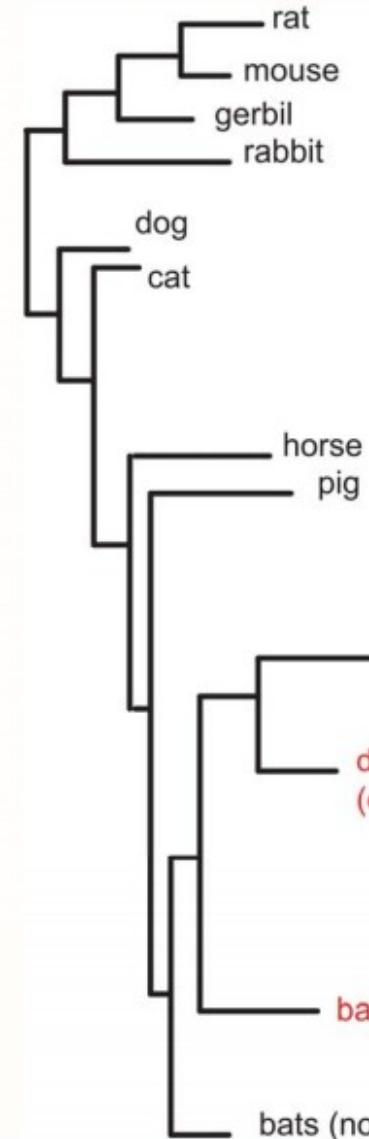
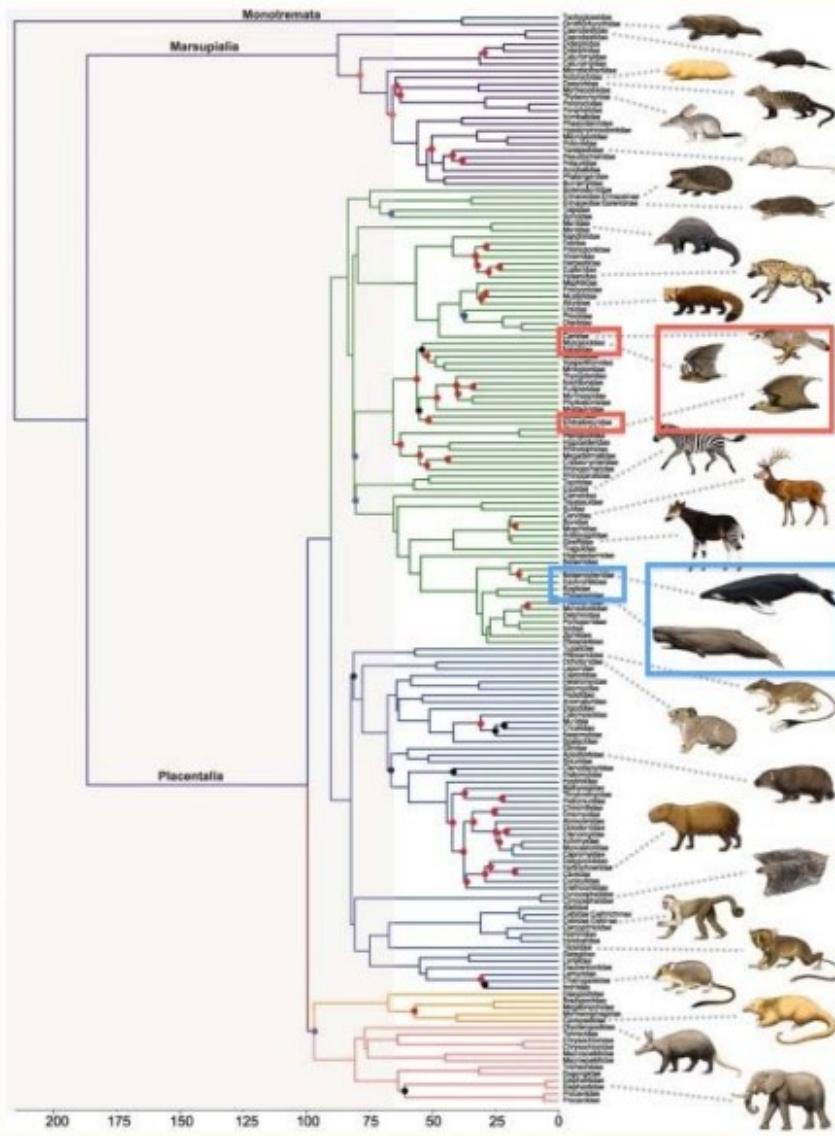
Incomplete lineage sorting in the primate lineage:

| | Informative Sites |
|--|-----------------------------------|
|  | 8,561 / 11,293 (~76%) |
|  | 1,302 / 11,293 (~11.5%) |
|  | 1,430 / 11,293 (~12.5%) |



Hidden paralogy caused by differential gene loss following duplication

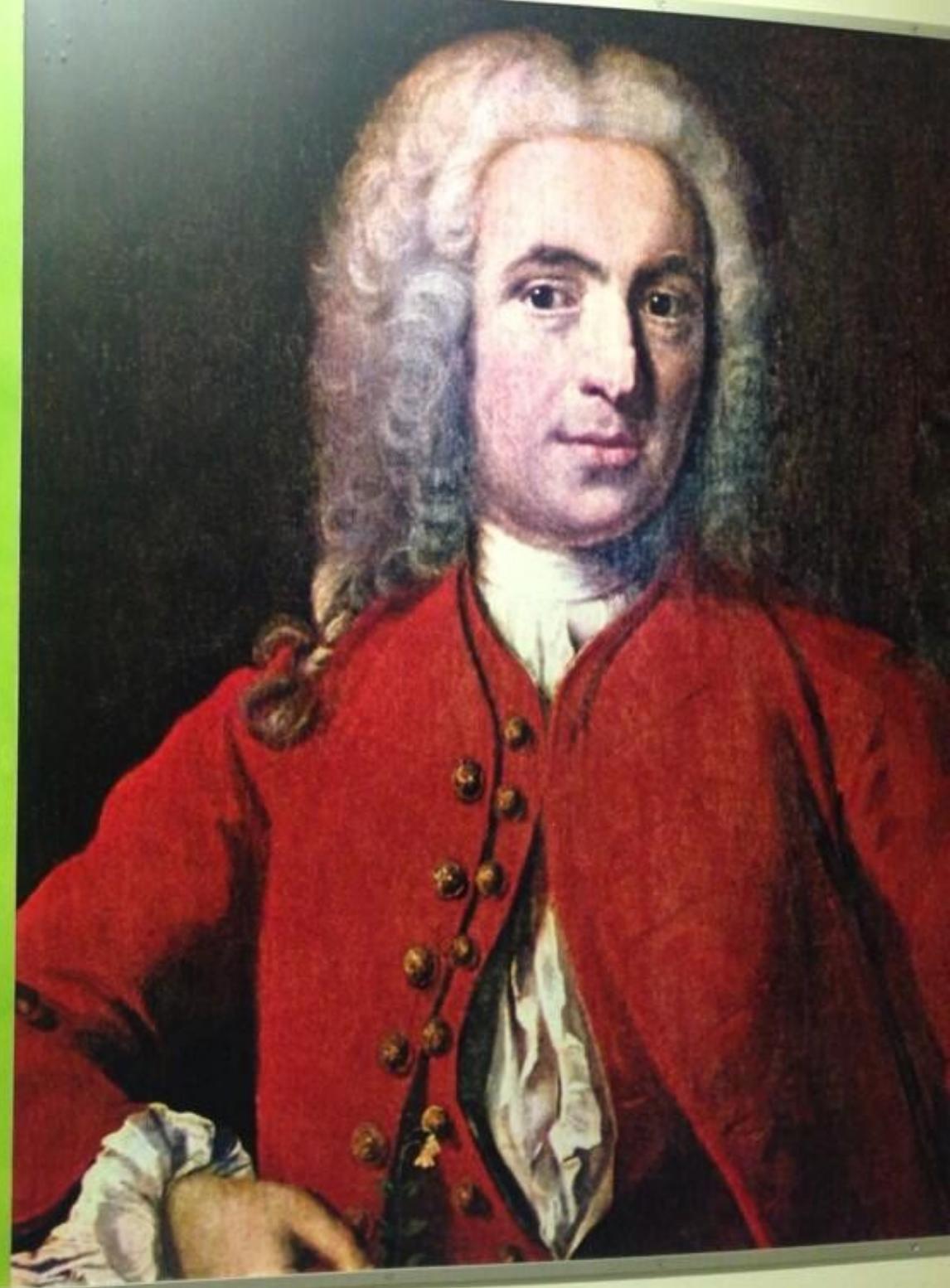
Positive selection can cause tree patterns different from the gene tree by convergent mutations



