# The Idea of Phylogenetic Trees



"I think" by Charles Darwin (1837) - One of the first evolutionary trees.



Ernst Haeckel, 1879

Ciccarelli, FD (2006). "Toward automatic reconstruction of a highly resolved tree of life.". Science; Letunic, I (2007). "Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation.". Bioinformatics
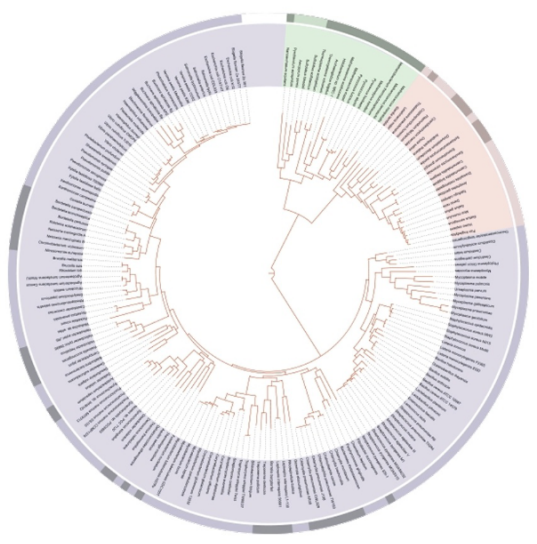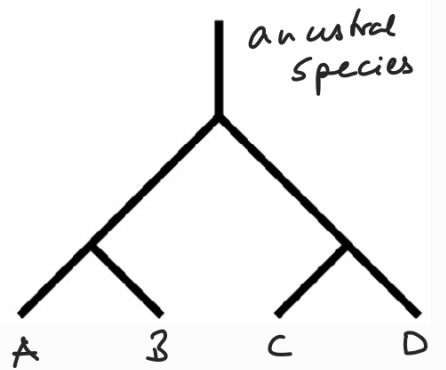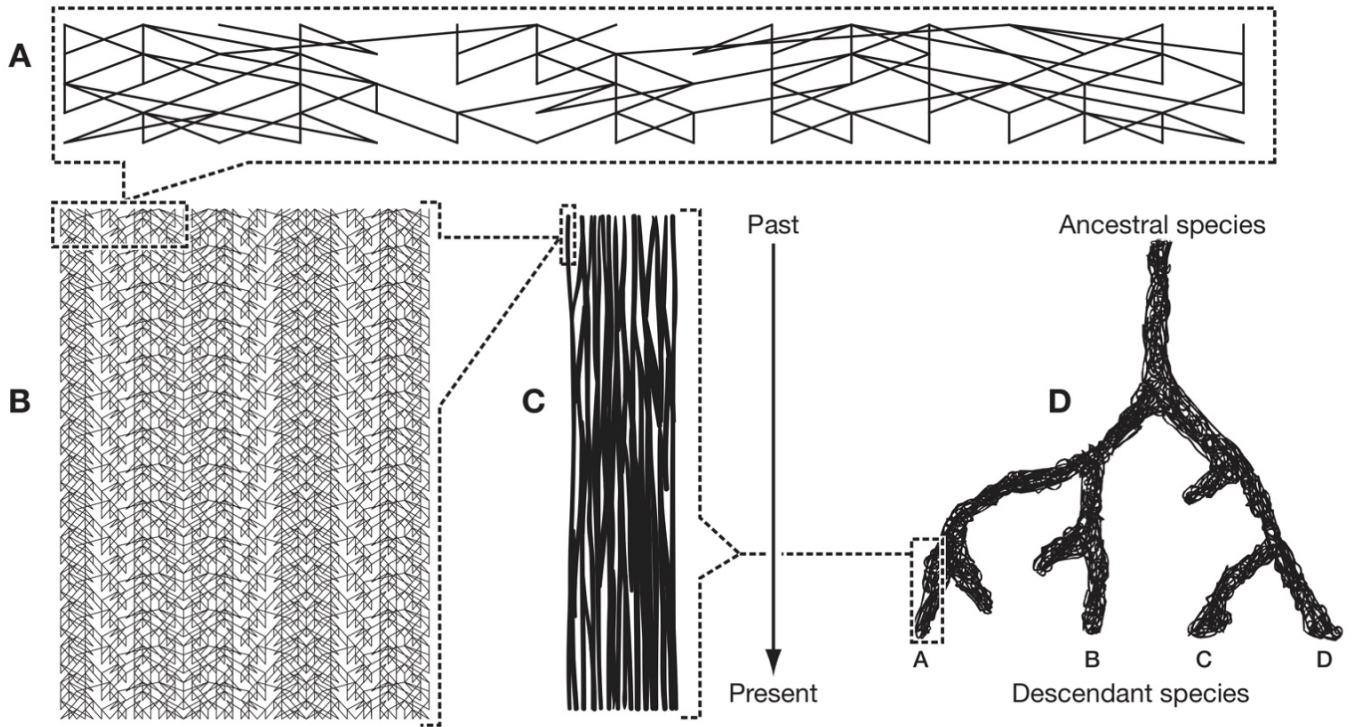
Relationship between species with sequenced genomes.

center = last universal ancestor of all life on earth.

three domains of life:

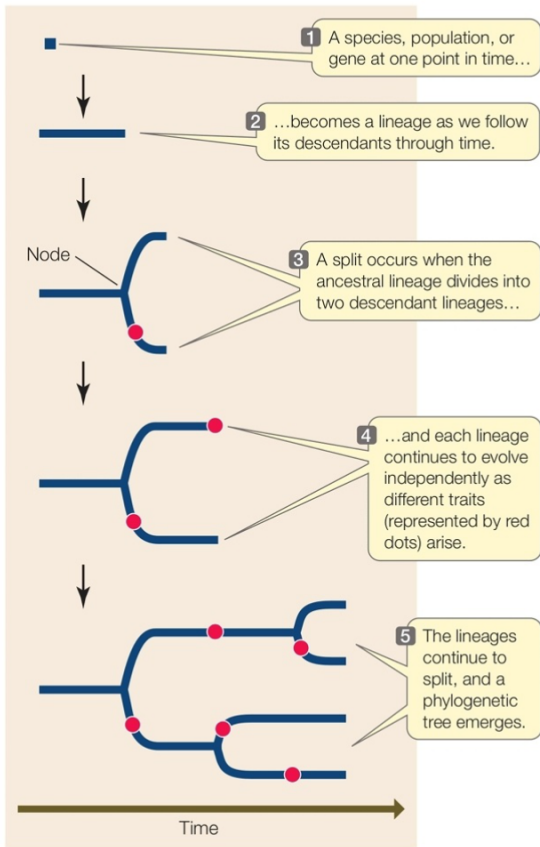eukaryota (animals, plants and fungi);

bacteria;

archaea.

# A general view



A

B

C

Past

Present

Ancestral species

D

A    B    C    D

Descendant species

A    B    C    D

ancestral species

A    B    C    D

Descendant species.

(A)

The splits in branches are called **nodes** and indicate a division of one lineage into two.

1 A species, population, or gene at one point in time…

2 …becomes a lineage as we follow its descendants through time.

Node

3 A split occurs when the ancestral lineage divides into two descendant lineages…

4 …and each lineage continues to evolve independently as different traits (represented by red dots) arise.

5 The lineages continue to split, and a phylogenetic tree emerges.

Time

Common ancestor

Chimpanzee

Human

Gorilla

Orangutan

15
Past

10

5

0
Present

Time (millions of years ago)

The positions of the nodes on the time scale (if present) indicate the times of the corresponding speciation events.

