

# Molecular Phylogeny and Phylogeny Tree Construction

## 分子系统发育学及系统发育树构建



孙田舒

指导老师： 罗静初 教授

特别感谢： 康菊清 博士

北京大学生命科学学院

# What is molecular phylogeny

- **Phylogeny** is the study of the **evolutionary history** of living organisms using tree like diagrams to represent pedigrees of these organisms.
- **Molecular phylogeny**: The study of evolutionary relationships of genes and other biological macromolecules by analyzing **mutations at various** positions in their sequences and developing hypotheses about the evolutionary relatedness of the biomolecules.



# Fossil Records Vs. Molecular Approach

- **Fossil records: hard evidence, often biased**

*Limitations:*

Available for certain species

Limited by abundance, habitat est.

Fragmentary and ambiguous

- **Molecular approach: DNA or Protein sequences**

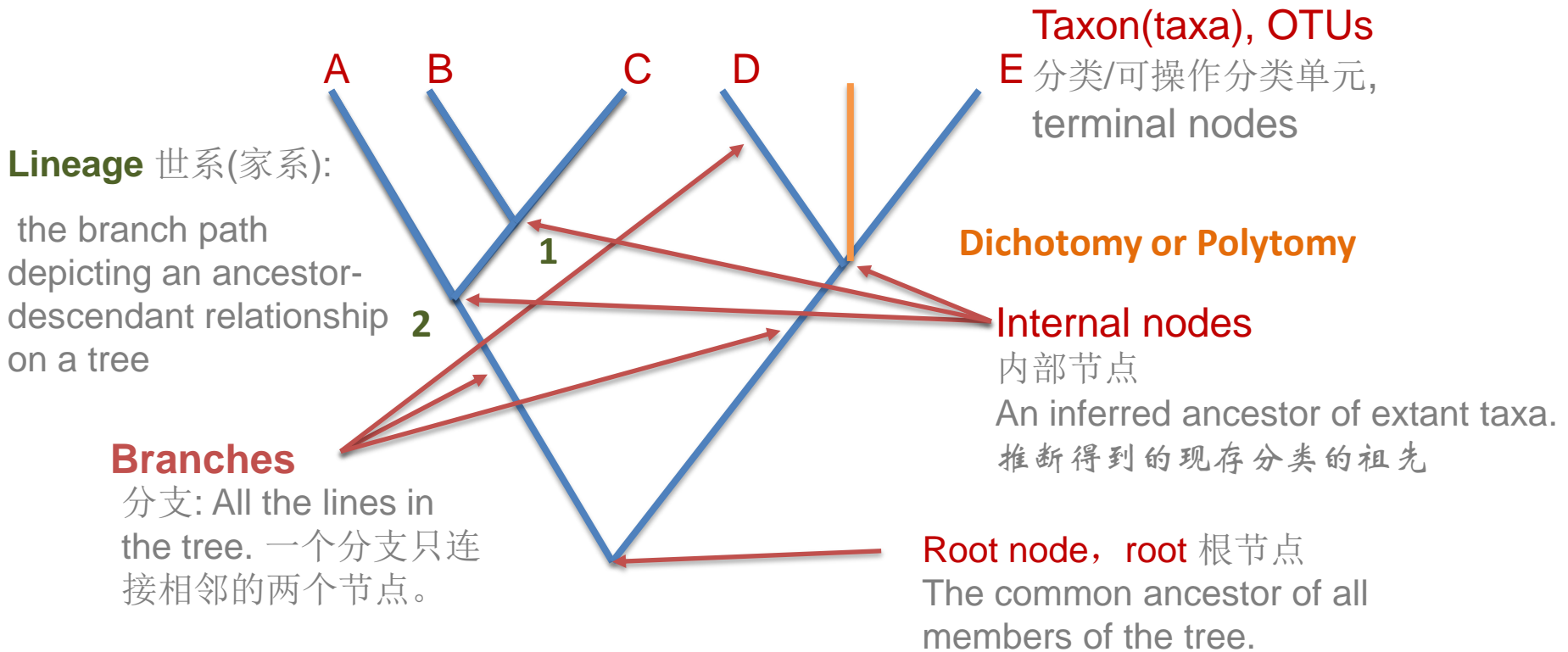
*Assumptions:*

Molecular sequences are homologous

Each position in a sequence evolved independently



# Terminology



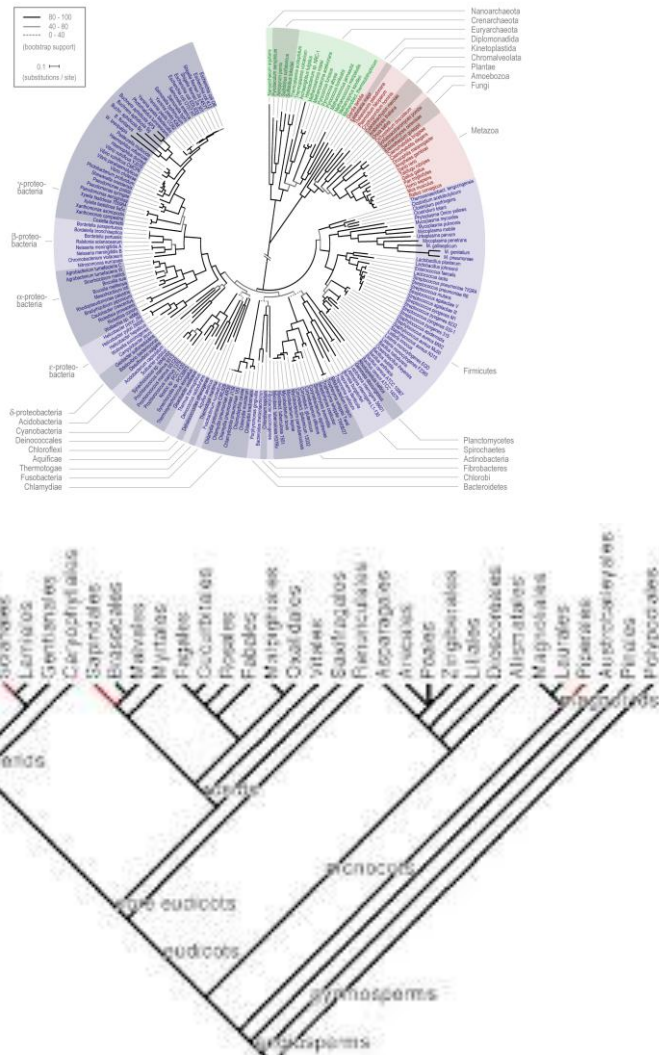
Any way, node is a taxonomic unit, it locates in the bifurcating branch point.

## Clades/monophyletics 进化枝/单系类群

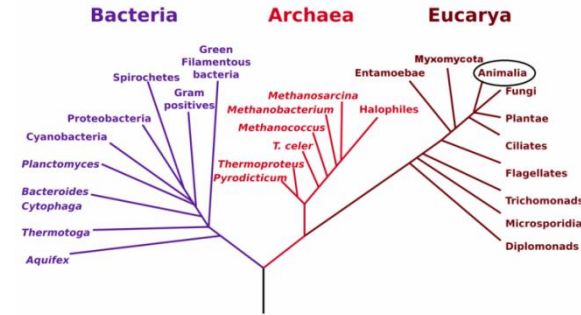
A group of taxon consists of a single common ancestor and all its descendants.



