# Sequence evolution & phylogenetics

Christophe Dessimoz

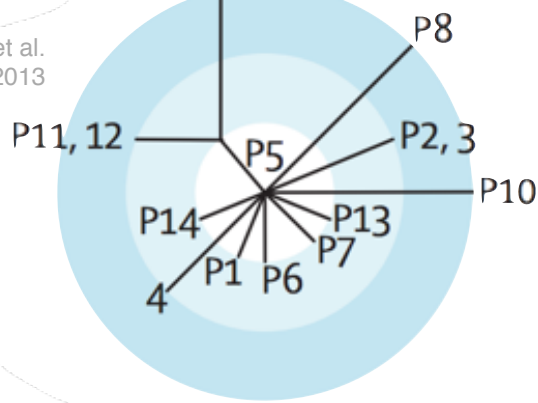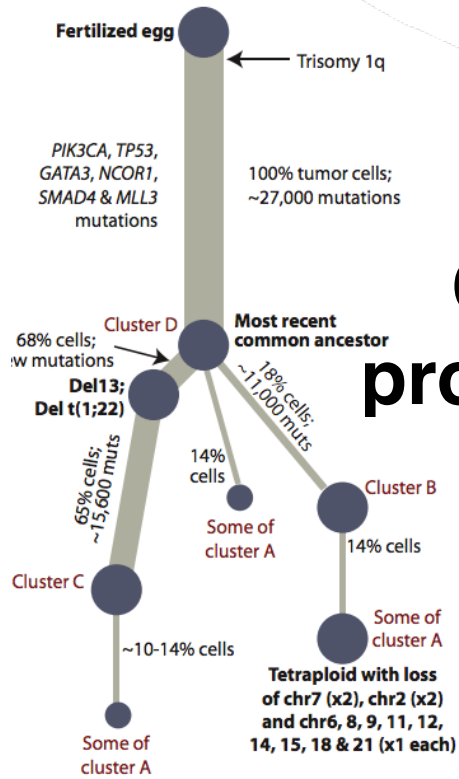http://lab.dessimoz.org 🐦@cdessimoz

