# Molecular Evolution and Phylogeny in Microbiology Research
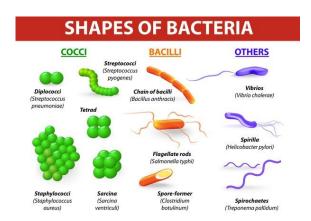
Sishuo WANG

# CONTENT

# Background: classification vs. phylogenetics

- Classification: the process of organizing and categorizing anything
  - Gram staining: Gram-negative, gram-positive
  - Shape: rod-shaped, spiral, coccoid
  - Pathogenic vs. non-pathogenic

- Phylogeny: purely based on the evolutionary relationships

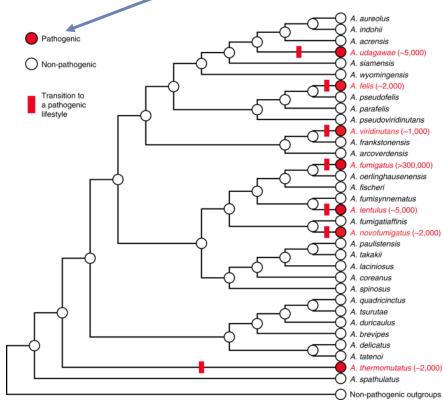s (from Greek φῦλον phylon "tribe" + γένεσις genesis "origin")

**SHAPES OF BACTERIA**

COCCI

BACILLI

OTHERS

Diplococci (Streptococcus pneumoniae)

Streptococci (Streptococcus pyogenes)

Chain of bacilli (Bacillus anthracis)

Vibrios (Vibrio cholerae)

Tetrad

Flagellate rods (Salmonella typhi)

Spirilla (Helicobacter pylori)

Staphylococci (Staphylococcus aureus)

Sarcina (Sarcina ventriculi)

Spore-former (Clostridium botulinum)

Spirochaetes (Treponema pallidum)

Charles Darwin's 1837 sketch

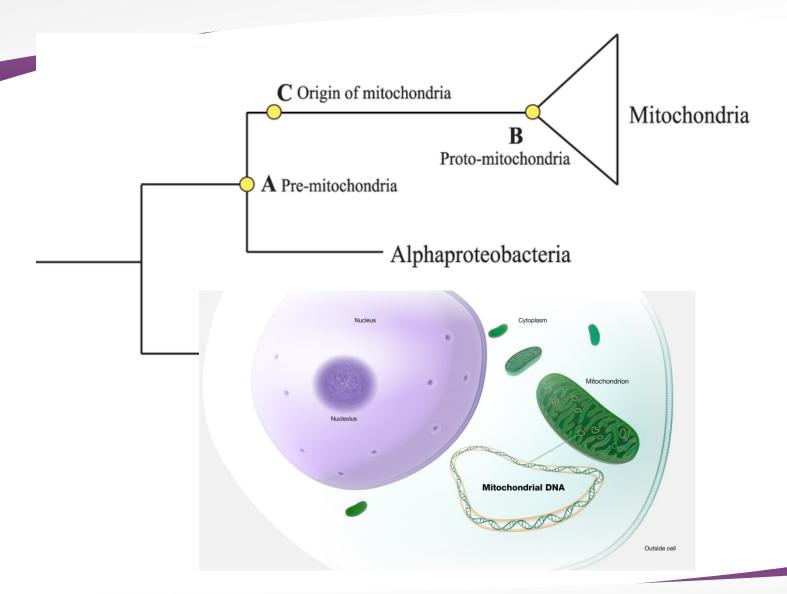https://www.shutterstock.com/image-vector/common-bacteria-infecting-human-vector-illustration-160035755

Independent transitions to pathogenic in the fungal genus *Aspergillus.*



Rokas, 2022

# Why study phylogeny:
# reveal intrinsic connections

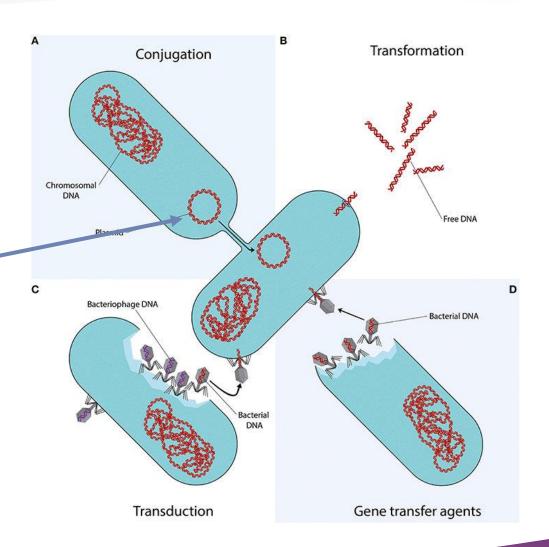https://www.genome.gov/genetics-glossary/Mitochondrial-DNA

# What's horizontal gene transfer (HGT)

Under antibiotics pressure, strains w/ AMR genes on the plasmid may have a higher opportunity to survive!