# Types of phylogenetic trees



## Phylogenetic Tree of Life



Papillomavirus Phylogeny 1249

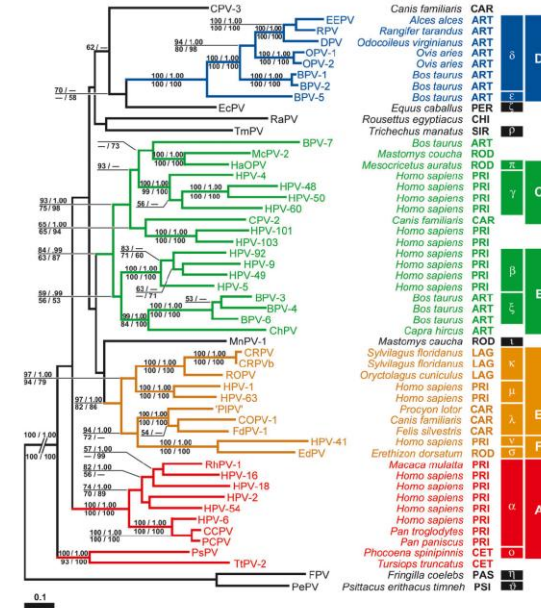


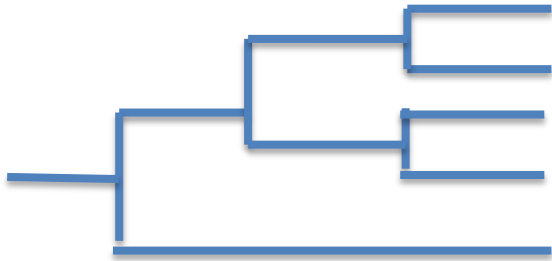
FIG. 2.—ML tree of 53 phylogenetically representative PVs as inferred from a combined E1-E2-L1 amino acid sequence analysis (1,082 parsimony-informative positions) justified by PHTs (table 2). All non-human PVs and 18 representative HPV types were used for analyses. PV genera (de Villiers et al. 2004) are indicated by Greek lettering; upper case lettering follow an alternative E1-E2 classification of Bravo and Alonso (2007). Higher order host taxa are abbreviated as follows: ART, Artiodactyla; CAR, Carnivora; CET, Cetacea; CHI, Chiroptera; LAG, Lagomorphs; PAS, Passeriformes; PER, Perissodactyla; PRI, Primates; PSI, Psittaciformes; ROD, Rodentia; and SIR, Sirenia. The supertaxa are colored blue (δ+ ε), ochre (κ+ μ+ ν+ ν+ σ), green (τ+ γ+ β+ ζ), and red (α+ θ), respectively. Branch lengths are drawn to scale, with the scale bar indicating the number of amino acid substitutions per site. Numbers on branches are bootstrap support values to clusters on the right of them (above: criteria = ML/Bayesian probabilities; below: criteria = MP/Distance; values under 50 are not shown).

Only different in the way its presented

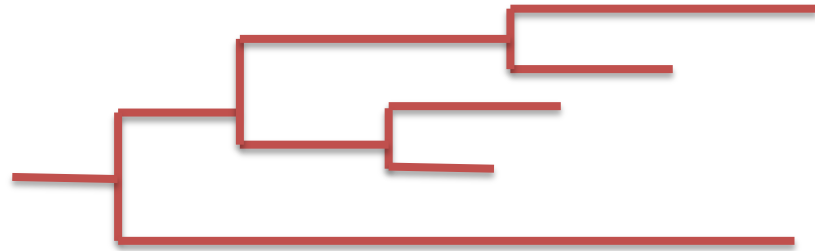


# Unscaled tree VS. Scaled tree

**Cladogram** 进化分支图



**Phylogram** 系统发育图



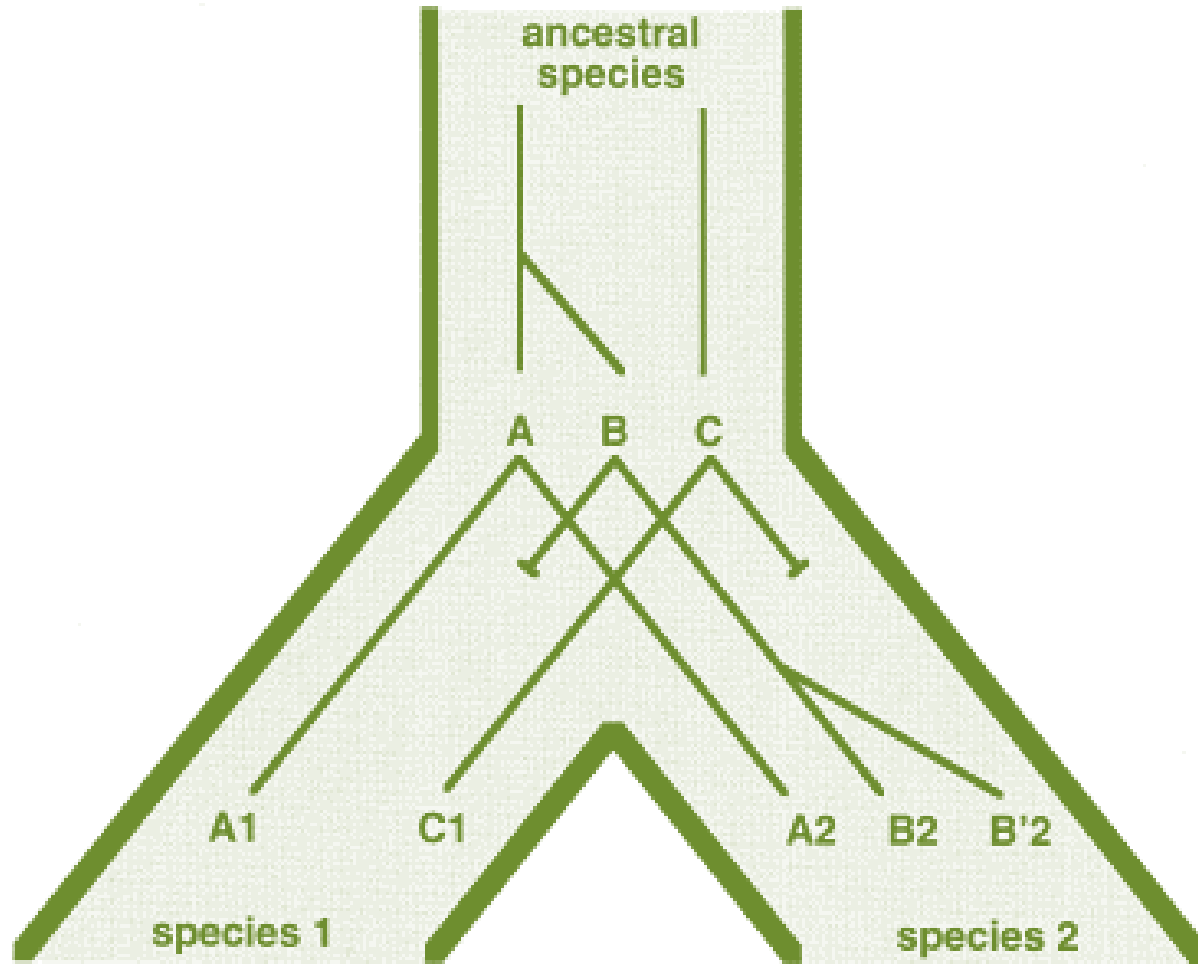
**Unscaled trees** - all branches in the tree are the same length. (only topology)

**Scaled trees** - branches will be different lengths based on the number of evolutionary changes or distance. (branch length & topology)