Phylogeny of *prestin*, a gene involved in echolocation

Non-vertical evolution

Welcome
to my hometown
Carl von Linné, the Father
of Taxonomy

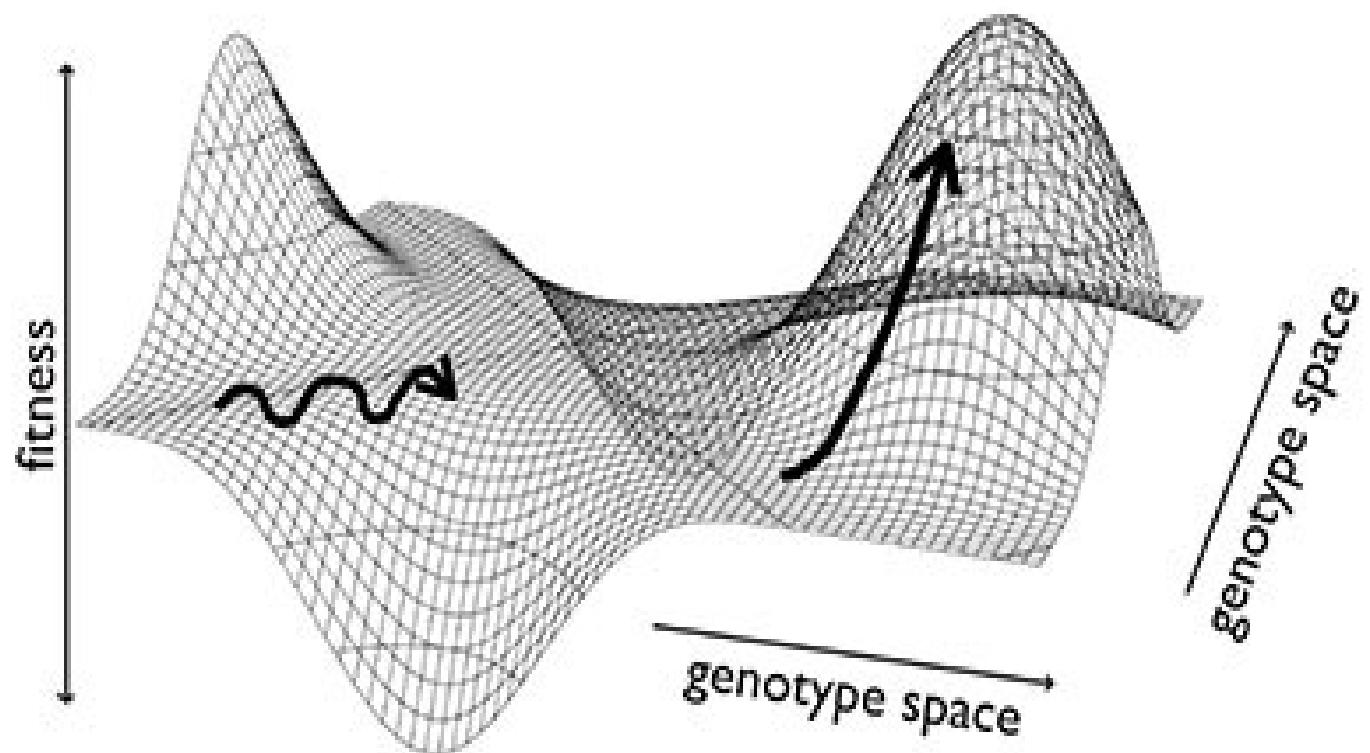




Huxley to Darwin (Nov 1859)

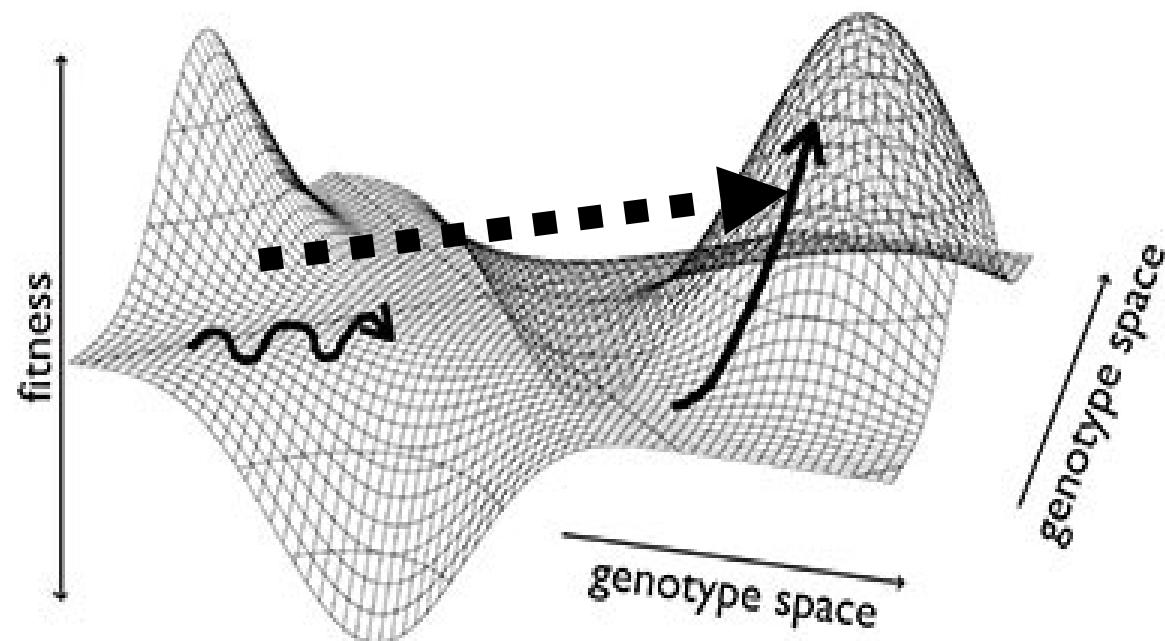
The only objections that have occurred to me are 1st that you have loaded yourself with an unnecessary difficulty in adopting 'Natura non facit saltum' so unreservedly. I believe she does make small jumps

Mutation, selection and drift:

