

# Using Phylogenies

## Assessing Robustness and Genomic Epidemiology

---

Finlay Maguire

April 1, 2020

FCS, Dalhousie

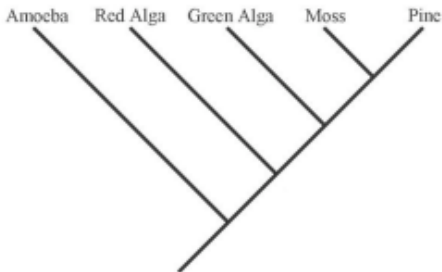
# Table of contents

1. Tree Thinking Refresher
2. Sequence Model Selection
3. Branch Support Testing
4. Comparing Trees
5. From A Single Gene to Many Genes
6. Genomic Epidemiology Phylogenetics
7. Conclusion

# Tree Thinking Refresher

---

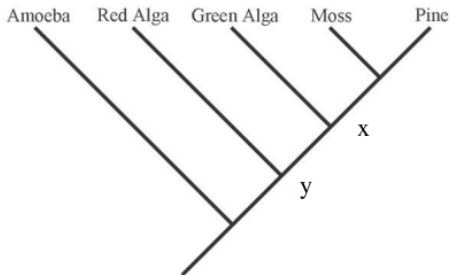
# Reading a tree



which of the following is an accurate statement of relationships?

1. A green alga is more closely related to a red alga than to a moss
2. A green alga is more closely related to a moss than to a red alga
3. A green alga is equally related to a red alga and a moss
4. A green alga is related to a red alga, but is not related to a moss

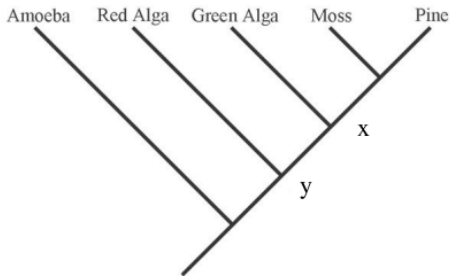
# Reading a tree



which of the following is an accurate statement of relationships?

1. A green alga is more closely related to a red alga than to a moss
2. A green alga is more closely related to a moss than to a red alga
3. A green alga is equally related to a red alga and a moss
4. A green alga is related to a red alga, but is not related to a moss

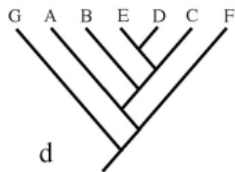
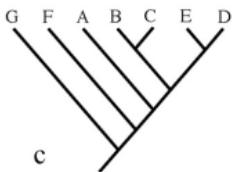
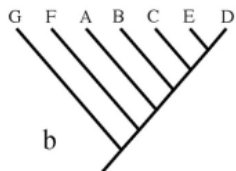
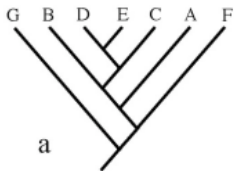
# Reading a tree



which of the following is an accurate statement of relationships?

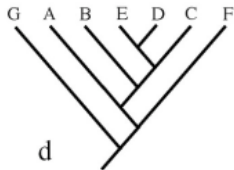
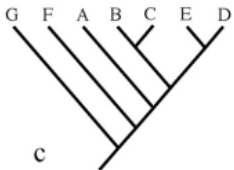
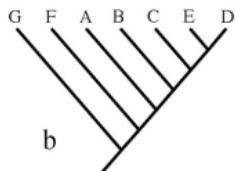
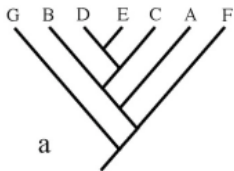
1. A green alga is more closely related to a red alga than to a moss
2. **A green alga is more closely related to a moss than to a red alga**
3. A green alga is equally related to a red alga and a moss
4. A green alga is related to a red alga, but is not related to a moss

# Comparing Topologies



Which of the four trees depicts a different pattern of relationships to the others?

# Comparing Topologies



Which of the four trees depicts a different pattern of relationships to the others?

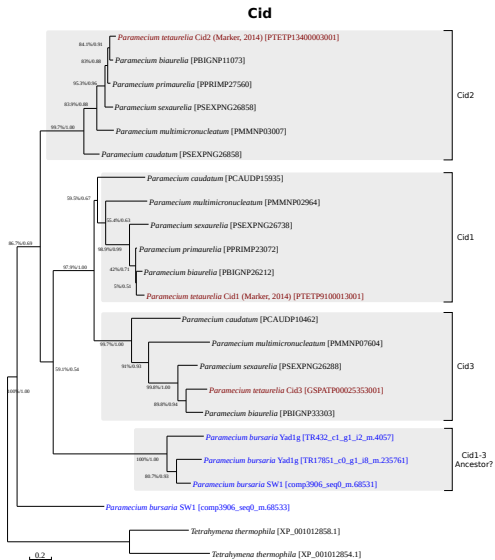
c: C is more closely related to E and D than to B in other trees.



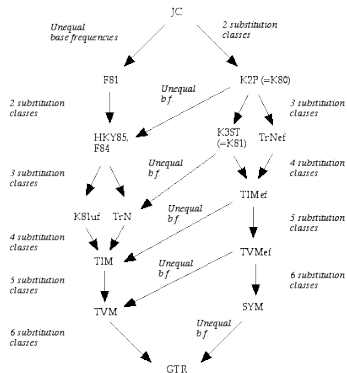
# Sequence Model Selection

---

# Phylogenies are hypotheses

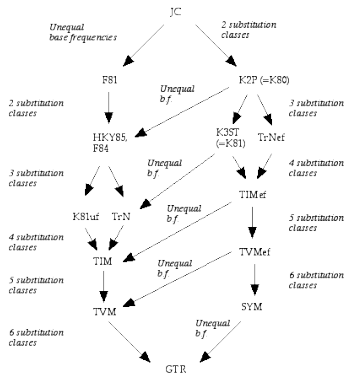


# Hypothesis testing



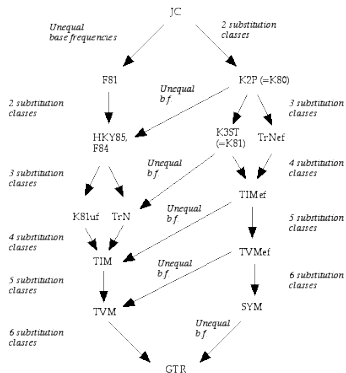
- Does another model of sequence evolution fit the data better?
- How well supported are individual branches in a tree?
- Does another tree explain the data better?

