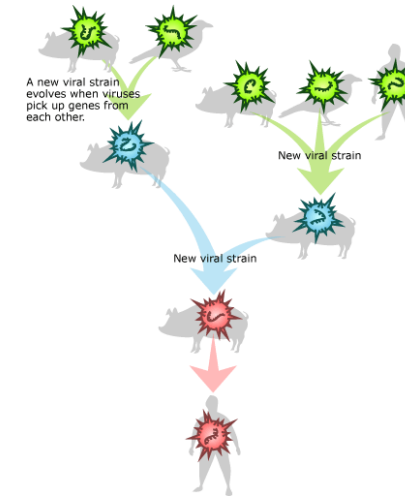# Molecular phylogeny -
# Using molecular sequences to infer evolutionary relationships

Tore Samuelsson Feb 2015

---

Molecular phylogeny is being used in the identification and characterization of new pathogens, like viruses and bacteria



A new viral strain evolves when viruses pick up genes from each other.

New viral strain

New viral strain

Bird flu virus evolution

---

## Molecular evidence of HIV-1 transmission in a criminal case

Michael L. Metzker*[†], David P. Mindell[‡], Xiao-Mei Liu*[§], Roger G. Ptak[¶||], Richard A. Gibbs*, and David M. Hillis**

*Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX 77030; [‡]Department of Ecology and Evolutionary Biology and Museum of Zoology, University of Michigan, Ann Arbor, MI 48109-1079; [¶]School of Dentistry, Biologic and Materials Sciences, University of Michigan, Ann Arbor, MI 48109; and **Section of Integrative Biology and Center for Computational Biology and Bioinformatics, University of Texas, Austin, TX 78712

A gastroenterologist was convicted of attempted second-degree murder by injecting his former girlfriend with blood or blood-products obtained from an HIV type 1 (HIV-1)-infected patient under his care. Phylogenetic analyses of HIV-1 sequences were admitted and used as evidence in this case, representing the first use of phylogenetic analyses in a criminal court case in the United States. Phylogenetic analyses of HIV-1 reverse transcriptase and *env* DNA sequences isolated from the victim, the patient, and a local population sample of HIV-1-positive individuals showed the victim's HIV-1 sequences to be most closely related to and nested within a lineage comprised of the patient's HIV-1 sequences. This finding of paraphyly for the patient's sequences was consistent with the direction of transmission from the patient to the victim.
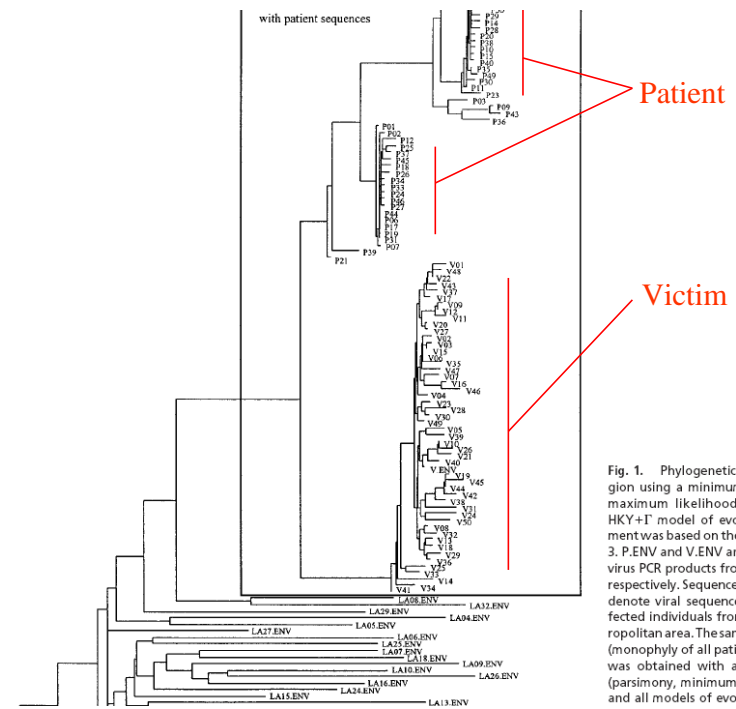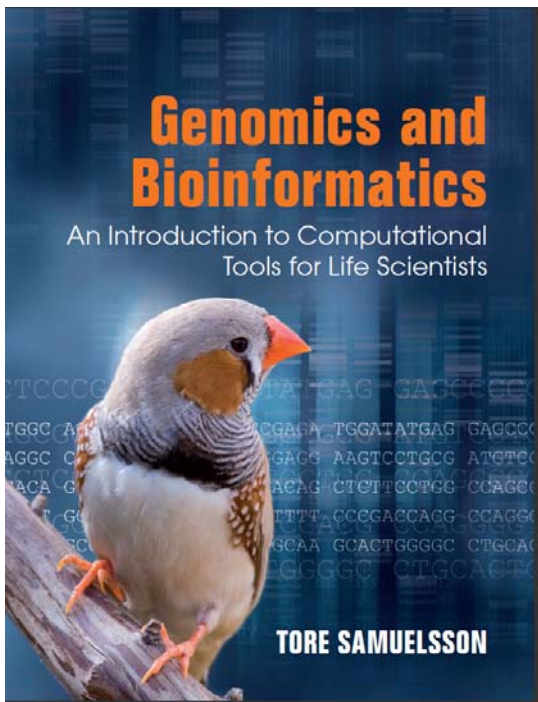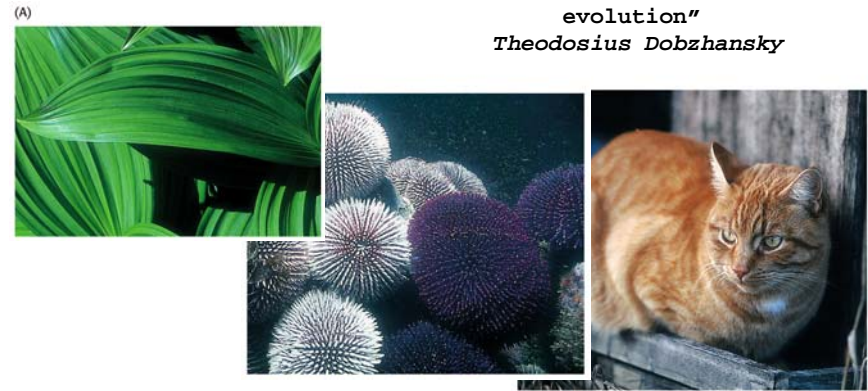
---



