

# Introducción a la filogenómica



CENTRO DE  
INVESTIGACIONES  
BIOLÓGICAS DEL  
NOROESTE, S.C.

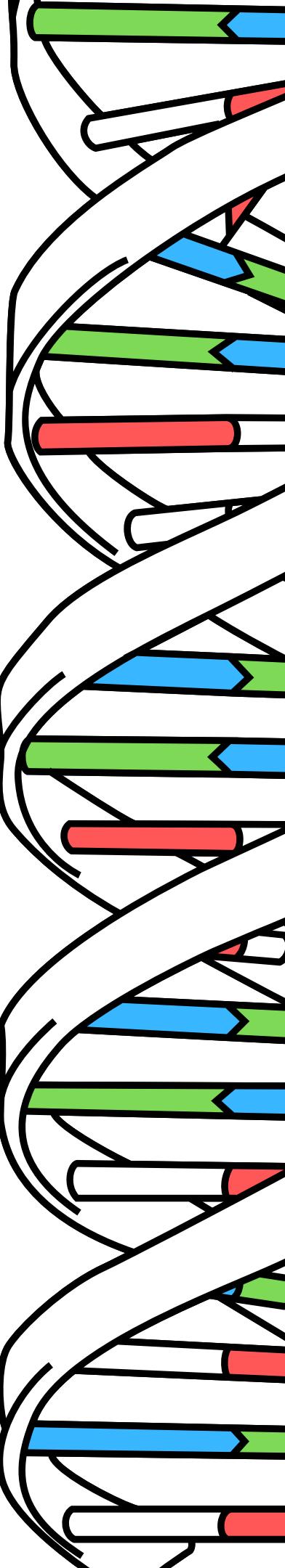


# Introducción

Día	Temas
23 de septiembre	<ul style="list-style-type: none"><li>• <b>Introducción a la Filogenómica. Conceptos básicos</b></li><li>• Introducción al sistema operativo Linux</li></ul>
24 de septiembre	<ul style="list-style-type: none"><li>• Métodos de construcción de librerías genómicas</li><li>• Control de calidad, ensamble y alineamiento de secuencias</li></ul>
25 de septiembre	<ul style="list-style-type: none"><li>• Método de Máxima Verosimilitud</li><li>• Método de Inferencia Bayesiana</li></ul>
26 de septiembre	<ul style="list-style-type: none"><li>• Modelo coalescente multiespecies</li><li>• Evaluación de árboles de especies</li></ul>
27 de septiembre	<ul style="list-style-type: none"><li>• Estimación de tiempos de divergencia</li><li>• Evaluación de flujo genético</li></ul>

## Introducción

Millones de años de evolución  
compartida con otras formas  
de vida



# Introducción

## Gran cadena del ser

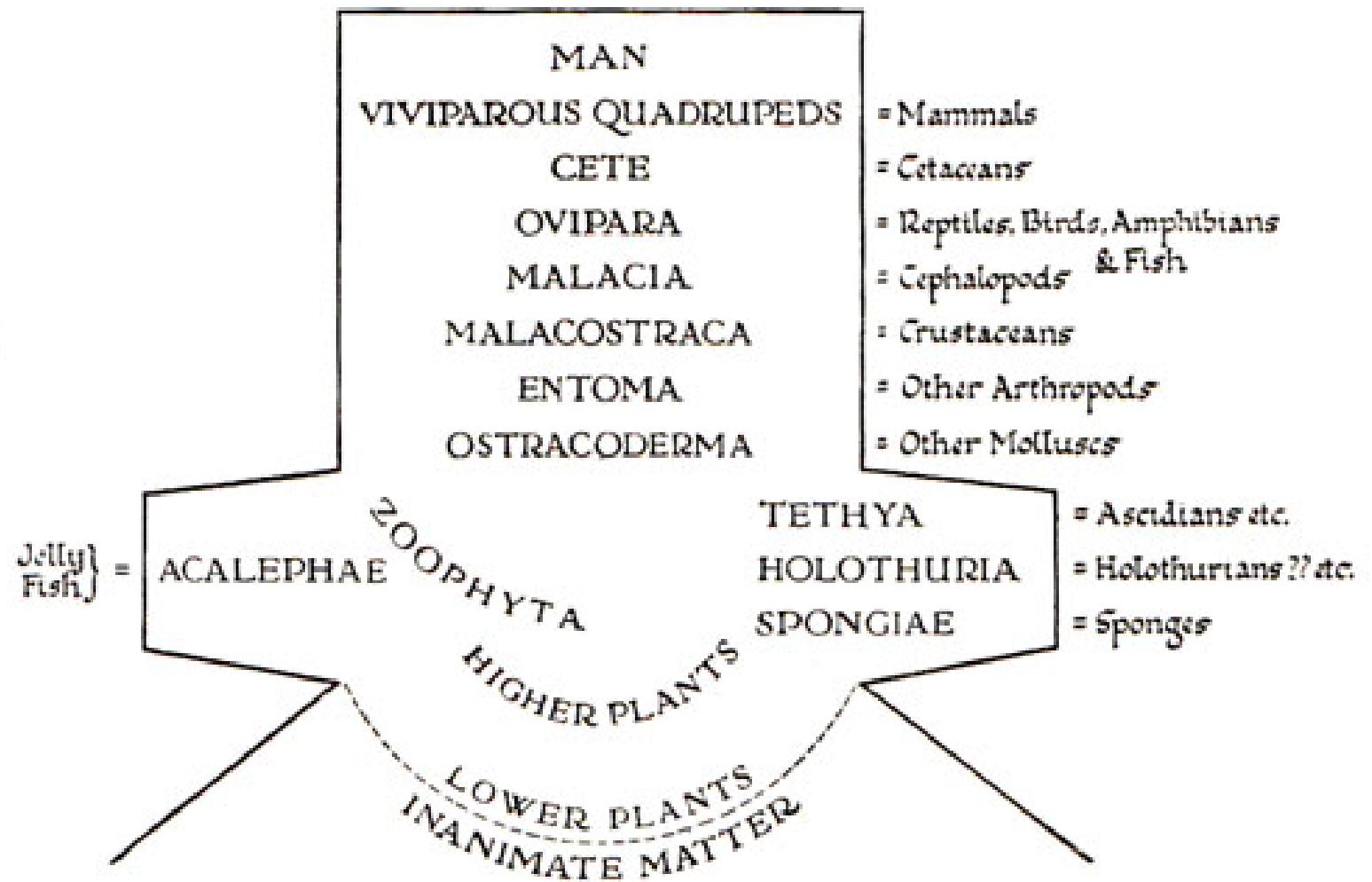


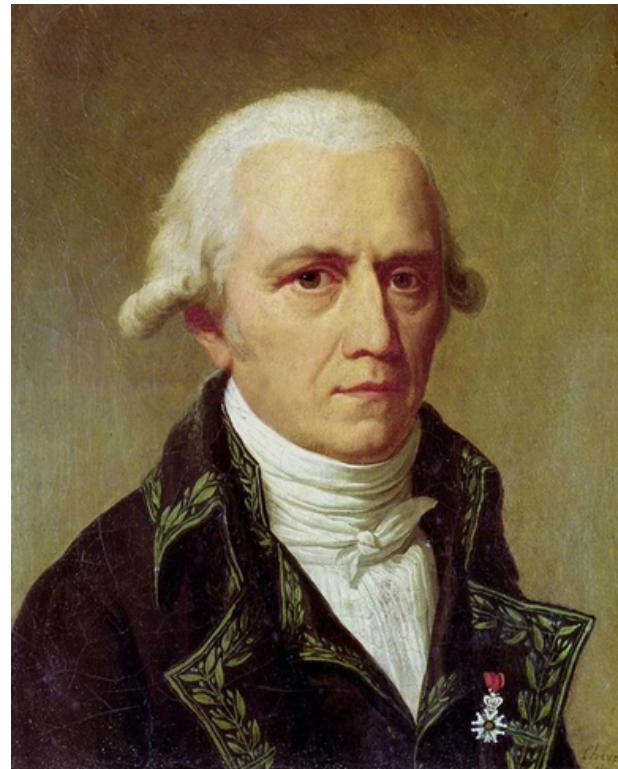
FIG. 18. The *Scala Naturae* or 'Ladder of Life' according to the descriptions of Aristotle.



Jerarquía lineal  
de la vida

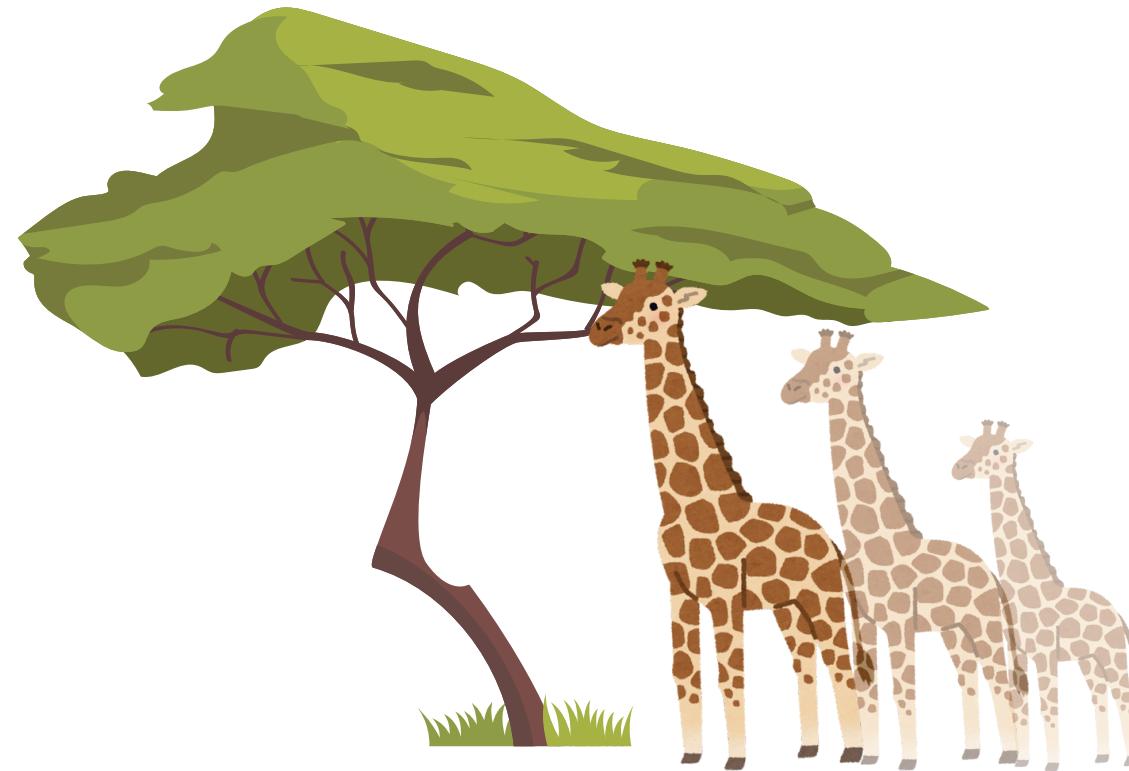
Mundo estático

# Introducción



**Jean-Baptiste  
Lamarck**

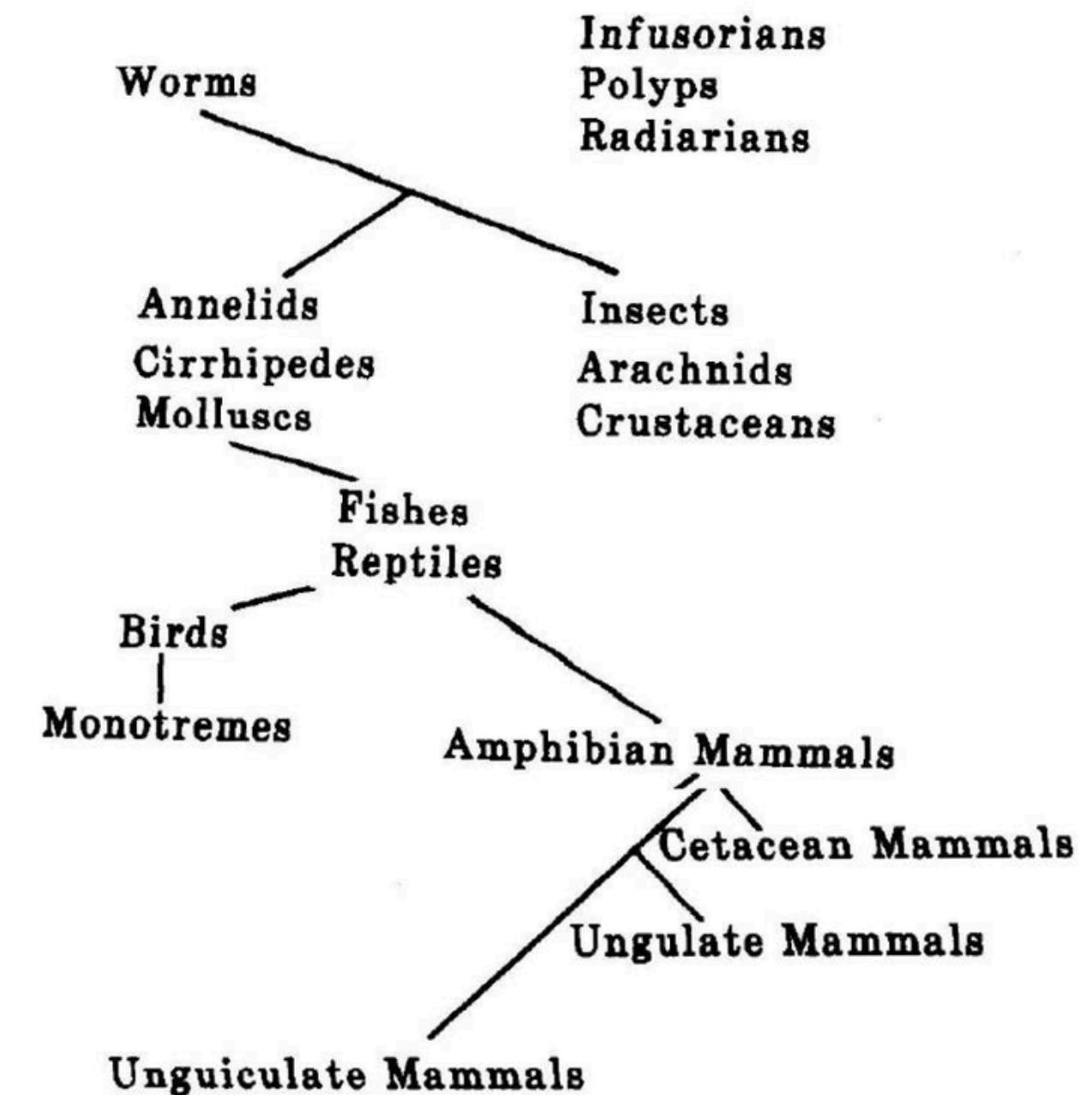
*Philosophie  
Zoologique* (1809)



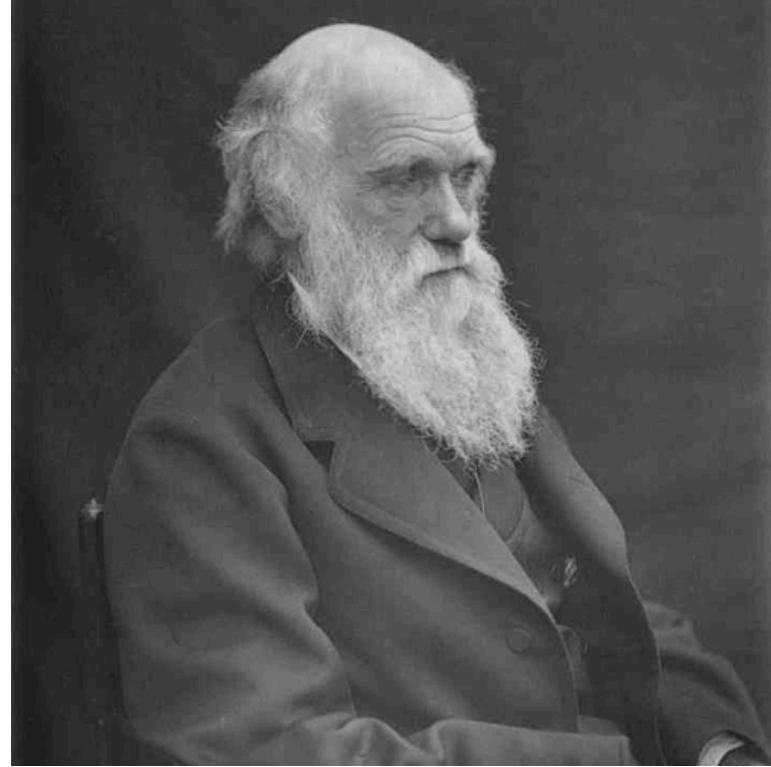
La vida evoluciona  
Árbol de animales  
Diferentes orígenes  
No lineal  
Progreso

## TABLE

### SHOWING THE ORIGIN OF THE VARIOUS ANIMALS

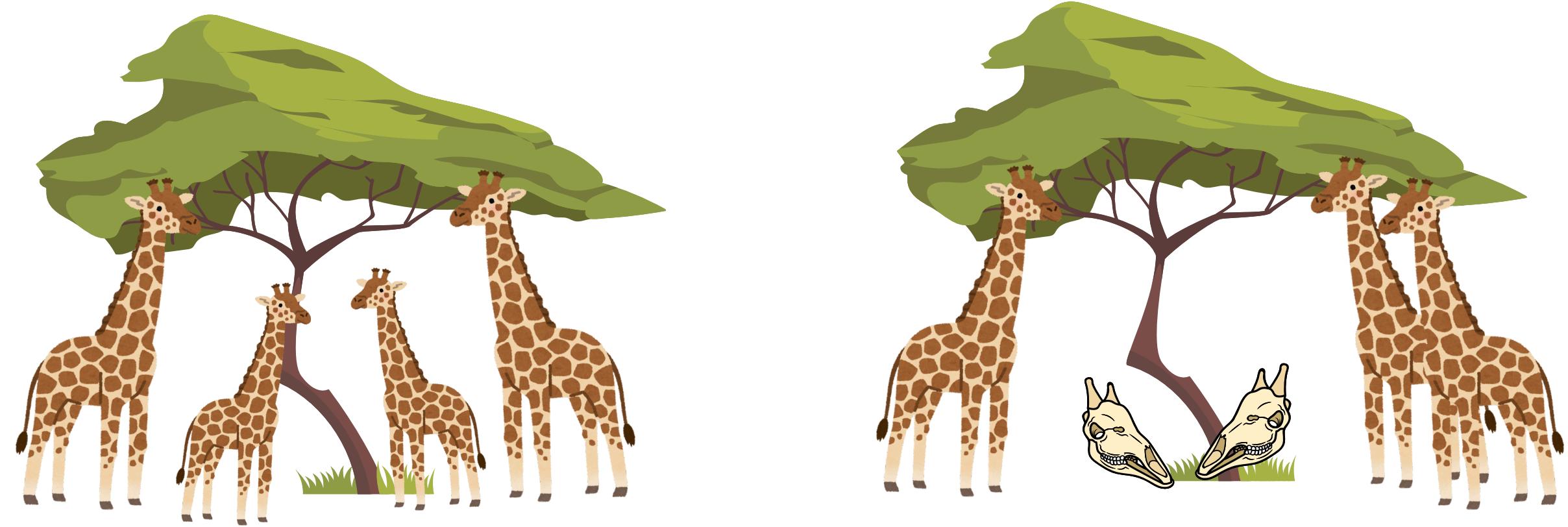


# Introducción



## Charles Darwin

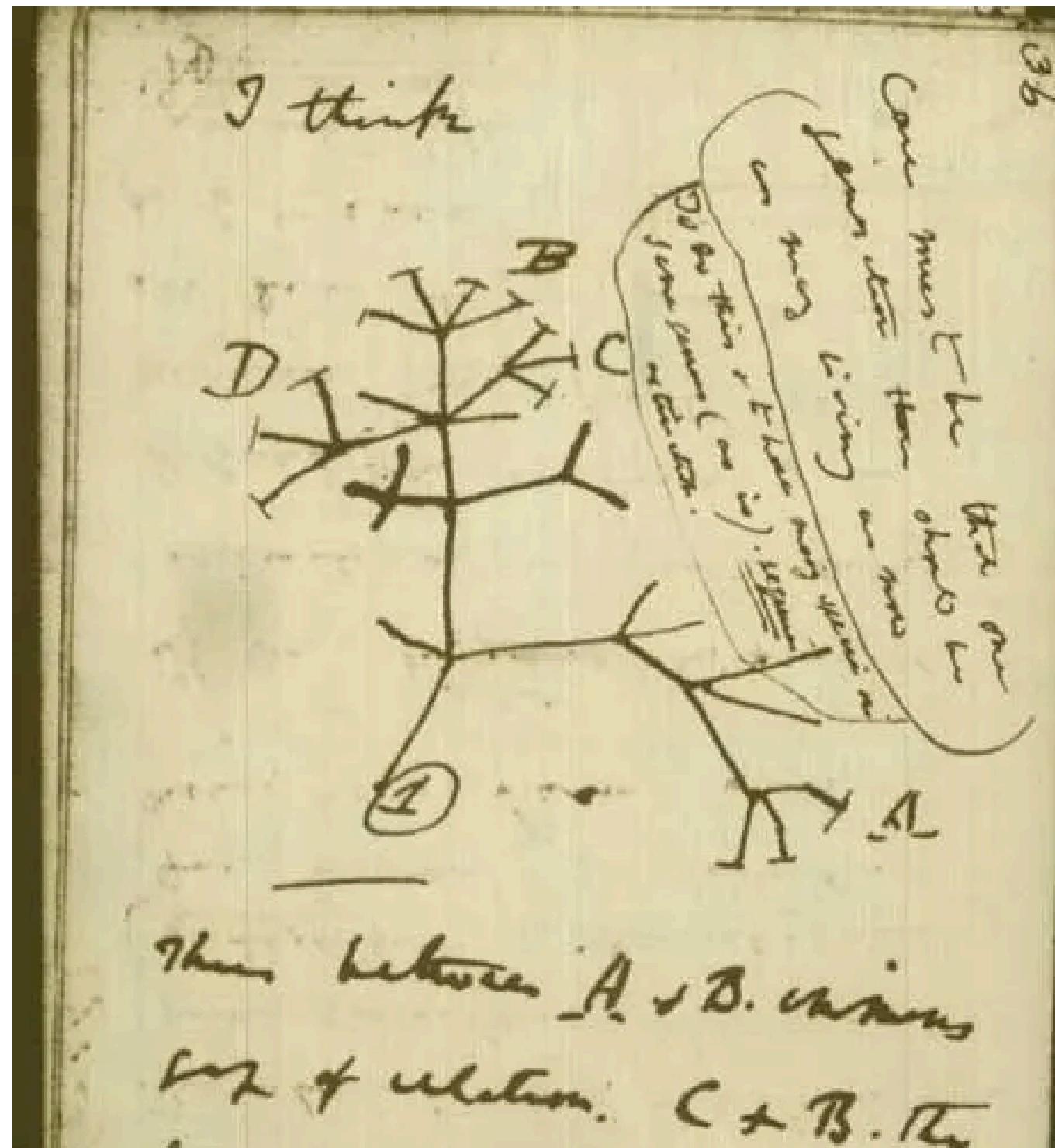
*Sobre el origen de las especies* por medio de la selección natural (1859).



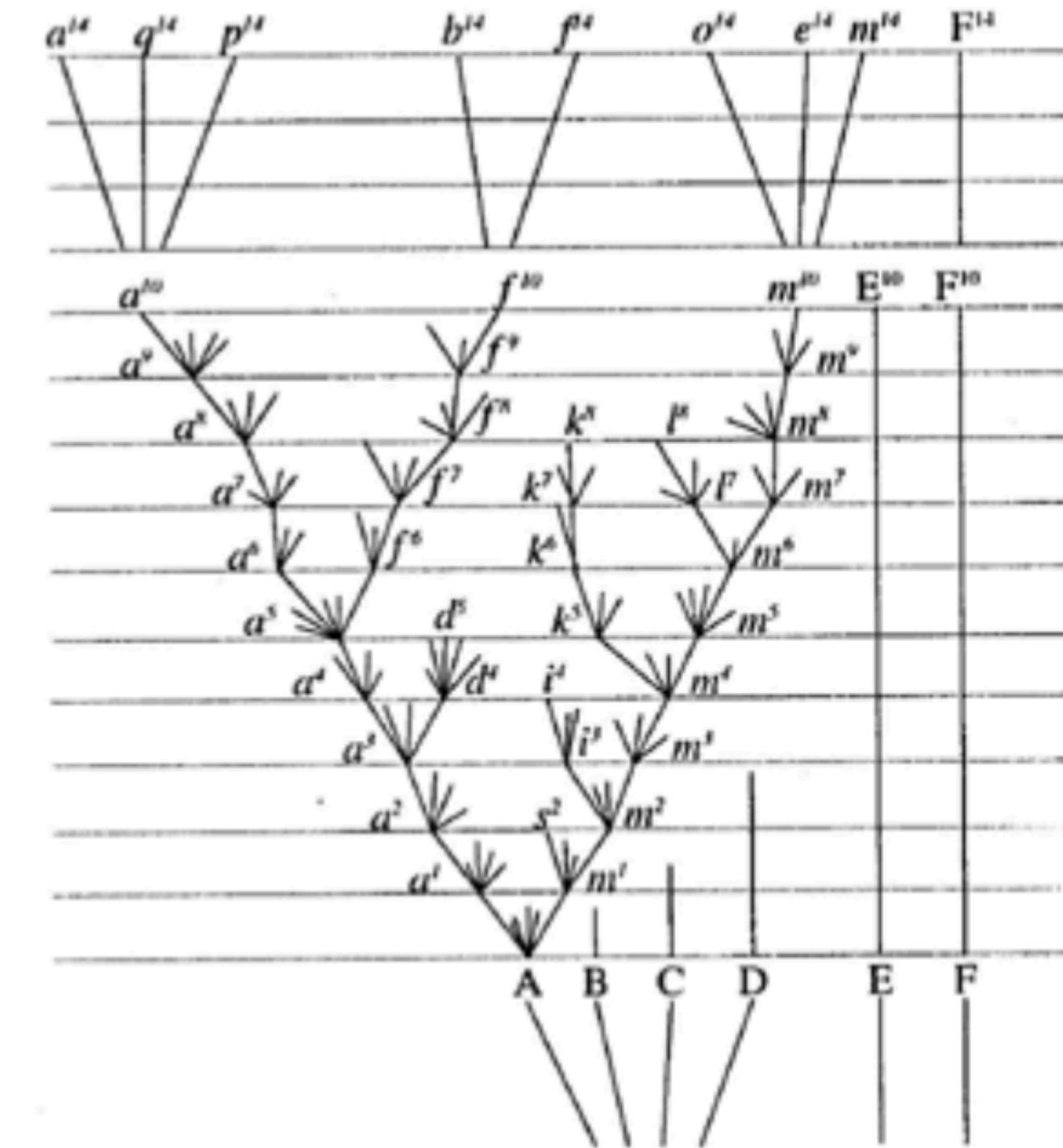
## Evolución por selección natural

Concluyó que las diversas especies vivientes se remontan a los mismos antepasados comunes.

# Introducción



Darwin, 1837



Darwin, 1859

...una generación sería tantas como las que viven ahora.