# How do we select a sequence model?



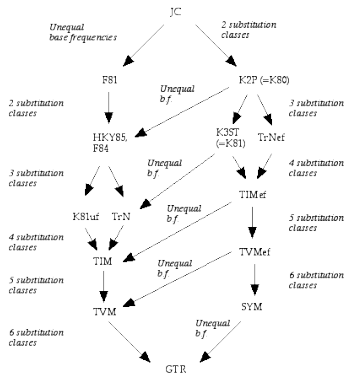
- Likelihood ratio test (LRT  $\delta$  i.e.  $p(\text{data} | \text{model})$ )

# How do we select a sequence model?



- Likelihood ratio test (LRT  $\delta$  i.e.  $p(\text{data} | \text{model})$ )
- $\delta = 2(\ln(L_1) - \ln(L_0))$

# How do we select a sequence model?



- Likelihood ratio test (LRT  $\delta$  i.e.  $p(\text{data} | \text{model})$ )
- $\delta = 2(\ln(L_1) - \ln(L_0))$
- Limitations: nested models (i.e. hLRT), order matters, no regularisation

- Akaike Information Criterion (**AIC**), penalising number of parameters:

- Akaike Information Criterion (**AIC**), penalising number of parameters:
- $AIC = -2\ln(L) + 2K$



# Information Criterion

- Akaike Information Criterion (**AIC**), penalising number of parameters:
- $AIC = -2\ln(L) + 2K$
- However, this penalises all high K models even if sample size is large too.

# Information Criterion

- Akaike Information Criterion (**AIC**), penalising number of parameters:
- $AIC = -2\ln(L) + 2K$
- However, this penalises all high K models even if sample size is large too.
- Corrected Akaike Information Criterion (**AICc**)

# Information Criterion

- Akaike Information Criterion (**AIC**), penalising number of parameters:
- $AIC = -2\ln(L) + 2K$
- However, this penalises all high  $K$  models even if sample size is large too.
- Corrected Akaike Information Criterion (**AICc**)
- $AICc = AIC + \frac{2K(K+1)}{n-K-1}$

# Information Criterion

- Akaike Information Criterion (**AIC**), penalising number of parameters:
- $AIC = -2\ln(L) + 2K$
- However, this penalises all high K models even if sample size is large too.
- Corrected Akaike Information Criterion (**AICc**)
- $AICc = AIC + \frac{2K(K+1)}{n-K-1}$
- Alternatively, there is the Bayesian Information Criterion (**BIC**):

# Information Criterion

- Akaike Information Criterion (**AIC**), penalising number of parameters:
- $AIC = -2\ln(L) + 2K$
- However, this penalises all high K models even if sample size is large too.
- Corrected Akaike Information Criterion (**AICc**)
- $AICc = AIC + \frac{2K(K+1)}{n-K-1}$
- Alternatively, there is the Bayesian Information Criterion (**BIC**):
- $BIC = -2\ln(L) + K\ln(n)$

# Information Criterion

- Akaike Information Criterion (**AIC**), penalising number of parameters:
- $AIC = -2\ln(L) + 2K$
- However, this penalises all high K models even if sample size is large too.
- Corrected Akaike Information Criterion (**AICc**)
- $AICc = AIC + \frac{2K(K+1)}{n-K-1}$
- Alternatively, there is the Bayesian Information Criterion (**BIC**):
- $BIC = -2\ln(L) + K\ln(n)$
- Decision Theory (**DT**) risk minimisation approach.

- What if everything fits poorly?

# Limitations

- What if everything fits poorly?
- Information criterion test relative goodness of fit instead of absolute



# Limitations

- What if everything fits poorly?
- Information criterion test relative goodness of fit instead of absolute
- Parametric Bootstrapping/Posterior Predictive Simulation

# Limitations

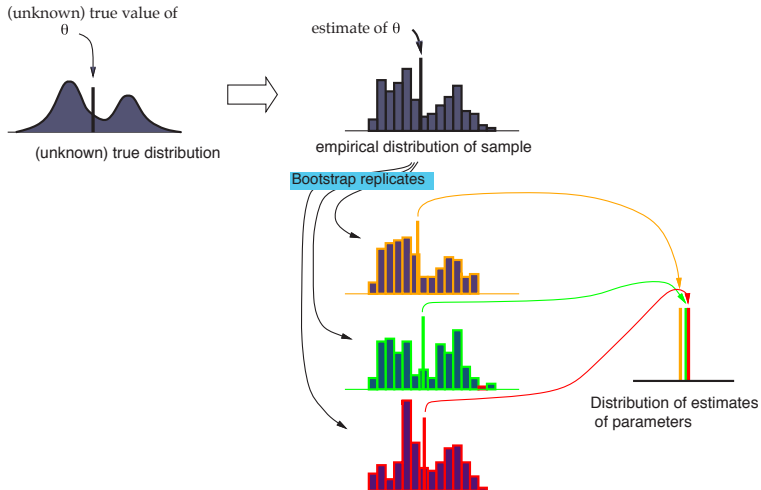
- What if everything fits poorly?
- Information criterion test relative goodness of fit instead of absolute
- Parametric Bootstrapping/Posterior Predictive Simulation
- If the model is reasonable then data simulated under should resemble the empirical data

# Branch Support Testing

---

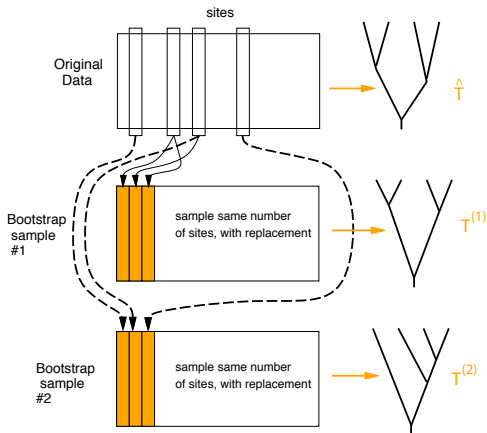
# Bootstrapping in General

## The bootstrap



# Bootstrapping Phylogenies

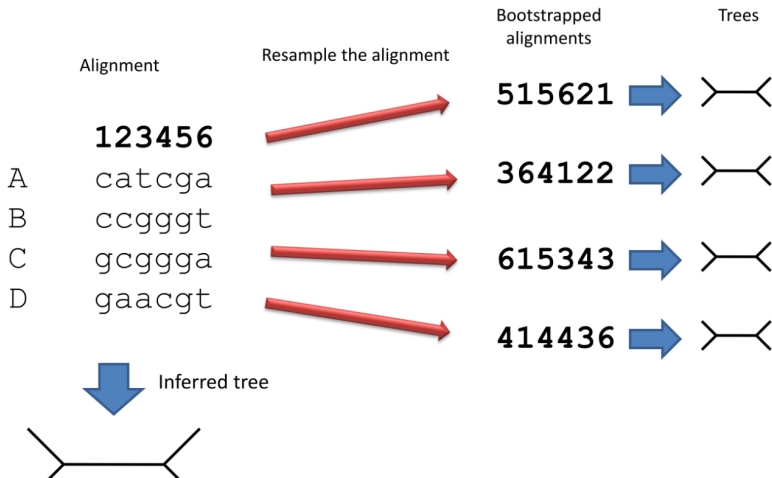
## The bootstrap for phylogenies



Slide from Joe Felsenstein

(and so on)

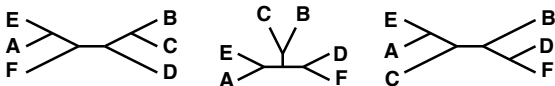
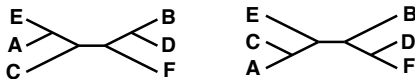
# Bootstrapping Phylogenies