# Today's lecture

- Pairwise distance estimation

- Tree thinking & terminology

- Tree inference

  - Methods

  - Confidence

  - Rooting

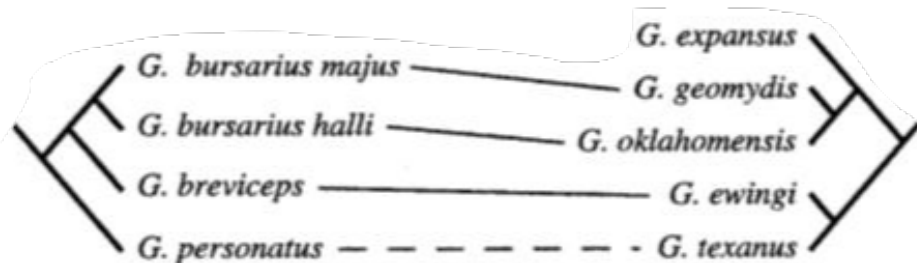Harris et al.
Lancet Infect Dis 2013

P8

P11, 12  P5  P2, 3
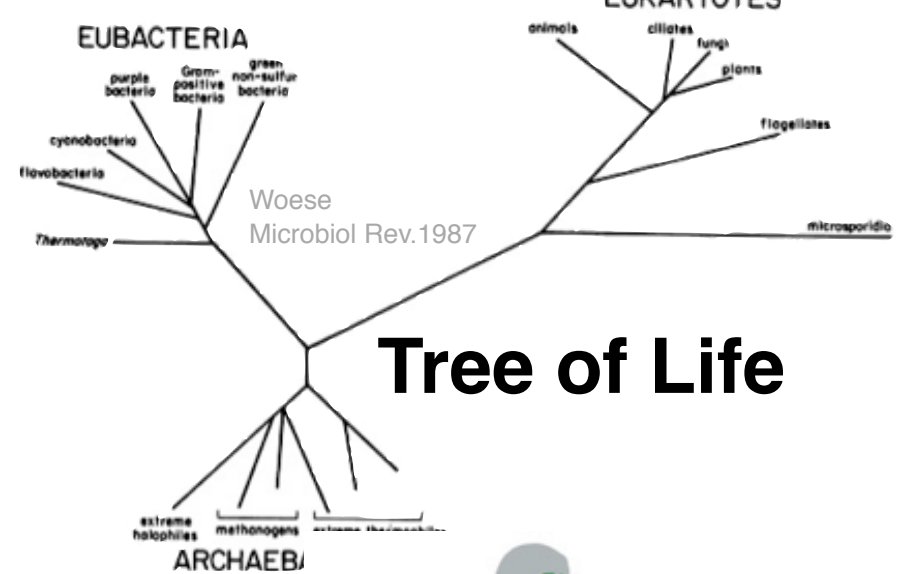
P14  P10

P13

P1  P7

P6

4

**Epidemiology**

EUBACTERIA

purple bacteria
Gram-positive bacteria
green non-sulfur bacteria

cyanobacteria

flavobacteria

Thermotoga

Woese
Microbiol Rev.1987

EUKARYOTES

animals
ciliates
fungi
plants

flagellates

microsporidia

extreme halophiles

methanogens
extreme thermophiles

ARCHAEBA

**Tree of Life**

**Fertilized egg**

Trisomy 1q

*PIK3CA, TP53,
GATA3, NCOR1,
SMAD4 & MLL3*
mutations

100% tumor cells;
~27,000 mutations

Cluster D

**Most recent
common ancestor**

68% cells;
w mutations

~18% cells;
~11,000 muts

**Del13;
Del t(1;22)**

14%
cells

65% cells;
~15,600 muts

Cluster B

14% cells

Cluster C

Some of
cluster A

Some of
cluster A

~10-14% cells

**Tetraploid with loss
of chr7 (x2), chr2 (x2)
and chr6, 8, 9, 11, 12,
14, 15, 18 & 21 (x1 each)**

Some of
cluster A

Nik-Zainal et al. Cell 2012

**Cancer
progression**

Katrin M. Weir/
Ed Marcotte Lab

**Model systems**

**Host/pathogen
co-evolution**

*G. bursarius majus*  *G. expansus*

*G. bursarius halli*  *G. geomydis*

*G. breviceps*  *G. oklahomensis*

*G. personatus*  *G. ewingi*

*G. texanus*

Legendre et al. Syst Biol 2002

**Cell lineage
trees**

NK cell

T cell

MPP, LMPP
or ELP

CLP

B cell

Ceredig et al. Nat Rev Immunol 2009
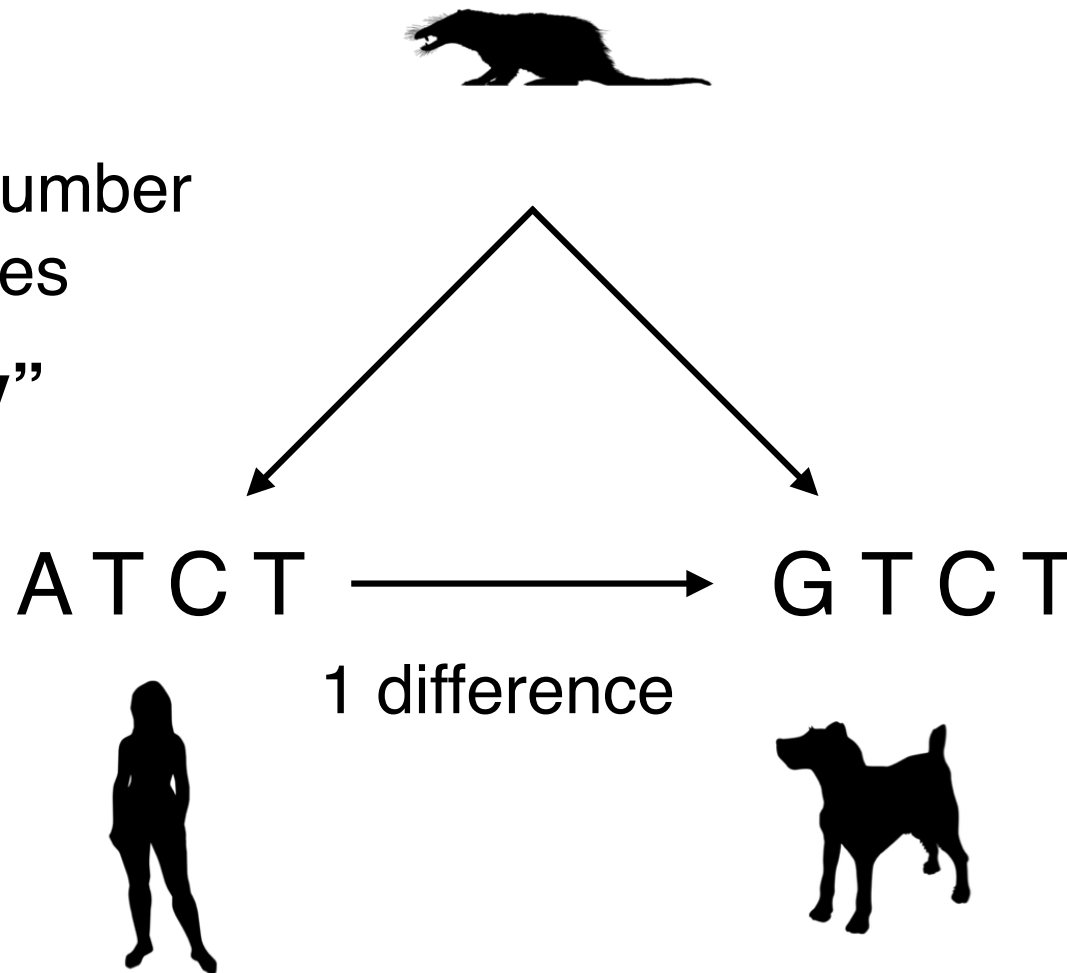
# Pairwise distance estimation & Markov models

