# PHYLOGENETIC ANALYSIS I.

# Definition of basic concepts:

phylogenetic tree = phylogeny (fylogenie): rooted, unrooted

branches = edges (větve): peripheral, internal, central

nodes = vertices (uzly): internal, terminal

dichotomy = bifurcation, polytomy = multifurcation

OTU = operational taxonomic unit, HTU = hypothetical taxonomic unit

tree topology

# Definition of basic concepts:

connects two
terminal nodes

c)

path

**path (dráha)**

connects
terminal node
with root

d)

lineage

**lineage (linie)**

# Definition of basic concepts:



star tree          partly resolved          fully resolved

network

Top network:

Snp39*
Smb27
Snp76
C
D
UND101
Sbr68
Smb−17
Sha161
Sty85
B
Sha154
Sha149,Snp34*
UND64,Snp128
Sty62
Sha158
Sty15*
Sha169
Sha183
E
San37
She7*
Sag129
Sha182
Sha147
Sca97, UND79
She12
A
Sha151,Sjo99
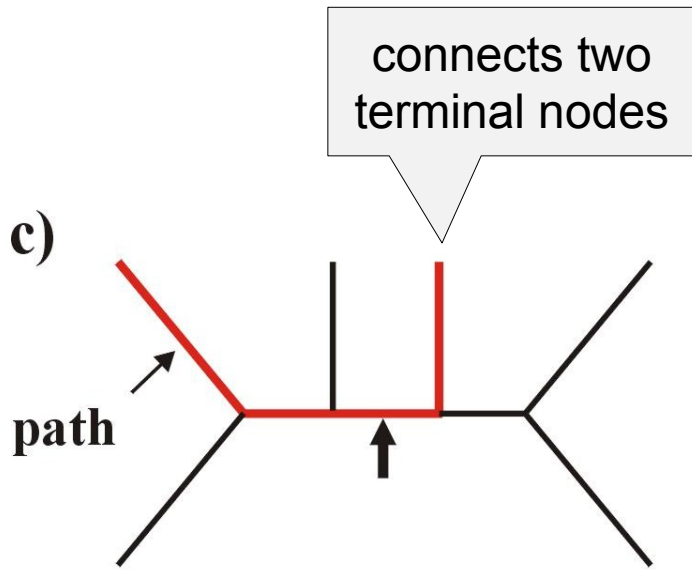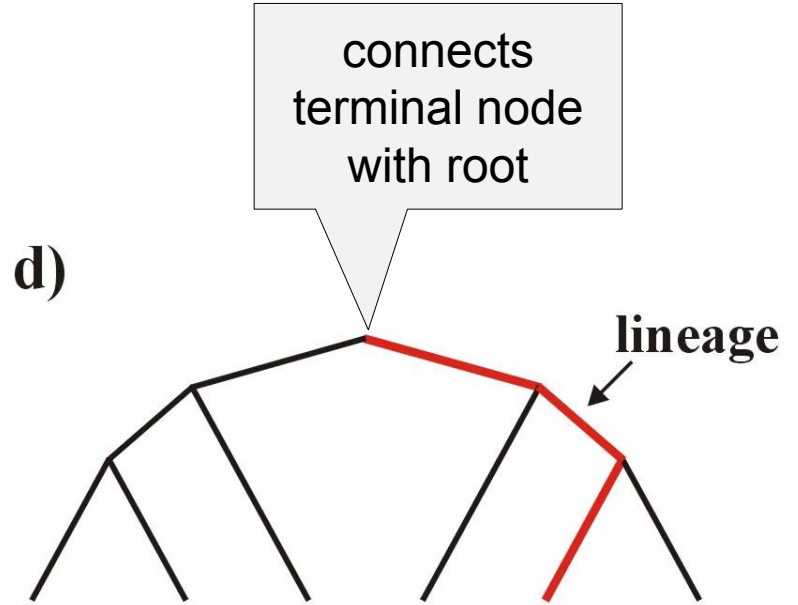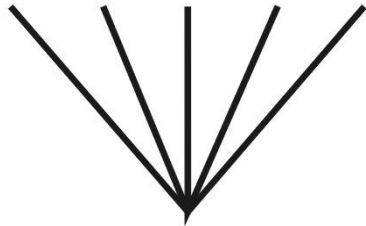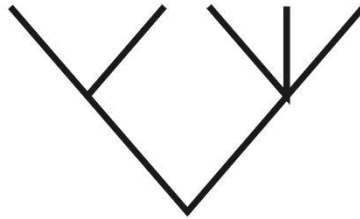Sse94
Sty90
Sre115
Sha135,Sha146
UND8
Sty19*
She49*
Sha184,Sen57*,Sha139,Sha60

0.01

Bottom network:

SHG7Y
969WV
N9VG5
SSYBZ
FQYN9
JTDBS
VCF3C  6XFVJ
VA6FD
87BA2
4AGBN
DUYF7
6PTZQ
9UYJR
68CPM
UXGSU
VZPTP
XYC78
EWCV8
3D6C5
DXATE
UDHAC
4U9UA
Y4ACV
E3C9W
4FS98
H7HZ3
EAZAV
M3MJU
MVRM9
QG3AJ
W8MF2
U4EE8
68DJZ
C4MB3
E78NV
UXE4F
HS8HH
B38TQ
PGCE4
ER8YT
XWHRK
RFPNW
C2MUB
QEDWY
GJMBJ  C27F2
9AFR6
RBQSR
SW6Z7
WVSF8
KKDC9
F6FNN
SPBSW
CRTJ5
JCEA9
KXNT5
3F536
CRNA3
F5CUZ
9W7J4
J5FX7
7468T  7JG3Y
NUKYN
F90132
F90130
F90133
ZAKSZ
U7U7Z
F90126
F90134  F90131
VJBX3
N9TPH
F39685
F90124
6PKCJ
F90128
F90129
F90127
F67866

# How many trees?

........................................................................

> Avogadro constant*)

number of electrons in visible universe (Eddington number)

*) 6,022 140 76×10$^{23}$ mol$^{-1}$

# What type of data can we use?



DATA

Distances

Discrete characters

Immunology
DNA-DNA hybridization

Binary

Multistate

11010010011

ABCDEF

unordered
ACGTTAGCT

ordered
A→B→C

# **Types of data**

## Nucleotide and protein sequences:

H_sapiens MTPMRKINPLMKLINHSFIDLPTPSNISAWWNFGS

base = character state

P_troglod ATGACCCCGACACGCAAAATTAACCCACTAATAAA

site = character

# **Types of data**

retroelements: SINE (*Alu*, B1, B2), LINE

microsatellites, SNP

# Problem with homology of sequences

# Problem with homology of sequences



Individual sites in DNA sequences may not be fully independent!

