

# Lecture:

# Algorithmic Bioinformatics

Doctoral School, Université Dauphine, 2022

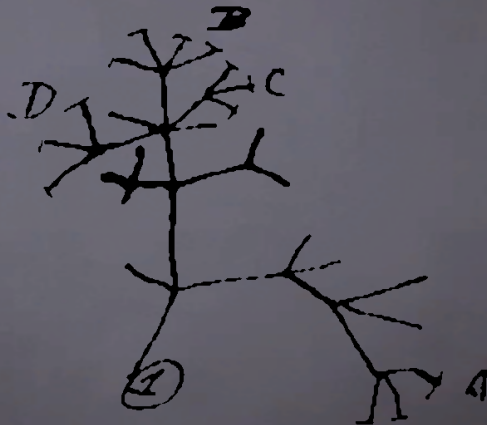
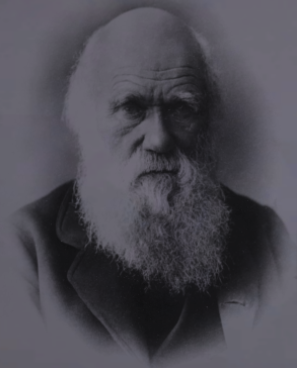


**Université  
Gustave Eiffel**

Some slides graciously provided by Daniel Huson & Celine Scornavacca

# Phylogenetic Trees - Motivation

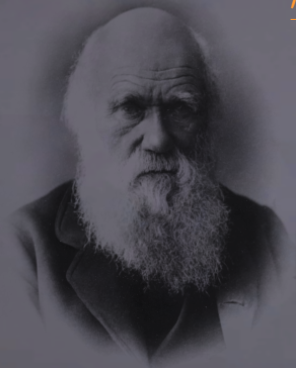
*I think*



# Phylogenetic Trees - Motivation

## Motivation

- study relation between species
- evolution of characteristics
- co-evolution (host-parasite)
- geological migration
- genetic development of viruses/diseases



# Phylogenetic Trees - Motivation

## Motivation

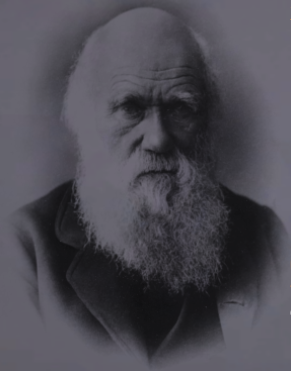
- study relation between species
- evolution of characteristics
- co-evolution (host-parasite)
- geological migration
- genetic development of viruses/diseases

## Evolution

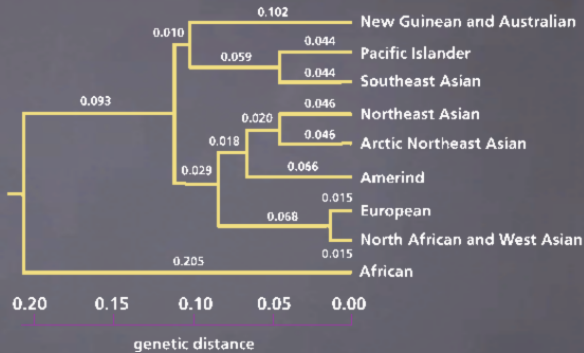
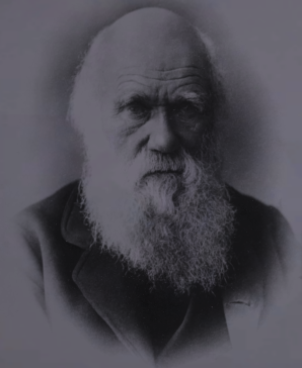
Genetic material changes over time

~> new species "Branch off"

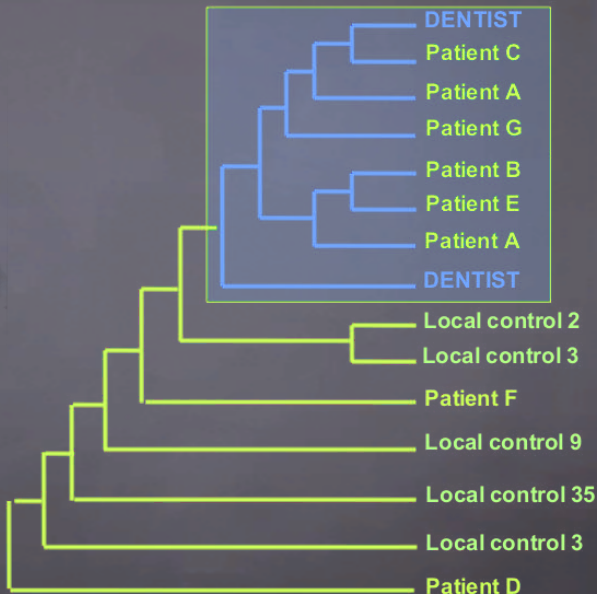
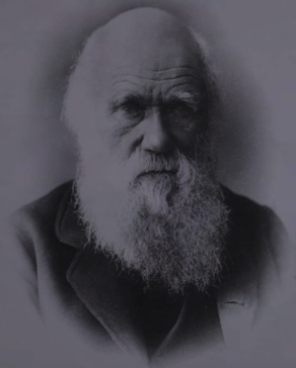
~> "tree of life"



# Phylogenetic Trees - Motivation

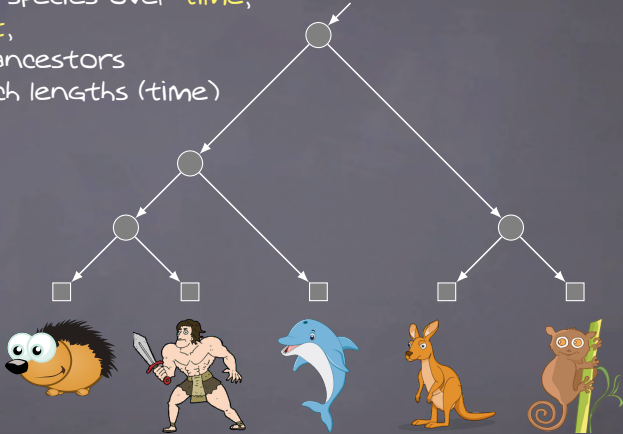


# Phylogenetic Trees - Motivation



# Rooted Phylogenetic Trees

evolution of species over **time**,  
leaves **extant**,  
hypothetical ancestors  
possibly branch lengths (time)



## Notation

taxon, cluster, triplet

# Rooted Phylogenetic Trees

evolution of species over **time**,  
leaves **extant**,  
hypothetical ancestors  
possibly branch lengths (time)



## Notation

taxon, cluster, triplet



# Rooted Phylogenetic Trees

evolution of species over **time**,  
leaves **extant**,  
hypothetical ancestors  
possibly Branch lengths (time)



## Notation

taxon, cluster, triplet

# Rooted Phylogenetic Trees

evolution of species over **time**,  
leaves **extant**,  
hypothetical ancestors  
possibly branch lengths (time)

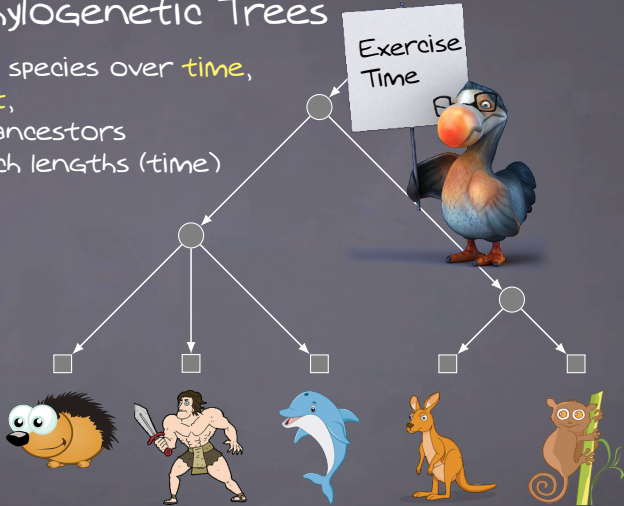


## Notation

taxon, cluster, triplet

# Rooted Phylogenetic Trees

evolution of species over **time**,  
leaves **extant**,  
hypothetical ancestors  
possibly branch lengths (time)



## Notation

taxon, cluster, triplet

## "Polytomies"

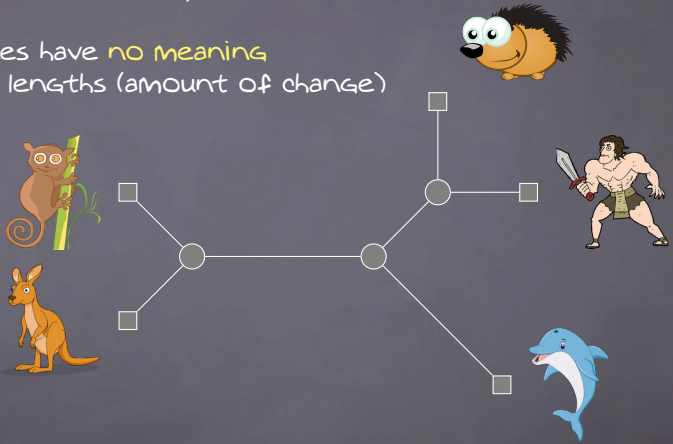
history not clear  $\rightsquigarrow$  "soft"  
known "fan out"  $\rightsquigarrow$  "hard"

## Exercise:

use  $xy|z \leftrightarrow LCA(xy) < LCA(xz) = LCA(xyz)$  to prove  $ab|c + b|cd \rightarrow ac|d$

# Unrooted Phylogenetic Trees

similarity between genomes,  
leaves **extant**,  
internal vertices have **no meaning**  
possibly branch lengths (amount of change)

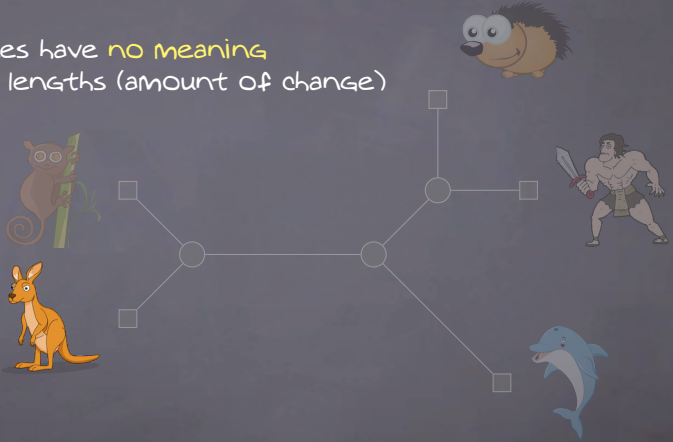


## Notation

taxon, split, quartet

# Unrooted Phylogenetic Trees

similarity between genomes,  
leaves **extant**,  
internal vertices have **no meaning**  
possibly branch lengths (amount of change)

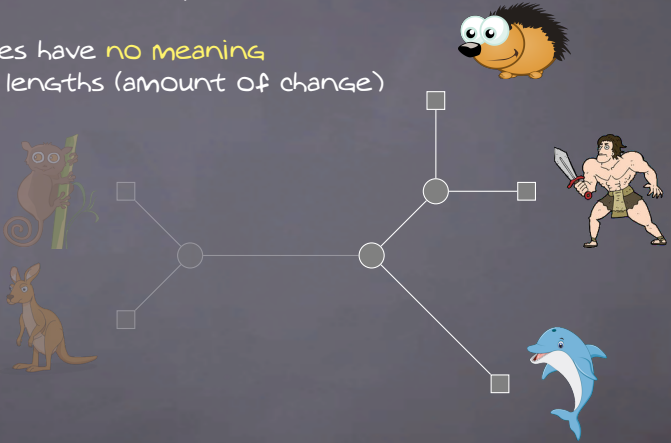


## Notation

taxon, split, quartet

# Unrooted Phylogenetic Trees

similarity between genomes,  
leaves **extant**,  
internal vertices have **no meaning**  
possibly branch lengths (amount of change)

