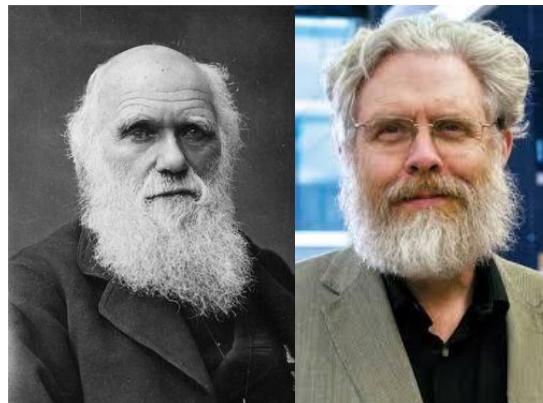


# 生物统计学： 生物信息中的概率统计模型

2019年秋



# 有关信息

- 授课教师: 宁康
  - Email: ningkang@hust.edu.cn
  - Office: 华中科技大学东十一楼504室
  - Phone: 87793041, 18627968927
- 课程网页
  - <http://www.microbioinformatics.org/teach/#>
  - QQ群: 882140516



2019生物统计学

扫一扫二维码，加入该群。

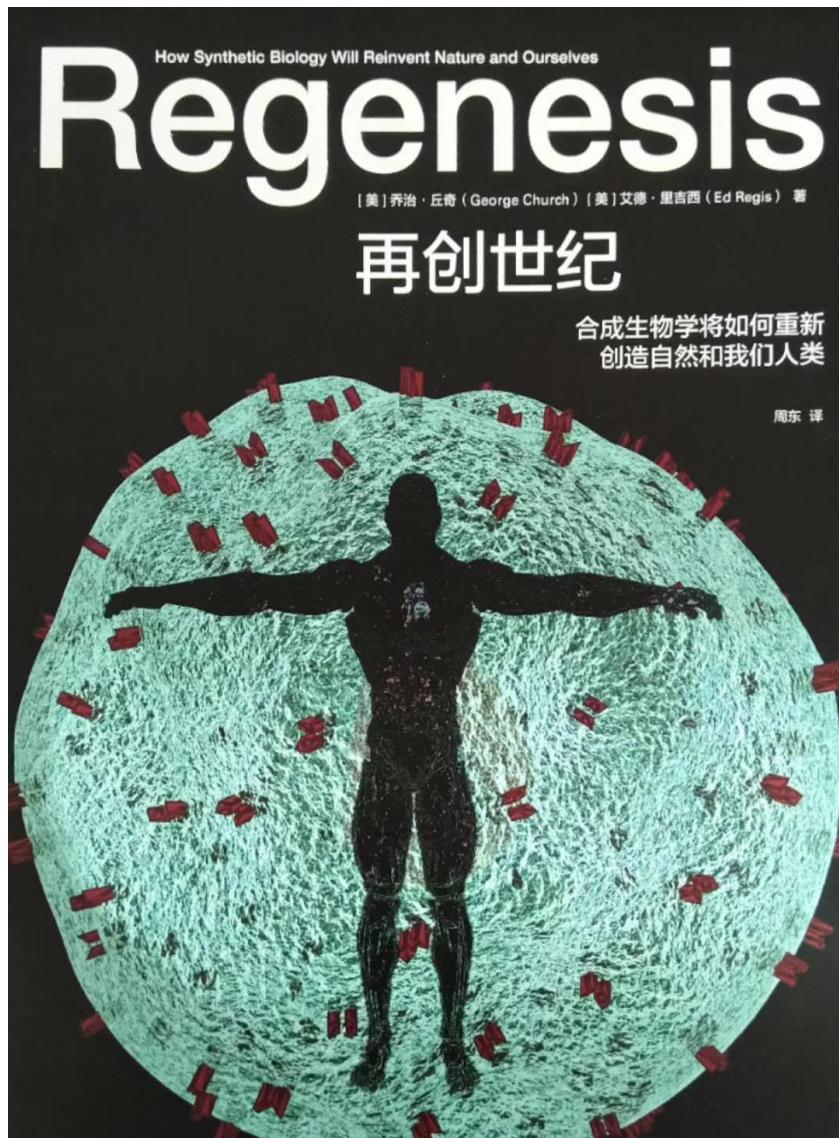
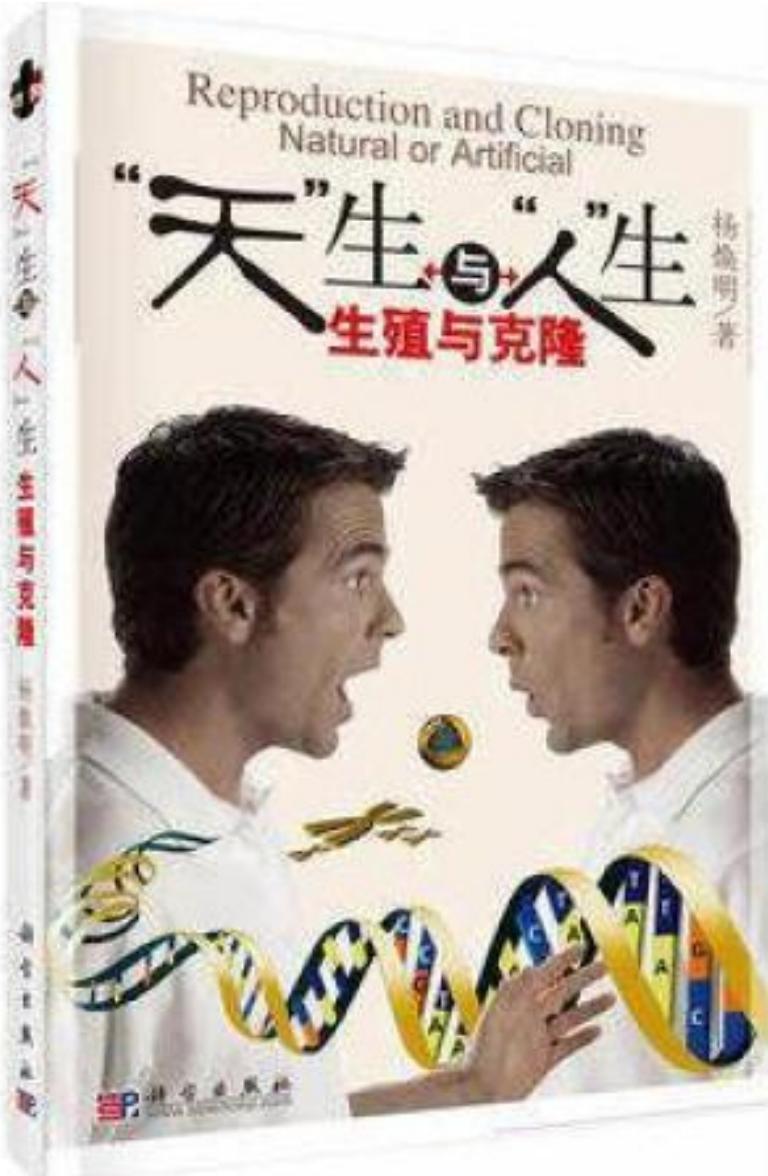
# 课程安排

- 生物背景和课程简介
- 传统生物统计学及其应用
- 生物统计学和生物大数据挖掘
  - Hidden Markov Model (HMM)及其应用
    - Markov Chain
    - HMM理论
    - HMM和基因识别 (Topic I)
    - HMM和序列比对 (Topic II)
  - 进化树的概率模型 (Topic III )
  - Motif finding中的概率模型 (Topic IV)
    - EM algorithm
    - Markov Chain Monte Carlo (MCMC)
  - 基因表达数据分析 (Topic V)
    - 聚类分析-Mixture model
    - Classification-Lasso Based variable selection
  - 基因网络推断 (Topic VI)
    - Bayesian网络
    - Gaussian Graphical Model
  - 基因网络分析 (Topic VII)
    - Network clustering
    - Network Motif
    - Markov random field (MRF)
  - Dimension reduction及其应用 (Topic VIII)
- 面向生物大数据挖掘的深度学习

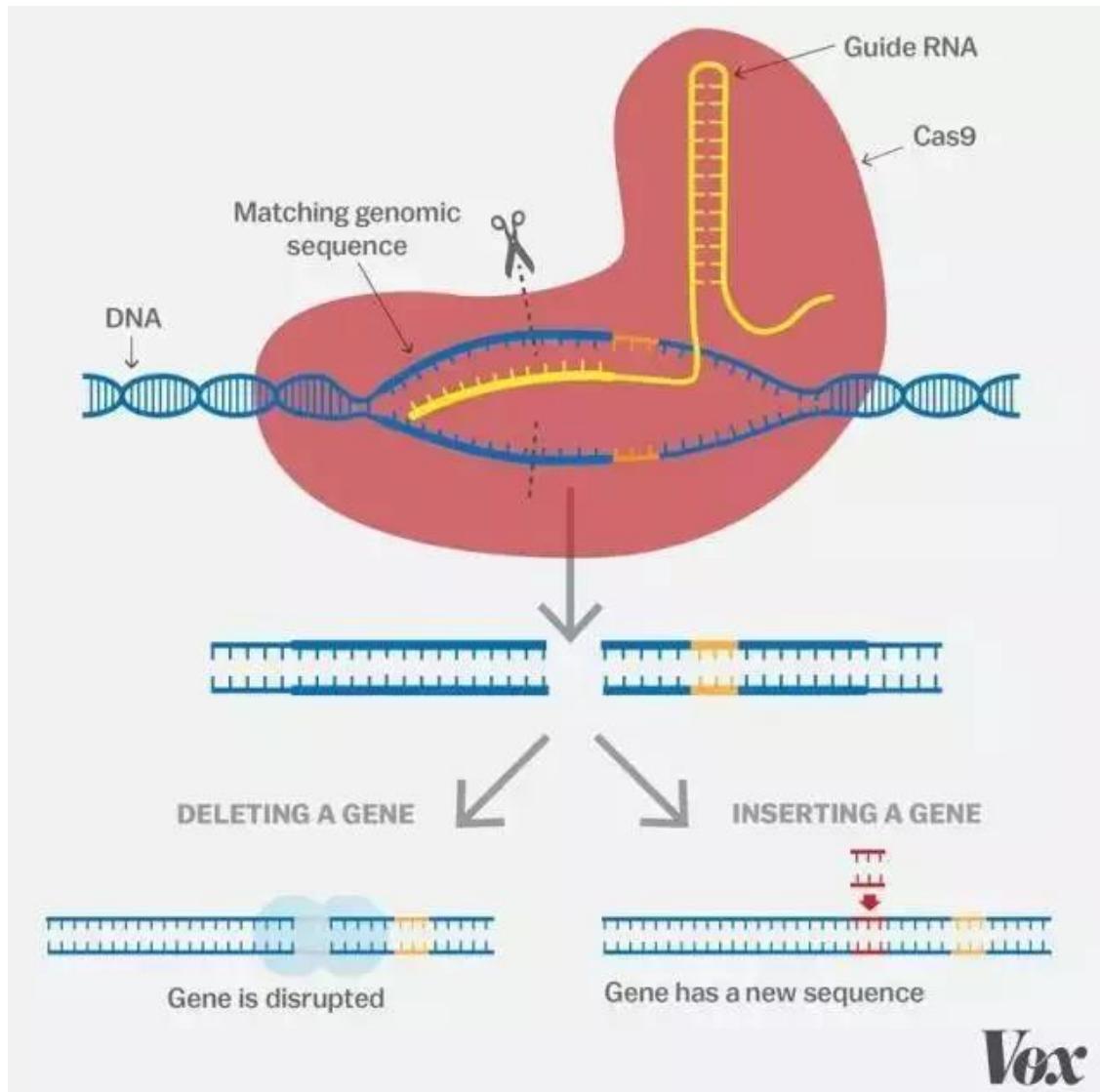
研究对象：  
生物序列，  
进化树，  
生物网络，  
基因表达  
...

方法：  
生物计算与生物统计

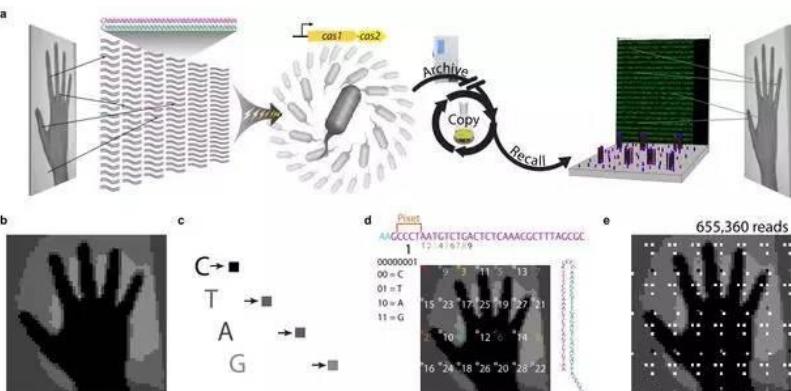
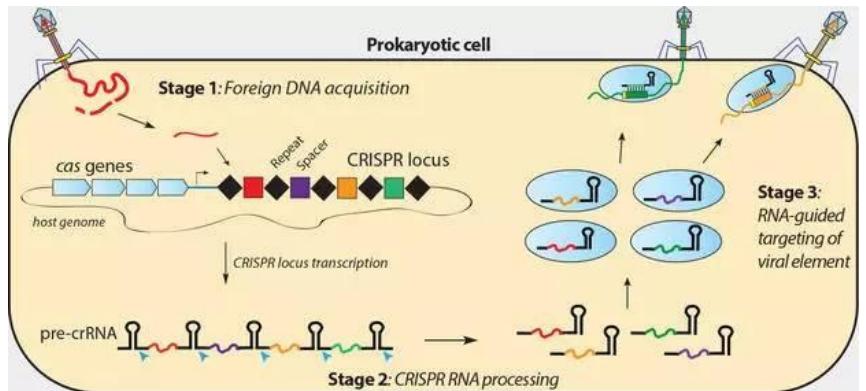
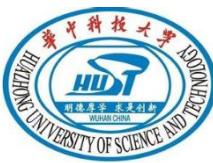
# 读物推荐



# CRISPR和基因编辑技术



# Understand it, create it!



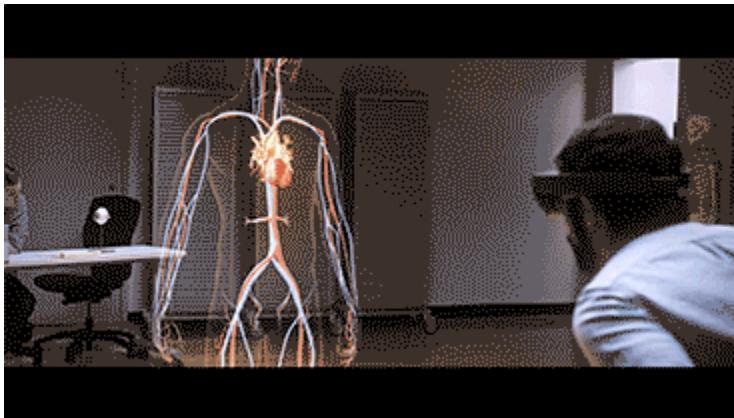
Original Image



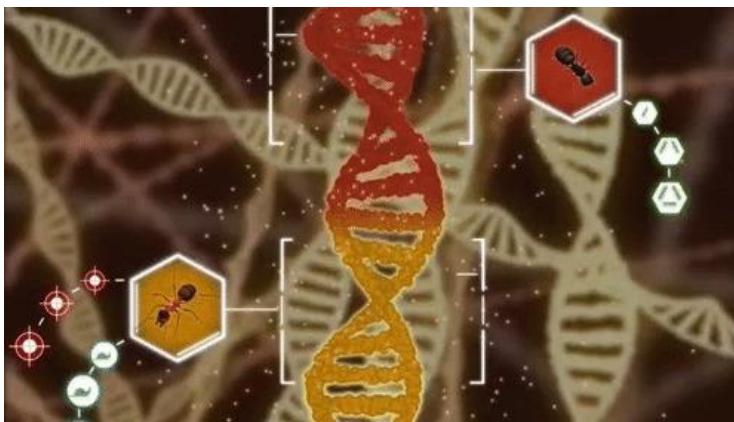
Image Reconstructed From Bacteria

原始图像

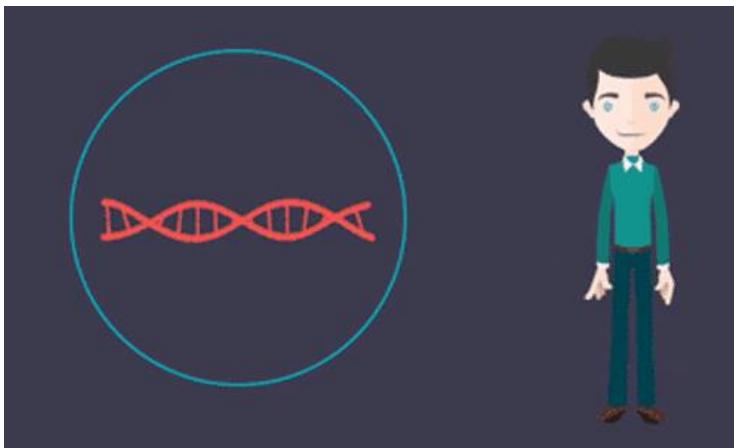
从细菌DNA还原的图像



See it!