# **Sequences**

## DNA databases:

EMBL (European Molecular Biology Laboratory) – European Bioinformatics Institute, Hinxton, UK: *http://www.ebi.ac.uk/embl/*

GenBank – NCBI (National Center for Biotechnology Information), Bethesda, Maryland, USA: *http://www.ncbi.nlm.nih.gov/Genbank/*

DDBJ (DNA Data Bank of Japan) – National Institute of Genetics, Mishima, Japan: *http://www.ddbj.nig.ac.jp/*

Database managment: usually packages Sybase or ORACLE

outputs: ASCII (*American Standard Code for Information Interchange*)

# **Sequences**

## Protein databases:

SWISS-PROT – University of Geneve & Swis Institute of Bioinformatics:
*http://www.expasy.ch/sprot/* a *http://www.ebi.ac.uk/swissprot/*

PIR (Protein Information Resource) – NBRF (National Biomedical Research Foundation, Washington, D.C., USA) & Tokyo University & JIPID (Japanese International Protein Information Database, Tokyo) & MIPS (Martinsried Institute for Protein Sequences, Martinsried, Germany): *http://www-nbrf.georgetown.edu/*

PRF/SEQDB (Protein Resource Foundation) – Ósaka, Japan:
*http://www.prf.or.jp/en/os.htm*

PDB (Protein Data Bank) – University of New Jersey, San Diego & Super-computer Center, University of California & National Institute of Standards and Technology:
*http://www.rcsb.org/pdb/*

# File formats:

## FASTA:

```
>H_sapiens
ATGACCCCAATACGCAAAATTAACCCCCTAATAAAATTAATTAACCACTCATTCATCGACCTCCCCACCC
CATCCAACATCTCCGCATGATGAAACTTCGGCTCACTCCTTGGCGCCTGCCTGATCCTCCAAATCACCAC
AGGACTATTCCTAGCCATACACTACTCACCAGACGCCTCAACCGCCTTTTCATCAATCGCCCACATCACT
CGAGACGTAAATTATGGCTGAATCATCCGCTACCTTCACGCCAATGGCGCCTCAATATTCTTTATCTGCC
TCTTCCTACACATCGGGCGAGGCCTATATTACGGATCATTTCTCTACTCAGAAACCTGAAACATCGGCAT
...
>P_troglod
ATGACCCCGACACGCAAAATTAACCCACTAATAAAATTAATTAATCACTCATTTATCGACCTCCCCACCC
CATCCAACATTTCCGCATGATGGAACTTCGGCTCACTTCTCGGCGCCTGCCTAATCCTTCAAATTACCAC
AGGATTATTCCTAGCTATACACTACTCACCAGACGCCTCAACCGCCTTCTCGTCGATCGCCCACATCACC
CGAGACGTAAACTATGGTTGGATCATCCGCTACCTCCACGCTAACGGCGCCTCAATATTTTTTATCTGCC
TCTTCCTACACATCGGCCGAGGTCTATATTACGGCTCATTTCTCTACCTAGAAACCTGAAACATTGGCAT
...
>P_paniscus
ATGACCCCAACACGCAAAATCAACCCACTAATAAAATTAATTAATCACTCATTTATCGACCTCCCCACCC
CATCCAATATTTCCACATGATGAAACTTCGGCTCACTTCTCGGCGCCTGCCTAATCCTTCAAATCACCAC
AGGACTATTCCTAGCTATACACTACTCACCAGACGCCTCAACCGCCTTCTCATCGATCGCCCACATTACC
CGAGACGTAAACTATGGTTGAATCATCCGCTACCTTCACGCTAACGGCGCCTCAATACTTTTCATCTGCC
TCTTCCTACACGTCGGTCGAGGCCTATATTACGGCTCATTTCTCTACCTAGAAACCTGAAACATTGGCAT
...
```

# File formats:

## GenBank:

```
ORIGIN
        1 tgaaatgaag atattctctt ctcaagacat caagaagaag gaactactcc ccaccaccag
       61 cacccaaagc tggcattcta attaaactac ttcttgtgta cataaattta catagtacaa
      121 tagtacattt atgtatatcg tacattaaac tattttcccc aagcatataa gcaagtacat
      181 ttaatcaatg atataggcca taaaacaatt atcaacataa actgatacaa accatgaata
      241 ttatactaat acatcaaatt aatgctttaa agacatatct gtgttatctg acatacacca
      301 tacagtcata aactcttctc ttccatatga ctatcccctt ccccatttgg tctattaatc
      361 taccatcctc cgtgaaacca acaacccgcc caccaatgcc cctcttctcg ctccgggccc
      421 attaaacttg ggggtagcta aactgaaact ttatcagaca tctggttctt acttcagggc
      481 catcaaatgc gttatcgccc atacgttccc cttaaataag acatctcgat ggtatcgggt
      541 ctaatcagcc catgaccaac ataactgtgg tgtcatgcat ttggtatttt tttattttgg
      601 cctactttca tcaacatagc cgtcaaggca tgaaaggaca gcacacagtc tagacgcacc
      661 tacggtgaag aatcattagt ccgcaaaacc caatcaccta aggctaatta ttcatgcttg
      721 ttagacataa atgctactca ataccaaatt ttaactctcc aaacccccca accccctcct
      781 cttaatgcca aaccccaaaa acactaagaa cttgaaagac atatattatt aactatcaaa
      841 ccctatgtcc tgatcgattc tagtagttcc caaaatatga ctcatatttt agtacttgta
      901 aaaattttac aaaatcatgc tccgtgaacc aaaactctaa tcacactcta ttacgcaata
      961 aatattaaca agttaatgta gcttaataac aaagcaaagc actgaaaatg cttagatgga
     1021 taattttatc cca
//
```

# File formats:

## PHYLIP ("interleaved" format):

```
6 1120
H_sapiens     ATGACCCCAA TACGCAAAAT TAACCCCCTA ATAAAATTAA TTAACCACTC
P_troglod     ATGACCCCGA CACGCAAAAT TAACCCACTA ATAAAATTAA TTAATCACTC
P_paniscus    ATGACCCCAA CACGCAAAAT CAACCCACTA ATAAAATTAA TTAATCACTC
G_gorilla     ATGACCCCTA TACGCAAAAC TAACCCACTA GCAAAACTAA TTAACCACTC
P_pygmaeus    ATGACCCCAA TACGCAAAAC CAACCCACTA ATAAAATTAA TTAACCACTC
H_lar         ATGACCCCCC TGCGCAAAAC TAACCCACTA ATAAAACTAA TCAACCACTC

              ATTCATCGAC CTCCCCACCC CATCCAACAT CTCCGCATGA TGAAACTTCG
              ATTTATCGAC CTCCCCACCC CATCCAACAT TTCCGCATGA TGGAACTTCG
              ATTTATCGAC CTCCCCACCC CATCCAATAT TTCCACATGA TGAAACTTCG
              ATTCATTGAC CTCCCTACCC CGTCCAACAT CTCCACATGA TGAAACTTCG
              ACTCATCGAC CTCCCCACCC CATCAAACAT CTCTGCATGA TGGAACTTCG
              ACTTATCGAC CTTCCAGCCC CATCCAACAT TTCTATATGA TGAAACTTTG
```