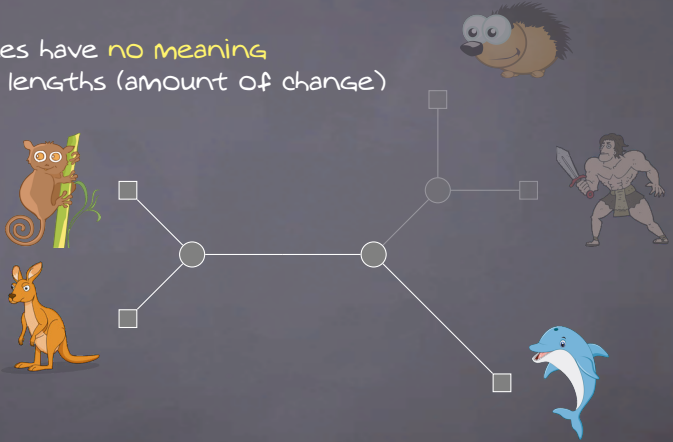


## Notation

taxon, split, quartet

# Unrooted Phylogenetic Trees

similarity between genomes,  
leaves **extant**,  
internal vertices have **no meaning**  
possibly branch lengths (amount of change)

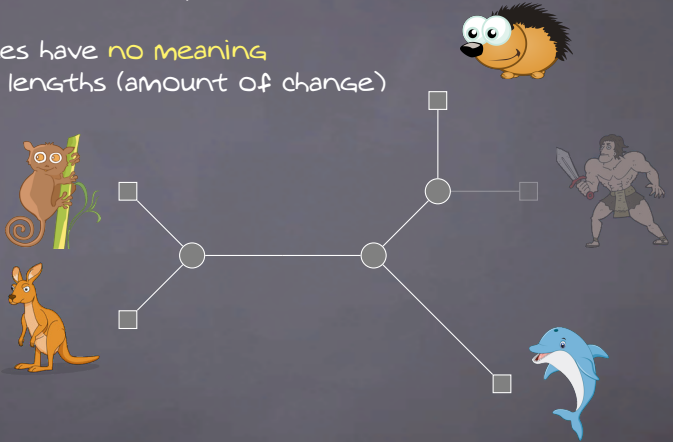


## Notation

taxon, split, quartet

# Unrooted Phylogenetic Trees

similarity between genomes,  
leaves **extant**,  
internal vertices have **no meaning**  
possibly branch lengths (amount of change)



## Notation

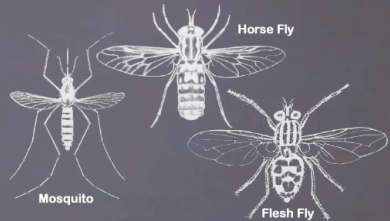
taxon, split, quartet



# Reconstructing Phylogenetic Trees

## Group Species By...

- morphology
- behavior
- geography

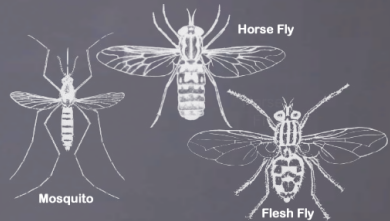


Diptera = two wings

# Reconstructing Phylogenetic Trees

## Group Species By...

- morphology
- behavior
- geography
- distance of sequences
- "genetic distance"
- etc.



Diptera = two wings

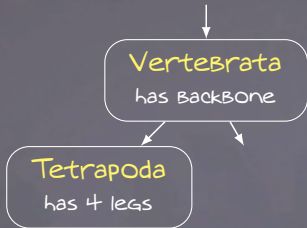
# Reconstructing Phylogenetic Trees



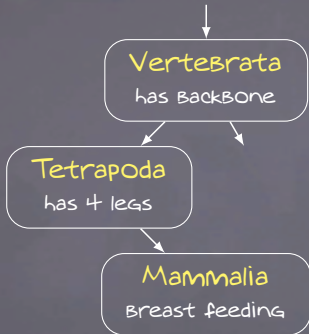
Vertebrata

has Backbone

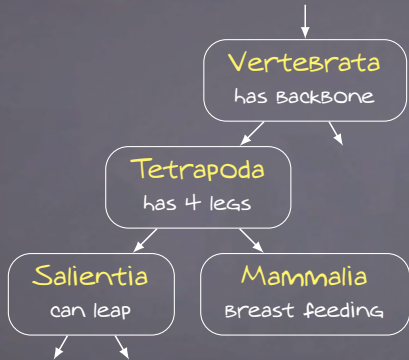
# Reconstructing Phylogenetic Trees



# Reconstructing Phylogenetic Trees



# Reconstructing Phylogenetic Trees



# Reconstructing Phylogenetic Trees



# Reconstructing Phylogenetic Trees





# Reconstructing Phylogenetic Trees

## Keep in Mind

Automatic reconstruction should be

- fast  
(deal with tons of species & genes)
- consistent  
(optimal data  $\leadsto$  correct tree)
- non-arbitrary



# Distance-Based Reconstruction

Idea: cluster hierarchically



# Distance-Based Reconstruction

Idea: cluster hierarchically



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters



# Distance-Based Reconstruction

Idea: cluster hierarchically

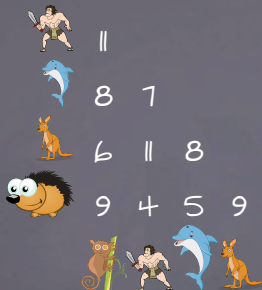


Idea: merge closest clusters



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

update matrix

$$d_{X \cup Y, Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric





# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

update matrix

$$d_{X \cup Y, Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

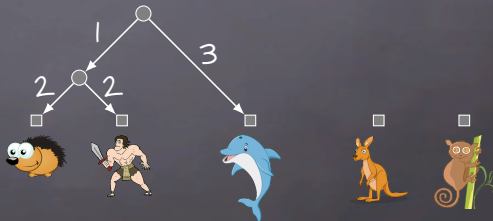
update matrix

$$d_{X \cup Y, Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

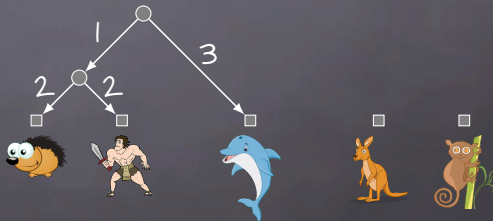
update matrix

$$d_{X \cup Y, Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

$\leadsto$  ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

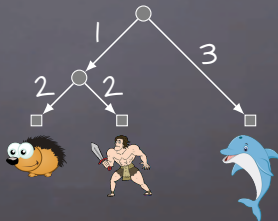
update matrix

$$d_{XUY,Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

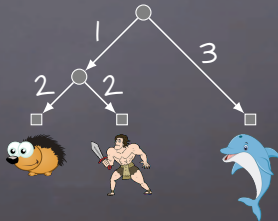
update matrix

$$d_{XUY,Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



Idea: merge closest clusters

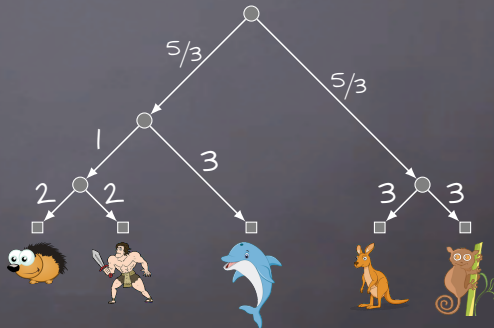
update matrix

$$d_{XUY,Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically

## Unweighted Pair Group Method w/ Avg.

- find "closest pair"
- "join" them
- update distances & recurse



Idea: merge closest clusters

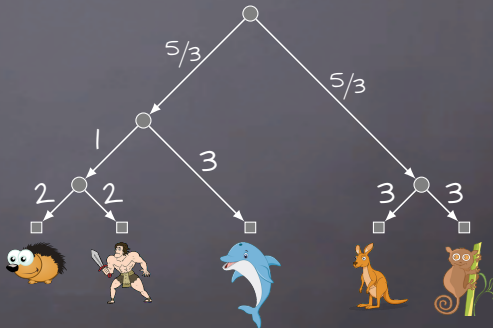
update matrix

$$d_{XUY,Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



## Unweighted Pair Group Method w/ Avg.

- find "closest pair"
- "join" them
- update distances & recurse

Idea: merge closest clusters

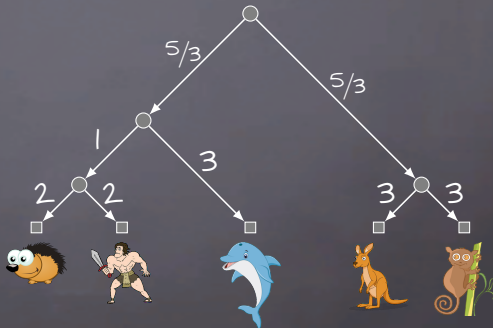
update matrix

$$d_{XUY,Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

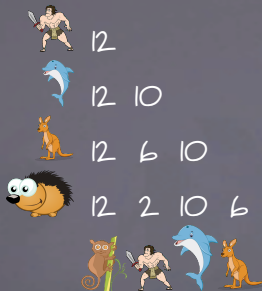
~> ultrametric





# Distance-Based Reconstruction

Idea: cluster hierarchically



## Unweighted Pair Group Method w/ Avg.

- find "closest pair"
- "join" them
- update distances & recurse

Idea: merge closest clusters

update matrix

$$d_{X \cup Y, Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

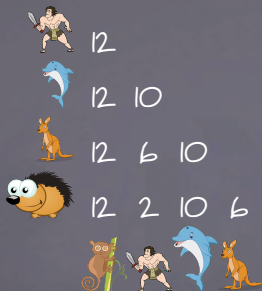
~> ultrametric

Exercise  
Time



# Distance-Based Reconstruction

Idea: cluster hierarchically



## Unweighted Pair Group Method w/ Avg.

- find "closest pair"
- "join" them
- update distances & recurse

Idea: merge closest clusters

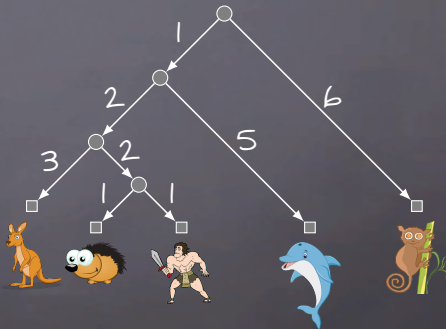
## update matrix

$$d_{XUY,Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

## Branch lengths ??

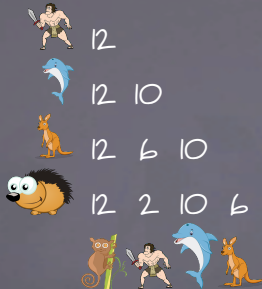
assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

Idea: cluster hierarchically



## Unweighted Pair Group Method w/ Avg.

- find "closest pair"
- "join" them
- update distances & recurse
- **only accurate if ultrametric**

Idea: merge closest clusters

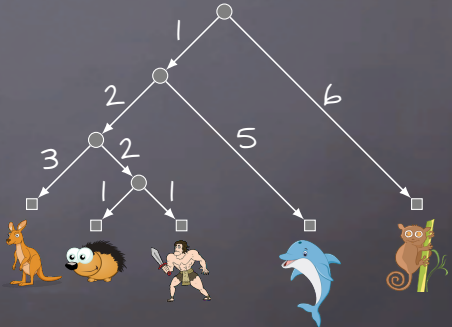
update matrix

$$d_{XUY,Z} = \frac{|X|d_{X,Z} + |Y|d_{Y,Z}}{|X| + |Y|}$$

Branch lengths ??

assume molecular clock

~> ultrametric



# Distance-Based Reconstruction

What about **unrooted** trees? No root  $\rightsquigarrow$  no molecular clock...

# Distance-Based Reconstruction

B	9		
C	11	8	
D	9	12	10
	A	B	C

What about **unrooted** trees? No root  $\rightsquigarrow$  no molecular clock...

# Distance-Based Reconstruction

Problem: correct pairs may not be closest

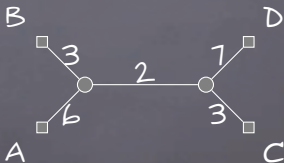
B	9		
C	11	8	
D	9	12	10
	A	B	C

What about **unrooted** trees? No root  $\leadsto$  no molecular clock...