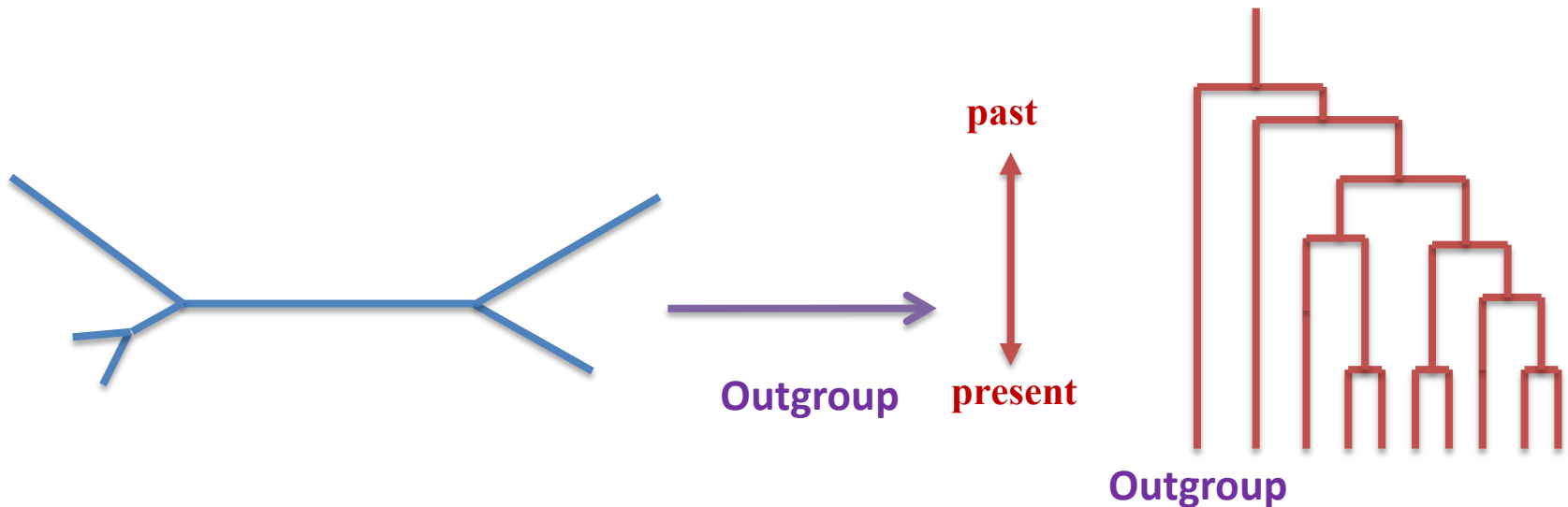
Difference: **phylogram** is scaled, but **cladogram** not.



# Gene tree VS. Species tree



# Unrooted tree Vs. rooted tree



## Unrooted tree

not assume knowledge of a common ancestor, but only positions the taxa to show their relative relationships.

## Rooted tree

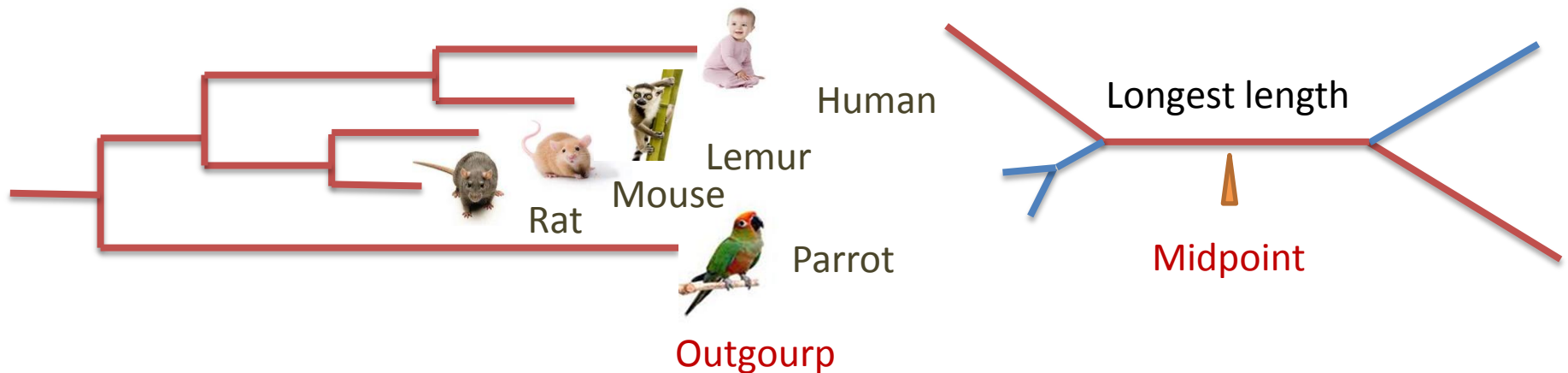
All sequences have a common ancestor or root node which a unique evolutionary path leads to all other nodes.





# Rooting Approach

- **Outgroup**, is a sequence that is homologous to the sequences under consideration, but separated from those sequences at an early evolutionary time.

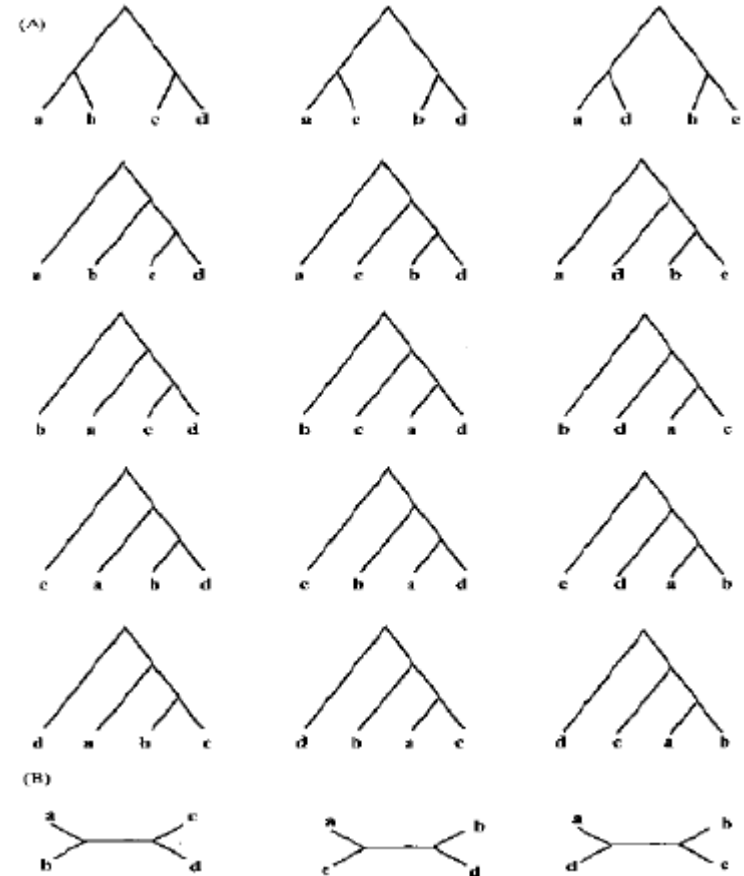


- **Midpoint**, the midpoint of the two most divergent groups judged by overall branch lengths is assigned as the root.



# Finding a true tree is difficult

Species	Unrooted	Rooted
Number	$N_U = (2n-5)! / 2^{n-3} (n-3)!$	$N_R = (2n-3)! / 2^{n-2} (n-2)!$
4	3	15
5	15	105
6	105	945
10	2027025	34459425
20	$2.22 \times 10^{20}$	$8.20 \times 10^{21}$
50	$2.84 \times 10^{74}$	$2.75 \times 10^{78}$



# Tree construction

①

- Choosing molecular marker

②

- Performing multiple sequence alignment

③

- Determining a tree building method

④

- Assessing tree reliability



# Protein or Nucleotide Sequence

The decision depends on the properties of the sequences and the purposes of the study. Generally, for studying very closed organisms, nucleotide sequences can be used. While if the phylogenetic relationships to be delineated are at the deepest level using protein sequences makes more sense.