# How to measure the amount of evolution?



Idea #1:
count the number of differences

"**similarity**"

A T C T $\longrightarrow$ G T C T
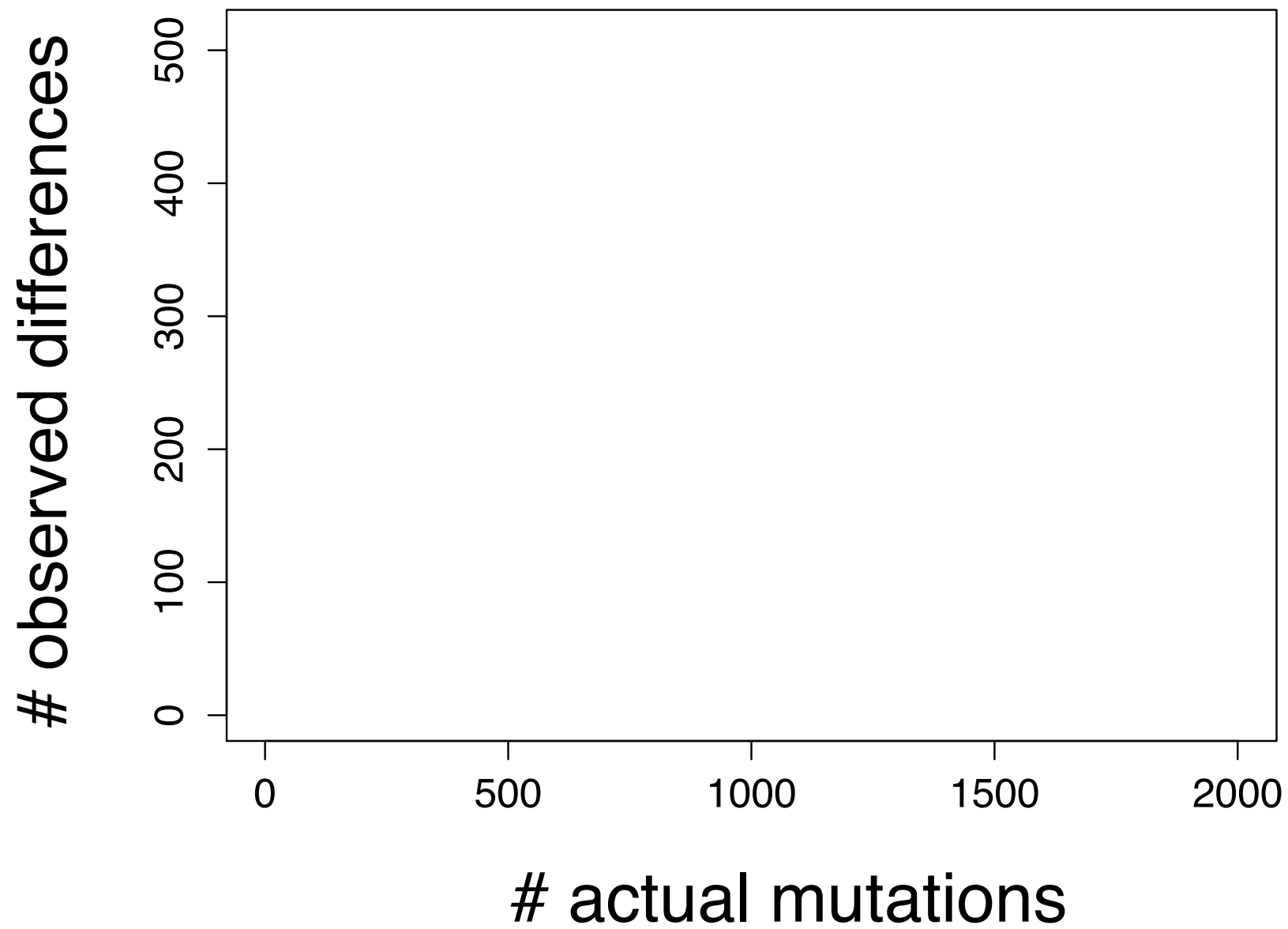
1 difference