Branches can be rotated around any node without changing the meaning of the tree.
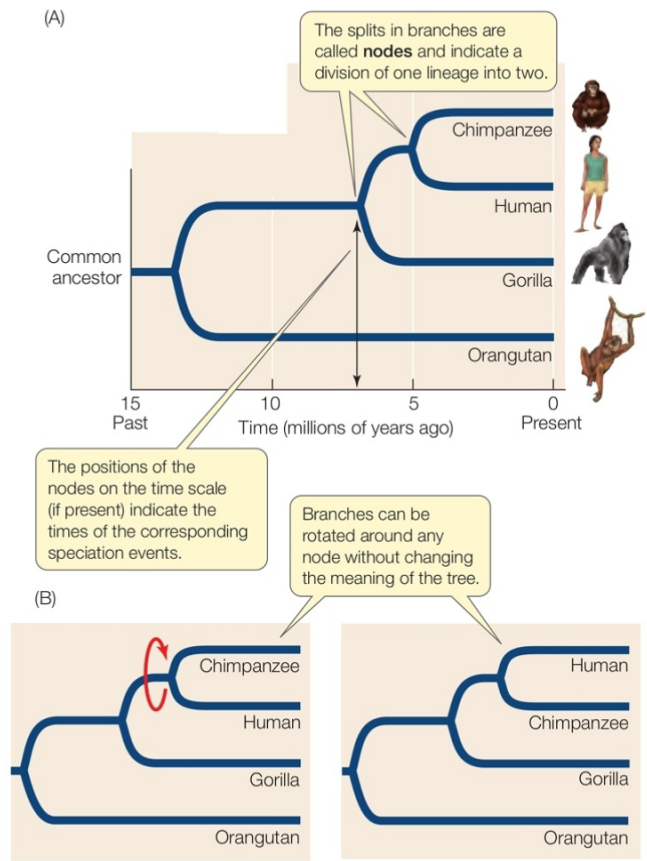
(B)

Chimpanzee

Human

Gorilla

Orangutan

Human

Chimpanzee

Gorilla

Orangutan

**22.1 The Components of a Phylogenetic Tree**   Evolutionary relationships among organisms can be represented in a treelike diagram.
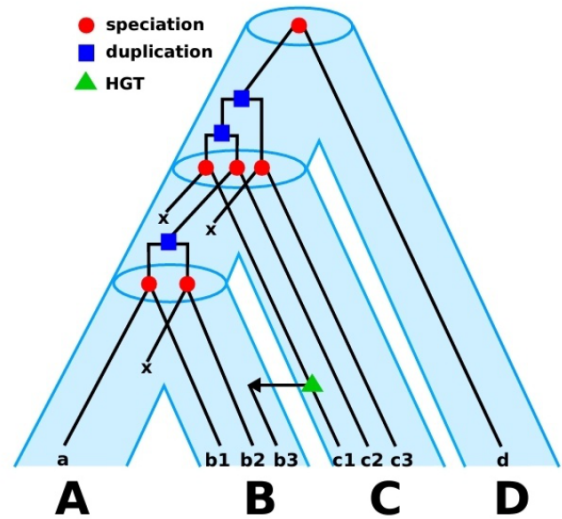
**22.2 How to Read a Phylogenetic Tree**   (A) Phylogenetic trees can be produced with time scales, as shown here, or with no indication of time. If no time scale is shown, then the trees are only meant to depict the relative order of divergence events. (B) Lineages can be rotated around a given node, so the vertical order of taxa is largely arbitrary.
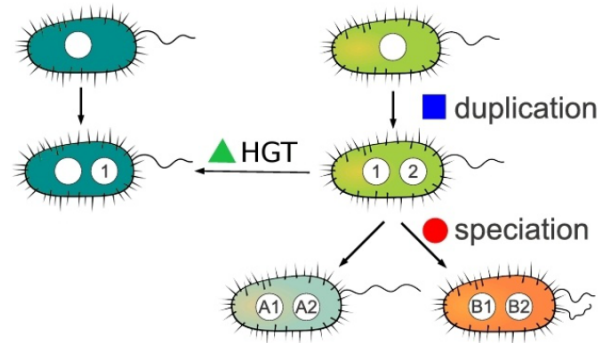
Sadava et al. (2012). "LIFE: The Science of Biology (10th edition)"
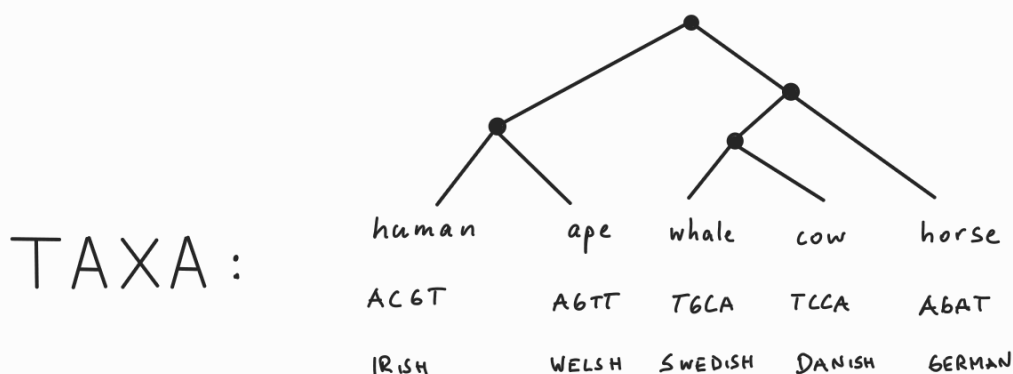
Applications:   ↗ beamer- slides.

- ▶ species are characterized by its genome:
  a "bag of genes"

- ▶ "Genes" evolve along a *rooted* tree with unique coloring
  $t : V^0 \rightarrow M = \{\bullet, \blacksquare, \blacktriangle\}$
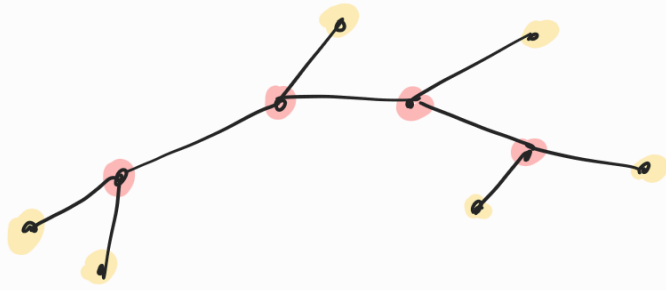
- ▶ "×" = gene loss



- ■ Gene duplication : an offspring has two copies of a single gene of its ancestor

- ● Speciation : two offspring species inherit the entire genome of their common ancestor

- ▲ HGT : transfer of genes between organisms in a manner other than traditional reproduction and across different species



- • All proposed phylogenetic trees are just HYPOTHESIS!

- • only leaves of trees are known & this knowledge must be used to infer the underlying trees

- • trees are not only about species evolution, but also of genes or other taxa as languages.

TAXA:



human    ape     whale   cow     horse

ACGT     AGTT    TGCA    TCCA    AGAT

IRISH    WELSH   SWEDISH DANISH  GERMAN

$T = (V, E)$ is  tree  if  connected & acyclic



- leaf
- internal vertices (inner)

$L(T)$ = leaf set of $T$

$T$  rooted if one vertex $\rho \in V$ is called root

unrooted



rooted



IF not stated differently $T$ is  phylogenetic , ie.

unrooted $T$ every inner vertex has at least degree 3
rooted $T$ —"— —"— 2 children

children of v



$T$ fully resolved (= binary ) if

$\forall v \in V \setminus L(T)$ : degree $v$ is 3   ($T$ unrooted)
                    $v$ has exactly 2 children ($T$ rooted)

in rooted trees we have $\boxed{\text{partial order } \leq_T}$ along vertices in $T$:

$v \leq_T w$ if $w$ lies on unique path from $v$ to $\varsigma$

($v = w$ possible) in this case, $v$ descendant of $w$

$w$ ancestor of $v$

write $v <_T w$ if $v \leq_T w$ & $v \neq w$.

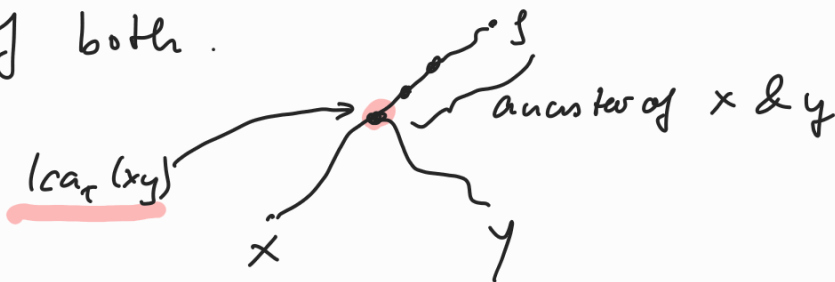last common ancestor $\boxed{lca_T(xy)}$ of $x, y \in V$ is $\leq_T$ - minimal vertex that is ancestor of both.

ancestor of $x$ & $y$

$lca_T(xy)$

$x$        $y$

Depending on the application, phylogenetic trees may:

- ▶ be rooted or unrooted
- ▶ have weighted or unweighted edges / vertices
- ▶ labeled vertices / edges
- ▶ have bounded degree
  (maximum nr of children of each internal node)
- ▶ . . .

---

- ▶ Inference of the gene or species tree $T$ is a classical problem of molecular phylogenetics.

  In practice it can only be solved approximately.

- ▶ Only leaves of tree corresponding to extant (currently "observable") taxa is available.