Understand it!



Create it!

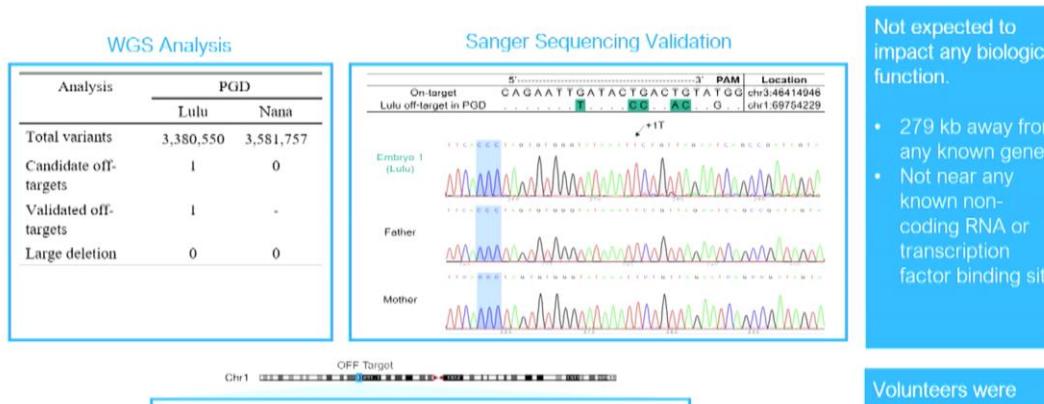
# CRISPR和基因编辑技术

要清楚什么可以做，什么不可以做！

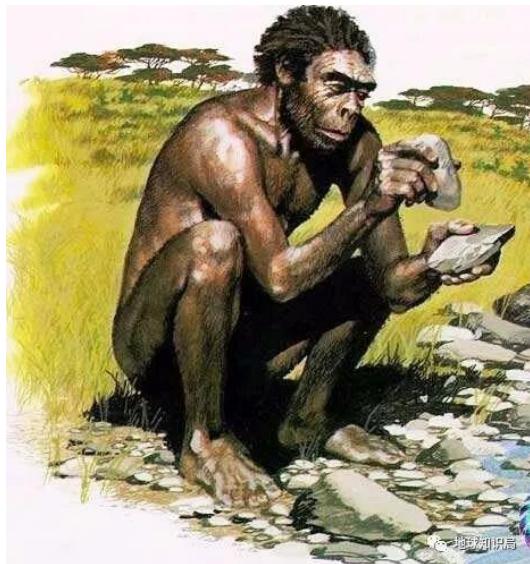
The slide title is "Embryo setting presents fewer editing events, but consequences pose very serious, whole-body risks". It includes a diagram showing that embryo editing targets 1 to 4 cells, while adults have 37 trillion cells. A table lists various cell types and their counts, with a note that 70% of cells are impacted by off-targets at 0.0005% frequency. A video player interface shows a speaker at the podium.

In humans, after advice in Feb last year; found could achieve with indications of safety within 14 days

PGD identified one potential intergenic off-target



# CRISPR和基因编辑技术



OR



# CRISPR和基因编辑技术

Are there compelling medical indications?

Disease prevention

- Huntington's
- Tay Sach's
- Cystic Fibrosis
- Sickle cell anemia

Consider alternatives...

- IVF, genetic diagnosis
- Somatic therapy

When no alternative...

- Couples, both affected
- Infertility

Modifying Disease Risk

- HIV resistance (CCR5)
- Heart disease (PCSK9)
- Alzheimer's (APP A673T/+)
- Cancer (BRCA1/2)
- Resistance to global pandemics...

"Enhancements"

- Muscularity (MSTN)
- Height, skin color
- Learning and memory  
<https://www.dnalc.org/view/1390-Genes-for-Learning-and-Memory.html>

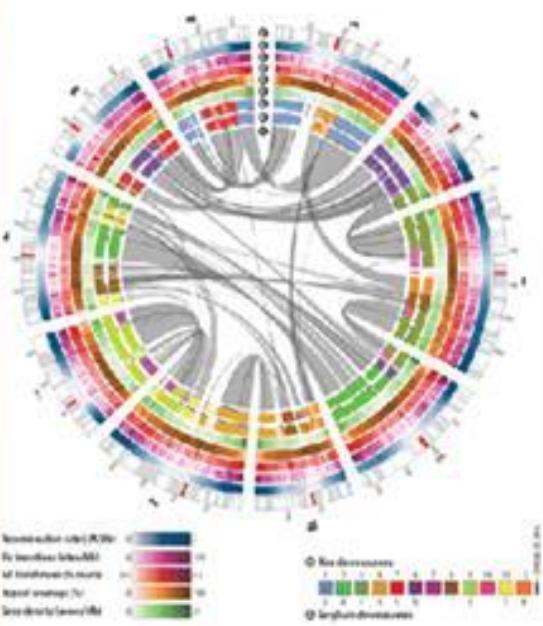
Permissible vs impermissible applications?

# 第4章：进化树构建的概率方法

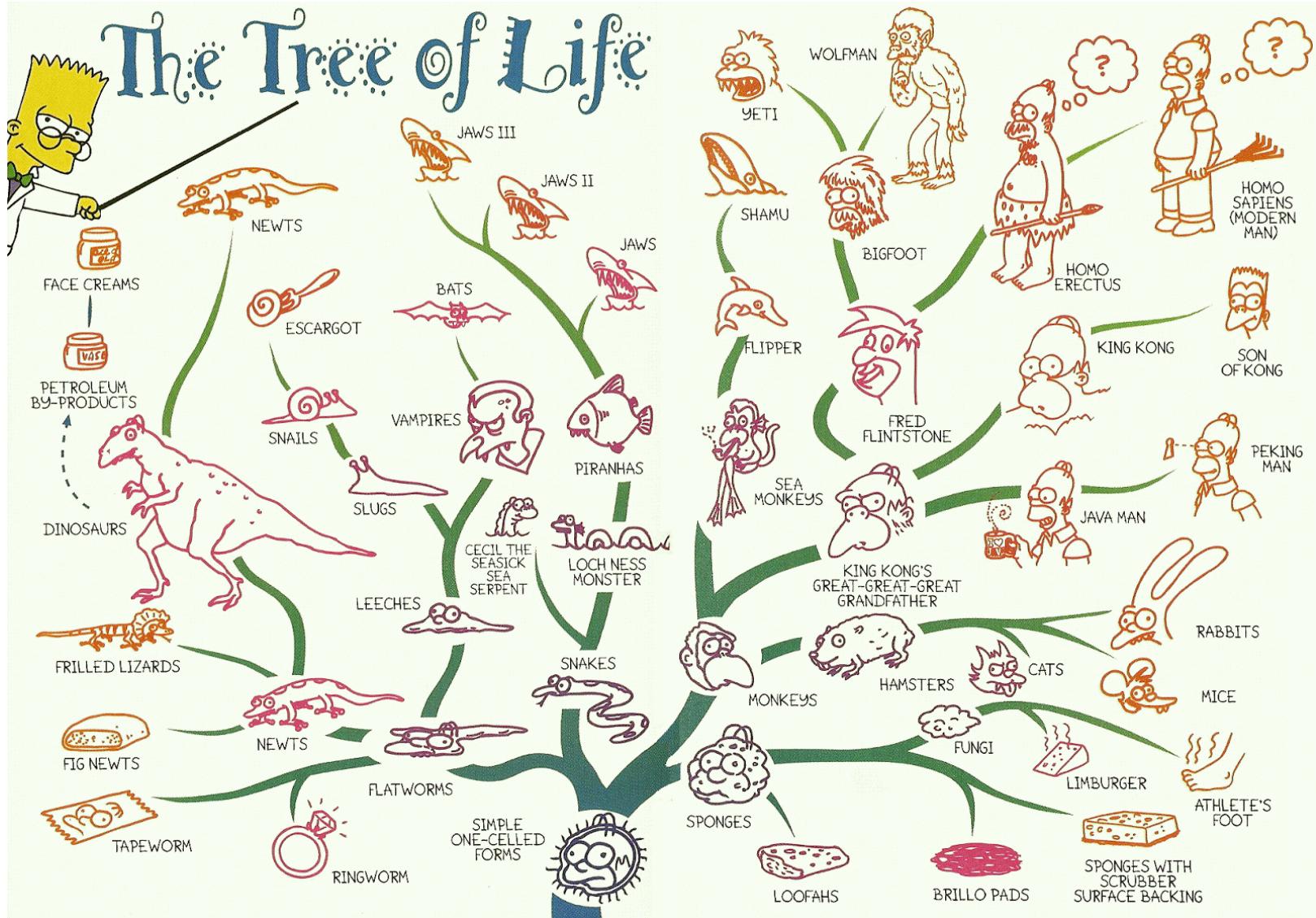
- 问题介绍
- 进化树构建方法的概率方法

# Phylogenetic Tree

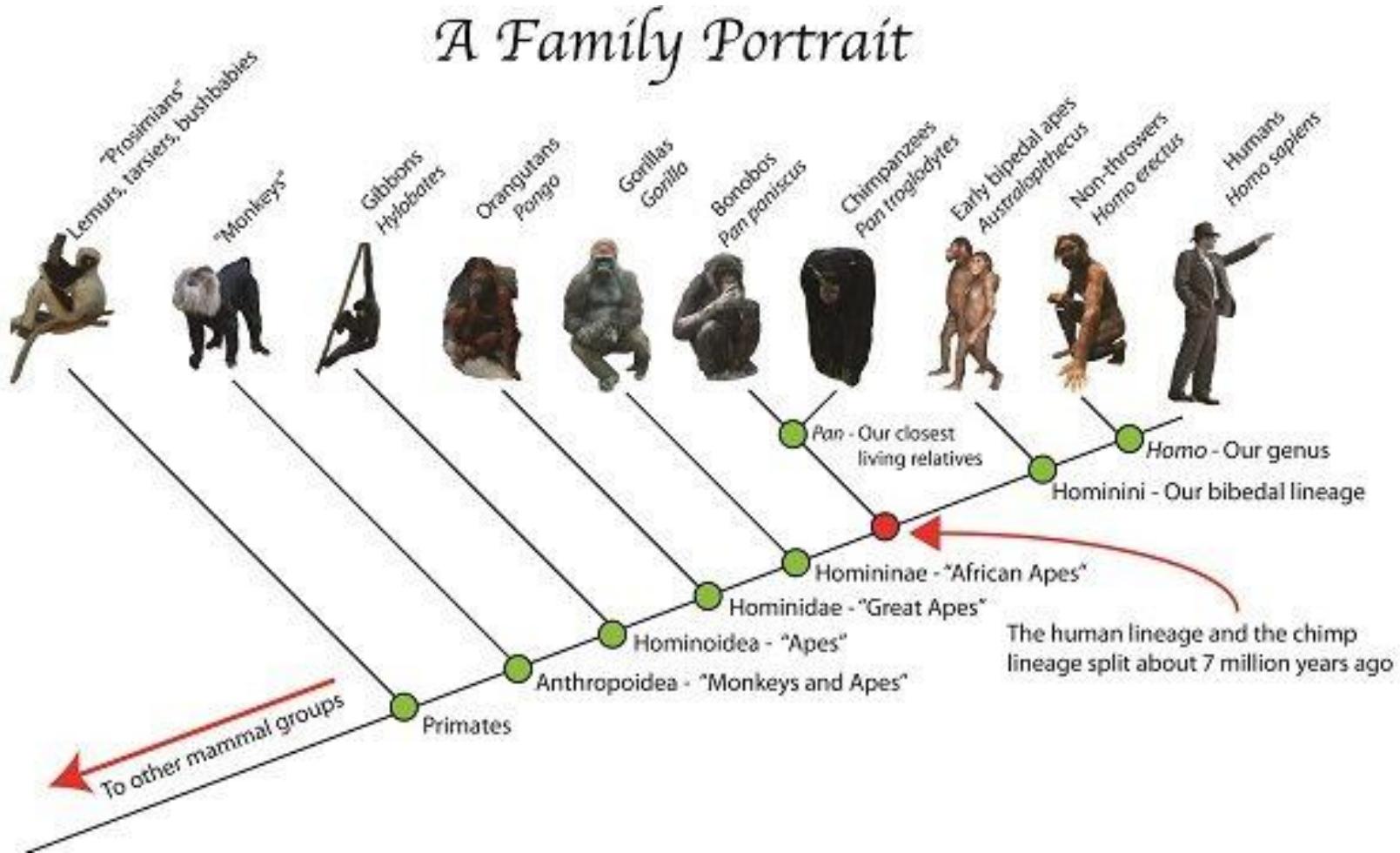
Everyone what to know how this tree  
look like...



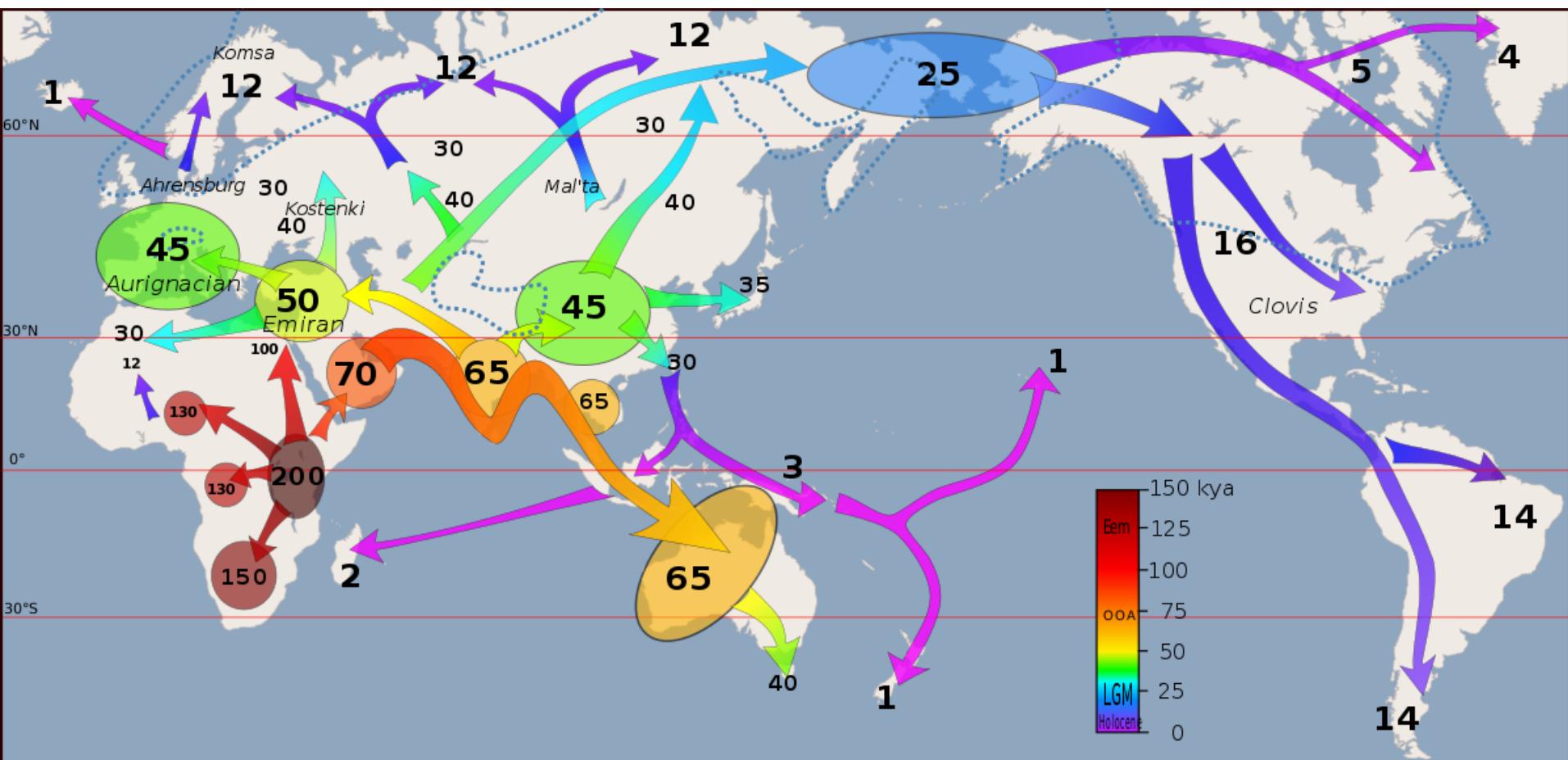
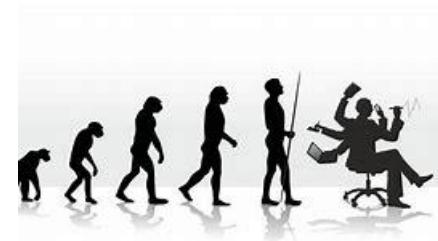
# Phylogenetic Tree