If we normalise with respect to length:

25% differences (75% "identity")

# Problem with counting differences

|  | actual #<br>of substitutions | # of observed<br>differences |
|---|---|---|
| → A T C T | | |
| → A T C T | 0 | 0 |
| G T C T | 1 | 1 |
| G T C T | 1 | 1 |
| C T C T | 2 | 1 |
| A T C T | 3 | 0 |

# Let's simulate sequence evolution

- Generate a random sequence of 0 and 1 of length 1000

- Repeat 2000 times:

  - Mutate a random position in the sequence

  - Compute the number of difference between the resulting sequence and the original sequence and store this number in a table

- Plot the values stored in the table

# Markov model of evolution

- Every site evolves independently, with prob. of mutation only depending on present state (no memory).

- Probabilities of mutation at a given distance are expressed by transition matrix.

$M^I =$

| ↙ | A | C | G | T |
|---|---|---|---|---|
| A | 0.900 | 0.033 | 0.033 | 0.033 |
| C | 0.033 | 0.900 | 0.033 | 0.033 |
| G | 0.033 | 0.033 | 0.900 | 0.033 |
| T | 0.033 | 0.033 | 0.033 | 0.900 |

After "one unit" of evolution, the probability that an A mutates into a C is given by the corresponding entry in the matrix:

$p(A \rightarrow C \mid d=1) = M^1[A \rightarrow C] = 0.033$

$M^2=$

| ↘ | A | C | G | T |
|---|---|---|---|---|
| A | **0.900** | **0.033** | **0.033** | **0.033** |
| C | 0.033 | 0.900 | 0.033 | 0.033 |
| G | 0.033 | 0.033 | 0.900 | 0.033 |
| T | 0.033 | 0.033 | 0.033 | 0.900 |

x

| ↙ | A | C | G | T |
|---|---|---|---|---|
| A | **0.900** | 0.033 | 0.033 | 0.033 |
| C | **0.033** | 0.900 | 0.033 | 0.033 |
| G | **0.033** | 0.033 | 0.900 | 0.033 |
| T | **0.033** | 0.033 | 0.033 | 0.900 |

$M^2=$

| ↙ | A | C | G | T |
|---|---|---|---|---|
| A | **0.813** | 0.062 | 0.062 | 0.062 |
| C | 0.062 | 0.813 | 0.062 | 0.062 |
| G | 0.062 | 0.062 | 0.813 | 0.062 |
| T | 0.062 | 0.062 | 0.062 | 0.813 |

$M^2[A,A] = M^1[A{\rightarrow}A] * M^1[A,{\rightarrow}A] +$
$M^1[A{\rightarrow}C] * M^1[C{\rightarrow}A] +$
$M^1[A{\rightarrow}G] * M^1[G{\rightarrow}A] +$
$M^1[A{\rightarrow}T] * M^1[T{\rightarrow}A]$

$M^\infty=$

| ↙ | A | C | G | T |
|---|---|---|---|---|
| A | 0.250 | 0.250 | 0.250 | 0.250 |
| C | 0.250 | 0.250 | 0.250 | 0.250 |
| G | 0.250 | 0.250 | 0.250 | 0.250 |
| T | 0.250 | 0.250 | 0.250 | 0.250 |

# Distance estimation

- Now that we have a model and data,
  how can we estimate the distance?

A T C T $\longrightarrow$ G T C T

$M^I =$

|   | A | C | G | T |
|---|---|---|---|---|
| A | 0.900 | 0.033 | 0.033 | 0.033 |
| C | 0.033 | 0.900 | 0.033 | 0.033 |
| G | 0.033 | 0.033 | 0.900 | 0.033 |
| T | 0.033 | 0.033 | 0.033 | 0.900 |

Let's assume that the distance is 1 "unit" of evolution.
Given the original sequence ATCT, the likelihood of
observing GTCT is:

$M^1[A,G] * M^1[T,T] * M^1[C,C] * M^1[T,T] = 0.0243$

**is that likely?!**

# A brief extrapolation on maximum likelihood (ML) parameter estimation



Ronald A. Fisher

## The likelihood function

L(parameter ; data) = p(data l parameter)

ML estimator: choose parameter that maximises the likelihood function!