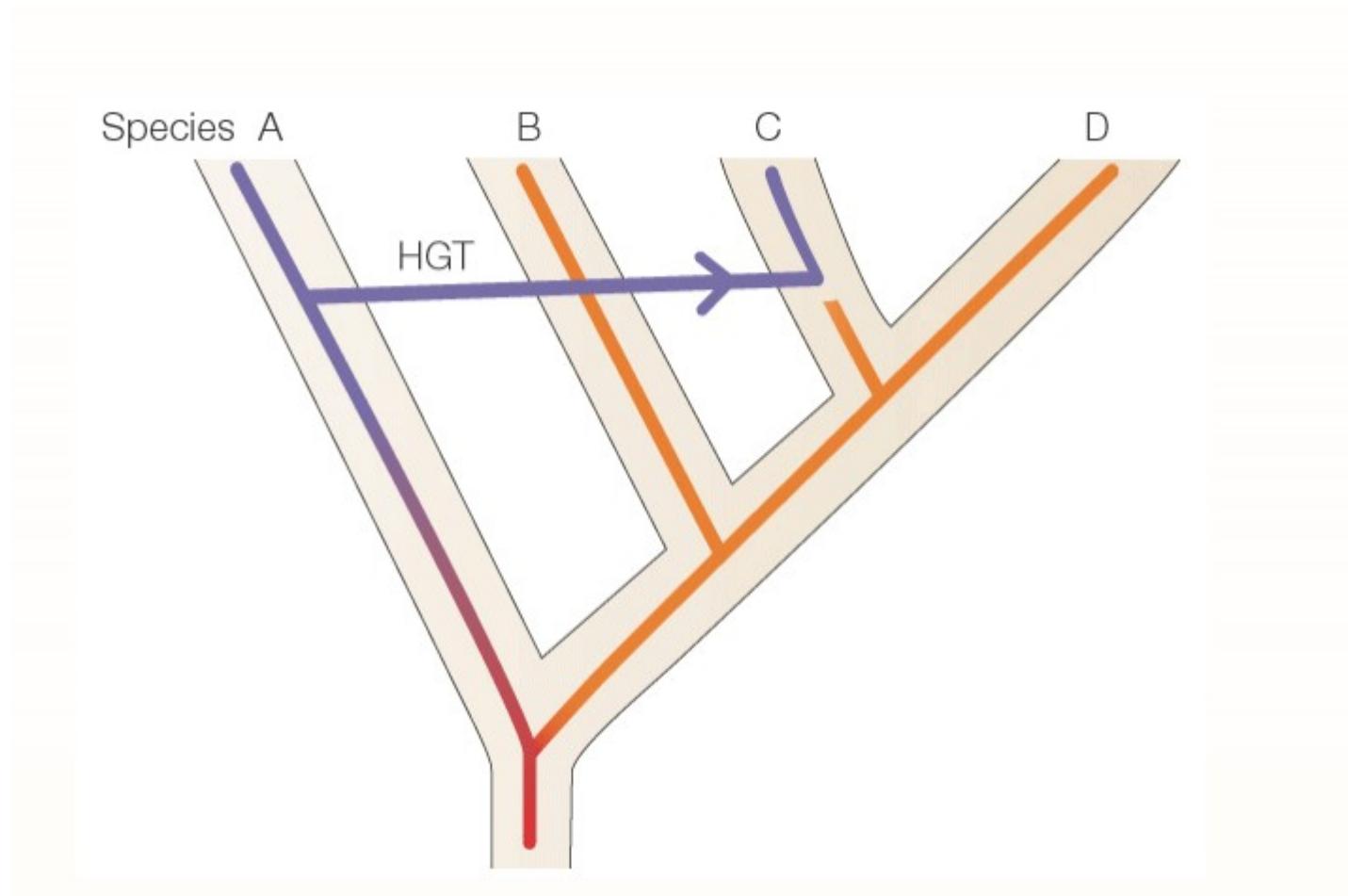


“Quantum leaps” in evolution:

- Gene (**Genome**) duplication
- (endo) symbiosis
- **hybridization**
- Lateral gene transfer

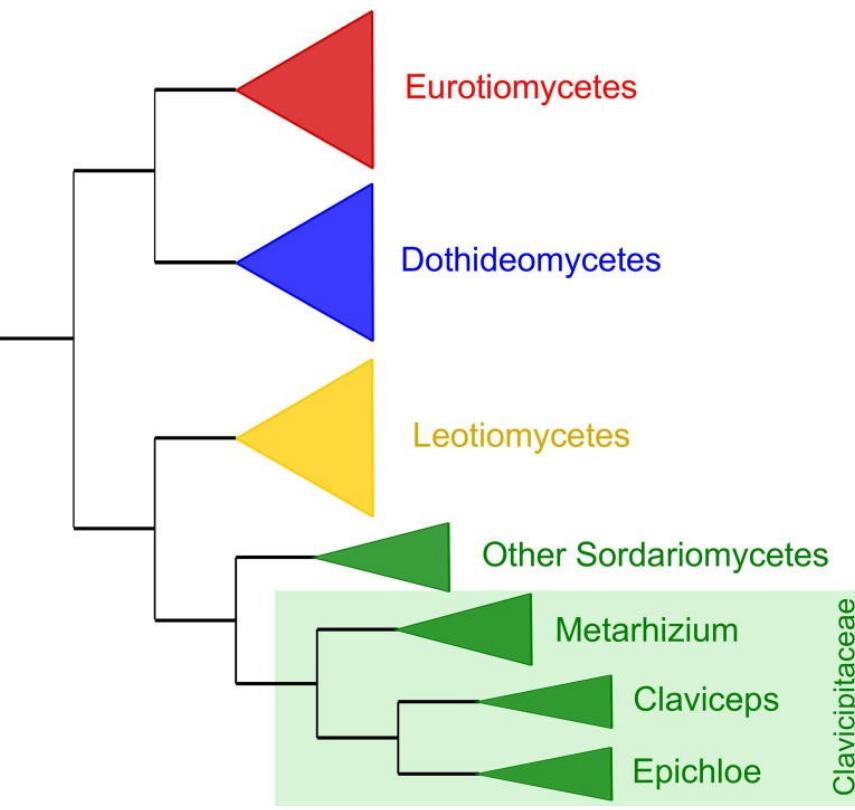


Horizontal Gene Transfer

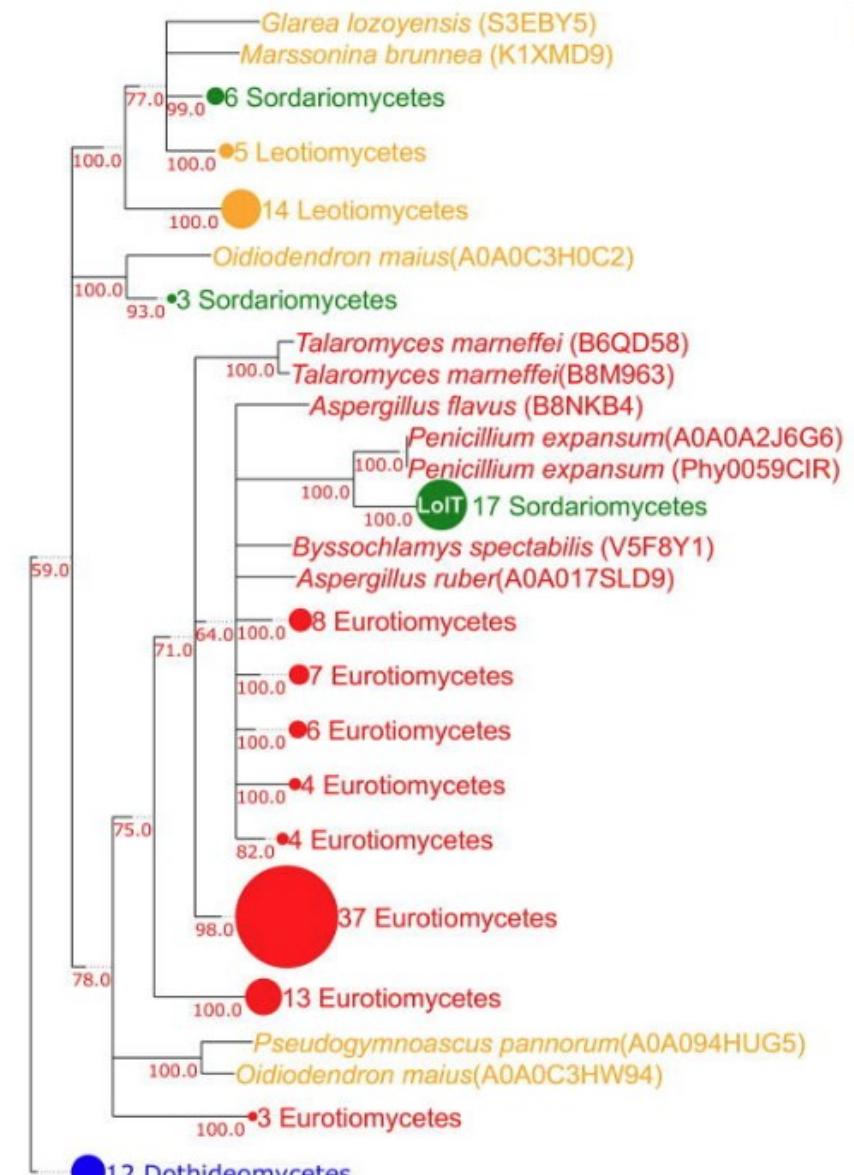


Horizontal acquisition of toxic alkaloid synthesis in a clade of plant associated fungi