- **Protein sequences' characters:**

- ① More conserved: 61 codons → 20 AAs
- ② No different evolutionary rates
- ③ No preferential codon usage
- ④ More sensitive alignment
- ⑤ Gap doesn't cause frameshift errors

- **Nucleotide sequences' advantages:**

- ① Rapid evolutionary rates can be informative for closely related sequences.
- ② depict synonymous and nonsynonymous substitutions, revealing evidence of positive or negative selection.



# Origin of Domestic Dogs



- mitochondrial DNA (mtDNA)
- 654 worldwide domestic dogs
- A larger genetic variation in East Asia than in other regions and the pattern of phylogeographic variation suggest an **East Asian** origin for the domestic dog, ~15,000 years ago.

# Controversy and Resampling

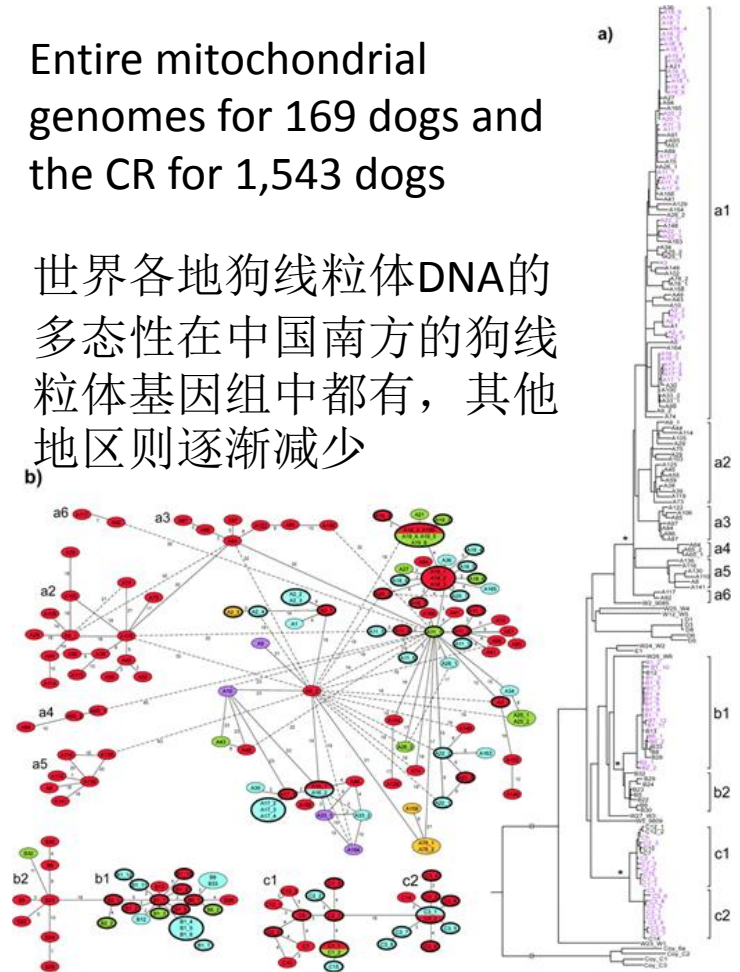


Surprisingly, we find similar mtDNA haplotype diversity in African and East Asian village dogs, potentially calling into question the hypothesis of an East Asian origin for dog domestication.

Boyko et al, 2009, PNAS

Entire mitochondrial genomes for 169 dogs and the CR for 1,543 dogs

世界各地狗线粒体DNA的多态性在中国南方的狗线粒体基因组中都有，其他地区则逐渐减少

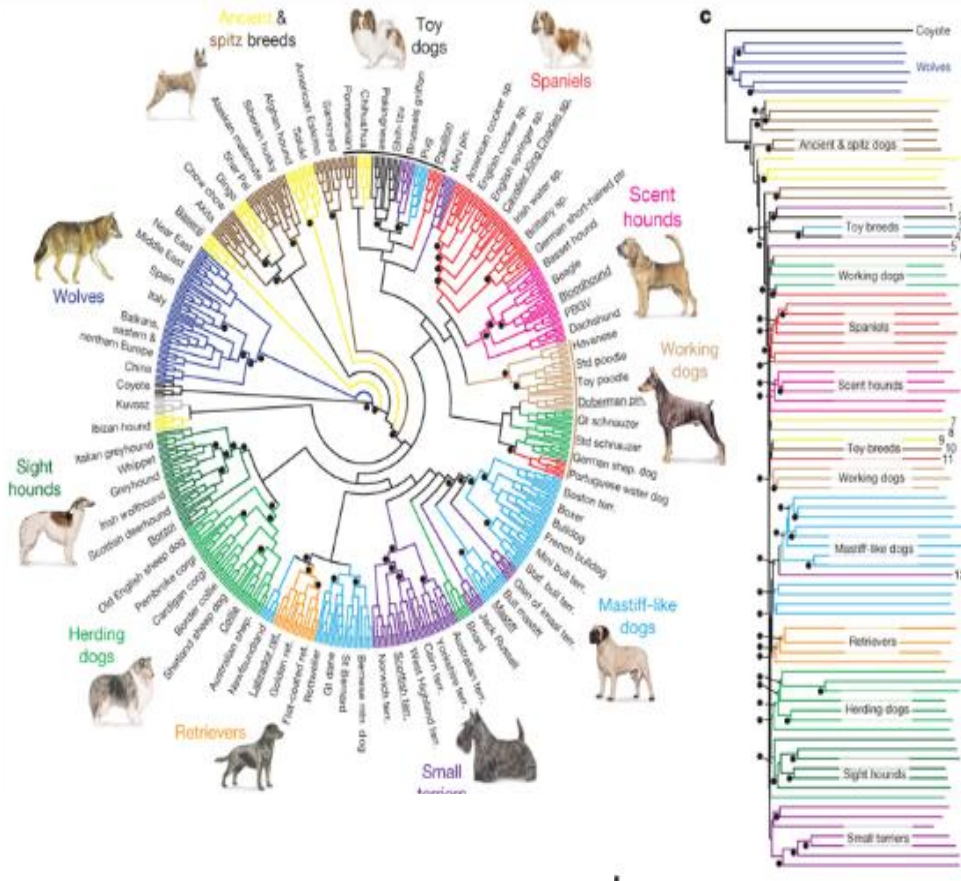


Pang et al, 2009, Molecular Biology and Evolution

狗起源于约一万多年前中国南方的几百只大灰狼



# Genome-wide SNP and haplotype analyses



Extensive genome-wide survey of more than 48,000 single nucleotide polymorphisms in dogs and their wild progenitor, the grey wolf.

**Middle Eastern** wolves were a critical source of genome diversity, although interbreeding with local wolf populations clearly occurred elsewhere in the early history of specific lineages.

Bridgett M. vonHoldt et al, 2010, Nature





# Alignment

Alignment may be the most critical step because it establishes positional correspondence in evolution.