*Note that L() is not a probability distribution (does not sum to 1)*

Another example:

Unfair coin. data: T T T T H T H

Model:     $p(x=T | \theta) = \theta$
           $p(x=H | \theta) = 1-\theta$

$L(\theta; data) = \theta*\theta*\theta*\theta*(1-\theta)*\theta*(1-\theta)$

Unfair coin. data: T T T T H T H

Model:    $p(x=T \mid \theta) = \theta$
          $p(x=H \mid \theta) = 1-\theta$

$L(\theta; \text{data}) = \theta*\theta*\theta*\theta*(1-\theta)*\theta*(1-\theta)$

$= \theta^5*(1-\theta)^2$

## Maximise L($\theta$; data)

$0.71 = 5/7$

# back to our example

A T C T  ⟶  G T C T



our "maximum likelihood"
distance estimate is 3 units!

Augustin Augier,
Arbre Botanique
(1801)

# Lamarck, Philosophie Zoologique , 1809

## TABLEAU

*Servant à montrer l'origine des différens animaux.*

Vers.

Infusoires.
Polypes.
Radiaires.

Insectes.
Arachnides.
Crustacés.

Annelides.
Cirrhipèdes.
Mollusques.

Poissons.
Reptiles.

Oiseaux.

Monotrèmes.

M. Amphibies.

M. Cétacés.

M. Ongulés.

M. Onguiculés.

Cette série d'animaux commençant par deux

# Darwin, Notebook B, 1837

Taf. I.

## Plantae

**Cormophyta**

Anthophyta
*Angiospermae*
*Gymnospermae*

Pteridophyta
*Lepidophyta*
*Rhizocarpeae*
*Filices*
*Calamophyta*

**Bryophyta**
*Phyllobrya*
*Thallobrya*

**Fucoideae**
*Sargassaceae*
*Laminariaceae*
*Chordariaceae*
11.

**Florideae**
*Sphaerococceae*
*Ceramiaceae*
10.

**Characeae**
12.

**Inophyta**
*Lichenes*
*Fungi*
13.
14.

**Archephyta**
*Ulva*
*Conferva*
*Desmidium*
*Nostoc*
*Codiolum*
9.

## Protista

**Myxomycetes**
*Physarum*
*Stemonitis*
*Lycogala*
*Trichia*
5.

**Spongiae**
Petrospongiae
*Siphonida*
*Ocellarida*
*Lymnorida*
*Bathrocmida*
*Turonida*

Autospongiae
*Calcispongiae*
*Silicispongiae*
*Ceratospongiae*
*Myxospongiae*
8.

**Rhizopoda**
*Radiolaria*
*Actinophryida*
*Acyttaria*
*Polythalamia*

**Flagellata**
*Peridinium*
*Euglena*
*Volvox*
3.
7.

**Myxocytoda**
*Noctilucae*
6.

**Protoplasta**
*Arcellae*
*Gregarinae*
*Aulamoebae*
*Amoebae*
2.

**Diatomeae**
*Areolatae*
*Vittatae*
*Striatae*
4.

**Moneres**
*Protogenes*
*Protamoeba*
*Vampyrella*
*Protomonas*
*Vibrio*
1.

## Animalia

**Articulata**

Arthropoda
*Trachaeata*
*Crustacea*

Vermes
*Annelida*
*Rotatoria*
*Scolecida*
*Infusoria*

**Echinodermata**
*Holothuriae*
*Echinida*
*Crinoida*
*Asterida*
16.
17.

**Coelenterata**
*Noctucalephae*
*Petracalephae*
15.

**Vertebrata**

Amniota
*Aves*
*Mammalia*
*Reptilia*

Anamnia
*Amphibia*
*Pisces*

Amphirhina
Monorhina
*Leptocardia*
19.

**Mollusca**
*Otocardia*
*Himatega*
18.

---

m  14 13 12 10 11 9 e  4 5 1 2 7 6 3 8  6 15 16 17 18 19 n

Pteridophyta  Fungi  Archephyta 11. 10. 9.
Characeae  Inophyta 12.
Cormophyta 14. 13.