A. Conjugation
B. Transformation
C. Transduction
D. Gene transfer agents

English: I eat dim sum.

English: He enjoys sushi.

# Can you identify "horizontally transferred" words?

- English: I eat dim sum.
- German: Ich esse dim sum.
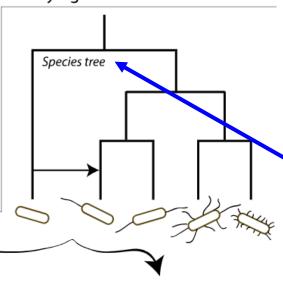  - Source: Cantonese

- English: He enjoys sushi.
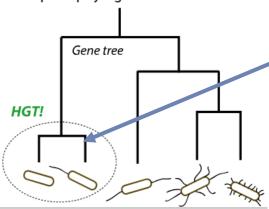- German: Er genießt sushi.
  - Source: Japanese

# Why studying phylogeny: Horizontal gene transfer (HGT)
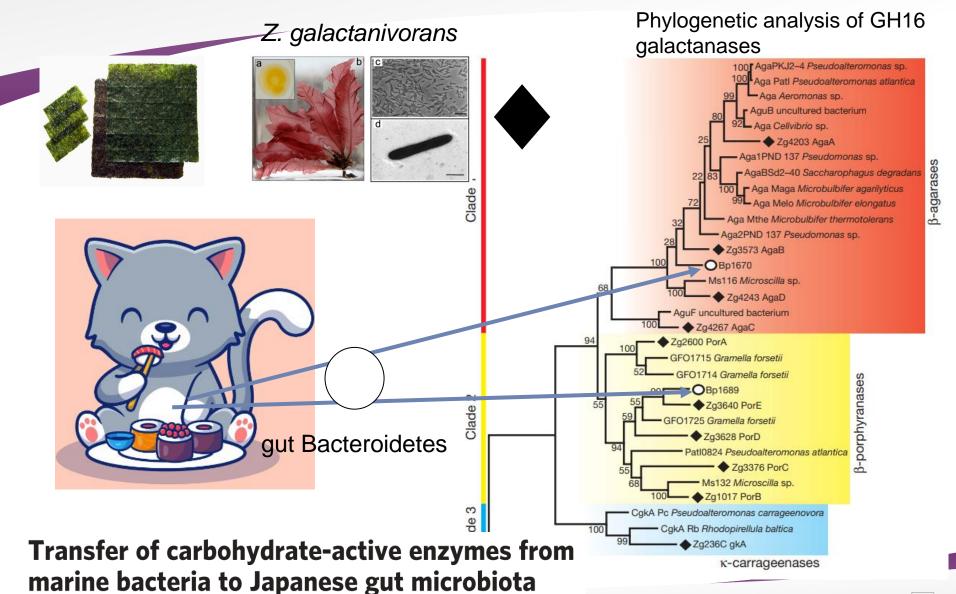
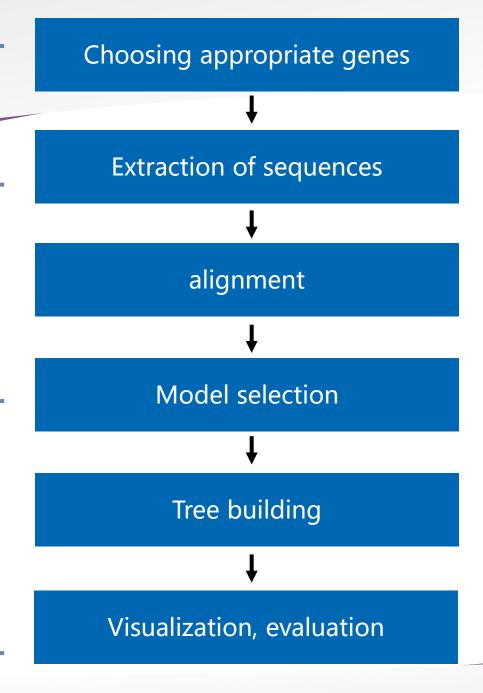

2. Phylogenetic methods

Species tree

2b. Explicit phylogenetic methods

Gene tree

HGT!