# Bootstrapping Phylogenies

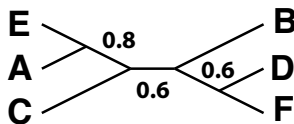
## The majority-rule consensus tree

Trees:

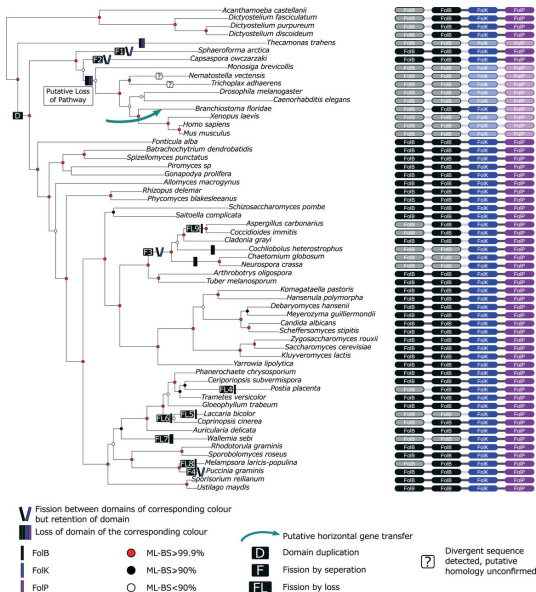


How many times each partition of species is found:

AE   BCDF	4
ACE   BDF	3
ACEF   BD	1
AC   BDEF	1
AEF   BCD	1
ADEF   BC	2
ABCE   DF	3



# Combining the results





# What is the bootstrap doing?

- Randomly reweighing the sites in an alignments

# What is the bootstrap doing?

- Randomly reweighing the sites in an alignments
- Probability of a site being excluded  $1 - \frac{1}{n}$

# What is the bootstrap doing?

- Randomly reweighing the sites in an alignments
- Probability of a site being excluded  $1 - \frac{1}{n}$
- Asymptotically approximately 0.36

# What is the bootstrap doing?

- Randomly reweighing the sites in an alignments
- Probability of a site being excluded  $1 - \frac{1}{n}$
- Asymptotically approximately 0.36
- Goal to simulate an infinite population (number of alignment columns)

- Typically underestimates the true probabilities

# Limitations

- Typically underestimates the true probabilities
- i.e biased but conservative

# Limitations

- Typically underestimates the true probabilities
- i.e biased but conservative
- Computationally demanding

# Limitations

- Typically underestimates the true probabilities
- i.e biased but conservative
- Computationally demanding
- Assumes independence of sites



# Limitations

- Typically underestimates the true probabilities
- i.e biased but conservative
- Computationally demanding
- Assumes independence of sites
- Relies on good input data

# Limitations

- Typically underestimates the true probabilities
- i.e biased but conservative
- Computationally demanding
- Assumes independence of sites
- Relies on good input data
- Only answers to what extent does input data support a given part of the tree

- Simulate data sets of this size assuming the estimate of the tree is the truth
- Key for many more sophisticated tests.
- Can be used to generate  $p$ -values, but non-trivial

- Resampling estimated log-likelihoods (**RELL**)

- Resampling estimated log-likelihoods (**RELL**)
- Instead of re-doing the full ML inference just re-sample the site  $\ln(L)$  values and sum

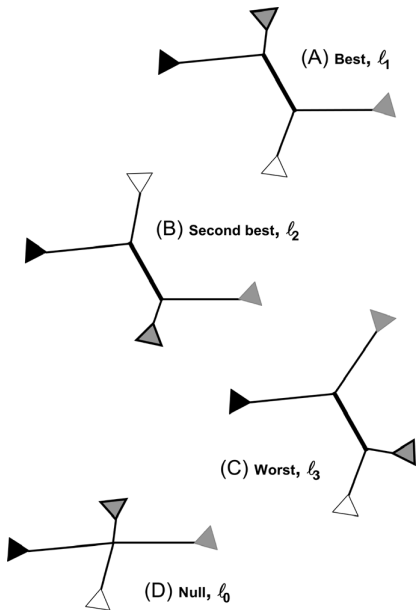
# Alternative Approaches

- Resampling estimated log-likelihoods (**RELL**)
- Instead of re-doing the full ML inference just re-sample the site  $\ln(L)$  values and sum
- Rapid Bootstraps (**RBS**)

# Alternative Approaches

- Resampling estimated log-likelihoods (**RELL**)
- Instead of re-doing the full ML inference just re-sample the site  $\ln(L)$  values and sum
- Rapid Bootstraps (**RBS**)
- Ultrafast Bootstraps (**UFBoot**)

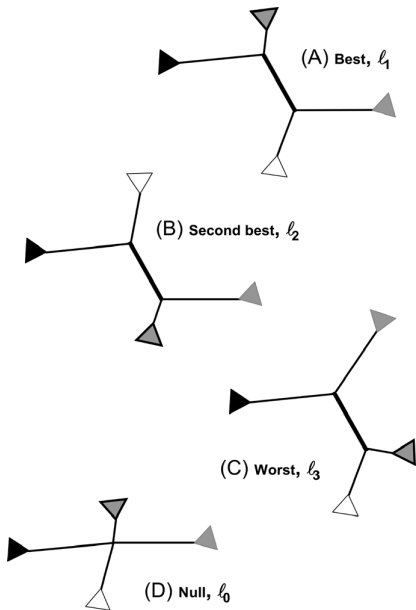
# Likelihood Tests



- Comparing the 3 nearest NNIs to a given branch:

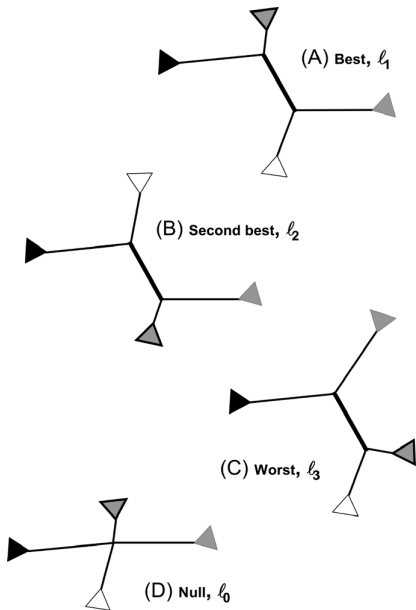


# Likelihood Tests



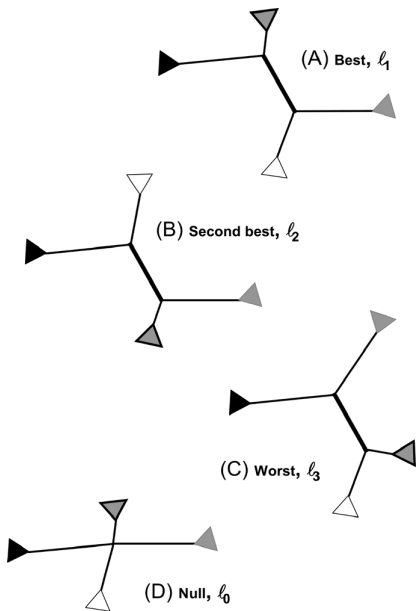
- Comparing the 3 nearest NNIs to a given branch:
- Parametric aLRT:  $\chi^2$  of  $\delta$  for branch vs. closest NNIs