# Introducción

Al considerar el origen de las especies se concibe perfectamente que un naturalista, reflexionando sobre las afinidades mutuas de los seres orgánicos, sobre sus relaciones embriológicas, su distribución geográfica, sucesión geológica y otros hechos semejantes, puede llegar a la conclusión de que las especies no han sido independientemente creadas, sino que han descendido, como las variedades, de otras especies. Sin embargo, esta conclusión, aunque estuviese bien fundada, no sería

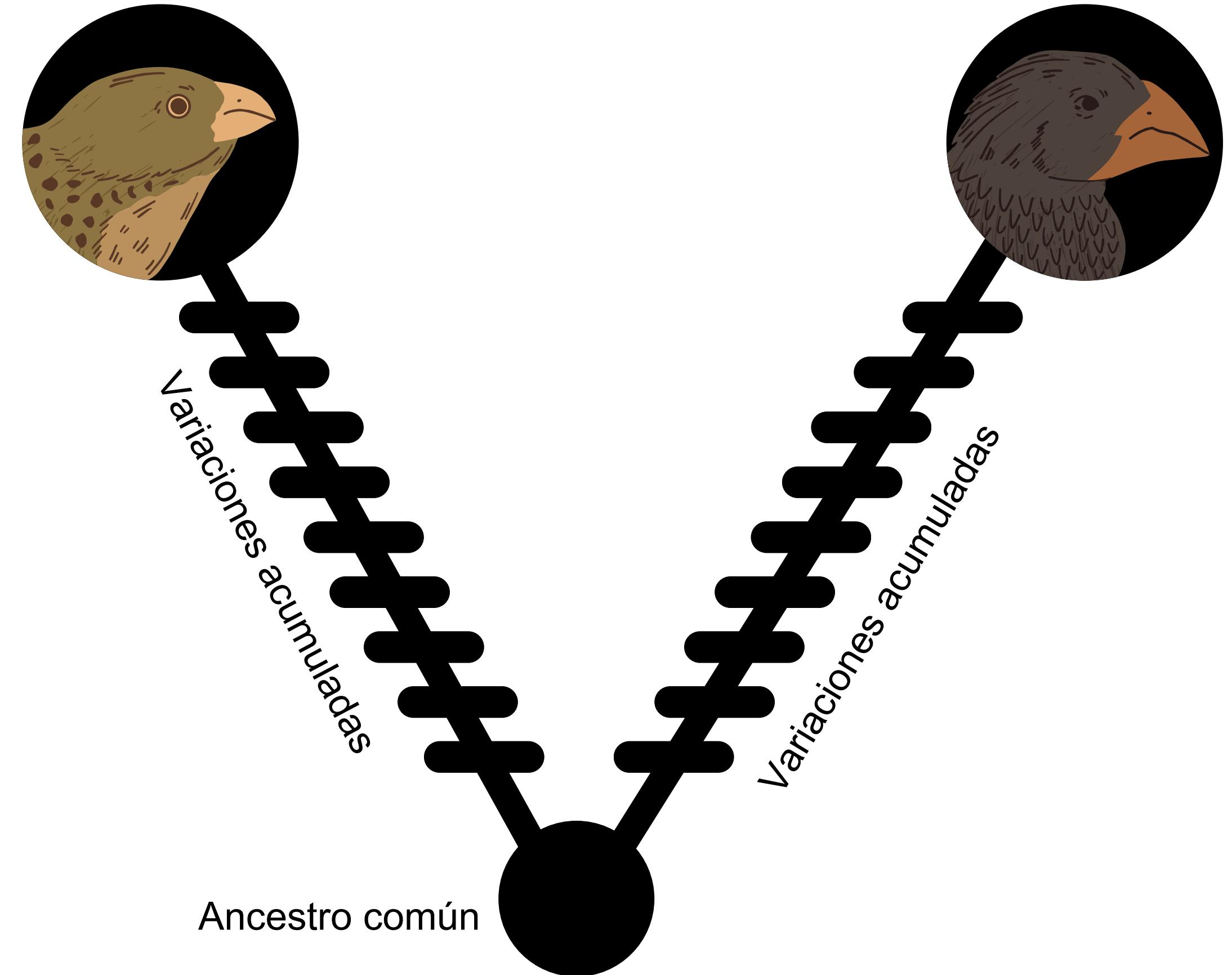
# Introducción

De las muchas ramitas que florecieron cuando el árbol era un simple arbollito, sólo dos o tres, convertidas ahora en ramas grandes, sobreviven todavía y llevan las otras ramas; de igual modo, de las especies que vivieron durante períodos geológicos muy antiguos, poquísimas han dejado descendientes vivos modificados. Desde el primer crecimiento del árbol, muchas ramas de todos tamaños se han secado

# Capítulo IV

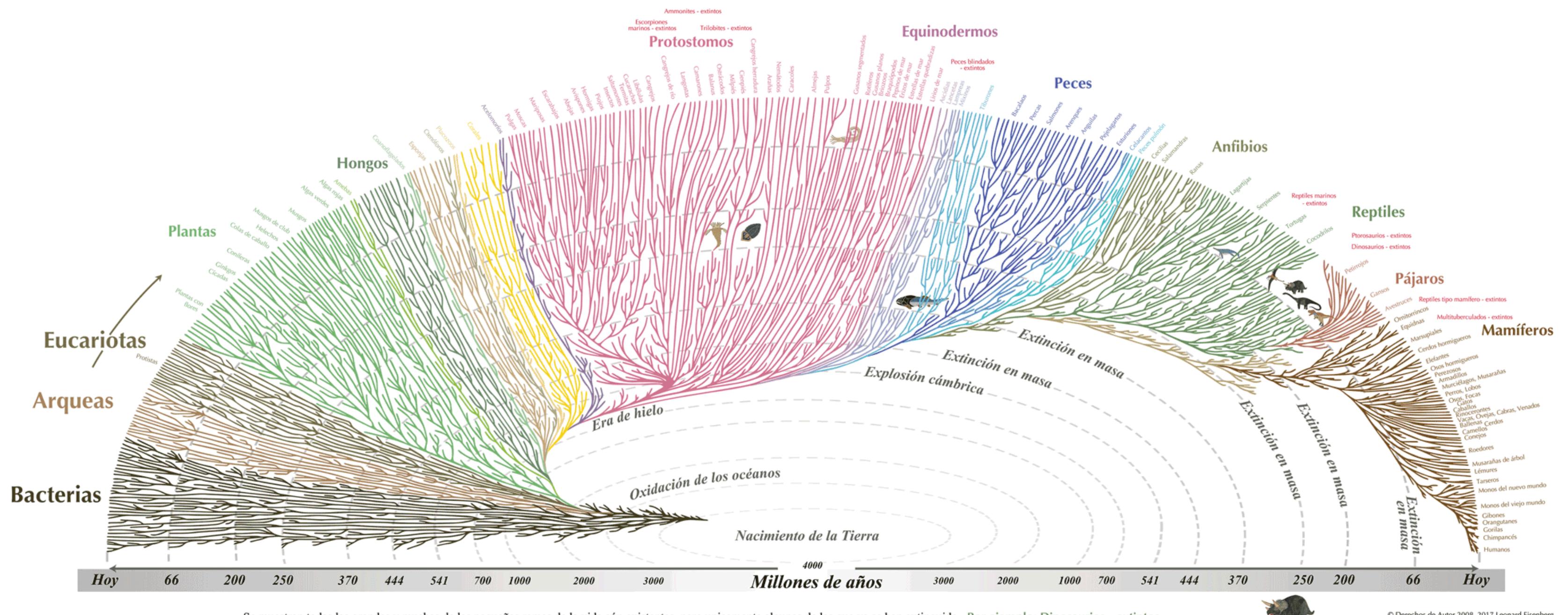
# Introducción

- Similitudes estructurales
- Similitudes moleculares
- Distribuciones geográficas
- Fósiles
- La distribución de ciertas características



# Introducción

## Representación de las relaciones de ancestría-descendencia de los linajes



# Introducción

Arbol filogenético

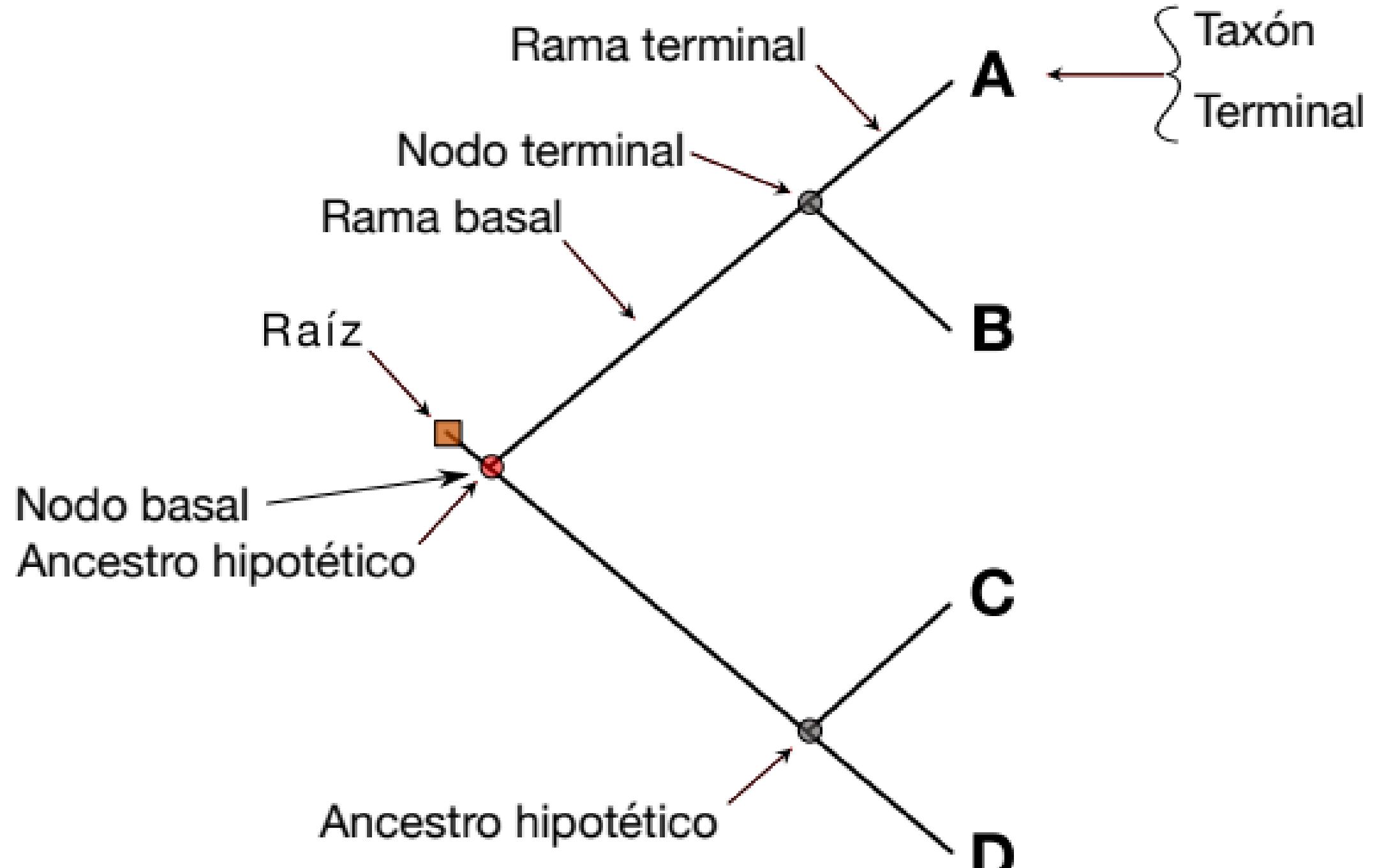
Filogenia

Árbol evolutivo

Diagrama ramificado que representa las relaciones entre organismos u otras entidades.

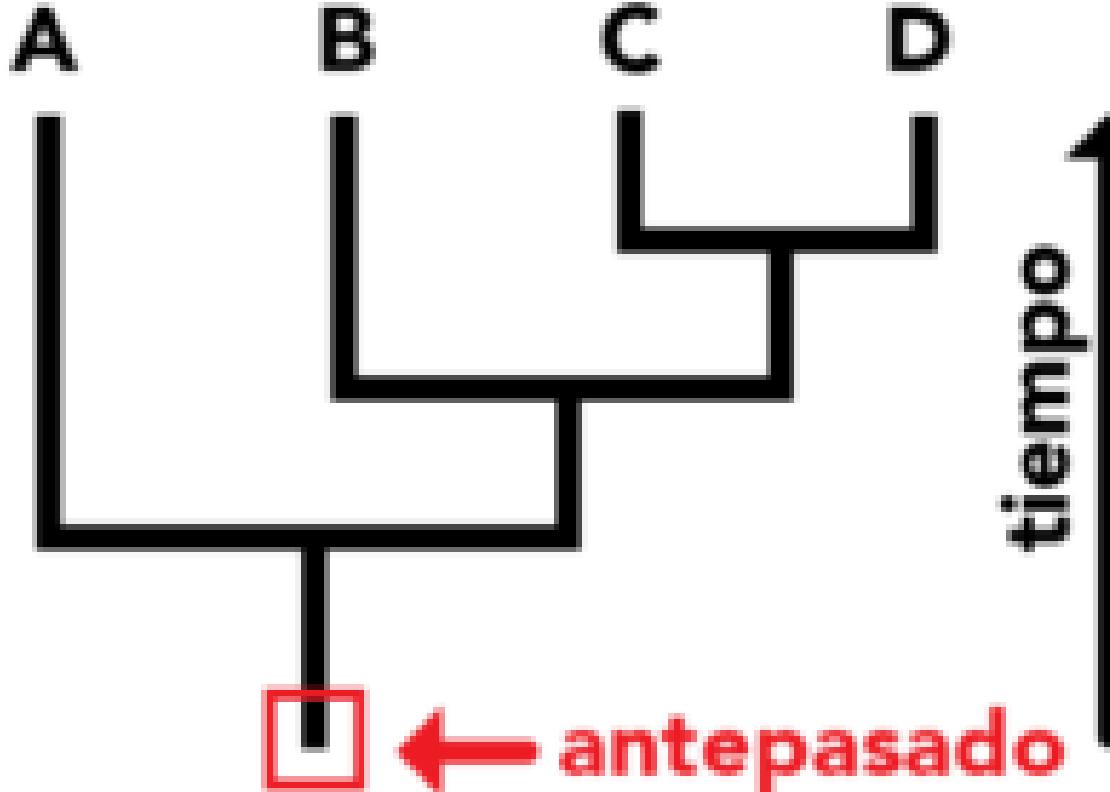
Es una hipótesis

Contexto

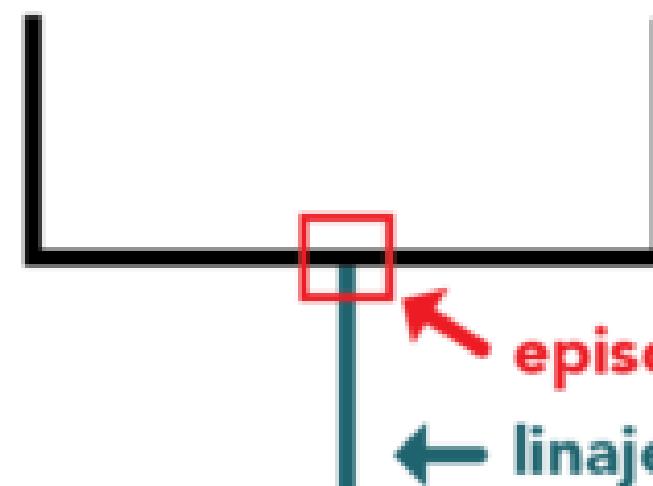


# Introducción

**descendientes → A**

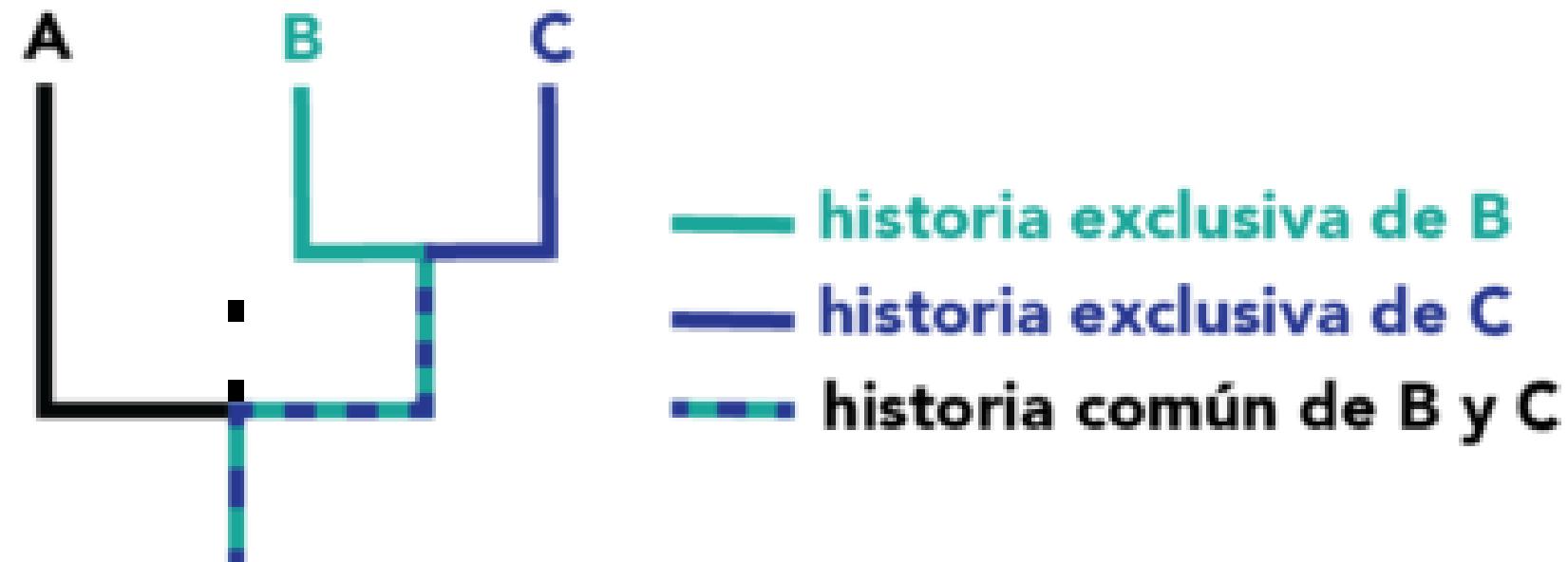


Descendientes que comparten  
un ancestro común.

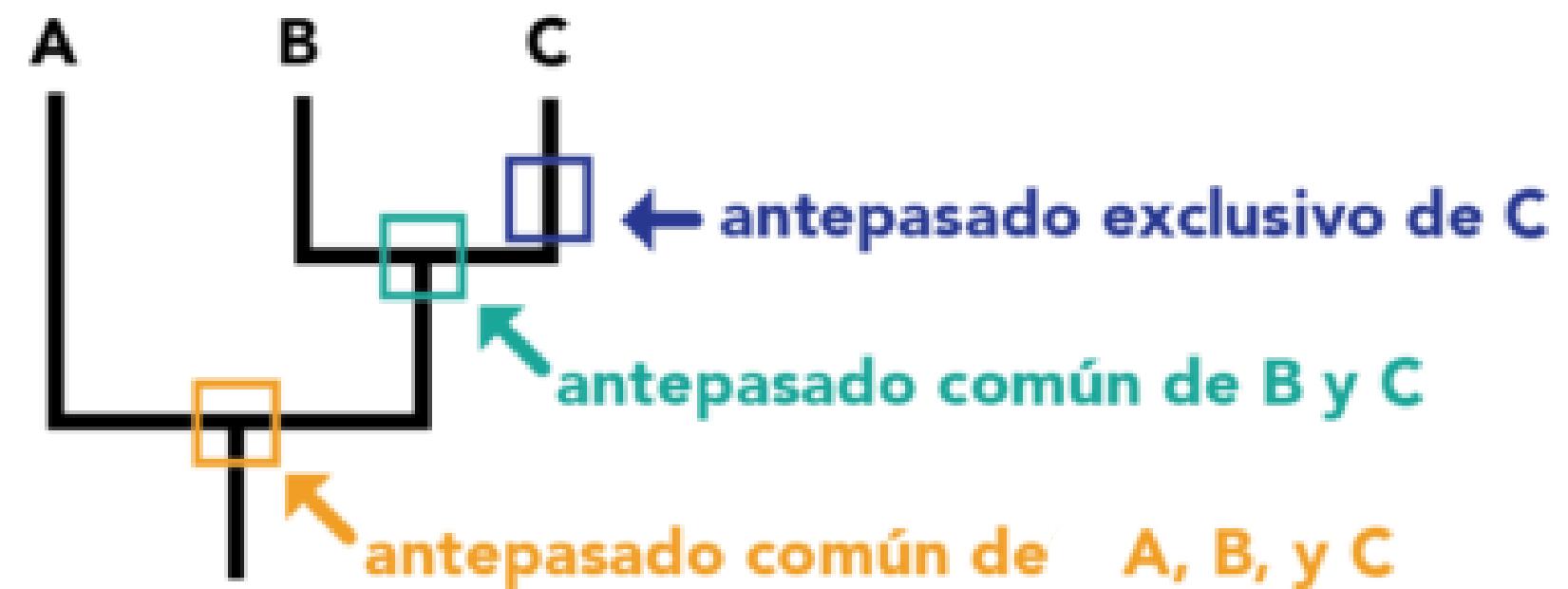


Cada ramificación indica un  
evento de divergencia.

# Introducción



B y C comparten historia evolutiva común antes de divergir

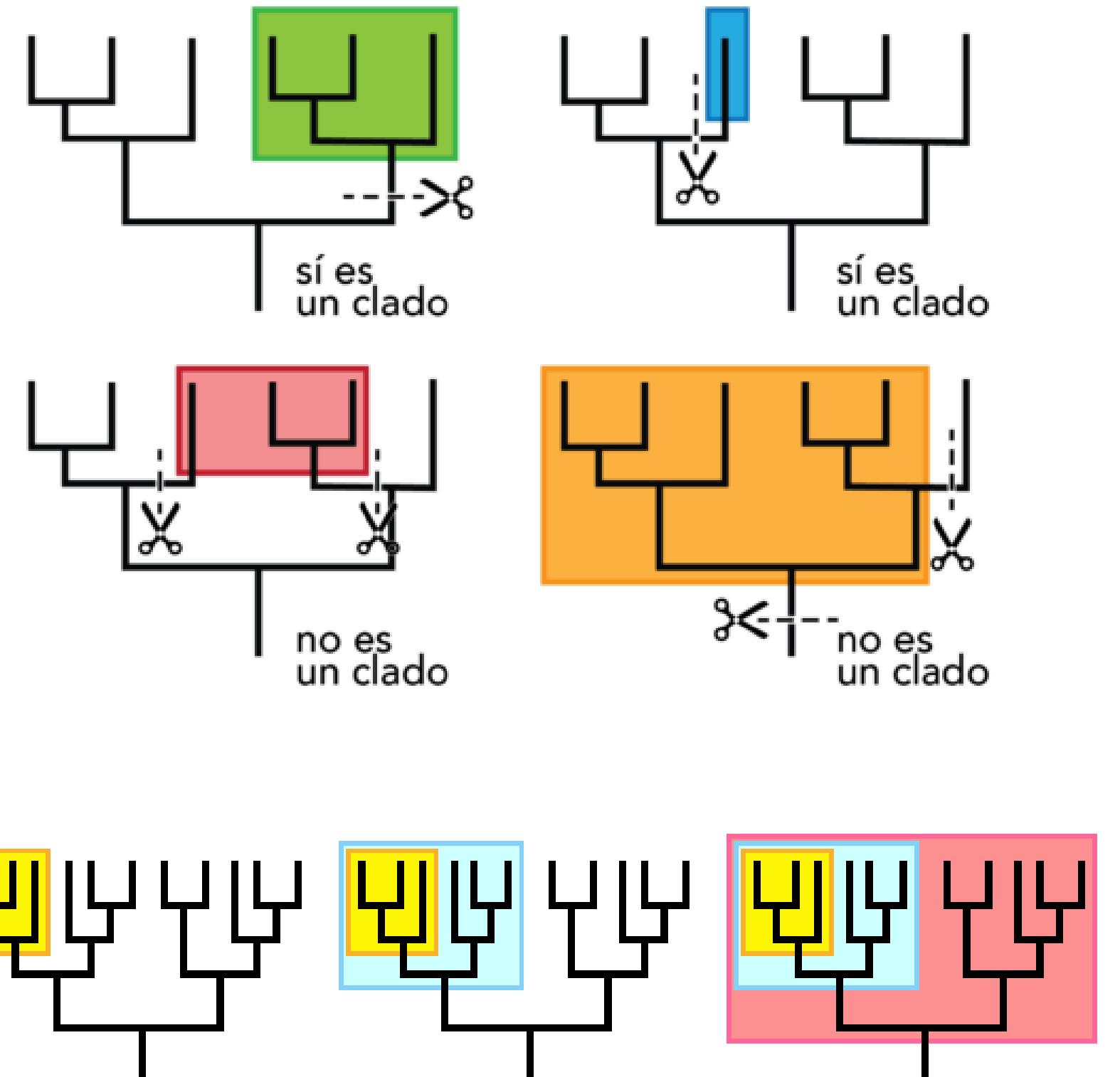
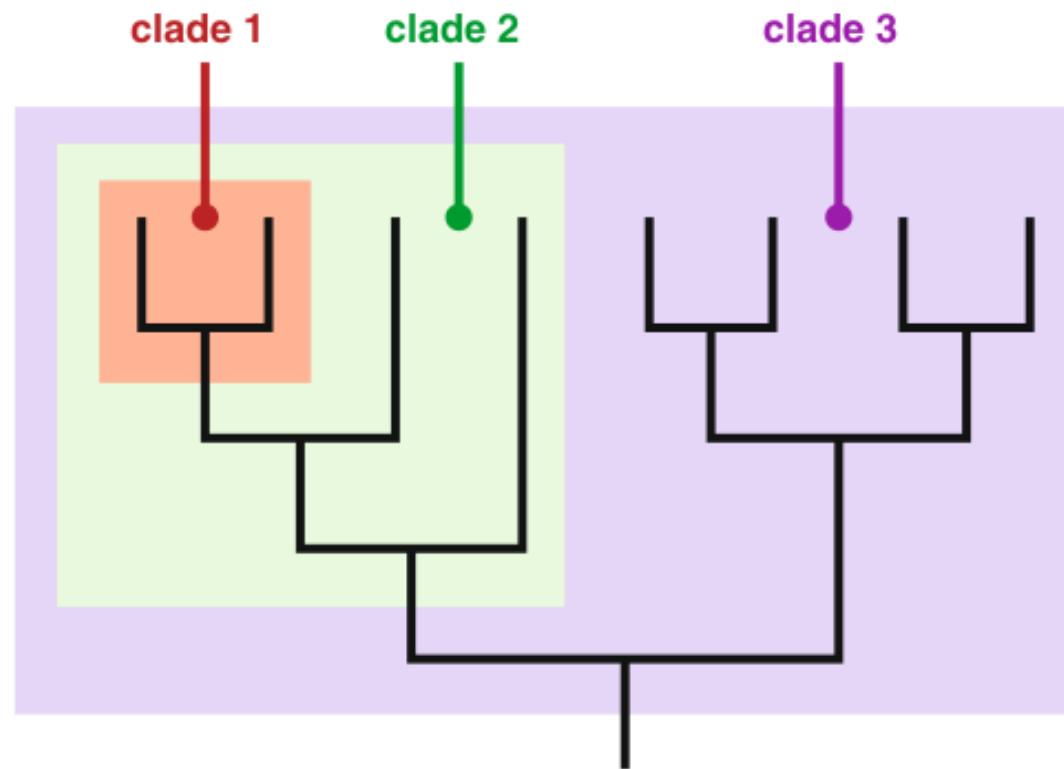


Divergencia continua  
cada linaje tiene una historia  
propia y una historia  
compartida.