Moneres
Flagellata 1.
Myxomycetes 3.
Rhizopoda 7.
Spongiae
Protoplasta
Myxocytoda 6.
Diatomeae 5.
8.

Echinodermata 16.
Articulata 17.
Mollusca 18.
Vertebrata 19.
Coelenterata 15.

x  Archephylum vegetabile  f  Archephylum protisticum  h  Archephylum animale  y

**Protista**

**Plantae**     **Animalia**

**Radix communis Organismorum**     **Moneres autogonum**

I. Feld: p m n q (*19 Stämme*)
II. Feld: p x y q (*3 Stämme*)
III. Feld: p s t q (*1 Stamm*)
stellen 3 mögliche Fälle der universalen Genealogie dar.

Monophyletischer
Stammbaum der Organismen
entworfen und gezeichnet von
Ernst Haeckel. Jena, 1866.

# Monophyletischer Stammbaum der Organismen

entworfen und gezeichnet von

Ernst Haeckel. Jena, 1866.

16S rRNA was used by Woese (1987) to group early life forms into three kingdoms



EUKARYOTES

EUBACTERIA

purple bacteria
Gram-positive bacteria
green non-sulfur bacteria
cyanobacteria
flavobacteria
Thermotoga

animals
ciliates
fungi
plants
flagellates
microsporidia

MICROBIOLOGICAL REVIEWS, June 1987, p. 221–271

**Bacterial Evolution**

CARL R. WOESE

extreme halophiles
methanogens
extreme thermophiles

ARCHAEBACTERIA

# Tree terminology and "tree thinking"

# Common Phylogenetic Tree Terminology

Terminal Nodes

Branches, edges, lineages

A

**Represent the TAXA (genes, populations, species, etc.) used to infer the phylogeny**

B

C

D

Ancestral Node or **ROOT** of the Tree

Internal Nodes or Divergence Points (represent hypothetical **ancestors** of the taxa)

E

# Phylogenetic trees diagram the *evolutionary relationships* between the taxa



No meaning to the spacing between the taxa, or to the order in which they appear from top to bottom.

This dimension either can have no scale (for 'cladograms'), can be proportional to genetic distance or amount of change (for 'phylograms' or 'additive trees'), or can be proportional to time (for 'ultrametric trees' or true evolutionary trees).

**((A,(B,C)),(D,E))** = **The above phylogeny as nested parentheses**

These say that B and C are more closely related to each other than either is to A, and that A, B, and C form a **clade** that is a **sister group** to the clade composed of D and E. If the tree has a time scale, then D and E are the most closely related.

# Which species are the closest living relatives of modern humans?

## *morphological tree*



Gorillas
Chimpanzees
Bonobos
Orangutans
**Humans**

15-30         0
**MYA**

**The pre-molecular view was that the great apes (chimpanzees, gorillas and orangutans) formed a clade separate from humans, and that humans diverged from the apes at least 15-30 MYA.**

## *mitochondrial DNA*



**Humans**
Chimpanzees
Bonobos
Gorillas
Orangutans

14         0
**MYA**

**Mitochondrial DNA, most nuclear DNA-encoded genes, and DNA/DNA hybridization all show that bonobos and chimpanzees are related more closely to humans than either are to gorillas.**

# Three types of trees

**Cladogram**

Taxon B

Taxon C

Taxon A

Taxon D

no meaning

**Phylogram**

6

1

1

Taxon B

Taxon C

3

1

Taxon A

5

Taxon D

genetic change

**Ultrametric tree**

Taxon B

Taxon C

Taxon A

Taxon D

time

**All show the same evolutionary relationships, or branching orders, between the taxa.**

# Evidence for the Molecular Clock: Cytochrome c



Slide credit: RAG

**Completely unresolved or "star" phylogeny**

A
B
C
D
E

Polytomy or **multifurcation (trifurchation)**

**Partially resolved phylogeny**

A
C
E
B
D

**Fully resolved, bifurcating phylogeny**

A
E
C
B
D

**A bifurcation**

All of these rearrangements show the same evolutionary relationships between the taxa

Rooted tree 1a

Slide credit: HK

# Limits of the tree representation

- Some events, such as hybridization, recombination, or combinations of lateral gene transfers, are poorly represented by trees.



concatenated alignment of 146 genes, and ML distances under a JTT + F + Γ model.