# Likelihood Tests



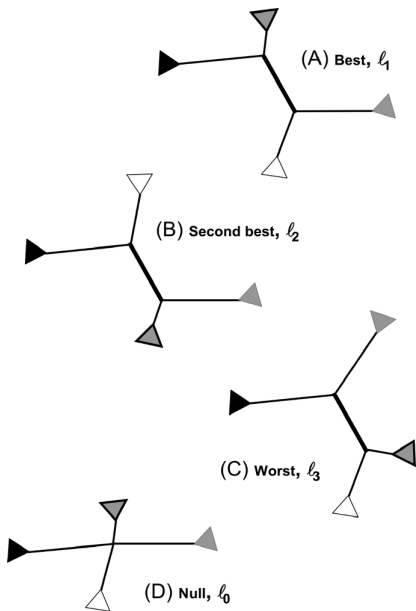
- Comparing the 3 nearest NNIs to a given branch:
- Parametric **aLRT**:  $\chi^2$  of  $\delta$  for branch vs. closest NNIs
- Non-parametric **SH-aLRT** based on RELL

# Likelihood Tests



- Comparing the 3 nearest NNIs to a given branch:
- Parametric **aLRT**:  $\chi^2$  of  $\delta$  for branch vs. closest NNIs
- Non-parametric **SH-aLRT** based on RELL
- **aBayes**:

# Likelihood Tests

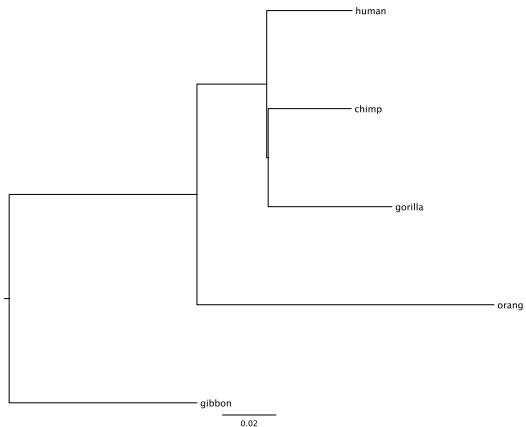


- Comparing the 3 nearest NNIs to a given branch:
- Parametric **aLRT**:  $\chi^2$  of  $\delta$  for branch vs. closest NNIs
- Non-parametric **SH-aLRT** based on RELL
- **aBayes**:
- $P(T_c | X) = \frac{P(X|T_c)P(T_c)}{\sum_{i=0}^2 P(X|T_i)P(T_i)}$  with flat prior  
 $P(T_0) = P(T_1) = P(T_2)$

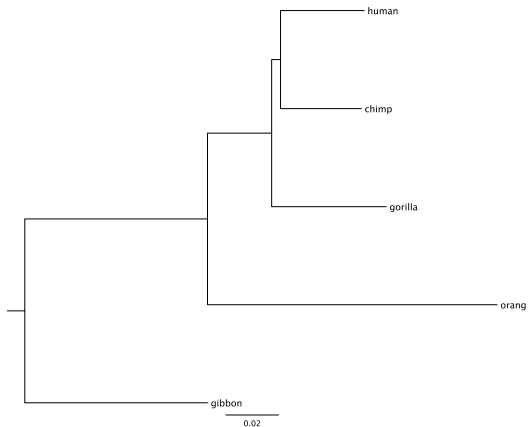
# Comparing Trees

---







# How to compare competing hypotheses?

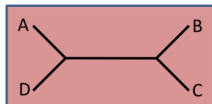


# How to compare competing hypotheses?



# Simplistic Comparison

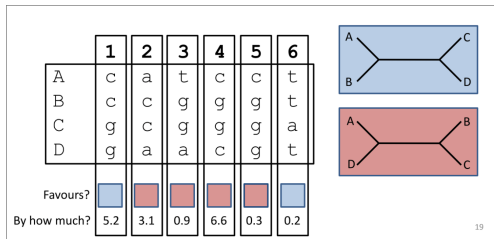
	1	2	3	4	5	6
A	c	a	t	c	c	t
B	c	c	g	g	g	t
C	g	c	g	g	g	a
D	g	a	a	c	g	t
Favours?						
By how much?	5.2	3.1	0.9	6.6	0.3	0.2



19

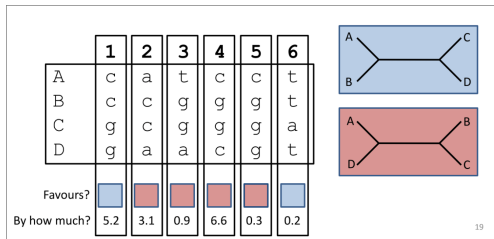


# Qualitative Comparison



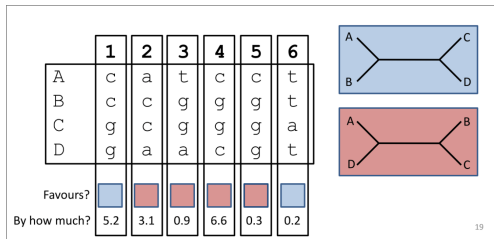
- 4 sites favour the red tree, 2 favour the blue

# Qualitative Comparison



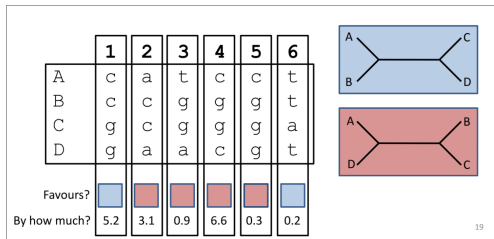
- 4 sites favour the red tree, 2 favour the blue
- $\binom{n}{k} p^k (1-p)^{n-k}$

# Qualitative Comparison



- 4 sites favour the red tree, 2 favour the blue
- $\binom{n}{k} p^k (1-p)^{n-k}$
- 4 out of 6  $p = 0.6875$

# Qualitative Comparison









- 4 sites favour the red tree, 2 favour the blue
- $\binom{n}{k} p^k (1-p)^{n-k}$
- 4 out of 6  $p = 0.6875$
- 40 out of 60  $p = 0.0124$

# Qualitative Comparison

	1	2	3	4	5	6
A	c	a	t	c	c	t
B	c	c	g	g	g	t
C	g	c	g	g	g	a
D	g	a	a	c	g	t
Favours?						
By how much?	5.2	3.1	0.9	6.6	0.3	0.2

- 4 sites favour the red tree, 2 favour the blue
- $\binom{n}{k} p^k (1-p)^{n-k}$
- 4 out of 6  $p = 0.6875$
- 40 out of 60  $p = 0.0124$
- 400 out of 600  $p = 2.3 * 10^{-16}$