- ▶ **Reconstructed trees do only provide a hypothesis about history!**

**Lemma**

*There are* $(2n - 3)!!$ *rooted trees* $(2n - 5)!!$ *unrooted trees with n leaves labeled from* $1, \ldots, n$.

$$(m)!! := \prod_{k=0}^{\lceil \frac{m}{2} \rceil - 1} (m - 2k) = m(m - 2)(m - 4) \cdots .$$

| $n$ | 3 | 4 | 5 | 6 | 10 | 20 |
|---|---|---|---|---|---|---|
| Exmpl: unrooted | 1 | 3 | 15 | 105 | 2'027'025 | $2.22 \cdot 10^{20}$ |
| rooted | 3 | 15 | 105 | 945 | 34'459'425 | $8.20 \cdot 10^{21}$ |

Enumeration / exhaustive search is no option!

---

**Aim:** Assemble a tree representing a hypothesis about the evolutionary history of a set of genes, species or other taxa.

**Methods:**

- ▶ Distance Based e.g.:
  - ▶ Ultrametric Tree Reconstruction (UPGMA)
  - ▶ Additive Tree Reconstruction (Neighbor-Joining)

- ▶ Character Based e.g.:
  - ▶ Parsimony Methods (Fitch- and Sankoff Algorithm)
  - ▶ Maximum Likelihood (not part here)

- ▶ Consensus Methods e.g.:
  - ▶ Supertree from subtrees (BUILD)
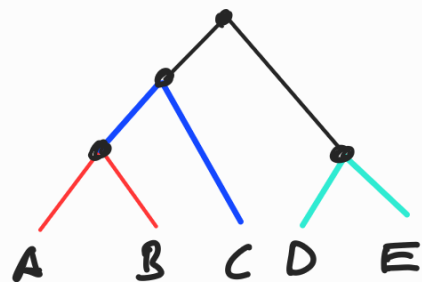
# DISTANCE-BASED METHOD

**UPGMA** (**u**nweighted **p**air **g**roup **m**ethod with **a**rithmetic **m**ean)

( $\equiv$ bottom-up hierarchichal clustering method)

**ALGO:** **IN:** Symmetric Distance matrix $D: X \times X \longrightarrow \mathbb{R}$ , $X = \{x_1 \dots x_n\}$
(or Similarities)

now, in each step merge two "closest" cluster starting with $C_1 = \{x_1\}, \dots, C_n = \{x_n\}$ as singleton clusters.



After merging 2 clusters $C_i$ & $C_j$ into new cluster $C_{new}$ distance as

$$D(C_{new}, C) = \frac{1}{|C_{new}||C|} \sum_{\substack{x \in C_{new} \\ y \in C}} D(x,y) \qquad \forall C \neq C_{new}.$$

$\equiv$ mean distance between objects $x \in C_{new}$ & $y \in C$

**REPEAT** until one cluster remains

$$C_1 = \{a\}, \; C_2 = \{b\}, \; C_3 = \{c\}, \; C_4 = \{d\}$$

D

|   | a | b | c | d |
|---|---|---|---|---|
| a | 0 | 8 | 5 | **3** |
| b |   | 0 | 8 | 8 |
| c |   |   | 0 | 5 |
| d |   |   |   | 0 |

↙ closest

⟹ merge $C_1$ & $C_4$ into $C_{new} = \{a, d\}$

⟹ new distances:

$$D(C_{new}, C_2) = \frac{1}{2 \cdot 1} \left( D(a,b) + D(db) \right)$$
$$= \frac{1}{2} (8 + 8) = 8$$
$$D(C_{new}, C_3) = \frac{1}{2} \left( D(a,c) + D(dc) \right)$$
$$= \frac{1}{2} (5 + 5) = 5$$

↓ update

|   | $C_{new}$ {a,d} | $C_2$ b | $C_3$ c |
|---|---|---|---|
| {a,d} | 0 | 8 | **5** |
| b |   | 0 | 8 |
| c |   |   | 0 |

⟹ merge $C_{new}$ & $C_3$ into $C'_{new} = \{a, c, d\}$

⟹ new distances

$$D(C'_{new}, C_2) = \frac{1}{3 \cdot 1} \left( D(a,b) + D(cb) + D(db) \right)$$
$$= \frac{1}{3} (8 + 8 + 8) = 8$$

↓ update

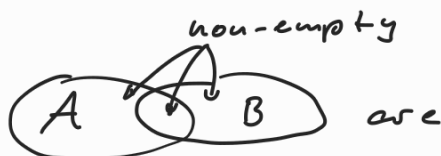|   | $C'_{new}$ {{a,d}, c}} | $C_2$ b |
|---|---|---|
| { | 0 | 8 |
| b |   | 0 |

finally merge $C'_{new}$ & $C_2$.

In this iterative process we obtained the set of clusters:

$$\mathcal{C} = \{ \{a\}, \{b\}, \{c\}, \{d\}, \{ad\}, \{acd\}, \{abcd\} \}.$$

**DEF**   2 sets $A, B$ s.t. $(A \overset{\text{non-empty}}{\cap} B)$ are said to **overlap**.

Hence, $A, B$ do **not** overlap if $A \cap B \in \{A, B, \emptyset\}$.

A set $\mathcal{C}$ of clusters is a __hierarchy__ if no two elements of $\mathcal{C}$ overlap.

Given a rooted tree $T$: with leaf set $X$

let $\mathcal{L}(v) = \{ x \in X : x \leq_T v \}$

& put $\mathcal{C}(T) = \{ \mathcal{L}(v) \mid v \in V(T) \}$

$\Rightarrow$ $\mathcal{C}(T)$ is a hierarchy [Exercise]

__Thm:__ Let $\mathcal{C}$ be a collection of non-empty

[without proof]

subsets of $X$. Then, there is a phylogenetic rooted tree on $X$ s.t.
$\mathcal{C}(T) = \mathcal{C} \iff \mathcal{C}$ is hierarchy on $X$

By construction, UPGMA gives us a hierarchy $\mathcal{C}$ & thus a tree!

$\mathcal{C} = \{ \{a\}, \{b\}, \{c\}, \{d\}, \{ad\}, \{acd\}, \{abcd\} \}$.

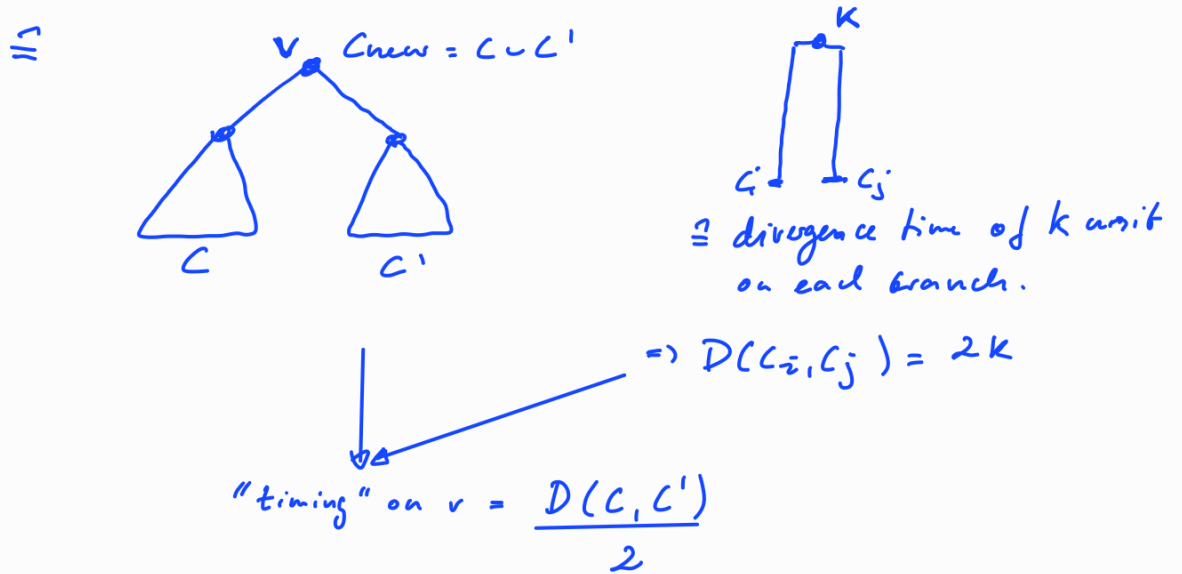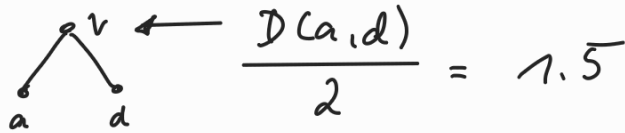## Can keep track of branch length $\delta$

merging $C \& C'$ into $C_{new}$:

$\cong$

$V$ $C_{new} = C \cup C'$

$C$ $C'$

$K$

$C_i$ $C_j$

$\underline{\cong}$ divergence time of $k$ unit on each branch.