Conflict btwn gene tree and species tree!

Ravenhall et al. 2015

# HGT: a real example

*Z. galactanivorans*

Phylogenetic analysis of GH16 galactanases



**Transfer of carbohydrate-active enzymes from marine bacteria to Japanese gut microbiota**

gut Bacteroidetes

Hehemann, et al., 2010, Nature
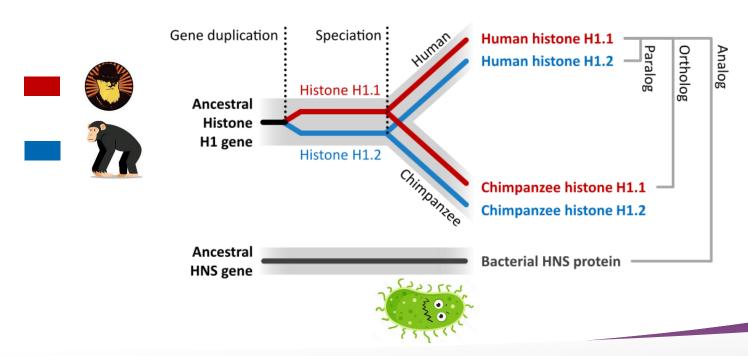
- For inferring species relationship: orthologs (resistant to HGT)
  - 16S
  - Other marker genes (house-keeping)
  - Single-copy genes

- For inferring gene trees: all genes within the same family
  - ❑ Gene family: homologs, or genes that are evolutionarily related (share a common origin)

# Orthologs vs. Paralogs

- Ortholog: genes separated by **speciation**
- Paralog: genes separated by gene **duplication**.
- Analog: genes not evolutionarily related but with similar functions

- English: I eat dim sum.
- German: Ich esse dim sum.

- Ortholog:
  - I – Ich
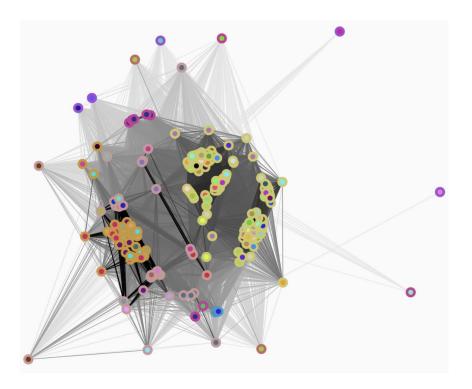  - eat – esse

- HGT:
  - dim sum

**ChatGPT**

Yes, the German word "esse" and the English word "eat" share a common origin. Both words are derived from the same Germanic root, which can be traced back to the Proto-Germanic word "*etaną." Over time, this root evolved into "essen" in German and "eat" in English, among other forms in various Germanic languages. So, "esse" and "eat" are cognates, meaning they have a common linguistic ancestor.

# How to identify orthologs (and paralogs)



- Clustering based on sequence similarity (BLAST) into families
  - Inparanoid
  - OrthoMCL
  - OrthoFinder
  - MMseqs2

- Typically, you can take **single-copy** genes as orthologs.
- If higher resolution needed, SNP or genome-wide nucleotide tree is needed.

https://orthomcl.org/orthomcl/app

# Building a tree

- Most parsimony (MP)
- Distance-based methods (not introduced)
- Maximum likelihood (ML)
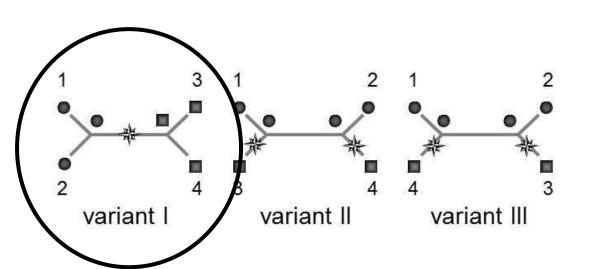- Bayesian approaches