- **What can you do?**

- ① Manual editing: correcting mismatching of key cofactor residues and residues of similar physicochemical properties
- ② Full alignment or parts of it (domain only)
- ③ Remove ambiguously aligned regions (subjective process)
- ④ Automatic approach: Rascal, NorMD and Gblocks
- ⑤ Statistical models to correct homoplasy
- ⑥ Using a  $\gamma$  correction factor to correct site-dependent rate variation

- **Choosing substitution models**

- ① Nucleotide: Juke-Cantor Model and Kimura Model
- ② Protein: PAM or JTT amino acid substitution matrix





# Phylogenetic Tree Construction Method

NJ

邻接法 Neighbor-joining

- 根据所有序列的两两比对结果，通过序列两两之间的差异决定进化树的拓扑结构和树枝长度

MP

最大简约法(Maximum Parsimony)

- 最大简约法根据序列的多重比对结果，对所有可能正确的拓扑结构进行计算并挑选出所需替代数最少的拓扑结构作为最优树，即能够利用最少的步骤去解释多重比对中的碱基差异。

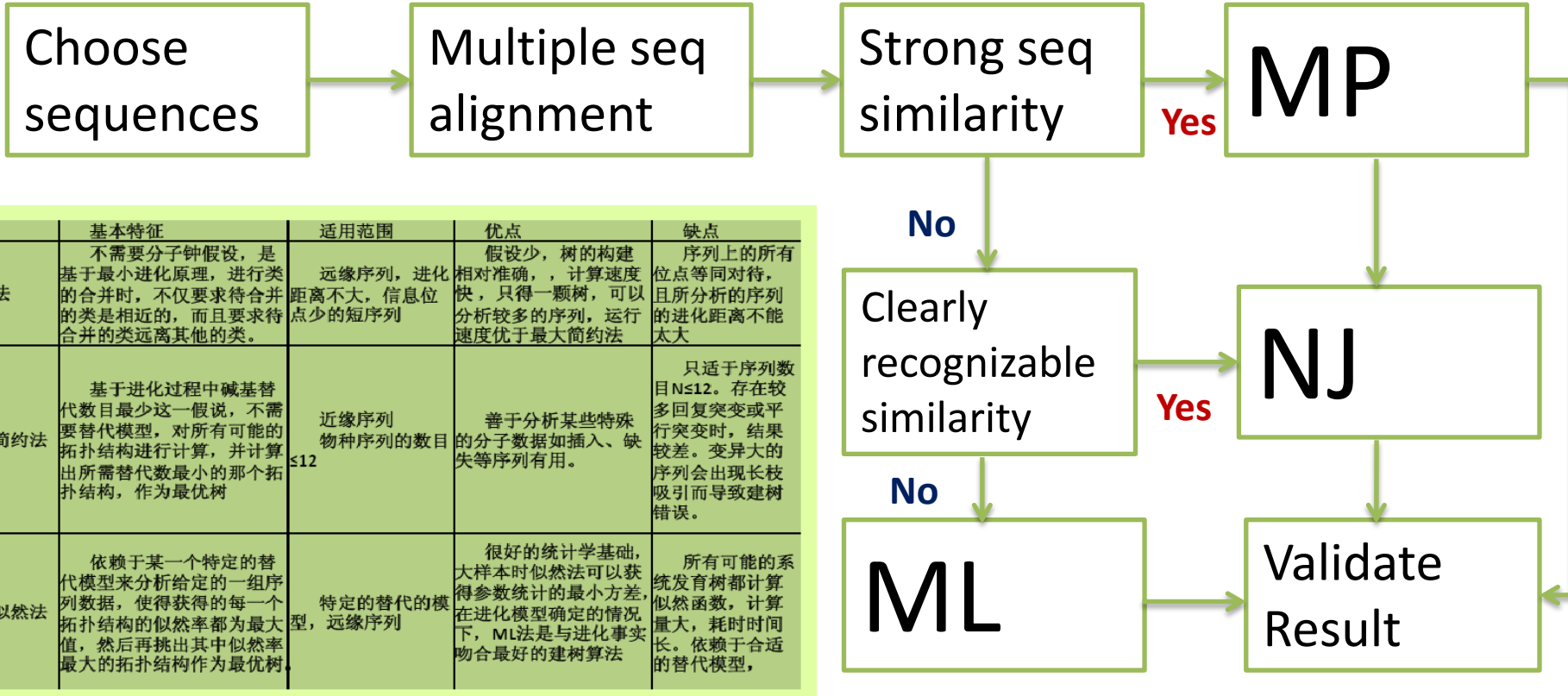
ML

最大似然法 Maximum Likelihood

- 最大似然法以一个特定的替代模型分析一组序列数据的多重比对结果，优化出拥有一定拓扑结构和树枝长度的进化树，使所获得的每一个拓扑结构的似然率均为最大，挑选似然率最大的拓扑结构作为最优树。



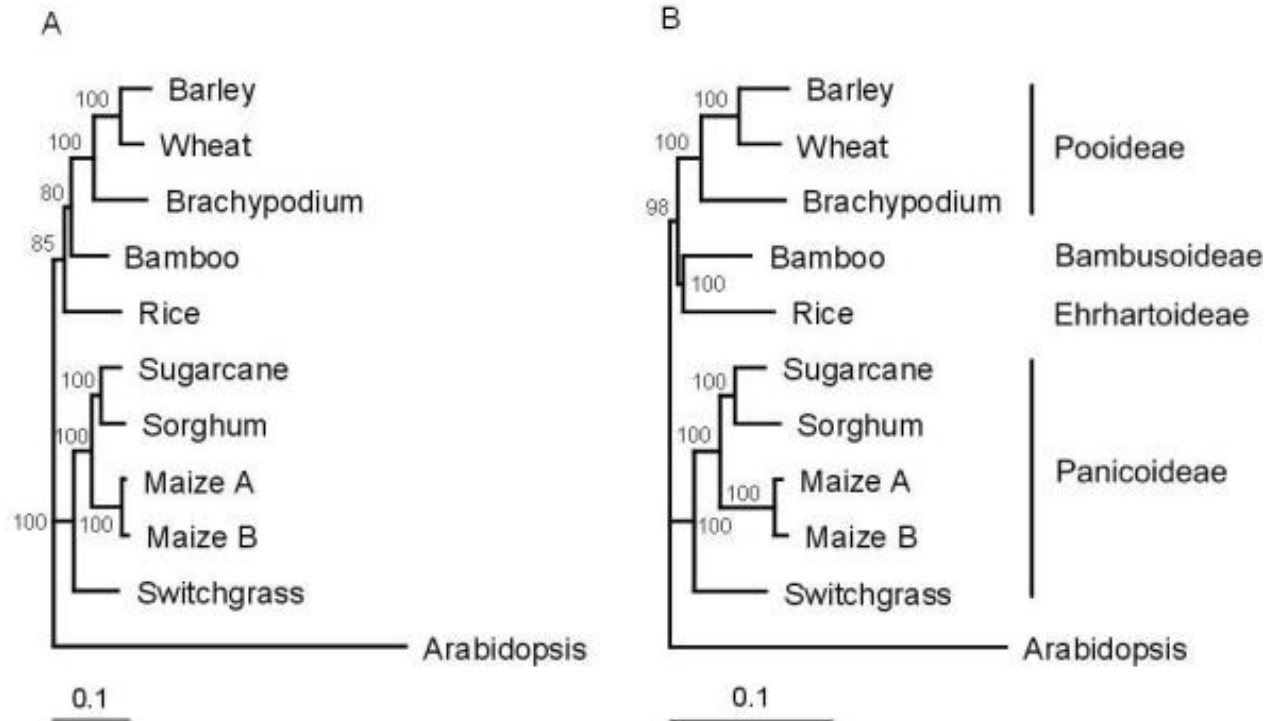
# General method



Combine Different Methods for Consensus



# Genome-wide characterization of the biggest grass, bamboo, based on 10,608 putative full-length cDNA sequences



**Phylogeny of grasses inferred from concatenated alignment of 43 putative orthologous cDNA sequences.** (A) Tree inferred from maximal likelihood method. Bayes inference yielded the same topology. (B) Tree inferred from neighbor joining method. Branch length ...