$\Rightarrow D(C_i, C_j) = 2k$

"timing" on $v = \dfrac{D(C, C')}{2}$

---

$\underline{\text{In Exmpl:}}$ merged first $\{a\} \& \{d\}$

$v \leftarrow \dfrac{D(a,d)}{2} = 1.5$

$a$ $d$

& updated distances & merged $C_{new} = \{ad\} \& \{b\}$ into $C_{new}$

$1.5$

$a$ $d$ $b$

$\dfrac{D(C_{new}, \{b\})}{2} = \dfrac{5}{2} = 2.5.$

& so on ....

---

Branch-length

$1.5$

$1$ $a$

$1.5$ $d$

$2.5$ $c$

$5$ $4$ $b$

4 3 2 1 0

Perfectly represented by tree

|   | a | b | c | d |
|---|---|---|---|---|
| a | 0 | 8 | 5 | 3 |
| b |   | 0 | 8 | 8 |
| c |   |   | 0 | 5 |
| d |   |   |   | 0 |

IF "branch length" not needed => we a have a heuristic
to build tree.


Q: IF "branch length" needed .... does it always work?
A: NO!

|   | a | b | c |
|---|---|---|---|
| a | 0 | 1 | 2 |
| b |   | 0 | 3 |
| c |   |   | 0 |

=

|       | {a,b} | c |
|-------|-------|---|
| {a,b} | 0     | 2,5 |
| c     |       | 0 |

=> tree



$$\gamma(ac) = 2 \cdot 1,25 = 2,5 \neq D(a,c) = 2$$

[in this case, at least a heuristic to get some tree]

**DEF:** rooted tree $T$ with branch-length $\delta$

at all leaves have same distance

to root $f$

&  $\Rightarrow \delta(u) \leq \delta(v)$

is called <u>ultrametric tree</u>

**DEF:** Distance $D: X \times X \to \mathbb{R}_{\geq 0}$ is <u>ultrametric</u>

if

1) $D(x,y) = 0 \iff x = y$

2) $D(xy) = D(yx)$

3) instead of usual $\Delta$ - inequ.:

$$D(x,y) \leq \max \{ D(x,z), D(y,z) \} \quad \forall x, y, z \in X$$

Let $D: X \times X \to \mathbb{R}_{\geq 0}$ be a map that satisfies D1/D2.

<u>Lemma:</u>

$\begin{bmatrix} 3 \text{ point} \\ \text{condition} \end{bmatrix}$

Then, $D$ is ultrametric

$\iff$ the two largest dist. among

$D(x,y), D(x,z), D(y,z)$ are equal
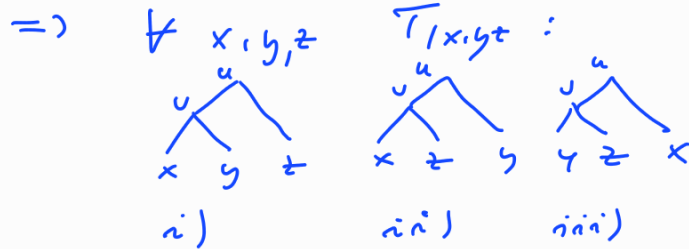
$$\forall x, y, z \in X$$

<u>proof</u> : <span style="color:red">Exercise</span>

# Thm

∃ ultrametric tree $T$ with branch length $\delta$
that represents $D: X \times X \to \mathbb{R}_{\geq 0}$
$\iff$ $D$ is ultrametric.

## proof:

"$\Rightarrow$" $(T, \delta)$ ultrametric tree

$\Rightarrow$ $\forall x, y, z$ $T_{|x,y,z}$ :



i)      ii)      iii)

case i)    $lca(xy) = v <_T lca(xz) = lca(yz) = u$

Since $(T, \delta)$ is ultrametric
& it __represents__ $D$ we have:   $D(xy) = \delta(v) \leq \delta(u) = D(xz) = D(y,z)$
$\Rightarrow$ 2 largest Dist are equal .

$\xrightarrow{\text{3 point cond}}$ $D$ ultrametric. $\quad$ [analog case ii/iii]

"$\Leftarrow$"

$X = \{1 .. n\}$ $\qquad\qquad$ $D \begin{pmatrix} 1224 1 4 4 1 \end{pmatrix} i$ $\quad D_1' = 1 < D_2' = 2 < D_3' = 4$ $\searrow$ Exmpl.

Take $i$-th row of $D$ $\quad (\hat{=} i$-th leaf in $T)$

Assume there are $m \geq 1$ diff. values
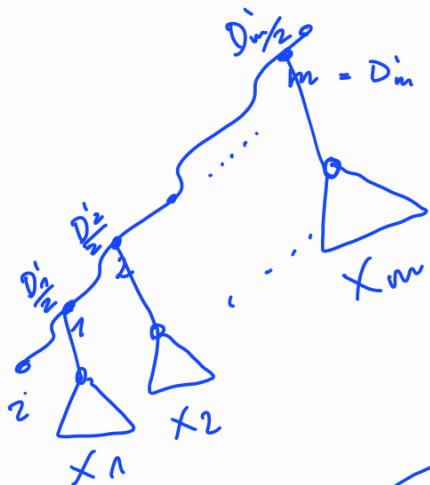$D_1' < D_2' < .. < D_m'$ in $i$th row of $D$

$\Rightarrow$ add path to $T$

$$1 \ldots j \ldots \ldots n$$

$i$-row — $D(i_{ij})$ _____

$D_{ij} \stackrel{!}{=} D_\ell'$

$\Rightarrow$ we can partition $X$ into $X_1 \ldots X_m$

with $X_\ell = \{ j \mid D(i,j) = \ell \}$

with $\ell \in \{ D_1' \ldots D_m' \}$

$\Rightarrow$



$\Rightarrow$ all $D(i,j)$, $1 \leq j \leq n$ are represented in $T$ constructed so-far.

$\Rightarrow D(xy) = D(kz)$ since $x, y \in X_\ell$

Rest for $D(yz)$ still to be computed. $\Rightarrow$ correctly represented

$\Rightarrow D_{\ell'} > D_\ell$

& $D(xy) < D(xz)$

$\stackrel{3p.cond}{\Rightarrow} D(yz) = D(xz) = \dfrac{D_{\ell'}}{2}$

now recurse on each class $X_\ell$

$\Rightarrow$ correctly represented

by ind. on steps

$\Rightarrow$ we get tree „$T + \delta$" ot $D$ is represented by it.

we never obtain:



s.t. $\delta' > \delta$

$\Rightarrow$ get ultrametric tree $T$ that represents $D$ $\quad \square$

# Drawbacks

constant - molecular - clock assumption:
"speed of evolution" $\hat{=}$ mutation rate is constant,
& thus, the same along all branches, ie.
path dist. from every leaf to root is the same.

Different rates ( reflected as branch length):

possible
true history:

branch length $\hat{=}$ rates

D

| | a | b | c | d |
|---|---|---|---|---|
| a | 0 | 7 | 5 | 6 |
| b | | 0 | 4 | 9 |
| c | | | 0 | 7 |
| d | | | | 0 |

UPGMA:

$\longrightarrow$

b c   a   d

# Neighbor - Joining (NJ)

- no "const. mol. clock" assumption

- Based on concept of minimum evolution, i.e., resulting tree will have min total branch length.

- quite fundamental approach!

IDEA:   start with "star tree"   ✳

&   stepwise seperate vertices that are
"quite" close to each other
& "quite" far away from rest

until **fully resolved** unrooted tree has been built.

$\hat{=}$ binary, ie each
inner vertex has degree 3



$i, j$ neighbors but

$D_{ij} = 13 > D_{jk} = 12$

but $i, j$ are together "farer"
away from rest.

**DEF:** given $n \times n$ dist. matrix $D$. Then $D^*$ denotes neighbor-joining matrix defined by

$$D^*_{i,j} = (n-2) \cdot D_{ij} - TotalDist_D(i) - TotalDist_D(j)$$

*"degree of freedom"* where $TotalDist(x) =$ sum of distances from $x$ to all other $n$ taxa, 2 taxa in $D_{ij}$

| D | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 13 | 21 | 22 |
| 2 |   | 0 | 12 | 13 |
| 3 |   |   | 0 | 13 |
| 4 |   |   |   | 0 |

> $D_{23}$ min but not neighbors in tree !