# Quantitative Comparison


	1	2	3	4	5	6
A	c	a	t	c	c	t
B	c	c	g	g	g	t
C	g	c	g	g	g	a
D	g	a	a	c	g	t
Favours?						
By how much?	5.2	3.1	0.9	6.6	0.3	0.2

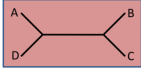
19

$$\bullet \mu = (-5.2 + 3.1 + 0.9 + 6.6 + 0.3 - 0.2)/6 = 0.916$$

# Quantitative Comparison

	1	2	3	4	5	6
A	c	a	t	c	c	t
B	c	c	g	g	g	t
C	g	c	g	g	g	a
D	g	a	a	c	g	t
Favours?						
By how much?	5.2	3.1	0.9	6.6	0.3	0.2

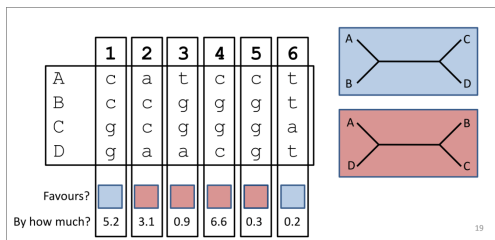




19

- $\mu = (-5.2 + 3.1 + 0.9 + 6.6 + 0.3 - 0.2)/6 = 0.916$
- $\sigma^2 = 15.22$







# Quantitative Comparison





- $\mu = (-5.2 + 3.1 + 0.9 + 6.6 + 0.3 - 0.2)/6 = 0.916$
- $\sigma^2 = 15.22$
- $t = \frac{\mu}{\sigma^2} * \sqrt{N} = 0.148$



# Quantitative Comparison

	1	2	3	4	5	6
A	c	a	t	c	c	t
B	c	c	g	g	g	t
C	g	c	g	g	g	a
D	g	a	a	c	g	t
Favours?						
By how much?	5.2	3.1	0.9	6.6	0.3	0.2



19

- $\mu = (-5.2 + 3.1 + 0.9 + 6.6 + 0.3 - 0.2)/6 = 0.916$
- $\sigma^2 = 15.22$
- $t = \frac{\mu}{\sigma^2} * \sqrt{N} = 0.148$
- therefore:  $p = 0.888$  under  $5d.f.$

- Null: if no sampling error (infinite data)  $T_1$  and  $T_2$  would explain the data equally well.

## More robust approaches

- Null: if no sampling error (infinite data)  $T_1$  and  $T_2$  would explain the data equally well.
- $\delta(X | T_1, T_2) = 2 [\ln L(X | T_1) - \ln L(X | T_2)]$

## More robust approaches

- Null: if no sampling error (infinite data)  $T_1$  and  $T_2$  would explain the data equally well.
- $\delta(X | T_1, T_2) = 2 [\ln L(X | T_1) - \ln L(X | T_2)]$
- Expectation under null  $\mathbb{E}[\delta(X | T_1, T_2)] = 0$

## More robust approaches

- Null: if no sampling error (infinite data)  $T_1$  and  $T_2$  would explain the data equally well.
- $\delta(X | T_1, T_2) = 2 [\ln L(X | T_1) - \ln L(X | T_2)]$
- Expectation under null  $\mathbb{E}[\delta(X | T_1, T_2)] = 0$
- Why can't we just use  $\chi^2$  to get a critical value for  $\delta$ ?

## More robust approaches

- Null: if no sampling error (infinite data)  $T_1$  and  $T_2$  would explain the data equally well.
- $\delta(X | T_1, T_2) = 2 [\ln L(X | T_1) - \ln L(X | T_2)]$
- Expectation under null  $\mathbb{E}[\delta(X | T_1, T_2)] = 0$
- Why can't we just use  $\chi^2$  to get a critical value for  $\delta$ ?
- Tree space is difficult.

- Many avenues:
- Non-parametric bootstrapping
- Parametric bootstrapping
- Related approaches.

- Shimodaira-Hasegawa Test



## Alternative tests

- Shimodaira-Hasegawa Test
- Compares candidate tree sets

## Alternative tests

- Shimodaira-Hasegawa Test
- Compares candidate tree sets
- $H_0$  = all topologies equally good

## Alternative tests

- Shimodaira-Hasegawa Test
- Compares candidate tree sets
- $H_0$  = all topologies equally good
- Very conservative when the number of candidate trees is large

## Alternative tests

- **Shimodaira-Hasegawa Test**
- Compares candidate tree sets
- $H_0 =$  all topologies equally good
- Very conservative when the number of candidate trees is large
- Can be corrected with weighted SH-test overcomes.

- **Shimodaira-Hasegawa Test**
- Compares candidate tree sets
- $H_0$  = all topologies equally good
- Very conservative when the number of candidate trees is large
- Can be corrected with weighted SH-test overcomes.
- **Approximately Unbiased Test**

- **Shimodaira-Hasegawa Test**
- Compares candidate tree sets
- $H_0$  = all topologies equally good
- Very conservative when the number of candidate trees is large
- Can be corrected with weighted SH-test overcomes.
- **Approximately Unbiased Test**
- Achieves weighted by varying bootstrap size for each tree.

## Alternative tests

- **Shimodaira-Hasegawa Test**
- Compares candidate tree sets
- $H_0$  = all topologies equally good
- Very conservative when the number of candidate trees is large
- Can be corrected with weighted SH-test overcomes.
- **Approximately Unbiased Test**
- Achieves weighted by varying bootstrap size for each tree.
- Better for larger comparisons, can have issues with P-space curvature.

## Alternative tests

- **Shimodaira-Hasegawa Test**
- Compares candidate tree sets
- $H_0$  = all topologies equally good
- Very conservative when the number of candidate trees is large
- Can be corrected with weighted SH-test overcomes.
- **Approximately Unbiased Test**
- Achieves weighted by varying bootstrap size for each tree.
- Better for larger comparisons, can have issues with P-space curvature.
- **Swofford–Olsen–Waddell–Hillis** same idea but uses parametric bootstraps instead.



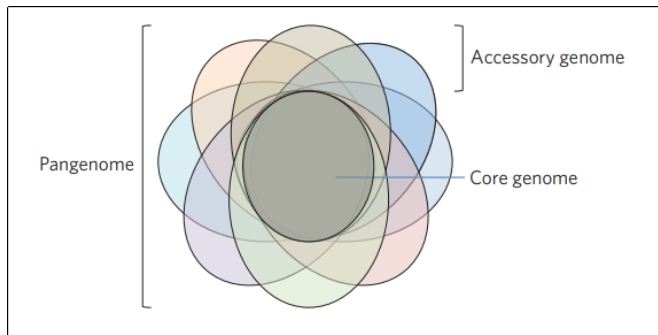
## Alternative tests

- **Shimodaira-Hasegawa Test**
- Compares candidate tree sets
- $H_0$  = all topologies equally good
- Very conservative when the number of candidate trees is large
- Can be corrected with weighted SH-test overcomes.
- **Approximately Unbiased Test**
- Achieves weighted by varying bootstrap size for each tree.
- Better for larger comparisons, can have issues with P-space curvature.
- **Swofford-Olsen-Waddell-Hillis** same idea but uses parametric bootstraps instead.
- Sensitive to model misspecification.