# File formats:

## NEXUS (PAUP*, "interleaved"):

```
#NEXUS
begin data;
dimensions ntax=6 nchar=1120;
format datatype=DNA interleave datatype=DNA missing=? gap=-;
matrix
P_troglod    ATGACCCCGACACGCAAAATTAACCCACTAATAAAATTAATTAATCACTC
P_paniscus   ATGACCCCAACACGCAAAATCAACCCACTAATAAAATTAATTAATCACTC
H_sapiens    ATGACCCCAATACGCAAAATTAACCCCCTAATAAAATTAATTAACCACTC
G_gorilla    ATGACCCCTATACGCAAAACTAACCCACTAGCAAAACTAATTAACCACTC
P_pygmaeus   ATGACCCCAATACGCAAAACCAACCCACTAATAAAATTAATTAACCACTC
H_lar        ATGACCCCCCTGCGCAAAACTAACCCACTAATAAAACTAATCAACCACTC

P_troglod    ATTTATCGACCTCCCCACCCCATCCAACATTTCCGCATGATGGAACTTCG
P_paniscus   ATTTATCGACCTCCCCACCCCATCCAATATTTCCACATGATGAAACTTCG
H_sapiens    ATTCATCGACCTCCCCACCCCATCCAACATCTCCGCATGATGAAACTTCG
G_gorilla    ATTCATTGACCTCCCTACCCCGTCCAACATCTCCACATGATGAAACTTCG
P_pygmaeus   ACTCATCGACCTCCCCACCCCATCAAACATCTCTGCATGATGGAACTTCG
H_lar        ACTTATCGACCTTCCAGCCCCATCCAACATTTCTATATGATGAAACTTTG

end;
```

# File formats:

## Clustal X:

```
P_troglod   ATGACCCCGACACGCAAAATTAACCCACTAATAAAATTAATTAATCACTCATTTATCGAC
P_paniscus  ATGACCCCAACACGCAAAATCAACCCACTAATAAAATTAATTAATCACTCATTTATCGAC
H_sapiens   ATGACCCCAATACGCAAAATTAACCCCCTAATAAAATTAATTAACCACTCATTCATCGAC
G_gorilla   ATGACCCCTATACGCAAAACTAACCCACTAGCAAAACTAATTAACCACTCATTCATTGAC
P_pygmaeus  ATGACCCCAATACGCAAAACCAACCCACTAATAAAATTAATTAACCACTCACTCATCGAC
H_lar       ATGACCCCCCTGCGCAAAACTAACCCACTAATAAACTAATCAACCACTCACTTATCGAC
            ********   *******  ***** ***   **** **** ** ****** * ** ***


P_troglod   CTCCCCACCCCATCCAACATTTCCGCATGATGGAACTTCGGCTCACTTCTCGGCGCCTGC
P_paniscus  CTCCCCACCCCATCCAATATTTCCACATGATGAAACTTCGGCTCACTTCTCGGCGCCTGC
H_sapiens   CTCCCCACCCCATCCAACATCTCCGCATGATGAAACTTCGGCTCACTCCTTGGCGCCTGC
G_gorilla   CTCCCTACCCCGTCCAACATCTCCACATGATGAAACTTCGGCTCACTCCTTGGTGCCTGC
P_pygmaeus  CTCCCCACCCCATCAAACATCTCTGCATGATGGAACTTCGGCTCACTTCTAGGCGCCTGC
H_lar       CTTCCAGCCCCATCCAACATTTCTATATGATGAAACTTTGGTTCACTCCTAGGCGCCTGC
            ** **   **** ** ** ** ** **      ****** ***** ** ***** ** ** ******
```

# File formats:

## FASTQ:

Line 1 begins with a '@' character and is followed by a sequence identifier and an *optional* description (like a <u>FASTA</u> title line).

Line 2 is the raw sequence letters.

Line 3 begins with a '+' character and is *optionally* followed by the same sequence identifier (and any description) again.

Line 4 encodes the quality values for the sequence in Line 2, and must contain the same number of symbols as letters in the sequence.

# Progressive alignment - ClustalX

3 phases:

1. Alignment of sequence pairs $\rightarrow$ pairwise distances
2. Construction of guide tree (eg. Neighbor-Joining)
3. Alignment of all sequences according to guide tree

# Problem with progressive alignment

## 6 species:

| | | | |
|---|---|---|---|
| gorilla | AGGTT | penguin | A-GTT |
| horse | AG-TT | chicken | A-GTT |
| panda | AG-TT | ostrich | AGGTT |



| | | | | |
|---|---|---|---|---|
| AGGTT | | gorilla | AGGTT | AGGTT |
| AG-TT | | horse | AG-TT | A-GTT |
| AG-TT | | panda | AG-TT | A-GTT |
| AG-TT | | penguin | A-GTT | A-GTT |
| AG-TT | | chicken | A-GTT | A-GTT |
| AGGTT | | ostrich | AGGTT | AGGTT |

Many other alignment programs: e.g. MAFFT, MUSCLE, Geneious...

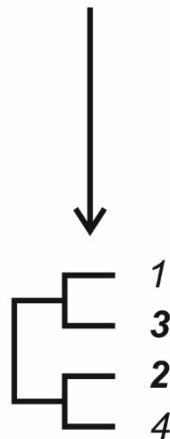# There are also methods without alignment:



homologní sekvence

referenční strom

seřazení všech sekvencí

metoda bez seřazení sekvencí

fylogenetický strom

# Methods

**Data types**

**distances**        **characters**

|  | distances | characters |
|---|---|---|
| **algorithms** | UPGMA<br><br>neighbor-joining | |
| **optimality criteria** | Fitch-Margoliash<br><br>minimum evolution | maximum parsimony<br><br>maximum likelihood<br><br>Bayesian a. |

**Methods of tree construction**

# How to assess the methods?

**Efficiency**:    how fast is the method?

**Power**:    how many characters we need?

**Consistency**:    does increasing characters result in true tree?

**Robustness**:    how does it work when assumptions are violated?

**Falsifiability**:    does it allow testing assumptions?

# MAXIMUM PARSIMONY, MP
## (maximální úspornost)

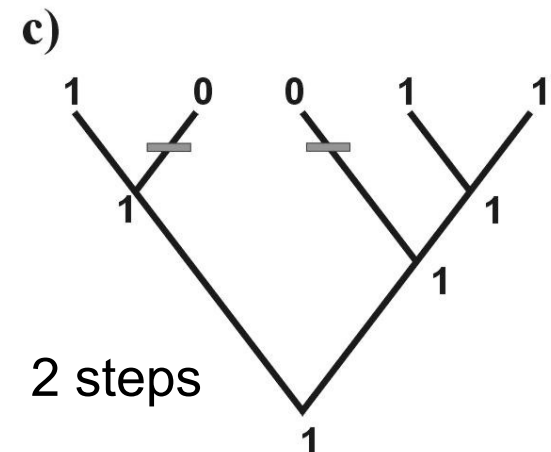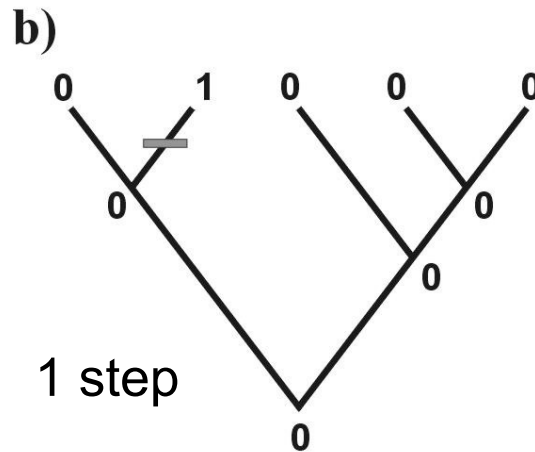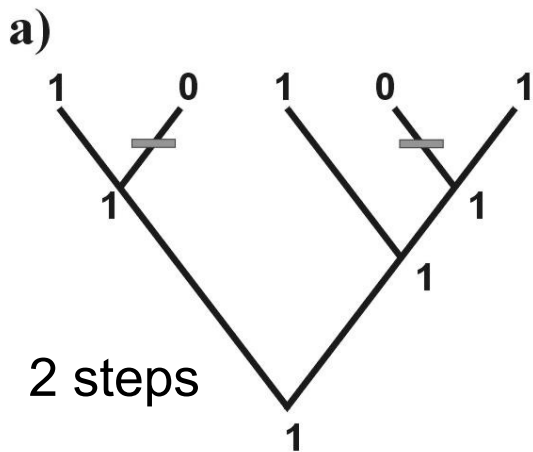William of Ockham (c. 1287 – 1347)