Sum branch-length between $i,j$
$$= D(i,j)$$

TotalDist $(1) = 56$
$(2) = 38$
$(3) = 46$
$(4) = 48$

$D^*$

| | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | -68 | -60 | -60 |
| 2 |   | 0 | -60 | -60 |
| 3 |   |   | 0 | -68 |
| 4 |   |   |   | 0 |

$$D^*_{12} = (4-2) \cdot 13 - 56 - 38 = -68$$

**"Intuition"** $D^*(i,j) \hat{=}$ "common net divergence"
$\longrightarrow$ TAKE lowest one.

<u>DEF:</u> $\Delta_{ij} = \left( \left| \text{Total Dist}_D(i) - \text{Total Dist}_D(j) \right| \right) \cdot \frac{1}{n-2}$

Take $i,j$ with min $D^*_{ij}$ & $\boxed{\text{adjust } D}$
by "joining" $i \& j$ column/row to new

m-th column/row $\rightarrow$ $D_{km} = D_{mk} = \dfrac{D_{ik} + D_{jk} - D_{ij}}{2}$

<u>"Intuition"</u> $(i,j)$ "joined as neighbors" in tree



$i$    $K$     $i$   $\delta_i$     $k$

$j$   $D_{ij}$   $D_{ik}$   $D_{ik'}$    $k'$

     $D_{jk}$   $D_{jk'}$    $j$   $\delta_j$   want to know $\delta_i, \delta_j$ = branch length.

$\Delta_{ij} = \left( D_{ij} + D_{ik} + D_{ik'} - \left( D_{ij} + D_{jk} + D_{jk'} \right) \right) \frac{1}{4-2}$    $// \; n=4$

$= \left( D_{ik} - D_{jk} + D_{ik'} - D_{jk'} \right) \cdot \frac{1}{2} = \delta_i - \delta_j$

$D_{ik} - D_{jk} = \delta_i + c - (\delta_j + c) = \delta_i - \delta_j$

$i \; \circ \overset{\delta_i}{\diagdown} \overset{c}{\underset{\phantom{}}{\diagup}} \circ \, k$
$j \; \circ \underset{\delta_j}{\diagup}$

$[\text{Analog: } D_{ik'} - D_{jk'} = \delta_i - \delta_j ]$

$D_{ij} = \delta_i + \delta_j + 2\delta_i - 2\delta_j$

$\boxed{\begin{array}{l} \frac{1}{2} \left( D_{ij} + \Delta_{ij} \right) = \frac{1}{2} \left( \underbrace{(\delta_i + \delta_j)}_{= D_{ij}} + (\delta_i - \delta_j) \right) = \delta_i \\[2mm] \frac{1}{2} \left( D_{ij} - \Delta_{ij} \right) = \frac{1}{2} \left( \underbrace{\delta_i + \delta_j}_{= D_{ij}} - \delta_i + \delta_j \right) = \delta_j \end{array}}$

# ALGO:

[runtime
$\mathcal{O}(n^3)$ ]

NeighborJoining (D)

IF D = 1×1 matrix stop
ELSE
1  construct $D^*$ from D
2  Take $i, j$ st $D_{ij}^* \xrightarrow{!} \min$
3  Compute $\Delta_{i,j}$
4  Compute $\delta_i$ & $\delta_j$
5  "Refine" tree    // Starting from star tree
6  D ← adjusted D [join i&j]
7  NeighborJoining (D)

## Exmpl:

first assume we dont know anything about tree => start with star tree



## Neighb). (D):

| D | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 13 | 21 | 22 |
| 2 |   | 0 | 12 | 13 |
| 3 |   |   | 0 | 13 |
| 4 |   |   |   | 0 |

**Step 1:**

$$\longrightarrow$$

| $D^*$ | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | -68 | -60 | -60 |
| 2 |   | 0 | -60 | -60 |
| 3 |   |   | 0 | -68 |
| 4 |   |   |   | 0 |

**Step 2**  may choose (1,2) or (3,4)
decide here for (1,2)

**Step 3**

$$\Delta(i,j) = (\text{Total Dist } D(i) - \text{Total Dist}(j)) \cdot \frac{1}{2}$$
$$= (56 - 38) \cdot \frac{1}{2} = 9$$

**Step 4**

$$\delta_i = \frac{1}{2}(D_{ij} + \Delta(i,j)) = \frac{1}{2}(13 + 9) = 11$$
$$\delta_j = \frac{1}{2}(D_{ij} - \Delta(i,j)) = \frac{1}{2}(13 - 9) = 2$$

**Step 5**



refine

**Step 6**

D adjusted:

$$\boxed{D_{km} = D_{mk} = \frac{D_{ik} + D_{jk} - D_{ij}}{2}}$$

| D | (12) | 3 | 4 |
|---|---|---|---|
| (12) | 0 | 10 | 11 |
| 3 |   | 0 | 13 |
| 4 |   |   | 0 |

$$D_{(12)3} = \frac{D_{13} + D_{23} - D_{12}}{2} = \frac{21 + 12 - 13}{2} = 10$$

**Step 7**  recurse on D.

**Def:** Distance $D: X \times X \to \mathbb{R}_{\geq 0}$ is <u>additive</u> if

1) $D(x,y) = 0 \iff x = y$
2) $D(xy) = D(yx)$
3) instead of usual $\Delta$-inequ.: $\forall\ x, y, a, b \in X$:

$$D(xy) + D(ab) \leq \max \{ D(xa) + D(yb),$$
$$D(xb) + D(ya) \}$$

[largest of must be equal]
$\triangleq$ 4 point condition

<u>if T looks like:</u>



$$\boxed{\textbf{Ex:} \quad D \text{ u metric} \overset{?}{\underset{\Leftarrow}{\Rightarrow}} \text{additiv}}$$

other example:

| D | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 1 | 0 | 3 | 4 | 3 |
| 2 | 3 | 0 | 4 | 5 |
| 3 | 4 | 4 | 0 | 2 |
| 4 | 3 | 5 | 2 | 0 |

UPGMA



NJ



No ultrametric - why?

$$|\{ D(1,3), D(1,4), D(3,4)\}| = 3$$

not additive:

$$D(1,2) + D(3,4) = 3 + 2 = 5$$
$$D(1,3) + D(2,4) = 4 + 5 = 9 \quad \Rightarrow 9 \neq \max(5,7).$$
$$D(1,4) + D(2,3) = 3 + 4 = 7$$

**Thm**

$\exists$ "additive" tree $T$ with branch length $\sigma$
that represents $D : X \times X \to \mathbb{R}_{\geq 0}$
[that is, $D(i,j) = \sum$ weights $\sigma$ along path connecting
$i, j$ in $T$]
$\iff$ $D$ additive metric

Drawback : trees may have negative branch-length.

# Summary - Distance based methods

- Distance base method work well on ultrametric or additive distances.
- in any other case quite useful as heuristics.

- when fetching sequ. alignments
  $ACGT$
  $ACGC$ ...

  $\longrightarrow$ get distances (but we loose inform in dist. matrix)

  $\implies$ we cannot say anything about the <u>ancestral state</u> !

# CHARACTER-BASED METHODS

Before "DNA-age", half a century ago, researchers constructed tree from morphological characters.

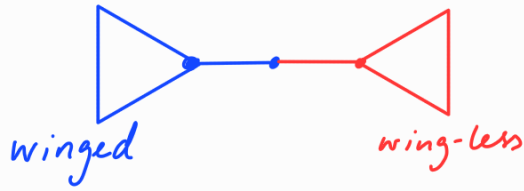| | Wings | # legs |
|---|---|---|
|  | yes | 6 |
|  | NO | 6 |
|  | NO | 42 |

stick insects

giant centipide →

Aim: Reconstruct phylogeny from characters

Input: n x m matrix ( n taxa, m characters)