## From A Single Gene to Many Genes

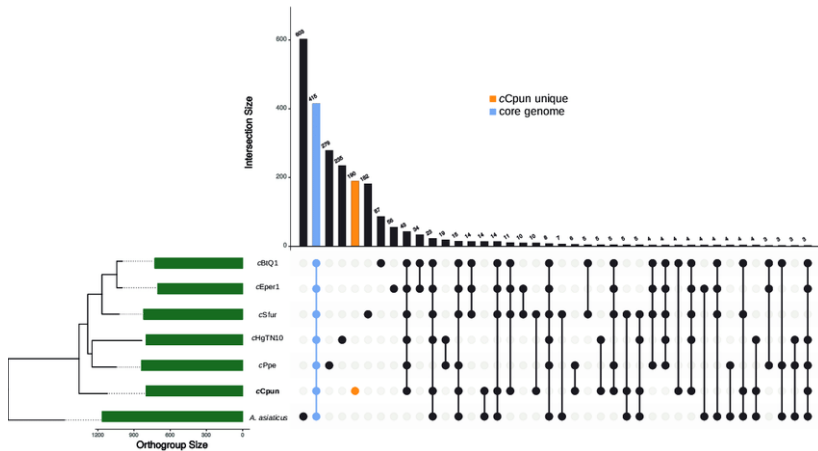
---

# Core vs Pan-Genome

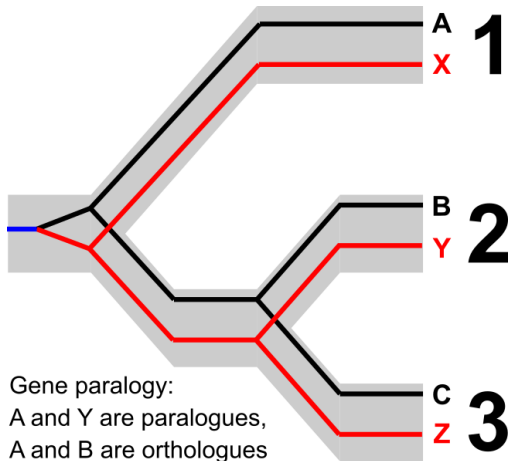


- How do we go from a bunch of individual genes to a species phylogeny?

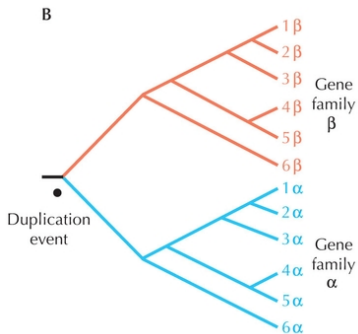
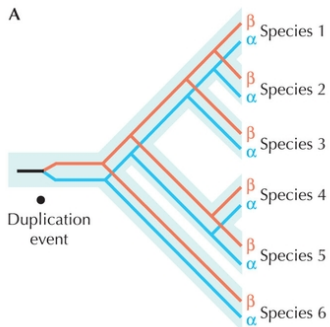
# Venn/Euler plots should be avoided



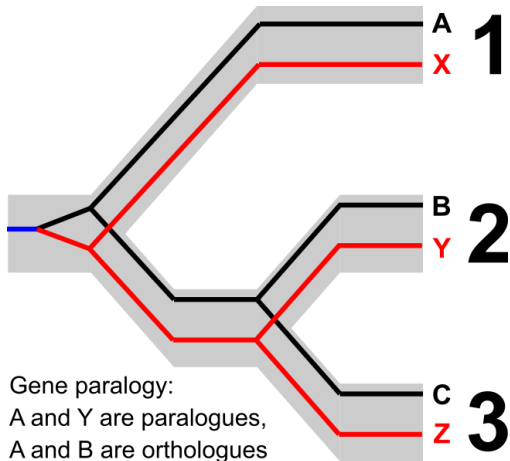
## Why can't just use a single gene tree?