# Introducción

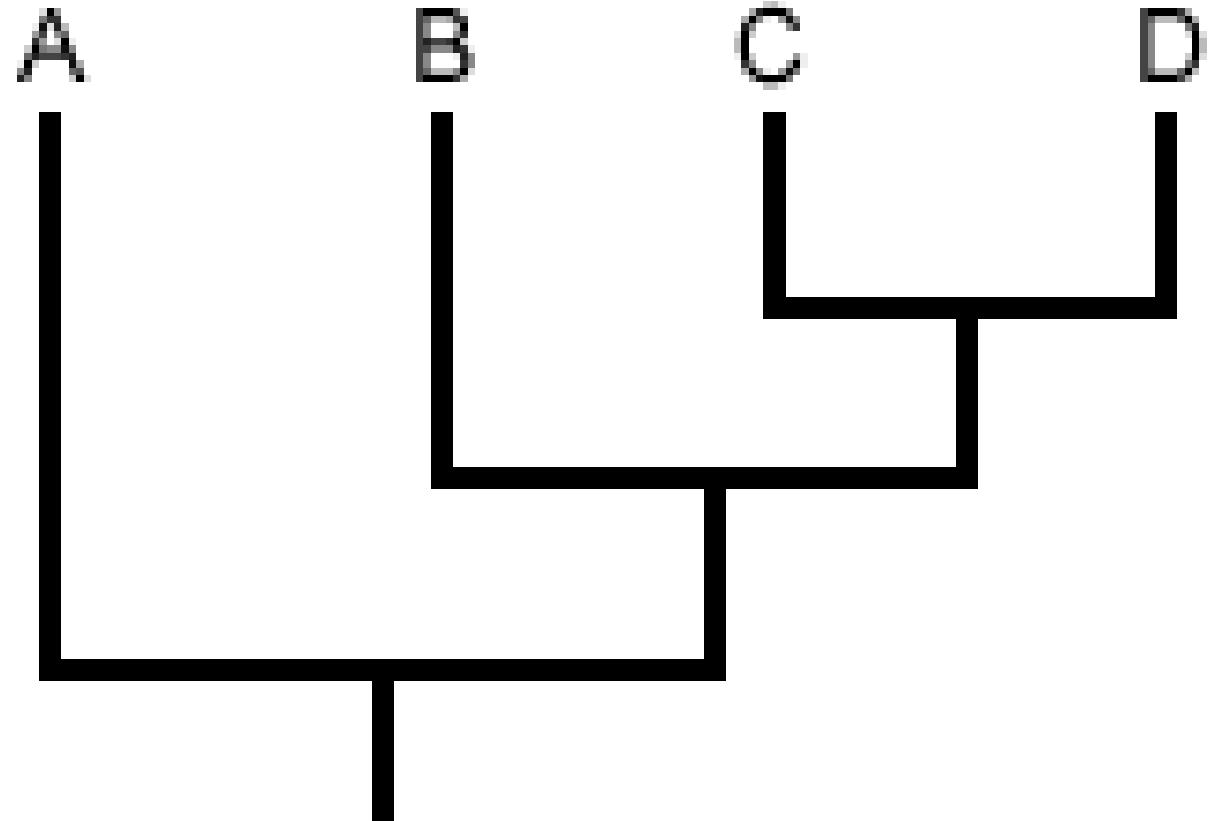
Clado. Grupo de organismos que incluye un ancestro común y a todos sus descendientes.

Se puede identificar al realizar un único corte en una rama, que separa un grupo completo sin dejar fuera a ningún descendiente de ese ancestro.



Thanukos, 2009

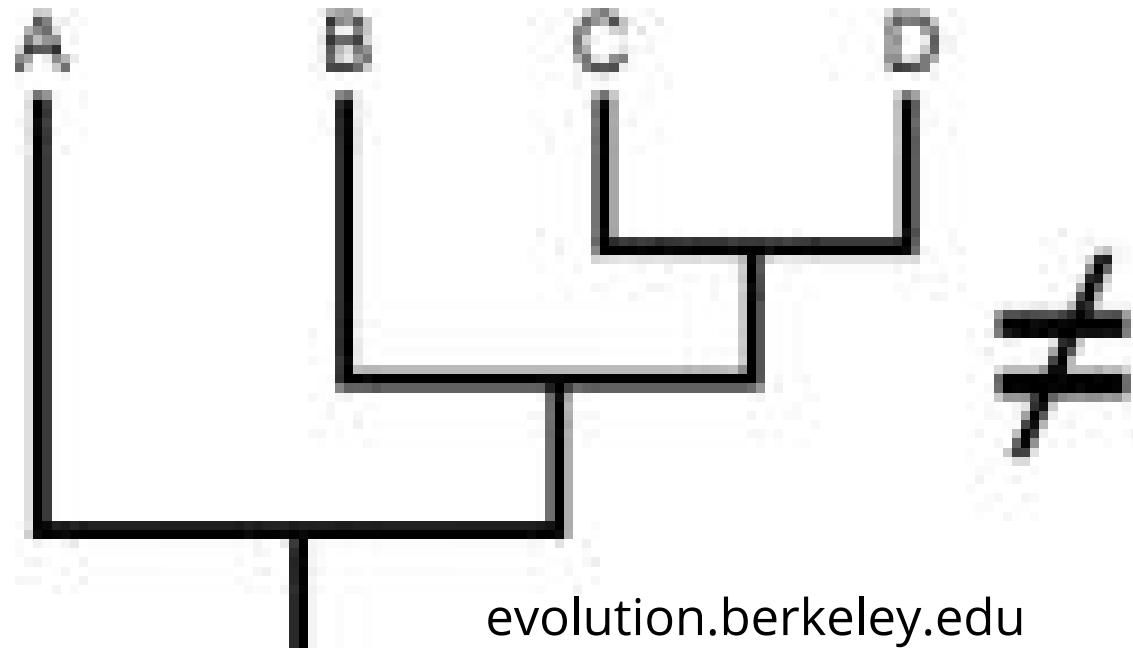
# Introducción



$\neq$



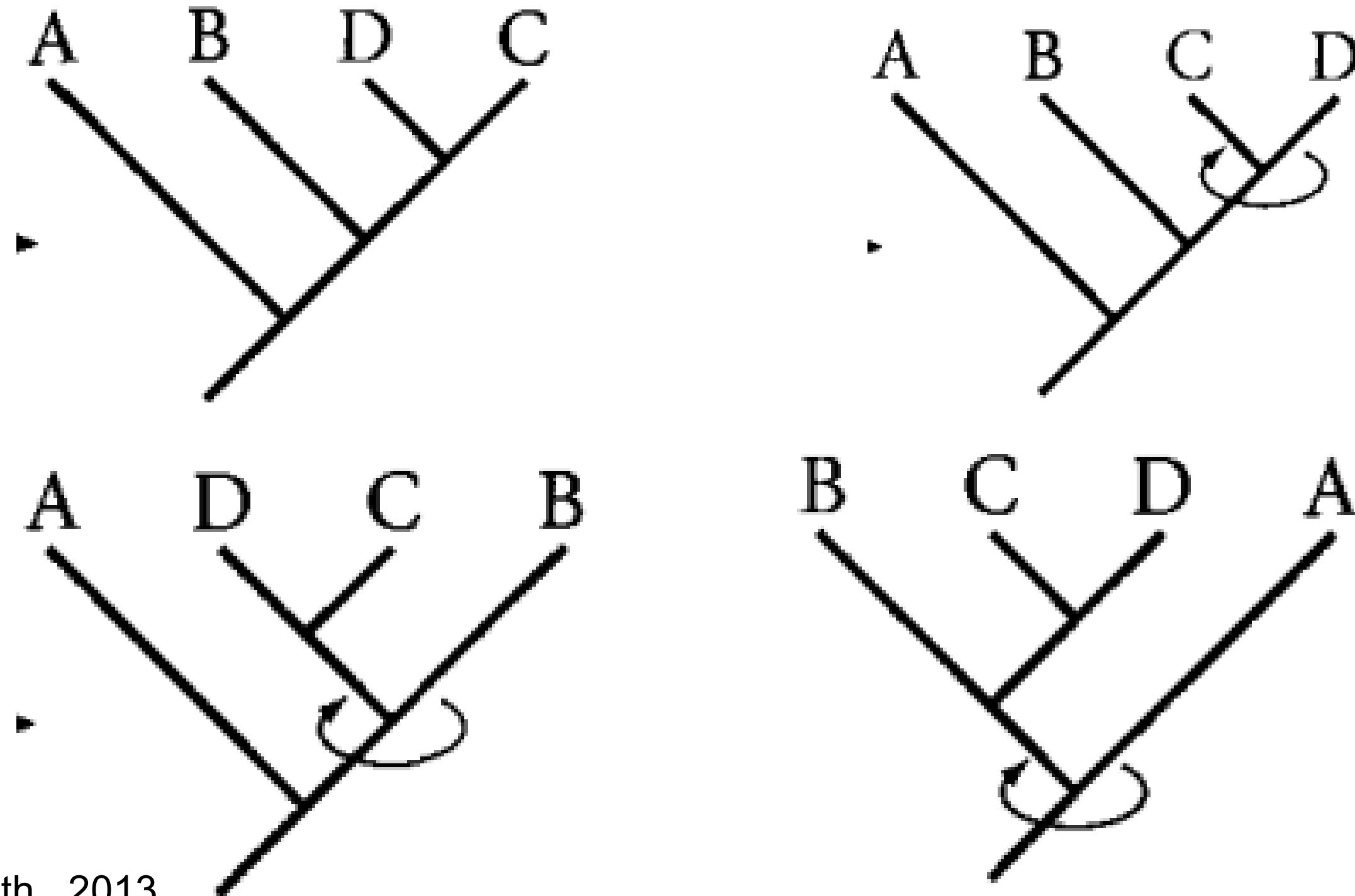
Orden no representa una secuencia lineal de evolución



A < B < C < D

Orden de las ramas no indica superioridad o progreso

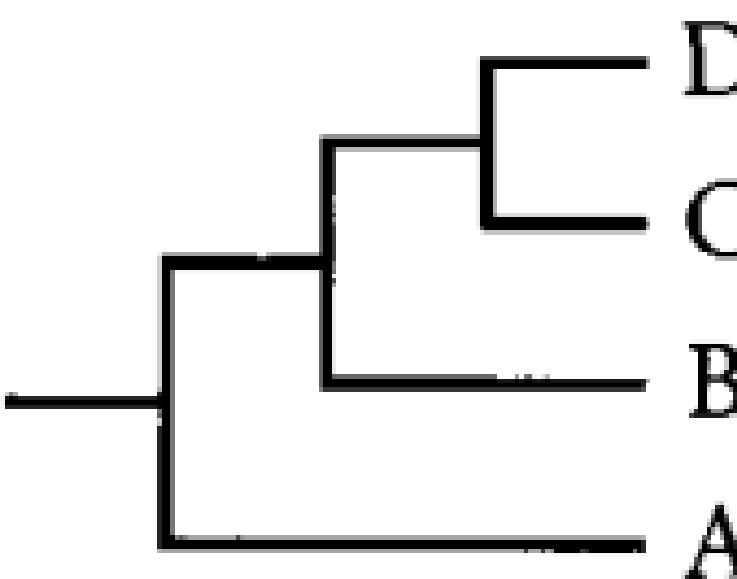
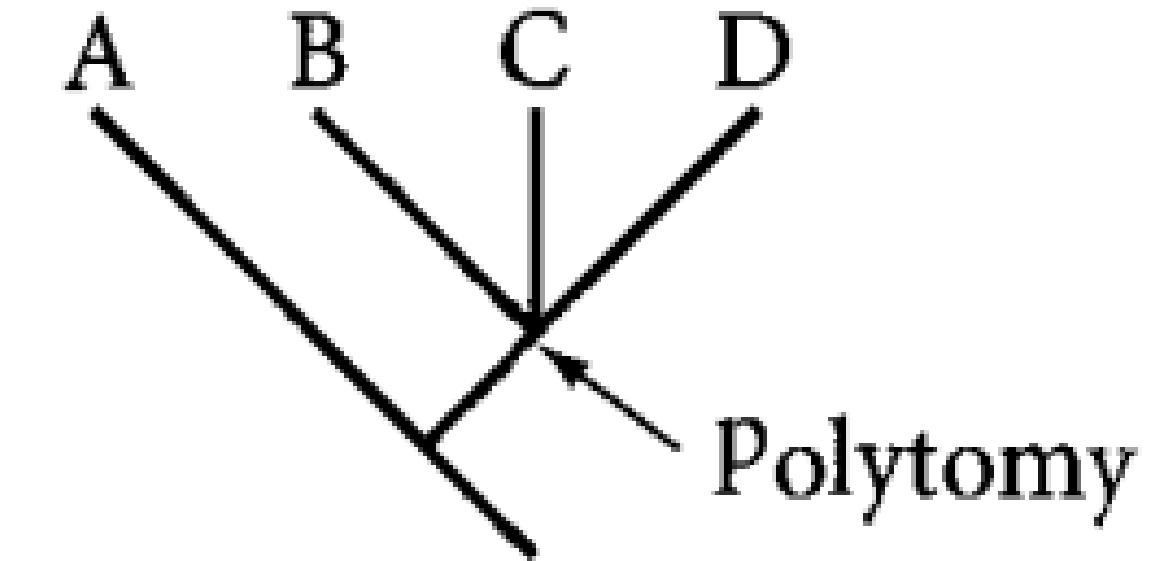
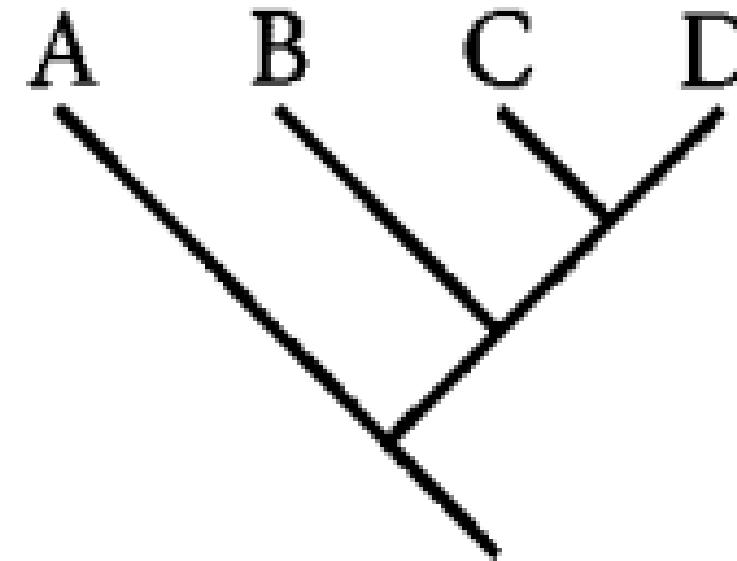
## Introducción



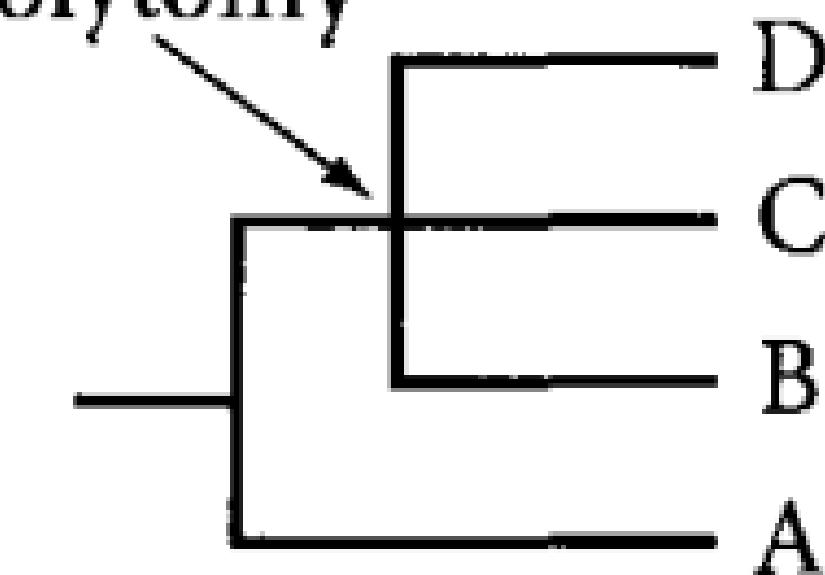
Baum & Smith, 2013

Las ramas pueden rotarse alrededor de los nodos sin alterar las relaciones evolutivas. El orden visual no importa; sino cómo los linajes se conectan.

## Introducción

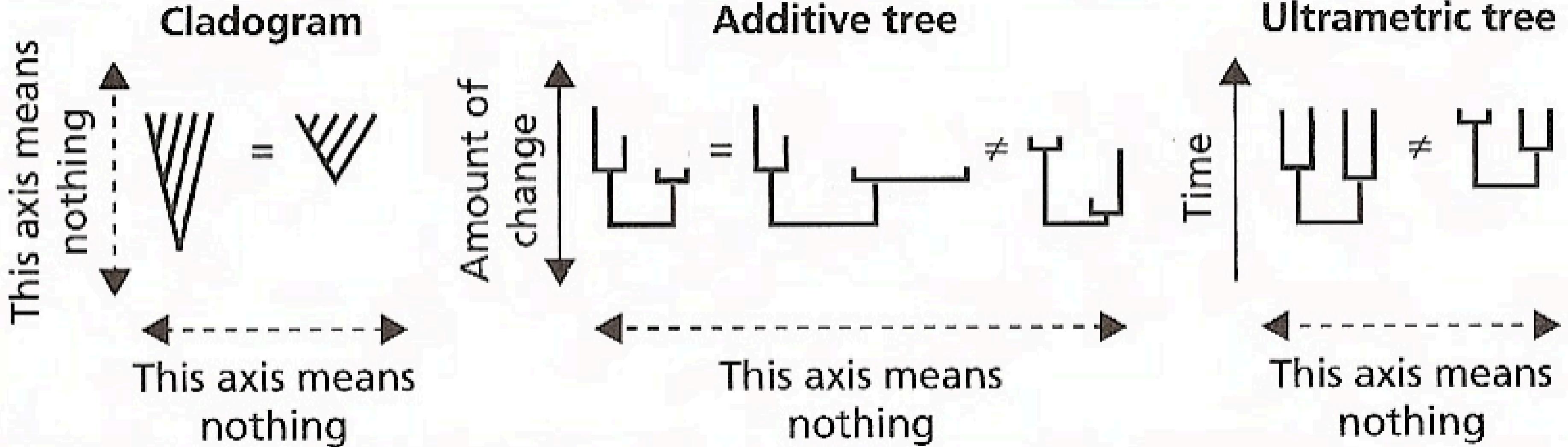


Baum & Smith, 2013



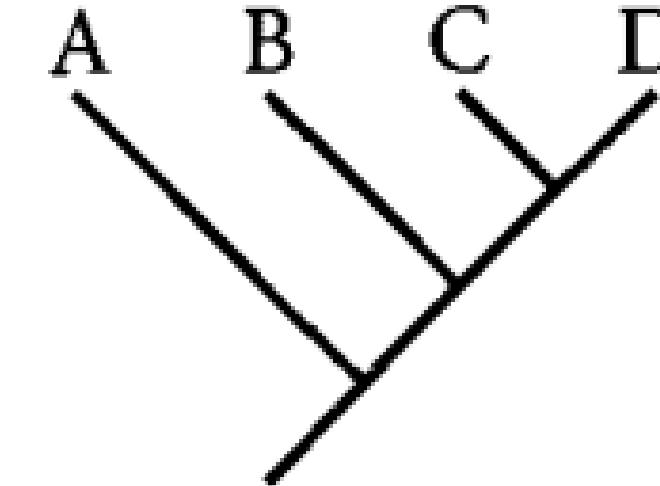
Politomías reflejan eventos evolutivos complejos o falta de información suficiente para resolver las relaciones exactas entre los linajes.