# Evolution everywhere



# Evolution everywhere



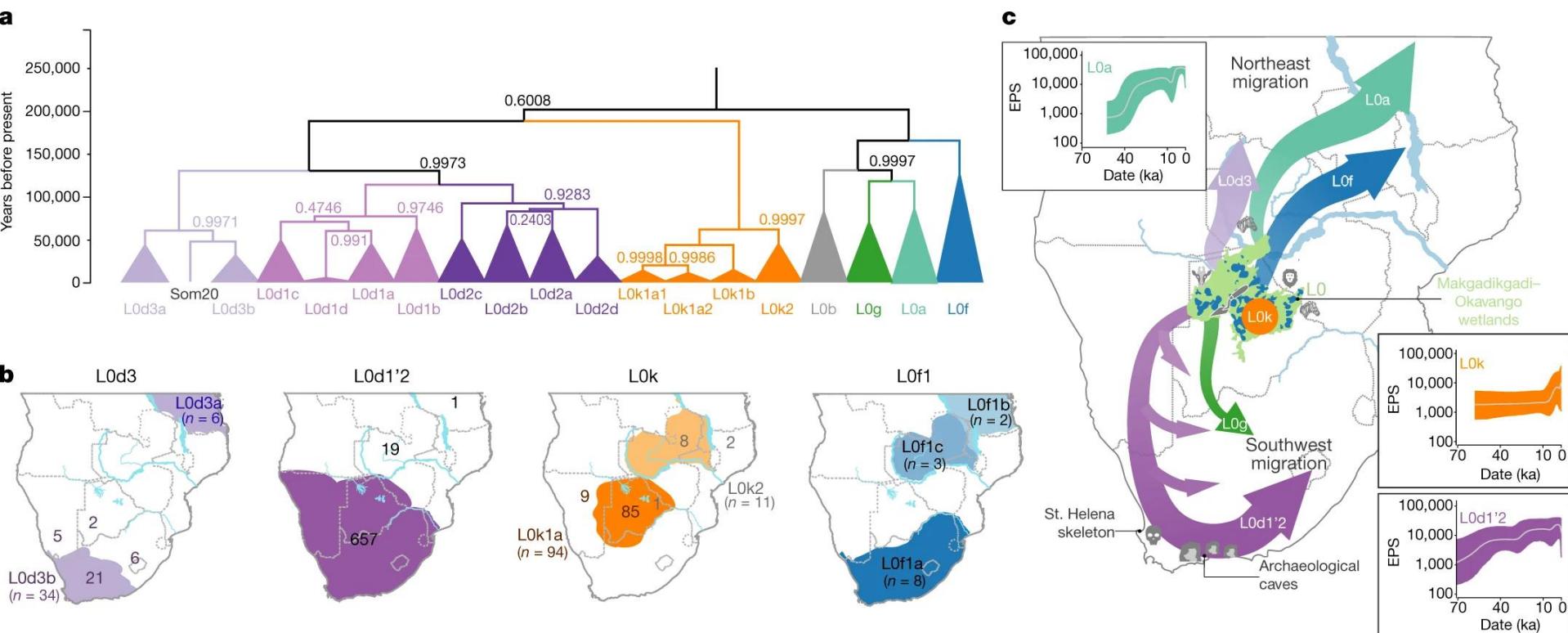
# Evolution everywhere

ISSUE 3185 | MAGAZINE COVER DATE: 7 July 2018



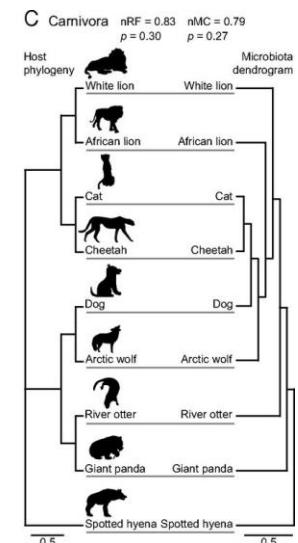
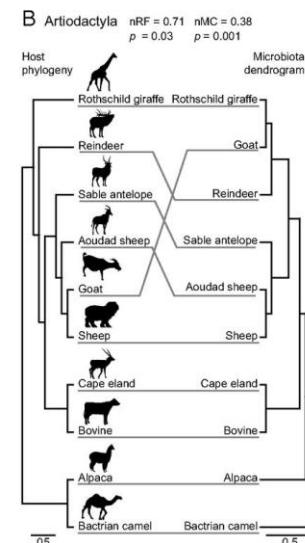
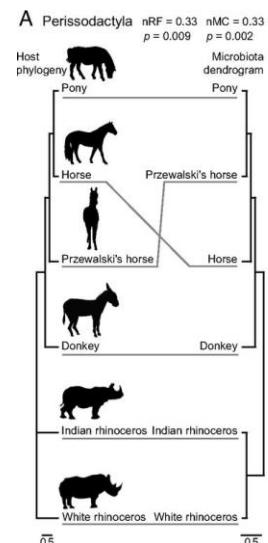
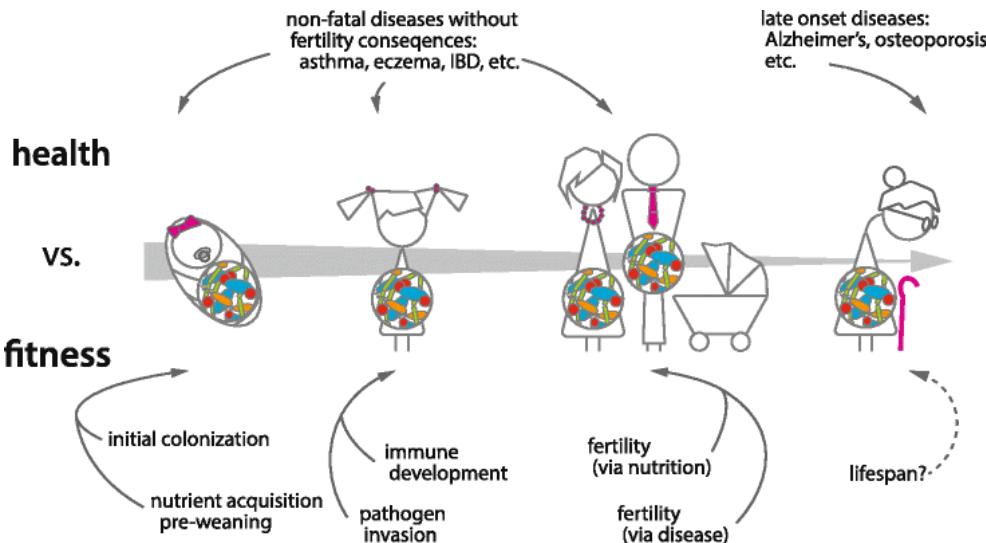
# Evolution everywhere

# Out of “homeland”?

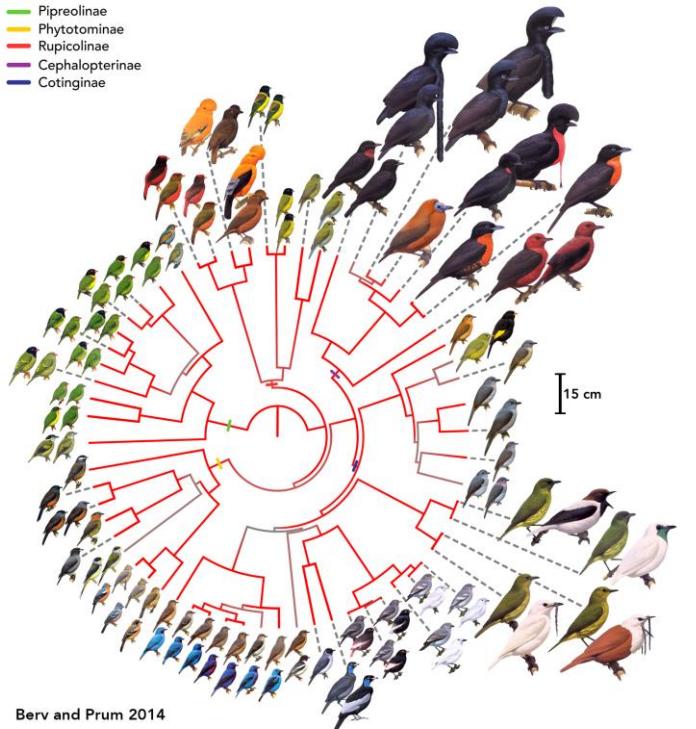


# Evolution everywhere

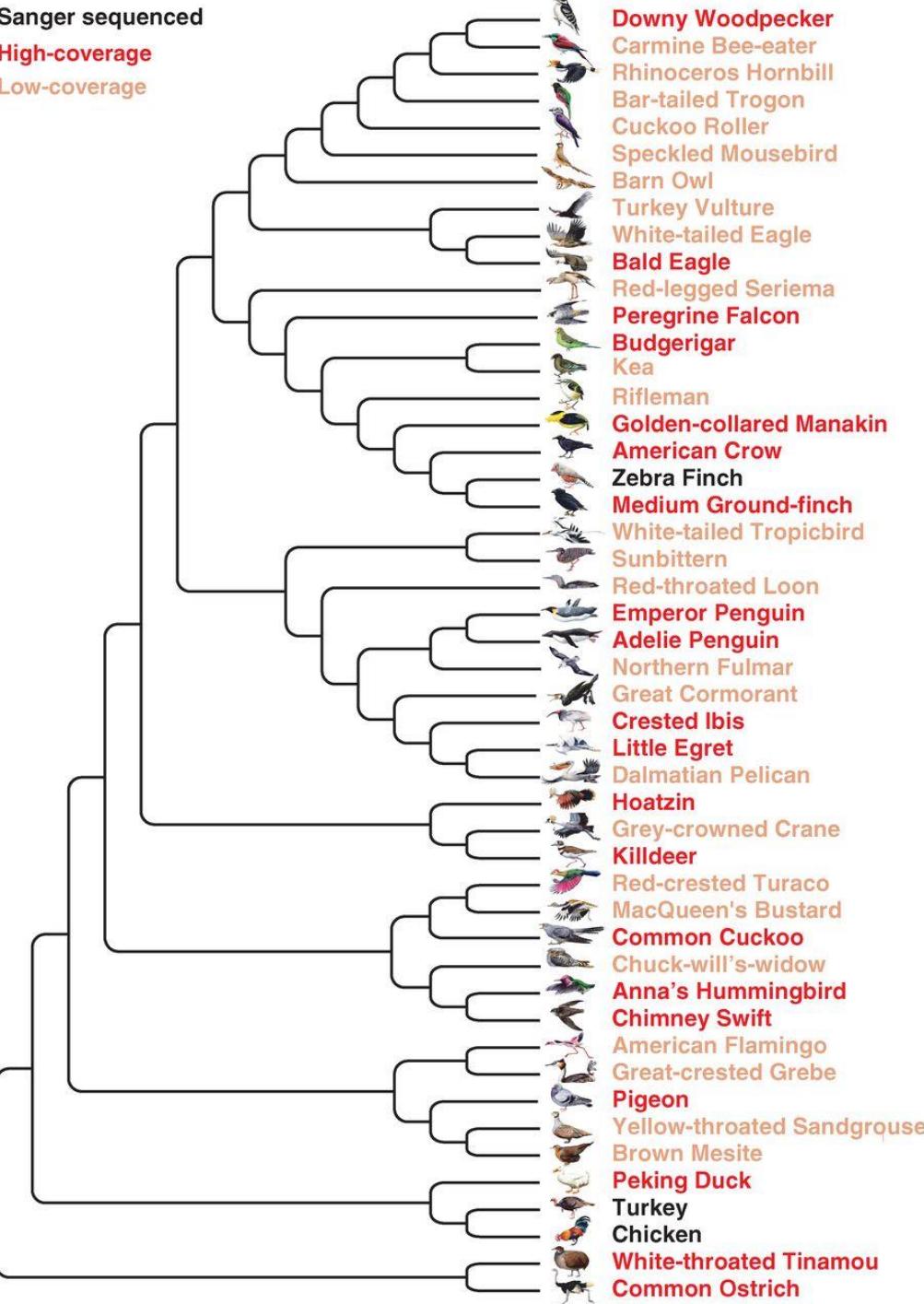
Microbiome evolution also exists



# Phylogenetic Tree



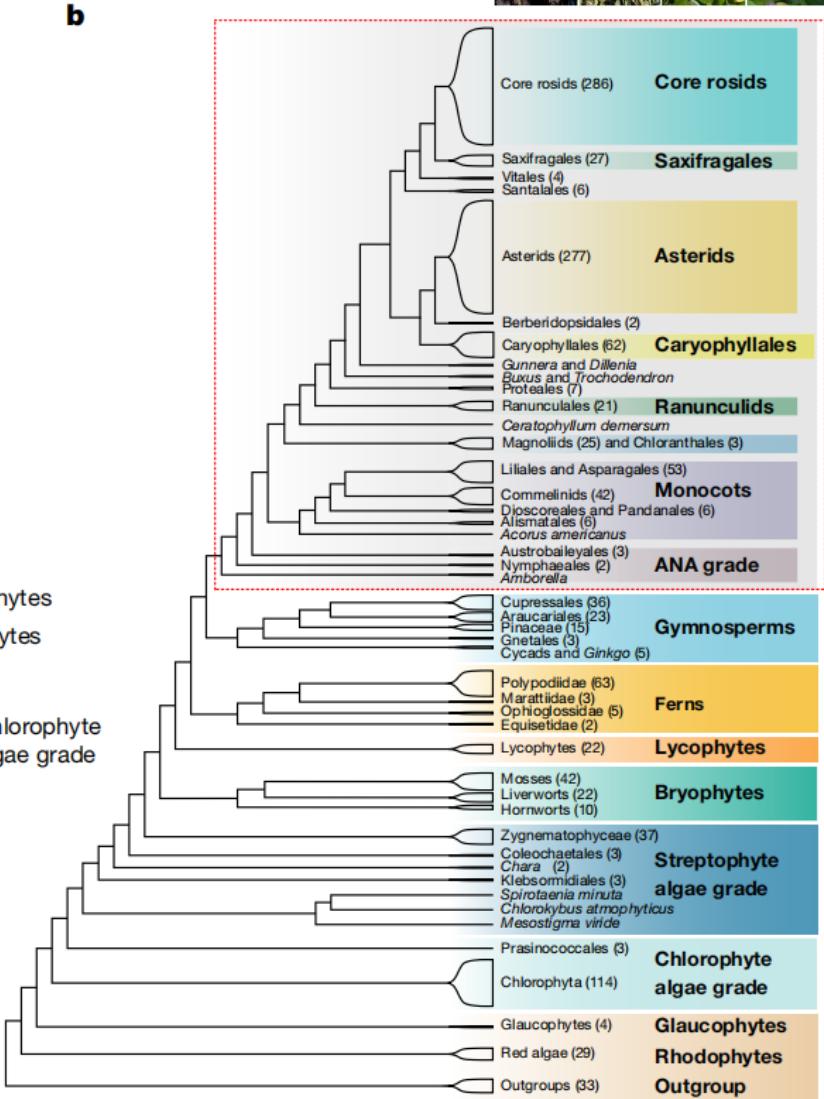
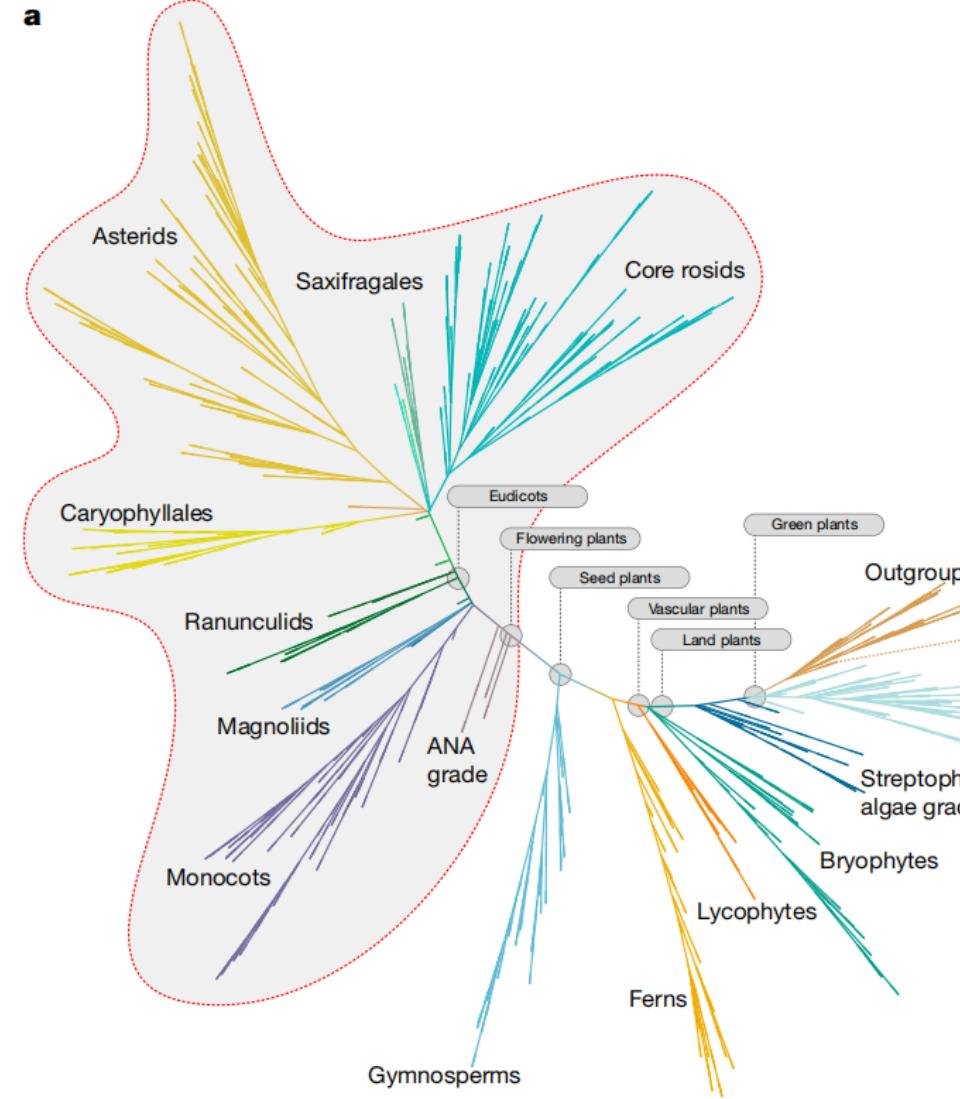
Reference:  
<https://b10k.genomics.cn/>