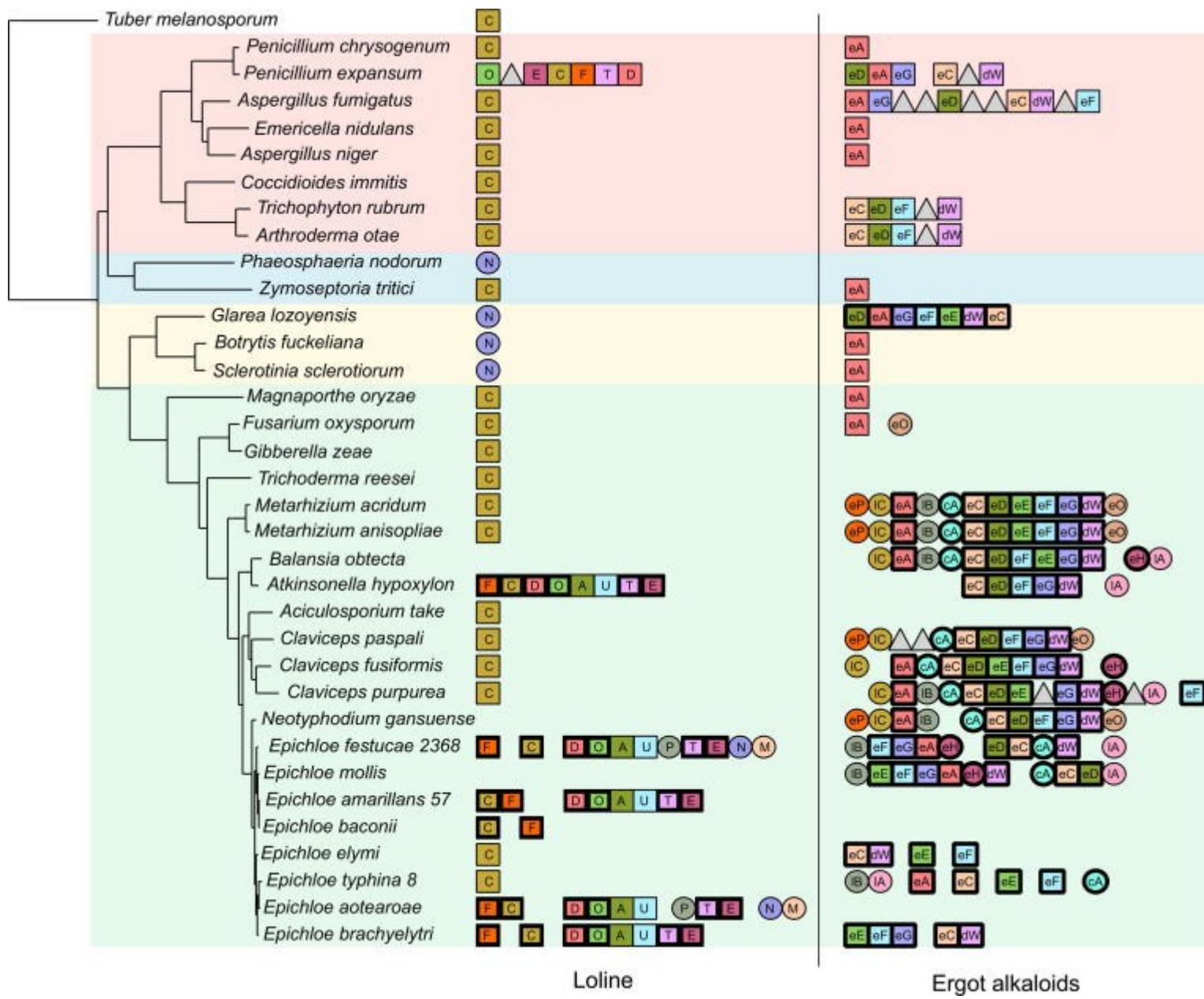
Marina Marçet-Houben^{a,b} and Toni Gabaldón^{a,b,c,*}



Species tree



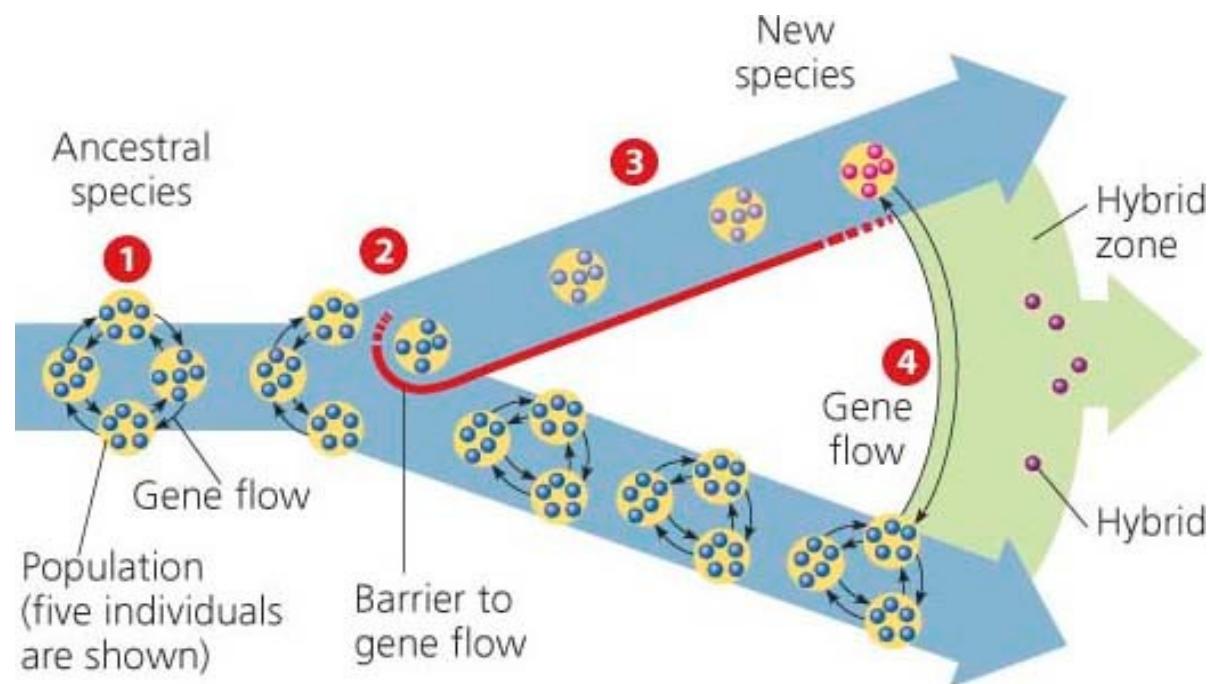
Gene tree

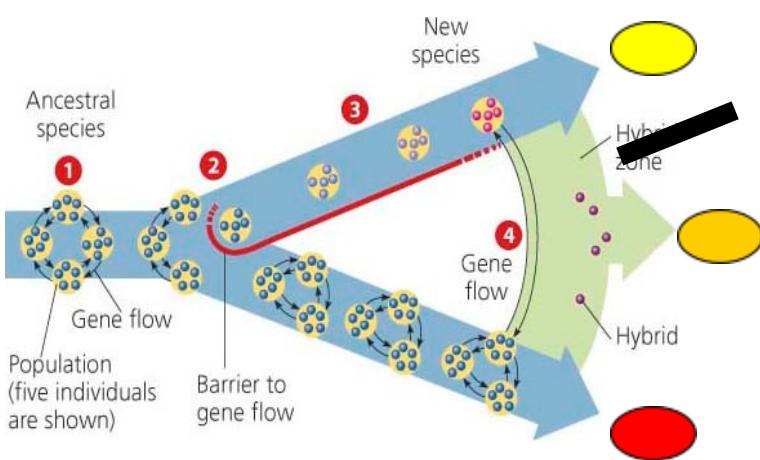


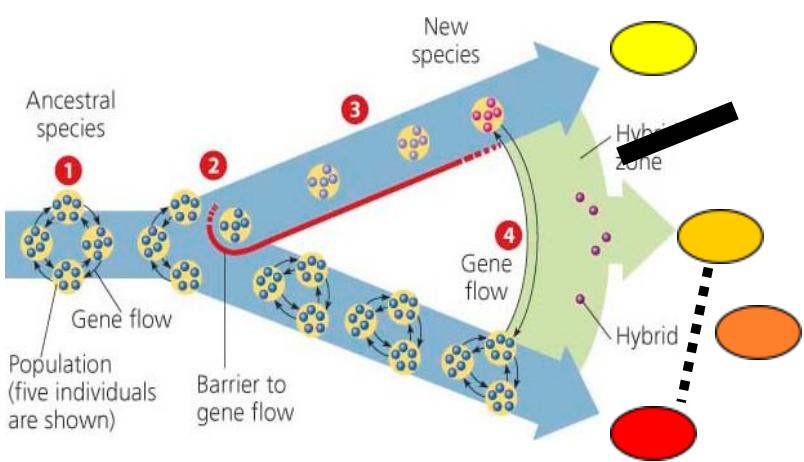
Hybrid concept is based on the concept of species (reproductive isolation)