with patient sequences

Patient

Victim

Fig. 1. Phylogenetic analysis of the gp120 region using a minimum evolution criterion and maximum likelihood distances assuming an HKY+Γ model of evolution. Nucleotide alignment was based on the protein alignment in Fig. 3. P.ENV and V.ENV are DNA sequences for provirus PCR products from the patient and victim, respectively. Sequence names beginning with LA denote viral sequences from control HIV-1 infected individuals from the Lafayette, LA, metropolitan area. The same pattern of relationships (monophyly of all patient and victim sequences) was obtained with all phylogenetic methods (parsimony, minimum evolution, and Bayesian) and all models of evolution examined. In addi-
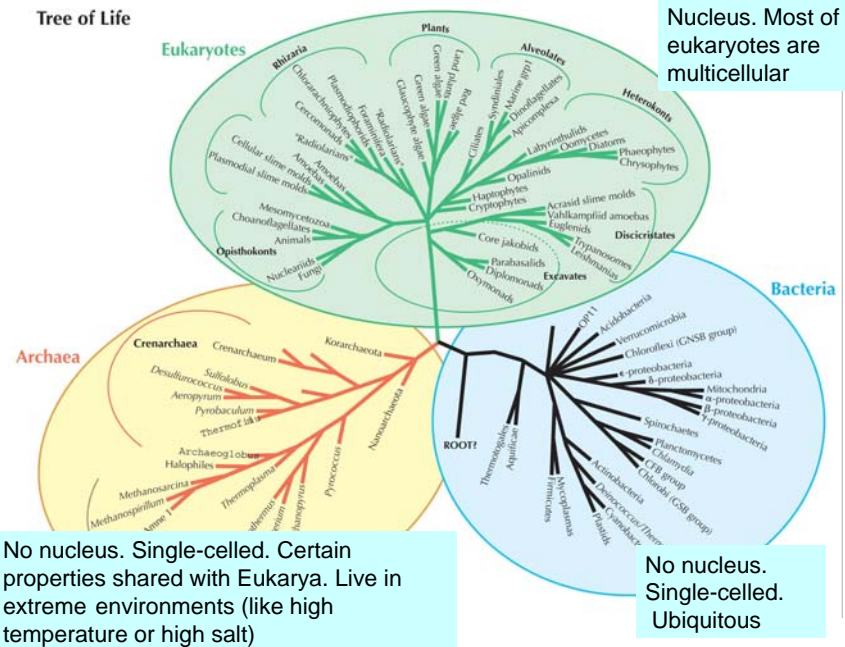
**Genomics and Bioinformatics**
An Introduction to Computational Tools for Life Scientists

**TORE SAMUELSSON**

---

**Evolution**

"Nothing in biology makes sense except in the light of evolution"
*Theodosius Dobzhansky*

(A)

Organisms are remarkably uniform at the molecular level

This uniformity reveals that organisms on Earth have arisen from a common ancestor

---

Tree of Life

Eukaryotes

Nucleus. Most of eukaryotes are multicellular

Archaea

Bacteria

No nucleus. Single-celled. Certain properties shared with Eukarya. Live in extreme environments (like high temperature or high salt)

No nucleus. Single-celled. Ubiquitous

---

**Principles of evolution**

At the molecular level evolution is a process of mutation with selection

* Reproduction
* Variation
* Competition/selective pressure

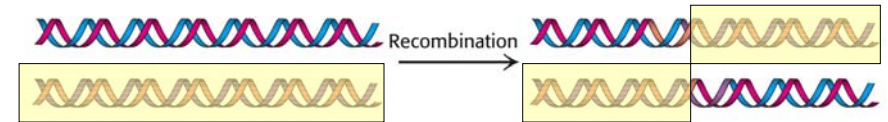## Mutations : changes in base sequence of DNA

1) single nucleotide change (point mutation)
   - **transition** (purine to purine or pyrimidine
       to pyrimidine
       C->T , T->C, A->G, G->A)
   - **transversion** (purine to pyrimidine
       or py to pu
       A->T, T->A, C->G, G->C etc

2) insertion / deletion of one or several
   nucleotides

Such mutations are the result of
   * **Replication errors**
   * **Chemicals & irradiation**

## Mutations : Homologous recombination cause large rearrangements in the genome



New gene families arise by
gene duplication and divergence

## Molecular phylogeny

**Phylogeny**
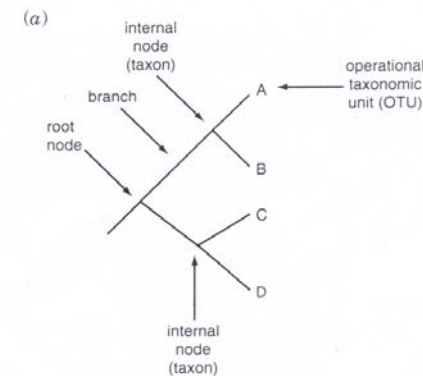   Inference of evolutionary relationships

**Molecular phylogeny**
   uses sequence information
   (as opposed to other characteristics
   frequently used in the past such as
   morphological features)

   **Goals**
   * Deduce trees to show how
   species/populations/inviduals/molecular sequences
   are related

## Nomenclature of trees



*nodes*
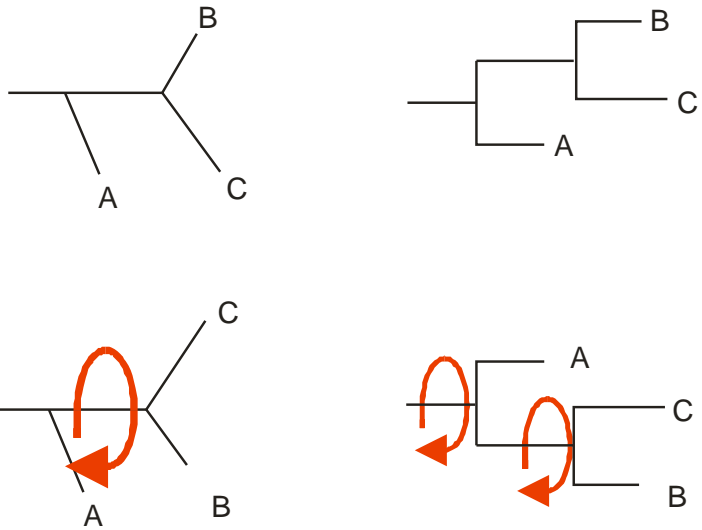   external (OTUs)
   internal
   root