# Distance-Based Reconstruction

Problem: correct pairs may not be closest

B	9		
C	11	8	
D	9	12	10
A	B	C	



# Distance-Based Reconstruction

Problem: correct pairs may not be closest

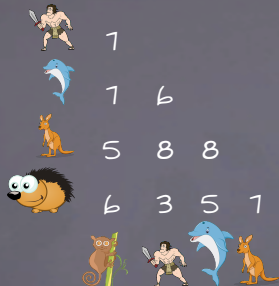
## Neighbor Joining (unrooted)

- Build **eccentricity** matrix:  
$$Q_{x,y} = \sum_z (d_{x,z} + d_{y,z} - d_{x,y}) + 2d_{x,y}$$
- find max in  $Q$
- join them
- update distances & recurse



# Distance-Based Reconstruction

Problem: correct pairs may not be closest



## Neighbor Joining (unrooted)

- Build **eccentricity** matrix:  
$$Q_{x,y} = \sum_z (d_{x,z} + d_{y,z} - d_{x,y}) + 2d_{x,y}$$
- find max in  $Q$
- join them
- update distances & recurse

# Distance-Based Reconstruction

Problem: correct pairs may not be closest



## Neighbor Joining (unrooted)

- Build **eccentricity** matrix:

$$Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$$

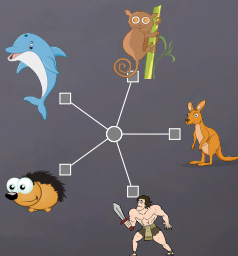
- find max in  $Q$

- join them

- update distances  $\neq$  recurse









## Theorem

$Q_{X,Y}$  max  $\Leftrightarrow$  any tree  $T$  yielding  $Q$  has "cherry"  $(X, Y)$



# Distance-Based Reconstruction

Problem: correct pairs may not be closest

	1			
	1	6		
	5	8	8	
	6	3	5	1
				
				
				
				

## Neighbor Joining (unrooted)

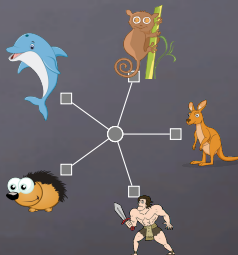
- Build **eccentricity** matrix:  
 $Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$
- find max in  $Q$
- join them
- update distances  $\neq$  recurse

## update distances

$$d_{X \cup Y, Z} = 1/2 (d_{X,Z} + d_{Y,Z} - d_{X,Y})$$









## Branch lengths

$$2b(X) = \frac{\sum_Z (d_{X,Z} - d_{Y,Z} + d_{X,Y})}{n-2}$$



# Distance-Based Reconstruction

Problem: correct pairs may not be closest

	1			
	1	6		
	5	8	8	
	6	3	5	1
				
				
				
				

## Neighbor Joining (unrooted)

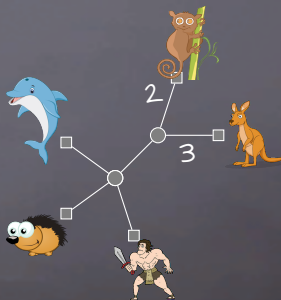
- Build **eccentricity** matrix:  
 $Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$
- find max in  $Q$
- join them
- update distances  $\neq$  recurse

## update distances

$$d_{X \cup Y, Z} = 1/2 (d_{X,Z} + d_{Y,Z} - d_{X,Y})$$

## Branch lengths

$$2b(X) = \frac{\sum_Z (d_{X,Z} - d_{Y,Z} + d_{X,Y})}{n-2}$$



# Distance-Based Reconstruction

Problem: correct pairs may not be closest



## Neighbor Joining (unrooted)

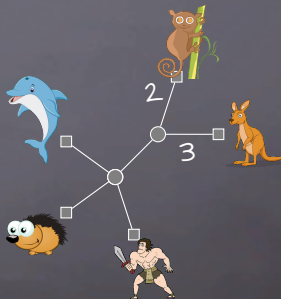
- Build **eccentricity** matrix:  
 $Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$
- find max in  $Q$
- join them
- update distances  $\neq$  recurse

## update distances

$$d_{X \cup Y, Z} = 1/2 (d_{X,Z} + d_{Y,Z} - d_{X,Y})$$

## Branch lengths

$$2b(X) = \frac{\sum_Z (d_{X,Z} - d_{Y,Z} + d_{X,Y})}{n-2}$$



# Distance-Based Reconstruction

Problem: correct pairs may not be closest



18



19 18



18 20 18



## Neighbor Joining (unrooted)

- Build **eccentricity** matrix:

$$Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$$

- find max in  $Q$

- join them

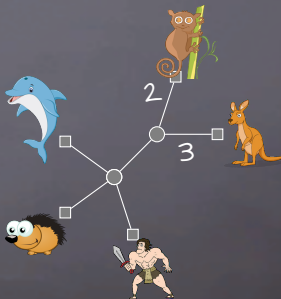
- update distances  $\neq$  recurse

## update distances

$$d_{X \cup Y, Z} = 1/2 (d_{X,Z} + d_{Y,Z} - d_{X,Y})$$

## Branch lengths

$$2b(X) = \frac{\sum_Z (d_{X,Z} - d_{Y,Z} + d_{X,Y})}{n-2}$$



# Distance-Based Reconstruction

Problem: correct pairs may not be closest



## Neighbor Joining (unrooted)

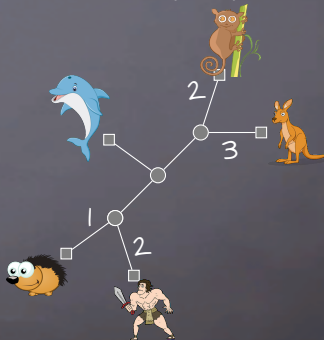
- Build **eccentricity** matrix:  
 $Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$
- find max in  $Q$
- join them
- update distances  $\neq$  recurse

## update distances

$$d_{X \cup Y, Z} = 1/2 (d_{X,Z} + d_{Y,Z} - d_{X,Y})$$

## Branch lengths

$$2b(X) = \frac{\sum_Z (d_{X,Z} - d_{Y,Z} + d_{X,Y})}{n-2}$$



# Distance-Based Reconstruction

Problem: correct pairs may not be closest



## Neighbor Joining (unrooted)

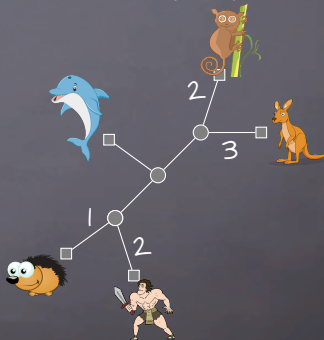
- Build **eccentricity** matrix:  
$$Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$$
- find max in  $Q$
- join them
- update distances  $\neq$  recurse

## update distances

$$d_{X \cup Y, Z} = 1/2 (d_{X,Z} + d_{Y,Z} - d_{X,Y})$$

## Branch lengths

$$2b(X) = \frac{\sum_Z (d_{X,Z} - d_{Y,Z} + d_{X,Y})}{n-2}$$





# Distance-Based Reconstruction

Problem: correct pairs may not be closest



## Neighbor Joining (unrooted)

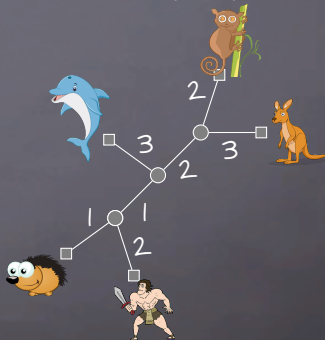
- Build **eccentricity** matrix:  
$$Q_{X,Y} = \sum_Z (d_{X,Z} + d_{Y,Z} - d_{X,Y}) + 2d_{X,Y}$$
- find max in  $Q$
- join them
- update distances  $\neq$  recurse

## update distances

$$d_{X \cup Y, Z} = 1/2 (d_{X,Z} + d_{Y,Z} - d_{X,Y})$$

## Branch lengths

$$2b(X) = \frac{\sum_Z (d_{X,Z} - d_{Y,Z} + d_{X,Y})}{n-2}$$



# Parsimony Reconstructing

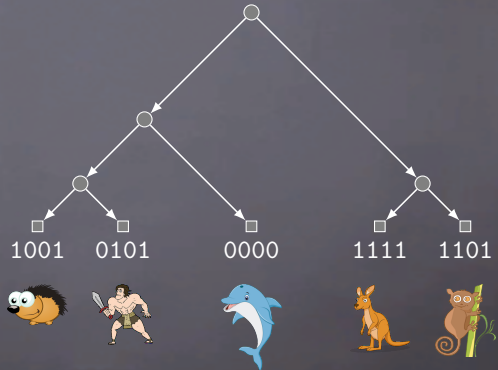


# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓

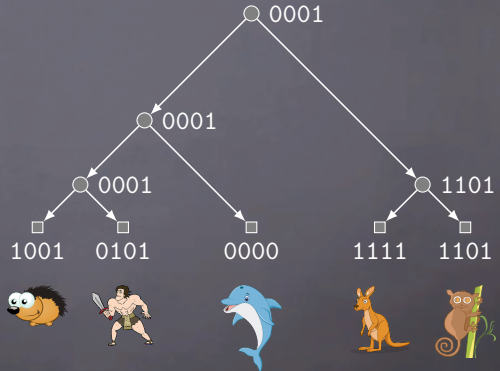
# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓



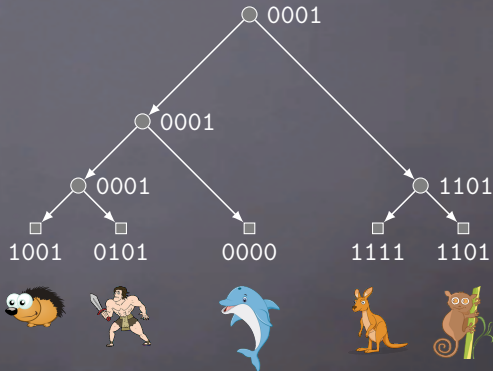
# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓



# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓



sum "distance" of endpoints of each edge



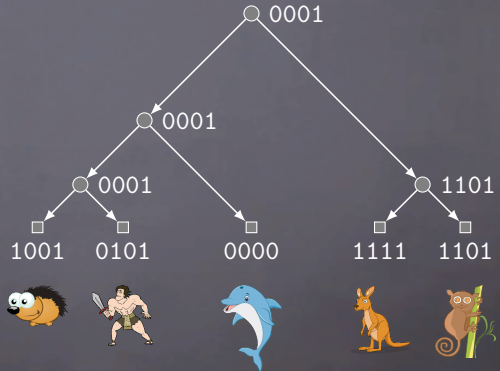
# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓

## Small Parsimony

Input: character state matrix  $M$ ,  
rooted tree  $T$

Task: assign characters to internal  
nodes minimizing **total cost**





# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓

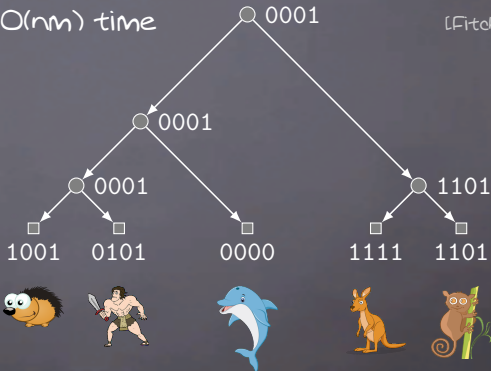
## Small Parsimony

Input: character state matrix  $M$ ,  
rooted tree  $T$

Task: assign characters to internal  
nodes minimizing total cost

$\leadsto O(nm)$  time

[Fitch, '71]



ALGO (simple DFS):  
Backtrack from  $v$ :  
if  $v$  is leaf:  
 $\leadsto R(v) = \text{label}(v)$   
else if  $R(u) \cap R(w) = \emptyset$ :  
 $\leadsto R(v) = R(u) \cup R(w)$   
else:  
 $\leadsto R(v) = R(u) \cap R(w)$

# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓

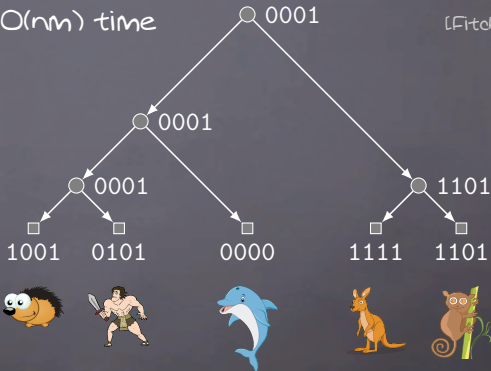
## Small Parsimony

Input: character state matrix  $M$ ,  
rooted tree  $T$

Task: assign characters to internal  
nodes minimizing **total cost**

$\leadsto O(nm)$  time

[Fitch, '71]