# Introducción

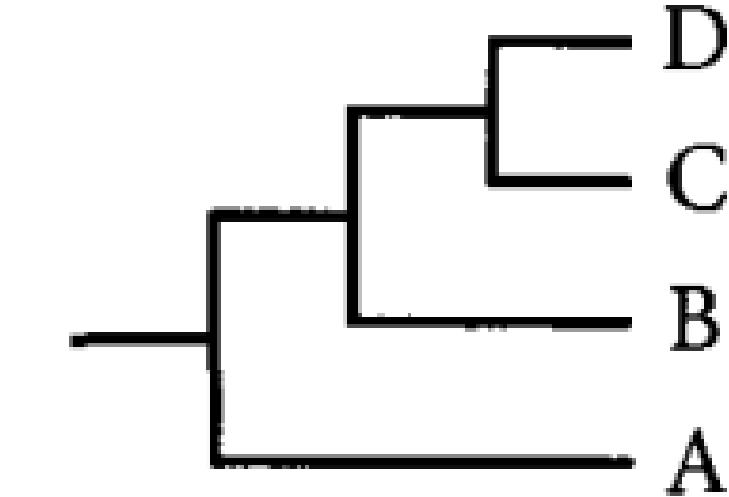


Interpretar los ejes y las longitudes de las ramas es clave

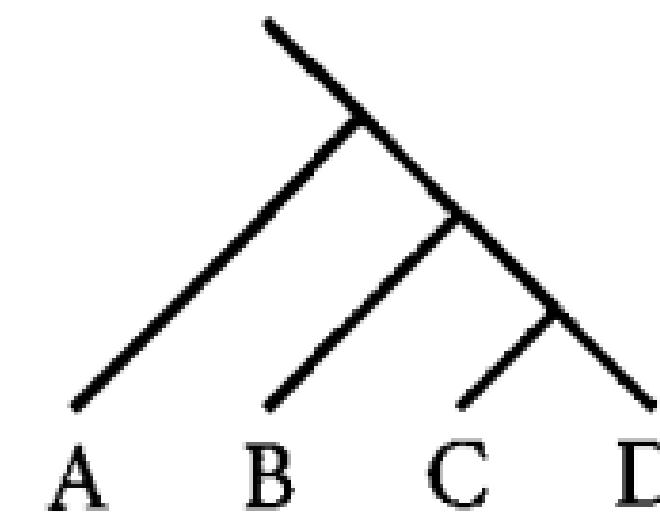
## Introducción



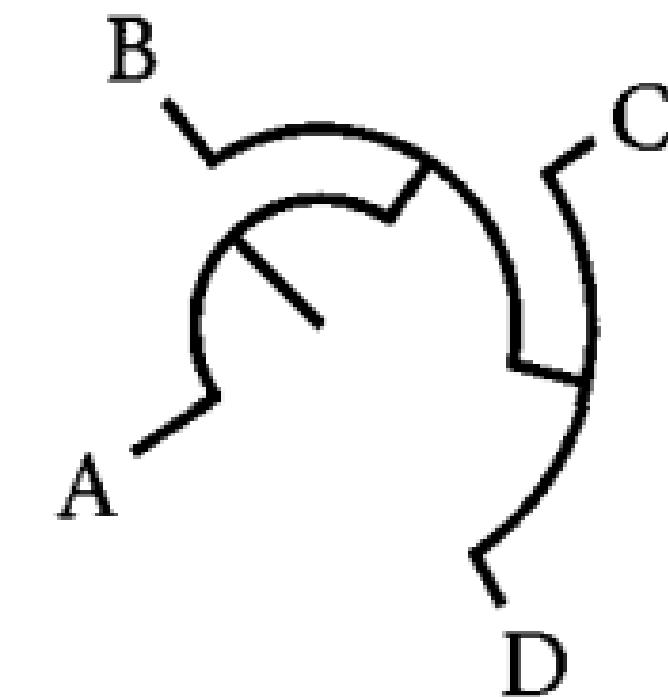
Diagonal-up



Rectangular-right



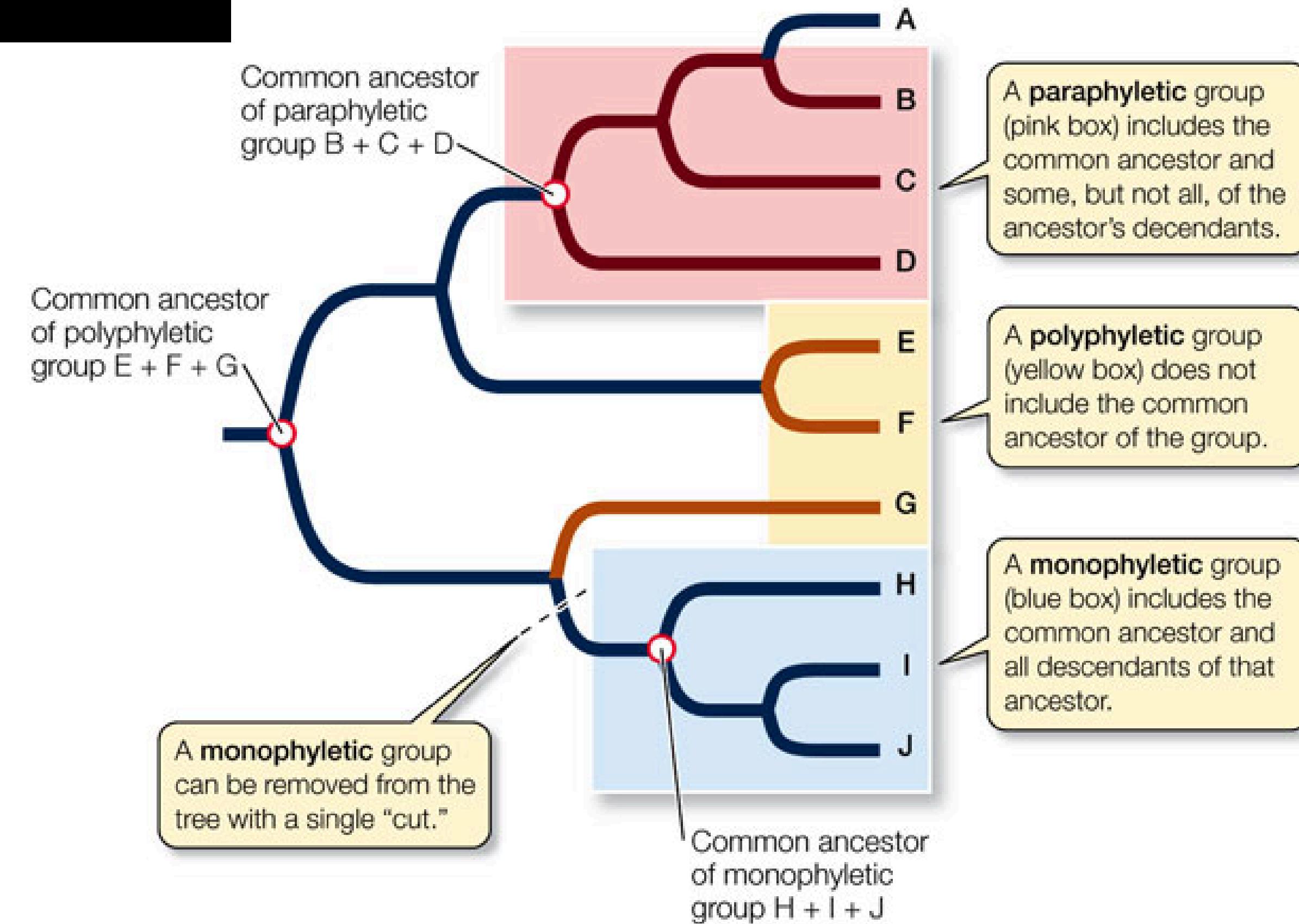
Diagonal-down  
Baum & Smith, 2013



Circle

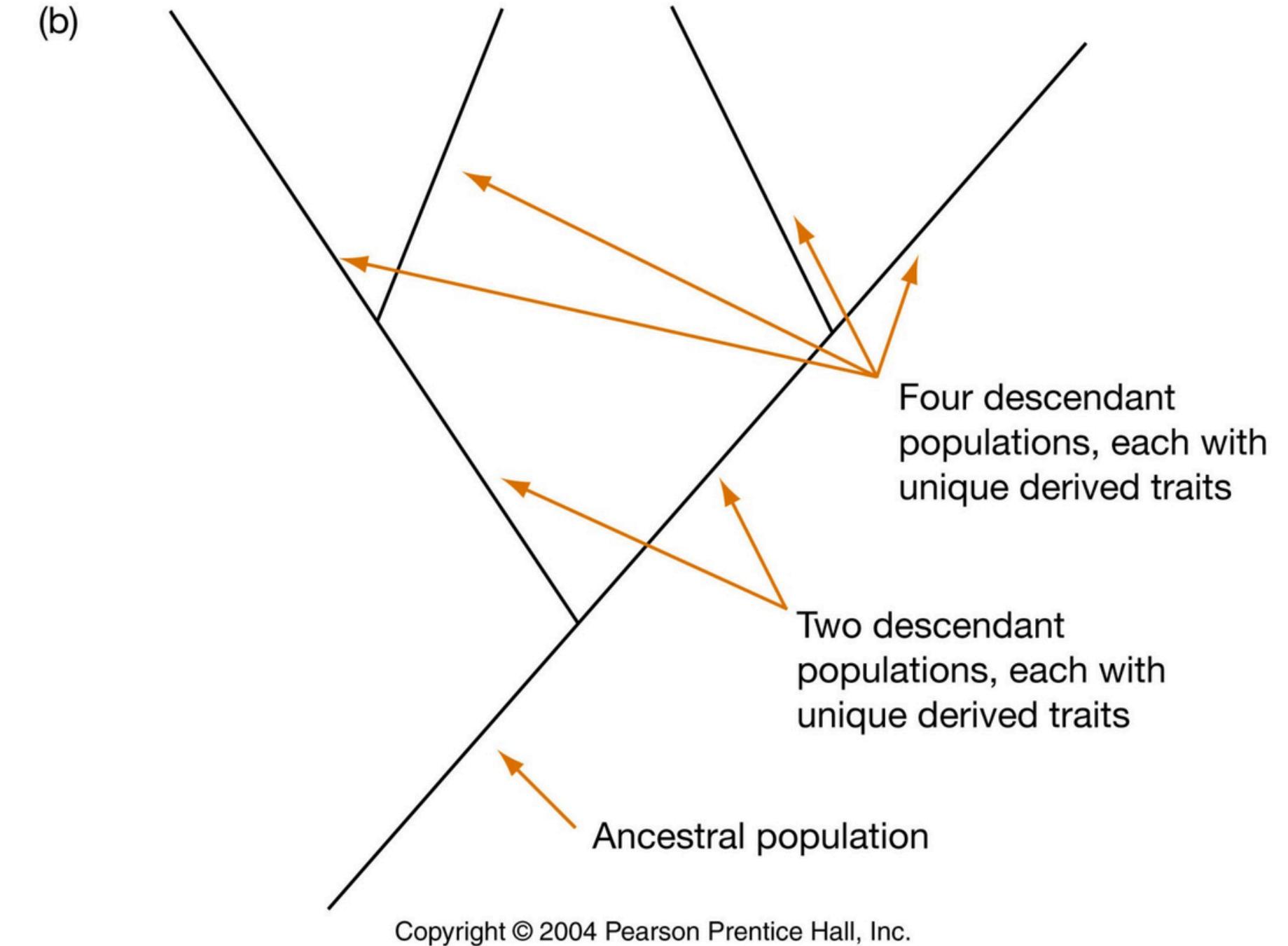
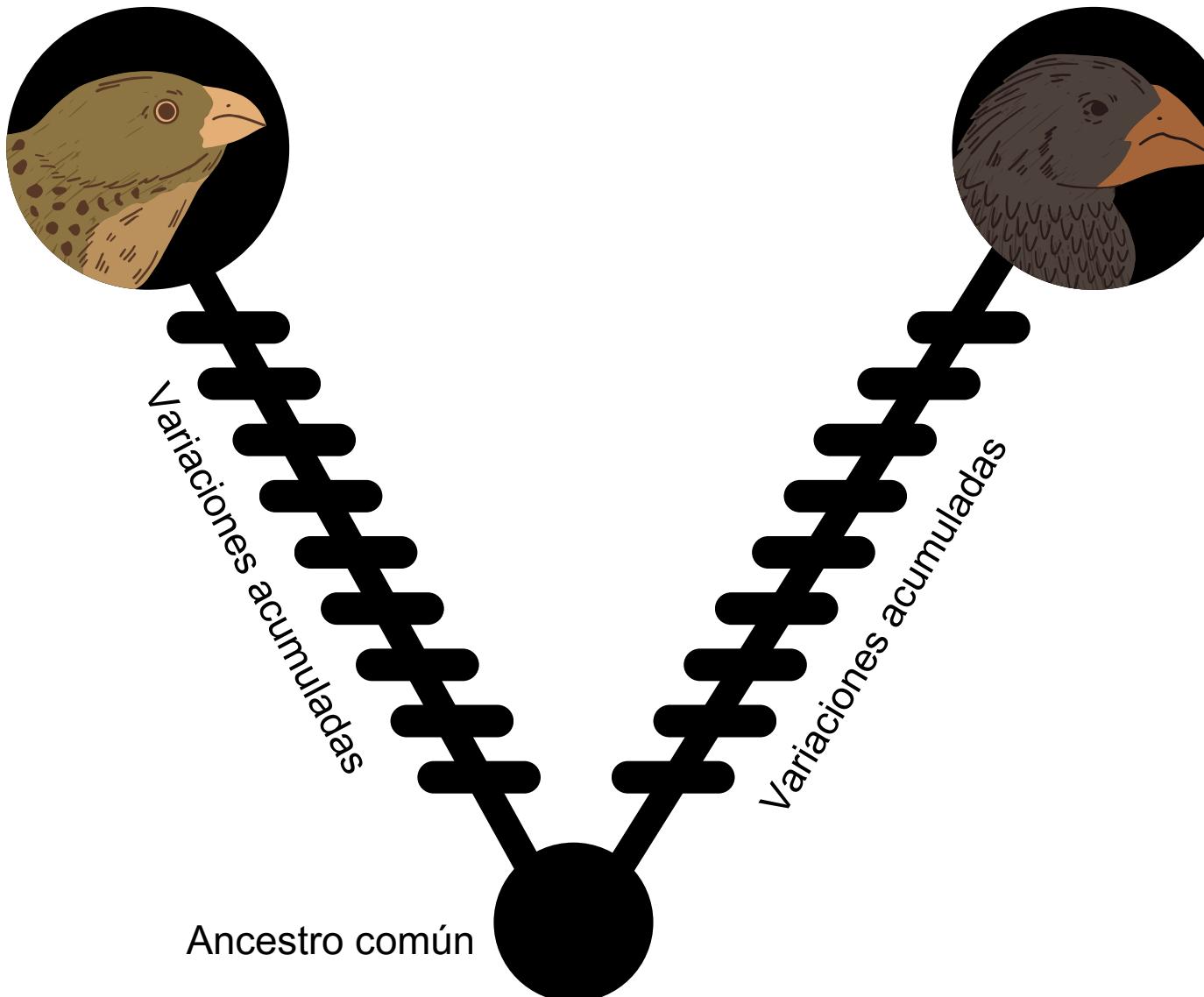
Representación visual cambia, relaciones evolutivas permanecen iguales.

# Introducción



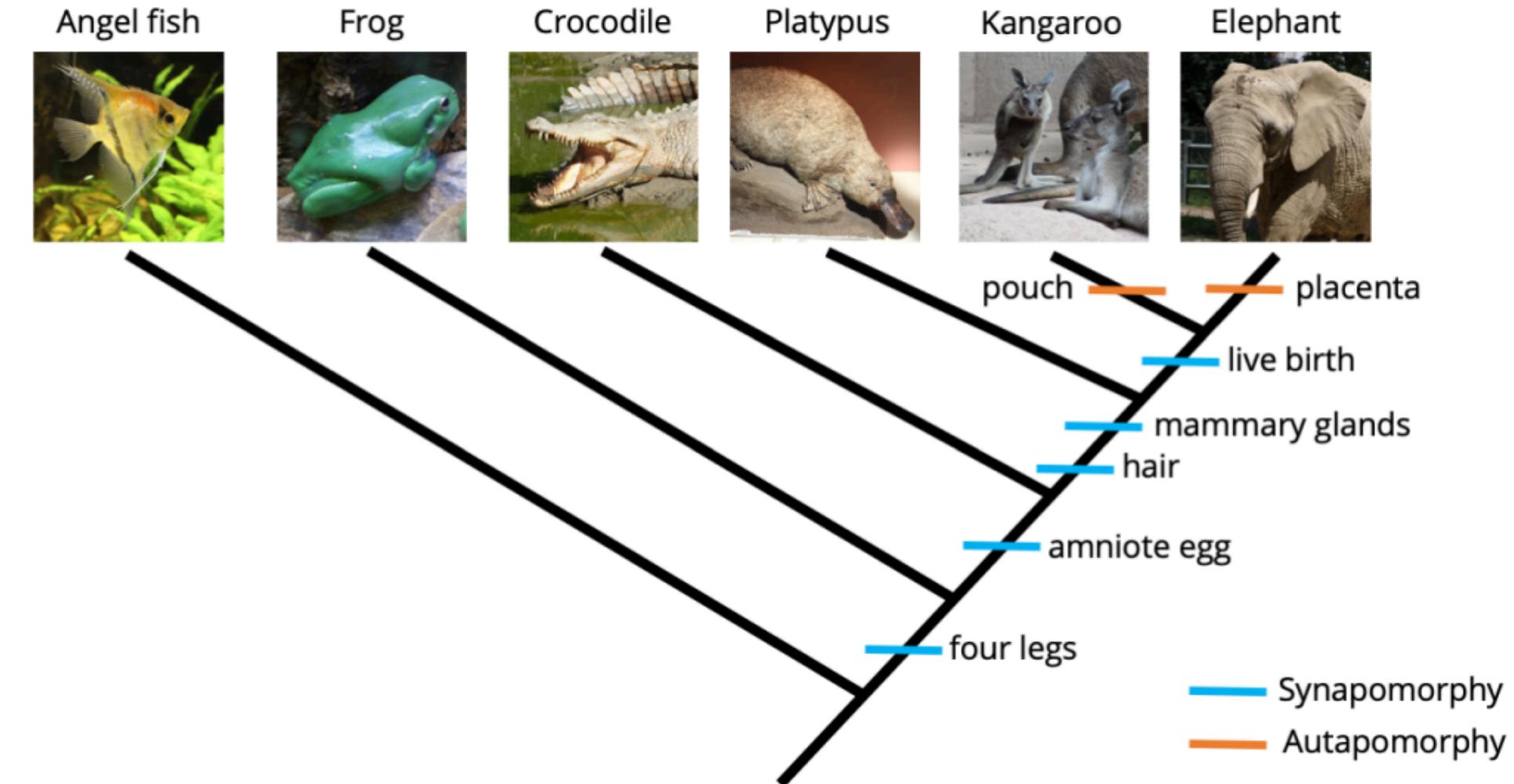
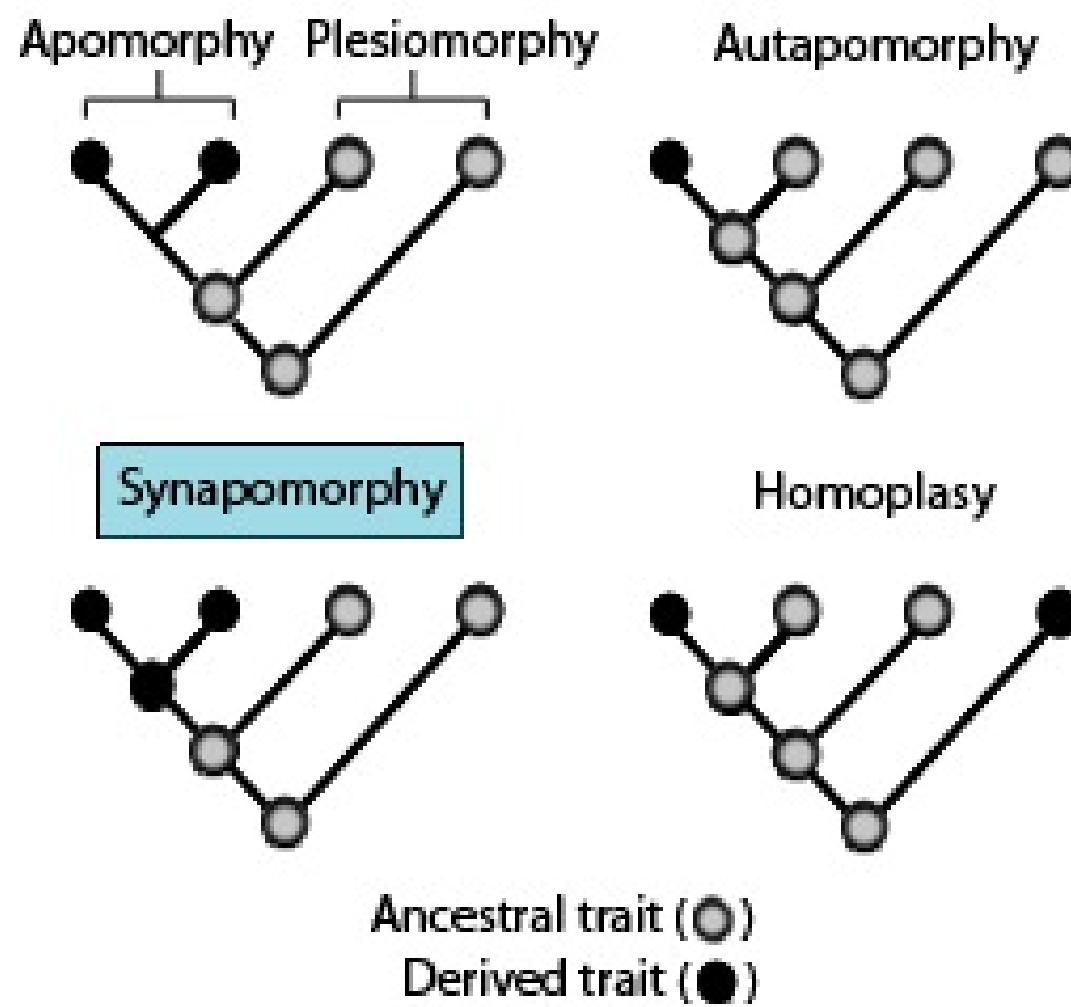
LIFE 8e, Figure 25.12

# Introducción



Los linajes resultantes del proceso de diversificación acumulan rasgos derivados únicos.

# Introducción



**Apomorfía:** Carácter derivado (nuevo) en un linaje.

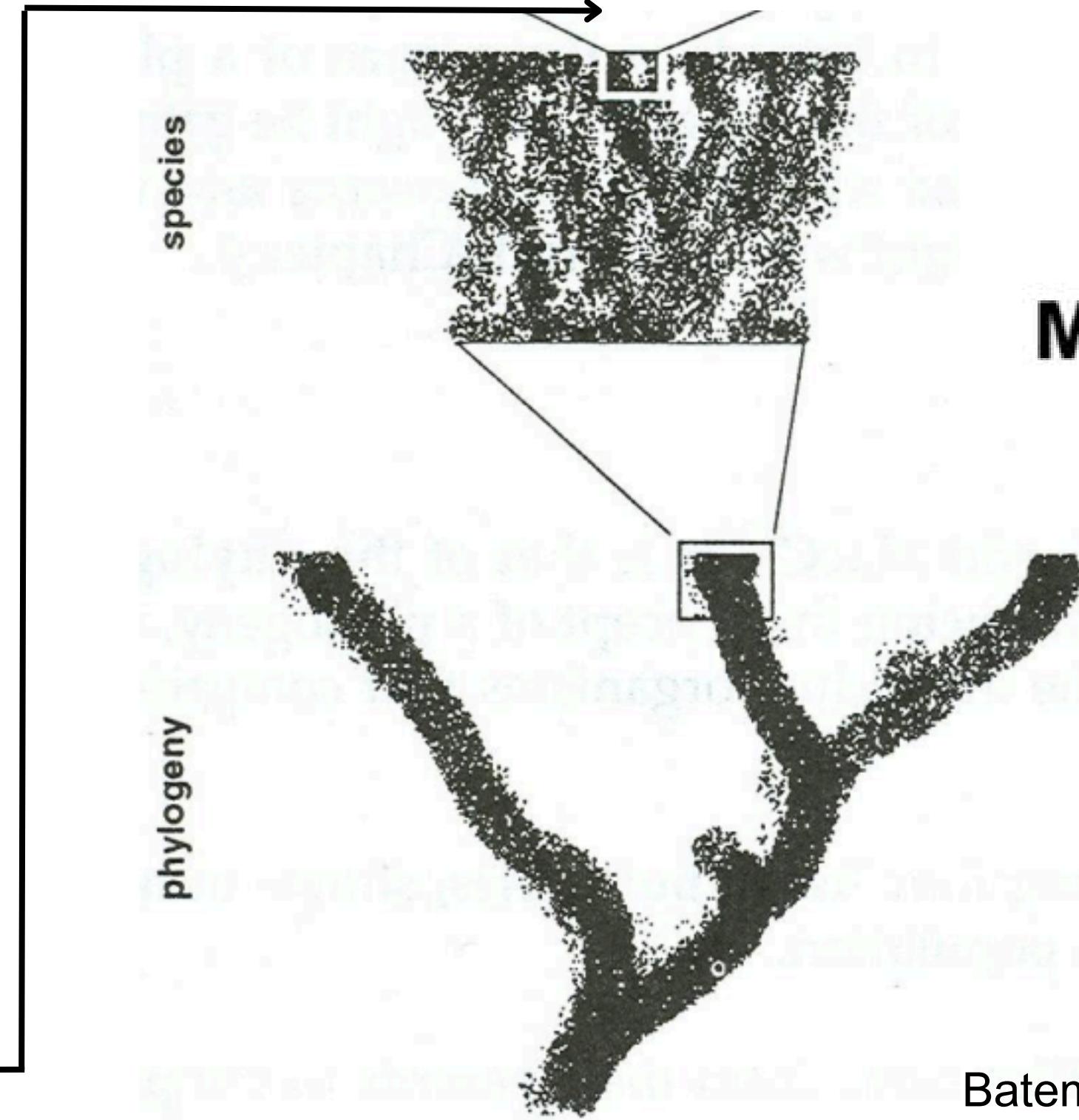
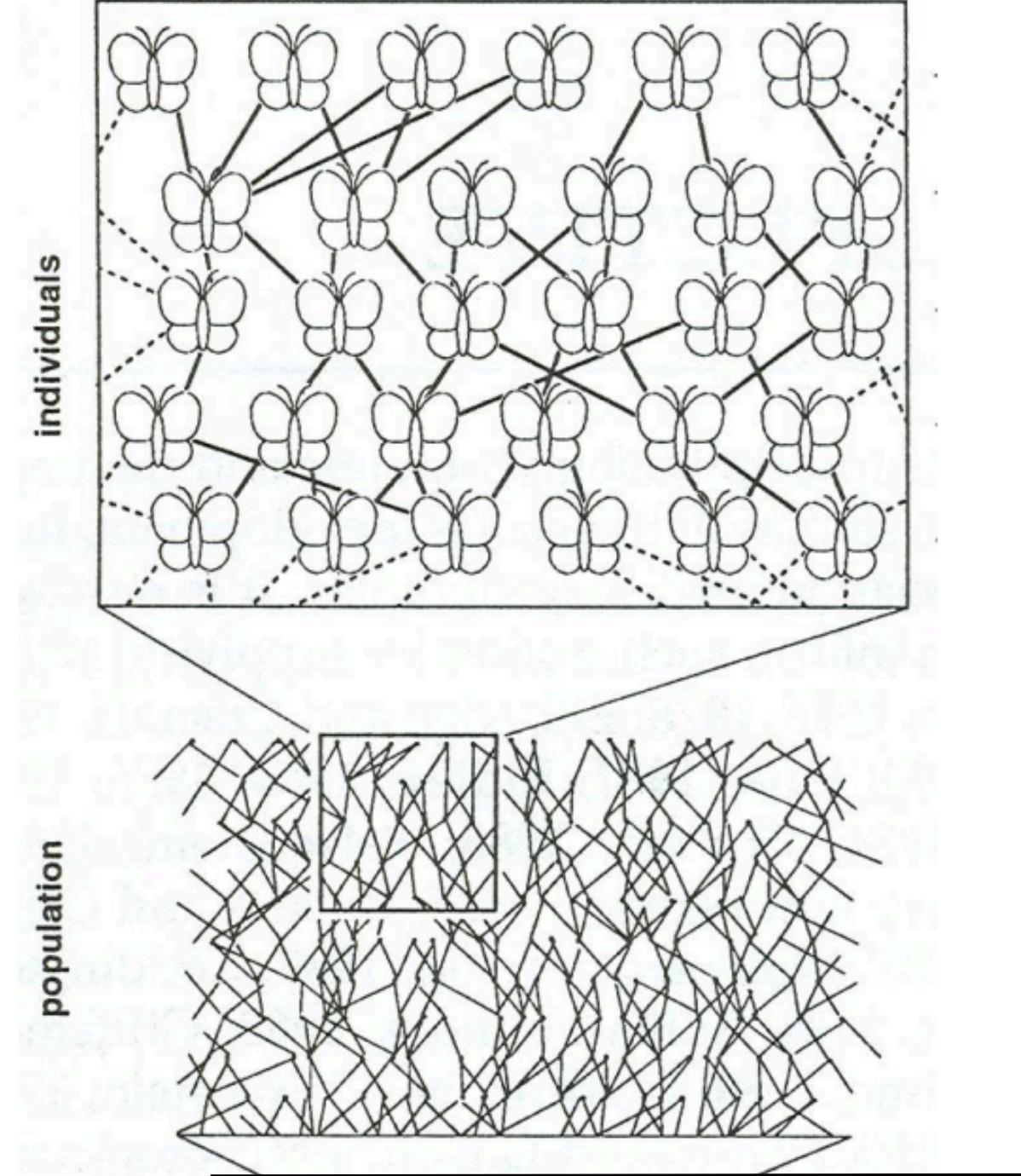
**Plesiomorfía:** Carácter ancestral presente en el ancestro común.

**Sinapomorfía:** Apomorfía compartida por varios linajes, indica un ancestro común.

**Autapomorfía:** Carácter derivado único de un solo linaje.

**Homoplasia:** Aparición independiente de un carácter en diferentes linajes (Convergencia).

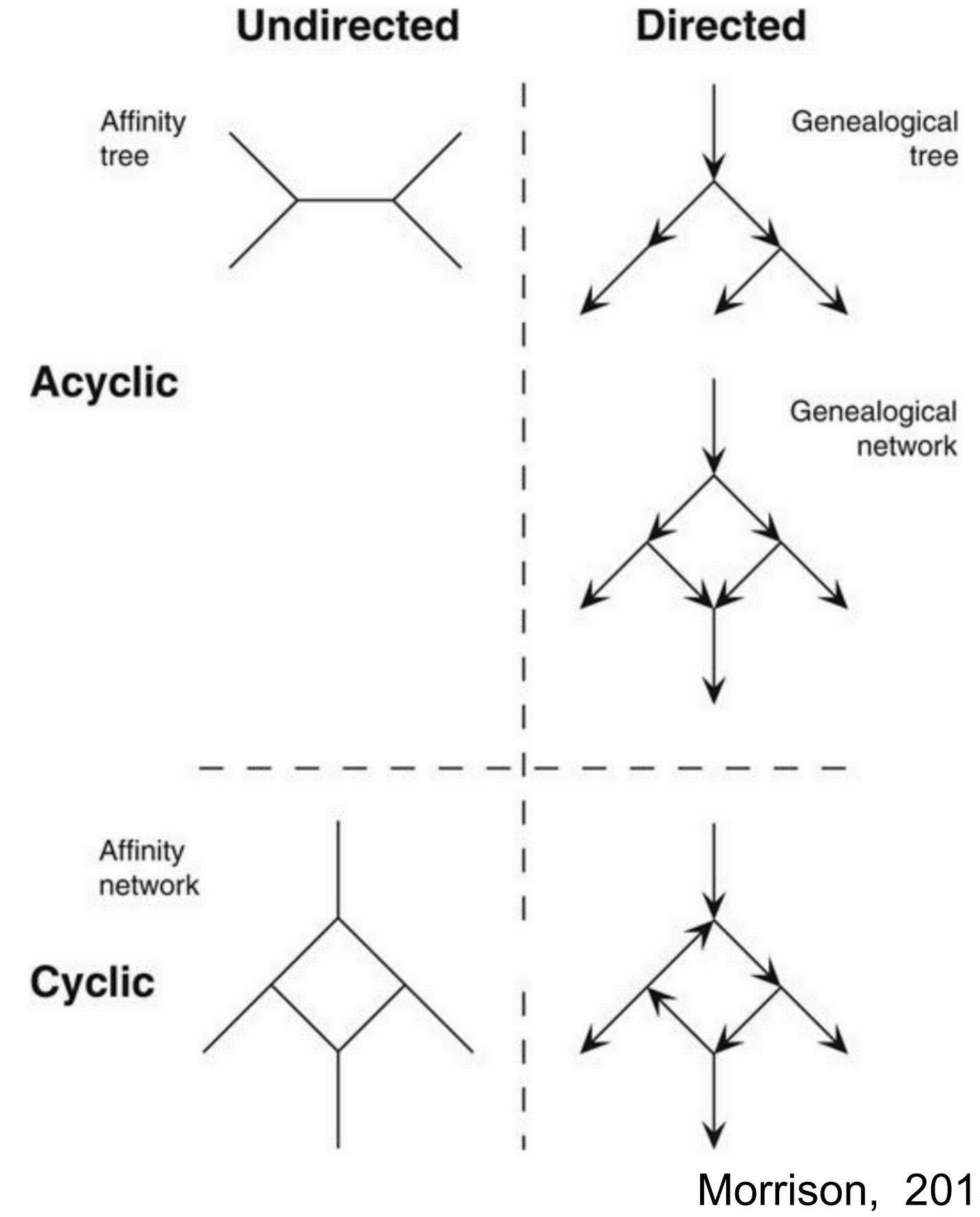
# Introducción



Bateman, 2009

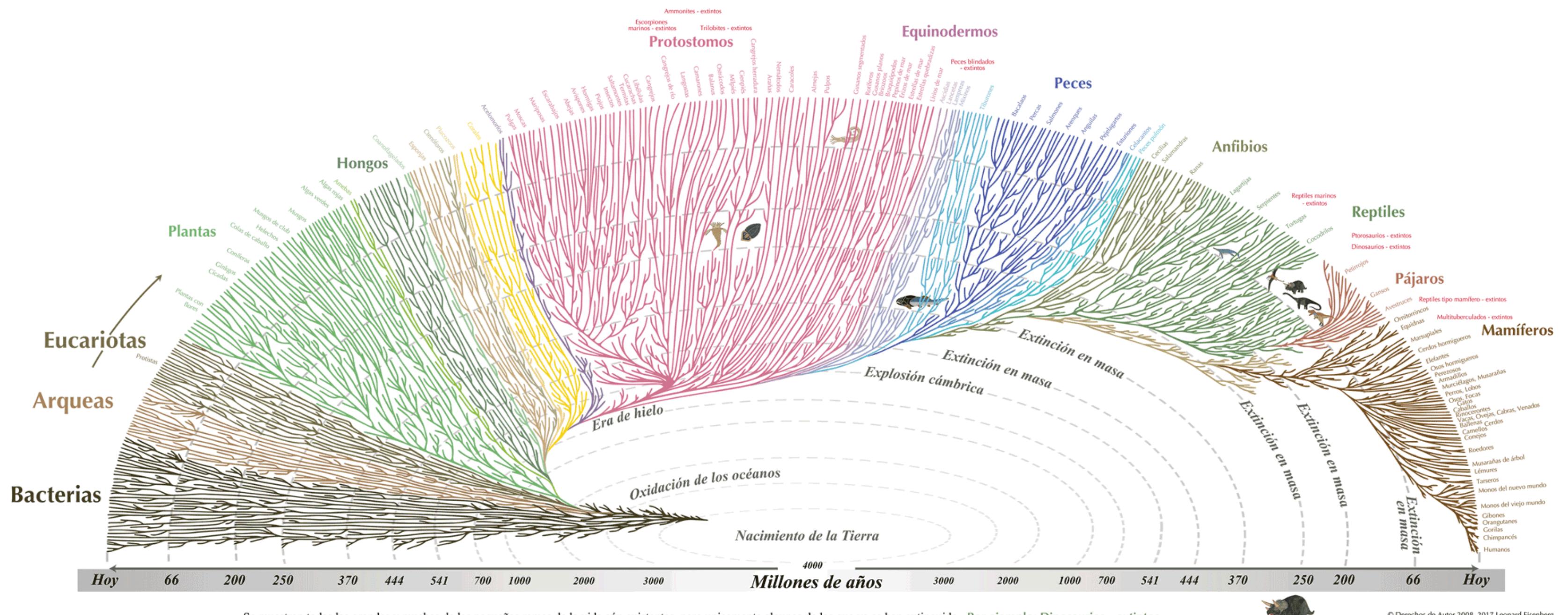
# Introducción

Procesos complejos, como el flujo genético o la hibridación, que no pueden ser representados solo con árboles bifurcados.