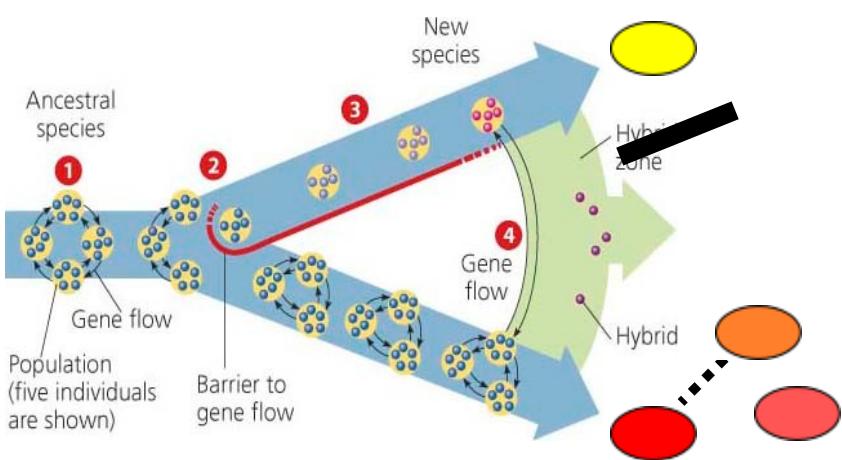


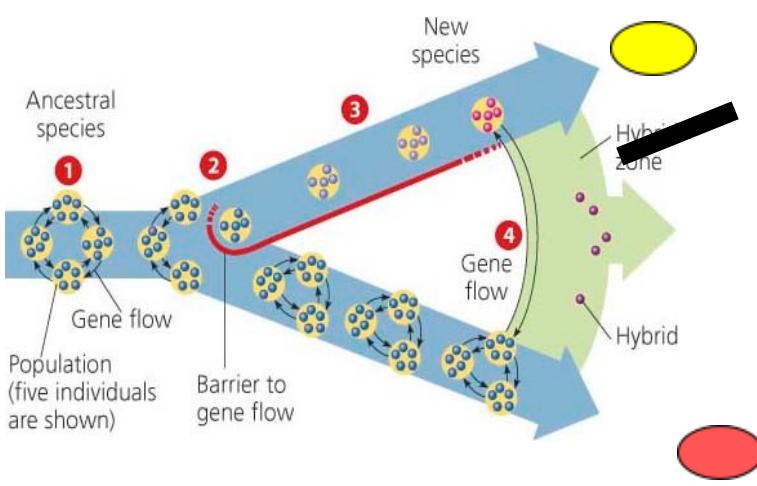
Hybridization





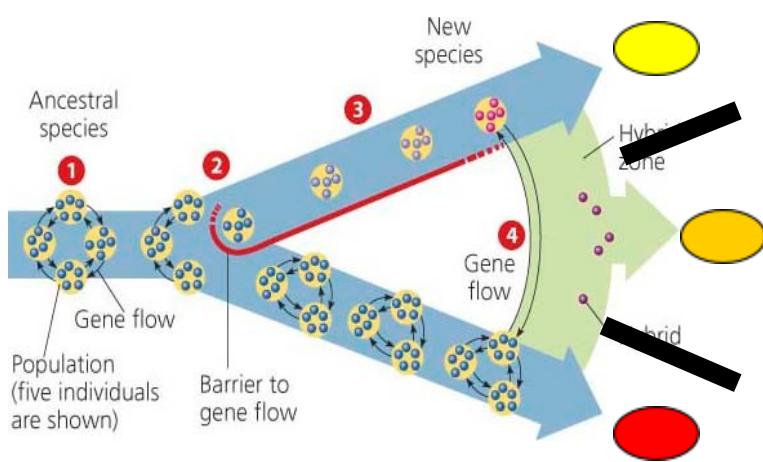


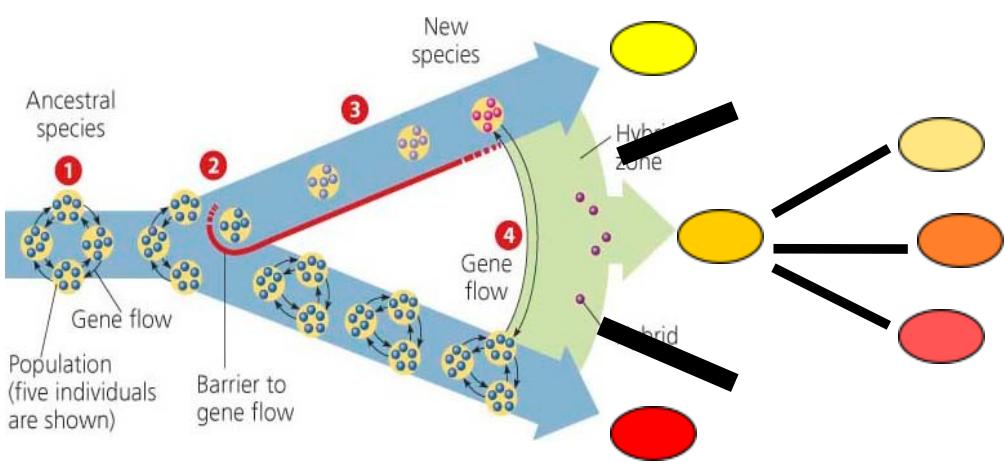


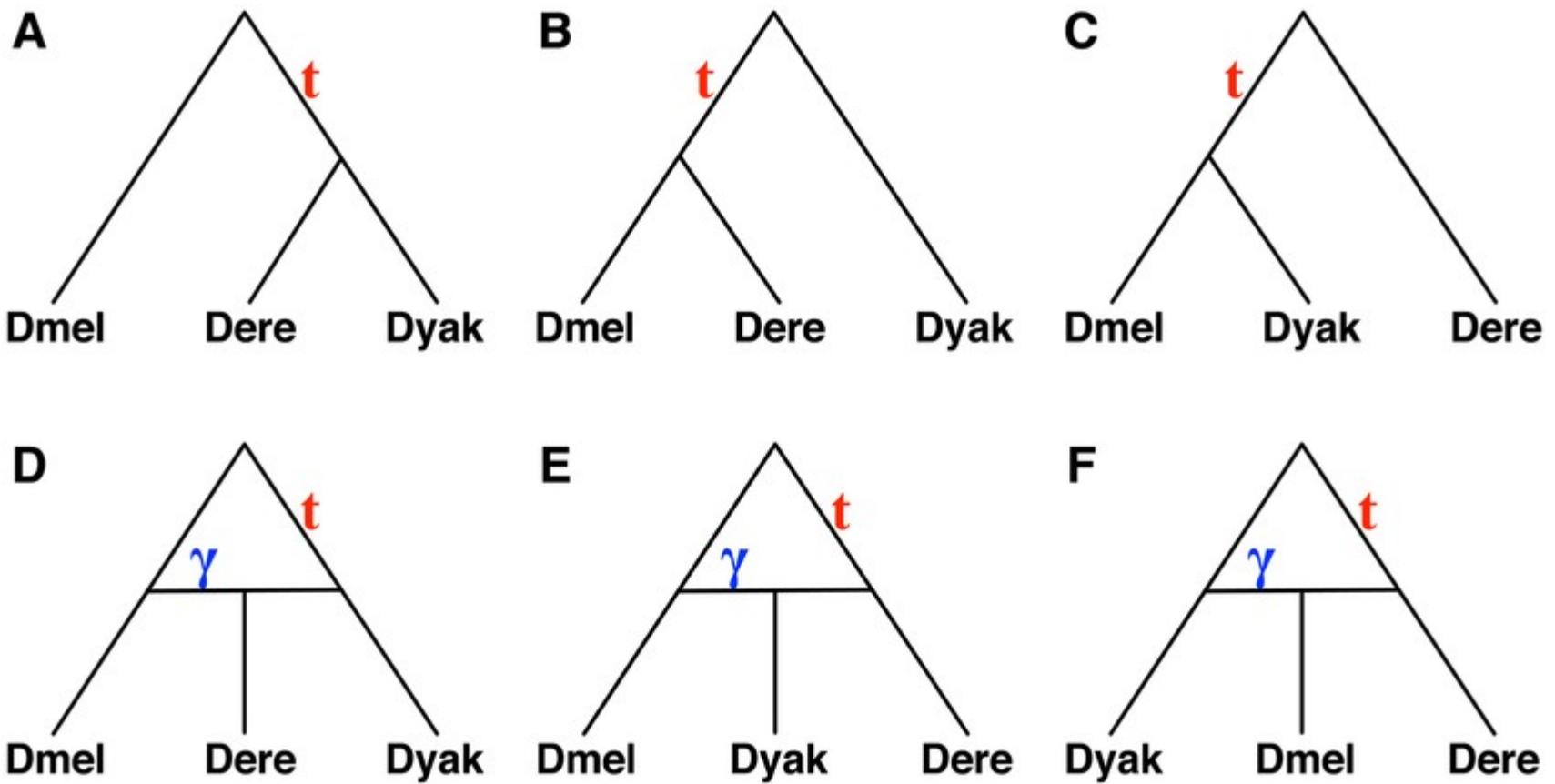