**ALGO** (simple DFS):  
Backtrack from  $v$ :  
if  $v$  is leaf:  
     $\leadsto R(v) = \text{label}(v)$   
else if  $R(u) \cap R(w) = \emptyset$ :  
     $\leadsto R(v) = R(u) \cup R(w)$   
else:  
     $\leadsto R(v) = R(u) \cap R(w)$

# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓

## Small Parsimony

Input: character state matrix  $M$ ,  
rooted tree  $T$

Task: assign characters to internal  
nodes minimizing **total cost**

$\leadsto O(nm)$  time

[Fitch, '71]

## Large Parsimony

Input: character state matrix  $M$

Task: find tree  $T$  & assign  
characters to internal nodes  
minimizing total cost

**ALGO** (simple DFS):  
Backtrack from  $v$ :  
if  $v$  is leaf:  
     $\leadsto R(v) = \text{label}(v)$   
else if  $R(u) \cap R(w) = \emptyset$ :  
     $\leadsto R(v) = R(u) \cup R(w)$   
else:  
     $\leadsto R(v) = R(u) \cap R(w)$

# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓

## Small Parsimony

Input: character state matrix  $M$ ,  
rooted tree  $T$

Task: assign characters to internal  
nodes minimizing **total cost**

~  $O(nm)$  time

[Fitch, '71]

## Large Parsimony

Input: character state matrix  $M$

Task: find tree  $T$  & assign  
characters to internal nodes  
minimizing total cost

~ NP-hard

**ALGO** (simple DFS):  
Backtrack from  $v$ :  
if  $v$  is leaf:  
    ~  $R(v) = \text{label}(v)$   
else if  $R(u) \cap R(w) = \emptyset$ :  
    ~  $R(v) = R(u) \cup R(w)$   
else:  
    ~  $R(v) = R(u) \cap R(w)$

# Parsimony Reconstructing

	fur	ALS	pouch	land
	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✗	✗
	✓	✓	✓	✓
	✓	✗	✗	✓

## Small Parsimony

Input: character state matrix  $M$ ,  
rooted tree  $T$

Task: assign characters to internal  
nodes minimizing **total cost**

~  $O(nm)$  time

[Fitch, '71]

## Large Parsimony

Input: character state matrix  $M$

Task: find tree  $T$  & assign  
characters to internal nodes  
minimizing total cost

~ NP-hard

Note: alignment is crucial!

**Algo** (simple DFS):  
Backtrack from  $v$ :  
if  $v$  is leaf:  
    ~  $R(v) = \text{label}(v)$   
else if  $R(u) \cap R(w) = \emptyset$ :  
    ~  $R(v) = R(u) \cup R(w)$   
else:  
    ~  $R(v) = R(u) \cap R(w)$

# Maximum Likelihood Reconstruction

Idea: find a tree (with branch lengths) under which evolution is most likely to have produced the observed characters/genomes

# Maximum Likelihood Reconstruction

Idea: find a tree (with branch lengths) under which evolution is most likely to have produced the observed characters/genomes  
~> need model of evolution

# Maximum Likelihood Reconstruction

Idea: find a tree (with branch lengths) under which evolution is most likely to have produced the observed characters/genomes  
~> need model of evolution

## Jukes & Cantor Model

- each Base evolves individually
- each Base occurs with equal frequency in the genome
- constant rate  $\mu$  of mutation
- each Base is equally likely to be result of mutation



# Maximum Likelihood Reconstruction

Idea: find a tree (with branch lengths) under which evolution is most likely to have produced the observed characters/genomes  
~> need model of evolution

## Jukes & Cantor Model

- each Base evolves individually
- each Base occurs with equal frequency in the genome
- constant rate  $\mu$  of mutation
- each Base is equally likely to be result of mutation

## Generalized Time Reversible Model

- each Base evolves individually
- each Base  $X$  has a frequency  $\pi_X$  to occur in the genome
- each Base-substitution has its own rate of occurrence

# Maximum Likelihood Reconstruction

Idea: find a tree (with branch lengths) under which evolution is most likely to have produced the observed characters/genomes  
~> need model of evolution

## Jukes & Cantor Model

- each base evolves individually
- each base occurs with equal frequency in the genome
- constant rate  $\mu$  of mutation
- each base is equally likely to be result of mutation

## Generalized Time Reversible Model

- each base evolves individually
- each base  $X$  has a frequency  $\pi_X$  to occur in the genome
- each base-substitution has its own rate of occurrence

compute likelihood, given tree & parameters ~>  $O(mn)$  time

# Maximum Likelihood Reconstruction

Idea: find a tree (with branch lengths) under which evolution is most likely to have produced the observed characters/genomes  
~> need model of evolution

## Jukes & Cantor Model

- each Base evolves individually
- each Base occurs with equal frequency in the genome
- constant rate  $\mu$  of mutation
- each Base is equally likely to be result of mutation

## Generalized Time Reversible Model

- each Base evolves individually
- each Base  $X$  has a frequency  $\pi_X$  to occur in the genome
- each Base-substitution has its own rate of occurrence

compute likelihood, given tree & parameters ~>  $O(mn)$  time  
find best tree & parameters ~> NP-hard

# Maximum Likelihood Reconstruction

Idea: find a tree (with branch lengths) under which evolution is most likely to have produced the observed characters/genomes  
~> need model of evolution

## Jukes & Cantor Model

- each Base evolves individually
- each Base occurs with equal frequency in the genome
- constant rate  $\mu$  of mutation
- each Base is equally likely to be result of mutation

## Generalized Time Reversible Model

- each Base evolves individually
- each Base  $X$  has a frequency  $\pi_X$  to occur in the genome
- each Base-substitution has its own rate of occurrence

compute likelihood, given tree & parameters ~>  $O(mn)$  time

find best tree & parameters ~> NP-hard

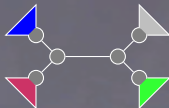
~> local search in the tree space

# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration  
of 4 (3) "neighboring"  
subtrees into another



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration  
of 4 (3) "neighboring"  
subtrees into another



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration  
of 4 (3) "neighboring"  
subtrees into another



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$





# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

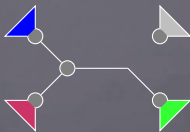
## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

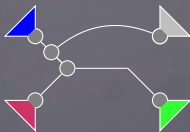
## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

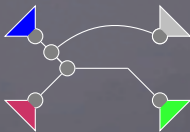
## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



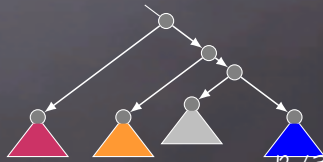
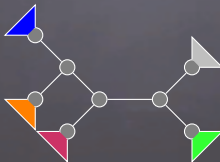
## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



## Tree Bisection & Reconnection

Break any edge  $\neq$   
insert a new  
reconnecting edge  
"Between" any 2 edges



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



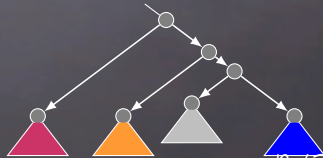
## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



## Tree Bisection & Reconnection

Break any edge  $\neq$   
insert a new  
reconnecting edge  
"Between" any 2 edges



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



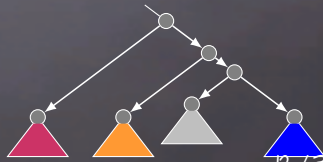
## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



## Tree Bisection & Reconnection

Break any edge  $\neq$   
insert a new  
reconnecting edge  
"Between" any 2 edges



# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

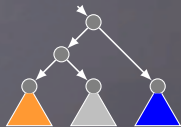
## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



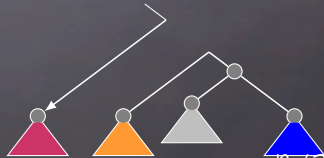
## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



## Tree Bisection & Reconnection

Break any edge  $\neq$   
insert a new  
reconnecting edge  
"Between" any 2 edges





# ML Reconstruction - Tree Spaces

Observe: a tree and a rearrangement operation span a space

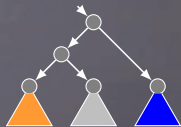
## Nearest Neighbor Interchange

change any configuration of 4 (3) "neighboring" subtrees into another



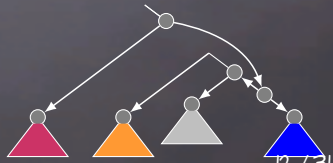
## Subtree Prune & Regraft

Break any edge  $uv \neq$   
connect  $v$  to any edge of  
the component of  $u$



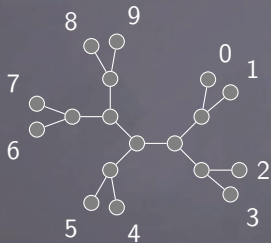
## Tree Bisection & Reconnection

Break any edge  $\neq$   
insert a new  
reconnecting edge  
"Between" any 2 edges



# ML Reconstruction - Tree Spaces

Exercise: turn into (any) caterpillar:



Exercise: how are the distances related?



# Checking ROBustness – Bootstrap Method

suppose: method  $X$  yields tree  $T$  from  $n \times m$  character-state matrix  $M$   
repeat  $k$  times the following experiment:

1. draw  $m$  columns from  $M$  (with repetition)
2. use  $X$  to compute  $T_i$

Finally, for each branch of  $T$ , check how often it occurs in the  $T_i$

~> "Bootstrap value" measures robustness ("support") of each branch

# Reconstruction By Gene Trees

# Reconstruction By Gene Trees

## A Common Method For Reconstructing Trees

1. get genomes of multiple species
2. extract "genes" using START & STOP codons
3. cluster genes in "**families**" of similar genes
4. within each family, infer a "**gene tree**" using dissimilarities
5. Build a consensus among the gene trees  $\rightsquigarrow$  "**species tree**"  
(**Note:** species tree may differ significantly from individual gene trees)
6. **reconcile** all gene trees with the species tree to learn the evolution of those genes

# Reconstruction By Gene Trees

## A Common Method For Reconstructing Trees

1. ~~get genomes of multiple species~~
2. ~~extract "genes" using START & STOP codons~~
3. cluster genes in "**families**" of similar genes
4. within each family, infer a "**gene tree**" using dissimilarities
5. Build a consensus among the gene trees  $\rightsquigarrow$  "**species tree**"  
(**Note:** species tree may differ significantly from individual gene trees)
6. **reconcile** all gene trees with the species tree to learn the evolution of those genes

# Reconstruction By Gene Trees

## A Common Method For Reconstructing Trees

1. ~~get genomes of multiple species~~
2. ~~extract "genes" using START & STOP codons~~
3. ~~cluster genes in "families" of similar genes~~
4. within each family, infer a "gene tree" using dissimilarities
5. Build a consensus among the gene trees  $\rightsquigarrow$  "species tree"  
(Note: species tree may differ significantly from individual gene trees)
6. **reconcile** all gene trees with the species tree to learn the evolution of those genes

# Reconstruction By Gene Trees

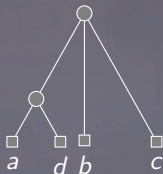
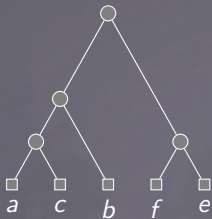
## A Common Method For Reconstructing Trees

1. ~~get genomes of multiple species~~
2. ~~extract "genes" using START & STOP codons~~
3. ~~cluster genes in "families" of similar genes~~
4. ~~within each family, infer a "gene tree" using dissimilarities~~
5. ~~Build a consensus among the gene trees  $\rightsquigarrow$  "species tree"~~  
(Note: species tree may differ significantly from individual gene trees)
6. ~~reconcile all gene trees with the species tree to learn the evolution of those genes~~



# Supertrees - "Build" Algorithm

Idea: find root partition  $\nVdash$  recurse (as long as there are  $\geq 3$  taxa)



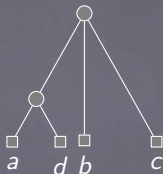
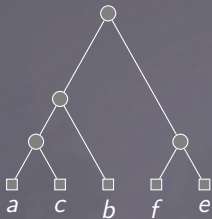
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



b

c

d

a

e

f

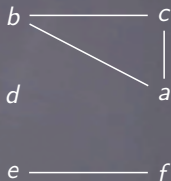
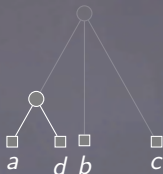
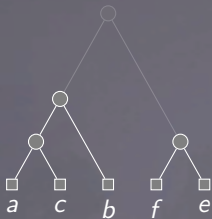
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