Occam's razor

|   | I | II | III |
|---|---|----|-----|
| A | 1 | 0  | 1   |
| B | 0 | 0  | 1   |
| C | 1 | 0  | 0   |
| D | 0 | 1  | 0   |
| E | 1 | 0  | 1   |

minimal number of steps = 3
real number of steps = 5
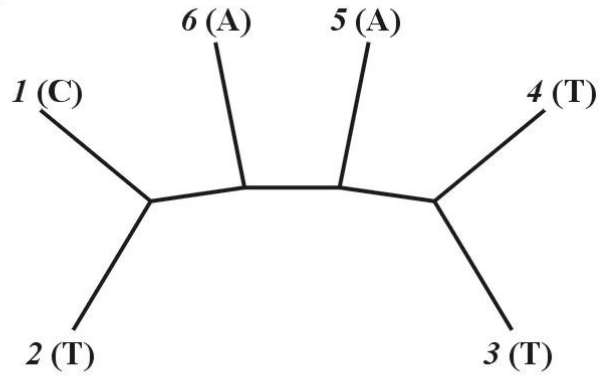$\Rightarrow$ 2 extra steps $\rightarrow$ <u>homoplasy</u>



a) 2 steps

b) 1 step

c) 2 steps

# Estimation of number of steps: Fitch algorithm



a)

6 (A)    5 (A)

1 (C)         4 (T)

2 (T)         3 (T)

1. arbitrary root

# Estimation of number of steps: Fitch algorithm



**a)**

6 (A)   5 (A)

1 (C)   4 (T)

2 (T)   3 (T)

**b)**

| 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|
| C | T | T | T | A |

[C,T] $w$   [T] $x$

[A,T] $y$

[T] $z$

6
A

1. arbitrary root

2. Downward:
   $w$ = C or T
   $x$ = T
   $y$ = A or T
   $z$ = T

# Estimation of number of steps: Fitch algorithm



**a)**

6 (A)   5 (A)

1 (C)                    4 (T)

2 (T)          3 (T)

**b)**

1  2  3  4  5
C  T  T  T  A

[C,T] w   [T] x

[A,T] y

[T] z

6
A

**c)**

1  2  3  4  5
C  T  T  T  A

T w   T x

A y

A z

6
A

**DELTRAN**
**(DELayed TRANsformation)**

**d)**

1  2  3  4  5
C  T  T  T  A

T w   T x

T y

T z

6
A

**ACCTRAN**
**(ACCelerated TRANsformation)**

1. arbitrary root

2. Downward:
   *w* = C or T
   *x* = T
   *y* = A or T
   *z* = T

3. Upward:
   *z* = T, nebo A

   total length = 3

# Problem of homoplasy:

parsimony-informative and non-informative characters (*sites*)
  - invariant sites (*symplesiomorphies*)
  - singletons (*autapomorphies*)

index of consistency, CI
retention index, RI
rescaled consistency index, RC
homoplasy index, HI

$$RC = CI \times RI$$
$$HI = 1 - CI$$

*m* = min. no. of possible steps
*s* = min. no. needed for explaining the tree
*g* = max. no. of steps for any tree

# Metods of parsimony:

Fitch:         X → Y a Y → X
               neseřazené znaky (A → T nebo A → G etc.)

Wagner:        X → Y a Y → X
               seřazené znaky (1 → 2 → 3)

Dollo:         X → Y a Y → X, potom nelze X → Y

*… restriction-site and
restriction-fragment data*

Camin-Sokal:   X → Y,
               not Y → X
… SINE, LINE



*"relaxed Dollo criterion"*

weighted = transversion p.

generalized p.: cost matrix = step matrix

**a)** Wagner

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

**b)** Fitch

|   | a | b | c | d |
|---|---|---|---|---|
| a | - | 1 | 1 | 1 |
| b | 1 | - | 1 | 1 |
| c | 1 | 1 | - | 1 |
| d | 1 | 1 | 1 | - |

**c)** Dollo

|   | a | b | c | d |
|---|---|---|---|---|
| a | - | *M*\*) | 2*M* | 3*M* |
| b | 1 | - | *M* | 2*M* |
| c | 2 | 1 | - | *M* |
| d | 3 | 2 | 1 | - |

**d)** transversion

|   | A | C | G | T |
|---|---|---|---|---|
| A | - | 5 | 1 | 5 |
| C | 5 | - | 5 | 1 |
| G | 1 | 5 | - | 5 |
| T | 5 | 1 | 5 | - |

\*) *M* is an arbitrarily large number, guaranteeing that only one transformation to each derived state will be permitted.

# Parsimony and consistency



"true"

A   $p$   $q$   $q$   C

$p>>q$

B   D

((A,B),(C,D))*
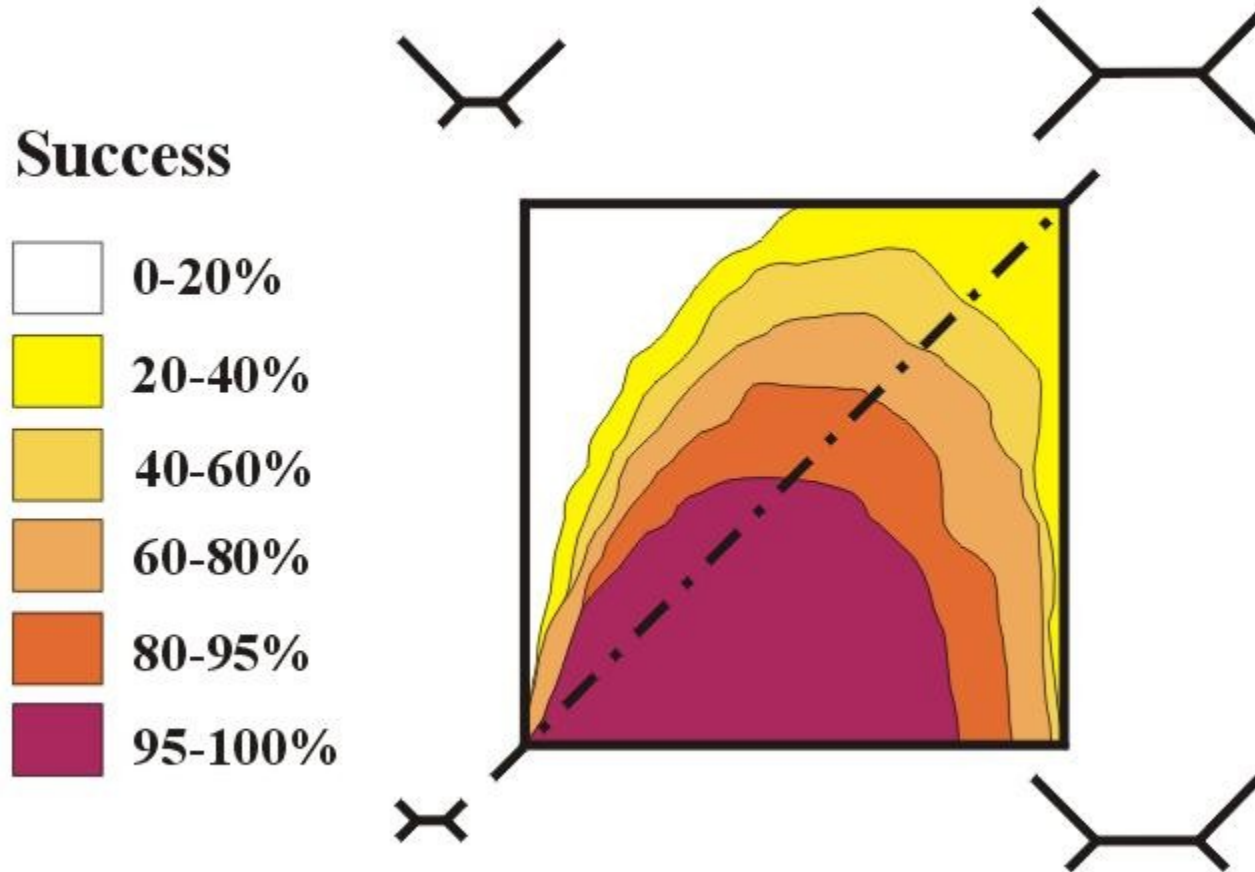
"wrong"