# Tree Evaluation

## Boot strapping

### 自举法

- 是对所比较序列上的替换位点作多次随机取样，根据每次取样的数据可以得到新的树形图，相同的组合出现在某一个节点上的次数占总取样次数的百分比就是该节点的bootstrap值。

## Jackknifing

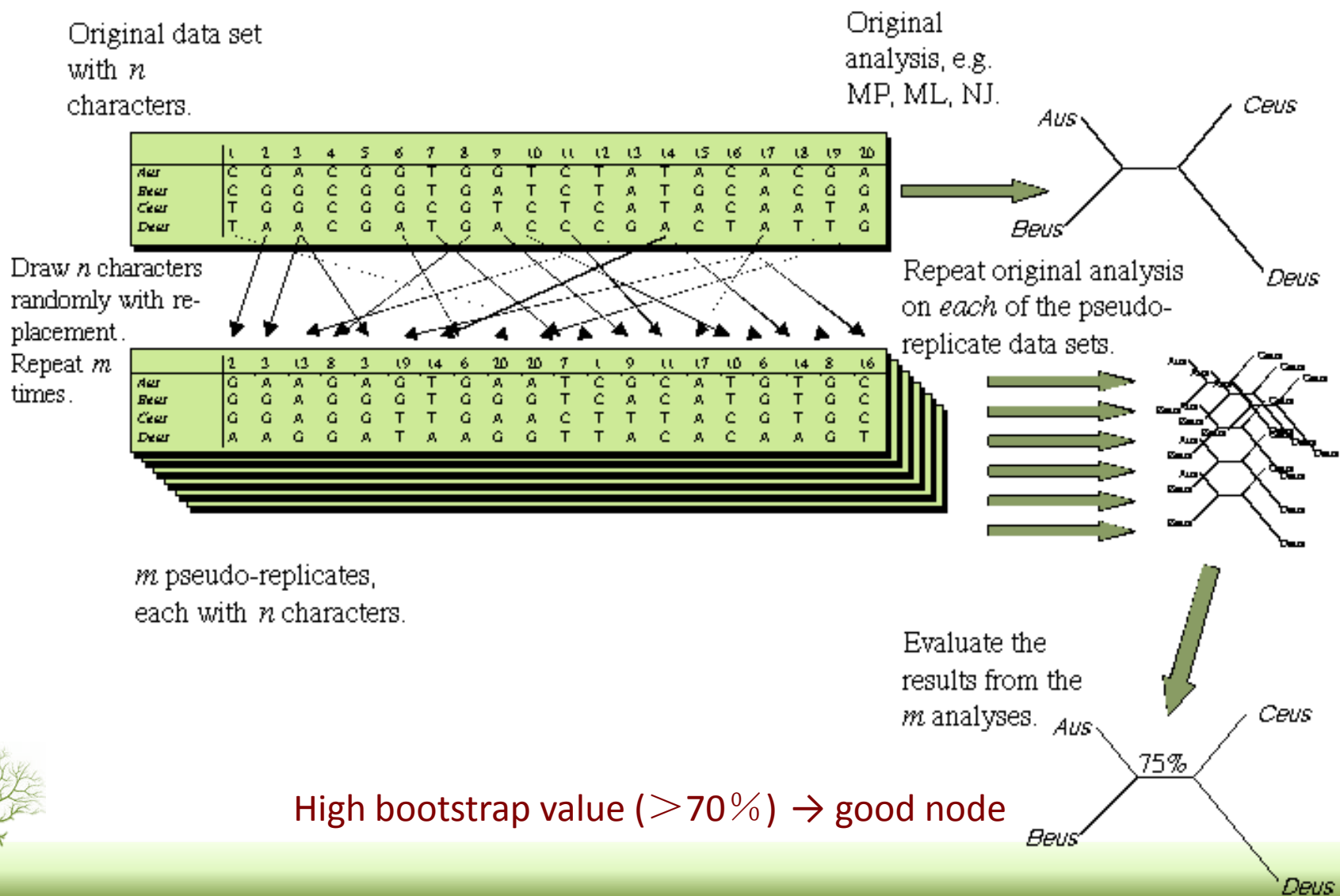
### 折刀法

- Half of the sites in a dataset are randomly deleted, each new dataset is subjected to phylogenetic tree construction using the same method as the original.