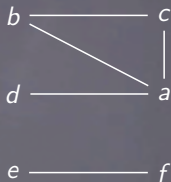
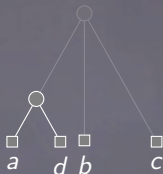
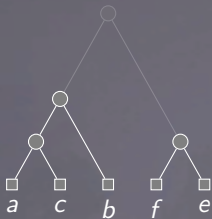
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



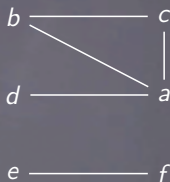
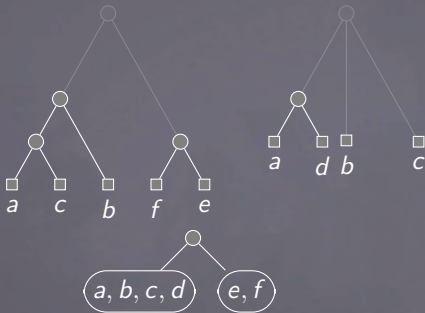
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



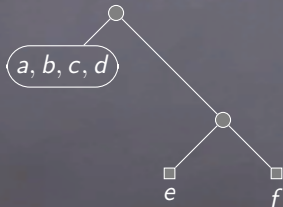
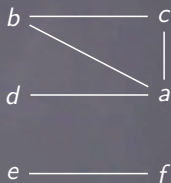
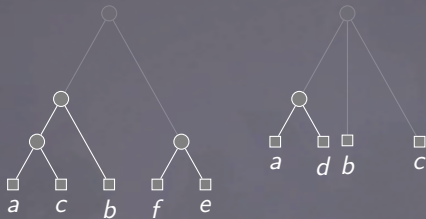
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



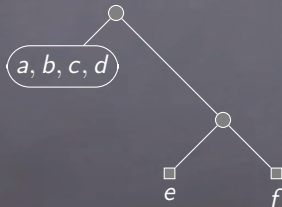
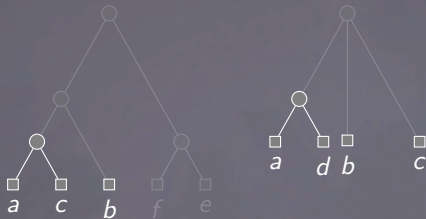
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



b

d

e

c

a

f

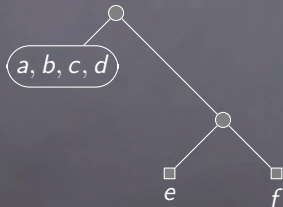
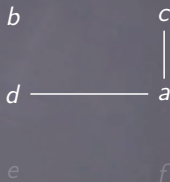
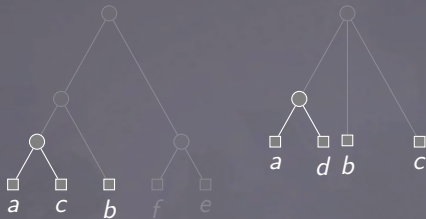
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



## ALGO

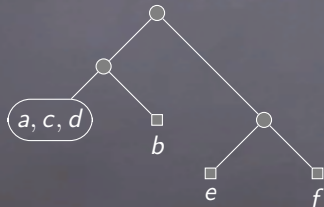
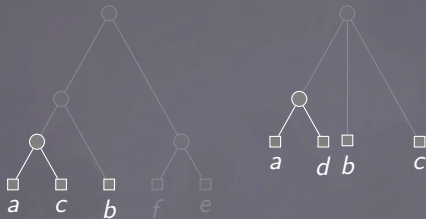
[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root



# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



b

c

d

a

e

f

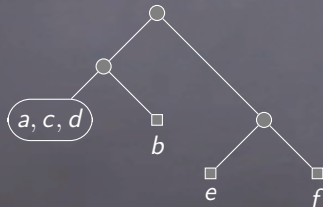
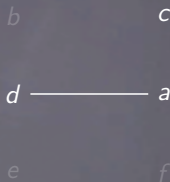
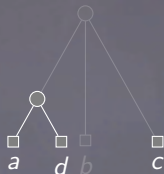
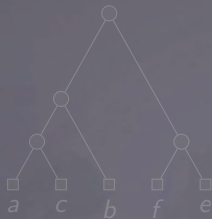
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



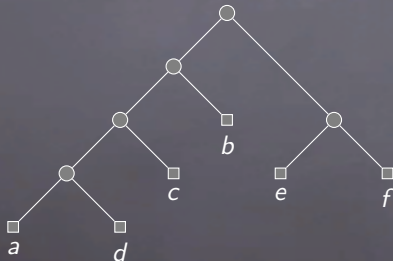
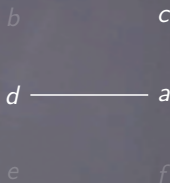
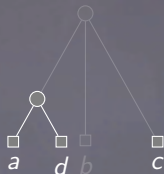
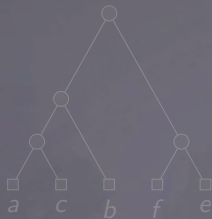
## ALGO

[Aho et al'81]

1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\nVdash$  recurse (as long as there are  $\geq 3$  taxa)



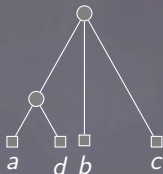
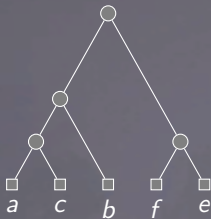
## ALGO

[Aho et al'81]

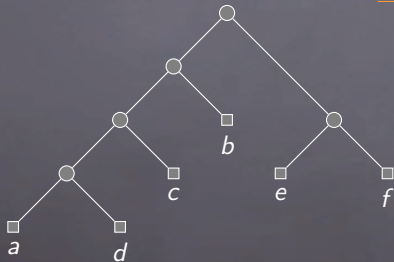
1. Build graph  $G$  with edge  $uv \Leftrightarrow \exists uv|x$
2. recurse for each component of  $G$
3. plug subtrees to root

# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)

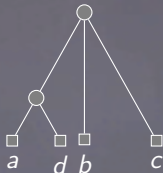
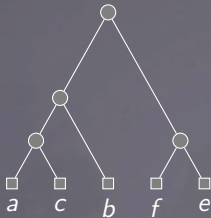


Note: always works if trees are **compatible**

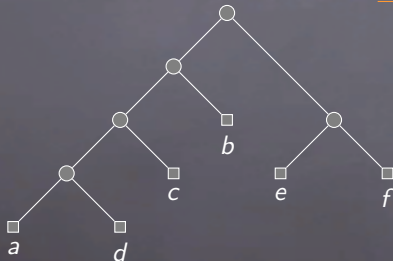


# Supertrees - "Build" Algorithm

Idea: find root partition  $\neq$  recurse (as long as there are  $\geq 3$  taxa)



Note: always works if trees are **compatible**



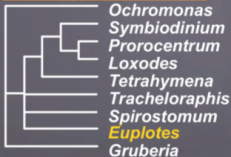
incompatible?

- largest compatible subset  
     $\rightsquigarrow$  NP-hard (even for triplets)
- voting schemes  
    (each tree votes for their clades)
- reinterpret clades as characters,  
    combine into matrix  $\neq$  reconstruct

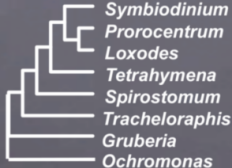
# Consensi of Non-Agreeing Trees



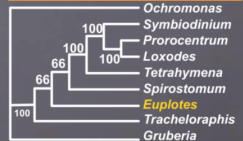
## strict consensus



## consensus subtree



## majority consensus



# Reconstruction By Gene Trees

## A Common Method For Reconstructing Trees

1. ~~get genomes of multiple species~~
2. ~~extract "genes" using START & STOP codons~~
3. ~~cluster genes in "families" of similar genes~~
4. ~~within each family, infer a "gene tree" using dissimilarities~~
5. ~~Build a consensus among the gene trees  $\rightsquigarrow$  "species tree"~~  
(Note: species tree may differ significantly from individual gene trees)
6. ~~reconcile all gene trees with the species tree to learn the evolution of those genes~~

# Reconstruction By Gene Trees

## A Common Method For Reconstructing Trees

1. ~~get genomes of multiple species~~
2. ~~extract "genes" using START & STOP codons~~
3. ~~cluster genes in "families" of similar genes~~
4. ~~within each family, infer a "gene tree" using dissimilarities~~
5. ~~build a consensus among the gene trees  $\rightsquigarrow$  "species tree"~~  
(~~Note:~~ species tree may differ significantly from individual gene trees)
6. **reconcile** all gene trees with the species tree to learn the evolution of those genes



# The History of a Gene Family

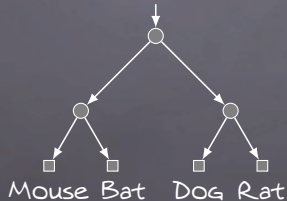
## Recall

gene = "functional element" of DNA, clustered into gene-families

Mouse	G	G	A	G	C	T	T	G	A	G	C	C	G	G	A	A	T	A	G	T	A	G	C	A	A	C	A	T	C	T	T	T	A	A	T	T	C	G	A	G	C									
Dog	G	G	A	A	T	C	T	G	A	A	C	A	G	G	C	T	T	A	G	T	A	G	C	C	A	C	T	A	G	A	A	T	A	A	G	A	C	T	T	T	T	A	A	T	T	C	G	A	G	C
Bat	G	G	A	A	T	T	T	G	A	A	C	A	G	G	T	T	T	A	G	T	A	G	C	C	A	C	T	A	G	A	A	T	A	A	G	A	C	T	C	T	T	A	A	T	T	C	G	A	G	C
Rat	G	G	A	A	T	T	T	G	A	A	C	C	G	G	C	C	T	C	G	T	A	G	C	A	A	C	A	A	G	A	A	T	A	A	G	C	T	T	A	T	T	A	A	T	C	C	G	T	G	C

each family yields a tree depicting its history  $\rightsquigarrow$  "gene tree"  
consensus of the gene trees yields "species tree"

But: what did really happen???



# The History of a Gene Family

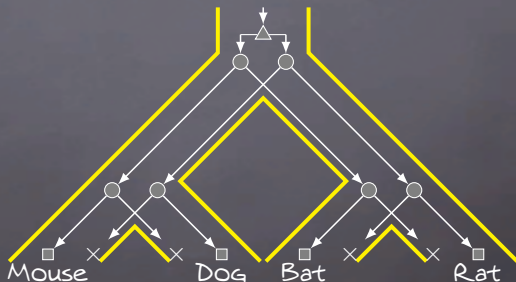
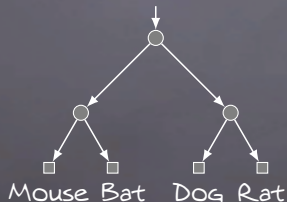
## Recall

gene = "functional element" of DNA, clustered into gene-families

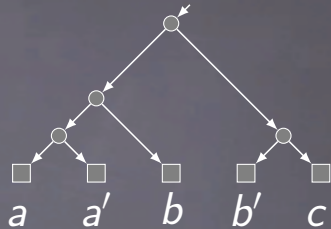
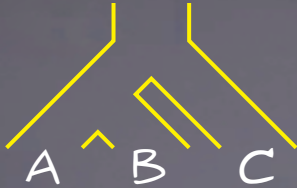
Mouse	G	G	A	G	C	T	T	G	A	G	C	C	G	G	A	A	T	A	G	T	A	G	C	A	A	C	A	T	C	T	T	T	A	A	G	A	A	T	T	T	A	A	T	T	C	G	A	G	C	
Dog	G	G	A	A	T	C	T	G	A	A	C	A	G	G	C	T	T	A	G	T	A	G	C	C	A	C	T	A	G	A	A	T	A	A	G	A	C	T	T	T	A	A	T	T	C	G	A	G	C	
Bat	G	G	A	A	T	T	T	G	A	A	C	A	G	G	T	T	T	A	G	T	A	G	C	C	A	C	T	A	G	A	A	T	A	A	G	A	C	T	C	T	T	A	A	T	T	C	G	A	G	C
Rat	G	G	A	A	T	T	T	G	A	A	C	C	G	G	C	C	T	C	G	T	A	G	C	A	A	C	A	A	G	A	A	T	A	A	G	C	T	T	A	T	T	A	A	T	C	C	G	T	G	C

each family yields a tree depicting its history  $\leadsto$  "gene tree"  
consensus of the gene trees yields "species tree"

But: what did really happen???