((A,C),(B,D))*

\* tree written in Newick format
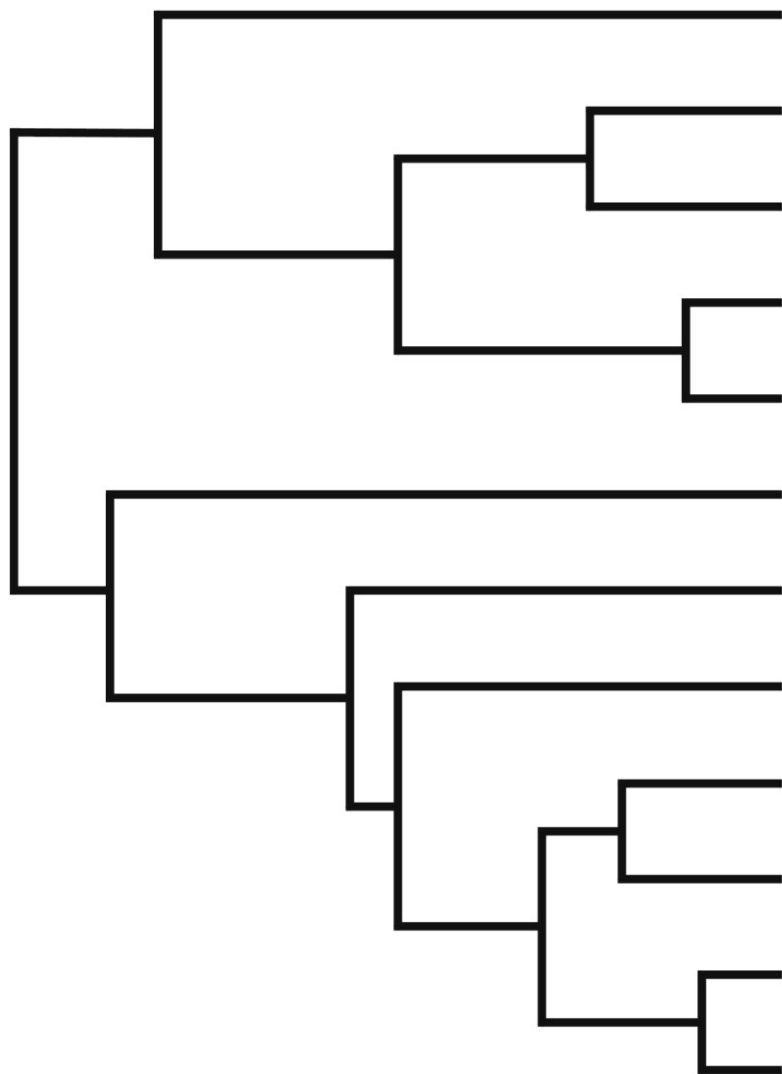
# Parsimony and consistency



In the Felsenstein zone, parsimony is inconsistent
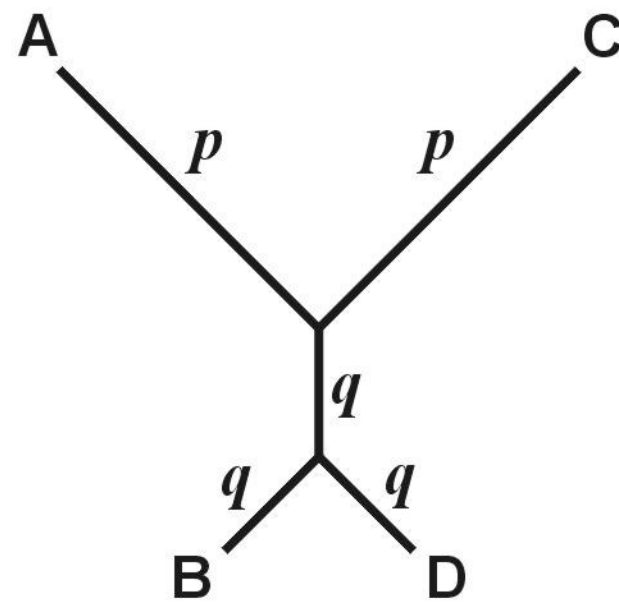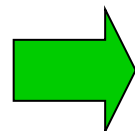
# Parsimony and consistency