*branch*
   connects 2 nodes

*OTUs* are existing (observable)
sequences / species /
 populations /individuals

an internal node is an inferred
ancestor (not observed)

(More ancient)          (More recent)

## Different ways of showing the same tree



## Nomenclature of trees

**Rooted tree**
  Root - Common ancestor of all sequences
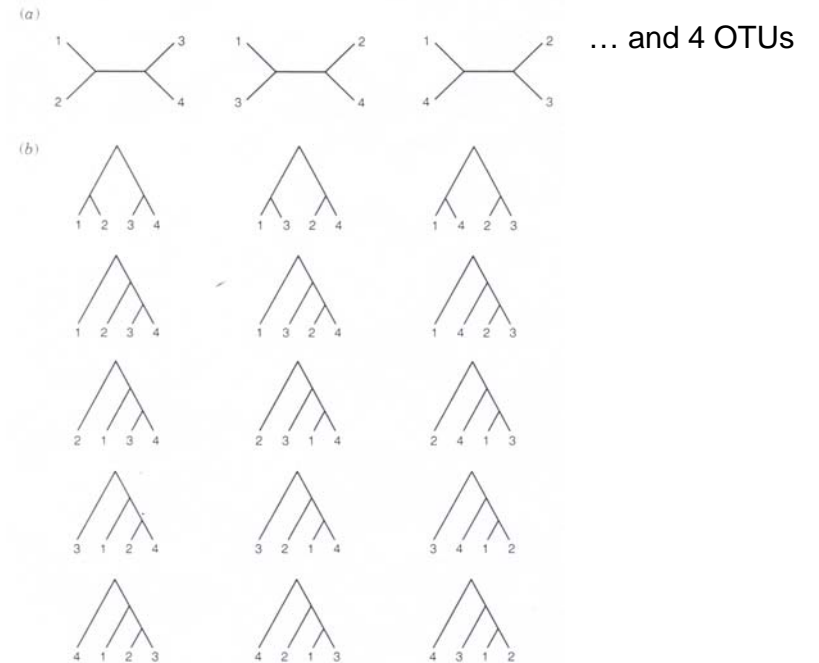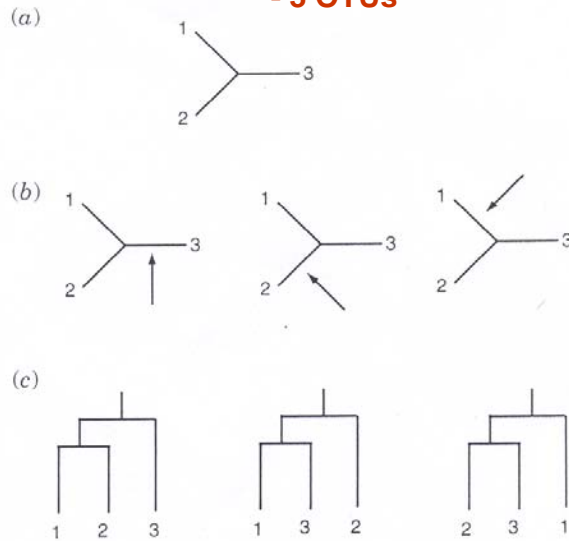  in the tree
  Unique path from the root to each of the other nodes
  Direction of each path corresponds to evolutionary time

**Unrooted tree**
  No root
  No complete definition of evolutionary path
  Direction of time is not determined

## Comparing the numbers of rooted and unrooted trees - 3 OTUs



… and 4 OTUs

## Goals of molecular phylogeny

Deduce the correct trees

* *Topology*
* *Branch lengths*

## Phylogenetic analysis

- *Selection of sequences for analysis*

DNA?
RNA?
protein?

- Multiple sequence alignment

- Construction of tree

---

Slowly changing sequences

* Protein
* ribosomal RNA, for instance 16S rRNA

Useful for comparing widely divergent species.
Ribosomal RNA database (rdp.cme.msu.edu)
> 50,000 aligned sequences

More rapidly changing sequences

* DNA
* Mitochondrial DNA

Useful for comparing more closely related
species or populations within a species.

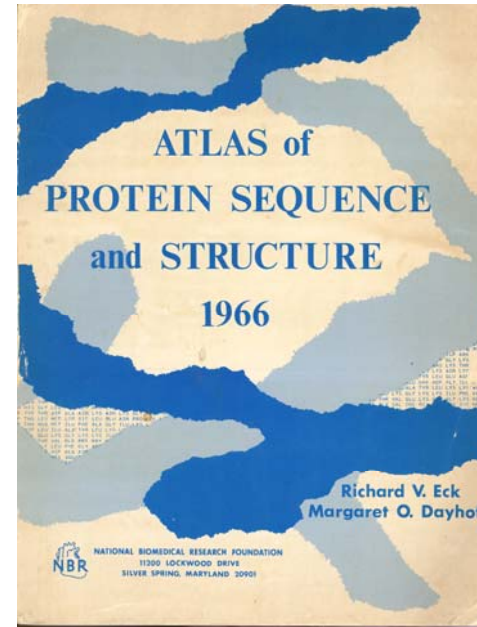---

DNA sequences evolve more rapidly than
protein sequences. This is to a large extent a result of
the genetic code degeneracy.

```
Seq 1 GGC AAG CGA AGT
Seq 2 GGA AGA CGT TCA
Seq 1 G   R   R   S
Seq 2 G   K   R   S
```

Approximate rates of substitution
(number of substitutions per site & billion years)

| | |
|---|---|
| rRNA | ~ 0.1 |
| protein | 0.01 - 10 |
| Hypervariable regions in mitochondria | 10 |
| HIV (RNA virus) | >1000 |

**Early days of molecular phylogeny**



ATLAS of PROTEIN SEQUENCE and STRUCTURE 1966

Richard V. Eck
Margaret O. Dayhoff

NATIONAL BIOMEDICAL RESEARCH FOUNDATION
11200 LOCKWOOD DRIVE
SILVER SPRING, MARYLAND 20901

**Margaret Dayhoff**



TABLE 2
CYTOCHROME C

NUMBER OF AMINO ACID DIFFERENCES BETWEEN SEQUENCES.