# Phylogenetic Tree

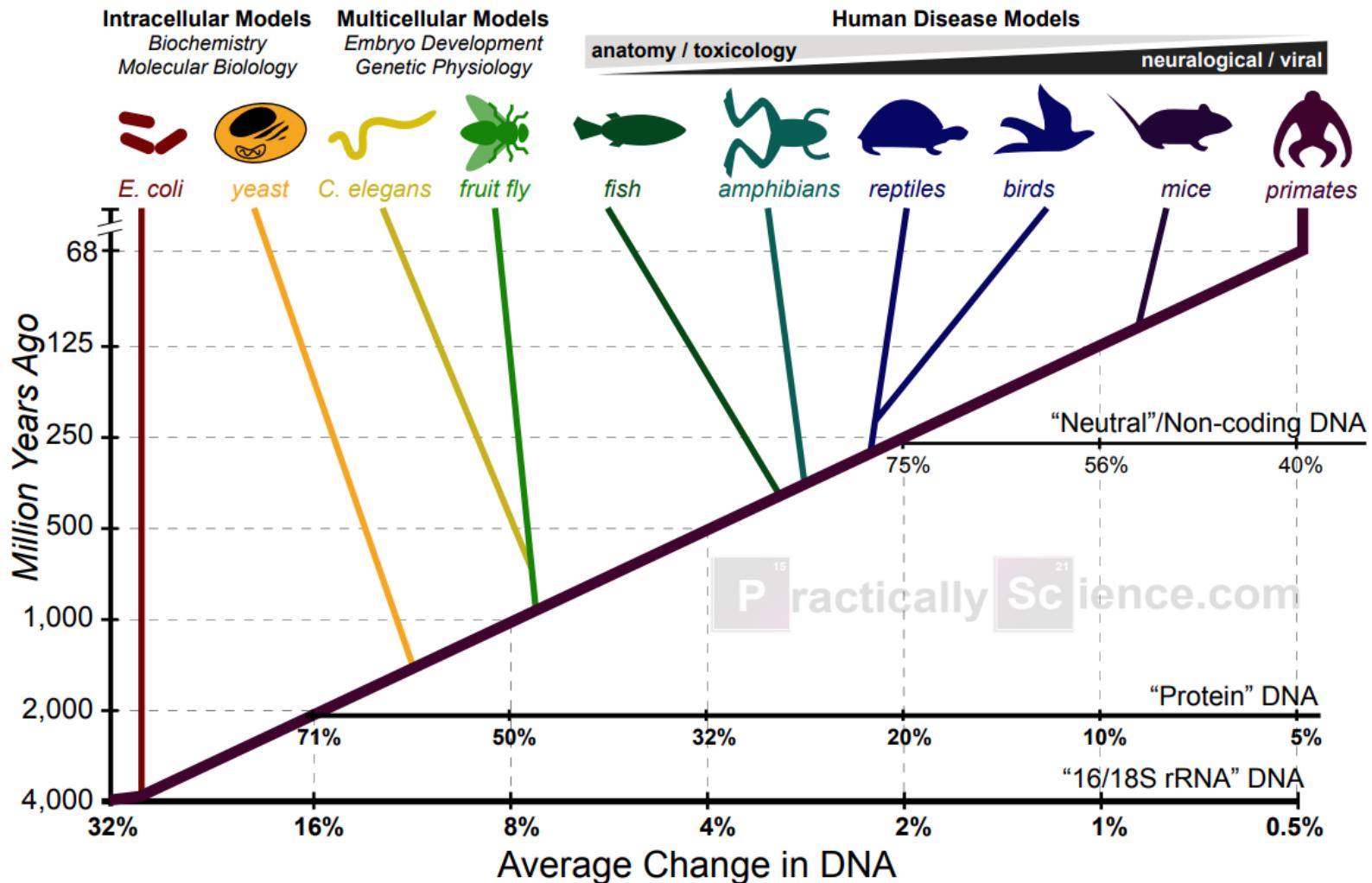
Reference:

One thousand plant transcriptomes and the phylogenomics of green plants,  
*Nature*, 2019

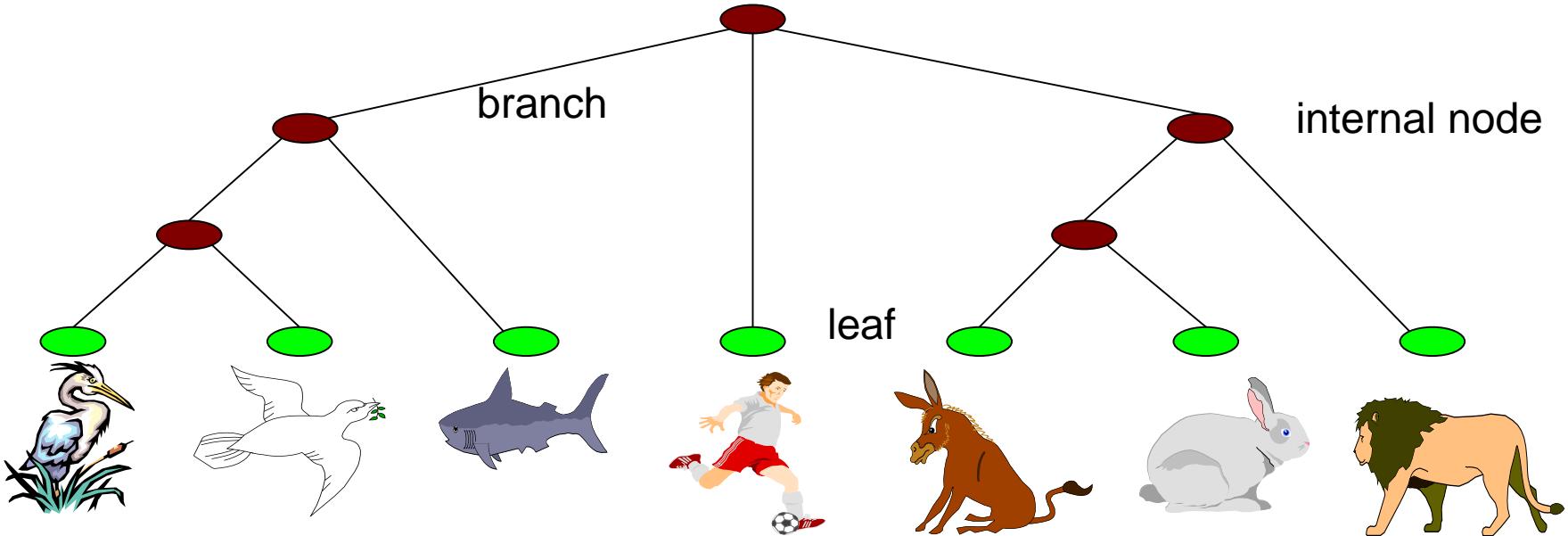


# Phylogenetic Tree

## Evolution of Model Organisms and the DNA Molecular Clock



# Phylogenetic Tree



- Topology: bifurcating
  - Leaves -  $1 \dots N$
  - Internal nodes  $N+1 \dots 2N-2$
- Branch length

Reference:  
<https://itol.embl.de/>

# 构建进化树算法

- 基于距离的构建方法：
  - UPGMA (Unweighted pair group method with arithmetic mean, 平均连接聚类法)
  - ME (Minimum Evolution, 最小进化法)
  - NJ (Neighbor-Joining, 邻接法)
- 基于特征的构建方法：
  - 最大简约法 (MP法)
  - 最大似然法 (ML法)
  - 进化简约法 (EP法)
  - 相容性方法

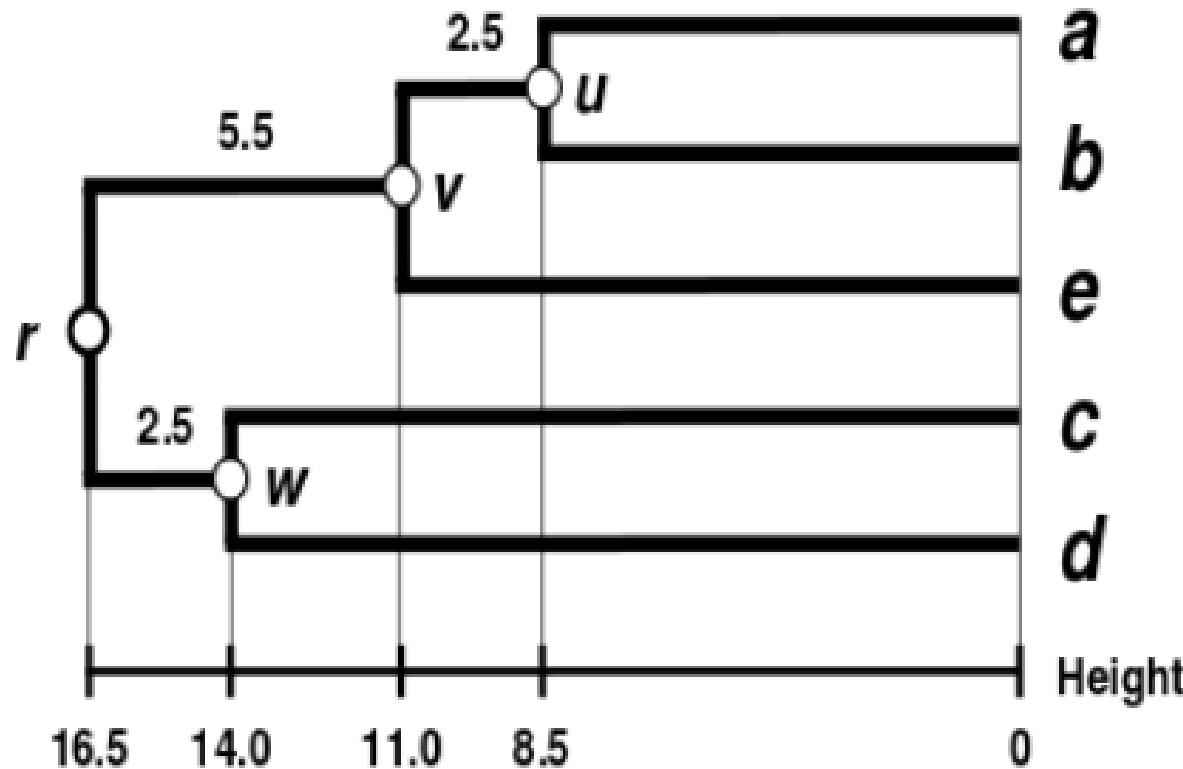
# 基于距离的构建方法

- UPGMA

- ①以已求得的距离系数,所有比较的分类单元的成对距离构成一个 $t \times t$ 方阵,即建立一个距离矩阵M。
- ②对于一个给定的距离矩阵,寻求最小距离值 $D_{pq}$ 。
- ③定义类群p和q之间的分支深度 $L_{pq}=D_{pq}/2$ 。
- ④若p和q是最后一个类群,侧聚类过程完成,否则合并p和q成一个新类群r。
- ⑤定义并计算新类群r到其他各分类群i( $i \neq p$ 和 $q$ )的距离 $D_{ri}=(D_{pi}+D_{qi})/2$ 。
- ⑥回到第一步,在矩阵中消除p和q,加入新类群r,矩阵减少一阶,重复进行直至达到最后归群。

# 基于距离的构建方法

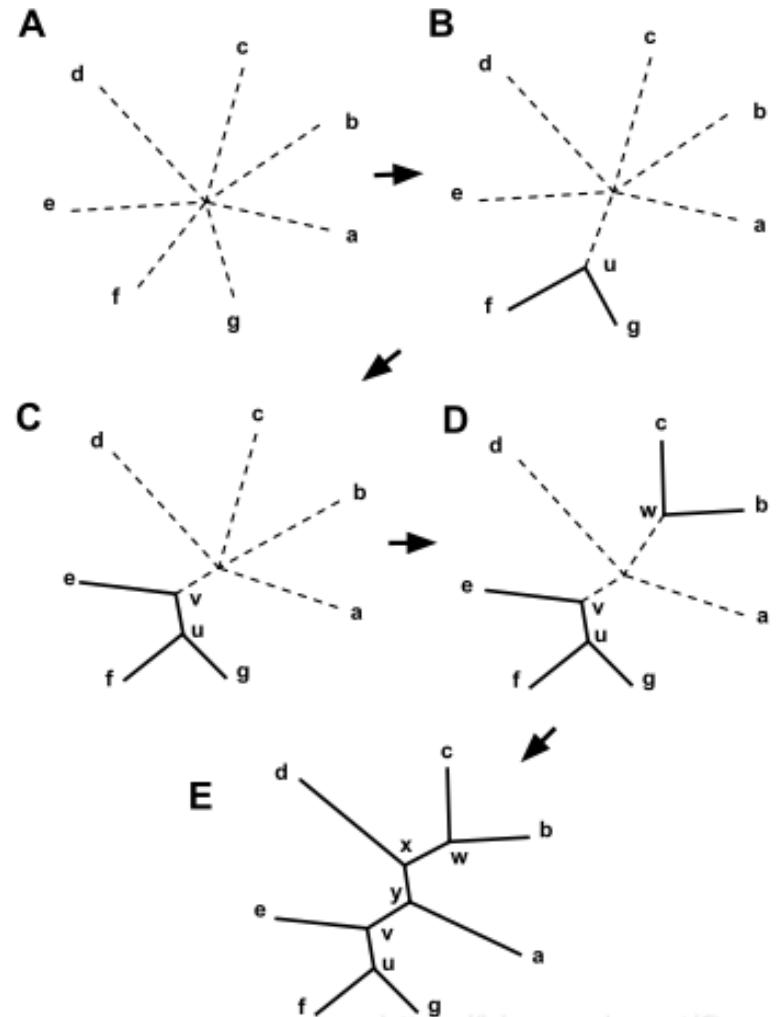
- UPGMA



# 基于距离的构建方法

## • Neighbor-Joining

- ① 对于给定距离矩阵中的每一端结*i*,用下式计算与其它分类单元之间的净趋异量( $R_i$ ) (*t*:矩阵中的分类单元数)  
$$R_i = \frac{D_{it} - D_{\text{min}}}{t - 2}$$
- ② 建立一个速率校正距离矩阵M,其元素由下式确定:  
$$M_{ij} = \frac{D_{ij}}{R_i + R_j}$$
- ③ 定义一个新节点,的三个分支分别与节点*i,j*和树的其余部分相连,并且Dij为矩阵中距离最小者,u到节点*i*和*j*的分支长度定义为  
$$d_{ui} = d_{uj} = \frac{D_{ij} - D_{\text{min}}}{2(R_i + R_j)}$$
- ④ 定义到树的其它节点*k*(*k*≠*i*和*j*外的所有节点)的距离:  
$$d_{uk} = D_{ik} + D_{jk} - D_{ij}$$
- ⑤ 从距离矩阵中删除*i*和*j*的距离,矩阵减少一阶。
- ⑥ 如果矩阵仍然多于两个的节点,重复第①-⑤步,否则删除最外两个节点的分支长度来确定外,树上其余节点都确定,最后是剩余的2个的分支长度 $Sy=Dij$



# 构建进化树算法

- 基于距离的构建方法：
  - UPGMA (Unweighted pair group method with arithmetic mean, 平均连接聚类法)
  - ME (Minimum Evolution, 最小进化法)
  - NJ (Neighbor-Joining, 邻接法)
- 基于特征的构建方法：
  - 最大简约法 (MP法)
  - 最大似然法 (ML法)
  - 进化简约法 (EP法)
  - 相容性方法

# 最大简约法 (Maximum Parsimony)

最大简约法的理论基础是奥卡姆 (Ockham) 哲学原则，这个原则认为：解释一个过程的最好理论是所需假设数目最少的那一个。

方法：

计算所有可能的拓扑结构；

计算出所需替代数最小的那个拓扑结构，作为最优树。

Occam's Razor

The simplest explanation is  
usually the correct one.