# Tree building: maximum parsimony (MP)



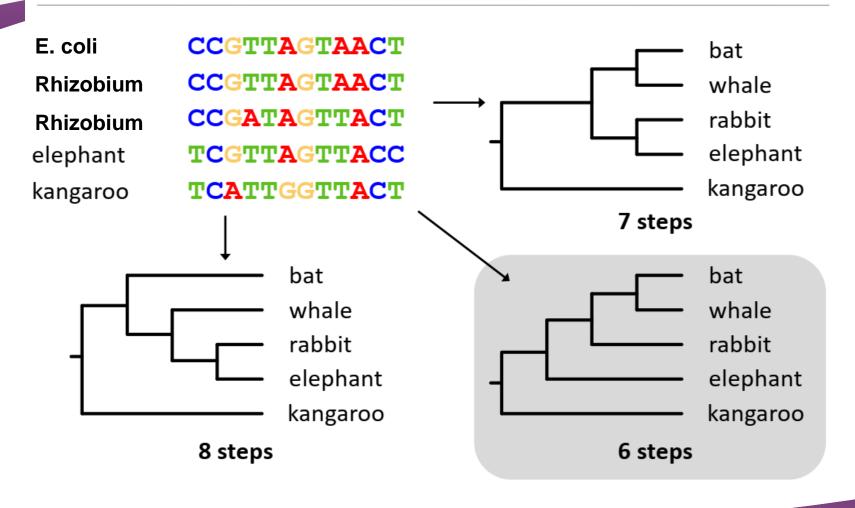| | character |
|---|---|
| species 1: | ● |
| species 2: | ● |
| species 3: | ■ |
| species 4: | ■ |

Occam's razor: the tree that involves the least changes in the character (nucleotides, amino acids, etc.), is the best.



**OCCAM'S RAZOR**

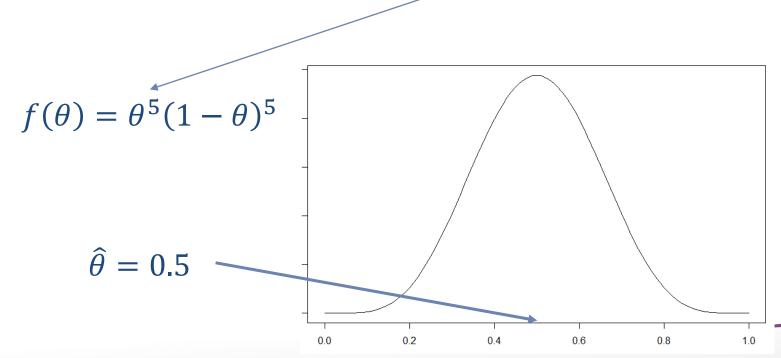PROBLEM SOLVING PRINCIPLE TO SIMPLE SOLUTIONS



variant I          variant II          variant III

https://www.frozenevolution.com/xxiii52-homologies-and-homoplasies-can-also-be-distinguished-using-maximum-parsimony-principle

# Maximum parsimony

Taken from Simon Ho's slides
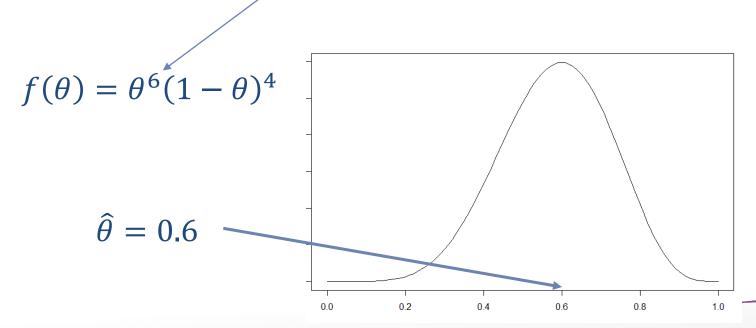
# Maximum likelihood estimation

- Maximum likelihood estimate (MLE) maximizes the likelihood.
- Consider you flip a coin for 10 times and 5 heads and 5 tails. What's the MLE of the probability the coin landing heads?

$$f(\theta) = \theta^5 (1 - \theta)^5$$

$$\hat{\theta} = 0.5$$

# Maximum likelihood estimation

- Maximum likelihood estimate (MLE) maximizes the likelihood.
- Consider you flip a coin for 10 times and 5 heads and 5 tails. What's the MLE of the probability the coin landing heads?
- What's if you 6 heads and 4 tails?

$$f(\theta) = \theta^6(1 - \theta)^4$$
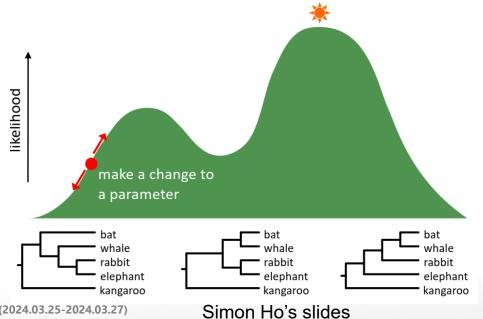
$$\hat{\theta} = 0.6$$
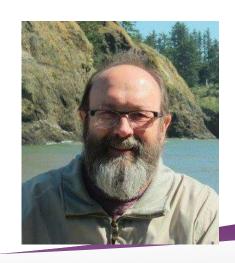
# Maximum likelihood tree