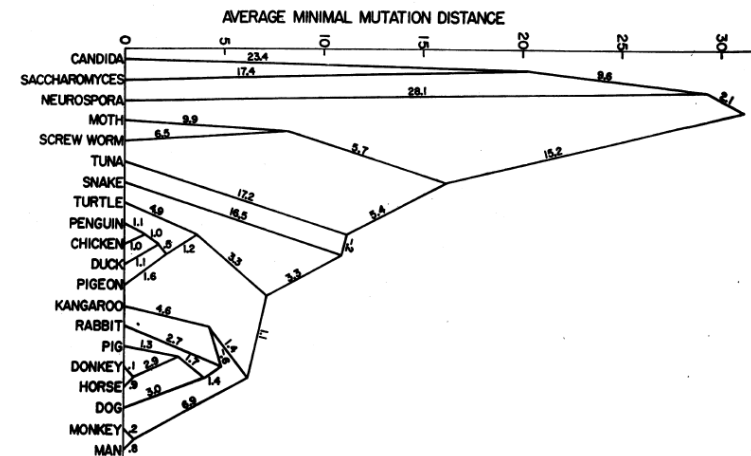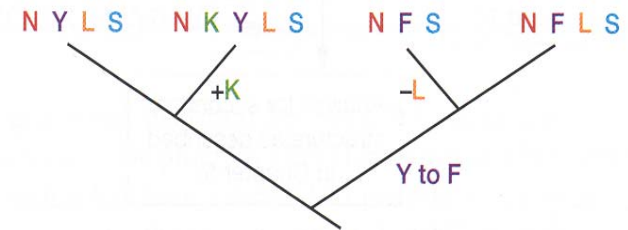| | Human | Monkey | Pig, Bovine, Sheep | Horse | Dog | Rabbit | Kangaroo | Chicken, Turkey | Duck | Rattlesnake | Turtle | Tuna Fish | Moth | Neurospora | Candida | Yeast |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Human | 0 | 1 | 10 | 12 | 11 | 9 | 10 | 13 | 11 | 14 | 15 | 21 | 31 | 48 | 51 | 45 |
| Monkey | 1 | 0 | 9 | 11 | 10 | 8 | 11 | 12 | 10 | 15 | 14 | 21 | 30 | 47 | 51 | 45 |
| Pig, Bovine, Sheep | 10 | 9 | 0 | 3 | 3 | 4 | 6 | 9 | 8 | 20 | 9 | 17 | 27 | 46 | 50 | 45 |
| Horse | 12 | 11 | 3 | 0 | 6 | 6 | 7 | 11 | 10 | 22 | 11 | 19 | 29 | 46 | 51 | 46 |
| Dog | 11 | 10 | 3 | 6 | 0 | 5 | 7 | 10 | 8 | 21 | 9 | 18 | 25 | 46 | 49 | 45 |
| Rabbit | 9 | 8 | 4 | 6 | 5 | 0 | 6 | 8 | 6 | 18 | 9 | 17 | 26 | 46 | 50 | 45 |
| Kangaroo | 10 | 11 | 6 | 7 | 7 | 6 | 0 | 12 | 10 | 21 | 11 | 18 | 28 | 49 | 51 | 46 |
| Chicken, Turkey | 13 | 12 | 9 | 11 | 10 | 8 | 12 | 0 | 3 | 19 | 8 | 17 | 28 | 47 | 51 | 46 |
| Duck | 11 | 10 | 8 | 10 | 8 | 6 | 10 | 3 | 0 | 17 | 7 | 17 | 27 | 46 | 51 | 46 |
| Rattlesnake | 14 | 15 | 20 | 22 | 21 | 18 | 21 | 19 | 17 | 0 | 22 | 26 | 31 | 47 | 51 | 47 |
| Turtle | 15 | 14 | 9 | 11 | 9 | 9 | 11 | 8 | 7 | 22 | 0 | 18 | 28 | 49 | 53 | 49 |
| Tuna Fish | 21 | 21 | 17 | 19 | 18 | 17 | 18 | 17 | 17 | 26 | 18 | 0 | 32 | 48 | 48 | 47 |
| Moth | 31 | 30 | 27 | 29 | 25 | 26 | 28 | 28 | 27 | 31 | 28 | 32 | 0 | 47 | 47 | 47 |
| Neurospora | 48 | 47 | 46 | 46 | 46 | 46 | 49 | 47 | 46 | 47 | 49 | 48 | 47 | 0 | 42 | 41 |
| Candida | 51 | 51 | 50 | 51 | 49 | 50 | 51 | 51 | 51 | 51 | 53 | 48 | 47 | 42 | 0 | 27 |
| Yeast | 45 | 45 | 45 | 46 | 45 | 45 | 46 | 46 | 46 | 47 | 49 | 47 | 47 | 41 | 27 | 0 |

Figure 4. Sequences of cytochrome c from 19 species. The amino acids common to all sequences and the allele groups at each position are shown. The sequences of the inferred common ancestors at each divergence point in the diagram are displayed below. Sites for which no single amino acid was most likely are left blank. The topology of the phylogenetic tree has been inferred from the sequences as explained in the text. The number of amino acid changes inferred between observed sequences and inferred ancestors are shown on the tree. The point of earliest time cannot be inferred directly from the sequences. We have placed it by assuming that on the average, species change at the same rate.



AVERAGE MINIMAL MUTATION DISTANCE

## Phylogenetic analysis

- Selection of sequences for analysis

- ***Multiple sequence alignment***
    *Alignment may be produced using methods
    such as CLUSTALW*

- Construction of tree

## Close relationship between multiple alignment and phylogenetic analysis



```
seqA   N  •  F     L  S
seqB   N  •  F  –  S
seqC   N  K  Y  L  S
seqD   N  •  Y  L  S
```

N Y L S        N K Y L S        N F S        N F L S

+K        –L

Y to F

## Inspecting the multiple alignment

Alignment should contain only homologous sequences.
Overall identity should ideally be significant ensuring that
the alignment is correct.

```
GGGCGGCGAGGCATTTATCGGGGGGGTTGCAAAAT
GGGCGGTGAGGCATTTATCGGGGGGGTTGCAAAAT
GGGCGGCGAAGCATAAATCGGGGAGTTGCAAAAT
GGGCGGCGAAGCATTTATCGGGGGGGTTGCGAAAT
GGGCGGCGAGGCATTTATCGGGGGGGCTGCAAAAT
```

## Phylogenetic analysis

- Selection of sequences for analysis

- Multiple sequence alignment

- ***Construction of tree***

    * ***Distance methods***

    * ***Character methods***
        *Maximum parsimony*
        *Maximum likelihood*

## Distance methods

Simplest distance measure:

Consider every pair of sequences in the multiple alignment and count the number of differences.