Huson and Bryant, MBE 2006

# How to infer trees?

# How many branches are there in an unrooted bifurcating tree of *n* taxa?



| # taxa | 2 | 3 | 4 | n |
|---|---|---|---|---|
| # branches | 1 | 3 | 5 | $2n-3$ |

# How many topologies?

| Number of "taxa" | Number of binary trees | |
| --- | --- | --- |
| | Unrooted | Rooted |
| 3 | 1 | |

# General Scheme to Build Phylogenetic Trees

# Distance matrix

Recall that have learned how to estimate the evolutionary distance between pairs of sequences:

|       | Mouse | Rat  | Human | Swine | Chimp |
|-------|-------|------|-------|-------|-------|
| Mouse |       | 0.12 | 0.38  | 0.28  | 0.38  |
| Rat   | 0.12  |      | 0.32  | 0.45  | 0.52  |
| Human | 0.38  | 0.32 |       | 0.38  | 0.08  |
| Swine | 0.28  | 0.45 | 0.38  |       | 0.33  |
| Chimp | 0.38  | 0.52 | 0.08  | 0.33  |       |

# UPGMA

## distance [mutations/site]

*(Unweighted pair group method using arithmetic averages)*

Recursively group the closest two remaining leaves

|  | 🟢 | 🔴 | 🔵 | 🟩 |
|---|---|---|---|---|
| 🟢 |  | 0.1 | 0.3 | 0.4 |
| 🔴 |  |  | 0.3 | 0.3 |
| 🔵 |  |  |  | 0.2 |
| 🟩 |  |  |  |  |

# Distance Trees as Scoring Scheme

distance [mutations/site]

| | 🟢 | 🔴 | 🔵 | 🟢 |
|---|---|---|---|---|
| 🟢 | | 0.1 | 0.3 | 0.4 |
| 🔴 | | | 0.3 | 0.3 |
| 🔵 | | | | 0.2 |
| 🟢 | | | | |

a+b = 0.1    b+c+d = 0.3
a+c+d = 0.3    b+c+e = 0.3
a+c+e = 0.4    d+e = 0.2

*5 unknown*
*6 equations*

→ overdetermined
→ minimise error

# Maximum Parsimony

- Occam's Razor: the simplest explanation is the likeliest.
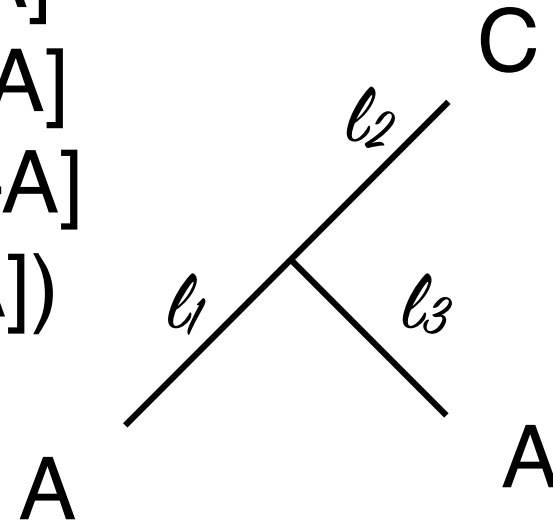- Scoring scheme: minimum number of change in (discrete) characters.

| | Number of feet | Tail? | Favorite food | Can fly? |
|---|---|---|---|---|
| Species 1 | 0 | Y | Carnivore | No |
| Species 2 | 4 | Y | Herbivore | No |
| Species 3 | 2 | Y | Herbivore | Yes |
| Species 4 | 6 | Y | Herbivore | Yes |
| Ancestor 1 | 2 | Y | Carnivore | No |
| Ancestor 2 | 2 | Y | Herbivore | No |
| Ancestor 3 | 2 | Y | Herbivore | Yes |

# Maximum Likelihood

The best tree is the one that maximizes the probability that the observed sequences would result **IF** the tree were correct.

likelihood = f(A)* (
  $M^{l1}$[A->A]*$M^{l2}$[A->C]*$M^{l3}$[A->A]
+ $M^{l1}$[A->C]*$M^{l2}$[C->C]*$M^{l3}$[C->A]
+ $M^{l1}$[A->G]*$M^{l2}$[G->C]*$M^{l3}$[G->A]
+ $M^{l1}$[A->T]*$M^{l2}$[T->C]*$M^{l3}$[T->A])



Note: here we arbitrarily start from the bottom left taxon (which has an A) as the "root" of the tree, but the likelihood is the same no matter what starting point we choose. This is a property of our Markov model, which is called "reversible".