Introgression

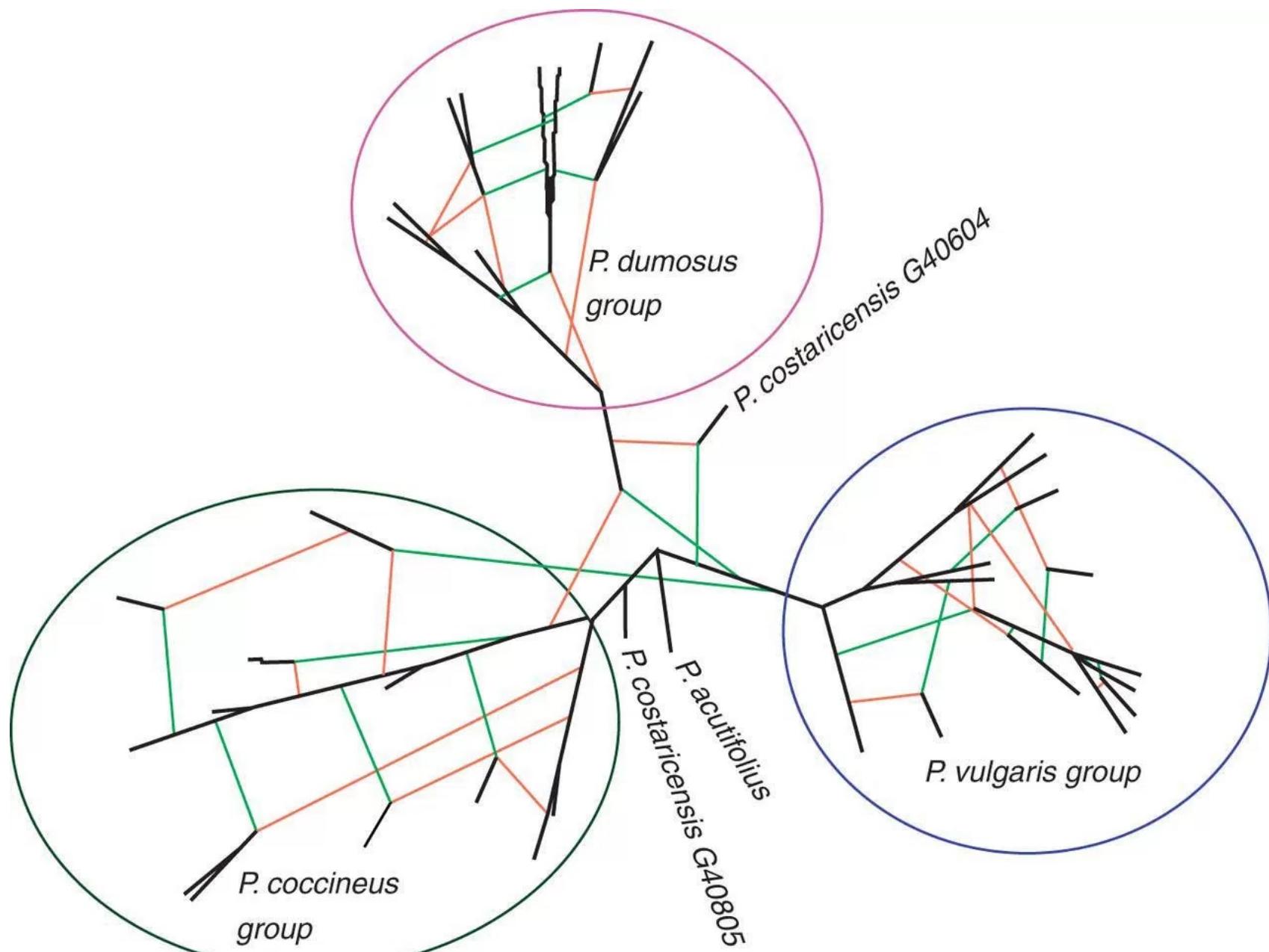
When part of the other separated species genome becomes part of the genome of the species



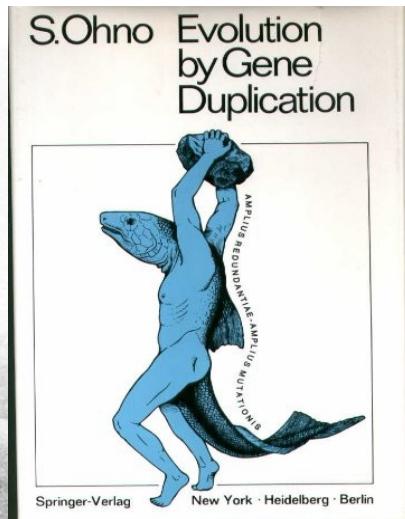
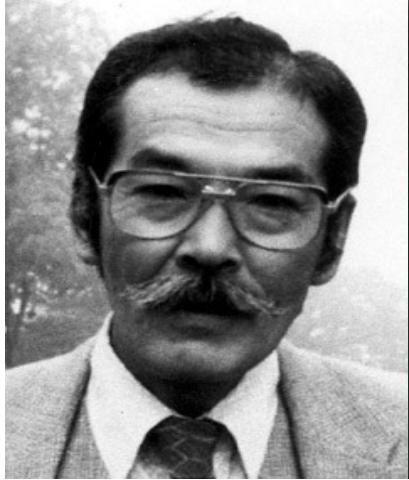