Degree of divergence = Hamming distance (D)

$D = n/N$
where N = alignment length
n = number of sites with differences

Example:
AG**G**CTT**T**TCA
AG**C**CTT**C**TCA
D = 2/10 = 0.2

## Generating a distance matrix

```
>A
GGACCACTACGAGCGCCTACGACGTA
>B
GGACCCCTACGAGCCCCTACGACGTA
>C
GGACCGCTGCGAGCTTCTACGACGTA
>D
GGACCTCTCCGGGCAGCTAGGACGTA
```
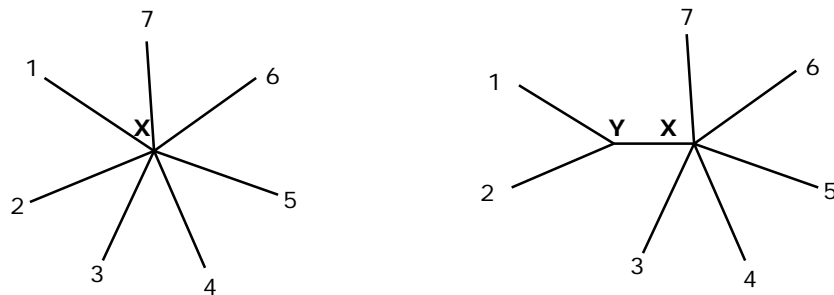
```
    A B C D
A   0 2 4 6
B   2 0 4 6
C   4 4 0 6
D   6 6 6 0
```

OR

```
      B C D
A  -  2 4 6
B  -  - 4 6
C  -  - - 6
```

## Distance methods - Neighbor joining (1987)
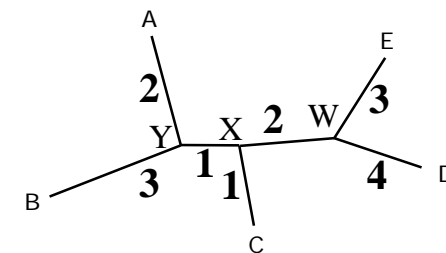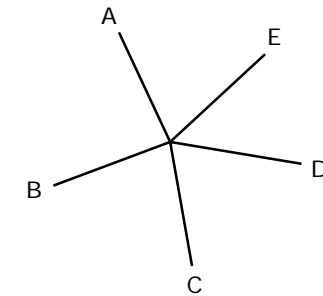


Uses *star decomposition* method

*Neighbors*: pair of nodes separated by one single node

*Minimal evolution*: minimizing total branch length.

Generates *unrooted tree*

Advantage: computationally *fast*

|   | B | C | D | E |
|---|---|---|---|---|
| A | 5 | 4 | 9 | 8 |
| B |   | 5 | 10 | 9 |
| C |   |   | 7 | 6 |
| D |   |   |   | 7 |

Character-based methods

* **_Maximum parsimony_**
* Maximum likelihood
* Bayesian statistics

---

Maximum parsimony

_parsimony_  - principle in science where the simplest answer
                is the preferred.

In phylogeny: The preferred phylogenetic tree is the one that
requires the fewest evolutionary steps.
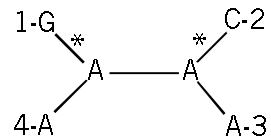
---
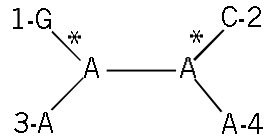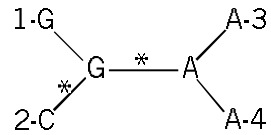
Maximum parsimony

1. Identify all _informative sites_ in the
    multiple alignment

2. For each possible tree, calculate the
    number of changes at each
    informative site.

3. Sum the number of changes for each
    possible tree.

4. Tree with the smallest number of changes
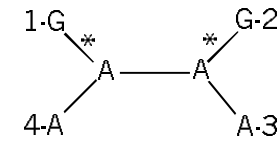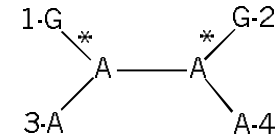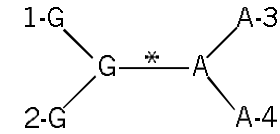    is selected as the most likely tree.

---

Maximum parsimony

Identify informative sites

```
                      Site
                1 2 3 4 5 6 7 8 9
Sequence    -------------------------
1               A A G A G T G C A
2               A G C C G T G C G
3               A G A T A T C C A
4               A G A G A T C C G
                    *       *       *
```

## Site 3 - non - informative

1-G   A-3
  G — * — A
2-C   A-4

1-G   C-2
  A — A
3-A   A-4

1-G   C-2
  A — A
4-A   A-3

## Site 5 - informative

1-G   A-3
  G — * — A
2-G   A-4

1-G   G-2
  A — A
3-A   A-4

1-G   G-2
  A — A
4-A   A-3

Summing changes:

|        | site 5 | site 7 | site 9 | Sum |
|--------|--------|--------|--------|-----|
| Tree I | 1 | 1 | 2 | 4 |
| Tree II | 2 | 2 | 1 | 5 |
| Tree III | 2 | 2 | 2 | 6 |

$\Rightarrow$Tree I most likely.