# Parsimony and consistency



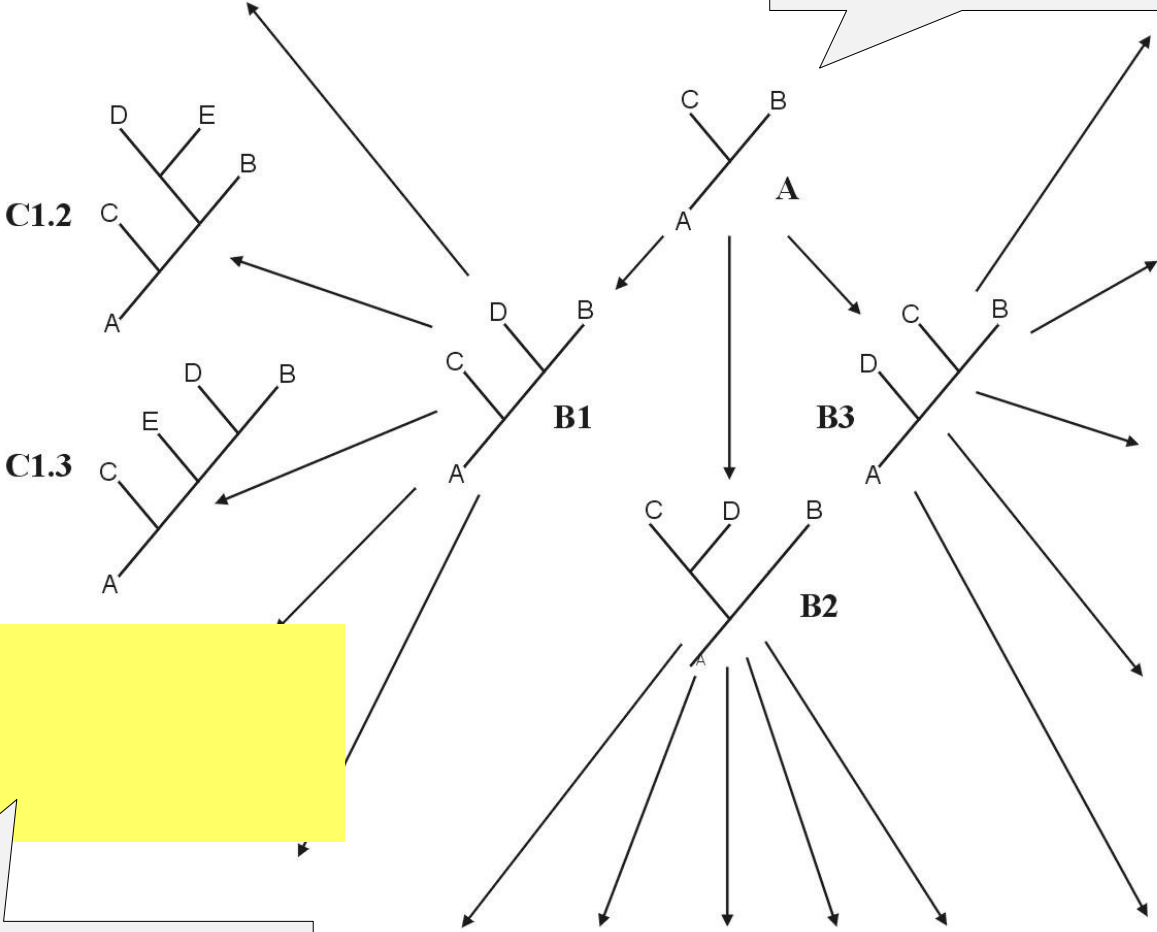long branches

long-branch attraction (LBA)

# Search for optimal tree

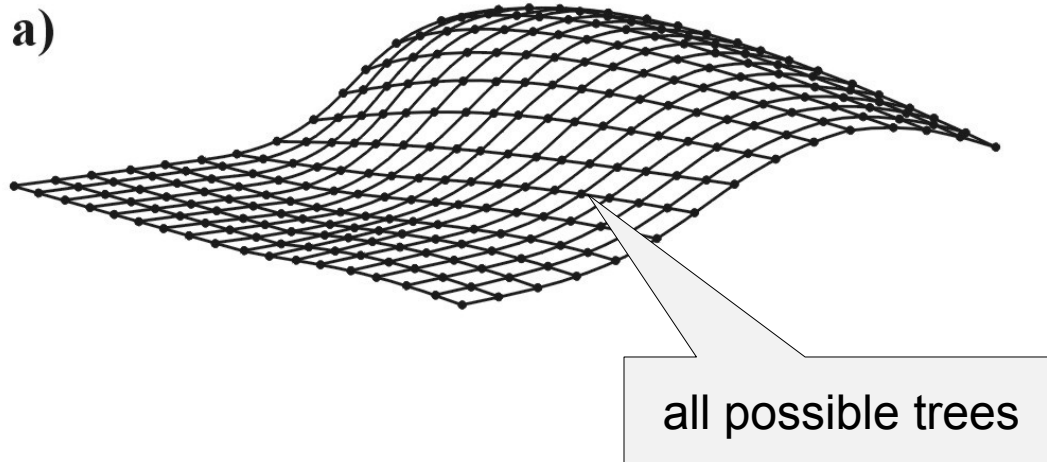1. Exact methods:

   a) exhaustive search

   b) branch-and-bound

# branch-and-bound

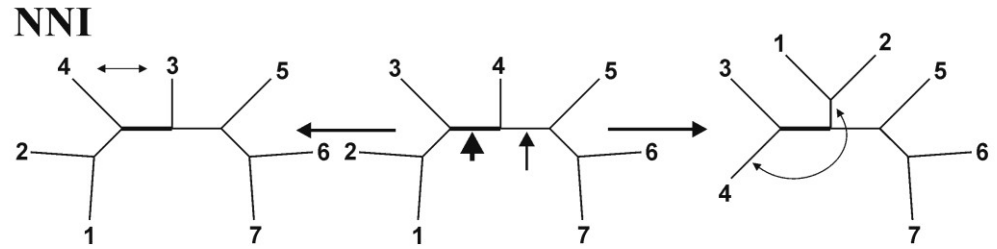C1.2

C1.3

B1

B3

B2

A

a)

all possible trees

stepwise addition

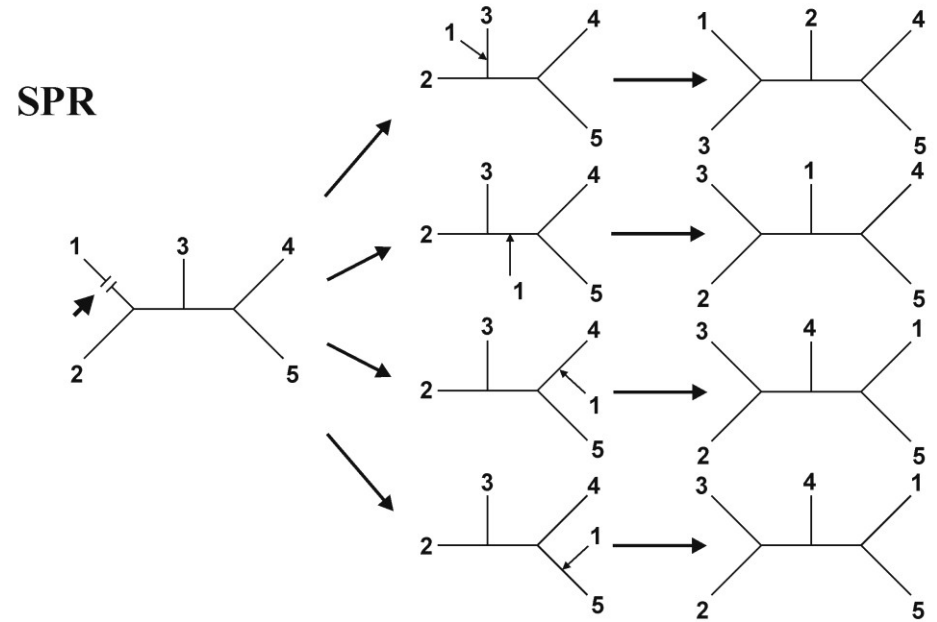star decomposition

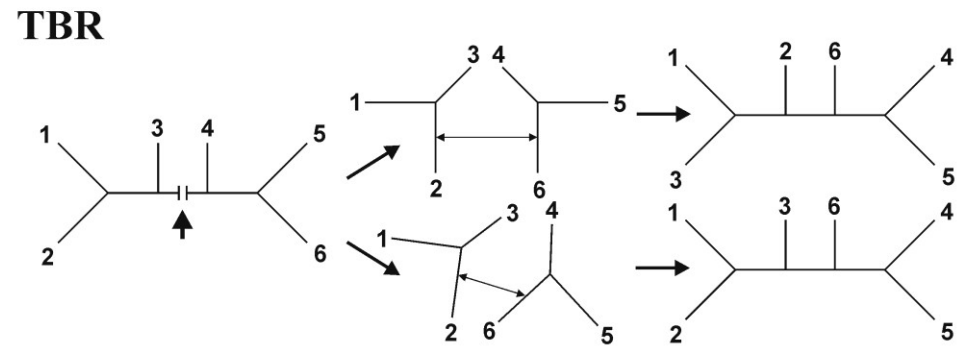branch swapping

a)