Two key points for a tree:
- Branch length
- Topology

So the ML tree consists of **topology + branch** lengths that maximize the likelihood given the observed DNA/protein sequences.
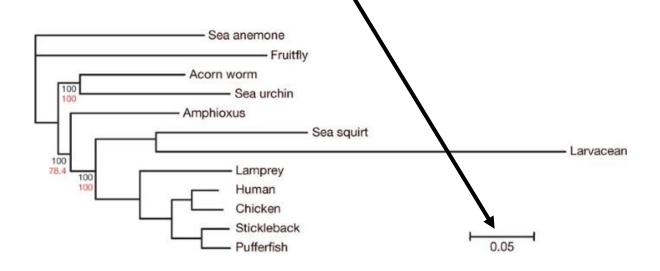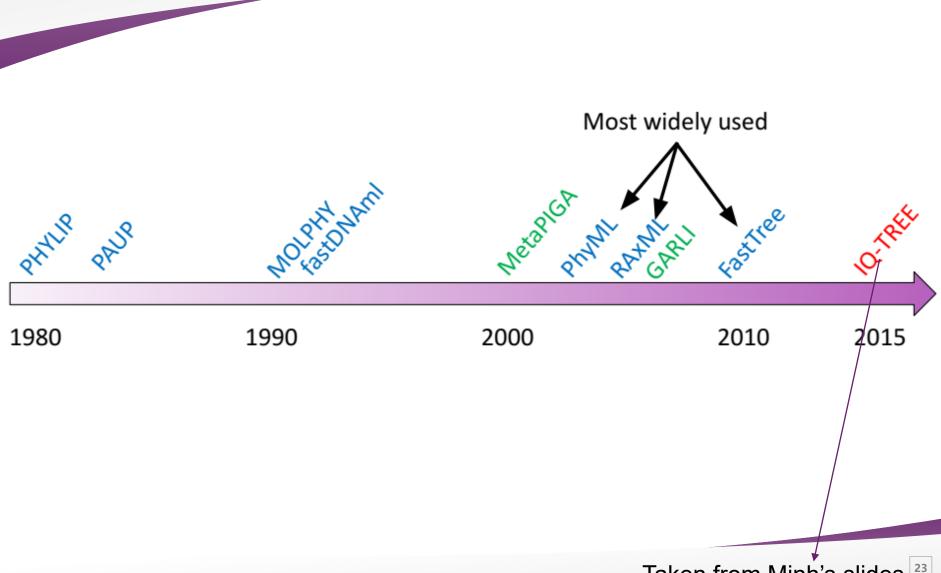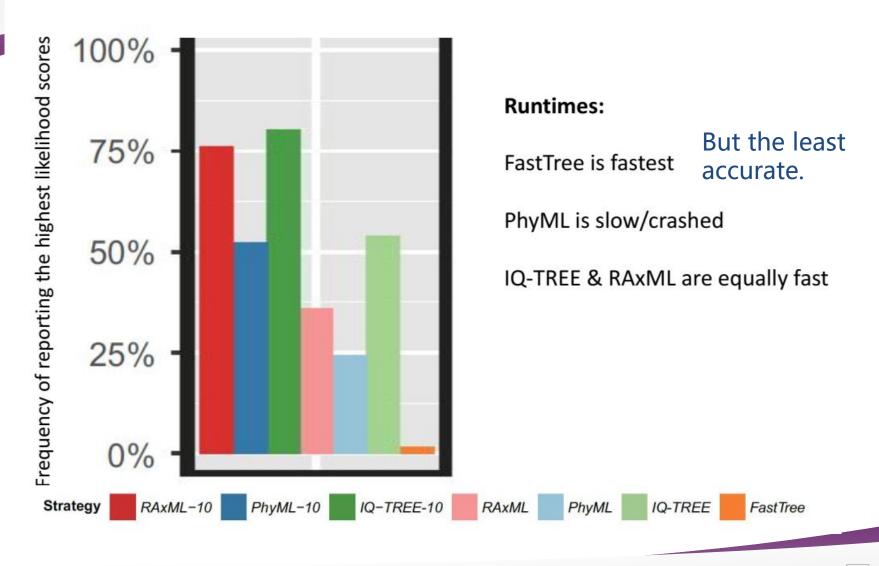


Simon Ho's slides

Felsenstein, J. 1981

In ML and Bayesian phylogenies, **branch length** represent substitutions/site (has nothing to do with time). So a longer branch means more changes compared to the ancestor.

https://www.nature.com/articles/nature06967

# maximum likelihood phylogenetics software



Most widely used

| PHYLIP | PAUP | MOLPHY fastDNAml | MetaPIGA | PhyML | RaxML | GARLI | FastTree | IQ-TREE |

1980          1990          2000          2010    2015

Taken from Minh's slides

# An independent benchmark by Zhou et al. (2018)



**Runtimes:**

FastTree is fastest

But the least accurate.