(In this case we are not considering branch lengths, only topology of tree is predicted)

Character-based methods

* Maximum parsimony
* **Maximum likelihood**
    **What is the probability that a particular tree generated the observed data under a specific model?**
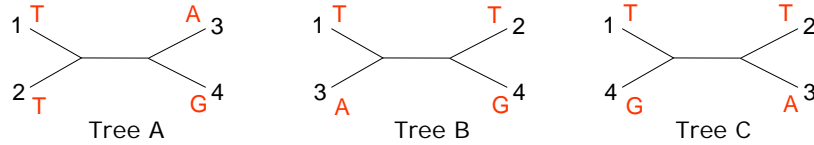* *Bayesian statistics*

Consider the following multiple alignment

```
1      A C T T
2      A C T T
3      A T A T
4      A T G C
```
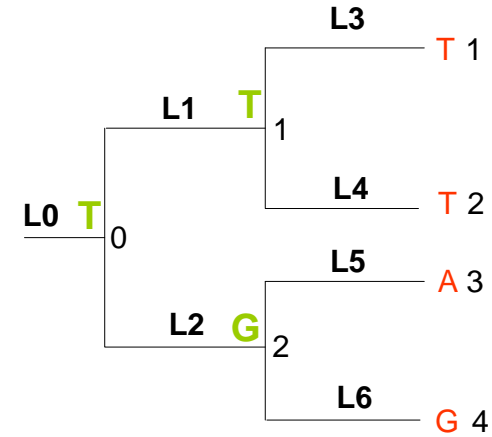
First , consider position 3 above (TTAG)
There are three possible unrooted trees for the OTUs 1-4:



Tree A          Tree B          Tree C

---

A rooted version of Tree A:



$$L(Tree1) = L0 * L1 * L2 * L3 * L4 * L5 * L6$$

---

Example of probability matrix for nucleotide substitutions

|   | A  | C  | T  | G  |
|---|----|----|----|----|
| A | ~ 1 | k  | k  | 2k |
| C | k  | ~1 | 2k | k  |
| T | k  | 2k | ~1 | k  |
| G | 2k | k  | k  | ~1 |

where  we here set k = 1E-6.

Transitions are more likely than transversions



Transitions          Transversions

---

A rooted version of Tree A:



$$L(Tree1) = L0 * L1 * L2 * L3 * L4 * L5 * L6 =$$
$$0.25 * 1 * 1E\text{-}6 * 1 * 1 * 2E\text{-}6 * 1 = 5E\text{-}13$$

A rooted version of Tree A:



$$L(Tree2) = L0 * L1 * L2 * L3 * L4 * L5 * L6 =$$
$$0.25 * 1E\text{-}6 * 2E\text{-}6 * 1E\text{-}6 * 1E\text{-}6 * 1E\text{-}6 * 1E\text{-}6 = 5E\text{-}37$$

$$L(Tree) = L(Tree1) + L(Tree2) + L(Tree3) .... L(Tree64)$$

Then we examine all positions of the alignment
in the same way. Probability of tree is the product of
probabilities for the different positions.

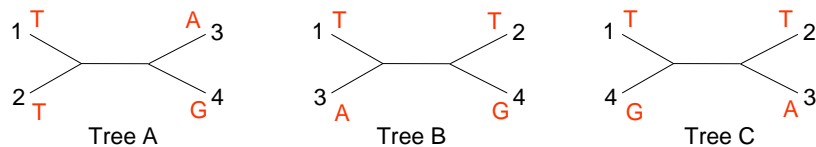$$L = L(Tree\ pos1) * L(Tree\ pos2) * L(Tree\ pos3) * L(Tree\ pos4)$$

$$lnL = ln\ L(Tree\ pos1) + ln\ L(Tree\ pos2)$$
$$+ ln\ L(Tree\ pos3) + ln\ L(Tree\ pos4)$$

Finally , the Trees B and C are handled the same way.
Tree with highest probability is preferred.

Consider the following multiple alignment

```
1     A C T T
2     A C T T
3     A T A T
4     A T G C
```

First , consider position 3 above (TTAG)
There are three possible unrooted trees for the OTUs 1-4:



Character-based methods

* Maximum parsimony
* *Maximum likelihood*
    *What is the probability of the data
    given the model?*
* **Bayesian statistics**
    **What is the probability of the
    model given the data?**

## Software for phylogenetic analysis

**PHYLIP**  (**Phyl**ogenetic **I**nference **P**ackage)
*Joe Felsenstein*
http://evolution.genetics.washington.edu/phylip.html

    DNADIST = create a distance matrix
    NEIGHBOR = neighbor joining / UPGMA
    DNAPARS = maximum parsimony
    DNAML = maximum likelihood

**PAUP** (**P**hylogenetic **A**nalysis **U**sing **P**arsimony)

**MrBayes**

---

## Applications of phylogenetic methods

- Reconstruction of evolutionary history /
  Resolving taxonomy issues
- Estimating divergence times
- Identification of gene duplication events
- Reconstructing ancient proteins
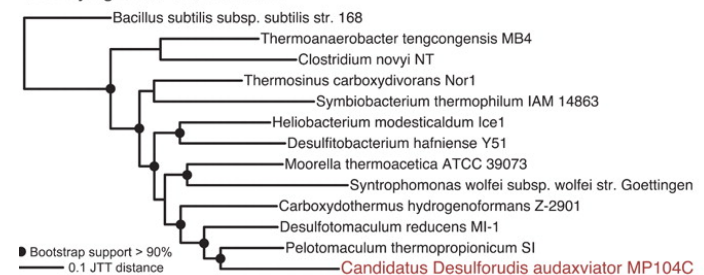- Identification of horizontal gene transfer

---



*Candidatus Desulforudis audaxviator*

---

## Life is Lonely at the Center of the Earth

*Environmental genomics reveals a single-species ecosystem deep within earth.*    Chivian et al.  Science 2008.
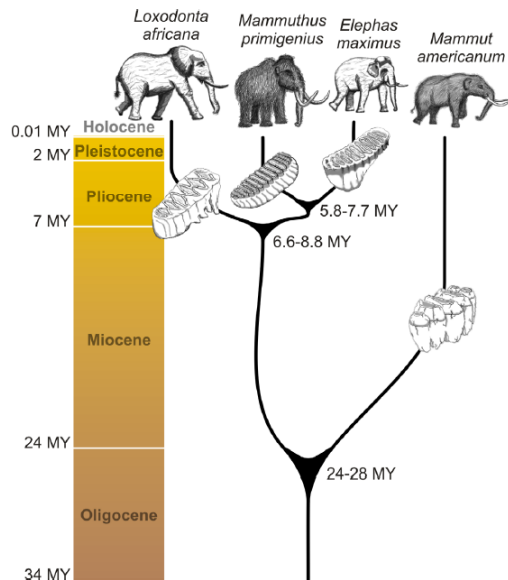