# Reconciliation

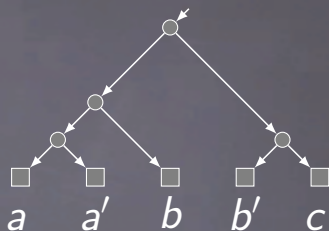
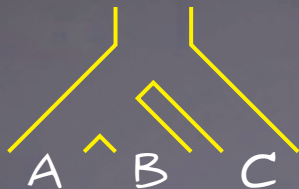


## Embedding Rules

Gene tree  $G$ , species tree  $S$

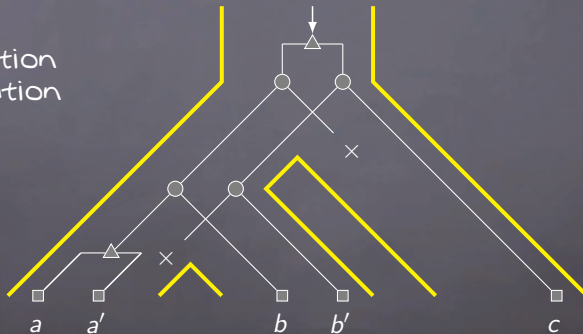
- mapping  $\rho: V(G) \rightarrow V(S)$
- $\ell$  is leaf in  $G \rightsquigarrow \rho(\ell)$  "corresponds" to  $\ell$  ( $a \rightarrow A$ ,  $a' \rightarrow A$ , etc.)
- $u \in V(G)$  is called duplication if  $\rho(u) = \rho(c)$  for any child  $c$  of  $u$  in  $G$
- all non-leaves of  $G$  that are not duplications are called speciations
- each edge  $uv$  of  $G$  incurs a loss-cost equal to the number of edges in the  $\rho(u)-\rho(v)$ -path in  $S$  minus 1 if  $v$  is a speciation or 0 if  $v$  is a duplication

# Reconciliation



## DL-model

- = speciation
- △ = duplication
- × = loss



# Reconciliation

Goal: embed gene tree into species tree  
(extant genes must map to their species)

## Max Likelihood

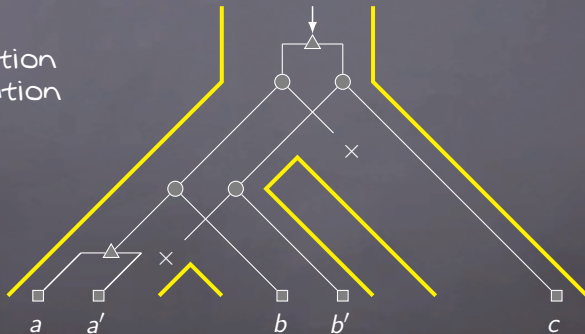
find most probable embedding  
(computationally expensive)

## Parsimony

find embedding minimizing  
#events (possibly weighted)

## DL-model

- = speciation
- △ = duplication
- × = loss





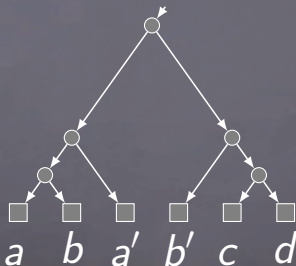
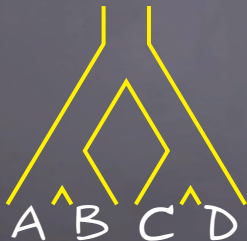
# Reconciliation

## Parsimonious Reconciliation

**Input:** species tree  $S$ , gene tree  $G$ ,  $\delta, \lambda, \tau \in \mathbb{N}$

**Task:** embed  $G$  in  $S$ , minimizing the weighted sum of events

**Result:** LCA-assignment solves this optimally in  $O(|S|+|G|)$



# Reconciliation

## Parsimonious Reconciliation

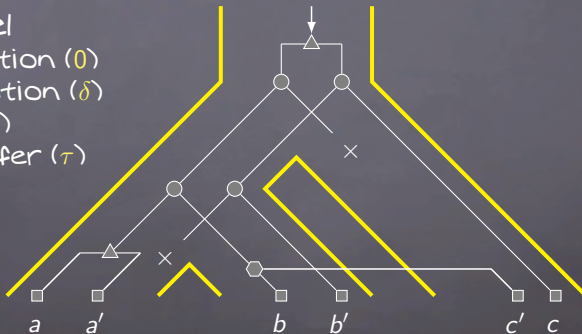
**Input:** species tree  $S$ , gene tree  $G$ ,  $\delta, \lambda, \tau \in \mathbb{N}$

**Task:** embed  $G$  in  $S$ , minimizing the weighted sum of events

**Result:** LCA-assignment solves this optimally in  $O(|S|+|G|)$

## DTL-model

- = speciation ( $\theta$ )
- △ = duplication ( $\delta$ )
- × = loss ( $\lambda$ )
- ⬡ = transfer ( $\tau$ )





# Reconciliation

## Parsimonious Reconciliation

Input: species tree  $S$ , gene tree  $G$ ,  $\delta, \lambda, \tau \in \mathbb{N}$

Task: embed  $G$  in  $S$ , minimizing the weighted sum of events

Result: LCA-assignment solves this optimally in  $O(|S|+|G|)$

Events only between co-existing species  $\leadsto$  time constraints  $\leadsto$  NP-hard

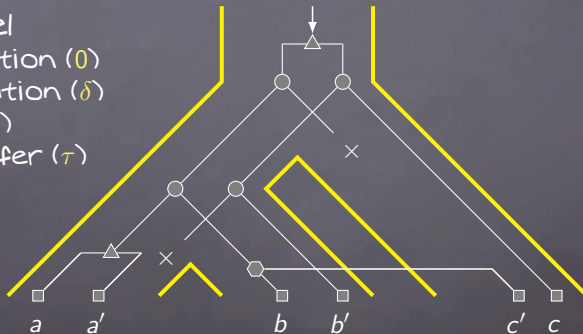
## DTL-model

○ = speciation ( $\theta$ )

△ = duplication ( $\delta$ )

× = loss ( $\lambda$ )

⬡ = transfer ( $\tau$ )



# Reconciliation

## Parsimonious Reconciliation

**Input:** species tree  $S$ , gene tree  $G$ ,  $\delta, \lambda, \tau \in \mathbb{N}$

**Task:** embed  $G$  in  $S$ , minimizing the weighted sum of events

**Result:** LCA-assignment solves this optimally in  $O(|S|+|G|)$

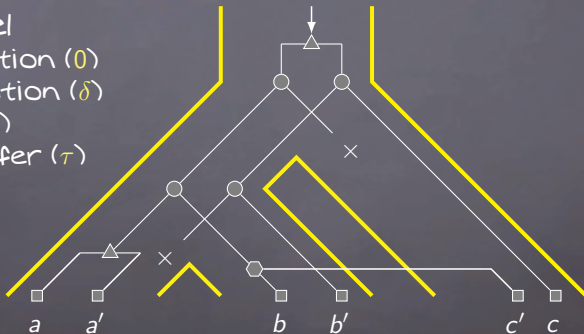
Events only between **co-existing** species  $\rightsquigarrow$  time constraints  $\rightsquigarrow$  NP-hard

**Idea:** take **dated** species tree  $\rightsquigarrow O(|S|^2|G|)$  time

[Doyon et al.'10]

## DTL-model

- = speciation (0)
- △ = duplication (δ)
- × = loss (λ)
- ⬡ = transfer (τ)



# Comparing Phylogenetic Trees

## Distance Measures

- Nearest Neighbor Interchange
- Subtree Prune & Regraft
- Tree Bisection & Reconnection

# Comparing Phylogenetic Trees

## Distance Measures

- Nearest Neighbor Interchange
- Subtree Prune & Regraft
- Tree Bisection & Reconnection
- **now**: via agreement-forests

# Comparing Phylogenetic Trees

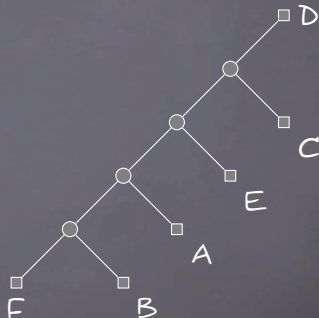
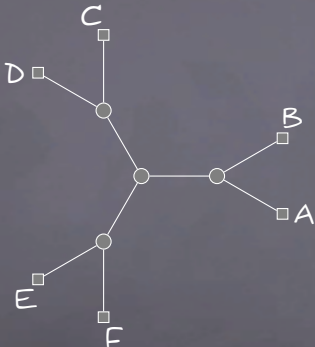
## Distance Measures

- Nearest Neighbor Interchange
- Subtree Prune & Regraft
- Tree Bisection & Reconnection
- **now**: via agreement-forests
- Robinson-Foulds distance
- Quartet/triplet distance

# Agreement Forests

## Definition

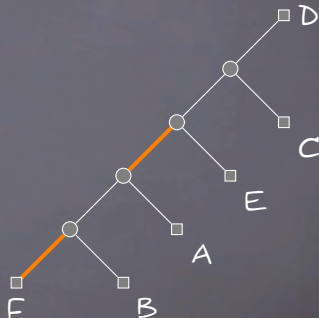
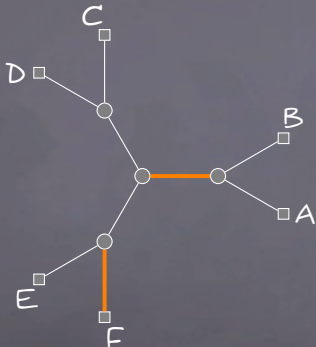
A forest  $F$  is called **agreement forest** of trees  $T_1$  and  $T_2$  if  $F$  can be obtained from  $T_1$  and  $T_2$  by removing edges.



# Agreement Forests

## Definition

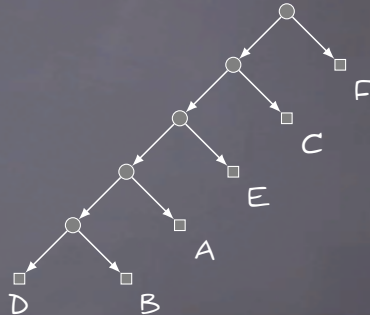
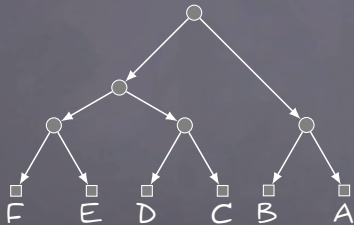
A forest  $F$  is called **agreement forest** of trees  $T_1$  and  $T_2$  if  $F$  can be obtained from  $T_1$  and  $T_2$  by removing edges.



# Agreement Forests

## Definition

A forest  $F$  is called **agreement forest** of trees  $T_1$  and  $T_2$  if  $F$  can be obtained from  $T_1$  and  $T_2$  by removing edges.

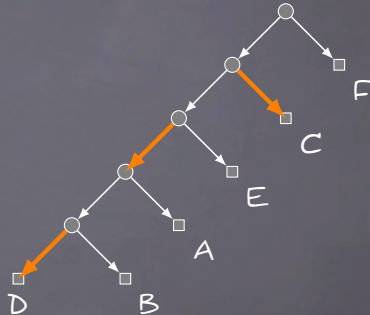
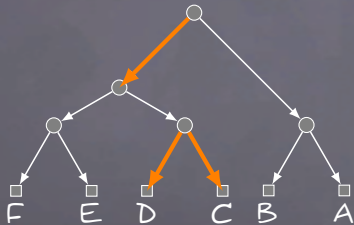




# Agreement Forests

## Definition

A forest  $F$  is called **agreement forest** of trees  $T_1$  and  $T_2$  if  $F$  can be obtained from  $T_1$  and  $T_2$  by removing edges.



# Agreement Forests

## Definition

A forest  $F$  is called **agreement forest** of trees  $T_1$  and  $T_2$  if  $F$  can be obtained from  $T_1$  and  $T_2$  by removing edges.