PhyML is slow/crashed

IQ-TREE & RAxML are equally fast

Strategy: RAxML–10, PhyML–10, IQ-TREE–10, RAxML, PhyML, IQ-TREE, FastTree

Modified on Minh's slides

# Bayesian inference

- Set priors for each parameter: branch length, topology, and others.

- The result will be an update of your priors based on the data.

Yang, ZH & Rannala, B., 1997

We can publish the work in Nature!!!

If a professor says confidently he can publish your study in Nature, how likely will the study be published in Nature?

The probability of your study being published in Nature, conditioned on your boss saying so.

Pr(your study being published in Nature | your boss saying so)

Let's write it in this way:
A: really publish in Nature
B: prof says "we'll publish the study in Nature"
P(A|B)

# A toy example

$$P(A|B) = P(A) \times \frac{P(B|A)}{P(B)}$$

posterior     prior     likelihood / marginal

A: really publish in Nature

B: **prof says** "we'll publish the study in Nature"

- Prior: $P(A) = 0.1$ ⟹ $P(\overline{A}) = 1 - 0.1 = 0.9$
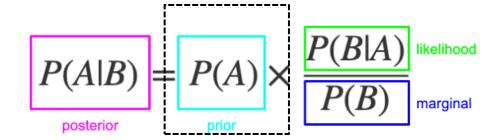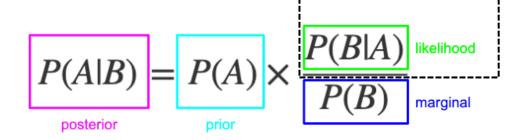
- Cond prob:
  - $P(B|A) = 1.0$
  - $P(B|\overline{A}) = 0.5$

- $P(B) = P(A)P(B|A) + P(B|\overline{A})P(\overline{A}) = 0.1 * 1 + 0.5 * 0.9 = 0.55$

$P(A|B) = 0.1*1.0/0.55 = 0.18$

# What if the journal has a higher acceptance rate?

$$P(A|B) = P(A) \times \frac{P(B|A) \text{ likelihood}}{P(B) \text{ marginal}}$$

posterior — prior

A: really publish in Nature

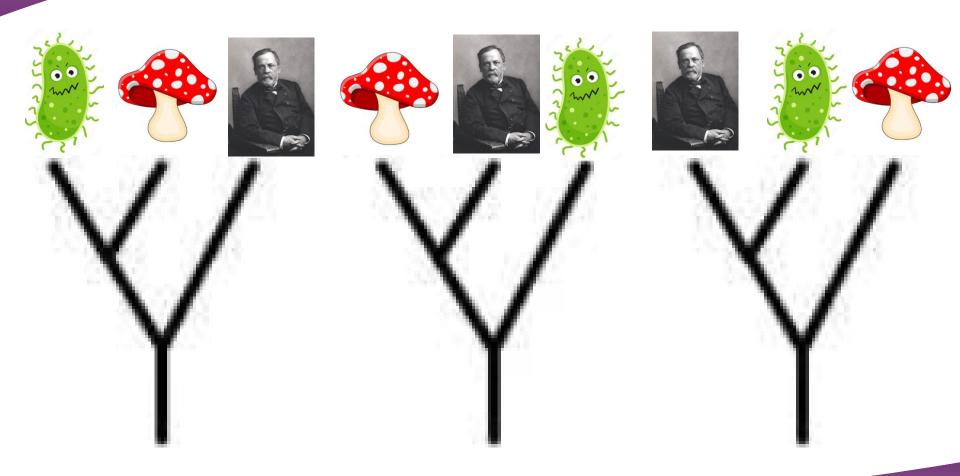B: **prof says** "we'll publish the study in Nature"

- Prior: P(A) = 0.2 ➡ $P(\overline{A}) = 1 - 0.2 = 0.8$

- Cond prob:
  - P(B|A)= 1.0
  - $P(B|\overline{A})$= 0.5

- $P(B) = P(A)P(B|A) + P(B|\overline{A})P(\overline{A}) = 0.2 * 1 + 0.8 * 0.5 = 0.6$