# occam's razor method

## The Role of Ockham's Razor: To Screen for Plurality at the Beginning of the Scientific Method

1. OBSERVATION
2. NECESSITY
3. INTELLIGENCE/AGGREGATION OF DATA (The Three Key Questions)
4. CONSTRUCT FORMULATION
5. SPONSORSHIP/PEER INPUT (Ockham's Razor)



False Skeptics seek to block this step at all costs.

6. HYPOTHESIS DEVELOPMENT
7. PREDICTIVE TESTING
8. COMPETITIVE HYPOTHESES FRAMING (ASKING THE RIGHT QUESTION)

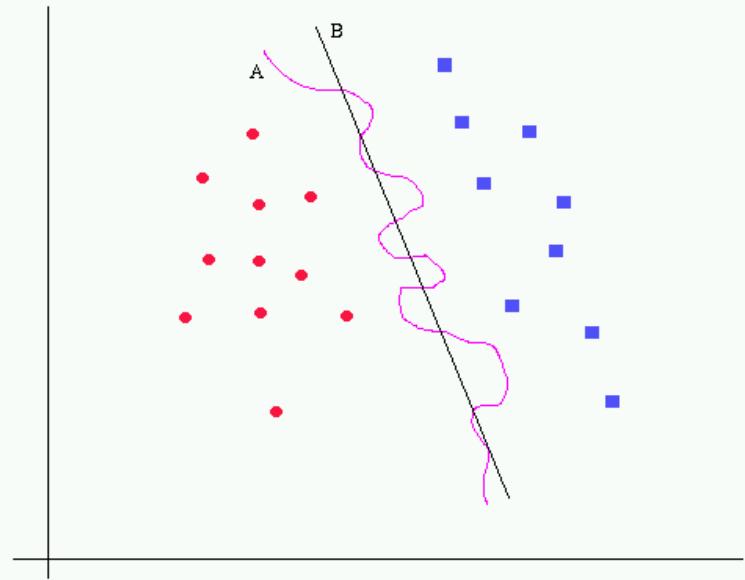
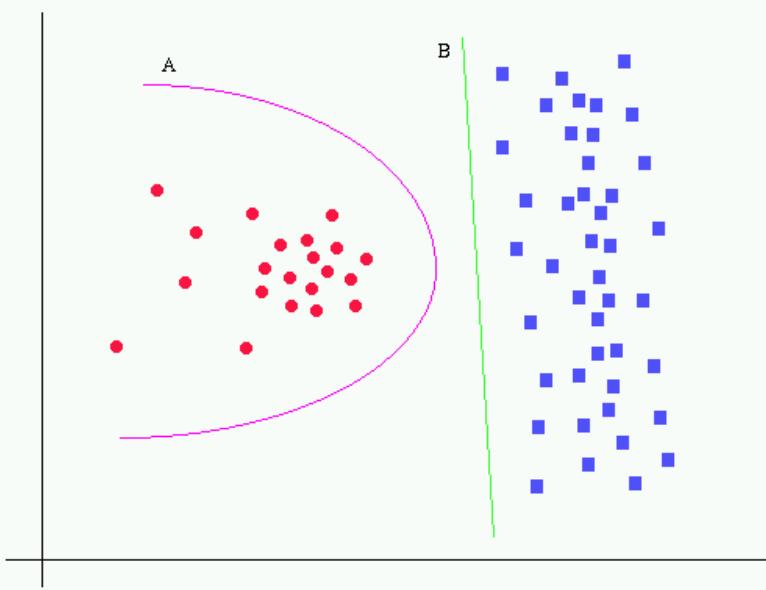
9. FALSIFICATION TESTING
10. HYPOTHESIS MODIFICATION
11. FALSIFICATION TESTING/REPEATABILITY
12. THEORY FORMULATION/REFINEMENT
13. PEER REVIEW (Community Vetting)
14. PUBLISH

15. ACCEPTANCE

Plurality

Proof

# occam's razor method



# 最大似然法 (Maximum Likelihood)

ML法对所有可能的系统发育树都计算似然函数，似然函数值最大的那棵树即为最可能的系统发育树。

利用最大似然法来推断一组序列的系统发生树，需首先确定序列进化的模型，如Jukes—Cantor模型、Kimura二参数模型及一般二参数模型等。在进化模型选择合理的情况下，ML法是与进化事实吻合最好的建树算法。其缺点是计算强度非常大，极为耗时。

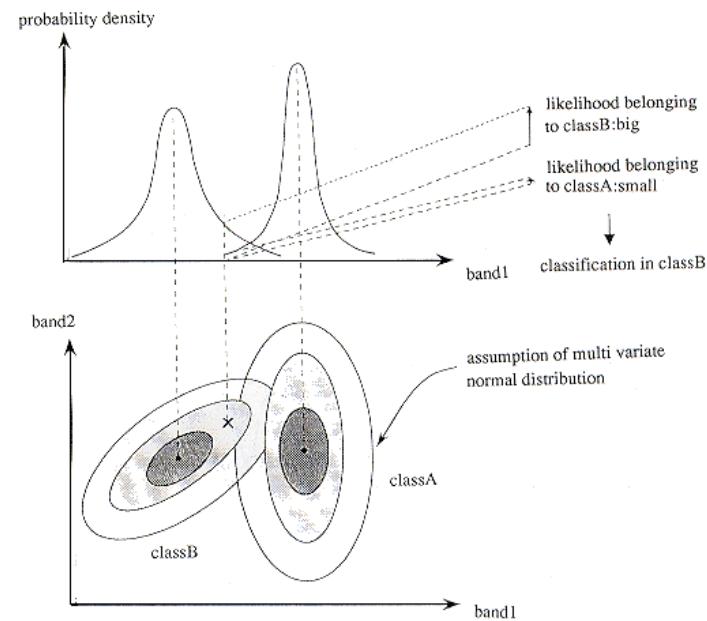


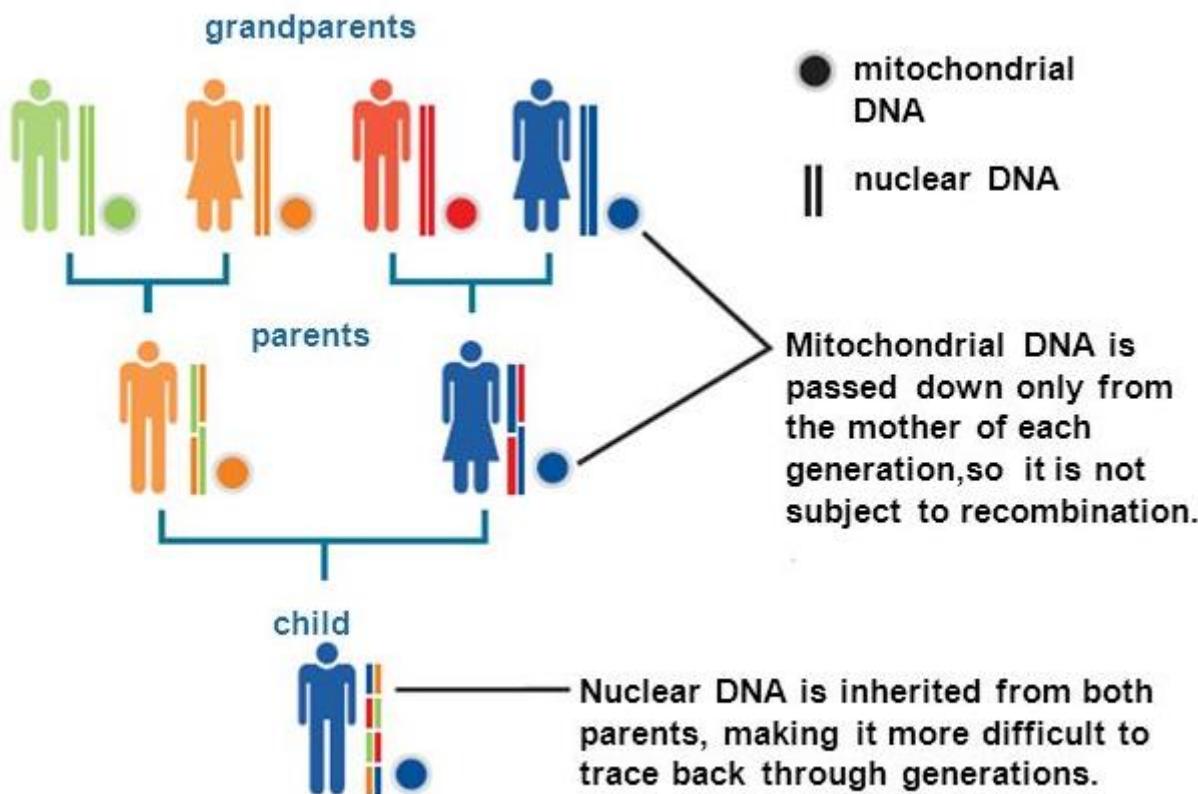
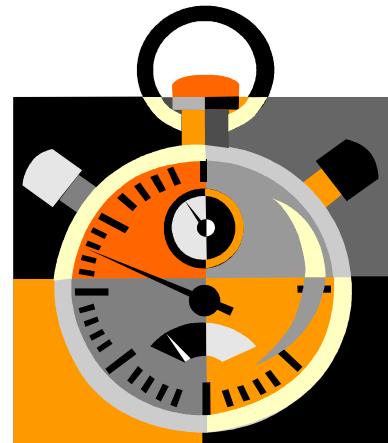
Figure 11.7.1 Concept of Maximum Likelihood Method

# Molecular Clock Hypothesis



- Amount of genetic difference between sequences is a function of time since separation.
- Rate of molecular change is constant (enough) to predict times of divergence

# Molecular Clock Hypothesis



# Likelihood of a Tree

- Given:
  - $n$  aligned sequences  $M = X_1, \dots, X_n$
  - A tree  $T$ , leaves labeled with  $X_1, \dots, X_n$
- Reconstruction  $t^*$ :
  - Labeling of internal nodes
  - Branch lengths

Goal: Find optimal reconstruction  $t^*$  : One maximizing the likelihood  $P(M | T, t^*)$

# Probabilistic Methods

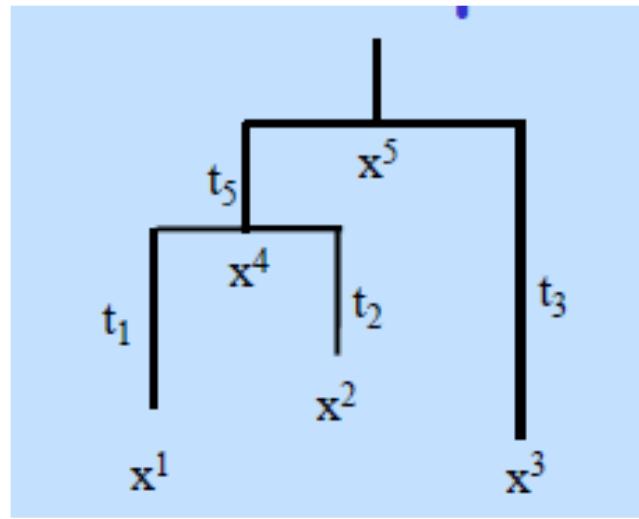
- The phylogenetic tree represents a generative probabilistic model (like HMMs) for the observed sequences.
- Background probabilities:  $q(a)$
- Mutation probabilities:  $P(a|b, t)$
- Models for evolutionary mutations
  - Jukes Cantor
  - Kimura 2-parameter model
- Such models are used to derive the probabilities

# Probabilistic Model

- Assumptions:
  - Each character is independent
  - The branching is a Markov process: The probability that a node  $x$  has a specific label is only a function of the parent node  $y$  and the branch length  $t$  between them
  - The probabilities  $P(x|y,t)$  are known

# Example

- Given the tree



$$\begin{aligned} & P(x_1, x_2, x_3, x_4, x_5 | T, t^*) \\ &= P(x_1 | x_4, t_1) P(x_2 | x_4, t_2) P(x_3 | x_5, t_3) P(x_4 | x_5, t_5) \end{aligned}$$

# Molecular Evolution

- Q: How can we model evolution on nucleotide level? (ignore gaps, focus on substitutions)
- A: Consider what happens at a specific position for small time interval  $\Delta t$
- $P(t)$  = vector of probabilities of {A,C,G,T} at time  $t$
  - $\mu_{AC}$  = rate of transition from A to C per unit time
  - $\mu_A = \mu_{AC} + \mu_{AG} + \mu_{AT}$  rate of transition out of A
  - $p_A(t+\Delta t) = p_A(t) - p_A(t) \mu_A \Delta t + p_C(t) \mu_{CA} \Delta t + \dots$

# Molecular Evolution

In matrix/vector notation, we get

$$P(t + \Delta t) = P(t) + QP(t)\Delta t$$

where  $Q$  is the substitution rate matrix

$$Q = \begin{pmatrix} -\mu_A & \mu_{AG} & \mu_{AC} & \mu_{AT} \\ \mu_{GA} & -\mu_G & \mu_{GC} & \mu_{GT} \\ \mu_{CA} & \mu_{CG} & -\mu_C & \mu_{CT} \\ \mu_{TA} & \mu_{TG} & \mu_{TC} & -\mu_T \end{pmatrix}$$

# Molecular Evolution

- This is a differential equation:

$$P'(t) = Q P(t)$$

- A substitution rate matrix  $Q$  implies a probability distribution over {A,C,G,T} at each position, including stationary (equilibrium) frequencies  $\pi_A, \pi_C, \pi_G, \pi_T$
- Each  $Q$  is an evolutionary model (some work better than others)

# Mutation Probabilities

$P(t)$  satisfy the following two property:

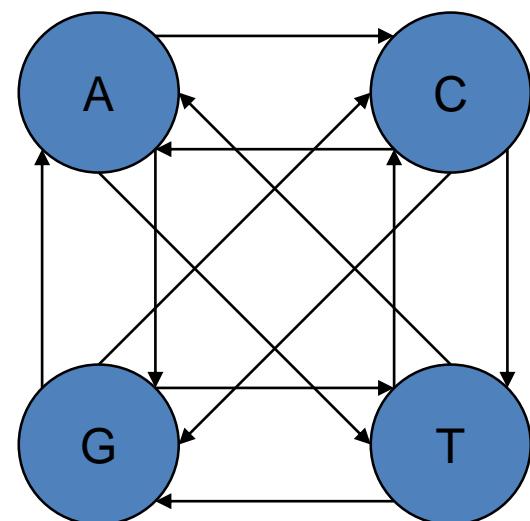
- **Lack of memory:**

$$- P_{a \rightarrow c}(t + t') = \sum_b P_{a \rightarrow b}(t) P_{b \rightarrow c}(t')$$