# Introducción

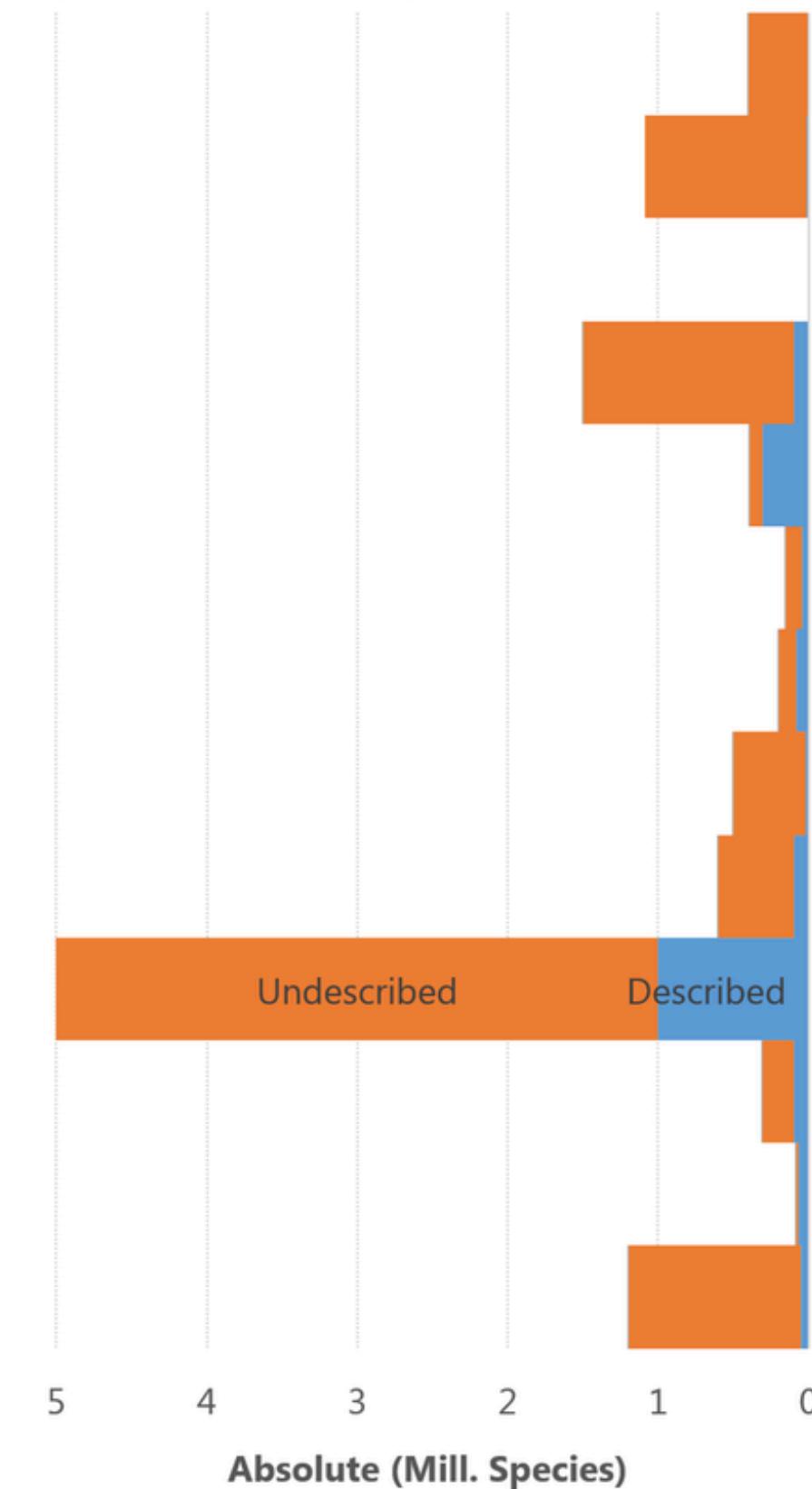
## Representación de las relaciones de ancestría-descendencia de los linajes



# Introducción

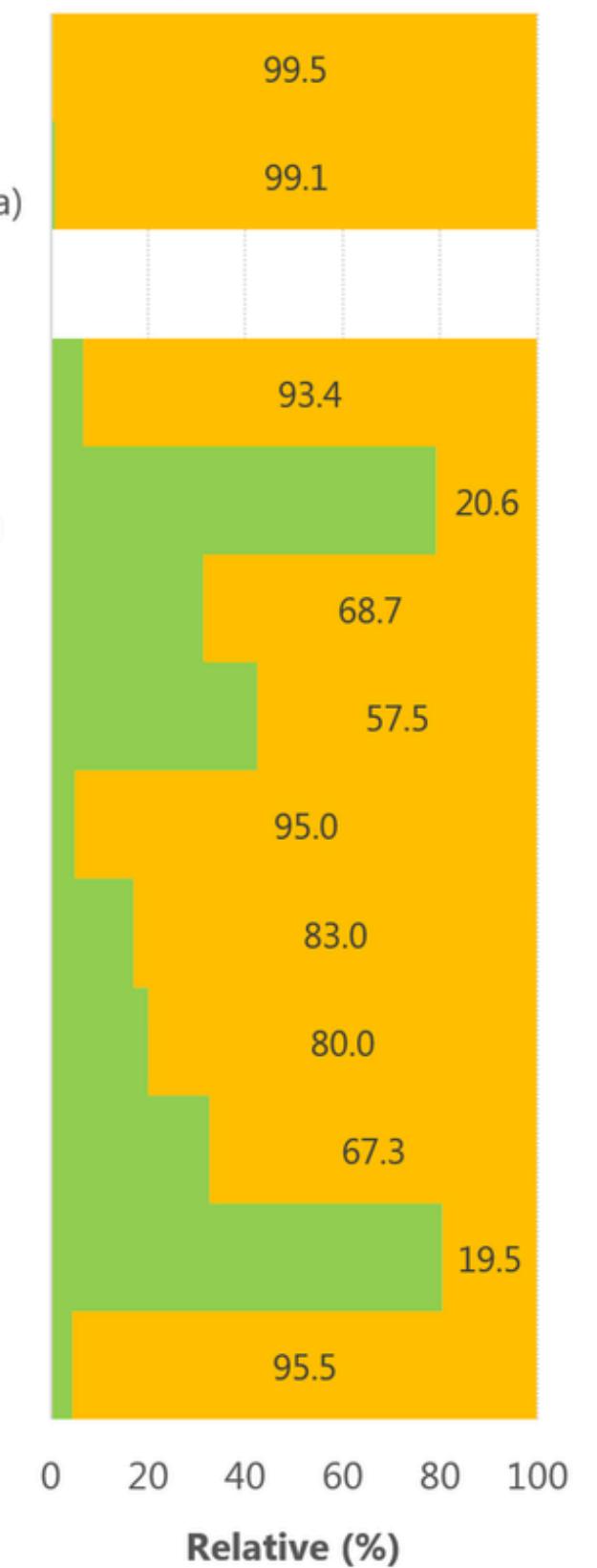


Species Richness by Taxonomic Groups



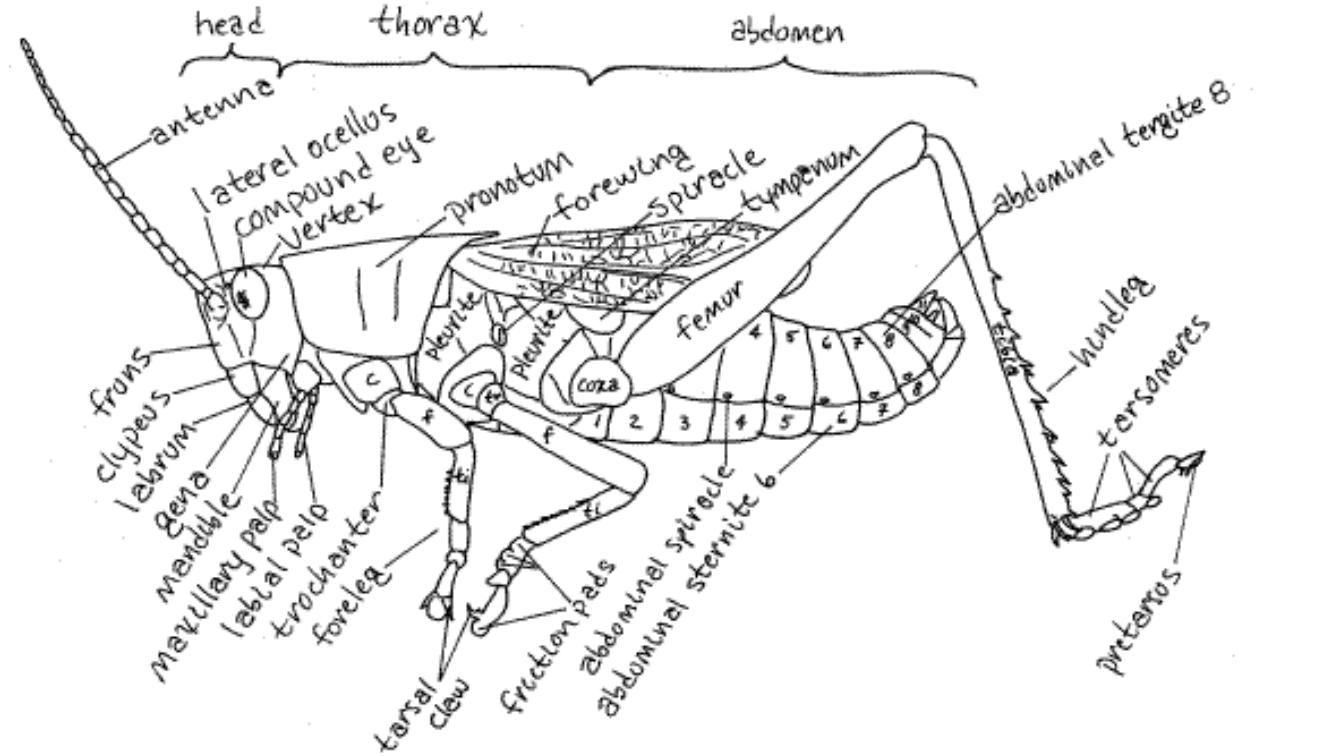
(Data: A. Chapham 2009)

Percentage Yet To Be Studied

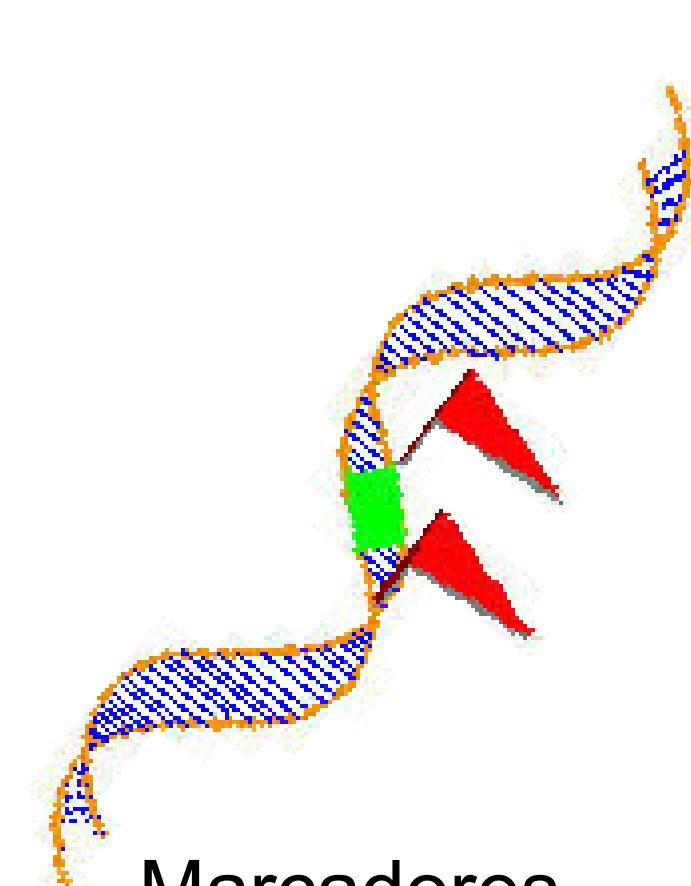


# Introducción

Carácter - Atributo que puede variar

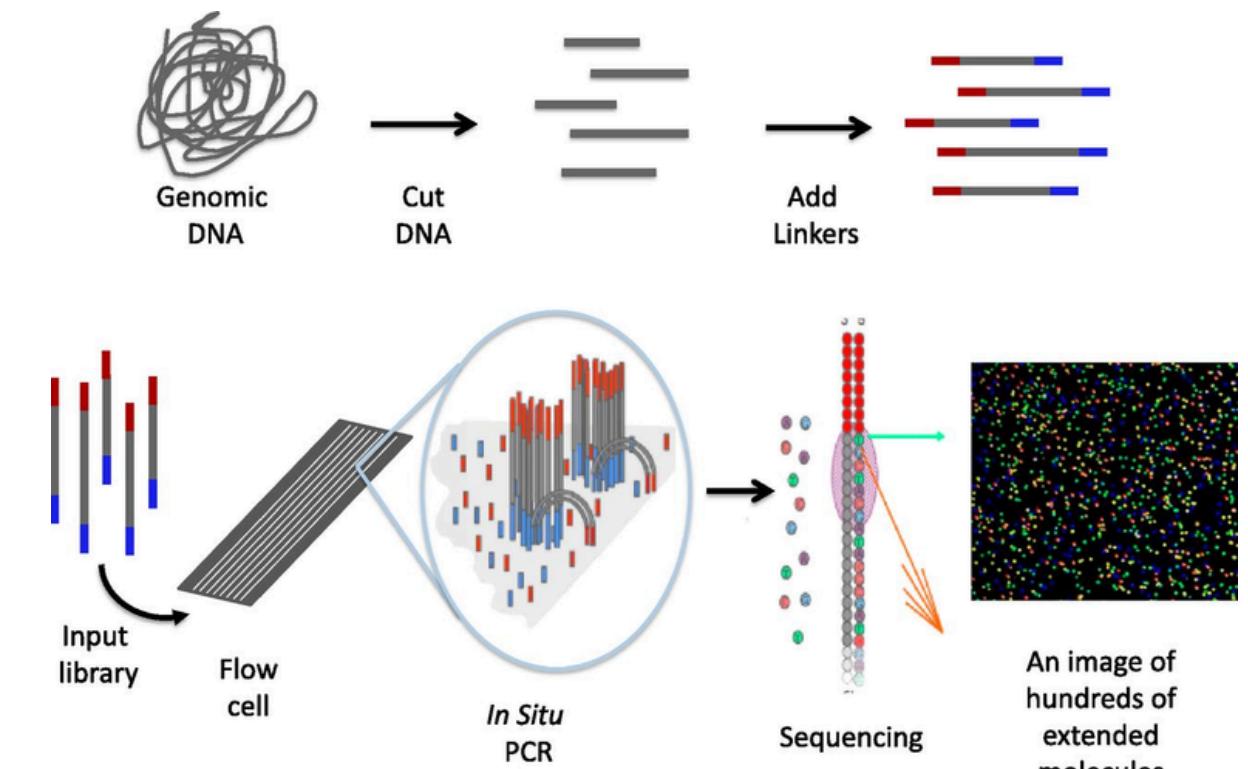


Morfología



Marcadores  
puntuales

Líneas de evidencia



Datos genómicos

# Introducción

## Marcadores puntuales

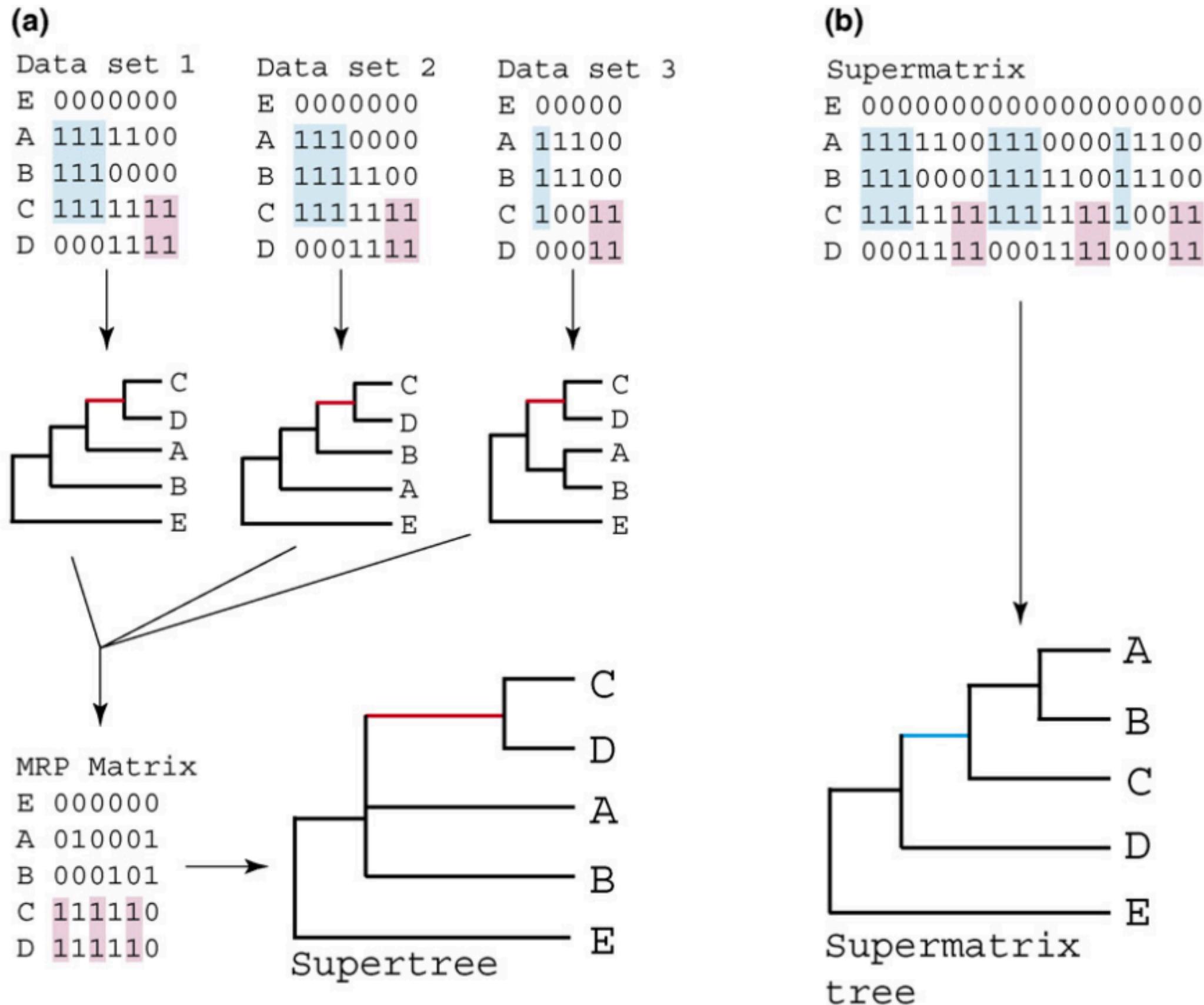
- Primers universales
- Tasas de evolución conocidas
- Métodos estandarizados
- Resolución
- Soporte
- Historia

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

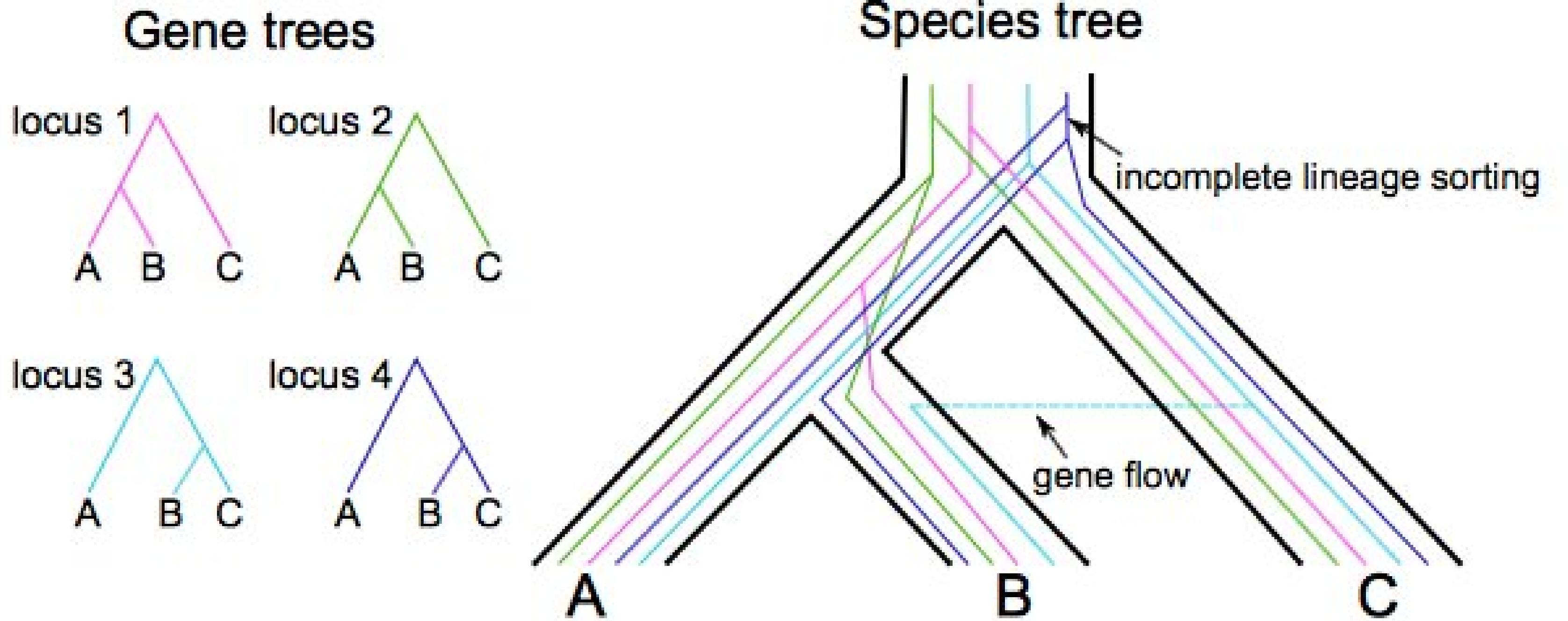
**Table 1:** List of some other molecular markers used in phylogeny research.

Patwardhan et al. 2014

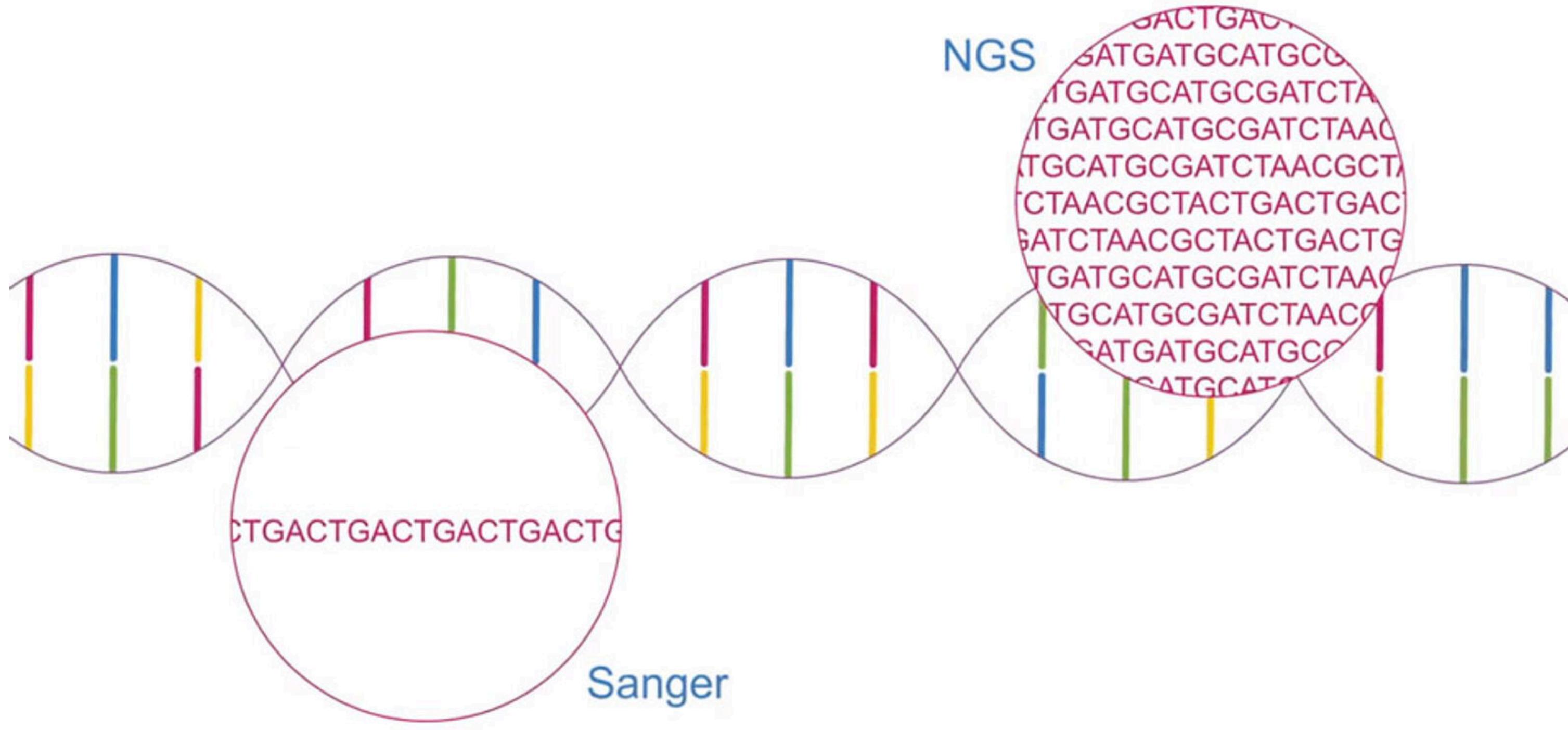
# Introducción



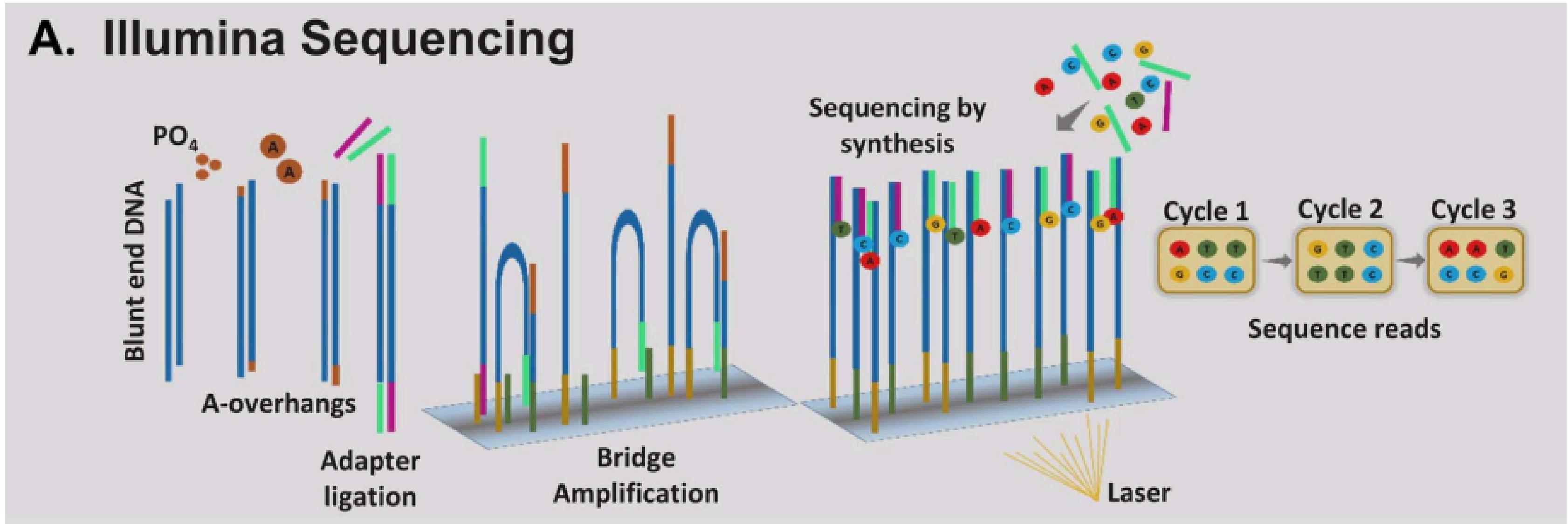
de Queiroz & Gatesy, 2006



# Introducción



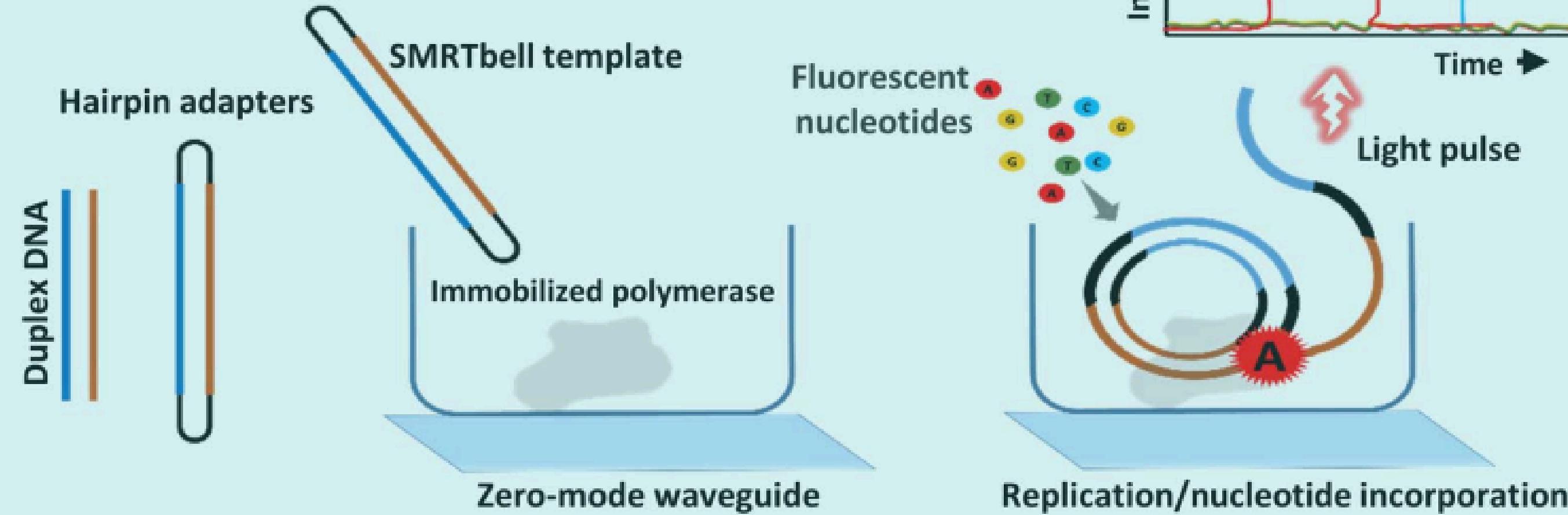
## A. Illumina Sequencing



- Lecturas Cortas: 50-300 pb
- Alta Precisión
- Millones de lecturas por corrida
- Bajo costo por base y tiempos cortos