b)

heuristic search

nearest-neighbor
interchanges (NNI)

subtree prunning
and regrafting (SPR)

tree bisection and
reconnection (TBR)
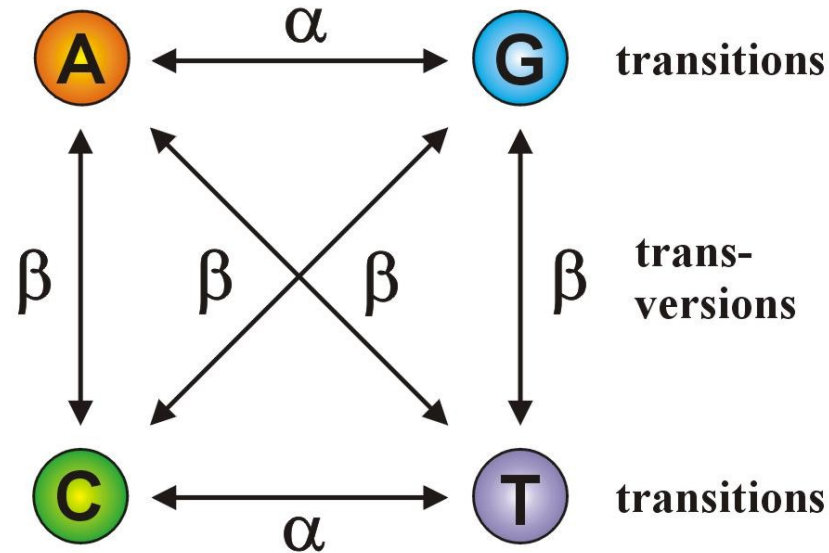
# Evolutionary models and distance methods

Base after substitution

| Original base | | A | C | G | T |
|---|---|---|---|---|---|
| | A | -¾ | ¼ | ¼ | ¼ |
| | C | ¼ | -¾ | ¼ | ¼ |
| | G | ¼ | ¼ | -¾ | ¼ |
| | T | ¼ | ¼ | ¼ | -¾ |

$$Q = \begin{pmatrix} - & \alpha & \alpha & \alpha \\ \alpha & - & \alpha & \alpha \\ \alpha & \alpha & - & \alpha \\ \alpha & \alpha & \alpha & - \end{pmatrix}$$

**Jukes-Cantor (JC):**    equal base frequencies
equal substitution rates

**Kimura 2-parameter (K2P):** transitions ≠ transversions

$$\mathbf{Q} = \begin{pmatrix} - & \beta & \alpha & \beta \\ \beta & - & \beta & \alpha \\ \alpha & \beta & - & \beta \\ \beta & \alpha & \beta & - \end{pmatrix}$$

If $\alpha = \beta$, K2P = JC

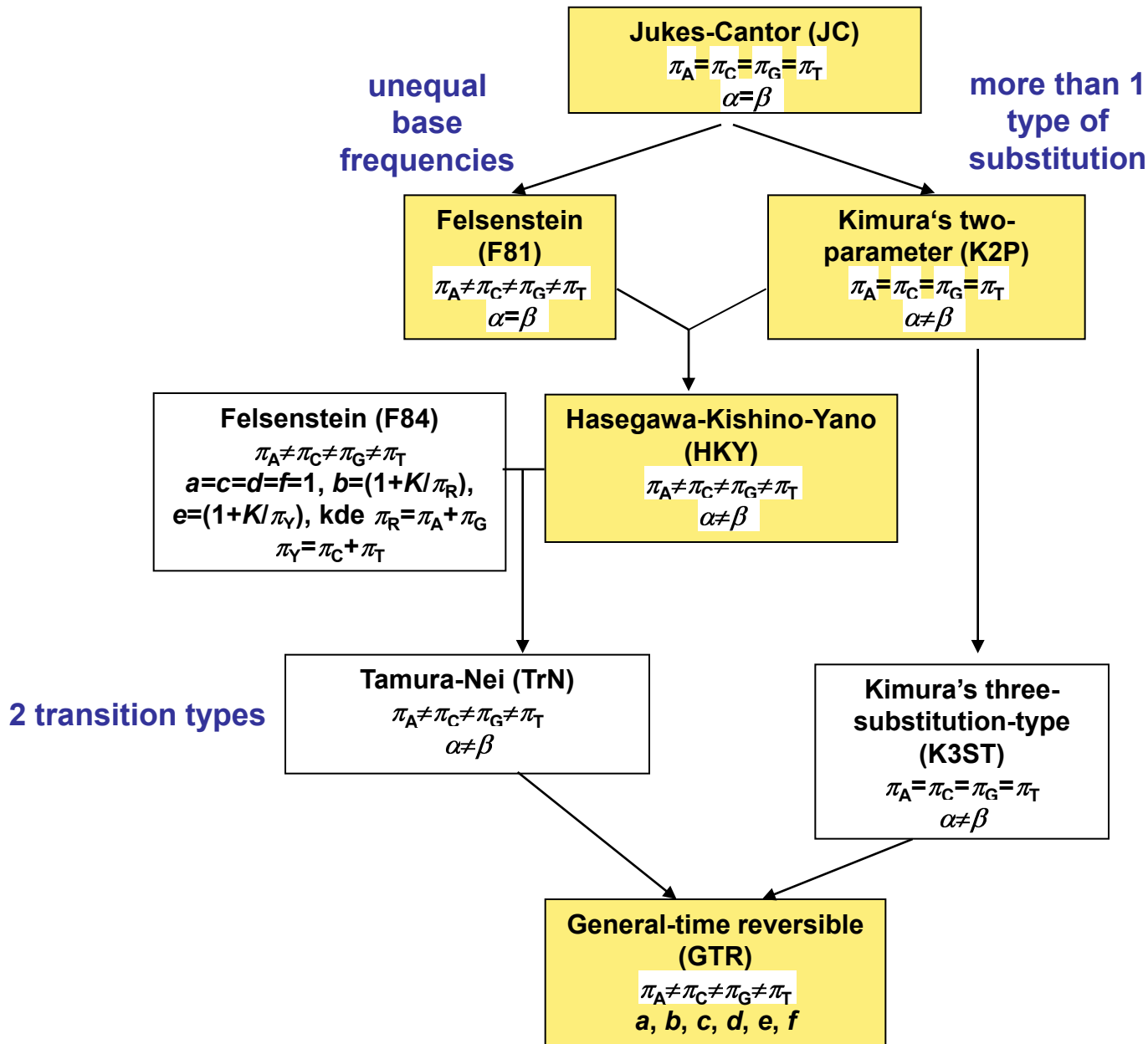**Felsenstein (F81):** different base frequencies

$$Q = \begin{pmatrix} - & \pi_C & \pi_G & \pi_T \\ \pi_A & - & \pi_G & \pi_T \\ \pi_A & \pi_C & - & \pi_T \\ \pi_A & \pi_C & \pi_G & - \end{pmatrix}$$