- **Reversibility:**

- Exist stationary probabilities  $\{P_a\}$  s.t.

$$P_a P_{a \rightarrow b}(t) = P_b P_{b \rightarrow a}(t)$$



# PAM矩阵

- Point accepted mutation (Dayhoff et al 1978)
- Given an tree of protein family, the frequence matrix  $A_{ab}$  counting the occurrence of an “a” in the ancestral sequence was replaced by a “b” in the descendant.
- Estimate the conditional probability  $p(b|a)$

$$P(b|a) = B_{a,b} = \frac{A_{ab}}{\sum_c A_{ac}}$$

# PAM矩阵

- Scaling B

$$C_{ab} = \sigma B_{ab}, C_{aa} = \sigma B_{aa} + (1 - \sigma)$$

- Such that the expected number of substitution is 1%, i.e.

$$\sum_{ab} q_a q_b C_{ab} = 0.01$$

- Then the PAM(1) matrix is given by

$$S(1) = (C_{ab})$$

# Calculating the Likelihood for Ungapped Alignments



$$P(x^1, x^2 | T, t_1, t_2) = \prod_{n=1}^N P(x_n^1, x_n^2 | T, t_1, t_2)$$

$$P(x_n^1, x_n^2 | T, t_1, t_2) = \sum_a q_a P(x_n^1 | a, t_1) P(x_n^2 | a, t_2)$$

Assuming Jukes-Cantor model &  $q_C = q_G = q_A = q_T = \frac{1}{4}$  :

$$P(C, C | T, t_1, t_2) = q_C r_{t_1} r_{t_2} + q_G s_{t_1} s_{t_2} + q_A s_{t_1} s_{t_2} + q_T s_{t_1} s_{t_2} = \frac{1}{4} (r_{t_1} r_{t_2} + 3s_{t_1} s_{t_2})$$

$$P(C, G | T, t) = P(G, C | T, t) = \frac{1}{4} (r_{t_1} s_{t_2} + s_{t_1} r_{t_2} + 2s_{t_1} s_{t_2})$$

$$\Rightarrow P(x^1, x^2 | T, t_1, t_2) = 16^{-(n1+n2)} \left(1 + 3e^{-4\alpha(t_1+t_2)}\right)^{n1} \left(1 - e^{-4\alpha(t_1+t_2)}\right)^{n2}$$

where n1=matches, n2=mismatches

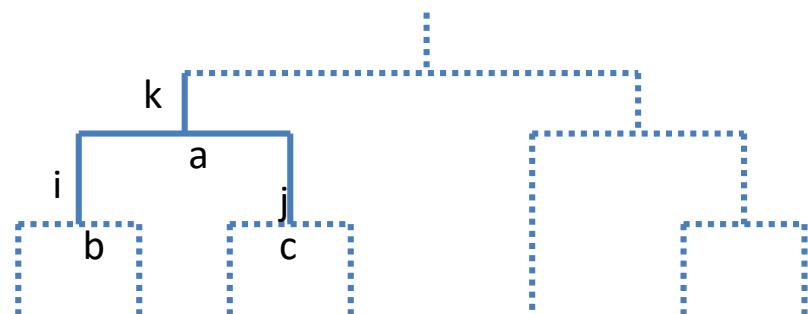
# Calculating the Likelihood for Ungapped Alignments

- $n$  sequences of length  $N$ , site  $u=1 \dots N$
- Given a rooted tree contains  $2n - 1$  nodes,  $1 \dots n$  being the leaf nodes,  $n+1 \dots 2n-1$  non-leaf, tree lengths  $t_1, \dots, t_{2n-1}$ .
- Let  $a(i)$  denote the ancestor of node  $a^i$

$$P(x^1, \dots, x^n | T, t) = \prod_{u=1}^N P(x_u^1, \dots, x_u^n | T, t)$$
$$P(x_u^1, \dots, x_u^n | T, t) = \sum_{a^{n+1}, \dots, a^{2n-1}} q_{a^{2n-1}} \prod_{i=n+1}^{2n-2} P(a^i | a^{\alpha(i)}, t_i)$$
$$\times \prod_{i=1}^n P(x_u^i | a^{\alpha(i)}, t_i)$$

# Felsenstein's Recursive Algorithm

- Let  $P(L_k | a)$  denote the probability of all the leafs below node  $k$  given that the residue at  $k$  is  $a$ .
- Then we compute  $P(L_k | a)$  from the probabilities  $P(L_i | b)$  and  $P(L_j | c)$  for all  $b$  and  $c$ , where  $i$  and  $j$  are the daughter nodes of  $k$ .



# Felsenstein's Recursive Algorithm

- Initialization: set  $k=2n-1$
- Recursion: Compute  $P(L_k | a)$  for all  $a$  as follows:
  - If  $k$  is leaf node:  $P(L_k | a)=1$  only if  $a = x_u^k$ .
  - If  $k$  is not a leaf node:
    - Compute  $P(L_i | a)$ ,  $P(L_j | a)$  for all  $a$  at the daughter nodes  $i, j$ , and set  $P(L_k | a) = \sum_{bc} P(b|a, t_i)P(L_i|b)P(c|a, t_j)P(L_j|c)$
- Termination: Likelihood at site  $u$ ,

$$P(x_u | T, t) = \sum_a P(L_{2n-1} | a) q_a$$

# Reversibility & Independence of Root Position

- The score of the optimal tree is independent of the root position if and only if:
  - the substitution matrix is **multiplicative**
  - the substitution matrix is **reversible**
- A substitution matrix is reversible if for all a,b and t:

$$P(b|a, t)q_a = P(a|b, t)q_b$$

# Maximum Likelihood (ML)

- Score each tree by
  - Assumption of independent positions “m”
- Branch lengths  $t$  can be optimized
  - Gradient Ascent
  - EM
- We look for the highest scoring tree
  - Exhaustive
  - Sampling methods (Metropolis)

# Computational Problem

- Such procedures are computationally expensive!
- Computation of optimal parameters, per candidate, requires non-trivial optimization step.
- Spend non-negligible computation on a candidate, even if it is a low scoring one.
- In practice, such learning procedures can only consider small sets of candidate structures

# 构建进化树算法

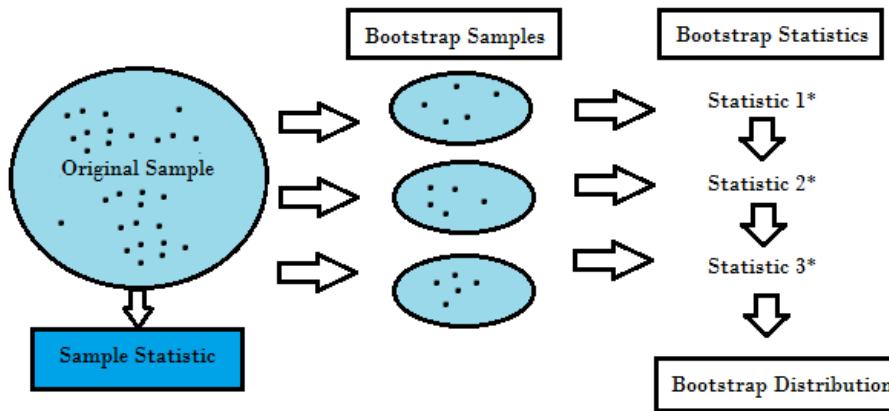
- 如果序列的相似性较高，各种方法都会得到不错的结果，模型间的差别也不大.
- 若有合适的分子进化模型可供选择，用最大似然法构树获得的结果较好.
- 对于近缘物种序列，通常情况下使用最大简约法.
- 而对于远缘物种序列，一般使用邻接法或最大似然法.
- 对于相似度很低的序列，邻接法往往出现长枝吸引(branch attraction)现象，有时严重干扰进化树的构建.

邻接法和最大似然法是需要选择模型的:

蛋白质序列的构树模型一般选择Poisson correction(泊松修正).

核酸序列的构树模型一般选择Kimura 2-parameter (Kimura-2参数).

# 构建进化树算法



在重建进化树过程中，均需选择**bootstrap**进行树的检验：

- 一般**bootstrap**的值>70，则认为重建的进化树较为可靠。
- 如果**bootstrap**的值太低，则有可能进化树的拓扑结构有错误，进化树是不可靠的。
- 一般推荐用两种以上不同的方法构建进化树，如果所得到的进化树类似，且**bootstrap**值总体较高，则得到的结果较为可靠。
- 通常情况下，只要选择了合适的方法和模型，构出的树均是有意义的，研究者可根据自己研究的需要选择最佳的树进行分析。

# 构建分子进化树相关的软件

## 软件 网址 说明

ClustalX <http://bips.u-strasbg.fr/fr/Documentation/ClustalX/>

**ClustalW** <http://www.cf.ac.uk/biosi/research/sequence-analysis/clustalw.html>

GeneDoc <http://www.psc.edu/biomed/genedoc/>

BioEdit <http://www.mbio.ncsu.edu/BioEdit/bioedit.html>

**MEGA** <http://www.megasoftware.net/>

PAUP <http://paup.csit.fsu.edu/>

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

**PHYML** <http://atgc.lirmm.fr/phym/>

**PAML** <http://abacus.gene.ucl.ac.uk/software/paml.html>

Tree-puzzle <http://www.tree-puzzle.de/>

**MrBayes** <http://mrbayes.csit.fsu.edu/>

MAC5 <http://www.agapow.net/software/mac5/>

**TreeView** <http://taxonomy.zoology.gla.ac.uk/rod/treeview.html>

图形化的多序列比对工具

**命令行格式的多序列比对工具**

多序列比对结果的美化工具

序列分析的综合工具

**图形化、集成进化分析工具，不包括ML**

商业软件，集成的进化分析工具

免费的、集成的进化分析工具

**最快的ML建树工具**

**ML建树工具**

较快的ML建树工具

**基于贝叶斯方法的建树工具**

基于贝叶斯方法的建树工具

**进化树显示工具**

# MEGA: 建树平台

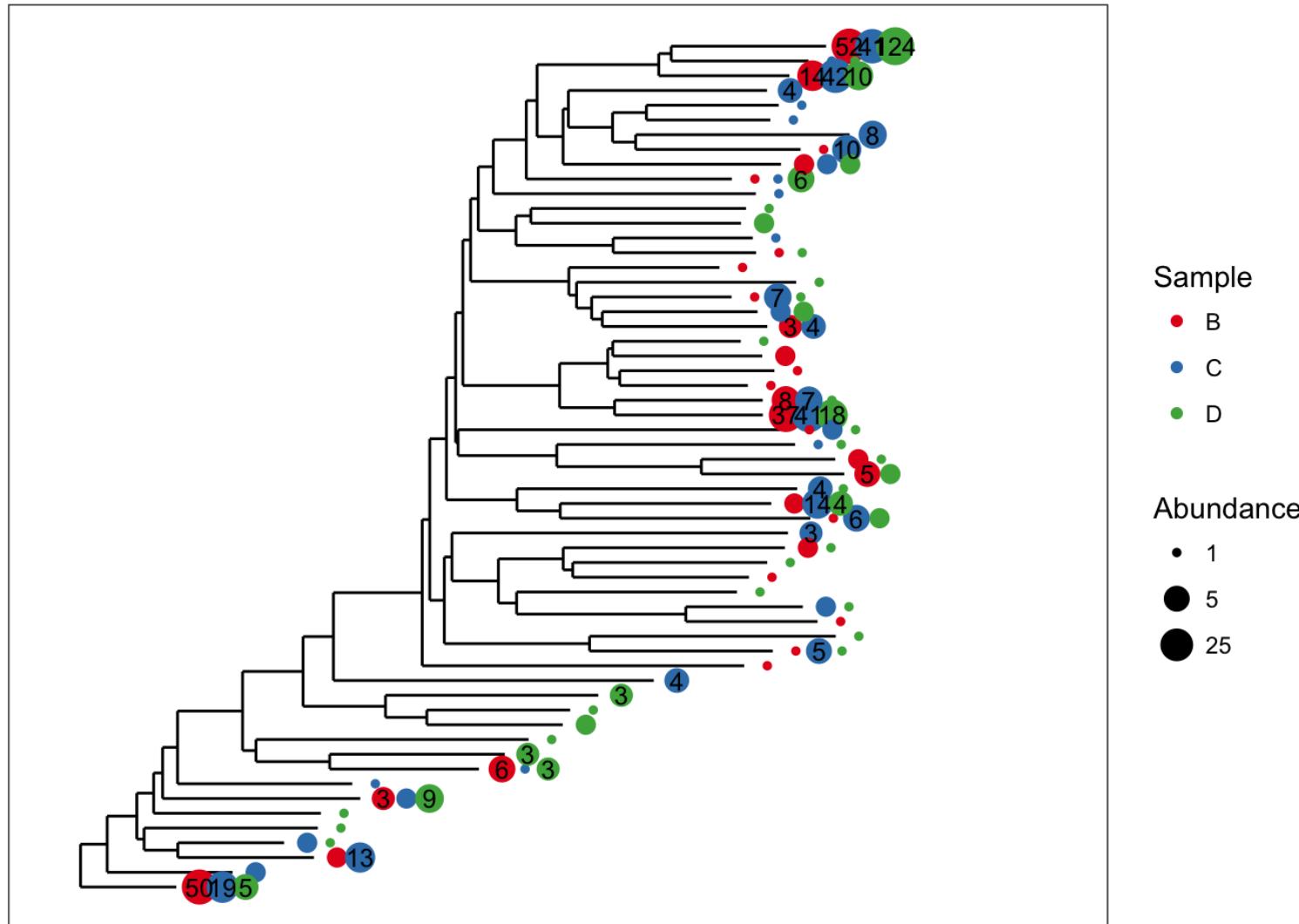
M | E | G | A  
Molecular Evolutionary  
Genetics Analysis

tutorial features documentation feedback

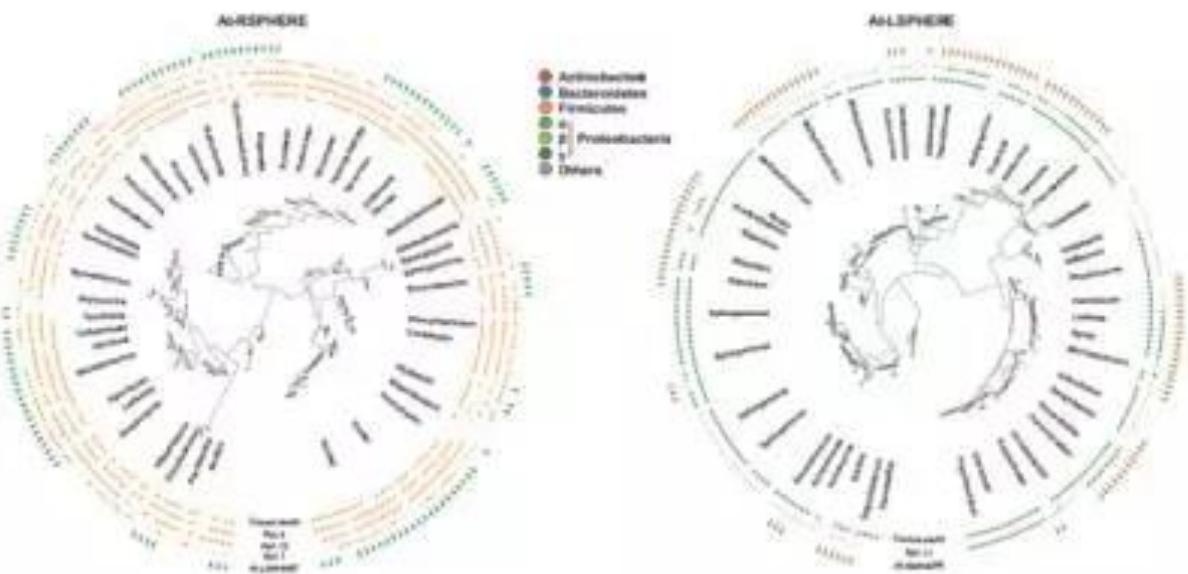


	Likelihood	Distance	Parsimony	Bayesian	Visual Explorer	Caption Expert
Phylogeny	✓	✓	✓		✓	✓
Bootstrap	✓	✓	✓		✓	✓
Distance/Diversity	✓	✓			✓	✓
Model Selection	✓					
Substitution Pattern	✓				✓ XL	✓
Rate Variation	✓				✓ XL	✓
Ancestral Sequence			✓	✓	✓	✓
Clock Test	✓	✓			✓ XL	✓
Time Tree	✓	✓			✓	✓
Selection Test	✓	✓			✓	✓
Disease Mutation	✓	✓			✓ XL	✓