P(A|B)= 1.0*0.2/0.6 = 1/3

# What if your boss is more conservative?

$$P(A|B) = P(A) \times \frac{P(B|A)}{P(B)}$$

posterior     prior     likelihood     marginal

A: really publish in Nature

B: **prof says** "we'll publish the study in Nature"

- Prior: P(A) = 0.1 ➡ P($\overline{A}$) = 1 − 0.1 = 0.9

- Cond prob:
  - P(B|A)= 1.0
  - P(B|$\overline{A}$)= 0.2

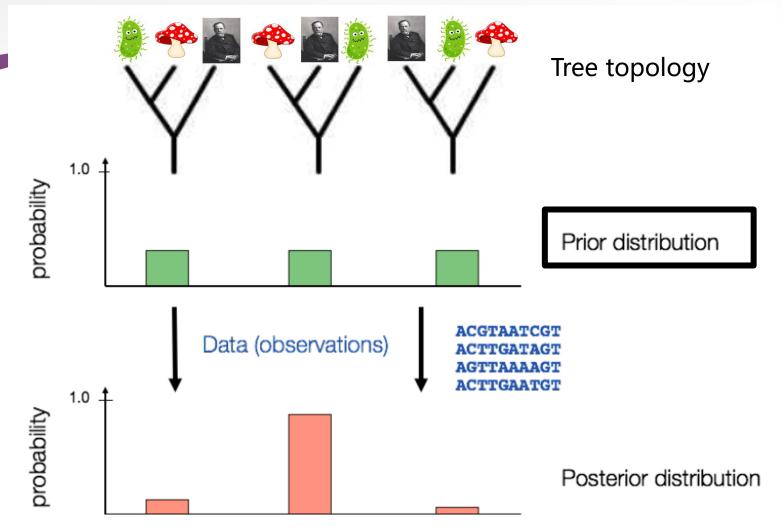- $P(B) = P(A)P(B|A) + P(B|\overline{A})P(\overline{A}) = 0.1 * 1 + 0.9 * 0.2 = 0.28$

P(A|B)= 1.0*0.1/0.28 = 0.36

Tree topology

Prior distribution

Data (observations)

ACGTAATCGT
ACTTGATAGT
AGTTAAAAGT
ACTTGAATGT

Posterior distribution

Will give you the posterior prob of each tree instead of showing a single "best" tree.

# Bayesian methods

- ## PhyloBayes
  - complex models

- ## BEAST
  - often combined with molecular clock
  - GUI but still not easy for beginners
  - GPU accelerated

- ## MrBayes
  - Not updated any more (but still useful)

- ## RevBayes (inheritance of MrBayes)
  - community efforts
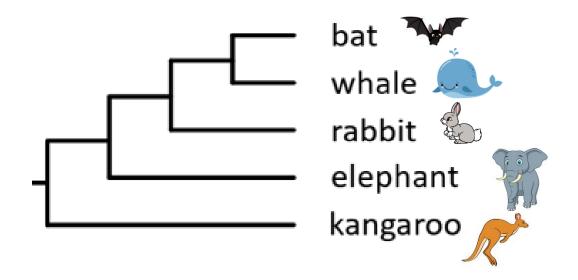  - so many functions
  - difficult to handle for beginners

Very user-friendly phylogenetics software with GUI, involving MP, distance-based, and ML phylogeny building, as well as many other functions.



Stecher et al., 2020

- You got a ML tree, but how confident are you in it?

# Evaluation: bootstrap

| bat | CCGTTAGTAACT |
| whale | CCGTTAGTAACT |
| rabbit | CCGATAGTTACT |
| elephant | TCGTTAGTTACC |
| kangaroo | TCATTGGTTACT |

**Randomly sample sites (with replacement)**

| bat | T |
| whale | T |
| rabbit | A |
| elephant | T |
| kangaroo | T |

| bat | CCGTTAGTAACT |
| whale | CCGTTAGTAACT |
| rabbit | CCGATAGTTACT |
| elephant | TCGTTAGTTACC |
| kangaroo | TCATTGGTTACT |

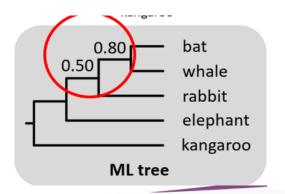| bat | TG |
| whale | TG |
| rabbit | AG |
| elephant | TG |
| kangaroo | TG |

Felsenstein, 1985
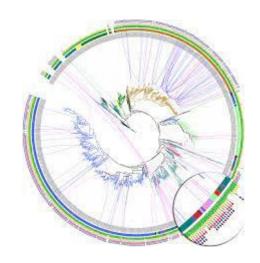
36

Bootstrapping

Taken from Ho's slides

- Maximum likelihood trees, the values on the nodes usually represent bootstrap so >= 80 are considered "reliable"
  - But IQ-Tree may use ultrafast bootstrap with 95 considered as "reliable"

- Bayesian trees, those values represent posterior probabilities, >= 95% considered as reliable.