If $\pi A = \pi C = \pi G = \pi T$, F81 = JC

**Hasegawa-Kishino-Yano (HKY):** different base frequencies
transitions ≠ transversions

$$Q = \begin{pmatrix} - & \pi_C\beta & \pi_G\alpha & \pi_T\beta \\ \pi_A\beta & - & \pi_G\beta & \pi_T\alpha \\ \pi_A\alpha & \pi_C\beta & - & \pi_T\beta \\ \pi_A\beta & \pi_C\alpha & \pi_G\beta & - \end{pmatrix}$$

**General time-reversible (GTR, REV):** different base frequencies
different substitution rates

**Jukes-Cantor (JC)**
$\pi_A=\pi_C=\pi_G=\pi_T$
$\alpha=\beta$

**unequal base frequencies**

**more than 1 type of substitution**

**Felsenstein (F81)**
$\pi_A\neq\pi_C\neq\pi_G\neq\pi_T$
$\alpha=\beta$

**Kimura's two-parameter (K2P)**
$\pi_A=\pi_C=\pi_G=\pi_T$
$\alpha\neq\beta$

**Felsenstein (F84)**
$\pi_A\neq\pi_C\neq\pi_G\neq\pi_T$
$a=c=d=f=1$, $b=(1+K/\pi_R)$,
$e=(1+K/\pi_Y)$, kde $\pi_R=\pi_A+\pi_G$
$\pi_Y=\pi_C+\pi_T$

**Hasegawa-Kishino-Yano (HKY)**
$\pi_A\neq\pi_C\neq\pi_G\neq\pi_T$
$\alpha\neq\beta$

**2 transition types**

**Tamura-Nei (TrN)**
$\pi_A\neq\pi_C\neq\pi_G\neq\pi_T$
$\alpha\neq\beta$

**Kimura's three-substitution-type (K3ST)**
$\pi_A=\pi_C=\pi_G=\pi_T$
$\alpha\neq\beta$

**General-time reversible (GTR)**
$\pi_A\neq\pi_C\neq\pi_G\neq\pi_T$
$a, b, c, d, e, f$

# Heterogenity of substitution rates in different parts of sequences

Gamma distribution:

shape parameter α

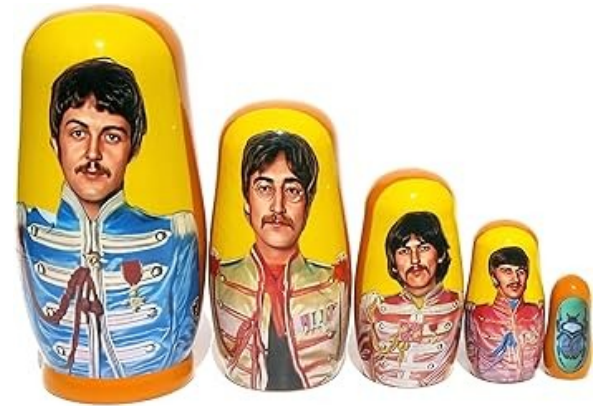discrete gamma model

invariant sites
  → GTR+Γ+I
  nebo GTR+G+I

the higher $\alpha$, the more homogeneous are substitutions

# Model comparison:

## Likelihood ratio test (LRT):

nested models

$LR = 2(\ln L2 - \ln L1)$

$\chi^2$ distribution, p2 – p1 degrees of freedom

## Akaike information criterion (AIC):

nonnested models

$AIC = -2\ln L + 2p$, kde $p$ = number of free parametres

better model $\rightarrow$ lower $AIC$

## Bayesian information criterion (BIC):

nonnested models

$BIC = -2\ln L + p\ln N$, where $N$ = sample size

# Model comparison:

## hierarchical LRT – ModelTest (Crandall and Posada), jModelTest

# Model comparison:

dynamic LRT:

# Model comparison:



More parametres $\Rightarrow$ more realism, but …

… also less confidence (estimates based on the same amount of data!)

# **Distances**

computed for each pair of taxa, from distance (or similarity) matrix
– tree inference

distance methods base on assumption that if we know true distances,
we can very easily infer the true phylogeny

advantage: very fast and simple (also with a calculator)

```
                1              10            20          30
sequence 1:   ACCCGTTAAGCTTAACGTACTTGGATCGAT
sequence 2:   ACCCGTTAGGCTTAATGTACGTGGATCGAT
```

*p*-distance:  $p = k/n = 3/30 = 0{,}10$

problem of
saturation:

Distances for some models:

# Cluster analysis - UPGMA

|              | chimp  | bonobo | gorilla | human  | orang. |
|--------------|--------|--------|---------|--------|--------|
| chimp (Š)    | --     |        |         |        |        |
| bonobo (B)   | 0,0118 | --     |         |        |        |
| gorilla (G)  | 0,0427 | 0,0416 | --      |        |        |
| human (Č)    | 0,0382 | 0,0327 | 0,0371  | --     |        |
| orangutan (O)| 0,0953 | 0,0916 | 0,0965  | 0,0928 | --     |

1. Find min $d(ij)$

2. Calculate new matrix (ŠB-$k$) = [$d$(B-$k$)+$d$(Š-$k$)]/2

3. Repeat 1 a 2.

|          | ŠB     | gorilla | human  | orang. |
| -------- | ------ | ------- | ------ | ------ |
| ŠB       | --     |         |        |        |
| gorilla (G) | 0,0422 | --      |        |        |
| human (Č)  | 0,0355 | 0,0371  | --     |        |
| orangutan (O) | 0,0935 | 0,0965  | 0,0928 | --     |

UPGMA (unweighted pair-group method using arithmetic means):

d[(BŠČ)G] = {d(BG)+d(ŠG)+d(ČG)}/3

WPGMA: d[(BŠČ)G] = {d[(BŠ)G] + d(ČG)}/2

single-linkage (metoda nejbližšího souseda)

complete-linkage (m. nejvzdálenějšího souseda)

# UPGMA and consistency

additive distances: $d_{AB} + d_{CD} \leq \max(d_{AC} + d_{BD}, d_{AD} + d_{BC})$

tj. distance between 2 taxa equals sum of branches connecting them



ultrametric distances: $d_{AC} \leq \max(d_{AB}, d_{BC})$





**additive tree**

**ultrametric tree**

# UPGMA and consistency

# Neighbor-Joining, NJ

Algorithmic method

Principle of minimal evolution $\rightarrow$ minimizes sum of branch lenghts $S$

Each pair of nodes adjusted according to its divergence from others

Single additive tree

**a)**

a) star tree

b) finding nearest neighbors

c) distance recalculation
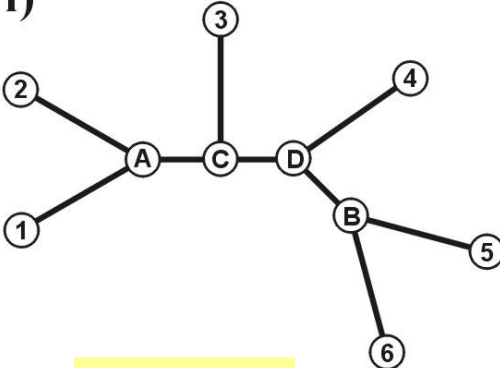
# Drawbacks of distance data:

1. loss of information during transformation

2. after transformation to distances, we cannot infer original data (different sequences may result in the same distance)

3. we cannot study the evolution in different parts of sequence

4. difficult biological interpretation of branch lengths

5. we cannot combine more distance matrices