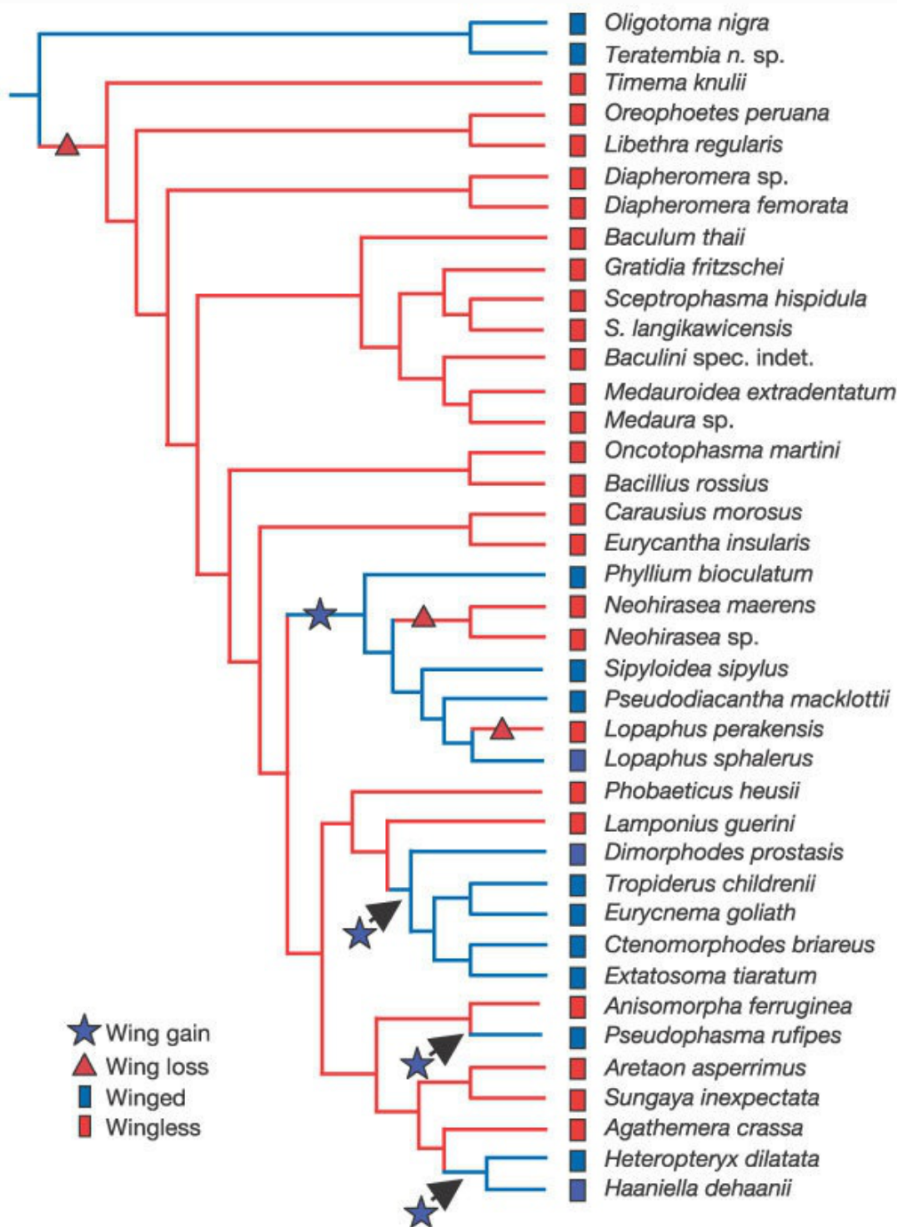output: tree in which taxa with similar character-values occur near each other.

stick-insects



winged          wing-less

# DOLLO'S principle of irrevesibility (1893):

Evolution doesn't reinvent the same organ (e.g. insect wings) [evolution is efficient]

## Stick-Insects Phylogeny:



☆/△  = 7 times where wings were gained or lost in stick insects alone!

Oligotoma nigra
Teratembia n. sp.
Timema knulii
Oreophoetes peruana
Libethra regularis
Diapheromera sp.
Diapheromera femorata
Baculum thaii
Gratidia fritzschei
Sceptrophasma hispidula
S. langikawicensis
Baculini spec. indet.
Medauroidea extradentatum
Medaura sp.
Oncotophasma martini
Bacillius rossius
Carausius morosus
Eurycantha insularis
Phyllium bioculatum
Neohirasea maerens
Neohirasea sp.
Sipyloidea sipylus
Pseudodiacantha macklottii
Lopaphus perakensis
Lopaphus sphalerus
Phobaeticus heusii
Lamponius guerini
Dimorphodes prostasis
Tropiderus childrenii
Eurycnema goliath
Ctenomorphodes briareus
Extatosoma tiaratum
Anisomorpha ferruginea
Pseudophasma rufipes
Aretaon asperrimus
Sungaya inexpectata
Agathemera crassa
Heteropteryx dilatata
Haaniella dehaanii

☆ Wing gain
▲ Wing loss
▌ Winged
▌ Wingless

What happened?

Evolution did not reinvent wings from scratch

"genes switched on/off

⟹ wings yes/no"

We can use genetic data as characters instead!

| Species | Alignment |
|---------|-----------|
| CHIMP | ACGTAGGCCT |
| HUMAN | ATGTAAGACT |
| SEAL | TCGAGAbCAC |
| WHALE | TCGAAAGCAT |

} n taxa

$\underbrace{\qquad\qquad}_{\text{m characters}}$

Given tree & reconstruct most-likely ancestral sequences.



?? ... ?

??? .. ?     ??? ... ?

ACGTAGGCCT        ATGTAAGACT        TCGAGAbCAC        TCGAAAGCAT
   CHIMP             HUMAN              SEAL              WHALE

parsimony - SCORE = sum of Hamming dist. along edges.



ACGAAAGCCT

ACGTAAGCCT          1          2          TCGAAAGCAT

1          2                    2          0

ACGTAGGCCT        ATGTAAGACT        TCGAGAbCAC        TCGAAAGCAT
   CHIMP             HUMAN              SEAL              WHALE

parsimony - SCORE = 8

Now we have info about ancestral states!

<u>Ockham's razor</u>    "simplest explanation is usually
(1287-1347)              best one"

## <u>SMALL PARSIMONY problem:</u>

Given rooted tree $T$, each leaf labeled by
                              string of length $m$
Find labeling (= string of length $m$) for
all internal nodes that minimize parsim. score.

IF "position" of string are independent THEN

Given rooted tree $T$, each leaf labeled by
                              single symbol
Find labeling (= single symbol ) for
all internal nodes that minimize parsim. score.
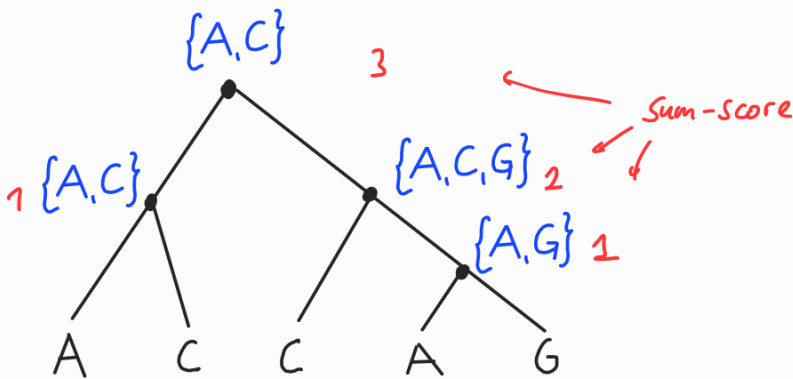
# FITCH - ALGO (Walter M Fitch 1971)

Given binary tree $T$ with leaf labels

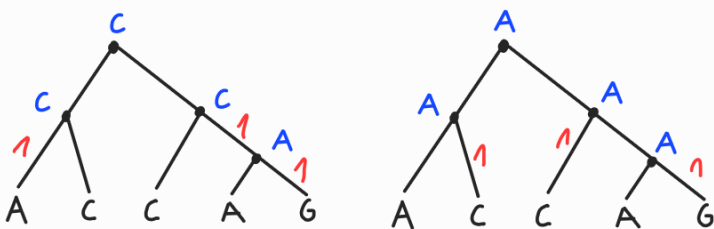let $X_v$ set of characters assigned to $v$

**bottom up:** assign to each internal vertex $v$ with children $u, w$ the state set

$$X_v := \begin{cases} X_u \cup X_w & , \text{ if } X_u \cap X_w = \emptyset \\ X_u \cap X_w & , \text{ else.} \end{cases}$$

until all vertices have been visited



Tree (bottom-up): root {A,C}  3 ← Sum-score

{A,C} 1, {A,C,G} 2, {A,G} 1

leaves: A   C   C   A   G

**Possible solutions:**



Solution 1: C / C, C / A, A leaves A C C A G (scores 1, 1, 1)

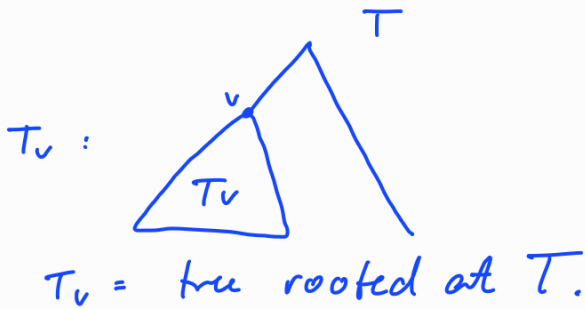Solution 2: A / A, A / A, A leaves A C C A G (scores 1, 1, 1)

**TOP DOWN:**

root $\rho$ take any character $l_\rho \in X_\rho$

Then for every internal node $v$

if $L_w \in X_v$ put $l_v = l_w$

else $l_v$ any of $X_v$

# SANKOFF - ALG  (David sankoff 1971)

## Dynamic Prog!

$T_v$ :

$T_v$ = tree rooted at $T$.