Bharti & Grimm, 2019

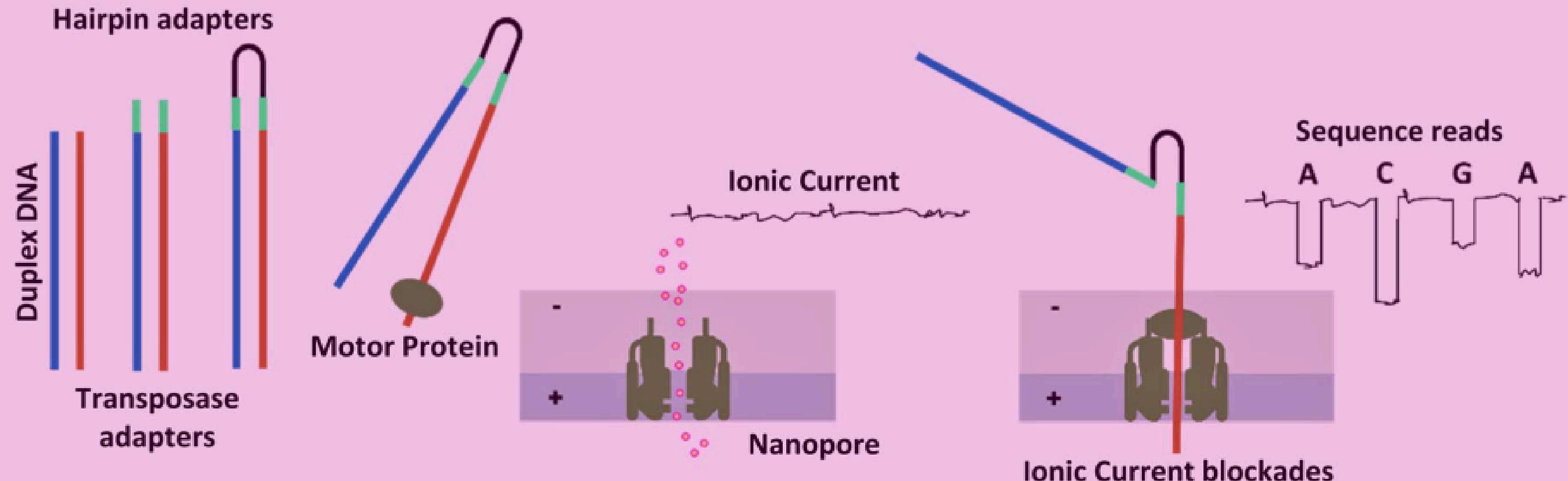
## B. PacBio Sequencing



- Lecturas largas: 10,000-100,000 pb.
- Secuenciación en tiempo real y detección epigenética.
- Errores inicial (10-15%) a >99.9% con consenso.
- Necesita más ADN y es más costosa

Bharti & Grimm, 2019

## C. Nanopore Sequencing



- Lecturas ultra largas: Hasta >1 Mb.
- Detecta cambios en la corriente eléctrica mientras las bases de ADN/ARN pasan por un poro.
- Portátil y escalable
- Error alto (~5-15%) que otras tecnologías

Bharti & Grimm, 2019

# Introducción

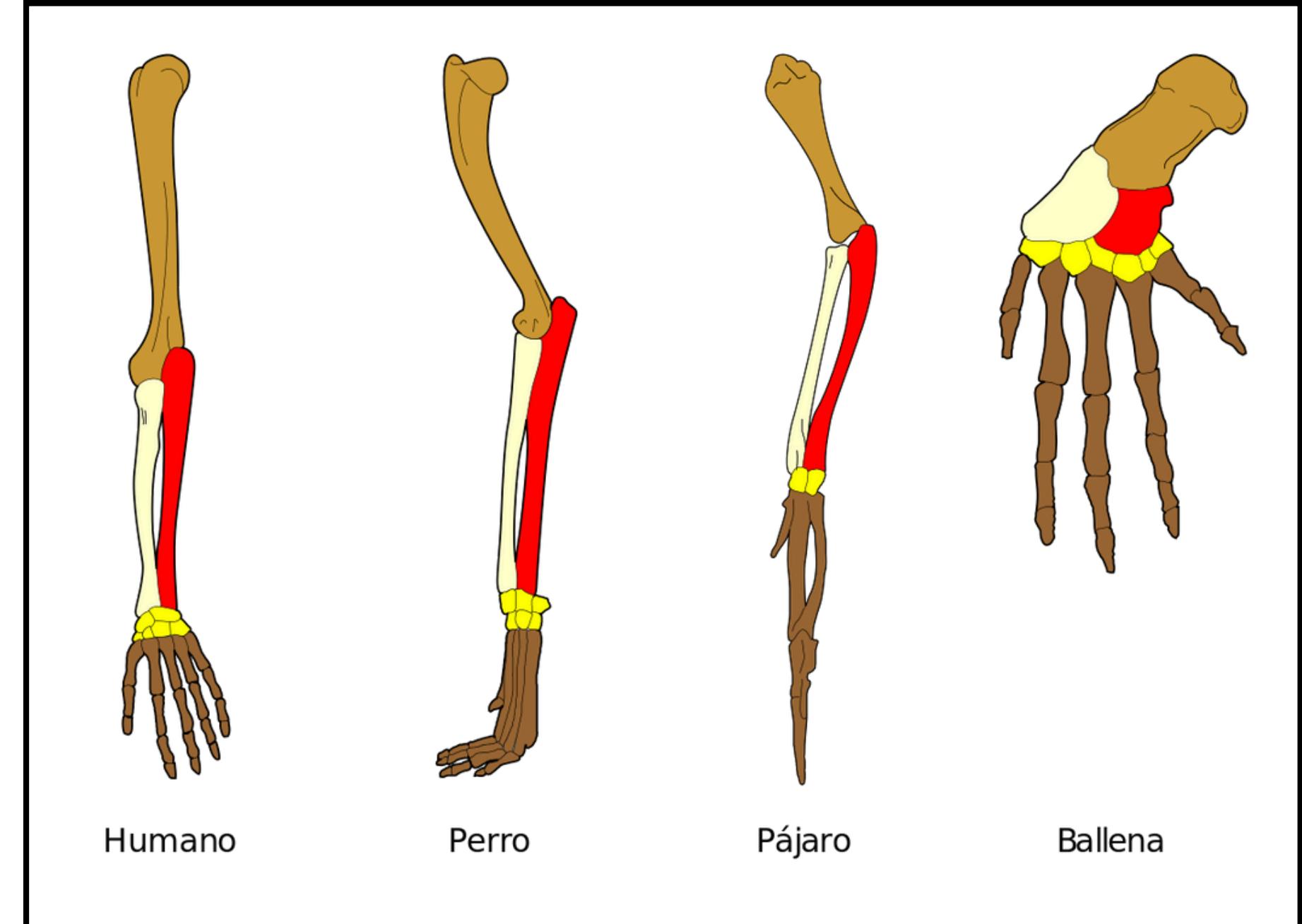
**Table 1.** Summary of the main advantages of long-reads sequencing over short-read sequencing.

Short-Read Technologies	Long-Read Technologies
<p>Fixed run time:</p> <ul style="list-style-type: none"><li>- Increased time to results and inability to identify workflow errors before completed sequencing</li><li>- Additional practical complexities associated with handling and storing large volumes of sequence data</li></ul>	<p>Real-time data acquisition:</p> <ul style="list-style-type: none"><li>- Achieve rapid turnaround with immediate access to results</li><li>- Enrich single targets during sequencing, with no additional sample prep using adaptive sampling</li><li>- Identify microbiome composition and resistance in real-time</li></ul>
<p>Limited flexibility:</p> <ul style="list-style-type: none"><li>- Sample batching often required for optimal efficiency</li><li>- Potentially leads to long turnaround times</li><li>- Benchtop devices confine sequencing to centralized locations</li></ul>	<p>Scalable and flexible:</p> <ul style="list-style-type: none"><li>- Scale to suit the throughput needs</li><li>- Decentralize sequencing</li><li>- No sample batching needed</li></ul>
Read length typically 50–300 bp	Unrestricted read length (>4 Mb achieved)
<p>Limited genomic characterization:</p> <ul style="list-style-type: none"><li>- Short reads do not span entire structural variants or important classes of genomic aberrations (repeat expansions and repeat-rich regions)</li><li>- Fragmented genome assemblies and ambiguous isoforms identification</li><li>- Short sequencing reads may not span complex genomic regions such as genes duplications, transposons and prophage sequences</li><li>- Potentially missing important genomic information</li></ul>	<p>Comprehensive genomic characterization:</p> <ul style="list-style-type: none"><li>- Identify mutations in complex and repetitive genomic regions</li><li>- Accurately phase single nucleotide variants, structural variants, and base modifications</li><li>- Can fully assemble genomes more easily</li><li>- Simplify de novo assembly and correct microbial reference genomes</li><li>- Possibility to completely assemble genomes and plasmids from metagenomic samples</li><li>- Resolving complex genomic regions and similar species</li></ul>
<p>Amplification required:</p> <ul style="list-style-type: none"><li>- Amplification can introduce bias reducing uniformity of coverage and removes base modifications</li><li>- Necessitating additional sample prep and sequencing runs</li></ul>	<p>Amplification-free protocols:</p> <ul style="list-style-type: none"><li>- Detect and phase base modifications as standard</li><li>- No additional preparation required</li></ul>
<p>Constrained to the lab:</p> <ul style="list-style-type: none"><li>- Traditional sequencing technologies are typically expensive and require substantial site infrastructure</li><li>- Usually limited its usage to well-resourced environments</li><li>- Delay in transmitting the results</li></ul>	<p>Sequence anywhere:</p> <ul style="list-style-type: none"><li>- Sequence in your lab or in the field</li><li>- Sequence at sample source and eliminate sample shipping delays</li><li>- Scale-up with high-throughput</li></ul>

# Introducción

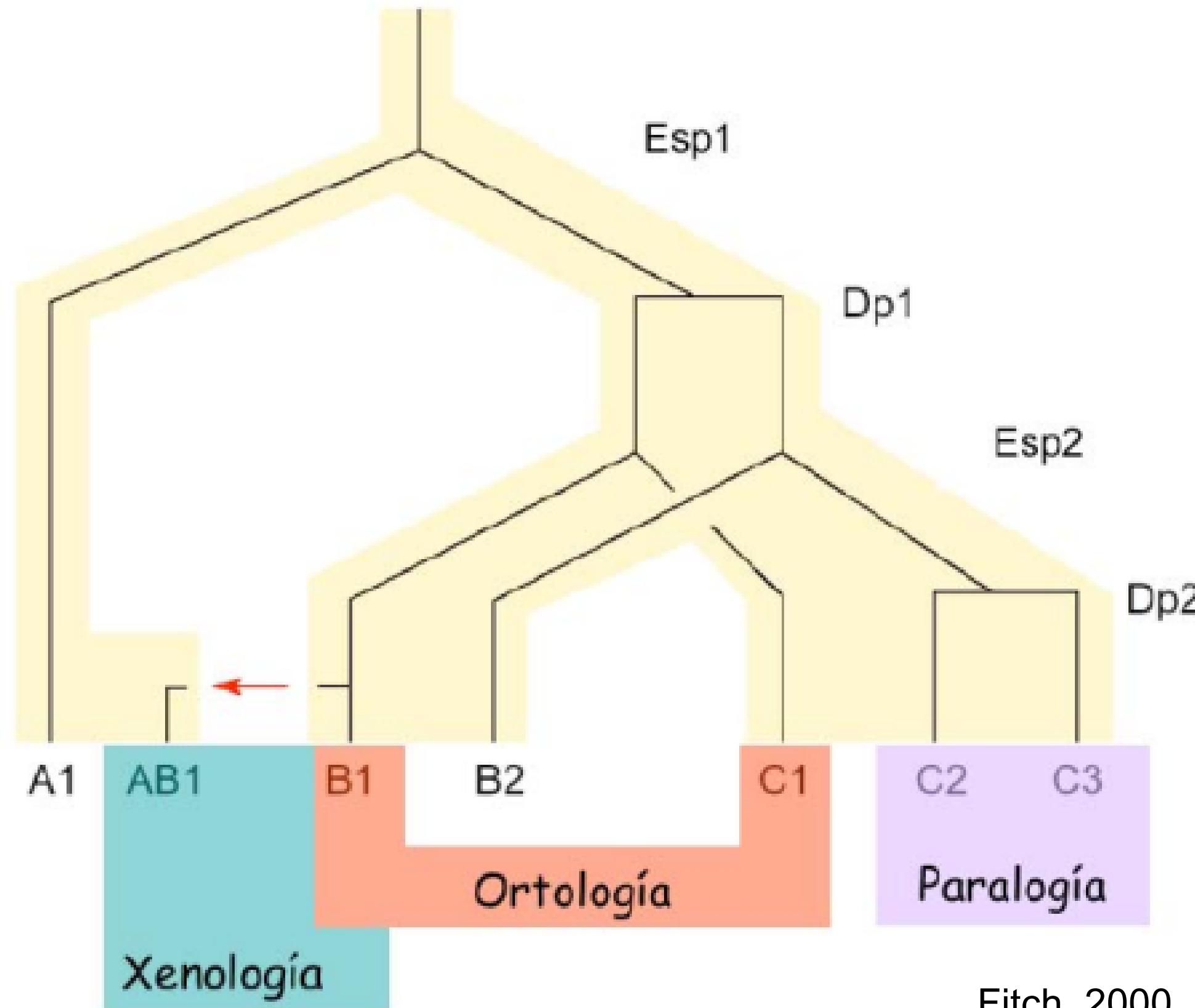


Analogía



Homología

# Introducción



## What is an Alignment?

Unaligned

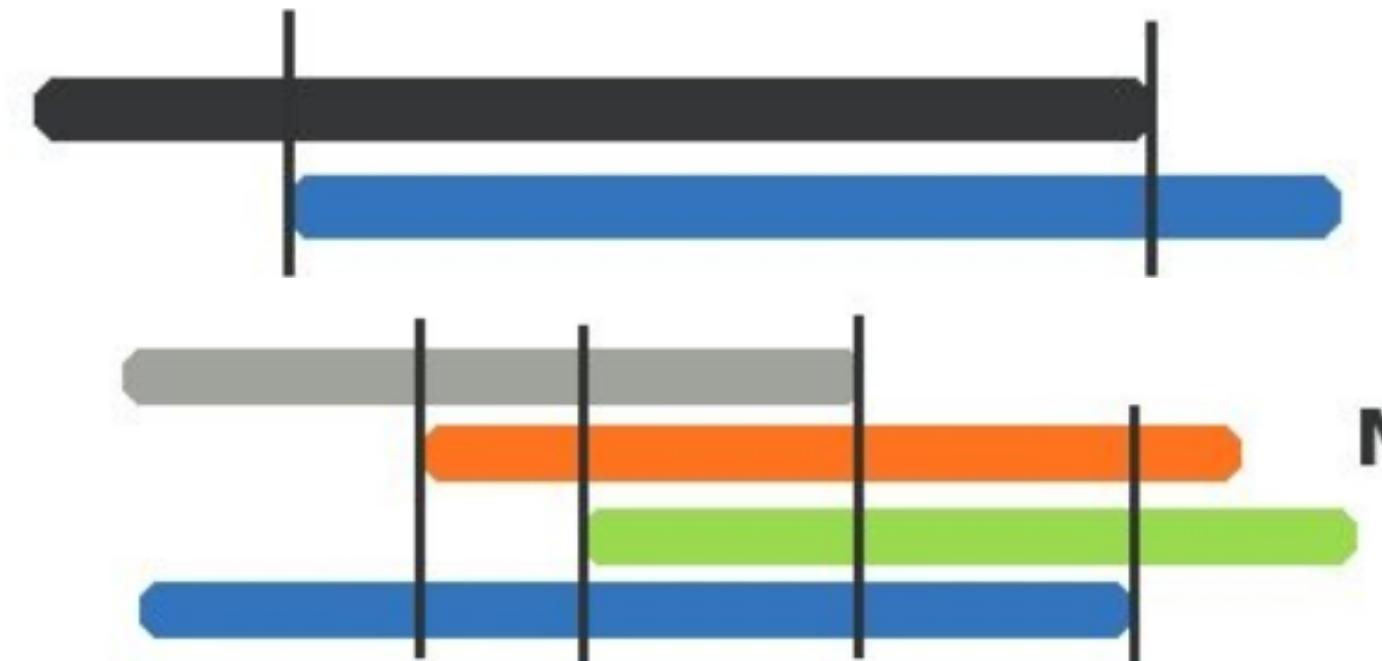
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.  
CTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGACAA.  
CACAGTACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAG.  
TAGAAAAAGAACATGTAACAGTAACACACACTCTGCTAACCTTCTAGAAGATAAG.  
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.  
AGAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGACAAGCAT.  
CTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACATTCTAGAAGACAA.  
GAAAAGAACATGTAACAGTAACACACATTCTGTTAACCTTCTAGAAGACAAGCA.  
CACAGTACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAG.  
AGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGACAAGC.

Aligned

TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.  
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.  
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.  
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGCTAACCTTCTAGAAGAT.  
TACTAGAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.  
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACATTCTAGAAGAC.  
TACTAGAAAAAGAACATGTAACAGTAACACACATTCTGTTAACCTTCTAGAAGAC.  
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.  
TACTAGAAAAAGAACATGTAACAGTAACACACACTCTGTTAACCTTCTAGAAGAC.

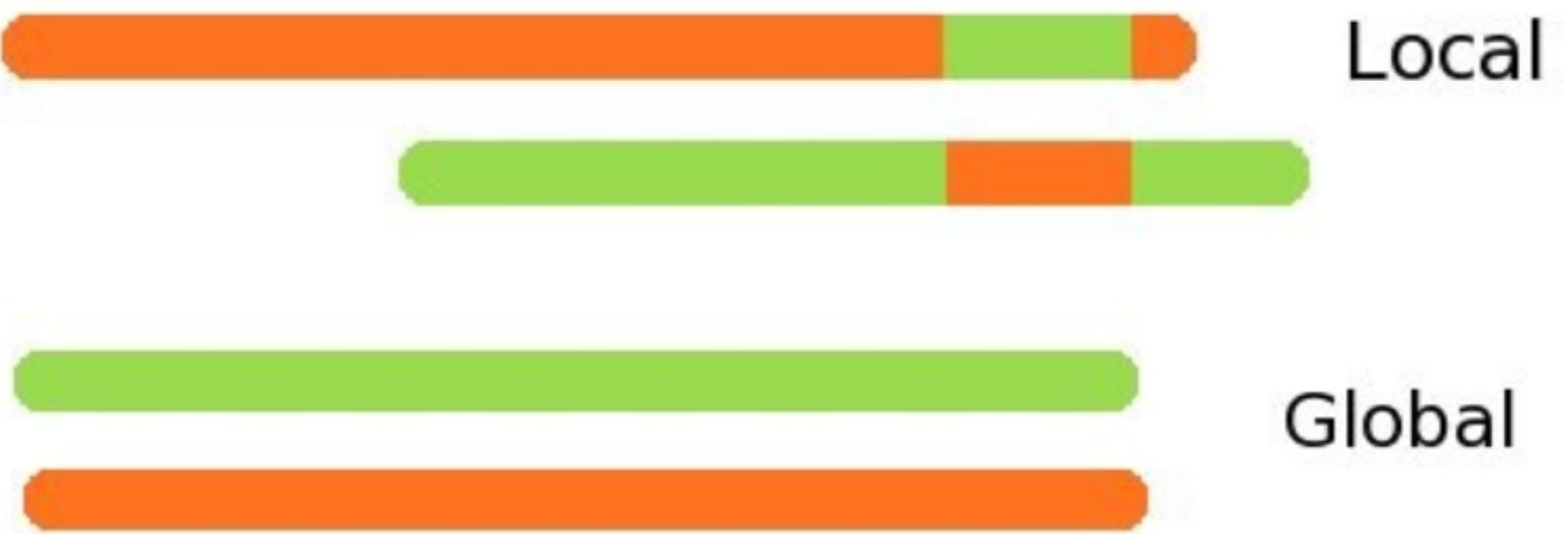
- Lining up related (homologous) positions
  - Allows comparison

# Introducción



Pareado

Múltiple



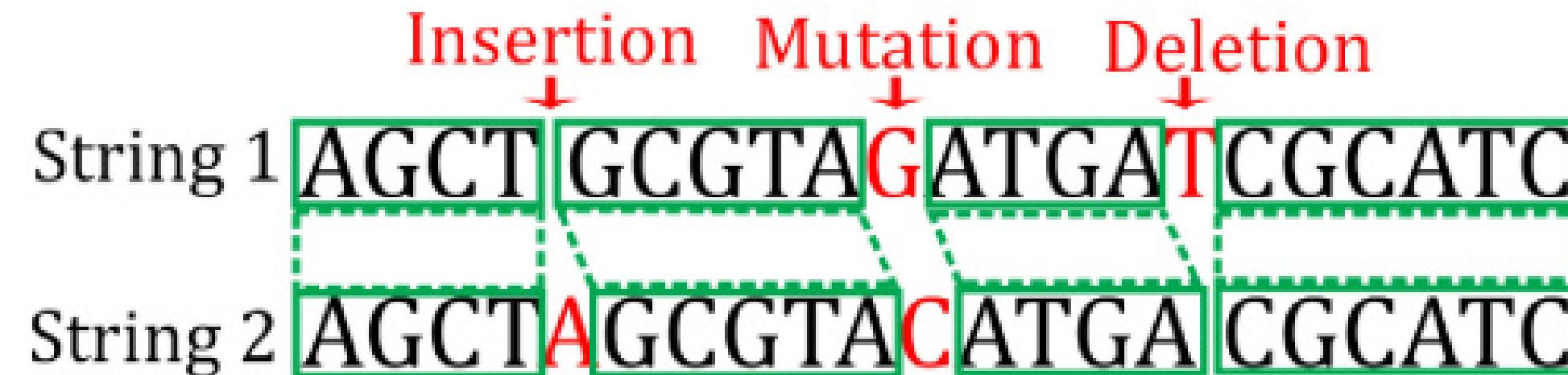
Local

Global

# Introducción

Gap (in/del)

GGAGCAATGCACAAGCTCTAGTTC	37218176
GGAGCAATGCACAAGCTCTAGTTGCAAGGTAAAGGTGCGACCTTCT 37250892	
Match	
CCCAAGGTAAAGGTGCGACCTTCTGGGG 37218206	
GGGGGCAGCAOGTCCAGTTCTCAAGGTAAAGGTGCGACCTTCTGGGG 37250942	
Mismatch	
GCAGCACGTOCAGTTCCCCAAGGTAAAGGTGCGACCTTCTAGGGTGG 37218256	
GCAGCACGTOCAGTTACCCAAGGTAAAGGTGCGACCTTCTAGGGTGG 37250992	
CACATCCCACCTTCTTCTGCTTCTCTGGACTCACTGGGGGACCAAAG 37218306	
CACATCCCACCTTCTTCTGCTTCTCTGGACTCACTGGGGGACCAAAG 37251042	
GGGTCCACTGGCAGGCAACTCTGCAGGGTGGAGGGTGGGGGCTGGCTC 37218356	
GGGTCCACTGGCAGGCAACTCTGCAGGGTGGAGGCTGGGGGCTGGCTC 37251092	



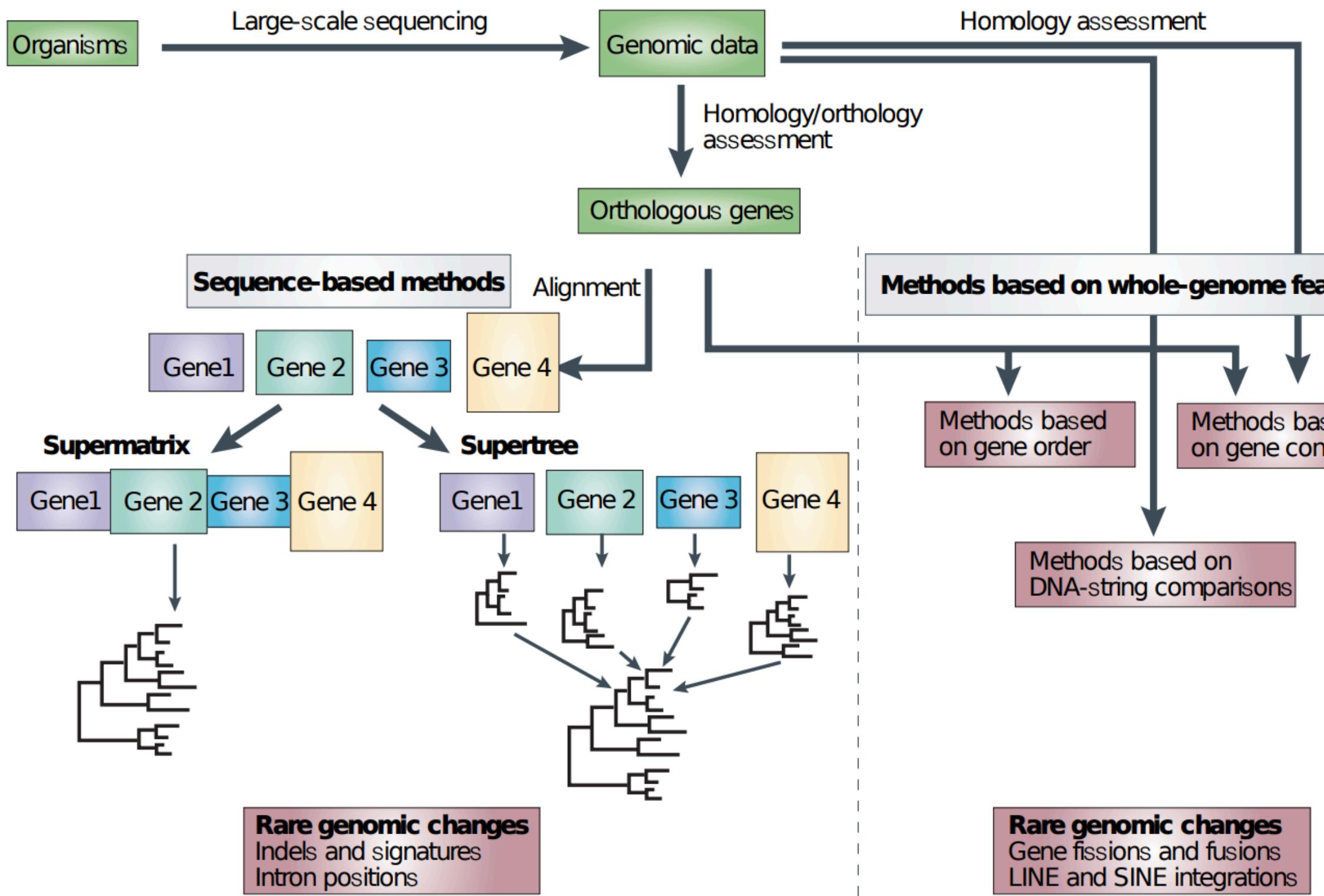
# FILOGENÓMICA

**Intersección de la filogenética y la genómica:  
El uso de grandes conjuntos de datos  
genómicos para resolver preguntas  
filogenéticas.**

Las relaciones entre los taxones se infieren en función de la homología (herencia de un ancestro común, comúnmente observada como patrones de similitud de secuencias) a lo largo de genomas completos, ya sea mediante un enfoque comparativo gen a gen, concatenado multigen o de genoma completo. (Xin-Chan & Ragan, 2013; Burki, 2014).