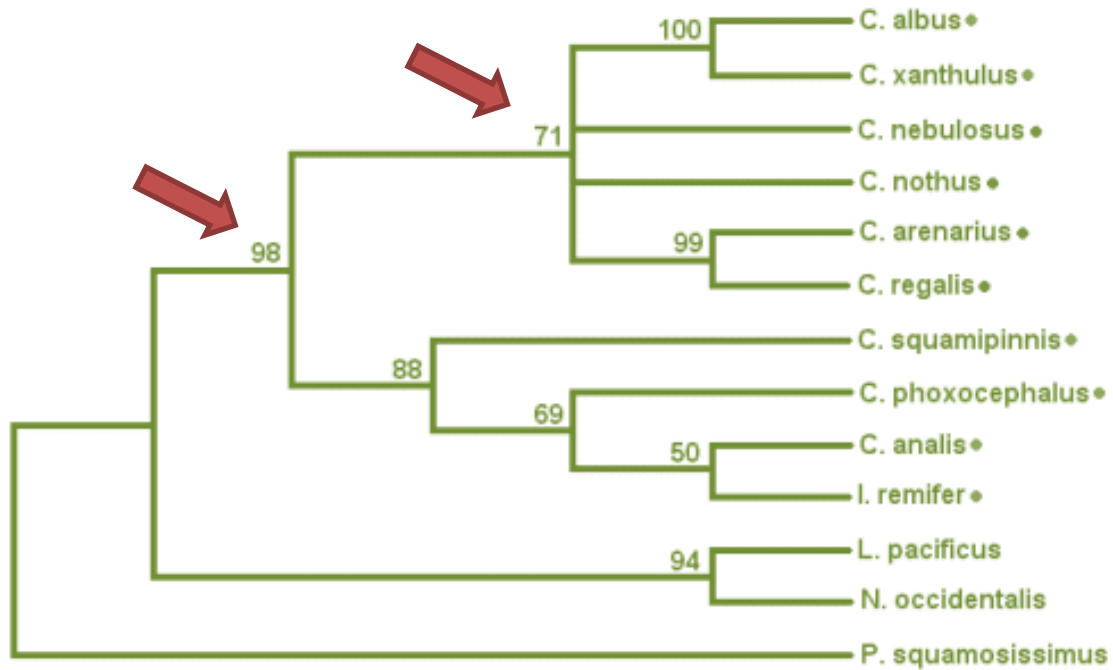


# Bootstrapping

To verify the solidity of each node



# Bootstrap Value



# Phylogenetic Tree Construction Programs

MEGA

- Molecular Evolutionary Genetics Analysis

PHYLIP

- PHYLogeny Inference Package

PAUP

- Phylogenetic Analysis Using Parsimony

PAML

- Phylogenetic Analysis by Maximum Likelihood

<http://evolution.genetics.washington.edu/phylip/software.html>

392 phylogeny packages and 54 free web servers



# Who needs molecular phylogeny tree





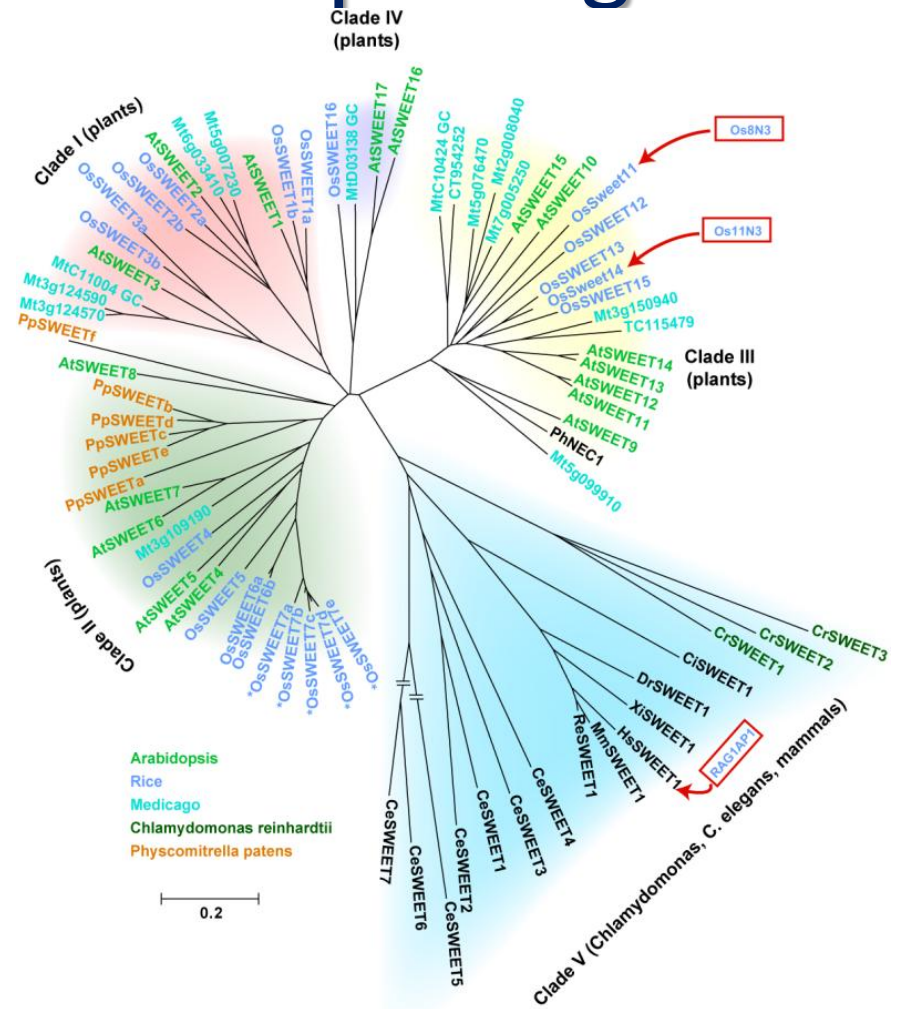
# Sugar transporters for intercellular exchange and nutrition of pathogens

Phylogenetic tree for SWEETs from different species.

Phylogenetic tree of the SWEET superfamily (PFAM PFO3083). Distances were calculated from a multiple sequence alignment (ClustalW) using the neighbor-joining method. The tree displays bootstrap values (percentage of 1000). SWEET genes fall into four clades. All sequences were obtained from NCBI or the Aramemnon database.

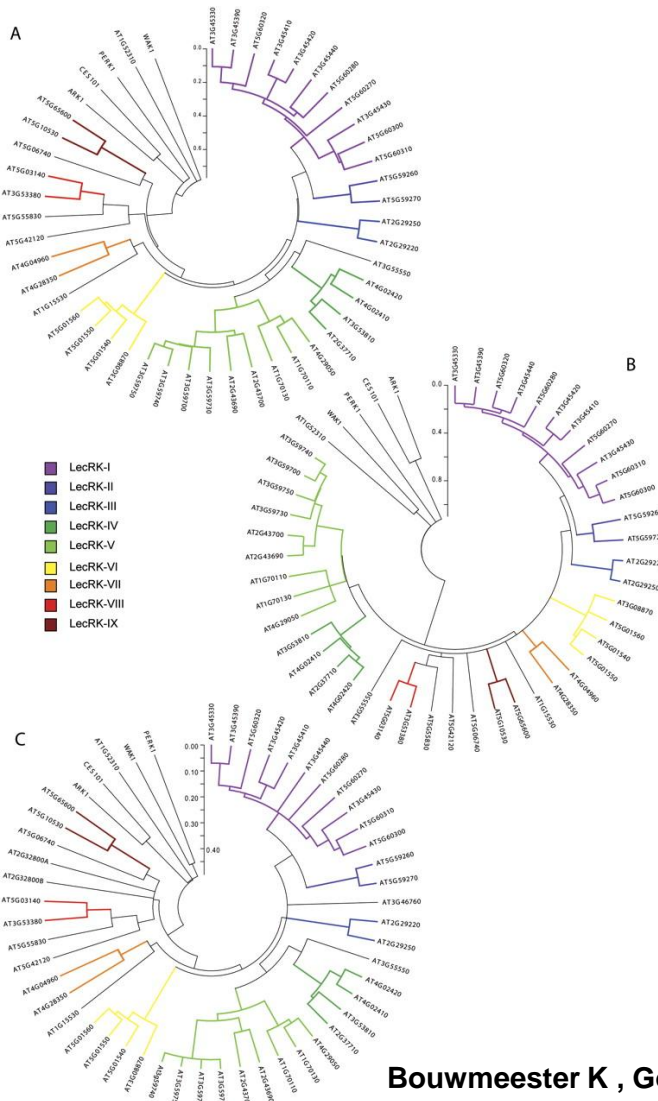
The family members were color-coded indicating the species:

- At, Arabidopsis thaliana (light green);
- Os, Oryza sativa (blue);
- Mt, Medicago trunculata (cyan);
- Chlamydomonas reinhardtii (darkgreen),
- Physcomitrella patens (orange).



Li-Qing Chen, et al. (2010)

# Phylogenetic analysis and classification of Arabidopsis LecRK proteins.

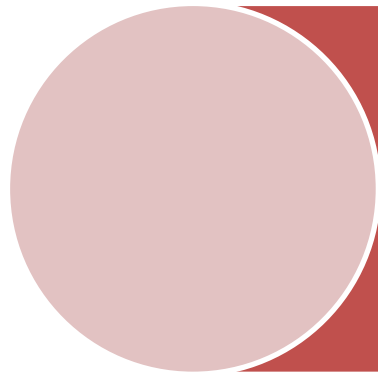


Phylogenetic analysis and classification of Arabidopsis LecRK proteins. Phylograms of (A) 43 full-length LecRK amino acid sequences, (B) 43 lectin domains and, (C) 46 kinase domains, including the kinase domain of LecRK-S.2 and the two kinase domains of LecRK-S.3. Each LecRK clade is depicted by a different colour.