$S_a(v)$ = min parsim. score of $T_v$ over all labelings of $T$ assuming $v$ is labeled with symbol $a$

⇒ min pars. score of $T$ = $S_a(root)$ over all symbols $a$

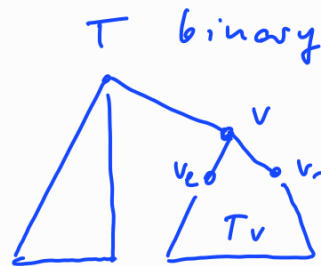To implement:  Def for symbols $a, b$:

$$\delta_{a,b} = \begin{cases} 0, & a = b \\ 1, & else \end{cases}$$

// also known as $\mathbb{1}_{ab}$ indicator fct
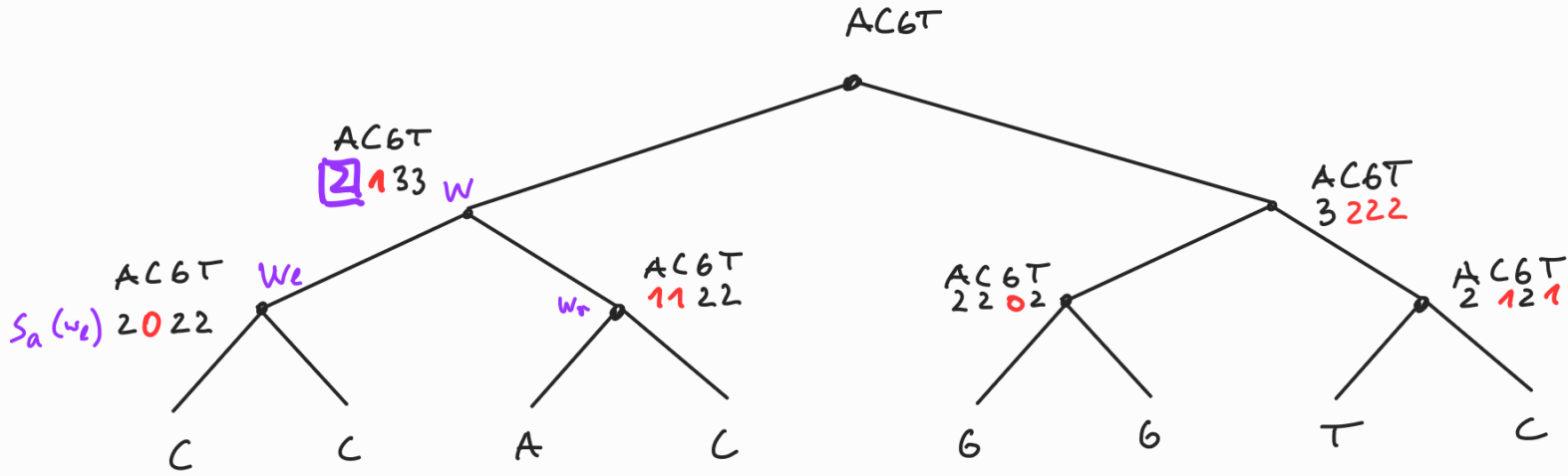
## recurrence relation:
(can be generalized to non-binary)

$T$ binary

left/right child of $v$ is $v_\ell / v_r$

$$S_a(v) = \min_{all\ symbols\ y} \left\{ S_y(v_\ell) + \delta_{ay} \right\} + \min_{all\ symbols\ y} \left\{ S_y(v_r) + \delta_{ay} \right\}$$

where ∀ leaves $\ell$:  $S_a(\ell) = \begin{cases} 0, & symbol\ of\ \ell\ is\ a \\ \infty, & else \end{cases}$

[detailed correctness as exercise]

Top tree:

ACGT (root)

ACGT (left), ACGT (right)

| a | ACGT |
|---|------|
| $S_a(v)$ | 2 0 2 2 |

v: ACGT 2 0 2 2

ACGT 11 22

ACGT 22 02

A ACGT 2 12 1

Leaves: C  C  A  C  G  G  T  C

Bottom tree:

ACGT (root)

ACGT [2] 133  w

ACGT We — $S_a(w_\ell)$ 2 0 2 2

$w_r$: ACGT 11 22

ACGT 3 222

ACGT 22 02

A ACGT 2 12 1

Leaves: C  C  A  C  G  G  T  C

__W = A__

$$S_A(w) = \min\{ S_A(w_\ell) + \delta_{AA},\ S_C(w_\ell) + \delta_{AC},\ S_G(w_\ell) + \delta_{AG},$$
$$S_T(w_\ell) + \delta_{TA} \}$$
$$+ \min\{ S_A(w_r) + \delta_{AA},\ S_C(w_r) + \delta_{AC},\ S_G(w_r) + \delta_{AG},$$
$$S_T(w_r) + \delta_{TA} \}$$

$$= \min(2+0,\ 0+1,\ 2+1,\ 2+1)$$
$$+ \min(1+0,\ 1+1,\ 2+1,\ 2+1)$$

$$= \quad 1 \quad + \quad 1 \qquad = \boxed{2}$$

analog:  $S_w(C) = 1$
$S_w(G) = 3$
$S_w(T) = 3$

min pars. score = 3

ACGT
5**3**44

ACGT
2**1**33 $W$

ACGT
$S_x(w_l)$ 2**0**22    $W_l$

ACGT
**1**1 22   $W_r$

ACGT
3 **222**

ACGT
22**0**2

A CGT
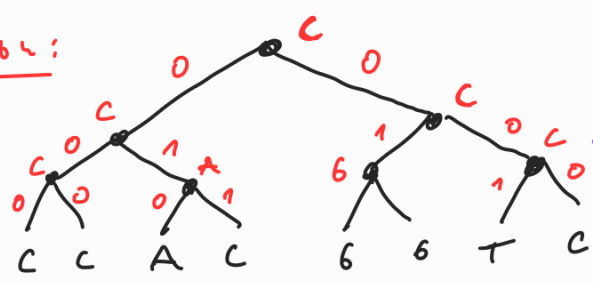2 **12 1**

C    C    A    C    G    G    T    C

To reconstruct ancestral state:    BACKTRACKING

possible solution:                              [exercise]


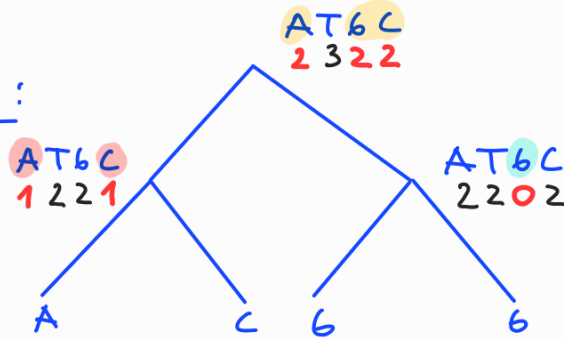
T not possible
since then

# Fitch vs Sankoff Alg (both O(nm) runtime)



Fitch:

A C G G → ACG (root), AC, G

Sankoff:

```
        ATGC
        2 3 2 2
       /        \
   ATGC          ATGC
   1 2 2 1       2 2 0 2
   /    \        /    \
  A      C      G      G
```

=> essentially "identical" in nature.

## LARGE PARSIMONY problem:

> **Find** rooted tree T, **for given** strings of length m
> **with** labeling (= string of length m) for all nodes that minimize parsim. score.

NP-hard!

=> heuristics needed!

[ not part here]

# CONSENSUS METHODS

**IDEA**      Given a collection of trees $T_1 .. T_k$
Find common "supertree" that
summarizes the information provided
by $T_1 ... T_k$ in a "best" way.

     **Why?** ▶ Different datasets or tree-finding
                                       methods

         ⟹    Different trees.

         ⟹    combine trees to get more
                reliable answer.

     ✏ comput. expensive methods can yield
highly accurate trees on small (overlapping)
data sets

              ⟹ Find 1 tree to represent
                   entire data set.

**Exmpl:** say we have only "partial" information about "similarities" between taxa $A, B, C, D, E$ in the form:     A & B are closer related than   A compared to C
                                                                                                                                & B compared to C
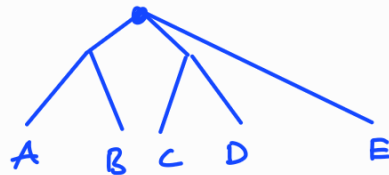
& C & D — " — than   C — " — E
                                                        D — " — E



$A$ $B$ $C$    &    $C$ $D$ $E$

**Q:** Is there a common tree that reflects both relationships?

**A:** yes:



$A$ $B$ $C$ $D$ $E$

**DEF:** (ROOTED) TRIPLE   $ab|c$ = binary rooted tree


$a$ $b$ $c$

$ab|c$ **displayed** by rooted tree if


$a$ $b$ $c$

$(\iff lca(ab) <_r lca(ac) = lca(bc))$

$ab$-path
   does not intersect
with $bc$-path

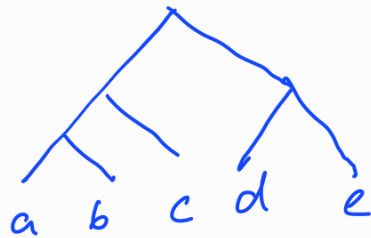Given set $R$ of triples compatible if exists tree that **displays** all triples in $R$.