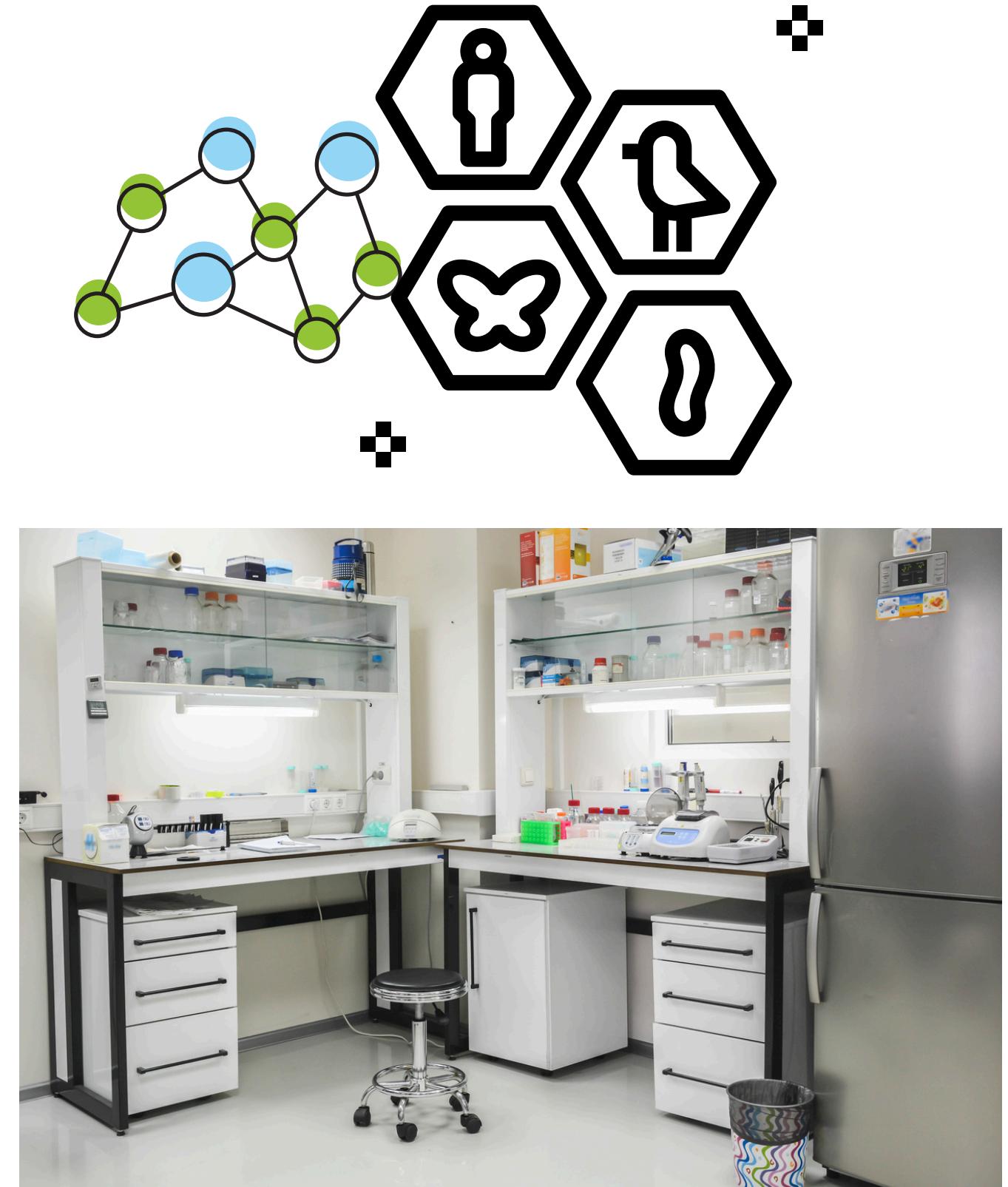
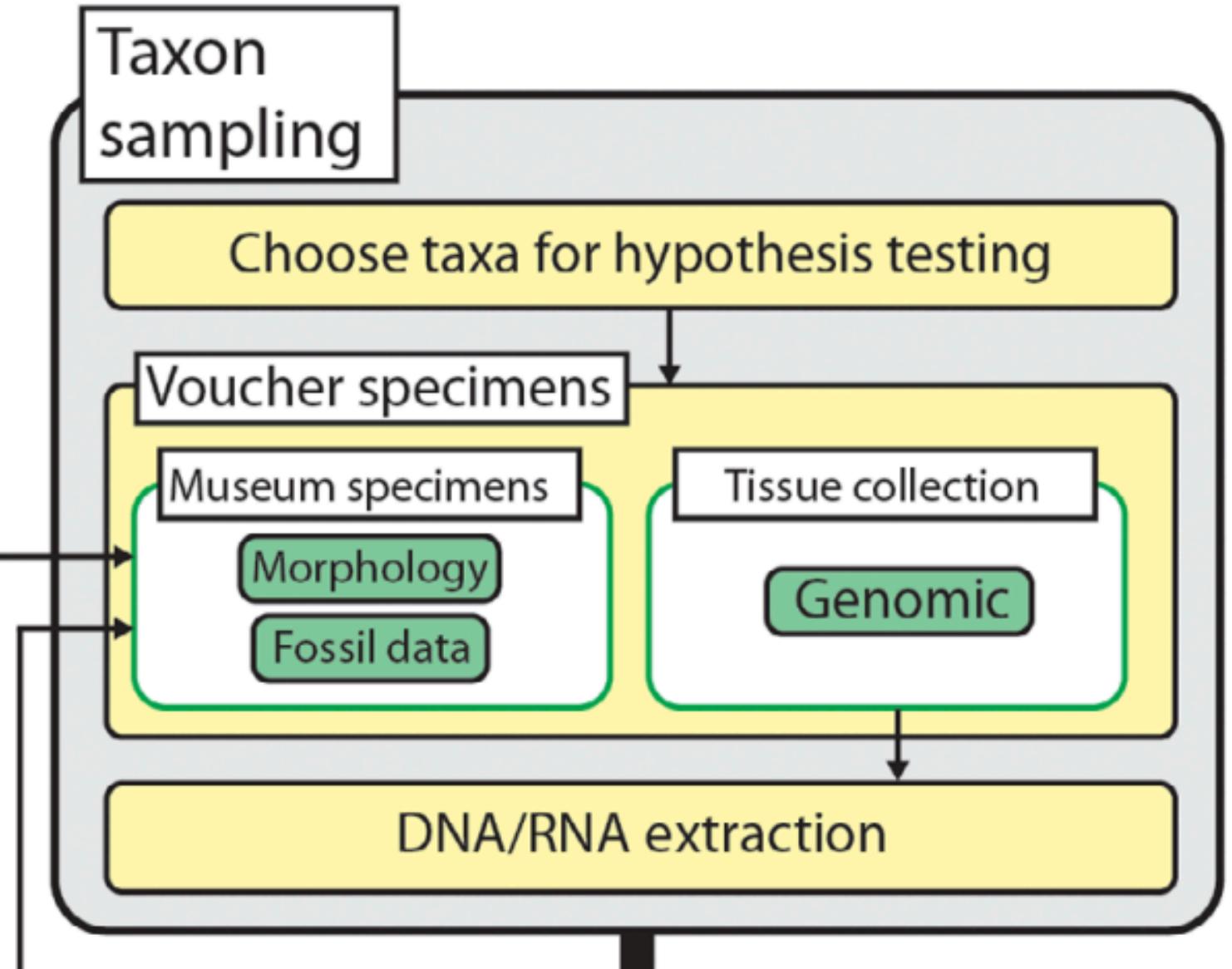
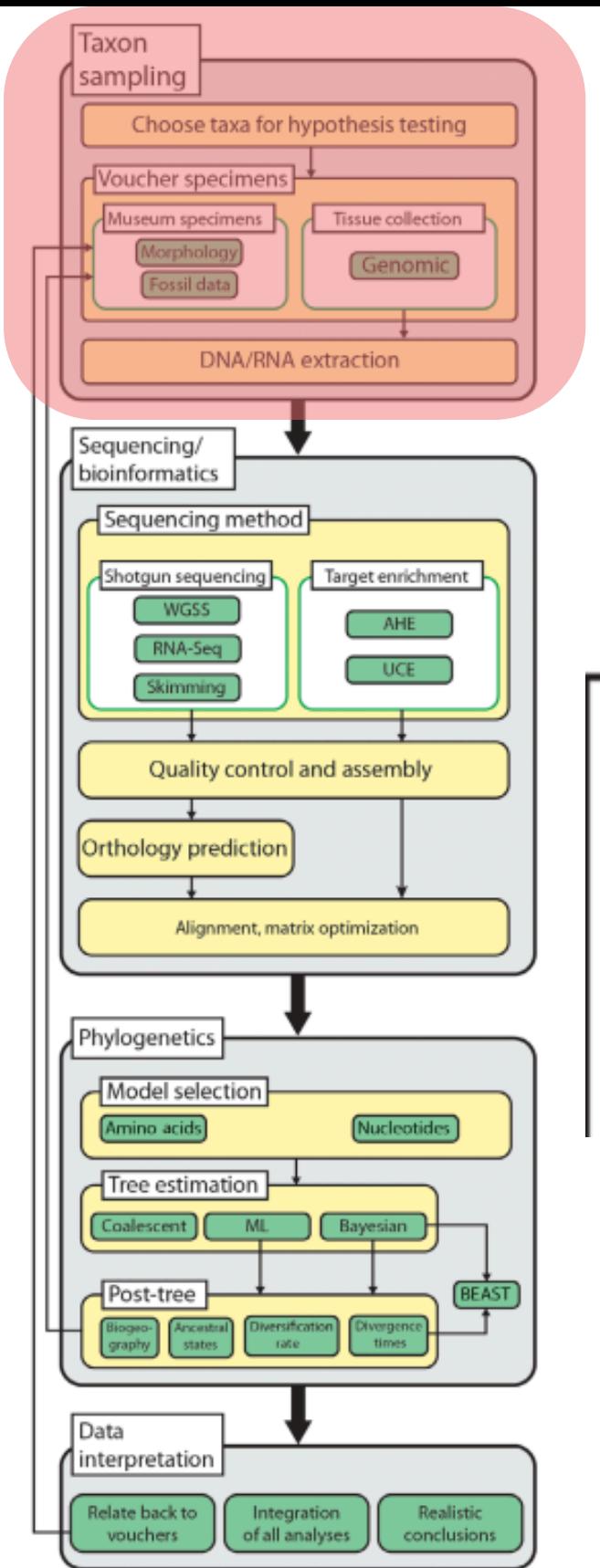
# Introducción