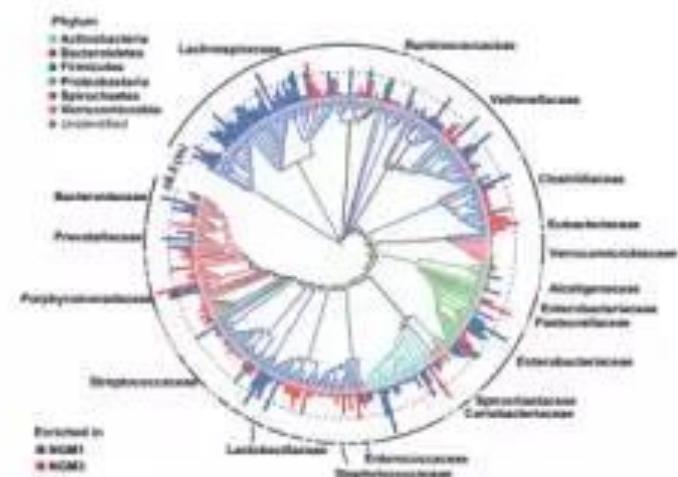
# ggplot: 物种进化关系分析



# iTOL: 物种进化关系分析

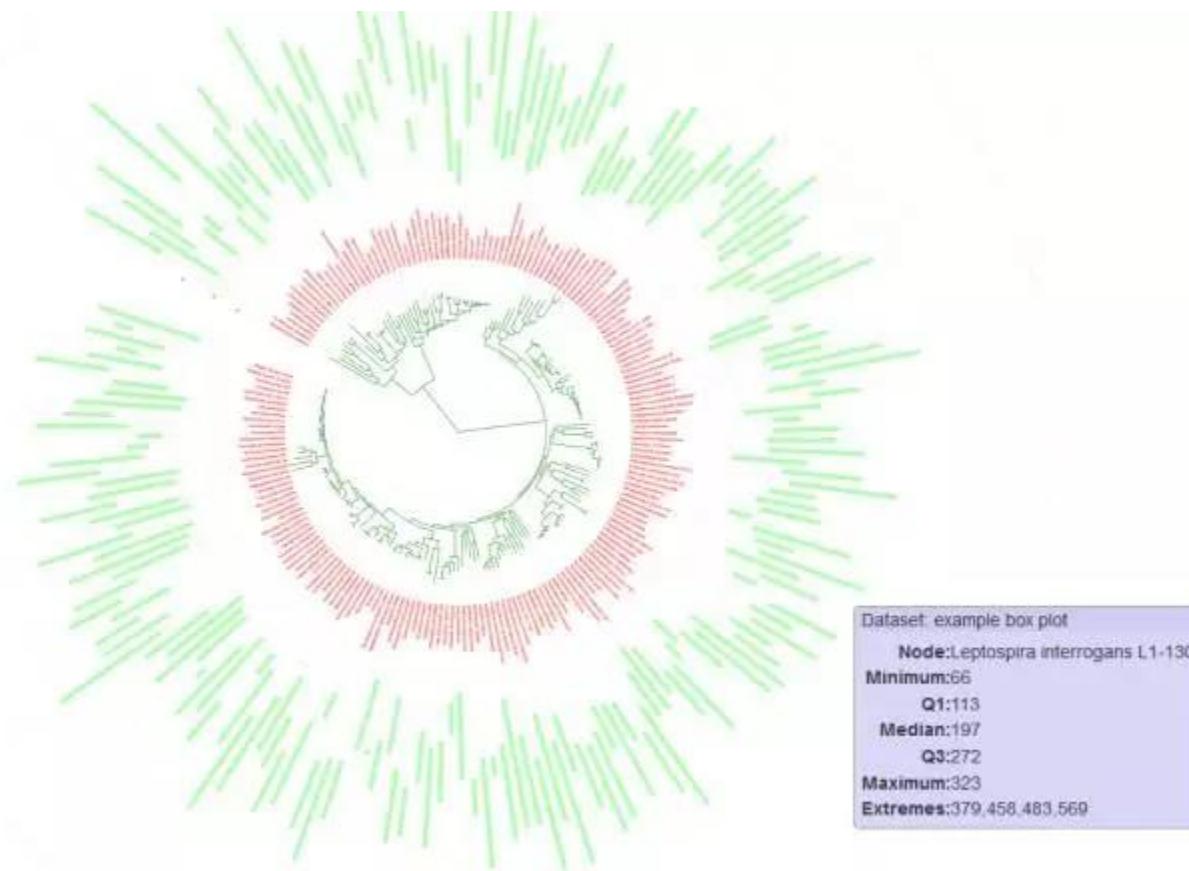


Nature 528, 364–369 (2015)

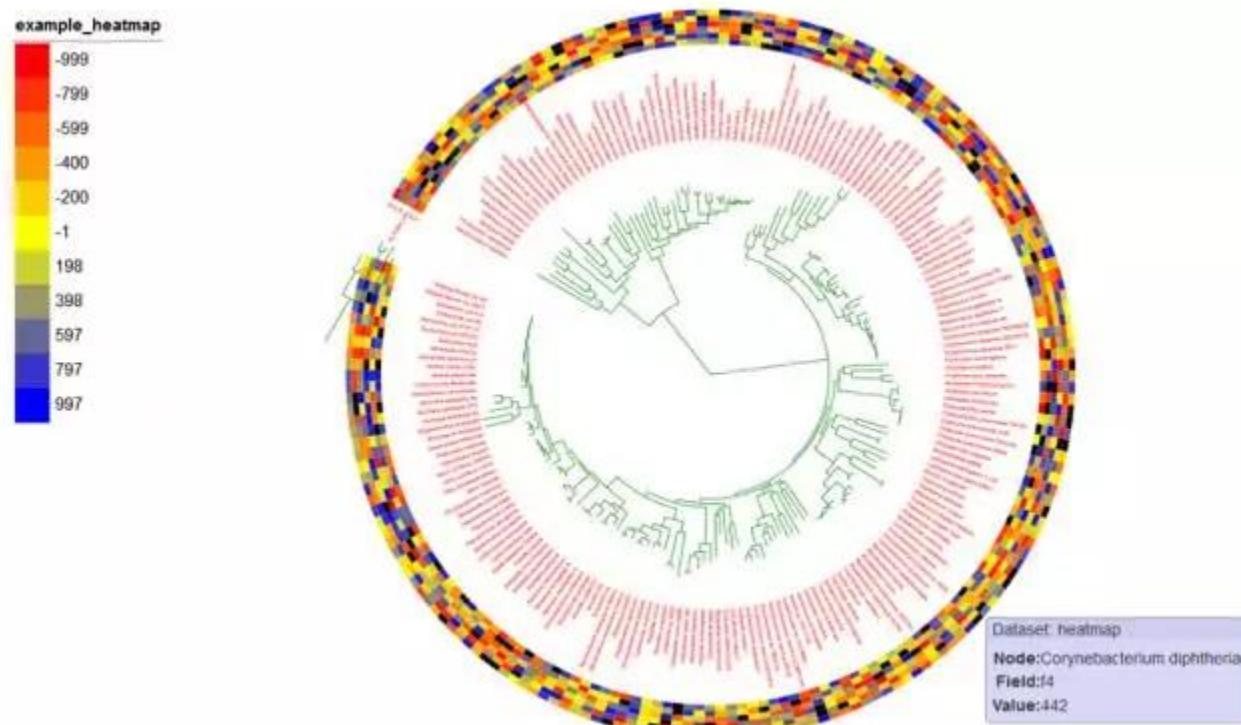


Nature Medicine 22, 1187–1191 (2016)