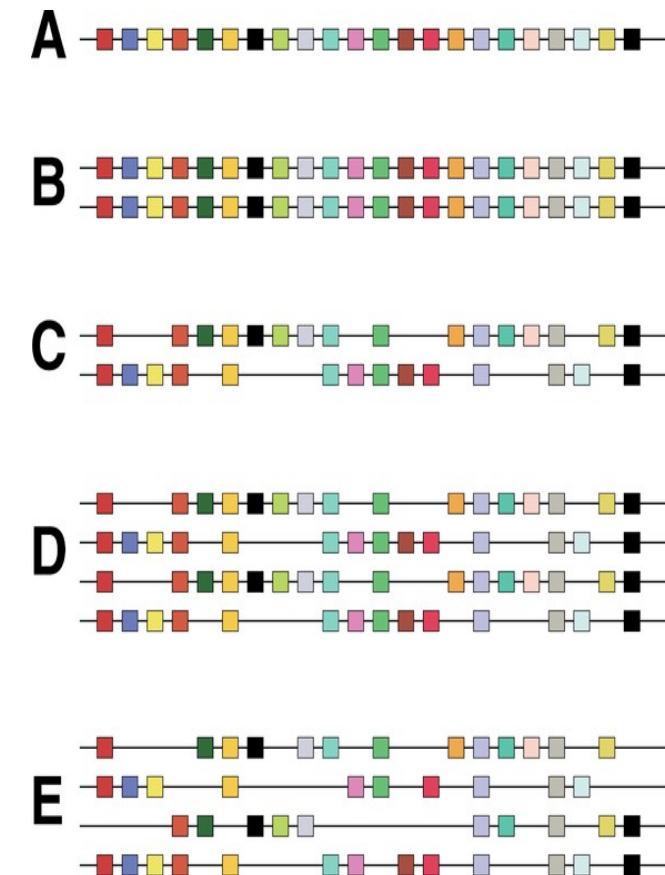
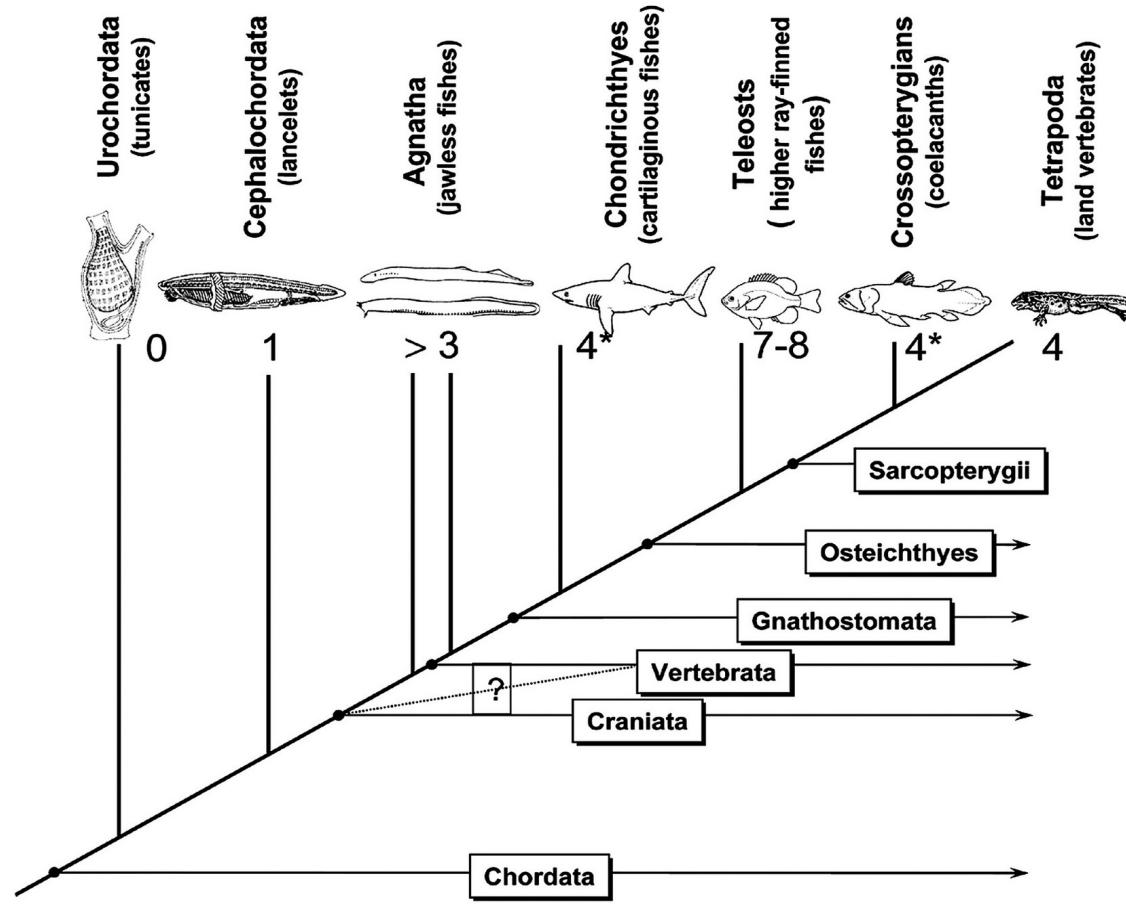
Hybridization can lead to networks rather than trees