## Theorem [Allen & Steel, '01]

$\text{TBR-distance}(T_1, T_2) = \# \text{trees in smallest agreement forest} - 1$   
NP-hard to compute

## Theorem [Bordewich & Semple, '04]

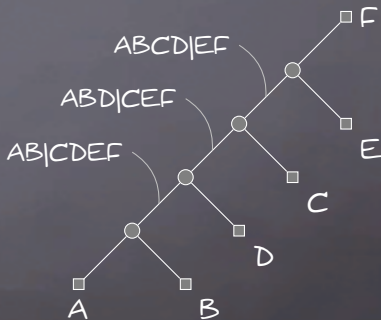
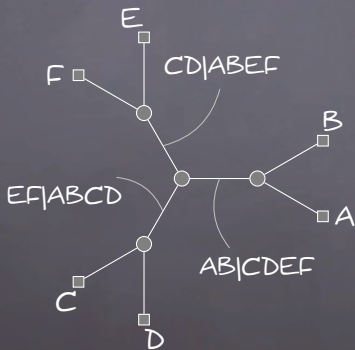
$\text{rSPR-distance}(T_1, T_2) = \# \text{trees in smallest rooted agreement forest} - 1$   
NP-hard to compute

# Robinson-Foulds Distance

## Definition

$RF(T_1, T_2) = \# \text{splits/clusters occurring in exactly one of } T_1 \text{ and } T_2$   
= edge-contraction distance a common tree

Note: observe relation to NNI:  $RF(T_1, T_2) \leq 2 \text{ NNI}(T_1, T_2)$

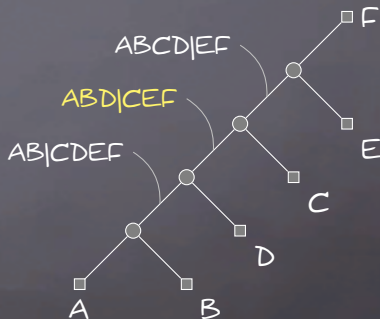
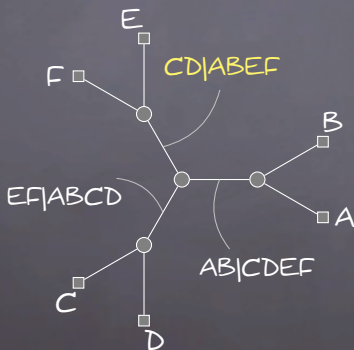


# Robinson-Foulds Distance

## Definition

$RF(T_1, T_2) = \# \text{splits/clusters occurring in exactly one of } T_1 \text{ and } T_2$   
= edge-contraction distance a common tree

Note: observe relation to NNI:  $RF(T_1, T_2) \leq 2 \text{ NNI}(T_1, T_2)$

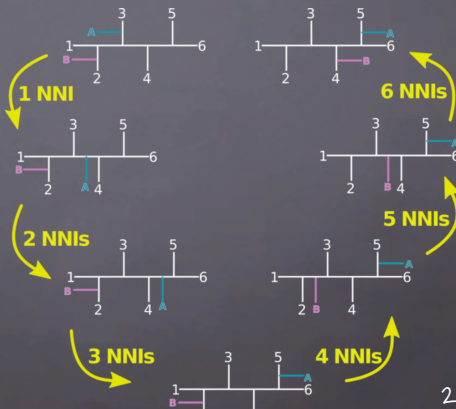
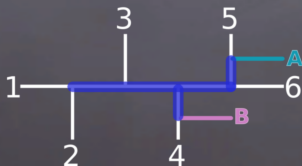
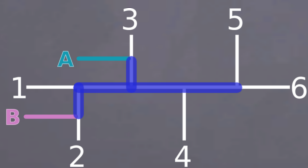


# Robinson-Foulds Distance

## Definition

$RF(T_1, T_2) = \# \text{splits/clusters occurring in exactly one of } T_1 \text{ and } T_2$   
= edge-contraction distance a common tree

Note: observe relation to NNI:  $RF(T_1, T_2) \leq 2 \text{ NNI}(T_1, T_2)$



# Robinson-Foulds Distance

## Definition

$RF(T_1, T_2) = \# \text{splits/clusters occurring in exactly one of } T_1 \text{ and } T_2$   
= edge-contraction distance a common tree

Note: observe relation to NNI:  $RF(T_1, T_2) \leq 2 \text{ NNI}(T_1, T_2)$

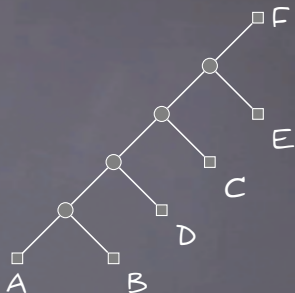
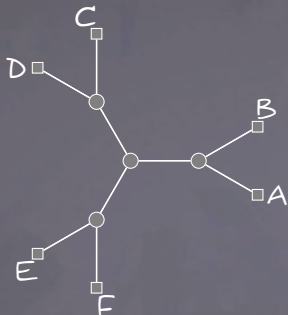
Note: splits correspond to clusters when rooted at **last** leaf

## Day's Algorithm (common clusters in $O(n)$ )

[Day'85]

1. relabel all leaves such that leaves continuous in  $T_1$
2. each node in  $T_1$  knows:  
 $L$  = smallest leaf in cluster       $\neq$        $R$  = largest leaf in cluster  
 $\leadsto$  note  $T_1$ 's clusters in hash-set
3. each node in  $T_2$  knows:  $L$ ,  $R$ , and size  $N$  of its cluster
4. each node in  $T_2$  checks  $[L, R]$  in table only if  $R - L = N - 1$   
(lookup in  $T_1$ 's cluster-set in  $O(1)$  (average) time)

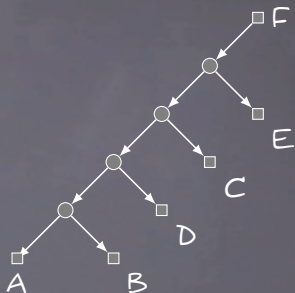
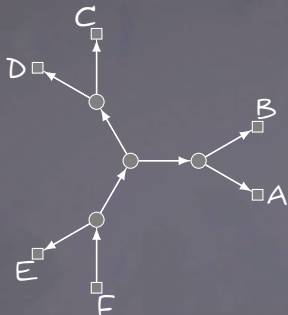
## Robinson-Foulds Distance



[Day'85]

1. relabel all leaves such that leaves continuous in  $T_1$
2. each node in  $T_1$  knows:  
 $L$  = smallest leaf in cluster       $R$  = largest leaf in cluster  
     $\rightsquigarrow$  note  $T_1$ 's clusters in hash-set
3. each node in  $T_2$  knows:  $L$ ,  $R$ , and size  $N$  of its cluster
4. each node in  $T_2$  checks  $[L, R]$  in table only if  $R - L = N - 1$   
(lookup in  $T_1$ 's cluster-set in  $O(1)$  (average) time)

# Robinson-Foulds Distance



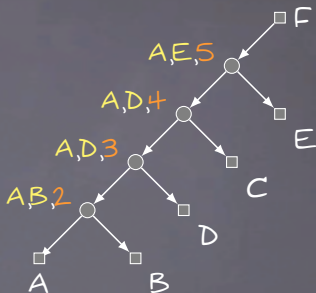
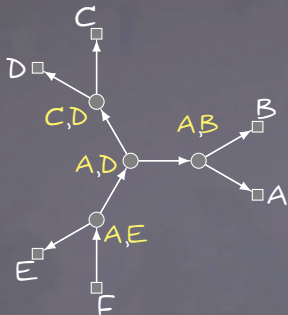
## Day's Algorithm (common clusters in $O(n)$ )

[Day'85]

1. relabel all leaves such that leaves continuous in  $T_1$
2. each node in  $T_1$  knows:  
 $L$  = smallest leaf in cluster  $\neq$   $R$  = largest leaf in cluster  
 $\leadsto$  note  $T_1$ 's clusters in hash-set
3. each node in  $T_2$  knows:  $L$ ,  $R$ , and size  $N$  of its cluster
4. each node in  $T_2$  checks  $[L, R]$  in table only if  $R - L = N - 1$   
 (lookup in  $T_1$ 's cluster-set in  $O(1)$  (average) time)



# Robinson-Foulds Distance

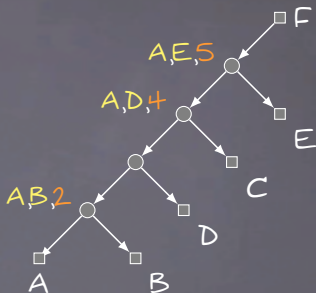
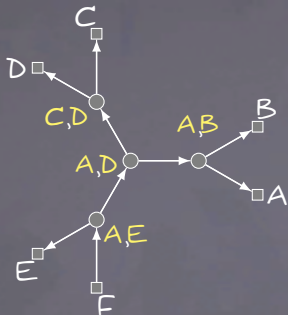


## Day's Algorithm (common clusters in $O(n)$ )

[Day'85]

1. relabel all leaves such that leaves continuous in  $T_1$
2. each node in  $T_1$  knows:  
 $L$  = smallest leaf in cluster  $\neq$   $R$  = largest leaf in cluster  
 $\leadsto$  note  $T_1$ 's clusters in hash-set
3. each node in  $T_2$  knows:  $L$ ,  $R$ , and size  $N$  of its cluster
4. each node in  $T_2$  checks  $[L, R]$  in table only if  $R - L = N - 1$   
 (lookup in  $T_1$ 's cluster-set in  $O(1)$  (average) time)

# Robinson-Foulds Distance



## Day's Algorithm (common clusters in $O(n)$ )

[Day'85]

1. relabel all leaves such that leaves continuous in  $T_1$
2. each node in  $T_1$  knows:  
 $L$  = smallest leaf in cluster  $\neq$   $R$  = largest leaf in cluster  
 $\leadsto$  note  $T_1$ 's clusters in hash-set
3. each node in  $T_2$  knows:  $L$ ,  $R$ , and size  $N$  of its cluster
4. each node in  $T_2$  checks  $[L, R]$  in table only if  $R - L = N - 1$   
 (lookup in  $T_1$ 's cluster-set in  $O(1)$  (average) time)

# Quartet/Triplet Distance

## Definition

$Q/T(T_1, T_2) = \# \text{quartets/triplets occur in exactly one of } T_1 \text{ and } T_2$

# Quartet/Triplet Distance

## Definition

$Q/T(T_1, T_2) = \# \text{quartets/triplets occur. in exactly one of } T_1 \text{ and } T_2$

## computing Q-distance (Binary trees) [Bryant et al.'00]

1. each edge  $uv$  has 4 sets (2 clusters for each of  $u \neq v$ )
2. quartet  $AB|CD$  "belongs" to edge  $e$  if  $e$  splits  $AB|CD$  and  $e$  touches  $AB$ -path  $\rightsquigarrow$  each quartet is owned exactly once
3.  $\forall uv \in T_1 \neq qr \in T_2$ : intersect 4 sets of  $uv$  with split of  $qr$  in  $T_2$
4. sizes of intersections can be precomputed bottom-up in  $O(n^2)$  time

# Quartet/Triplet Distance

## Definition

$Q/T(T_1, T_2) = \# \text{quartets/triplets occur in exactly one of } T_1 \text{ and } T_2$

## computing Q-distance (Binary trees) [Bryant et al.'00]

1. each edge  $uv$  has 4 sets (2 clusters for each of  $u \neq v$ )
2. quartet  $AB|CD$  "belongs" to edge  $e$  if  $e$  splits  $AB|CD$  and  $e$  touches  $AB$ -path  $\rightsquigarrow$  each quartet is owned exactly once
3.  $\forall uv \in T_1 \neq qr \in T_2$ : intersect 4 sets of  $uv$  with split of  $qr$  in  $T_2$
4. sizes of intersections can be precomputed bottom-up in  $O(n^2)$  time

## State of the Art

count conflict quartets/triplets  $\rightsquigarrow O(n \log n)$  time

[Brodal et al.'13]

enumerate conflict quartets  $\rightsquigarrow O(n^2 + d)$  time

[Bryant et al.'00]

enumerate conflict triplets  $\rightsquigarrow O(n + d)$  time

[Weller'17]