## Box 2 | Methods of phylogenomic inference

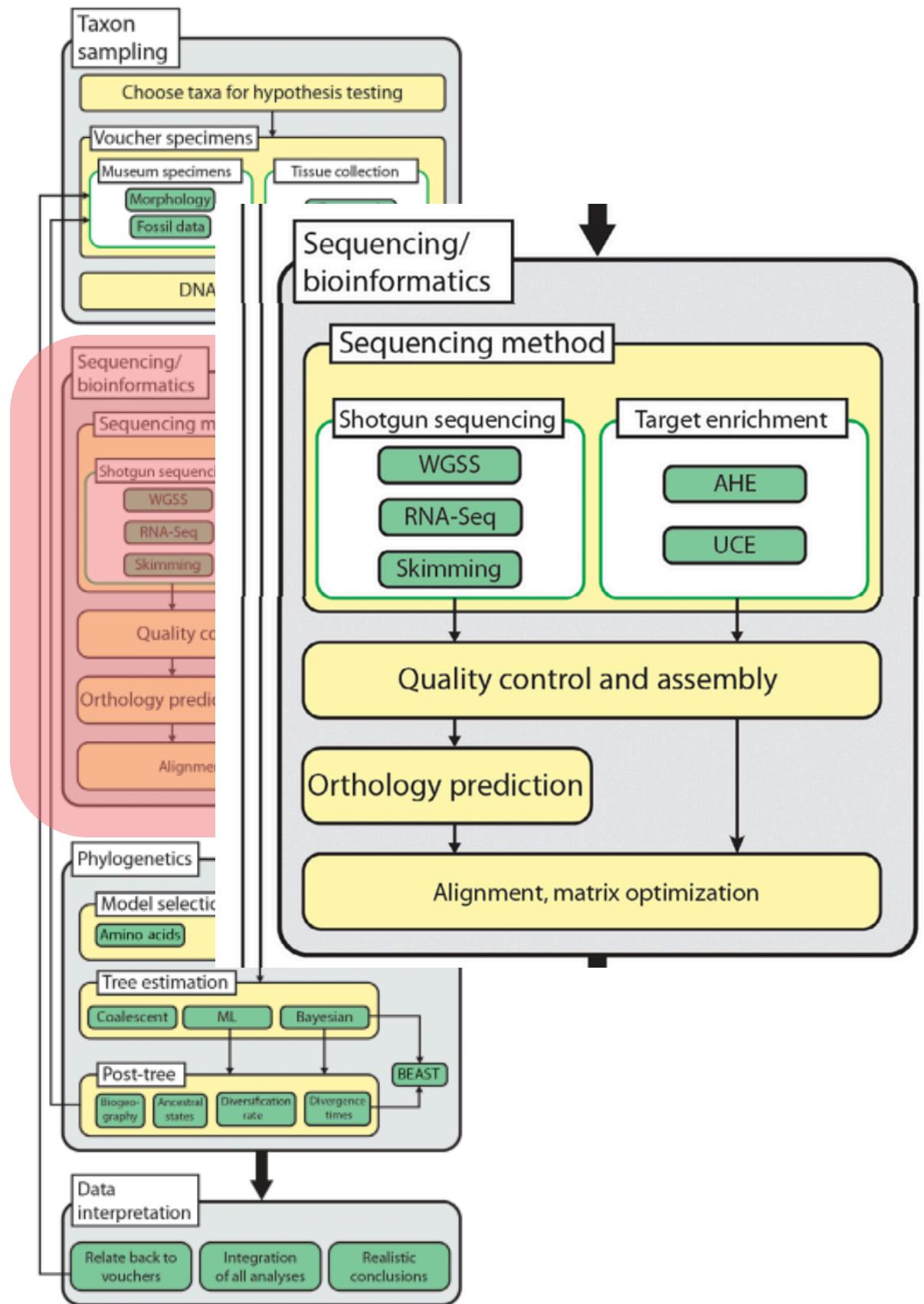


Delsuc et al., 2005

# Introducción



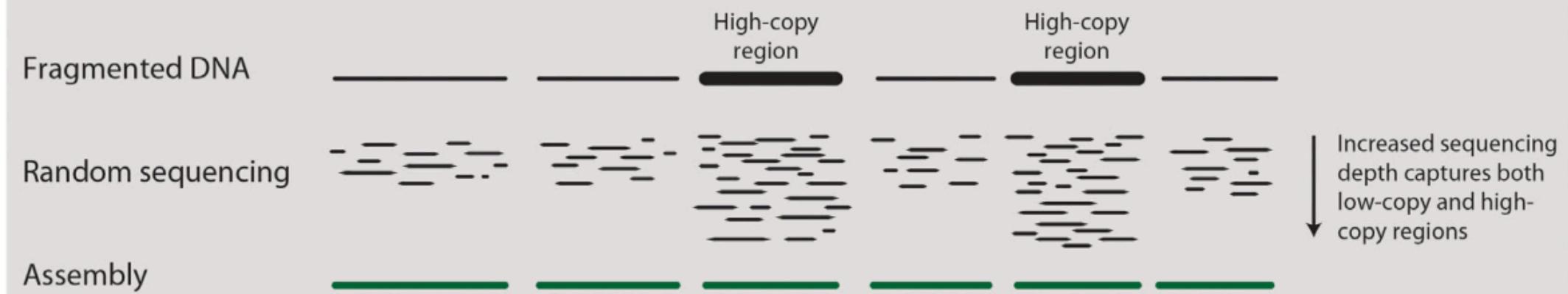
# Introducción



## (A) Genome skimming



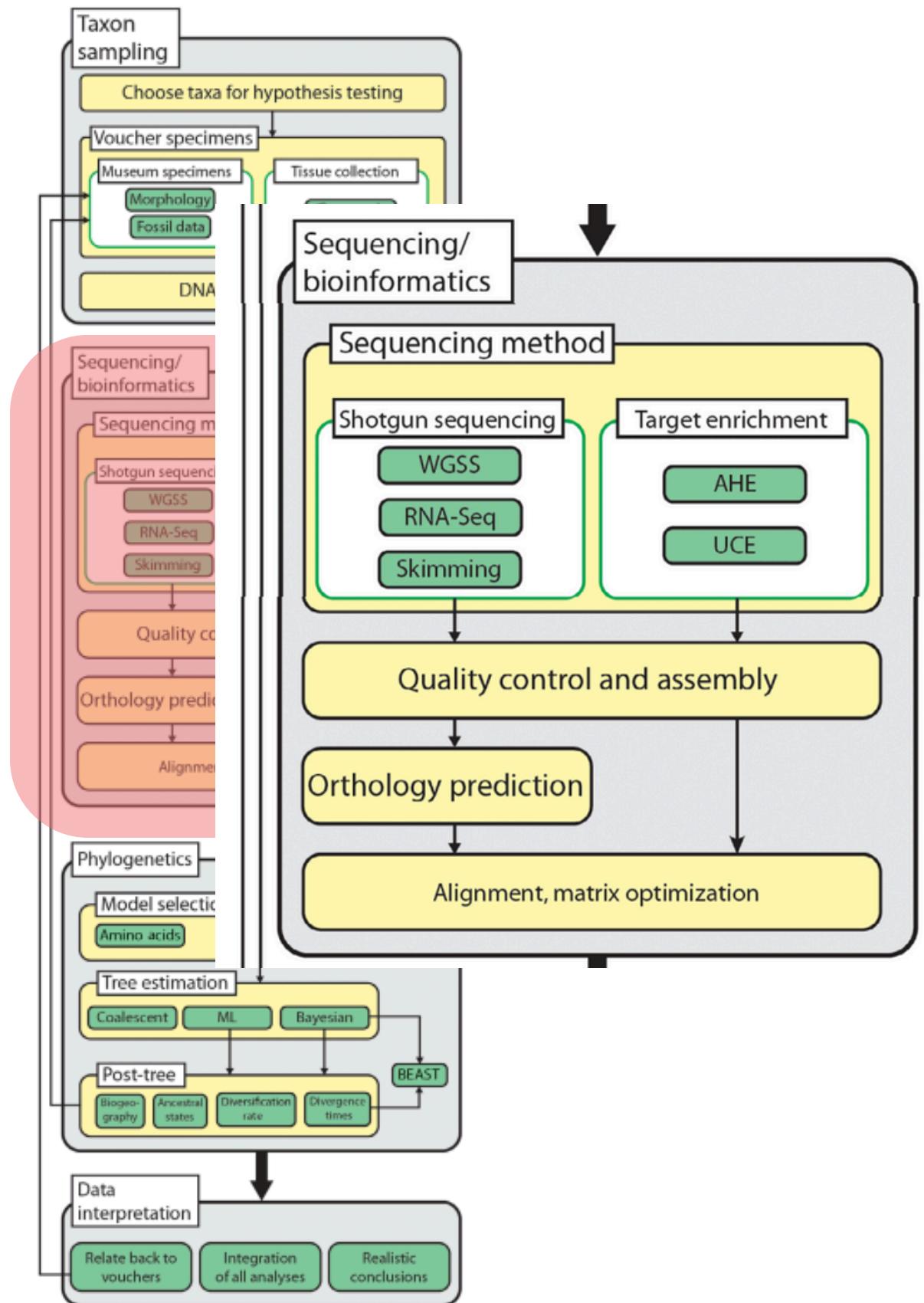
## (B) Whole-genome sequencing



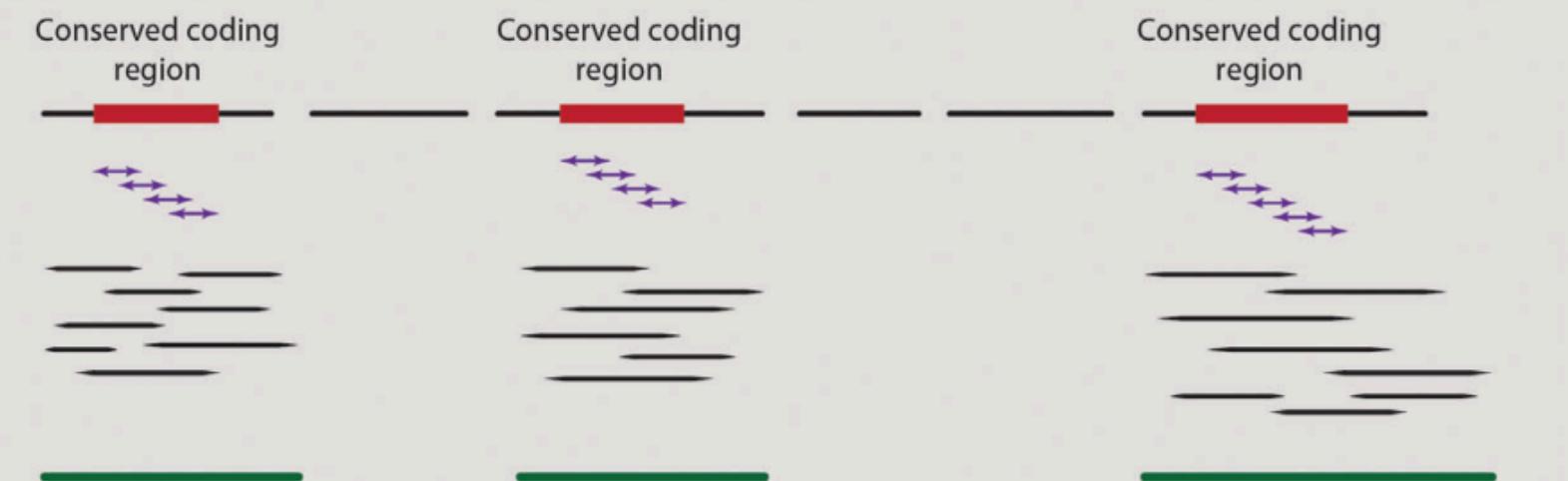
## (C) RNA sequencing



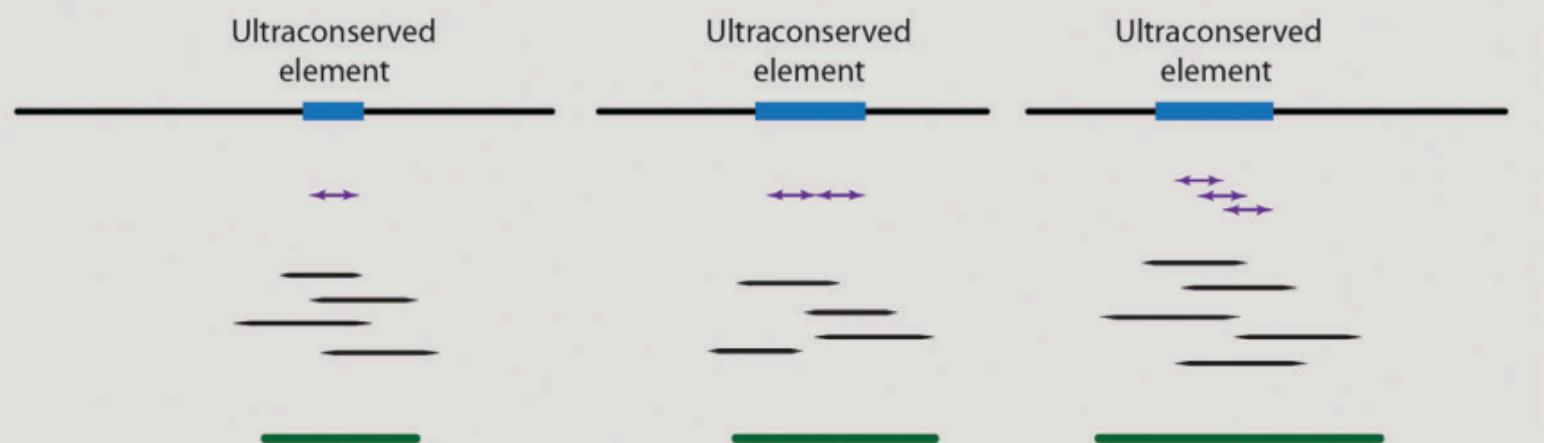
# Introducción



## (A) Anchored hybrid enrichment



## (B) Ultraconserved elements



# Introducción

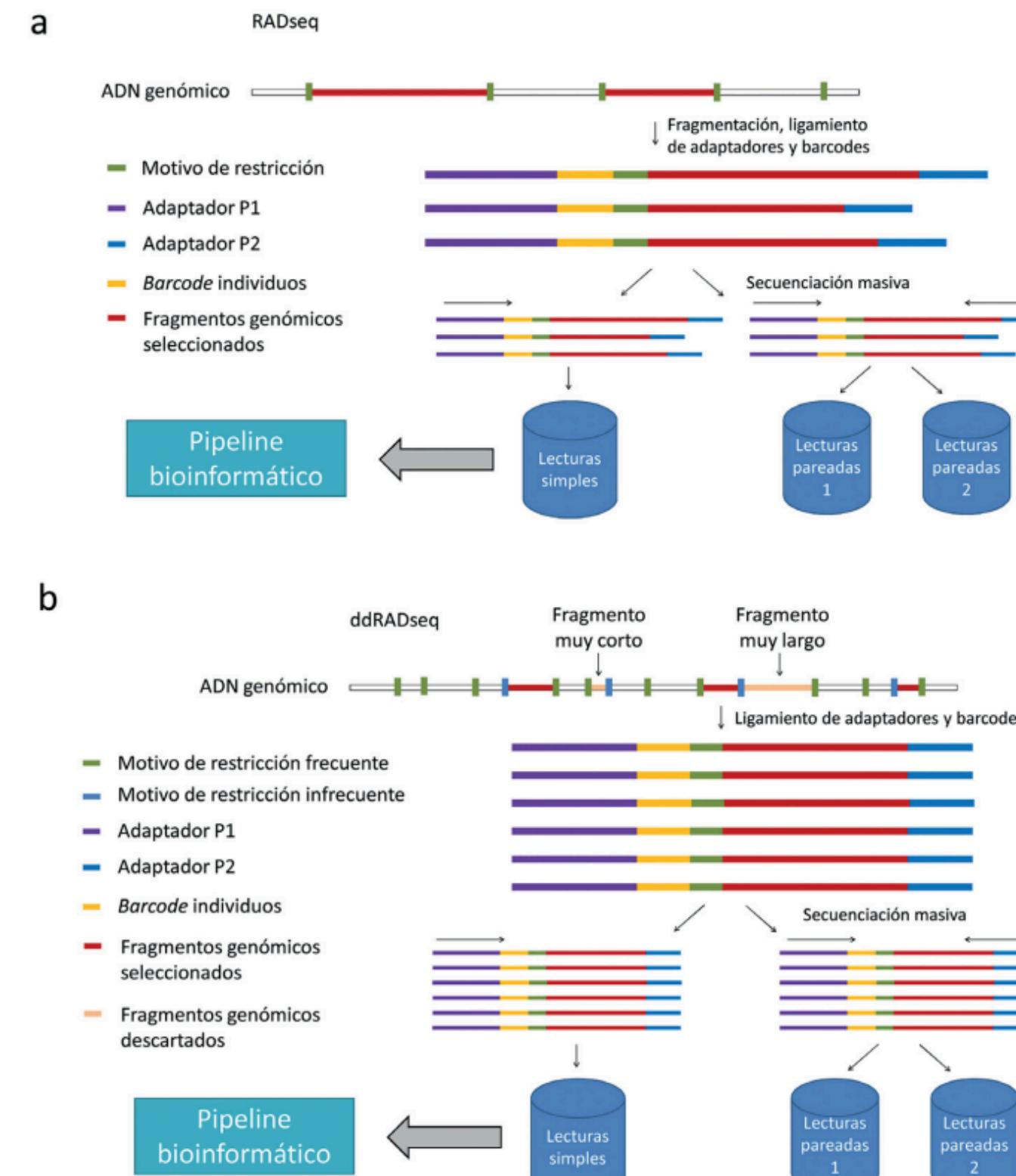
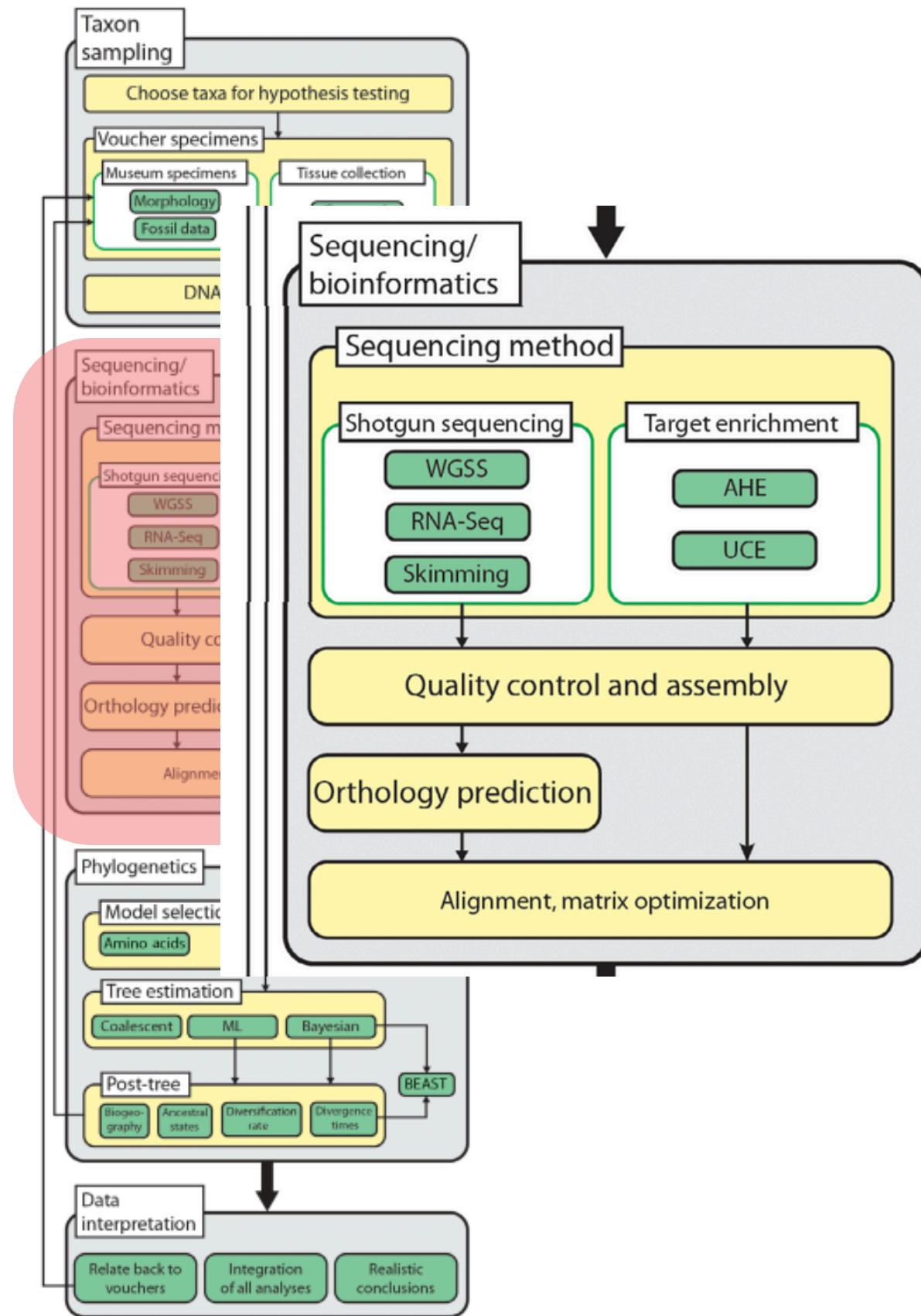
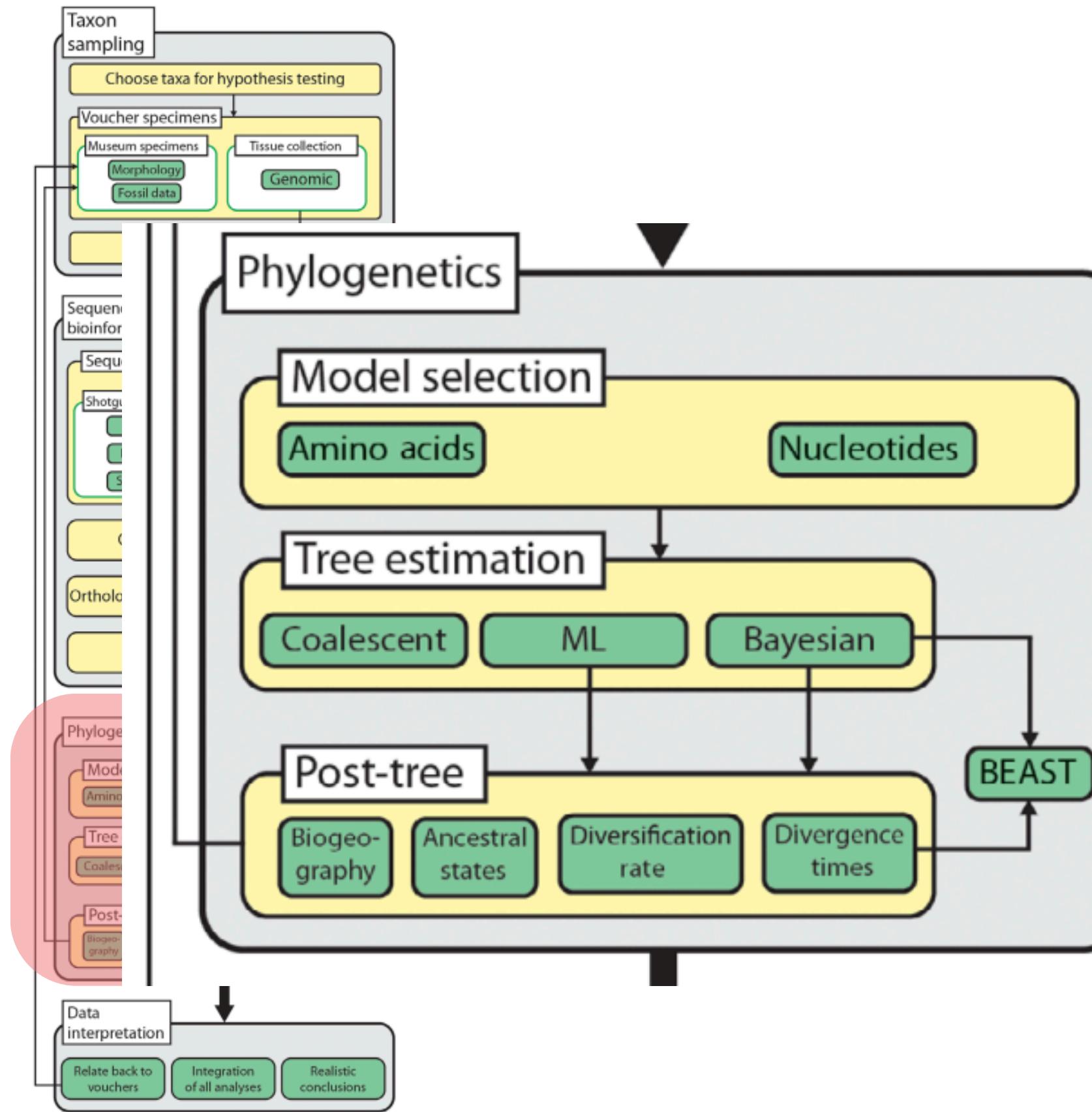


Fig. 2.- Esquema de las metodologías de genotipado RAD-seq (a) y ddRAD-seq (b).

Fig. 2.- Outline of RAD-seq (a) and ddRAD-seq (b) genotyping methodologies.

# Introducción



Young & Gillung, 2019

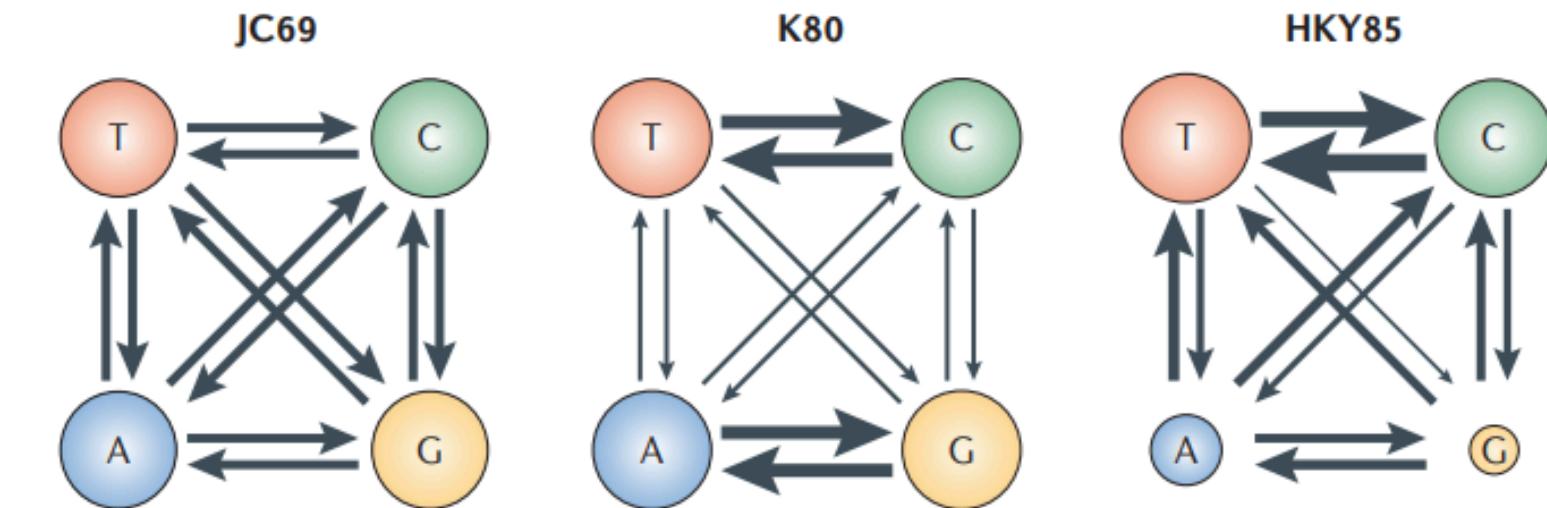
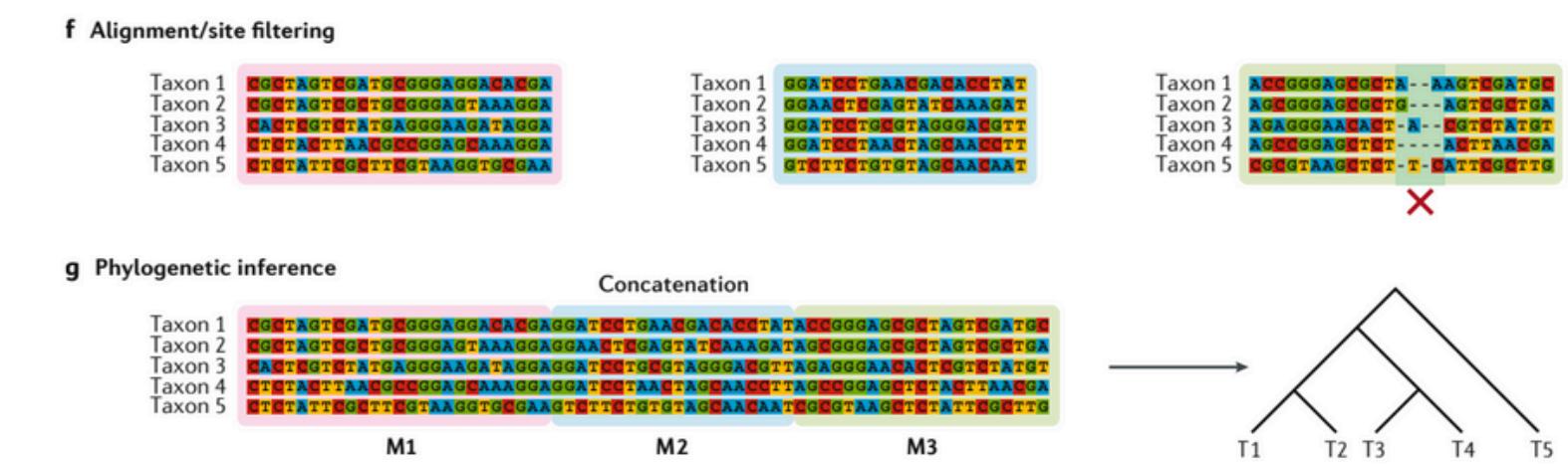


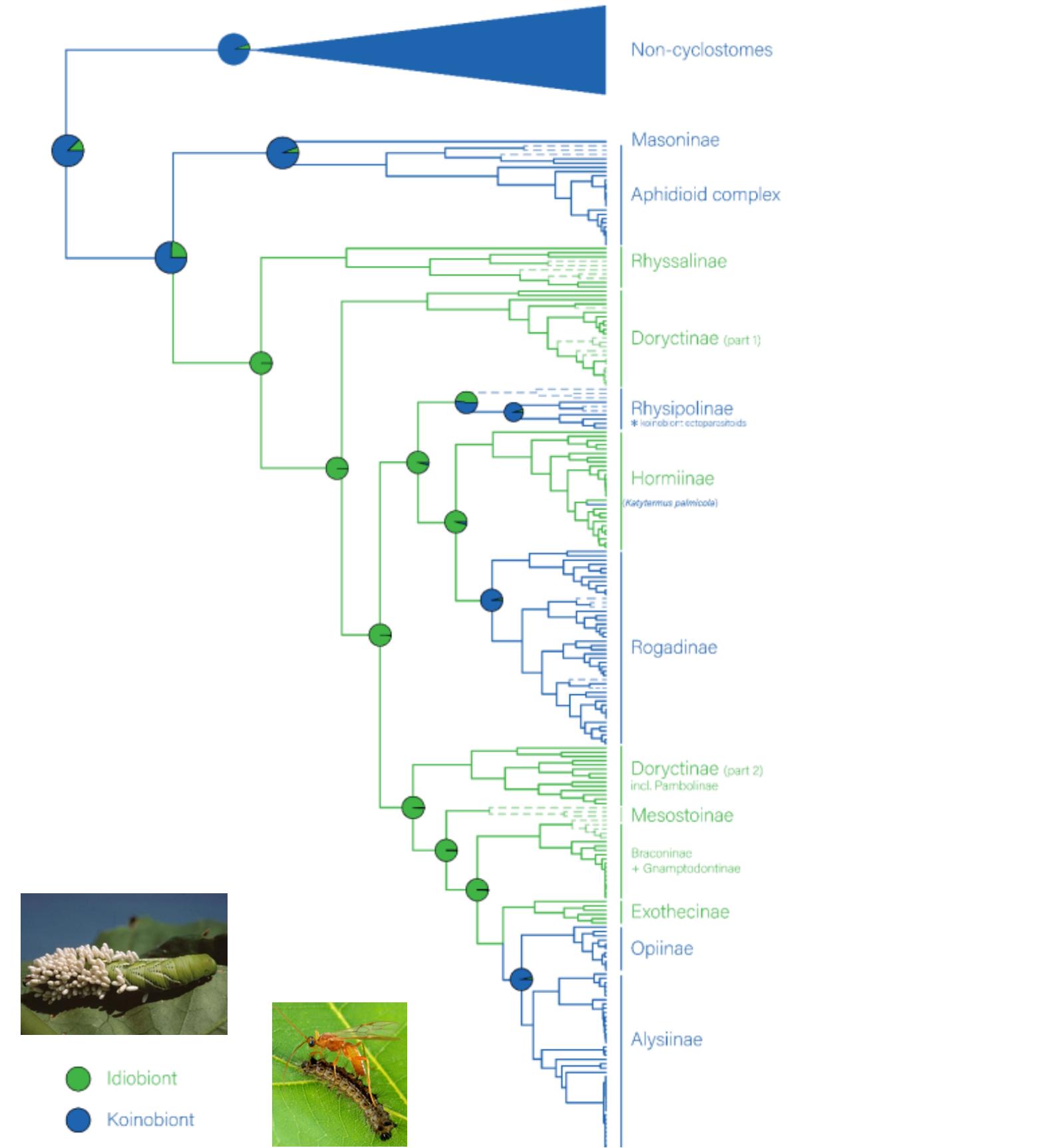
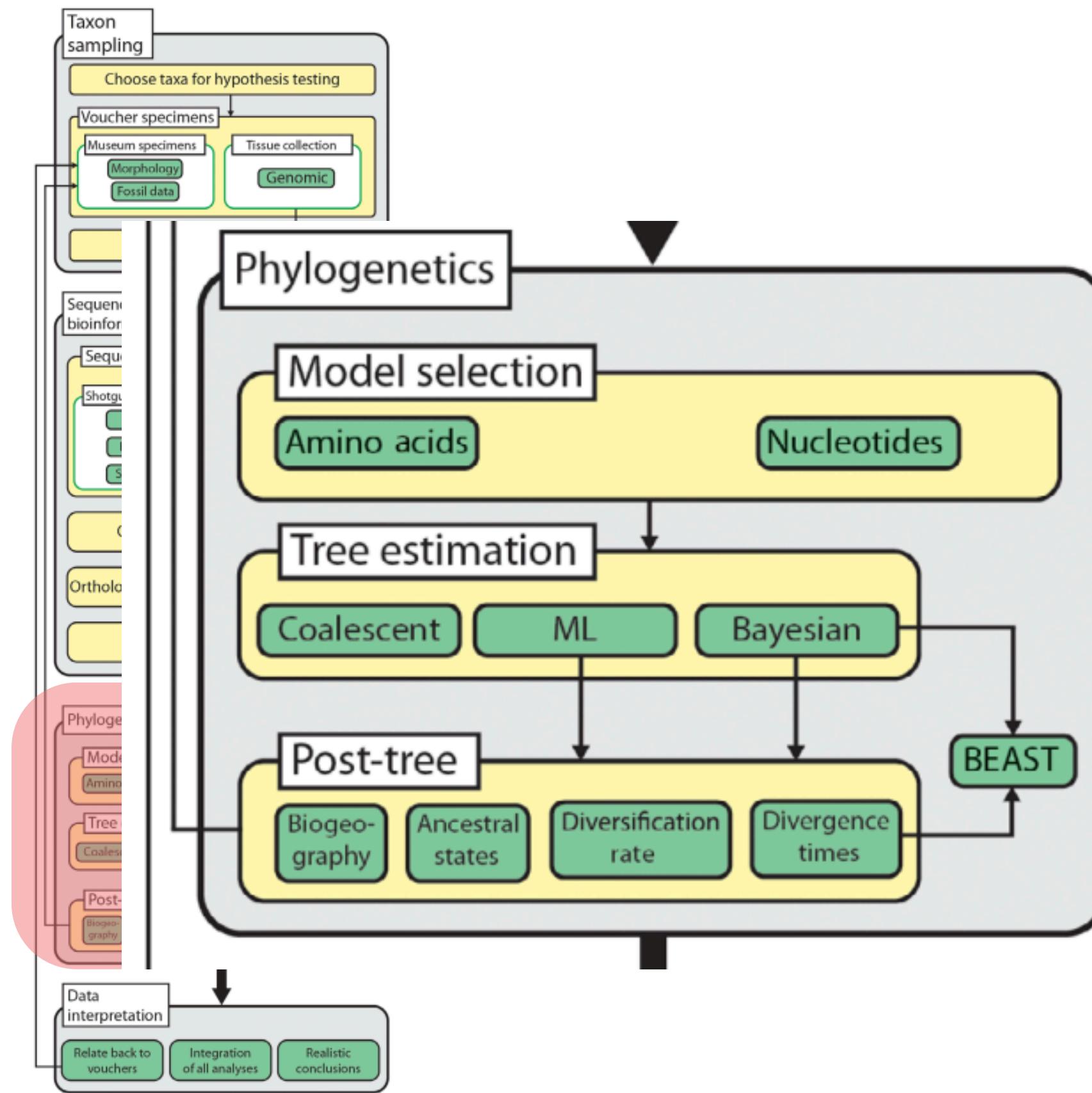
Figure 1 | **Markov models of nucleotide substitution.** The thickness of the arrows indicates the substitution rates of the four nucleotides (T, C, A and G), and the sizes of the circles represent the nucleotide frequencies when the substitution process is in equilibrium. Note that both JC69 and K80 predict equal proportions of the four nucleotides.

Yang & Rannala, 2012



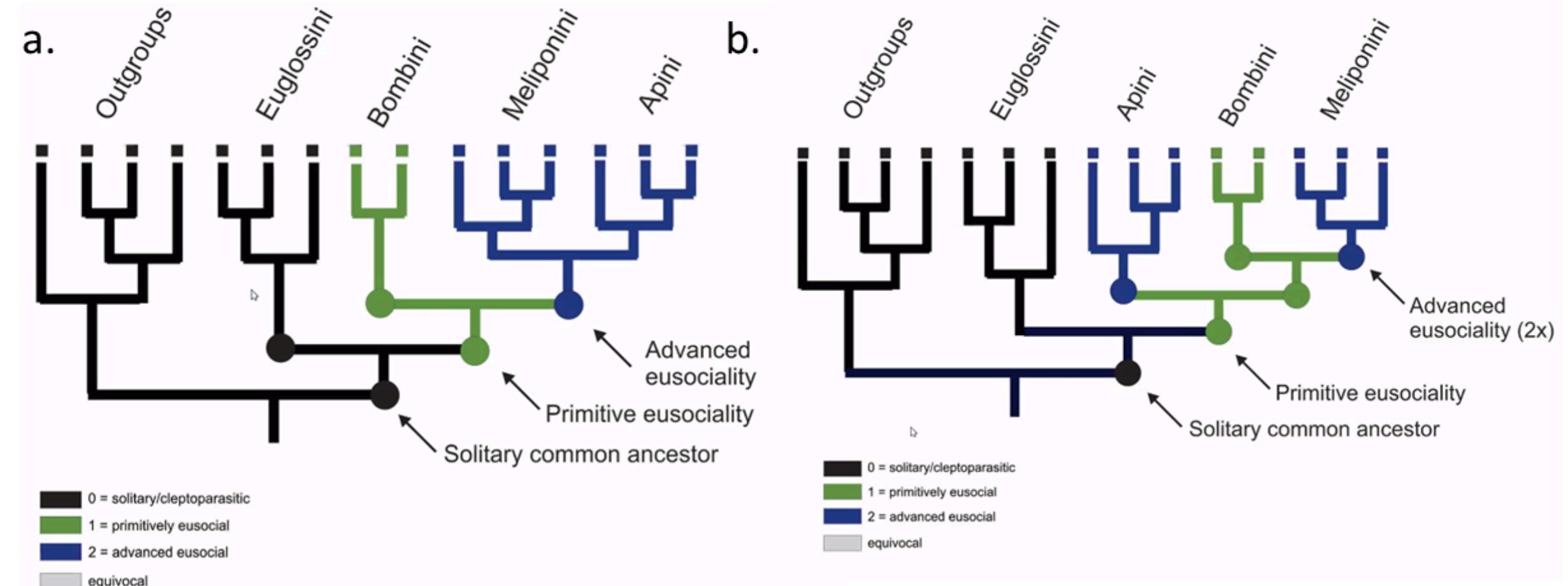
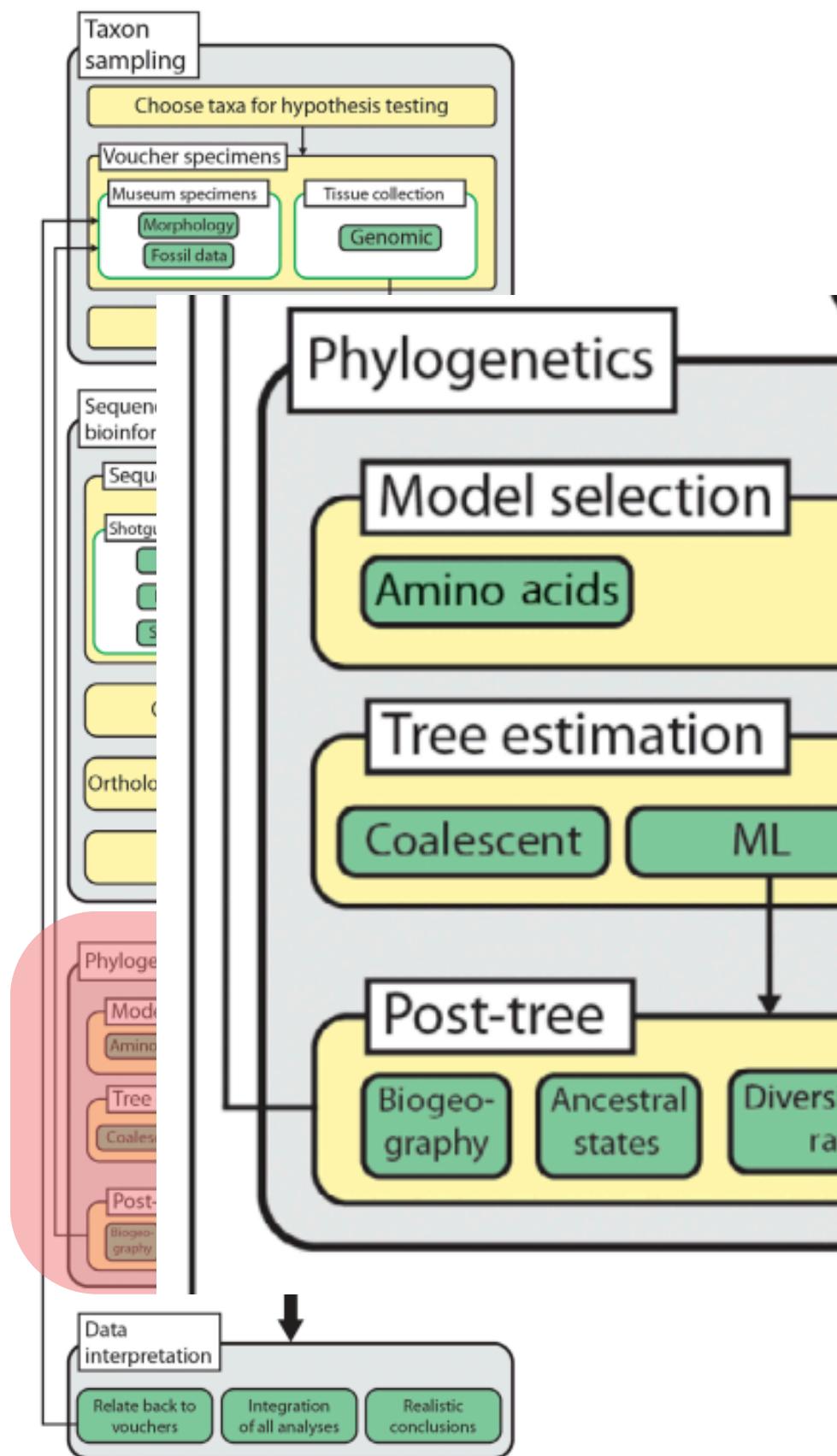
Kapli et al., 2020

# Introducción



Jasso-Martínez et al., 2022

# Introducción



Bossert et al., 2017

# Introducción