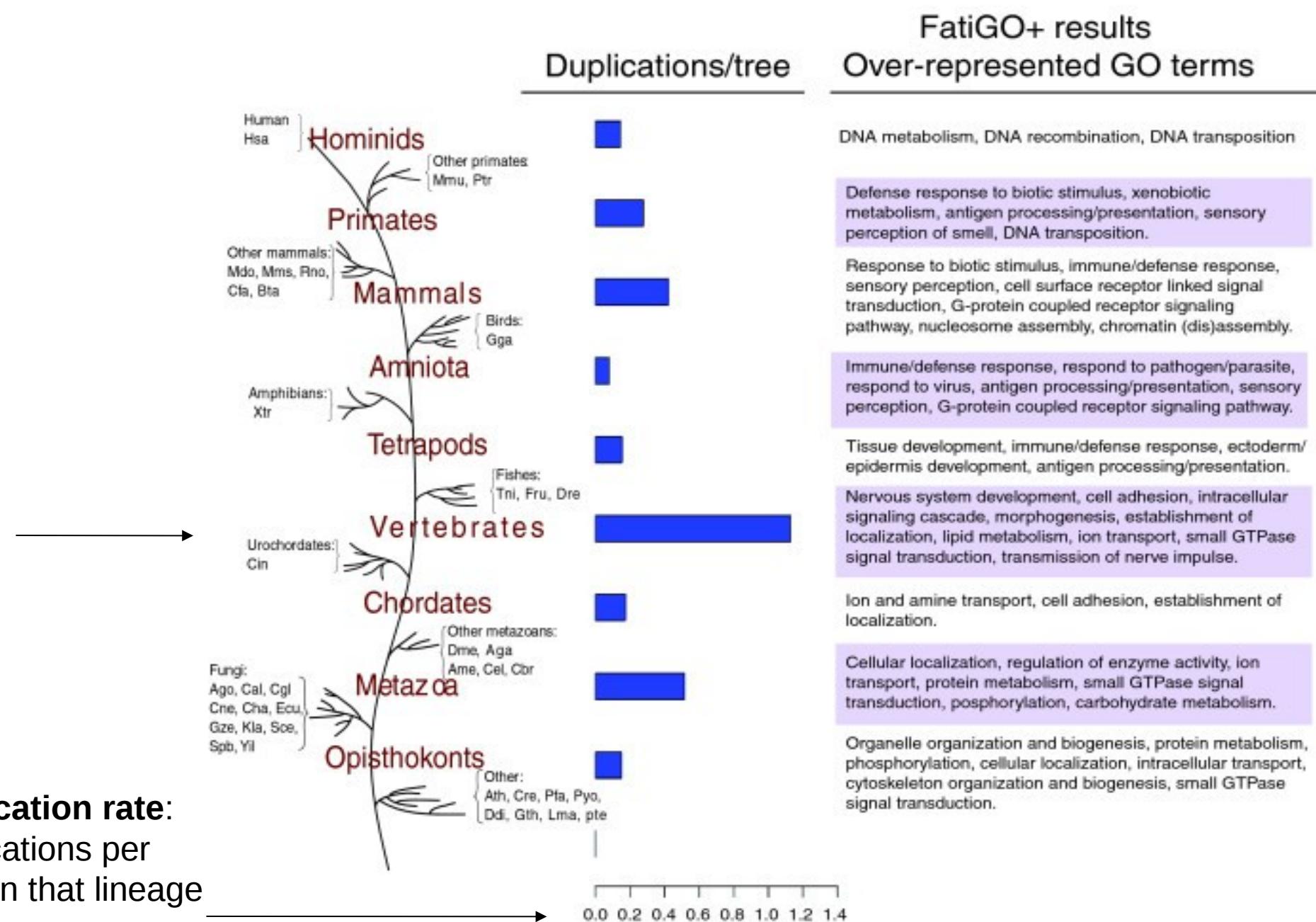


Whole genome duplications



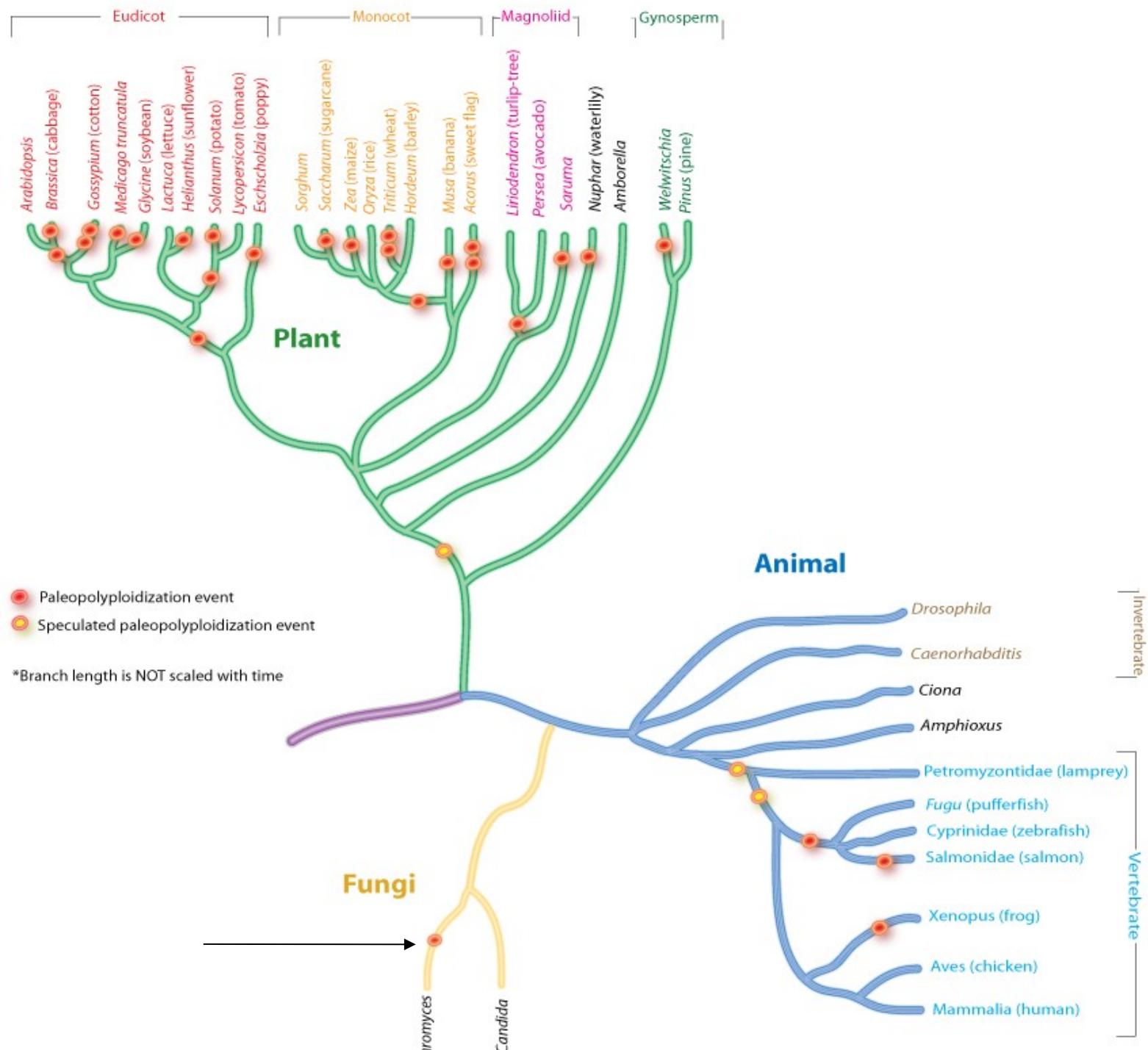
Susumu Ohno (1970)

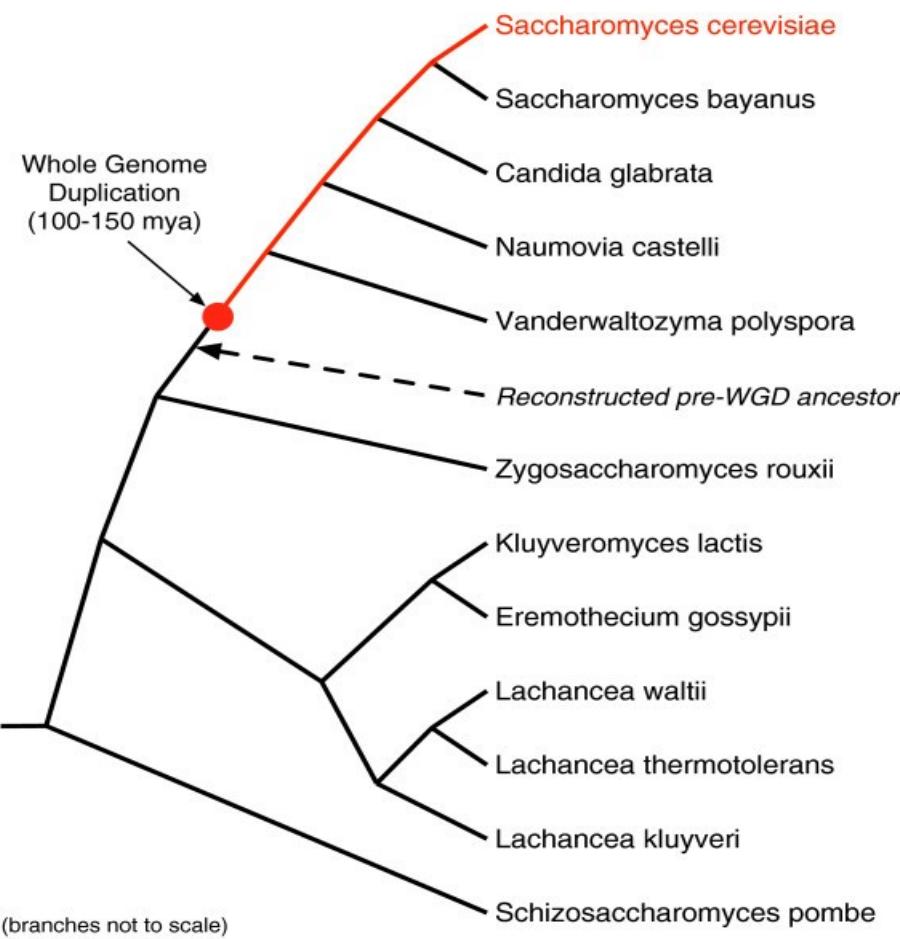




Huerta-Cepas et. al. (2007)

Known Paleopolyploidy in Eukaryotes





Article

Nature **428**, 617-624 (8 April 2004) | doi:10.1038/nature02424; Received 17 December 2003; Accepted 19 January 2004; Published online 7 March 2004

Proof and evolutionary analysis of ancient genome duplication in the yeast *Saccharomyces cerevisiae*

Manolis Kellis^{1,2}, Bruce W. Birren¹ & Eric S. Lander^{1,3}

Letters to Nature

Nature **387**, 708-713 (12 June 1997) | ; Received 17 December 1996; Accepted 20 April 1997

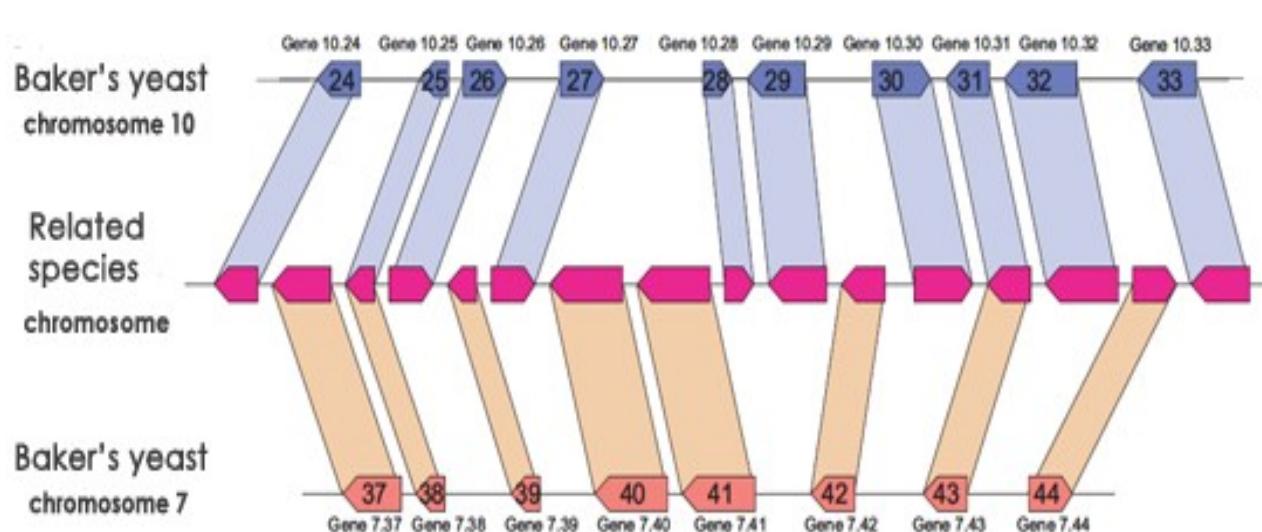
Molecular evidence for an ancient duplication of the entire yeast genome

Kenneth H. Wolfe¹ & Denis C. Shields¹

ARTICLE LINKS

▶ Figures and tables

ARTICLE TOOLS



Hybridizations can lead to whole genome duplications