Bouwmeester K , Govers F J. Exp. Bot. 2009;60:4383-4396



# Application to the real world



Satisfying One's  
Curiosity

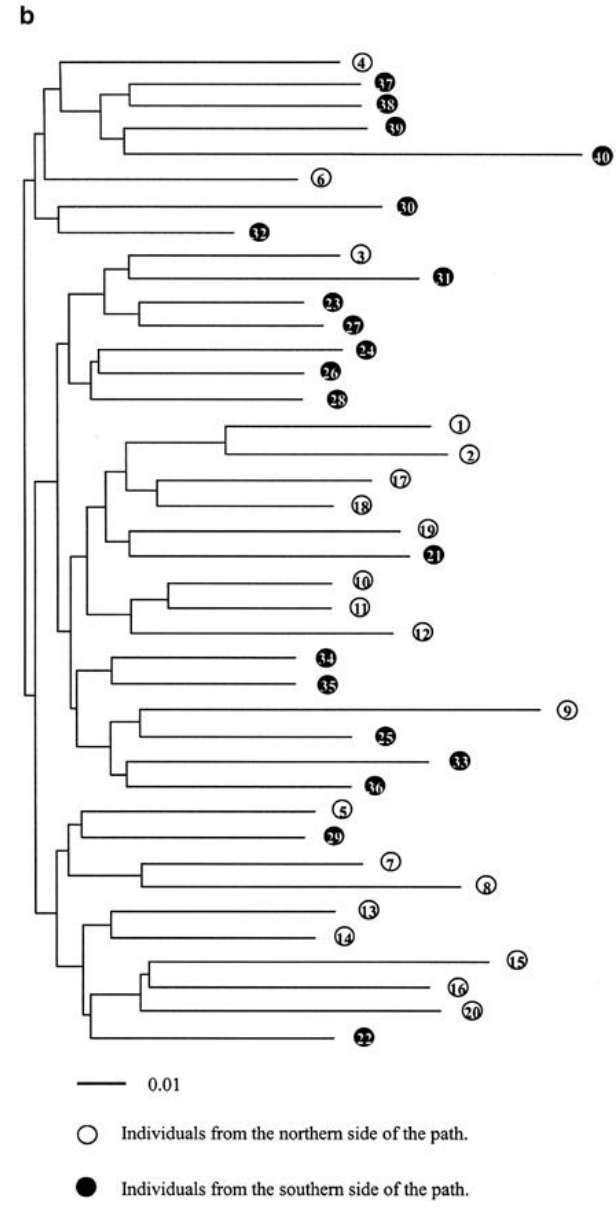
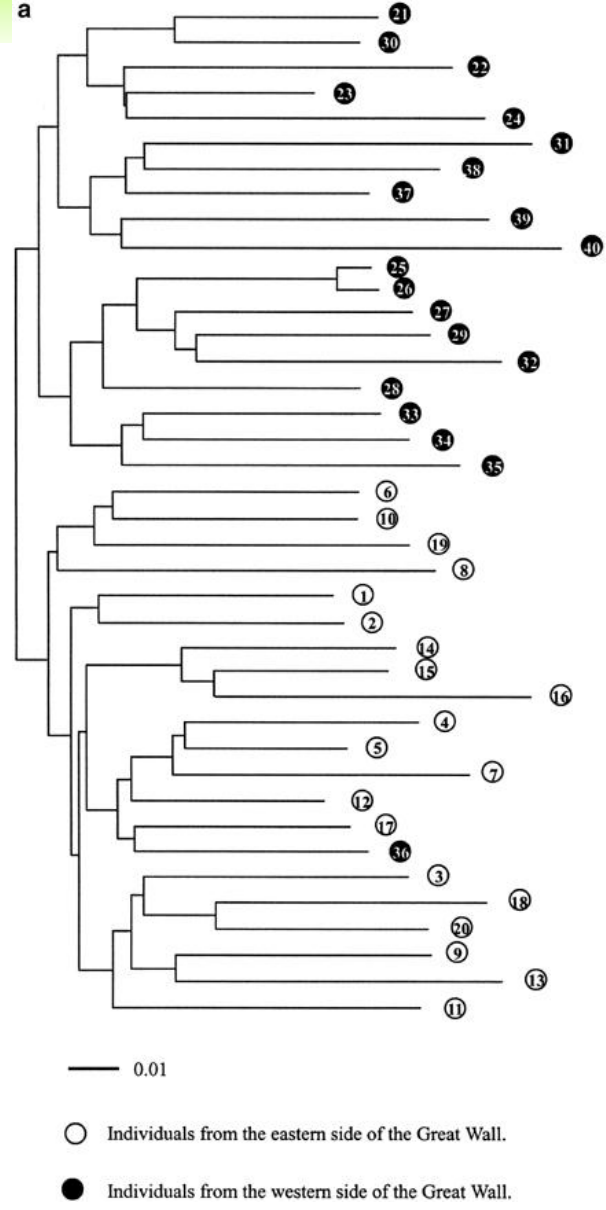


Solving Crimes



# The Great Wall of China: a physical barrier to gene flow?

Therefore, it is reasonable to deduce that the Juyong-guan Great Wall has served as a physical barrier to gene flow between subpopulations separated for more than 600 years.



(a) Neighbour-joining dendrogram of 40 *P. armeniaca* individuals from the Great Wall site. (b) Neighbour-joining dendrogram of 40 *P. armeniaca* individuals from the control site.



# 20世纪90年代美国的一件诉讼案

Florida的一位牙医80年代末得了艾滋病，他的一位病人90年代初被诊断为HIV阳性

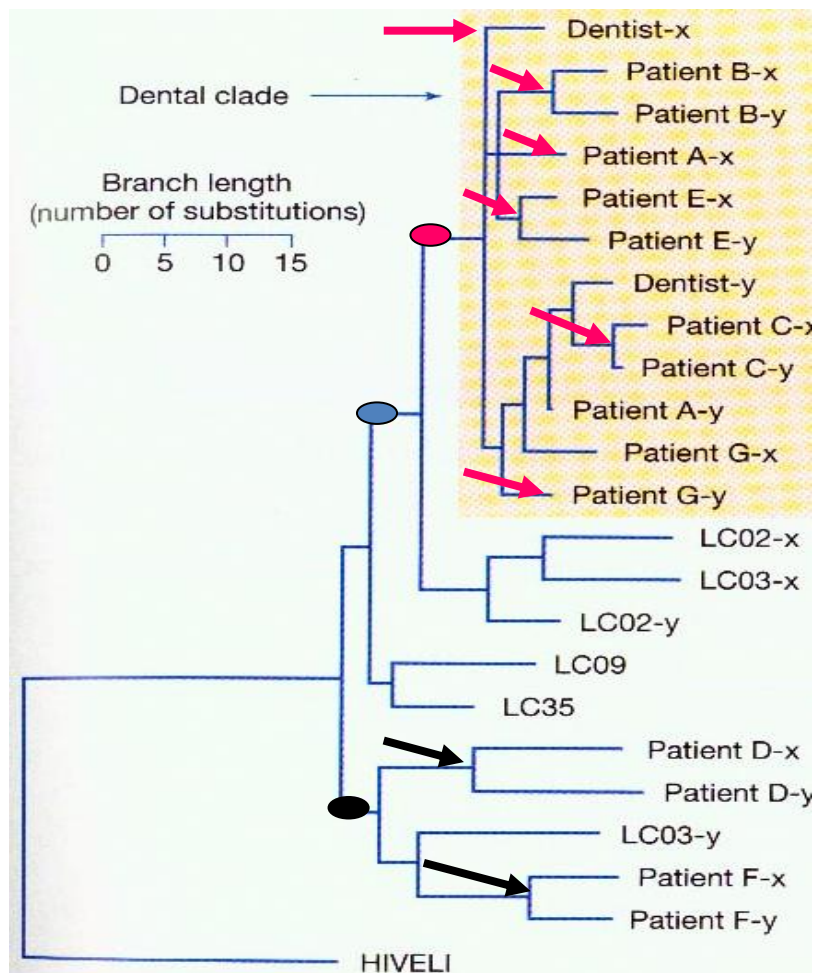
他的病人均去检查，又发现了几个HIV阳性，便起诉牙医

美国CDC用牙医及其病人体内不同时间分离的HIV基因序列建树

结论：

- 有两位病人从其他来源感染了HIV
- 有5位病人从牙医那感染了HIV（同时考虑了其他因素）

该事件促使医院加强保护措施





# The most important step: picking up sequences

- Be aware:
- 1.All the sequences are extracted from reliable source and correct
- 2.All the sequences are **orthologous**
- 3.Choose the right region
- 4.Random evolution
- 5.Each position evolve independently



*Thank  
You!*



*Sun Tianshu*