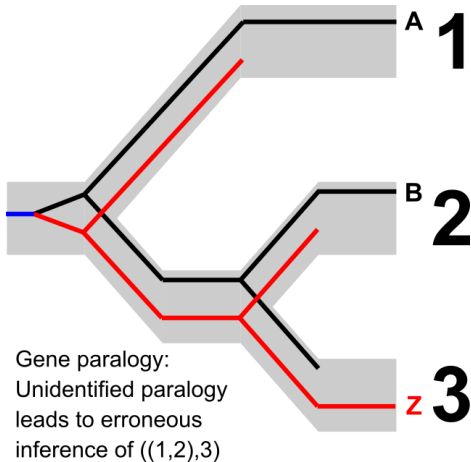
# Paralogy



# Paralogy

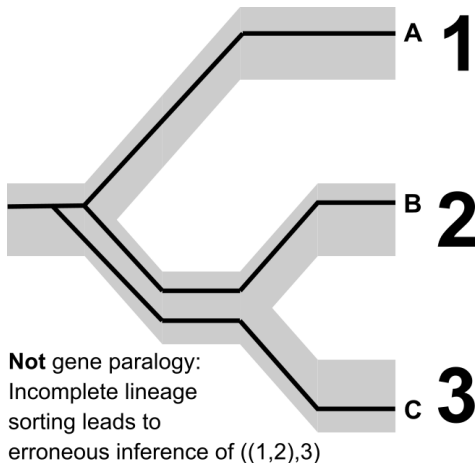


# Hidden Paralogy

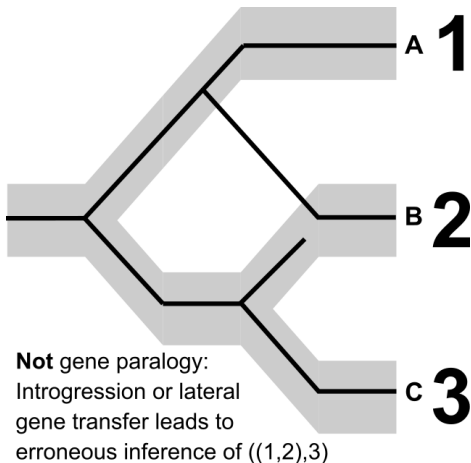




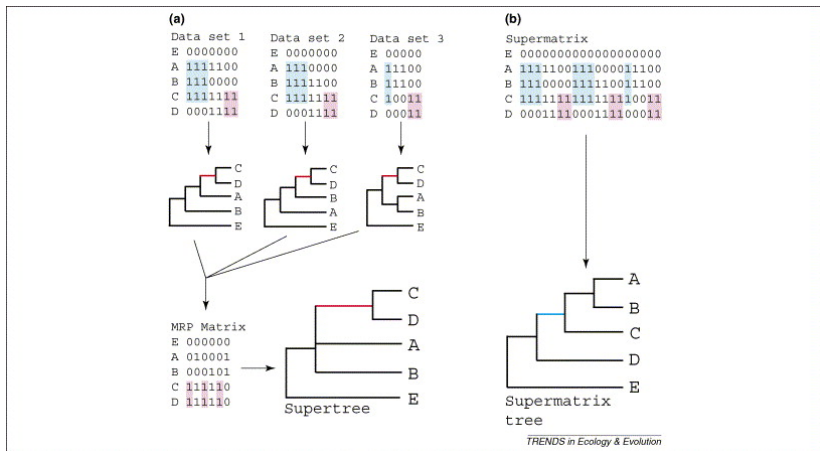
# Incomplete Lineage Sorting



# Lateral/Horizontal Gene Transfer



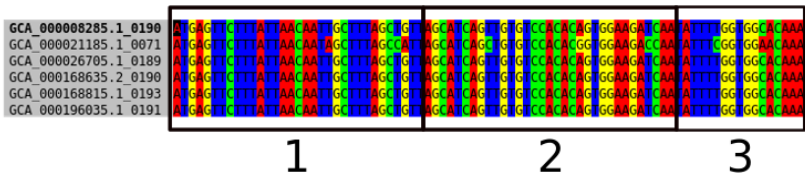
# Supermatrix and SuperTrees



# Supermatrix Evolution Models

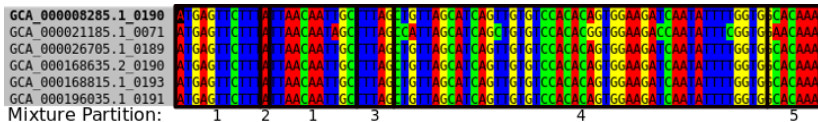
```
GCA_000008285.1_0190 ATGAGTTCCTTATTAAACAATGCCTTAGCTGTAGCAACAGTTGTGCCACAGTGGAAAGATCAATAATTTGGTGGCACAAA
GCA_000021185.1_0071 ATGAGTTCCTTATTAAACAATGCCTTAGCCATTAGCAACAGCTGTGCCACAGGTGGAAAGACCAATAATTCGGTGGAACAAA
GCA_000026705.1_0189 ATGAGTTCCTTATTAAACAATGCCTTAGCTGTAGCAACAGTTGTGCCACAGTGGAAAGATCAATAATTTGGTGGCACAAA
GCA_000168635.2_0190 ATGAGTTCCTTATTAAACAATGCCTTAGCTGTAGCAACAGTTGTGCCACAGTGGAAAGATCAATAATTTGGTGGCACAAA
GCA_000168815.1_0193 ATGAGTTCCTTATTAAACAATGCCTTAGCTGTAGCAACAGTTGTGCCACAGTGGAAAGATCAATAATTTGGTGGCACAAA
GCA_000196035.1_0191 ATGAGTTCCTTATTAAACAATGCCTTAGCTGTAGCAACAGTTGTGCCACAGTGGAAAGATCAATAATTTGGTGGCACAAA
```

# Gene Partitions





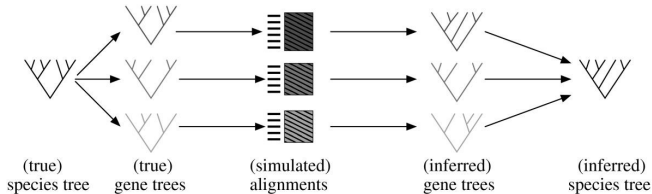
# Mixture Models



- Can be very computationally demanding
- Variable evolutionary rates/models combined (one-size fits all)
- CAT model does offer a solution but can overfit.
- More robust to random error when phylogenetic signal is consistent.

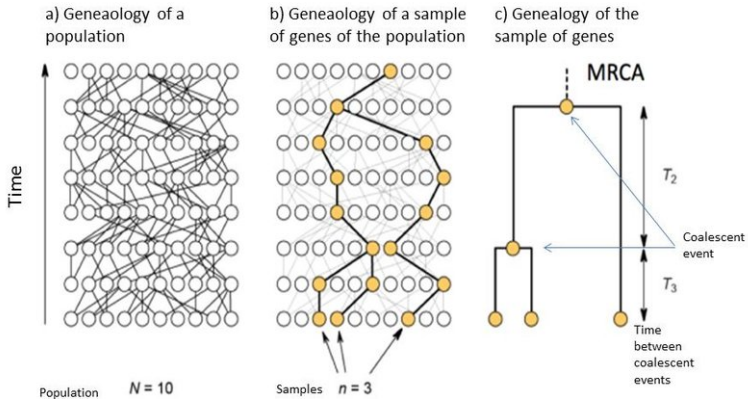


# SuperTree



- Allows reconciliation of partial overlaps (i.e. not just core genome)
- Faster/more tractable
- Observed to have lower accuracy generally but more robust to incongruent signal (i.e. frequent HGT).

# Coalescent Theory

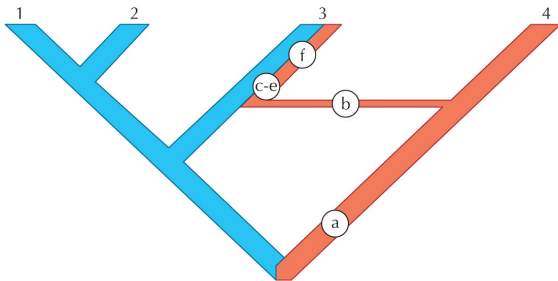


# Genomic Epidemiology

## Phylogenetics

---

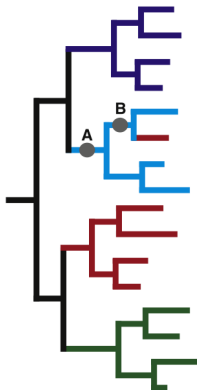
# Identifying HGT with phylogenies



**FIGURE 27.30.** Stages in lateral gene transfer (LGT). The evolution of four species and one example of LGT are shown. Some key steps in LGT are labeled: (a) Divergence of genomes of different lineages; (b) movement of DNA from one lineage to another; (c) maintenance and replication of the foreign DNA; (d) possible positive selection for the foreign DNA; (e) spread into the new species' population; (f) amelioration. (Modified from Penny D. and Poole A. *Curr. Opin. Genet. Dev.* **9**: 672–677, © 1999 Elsevier.)