<u>Exmpl:</u>    ablc,  acld,  delb



=)  common supertree



<u>Exmpl:</u>  For   ablc  &  cbla   no tree!


<u>OBSERVATION:</u>

given set R of triples  &  xylz ∈ R

IF  exists tree for R  =>  x & y  cannot be descendants
                          of  two different children
                          of root.



not possible!

=> Central idea:   determine for potential tree
                   the set of leaves that are descendents
                   of each child of root.
                   [then recurse on children]

$\Rightarrow$ Find partition $X_1 \dots X_\ell$ of $X$

st



$v_1 \quad \dots \quad v_\ell$

$X_1 \quad \dots \quad X_\ell$

($|X_i| = 1$ identify $v_i$ with $x \in X_i$ )

Example:

T

$a \quad b \quad c \quad d \quad e$

$X = X_1 \cup X_2 \cup X_3$ with $X_1 = \{a, b\}$
$X_2 = \{c, d\}$
$X_3 = \{e\}$.

| ∀ triples $xy|z$ such a partition must satisfy:
∘ $xy$ are in same set $X_i$ ( $\nearrow$ ⊛ )

___

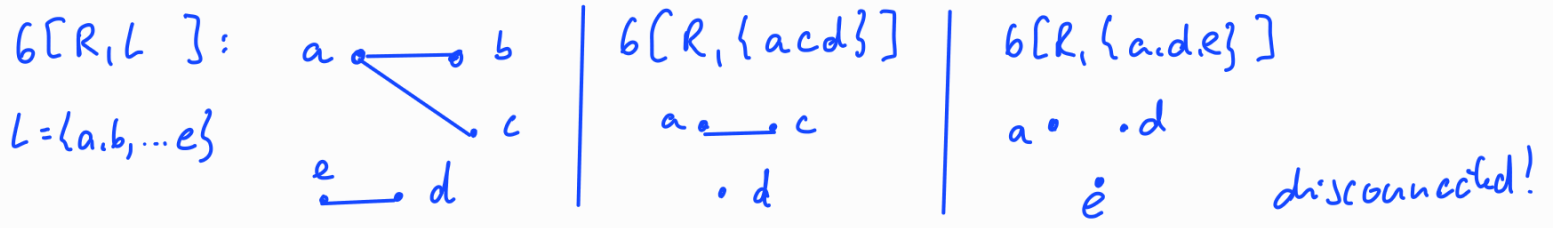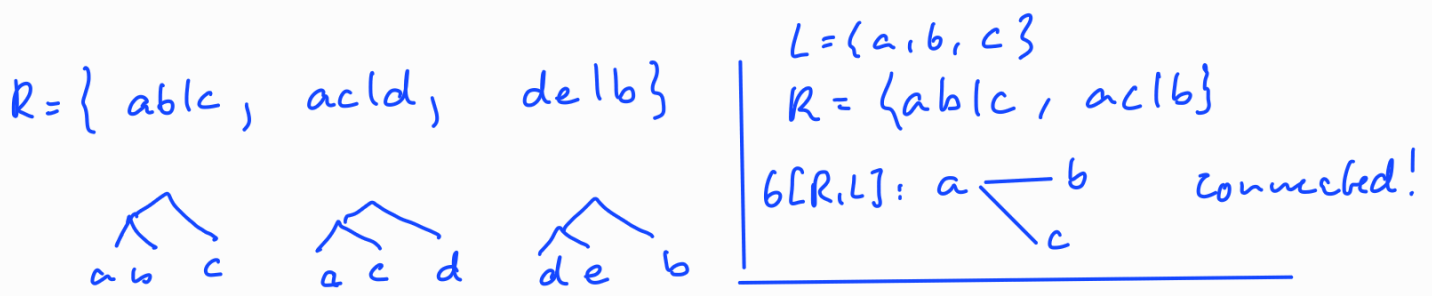DEF    set of triples $R$, $L$ set of leaves.

$$R_{|L} := \{ xy|z \in R : x, y, z \in L \}$$

___

DEF:    Compatibility-graph $G[R, L]$

$R$ = set of triples , $L$ = set of leaves.

Then $G[R, L]$ has leaf set $L$ &

$\{xy\}$ is an edge $\iff$ $\exists \, xy|z \in R_{|L}$

$R = \{ ab|c, \ ac|d, \ de|b \}$



$L = \{a, b, c\}$
$R = \{ab|c, \ ac|b\}$

$6[R,L]: \ a \diagdown^b_c$   connected!

$6[R, L]: \quad a \multimap b$
$L = \{a, b, ... e\}$   $\quad \diagdown c$
$\overset{e}{\multimap} d$

$6[R, \{acd\}]$
$a \multimap c$
$\cdot d$

$6[R, \{a, d, e\}]$
$a \bullet \quad \bullet d$
$\overset{\bullet}{e}$   disconnected!

## <u>ALGO</u>   (from Aho, Sagiv, Szymanski & Ullmann 1981)

Build $( R, \ v, \ T, \ L )$   // input: set of triples $R$,
                              vertex $v$ & tree $T$

   IF $(|L| = 1)$
      output: rooted tree $\overset{\bullet}{x}$   $(x \in L)$

   IF $(|L| = 2)$
      output: $\overset{v}{\underset{x \ y}{\wedge}}$   $(x, y \in L)$

   IF $(|L| \geq 3)$
      construct $6[R, L]$
      Let $L_1 ... L_k$ vertex set of conn. components of $6[R, L]$
      IF $(k = 1)$ stop & output "$R$ not compatible"
      FOR $(i = 1 ... k)$
         call   BUILD$( R, v_i, T_i, L_i )$
         IF (BUILD$(R, v_i, T_i, L_i)$ outputs tree $T_i$)
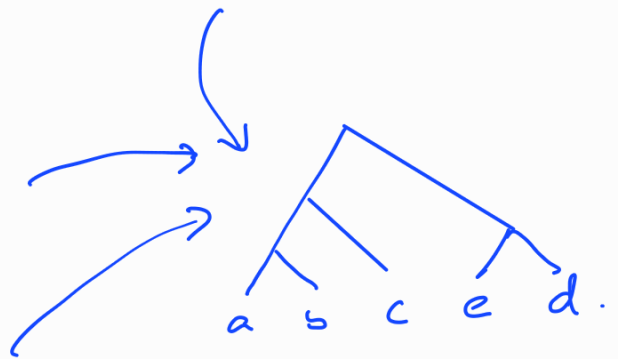            attach $T_i$ to $v$ via edge $\{v_i, v\}$

**Exmpl:** $R = \{\, ab|c, \; ac|d, \; de|b\,\}$

call of BuilD: $6[R,L]$, $L = \{a,b,\dots e\}$



$L_1$    $L_2$

$BuilD(R|_{L_1}\cdots)$



a b c

$BuilD(R|_{L_2},\cdots)$

e d     a b c e d.

**theorem**    BuilD runs in $O(|L||R|)$ time
& is correct

**proof [SKETCH]**

R comp. $\Rightarrow$ R' comp. $\forall R' \subseteq R$.

$\Rightarrow \exists \tilde{T}'$ for R'

if $6[R',L']$ conn. $\Rightarrow \exists T'$ 



& $ab|c \in R'$

but $ab|c$ not displ. by $T'$ ⚡
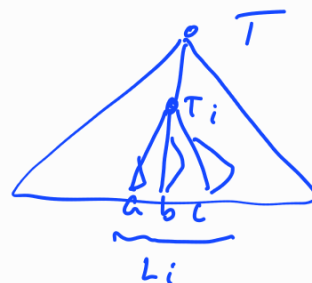
$\Rightarrow 6[R',L']$ disconn. in each step.

remains to show  T displays each triple in R.

Let  ablc ∈ R    &    $T_i$ "min. subtree" in T that
                                    contain a,b,c

6 [$R_1 L_i$] => a ⟍ᵇ
                      c

    => ab ir same comp.

    ⟶ 

$T_i$
a b   c

T
$T_i$
a b c
$L_i$

[runtime: Exercise]
□