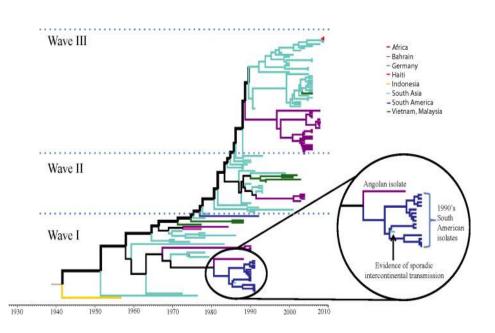

ML tree

# Visualization

- ITOL: https://itol.embl.de/

- Figtree:
https://github.com/rambaut/figtree

- R package ggtree:
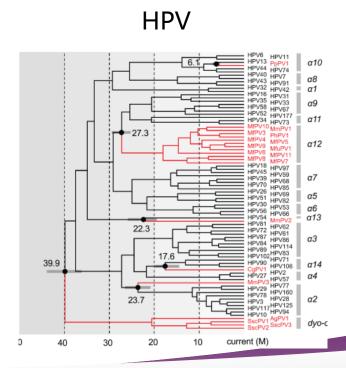https://github.com/YuLab-SMU/ggtree

# Timetree

- If you have some info about time or rates, you can convert your phylogeny into a timetree.
- Note here the branch length means substitutions/site/unit time.



Vibrio cholerae 7th pandemic lineage



HPV

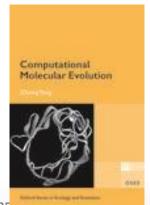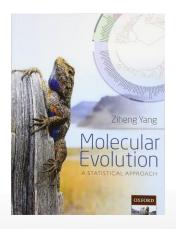Chen Zigui, et al., 2018, plos patho

- Rokas, A. Evolution of the human pathogenic lifestyle in fungi. *Nat Microbiol* **7**, 607–619 (2022).
- Wang, Zhang, and Martin Wu. "Phylogenomic reconstruction indicates mitochondrial ancestor was an energy parasite." PLoS One 9.10 (2014): e110685.
- Ravenhall M, Škunca N, Lassalle F, Dessimoz C (2015) Inferring Horizontal Gene Transfer. *PLoS Comput Biol 11(5): e1004095.*
- Felsenstein, Joseph. "Evolutionary trees from DNA sequences: a maximum likelihood approach." Journal of molecular evolution 17 (1981): 368-376.
- Yang Z, Rannala B. (1997) Bayesian phylogenetic inference using DNA sequences: a Markov chain Monte Carlo method. Molecular Biology and Evolution, 14(7):717-724.
- Felsenstein, Joseph. "Confidence limits on phylogenies: an approach using the bootstrap." evolution 39.4 (1985): 783-791.
- Zhou, X, et al. "Evaluating fast maximum likelihood-based phylogenetic programs using empirical phylogenomic data sets." *Molecular biology and evolution* 35.2 (2018): 486-503.
- Moore, S., et al. "Widespread epidemic cholera caused by a restricted subset of Vibrio cholerae clones." Clinical Microbiology and Infection 20.5 (2014): 373-379.
- Chen Z, DeSalle R, Schiffman M, Herrero R, Wood CE, et al. (2018) Niche adaptation and viral transmission of human papillomaviruses from archaic hominins to modern humans. PLOS Pathogens 14(11): e1007352.
- Glen Stecher, Koichiro Tamura, Sudhir Kumar, Molecular Evolutionary Genetics Analysis (MEGA) for macOS, Molecular Biology and Evolution, Volume 37, Issue 4, April 2020, Pages 1237–1239
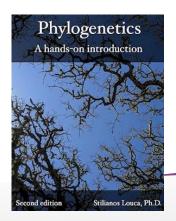
# Resources

- Yang, Ziheng. Computational molecular evolution. OUP Oxford, 2006.
- Yang, Ziheng. Molecular evolution: a statistical approach. Oxford University Press, 2014.
  - solutions: https://github.com/sishuowang/
- Louca, Stilianos. Phylogenetics, a **hands-on** introduction. 2023.

- Sydney phylogenetics workshop. https://github.com/simon-ho/SydneyPhyloWorkshop/ (organized by **Simon Ho**)

# lottery

https://numbergenerator.org/randomnumbergenerator/1-40

**Sishuo WANG**

**Thanks**
For Your Listening