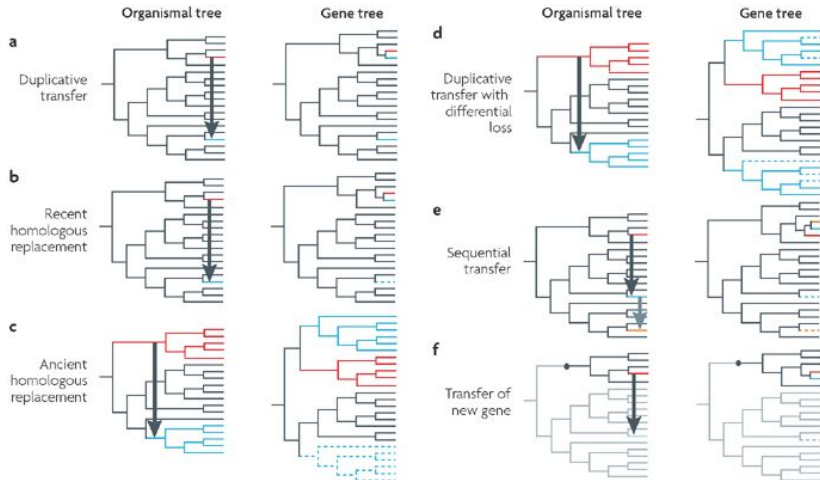
*Evolution* © 2008 Cold Spring Harbor Laboratory Press

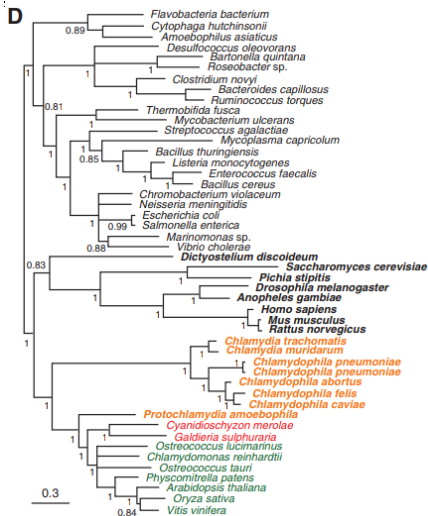
# Schematic HGT Tree



- (A) Need strong branch support for recipient branching with donor lineages
- (B) Need strong support for recipient branching within donor lineages

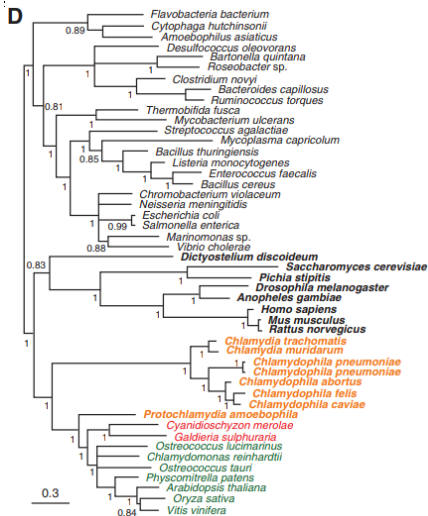
# Different Types result in different topologies



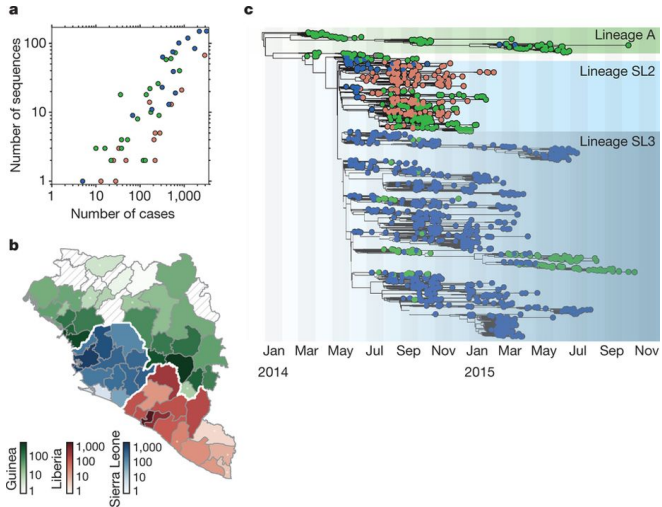




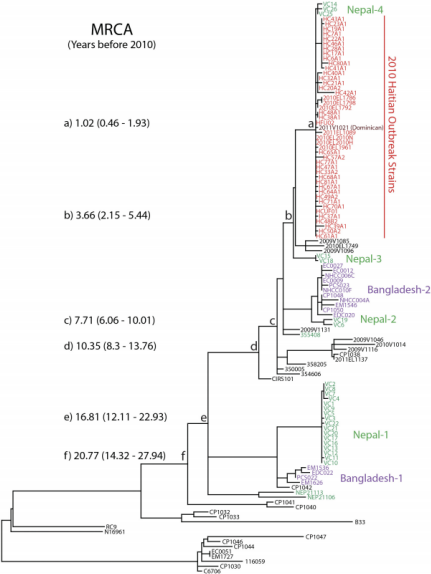




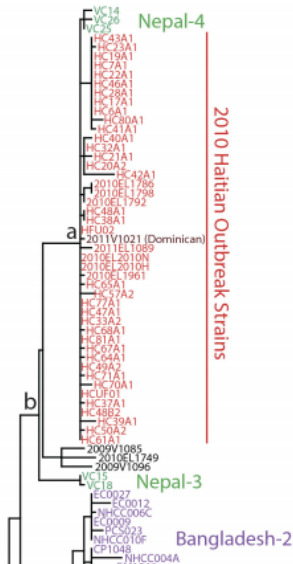
# West African Ebola Epidemic



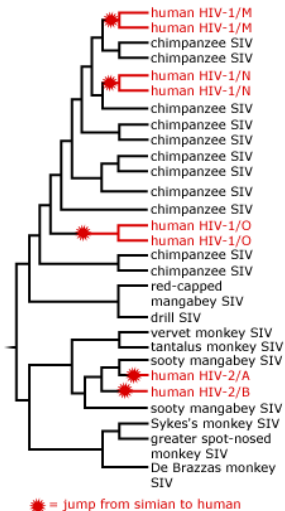
# Haitian Cholera Outbreak: Identifying the origin



# Haitian Cholera Outbreak: Identifying the origin



# Finding zoonoses



# Conclusion

---

# Summary

- Phylogenies are hypothesis and their inference includes assumptions that need testing.

# Summary

- Phylogenies are hypothesis and their inference includes assumptions that need testing.
- Bootstrapping is a slow, biased but conservative way to estimate the support for a given branch in your tree.



# Summary

- Phylogenies are hypothesis and their inference includes assumptions that need testing.
- Bootstrapping is a slow, biased but conservative way to estimate the support for a given branch in your tree.
- Comparing trees directly is non-trivial due to tree-space.

# Summary

- Phylogenies are hypothesis and their inference includes assumptions that need testing.
- Bootstrapping is a slow, biased but conservative way to estimate the support for a given branch in your tree.
- Comparing trees directly is non-trivial due to tree-space.
- Supermatrix and supertree approaches allow reconciliation of data from multiple genes.

# Summary

- Phylogenies are hypothesis and their inference includes assumptions that need testing.
- Bootstrapping is a slow, biased but conservative way to estimate the support for a given branch in your tree.
- Comparing trees directly is non-trivial due to tree-space.
- Supermatrix and supertree approaches allow reconciliation of data from multiple genes.
- Incongruence between the species tree and the gene tree might be evidence for HGT.

# Summary

- Phylogenies are hypothesis and their inference includes assumptions that need testing.
- Bootstrapping is a slow, biased but conservative way to estimate the support for a given branch in your tree.
- Comparing trees directly is non-trivial due to tree-space.
- Supermatrix and supertree approaches allow reconciliation of data from multiple genes.
- Incongruence between the species tree and the gene tree might be evidence for HGT.
- Phylogenies can be used to trace outbreak origins and parameters.

Questions?