A Phylogenetic classification

Bacillus subtilis subsp. subtilis str. 168
Thermoanaerobacter tengcongensis MB4
Clostridium novyi NT
Thermosinus carboxydivorans Nor1
Symbiobacterium thermophilum IAM 14863
Heliobacterium modesticaldum Ice1
Desulfitobacterium hafniense Y51
Moorella thermoacetica ATCC 39073
Syntrophomonas wolfei subsp. wolfei str. Goettingen
Carboxydothermus hydrogenoformans Z-2901
Desulfotomaculum reducens MI-1
Pelotomaculum thermopropionicum SI
Candidatus Desulforudis audaxviator MP104C
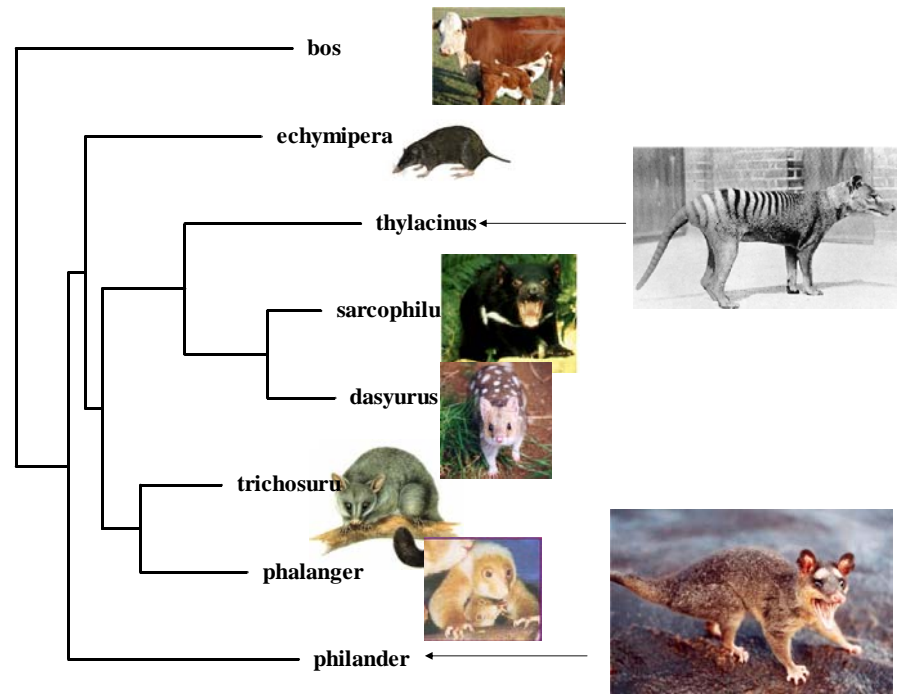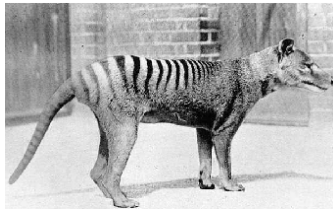
Bootstrap support > 90%
0.1 JTT distance

**Molecular phylogeny to examine extinct species - I**

Is the south american opossum

evolutionary related to the australian 'marsupial wolf' ?



bos
echymipera
thylacinus
sarcophilu
dasyurus
trichosuru
phalanger
philander



**Molecular phylogeny to examine extinct species - II**

*"Sequencing the nuclear genome of the extinct woolly mammoth". Miller et al. Nature Nov. 2008*

Loxodonta africana    Mammuthus primigenius    Elephas maximus    Mammut americanum

0.01 MY  Holocene
2 MY  Pleistocene
Pliocene
7 MY
5.8-7.7 MY
6.6-8.8 MY
Miocene
24 MY
24-28 MY
Oligocene
34 MY

**Molecular phylogeny to examine extinct species - III**

Phylogeny of Neanderthal individuals



Feldhofer    Mezmaiskaya

Svante Pääbo

A Complete Neandertal Mitochondrial Genome Sequence Determined by High-Throughput Sequencing

Richard E. Green,[1,*] Anna-Sapfo Malaspinas,[2] Johannes Krause,[1] Adrian W. Briggs,[1] Philip L.F. Johnson,[3] Caroline Uhler,[4] Matthias Meyer,[1] Jeffrey M. Good,[1] Tomislav Maricic,[1] Udo Stenzel,[1] Kay Prüfer,[1] Michael Siebauer,[1] Hernán A. Burbano,[1] Michael Ronan,[5] Jonathan M. Rothberg,[5] Michael Egholm,[5] Pavao Rudan,[7] Dejana Brajković,[8] Željko Kućan,[7] Ivan Gušić,[7] Mårten Wikström,[9] Liisa Laakkonen,[10] Janet Kelso,[1] Montgomery Slatkin,[2] and Svante Pääbo[1]

[1]Max-Planck Institute for Evolutionary Anthropology, D-04103 Leipzig, Germany

Labels in figure: 706 KYA human-Neanderthal nuclear TMRCA; 660 KYA human-Neanderthal mtDNA TMRCA; 370 KYA ancestral population splitting; 171 KYA mtDNA TMRCA; 38 KYA; Neanderthal; modern human



Mitochondrial genome variation and the origin of modern humans

Max Ingman[*], Henrik Kaessmann[†], Svante Pääbo[†] & Ulf Gyllensten[*]

* Department of Genetics and Pathology, Section of Medical Genetics, Rudbeck Laboratory, University of Uppsala, S-751 85 Uppsala, Sweden
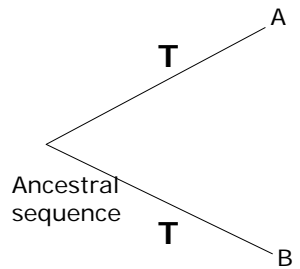† Max Planck Institute for Evolutionary Anthropology, Inselstrasse 22, D-04103 Leipzig, Germany



**"Out of Africa" hypothesis**
Modern humans evolved from archaic forms only in Africa.
Archaic humans living in Asia and Europe (like the Neanderthal) were replaced by
modern humans migrating out of Africa.