# Discussion on methods

| | 👍 | 👎 |
|---|---|---|
| Distance | - Fast & scalable<br>- Statistically consistent (converges to true tree if the model is correct) | - Does not use all information available optimally |
| Parsimony | - Fast<br>- Intuitive | - Statistically inconsistent<br>- Lowly regarded in the phylogenetic community |
| Likelihood | - Statistically consistent<br>- Statistically efficient<br>- Highly regarded in the phylogenetic community | - Slow |

# Confidence

# Measuring confidence with the Bootstrap

Our data is limited, representing (infinitesimally) small fraction of an ideal, infinite set
What is the uncertainty due to our limited data?

Scenario: take 100 same-size sets of data from infinite set

Calculate tree for each set

See what is consistent across trees

Bootstrap value: what fraction of all trees have a given node

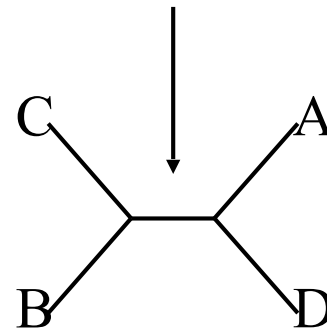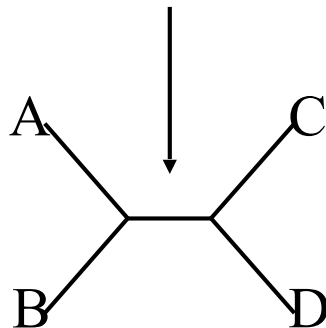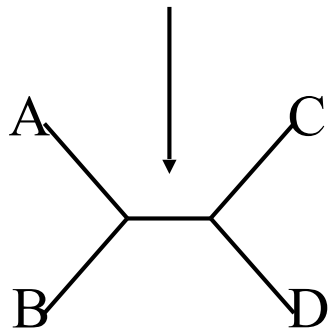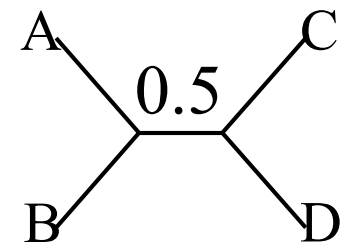# Measuring confidence with the Bootstrap

But we don't have access to infinite set!

Best we can do: use the set we have to represent the sets we don't have: **Bootstrapping**

| Data | Replicate 1 | Replicate 2 |
|------|-------------|-------------|
| A: ALTFCG | A: LCGCAL | A: ATFALF |
| B: NLTFCG | B: LCGCNL | B: NTFNLF |
| C: ALSFRG | C: LRGRAL | C: ASFALF |
| D: NLSFRG | D: LRGRNL | D: NSFNLF |

# Measuring confidence with the Bootstrap

But we don't have access to infinite set!

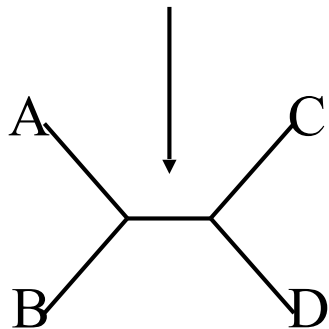Best we can do: use the set we have to represent the sets we don't have: **Bootstrapping**
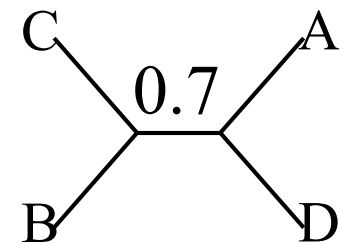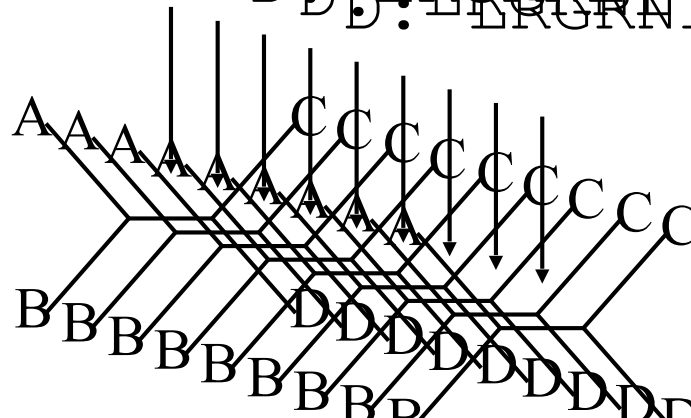


Data

```
A:  ALTFCG
B:  NLTFCG
C:  ALSFRG
D:  NLSFRG
```

# Questions?

🌐 http://lab.dessimoz.org

🐦 @cdessimoz