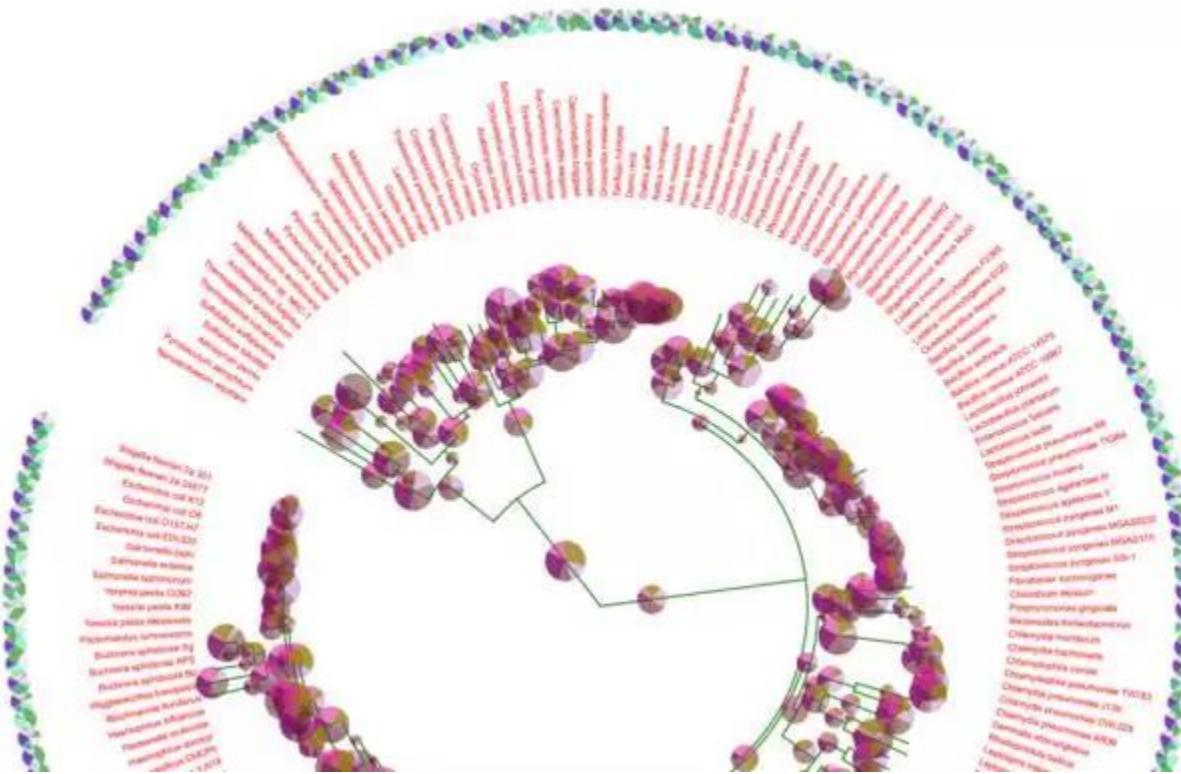
# iTOL: 物种进化关系分析



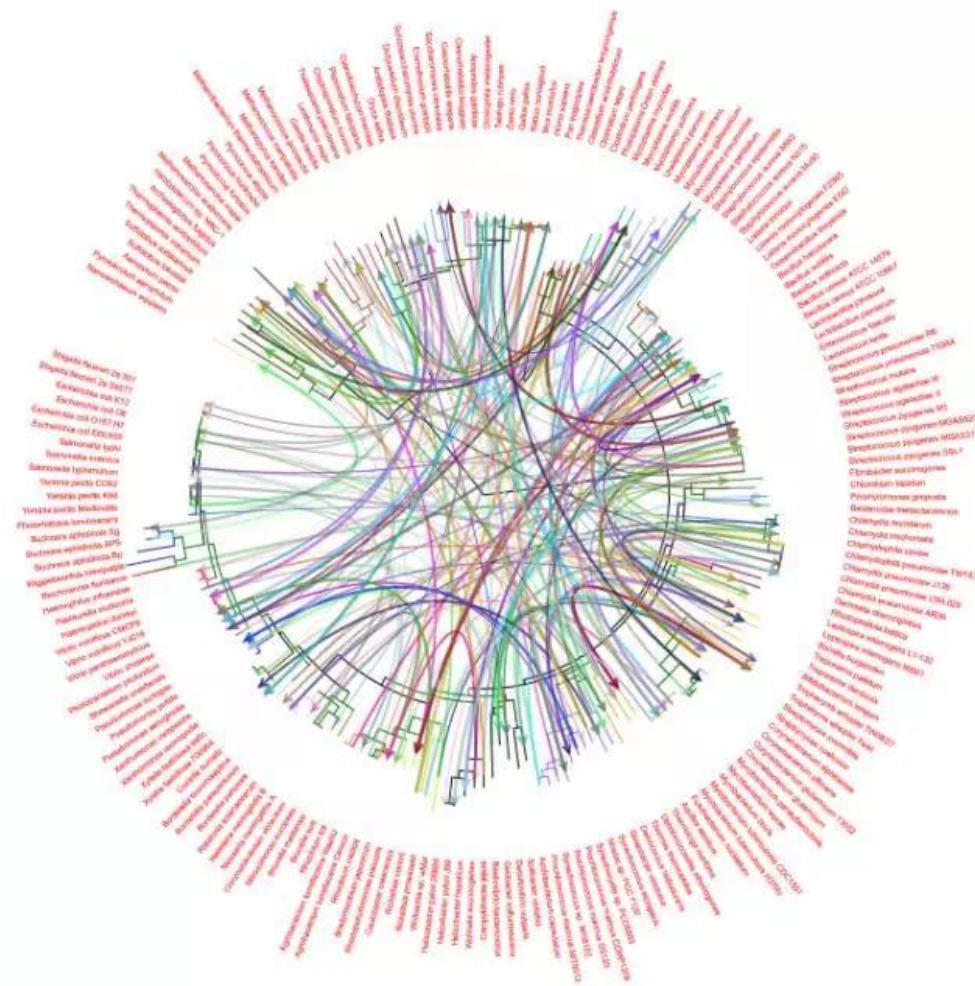
# iTOL: 物种进化关系分析



# iTOL: 物种进化关系分析



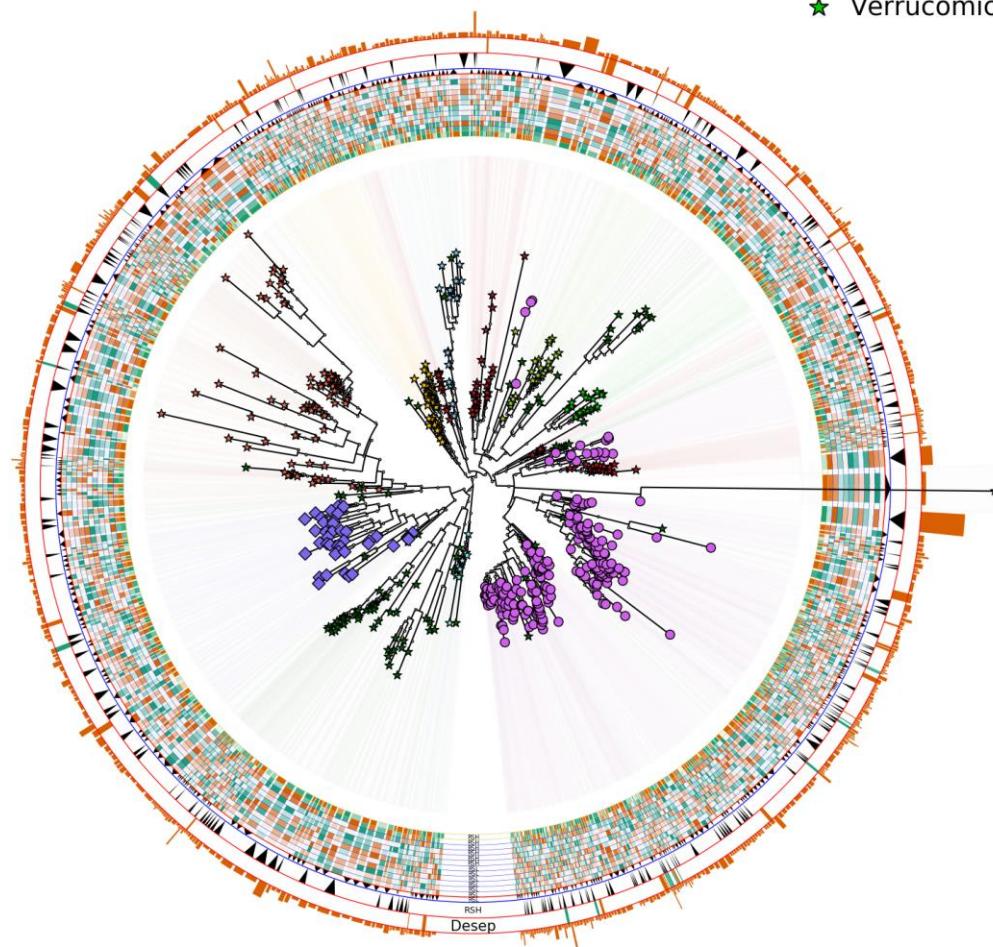
# iTOL: 物种进化关系分析



# GraPhIAn: 分类树分析

Metagenomic

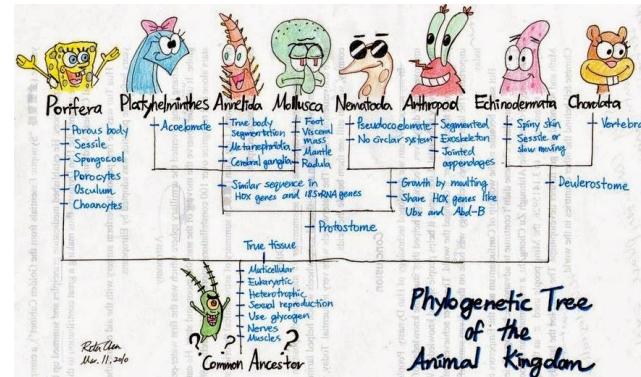
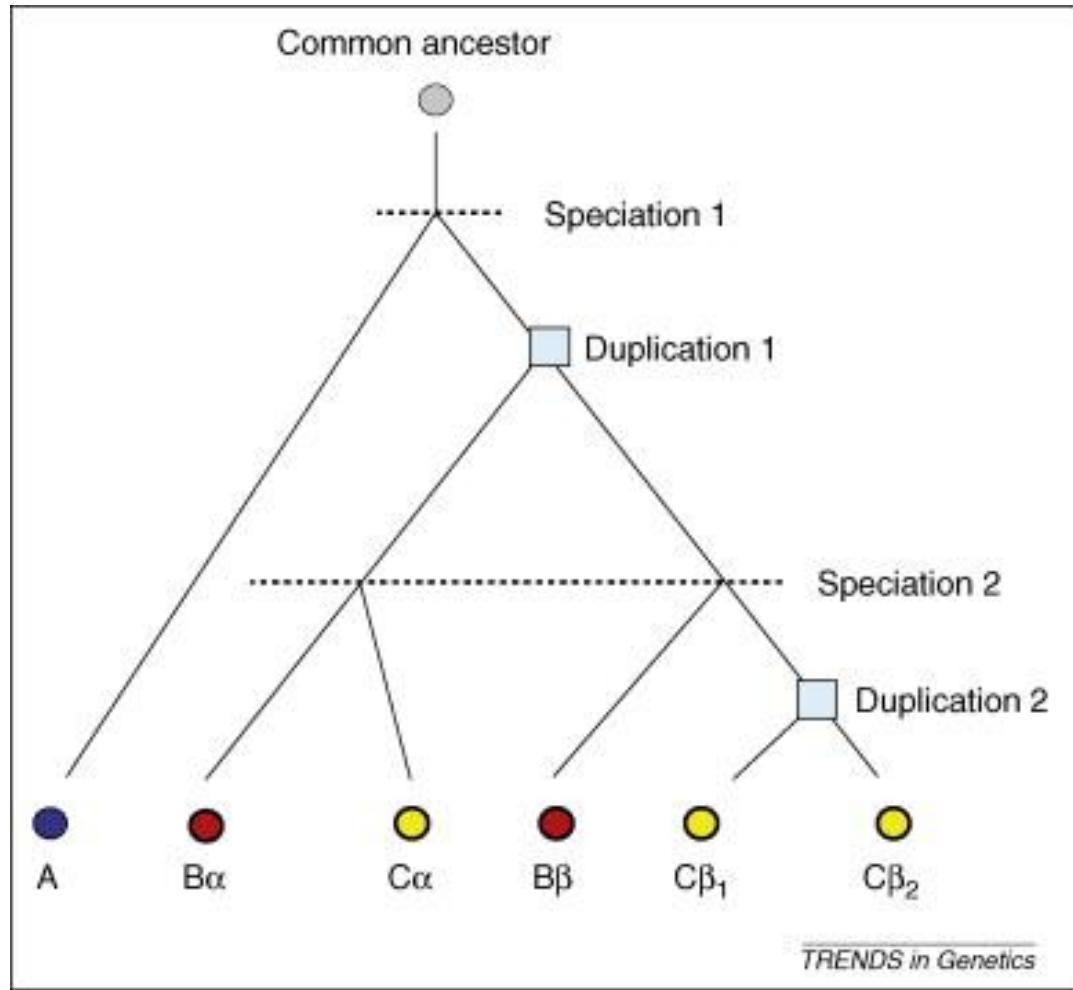
- ★ Acidobacteria
- ★ Actinobacteria
- ★ Bacteroidetes
- ◆ Chloroflexi
- ◆ Cyanobacteria
- ◆ Firmicutes
- ◆ Gemmatimonadetes
- ◆ Others
- ◆ Proteobacteria
- ◆ Verrucomicrobia



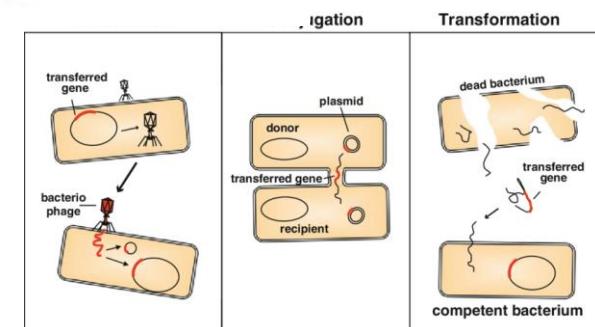
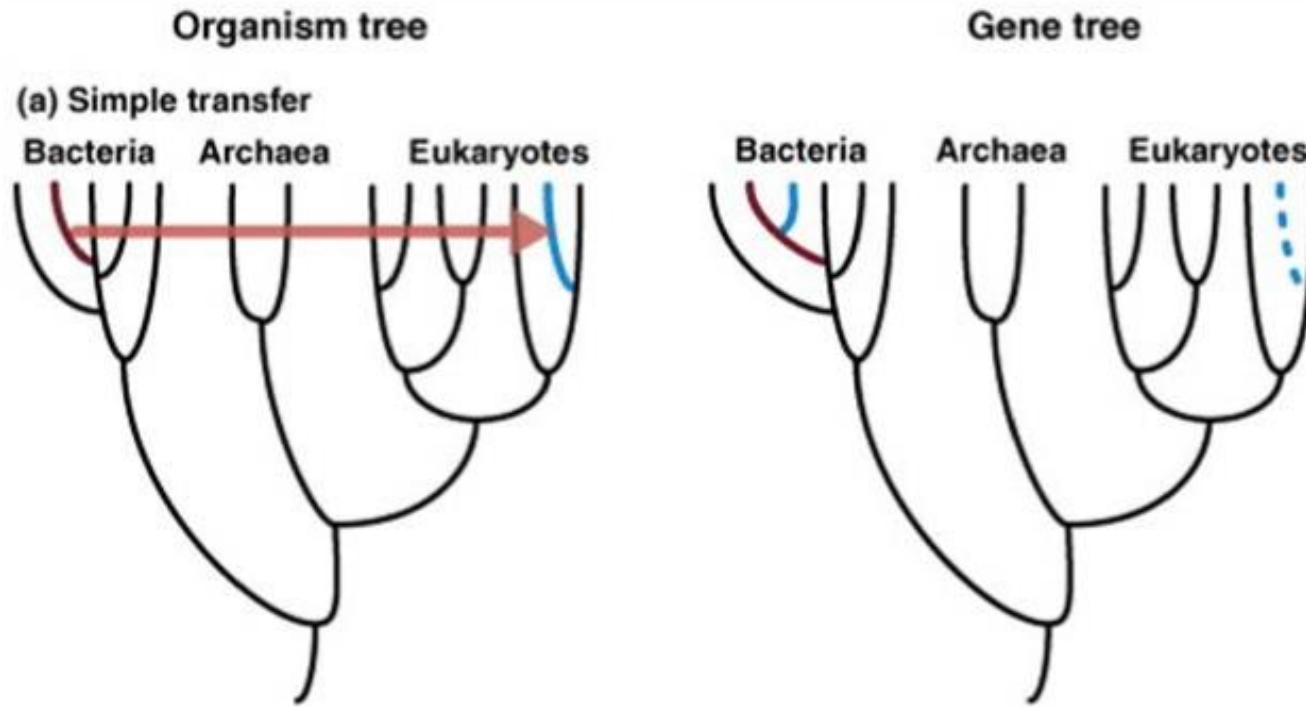
# 构建分子进化树相关的软件

- 就进化树而言， iTOL功能最为全面。iTOL无限制添加的数据集，外环可以制作各种图形包括箱线图等，可以使用多种符号填充外环。但Graphlan就不能这么随便了，只能使用两个符号填充外环。再多的环属性也就是设置环数量和颜色，透明度了。
- ggtree最容易上手，但是就一张圈图来说，它不能添加除了热图以外的其他图形，但是在非圈图的模式下，可以对多种数据进行合并，方法将更为简单，操作也容易一些。
- Graphlan可以制作分类树，是它不同于R包ggtree和iTOL的地方。

# Gene: ortholog and paralog



# Gene: HGT

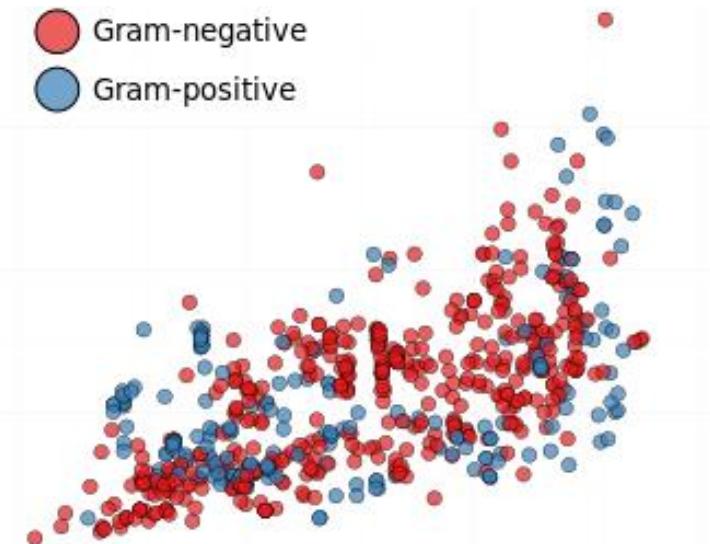


# Gene: HGT

Think about it:  
how does HGT affect the molecular clock calculation?

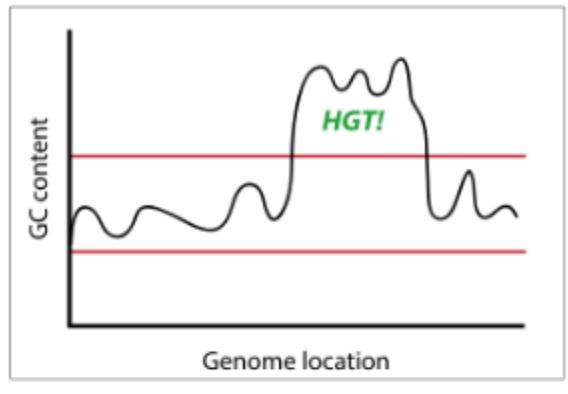


- Gram-negative
- Gram-positive

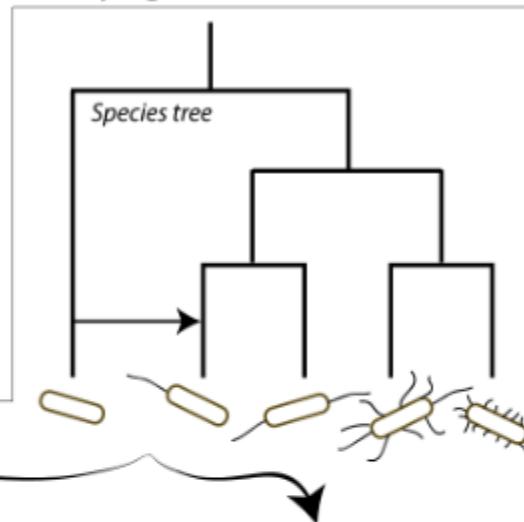


# HGT identification methods (phylogeny analysis)

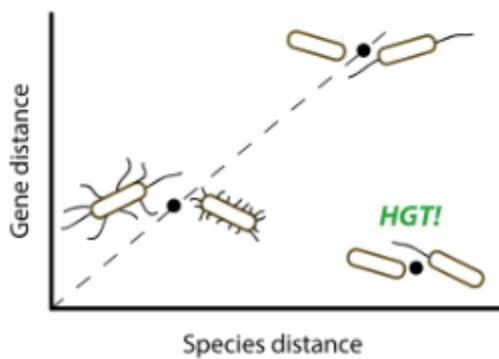
## 1. Parametric methods



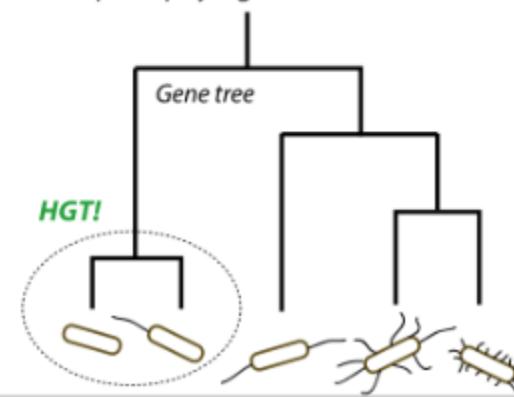
## 2. Phylogenetic methods



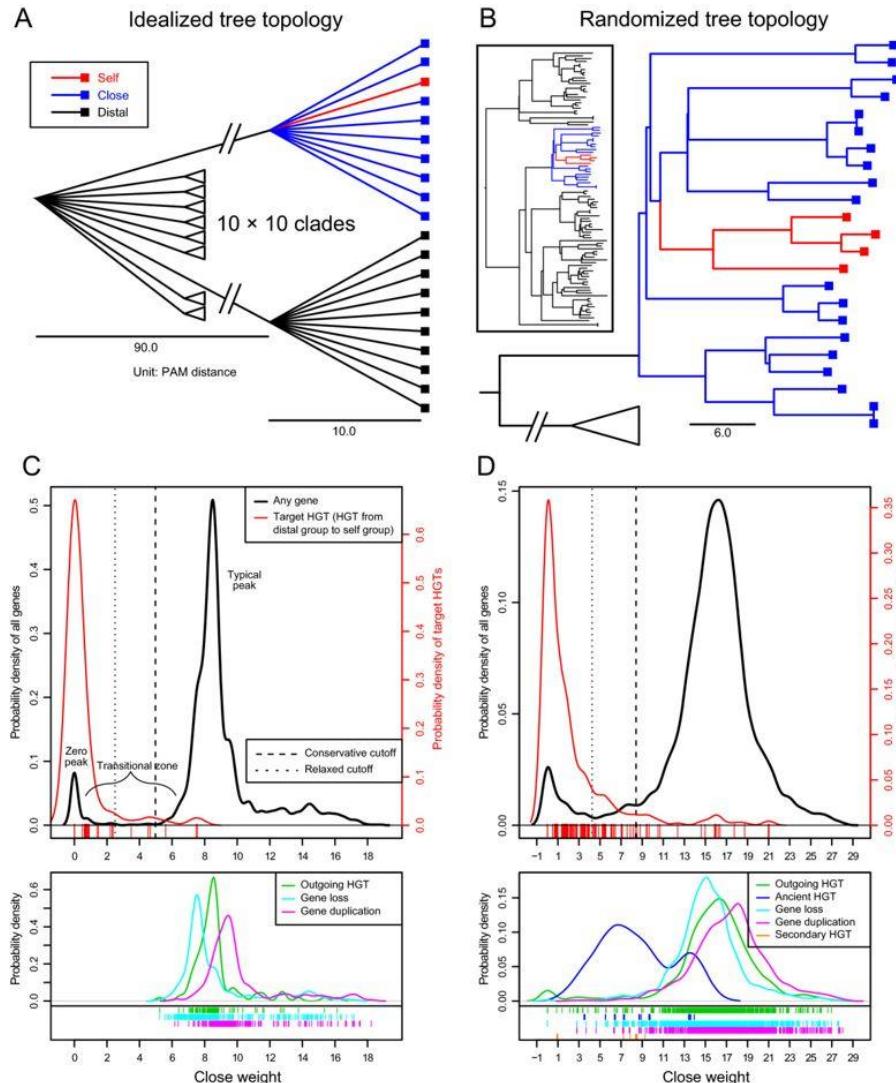
### 2a. Implicit phylogenetic methods



### 2b. Explicit phylogenetic methods



# HGT identification methods (phylogeny analysis)



# 参考文献

- R. Durbin, S. Eddy, A. Krogh and G. Mitchison.  
Biological Sequence Analysis—Probabilistic  
Models of Proteins and Nucleic Acids. 1998,  
Cambridge University Press.

# References