### Applications of phylogenetic methods

- Reconstruction of evolutionary history / Resolving taxonomy issues
- Estimating divergence times
- Identification of gene duplication events
- Reconstructing ancient proteins
- Identification of horizontal gene transfer

**Slide (top-left):**

A molecular clock may be used in the estimation of *time of divergence* between two species

r = K / 2T   or

T = K/2r

where

r = rate of nucleotide substitution (estimated from fossil records)

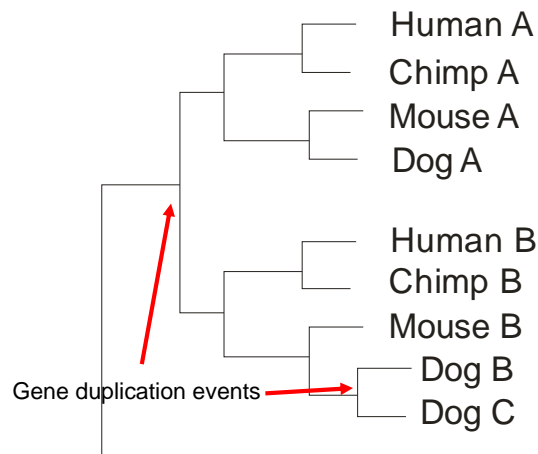K = number of substitutions K between the two homologous sequences

T = **Time of divergence between the two species**

A

T

Ancestral sequence

T

B

**Slide (top-right):**

## Applications of phylogenetic methods

- Reconstruction of evolutionary history / Resolving taxonomy issues
- Estimating divergence times
- Identification of gene duplication events
- Reconstructing ancient proteins
- Identification of horizontal gene transfer

**Slide (bottom-left):**

## Analysis of gene and protein evolution

Human A
Chimp A
Mouse A
Dog A

Human B
Chimp B
Mouse B
Dog B
Dog C

Gene duplication events

**Slide (bottom-right):**

## Applications of phylogenetic methods

- Reconstruction of evolutionary history / Resolving taxonomy issues
- Estimating divergence times
- Identification of gene duplication events
- Reconstructing ancient proteins
- Identification of horizontal gene transfer

## Slide 1

### Resurrecting ancestral proteins responsible for ethanol digestion
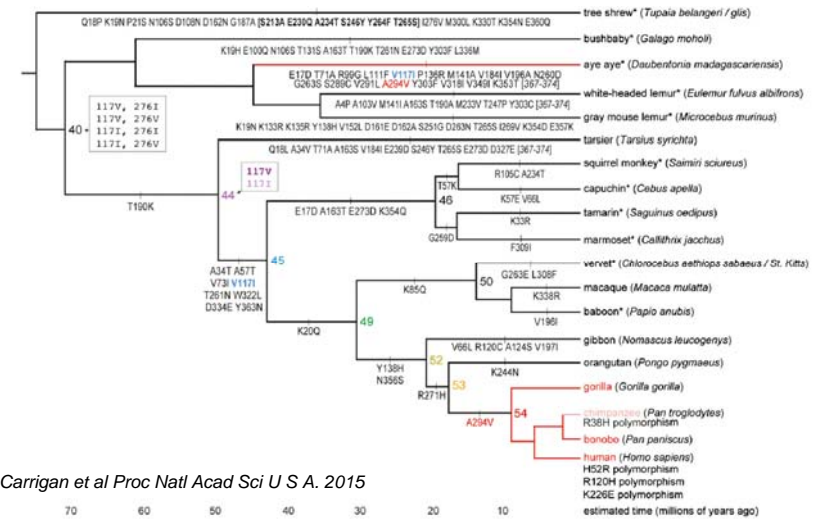
alcohol dehydrogenase

Ethanol => acetaldehyde

Modern humans can metabolize ethanol.
Many monkeys such as gibbon and orangutang cannot.
Was the ability to metabolize ethanol developed when humans started intentional fermentation of food?

## Slide 2

### Resurrecting ancestral proteins responsible for ethanol digestion
**Ancestral sequences are inferred from present sequences and proteins are then produced in the lab to examine their properties.**
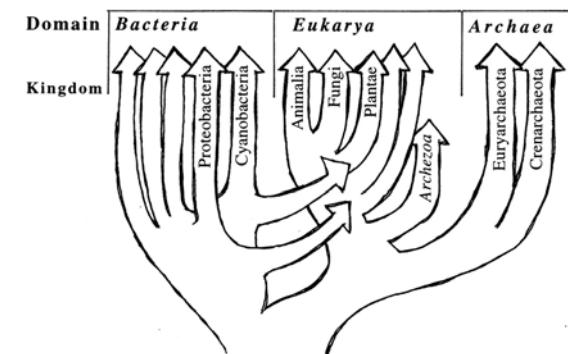


*Carrigan et al Proc Natl Acad Sci U S A. 2015*

Conclusion: The ability to digest ethanol appeared ~10 million years ago
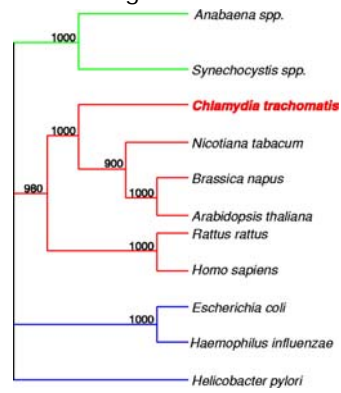
## Slide 3

### Applications of phylogenetic methods

• Reconstruction of evolutionary history / Resolving taxonomy issues
• Estimating divergence times
• Identification of gene duplication events
• Reconstructing ancient proteins
• Identification of horizontal gene transfer

## Slide 4

### Horizontal gene transfer - transfer of genes between species

## Phylogenetic analysis may be used to identify horisontal gene transfer.

Some Chlamydia (Eubacteria kingdom) proteins group
with plant homologs



Phylogeny of chlamydial
enoyl-acyl carrier protein
reductase as an example of
horizontal transfer.

From: Stephens RS, et al  Genome sequence of an obligate

intracellular pathogen of humans: Chlamydia trachomatis.

Science. 1998 Oct 23;282(5389):754-9.



Mitochondria and chloroplasts resulted from
bacteria that lived in symbiosis with a primitive
eukaryote. Eventually many genes were lost
or transferred to the
nuclear genome. Therefore, some nuclear
encoded proteins resemble bacterial proteins.