# Phylogenetic Networks

## Observation

Trees cannot capture hybridization



# Phylogenetic Networks

## Observation

Trees cannot capture hybridization  $\rightsquigarrow$  phylogenetic network



# Phylogenetic Networks

## Observation

Trees cannot capture hybridization  $\leadsto$  phylogenetic network

## Definition

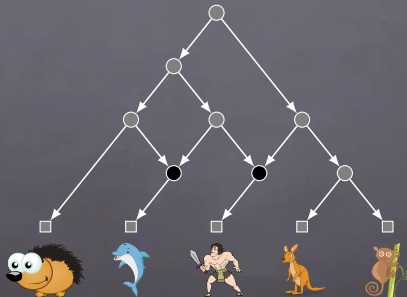
evolutionary network  $N$  = rooted DAG, leaves labeled (taxa)

reticulations  $R$  = vertices of in-degree  $\geq 2$

Binary = all inner vertices degree 3

Block = maximal Biconnected component

display  $T$  = subdivision of  $T$  is a subgraph





# Phylogenetic Networks

## Observation

Trees cannot capture hybridization  $\leadsto$  phylogenetic network

## Definition

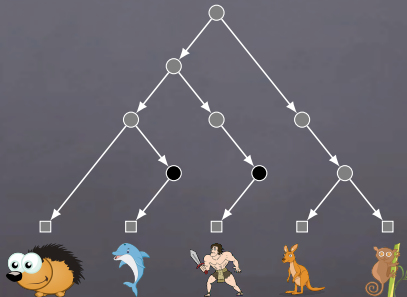
evolutionary network  $N$  = rooted DAG, leaves labeled (taxa)

reticulations  $R$  = vertices of in-degree  $\geq 2$

Binary = all inner vertices degree 3

Block = maximal Biconnected component

display  $T$  = subdivision of  $T$  is a subgraph

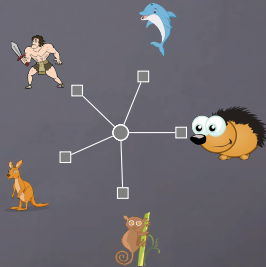


# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]



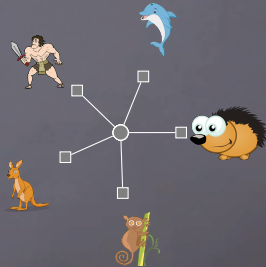
c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]



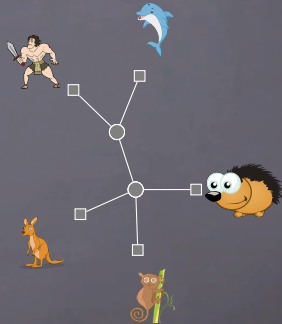
c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nmid X|Y$  **incompatible** if both  $A \nmid B$  intersect both  $X \nmid Y$

Convex Hull Algorithm [Holland et al., '04]



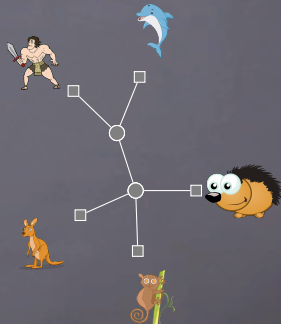
c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]



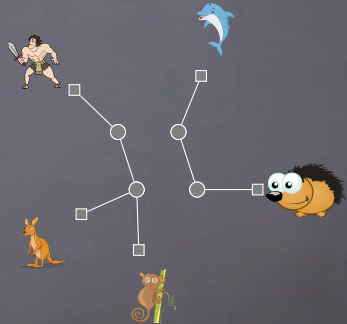
c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]



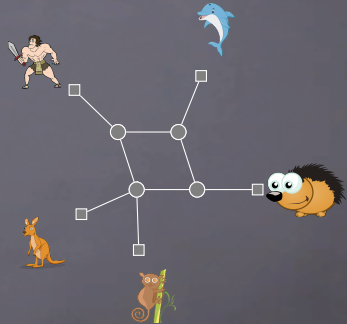
c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]



c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Split Networks

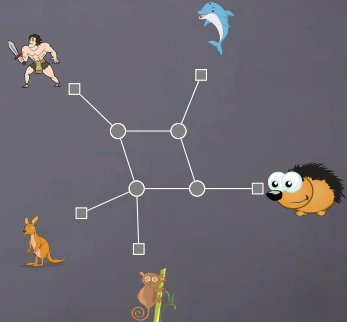
**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]

fox  
AUS  
pouch  
land

	✓	✓	✗	✓
	✗	✓	✗	✓
	✗	✗	✓	✗
	✓	✓	✓	✓
	✓	✗	✗	✓



c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

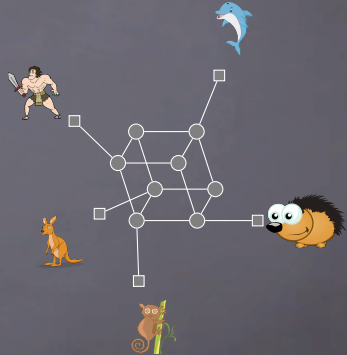


# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]



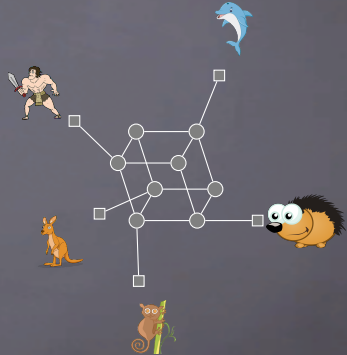
c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Split Networks

**split** = bipartition of set of taxa

splits  $A|B \nsubseteq X|Y$  **incompatible** if both  $A \nsubseteq B$  intersect both  $X \nsubseteq Y$

Convex Hull Algorithm [Holland et al., '04]



c.f.: Neighbor Net [Bryant & Moulton, '03]  
(for **circular splits**)

# Consensus Split Networks

## Strategy

1. list all splits of all input trees
2. extend splits to full taxa using "Z-closure"
3. Build consensus

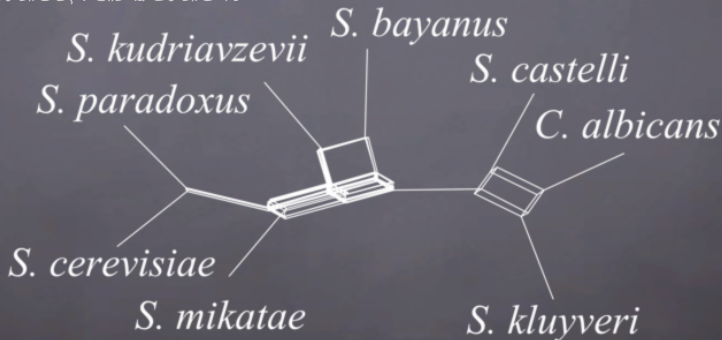
# Consensus Split Networks

## Strategy

1. list all splits of all input trees
2. extend splits to full taxa using "Z-closure"
3. Build consensus

## Experimental Study – 106 gene trees (yeast)

[Rokas et al.'03, Holland et al.'04]



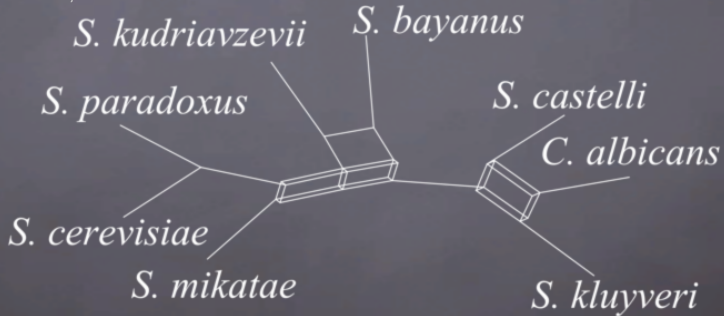
# Consensus Split Networks

## Strategy

1. list all splits of all input trees
2. extend splits to full taxa using "Z-closure"
3. Build consensus

## Experimental Study – 106 gene trees (yeast)

[Rokas et al.'03, Holland et al.'04]



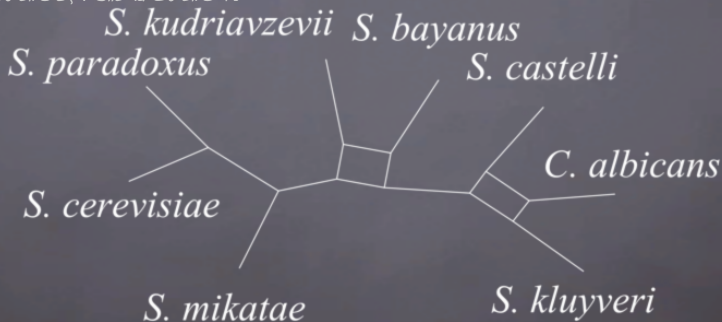
# Consensus Split Networks

## Strategy

1. list all splits of all input trees
2. extend splits to full taxa using "Z-closure"
3. Build consensus

## Experimental Study – 106 gene trees (yeast)

[Rokas et al.'03, Holland et al.'04]



# Consensus Split Networks

## Strategy

1. list all splits of all input trees
2. extend splits to full taxa using "Z-closure"
3. Build consensus

## Experimental Study – 106 gene trees (yeast)

[Rokas et al.'03, Holland et al.'04]



# Rooted Network Reconstruction

## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters



# Rooted Network Reconstruction

## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters

## Example

$\{a,b,c,d\}, \{c,d,e,f,g,h\}, \{c,d,e,f,g\}, \{e,f,g,h\}, \{c,d,e\}, \{e,f,g\}, \{a,b\}, \{c,d\}, \{f,g\}$



# Rooted Network Reconstruction

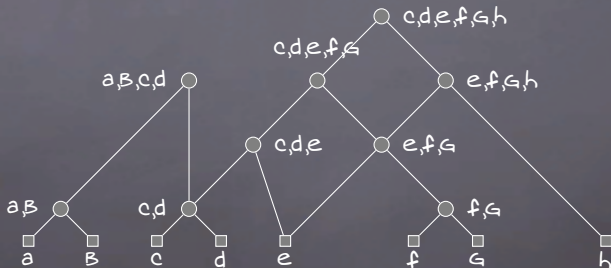
## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters

## Example

$\{a,b,c,d\}, \{c,d,e,f,g,h\}, \{c,d,e,f,g\}, \{e,f,g,h\}, \{c,d,e\}, \{e,f,g\}, \{a,b\}, \{c,d\}, \{f,g\}$



# Rooted Network Reconstruction

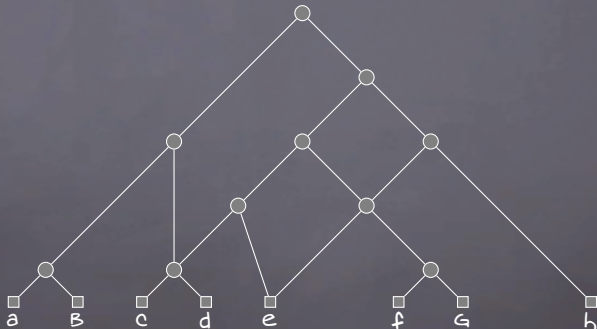
## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters

## Example

$\{a,b,c,d\}, \{c,d,e,f,g,h\}, \{c,d,e,f,g\}, \{e,f,g,h\}, \{c,d,e\}, \{e,f,g\}, \{a,b\}, \{c,d\}, \{f,g\}$



# Rooted Network Reconstruction

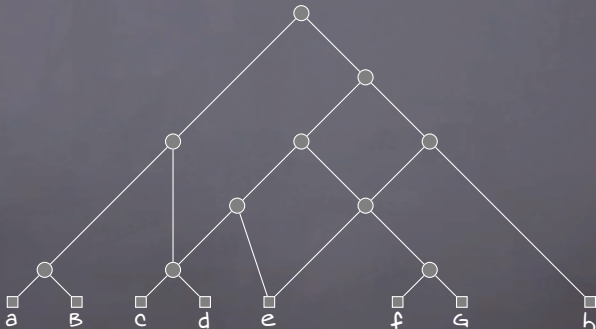
## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters

## Example

$\{a,b,c,d\}, \{c,d,e,f,g,h\}, \{c,d,e,f,g\}, \{e,f,g,h\}, \{c,d,e\}, \{e,f,g\}, \{a,b\}, \{c,d\}, \{f,g\}$



c.f. "cluster popping" [Huson & Rupp, 08]

# Rooted Network Reconstruction

## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters

## Problem

may produce more reticulations than necessary to explain the data

# Rooted Network Reconstruction

## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$   
 $\leadsto$  rooted network is **hasse diagram** of its clusters

## Problem

may produce more reticulations than necessary to explain the data

## Hybridization Number

**Input:** set of trees  $T$ , int  $k$

**Question:** Is there a network with  $\leq k$  reticulations displaying all trees in  $T$ ?

# Rooted Network Reconstruction

## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters

## Problem

may produce more reticulations than necessary to explain the data

## Hybridization Number

**Input:** set of trees  $T$ , int  $k$

**Question:** Is there a network with  $\leq k$  reticulations displaying all trees in  $T$ ?

$\leadsto$  NP-hard for 2 trees [Bordewich & Semple, '07]

# Rooted Network Reconstruction

## Observation

rooted network: cluster of  $u \subseteq$  cluster of  $v \Leftrightarrow u \leq v$

$\leadsto$  rooted network is **hasse diagram** of its clusters

## Problem

may produce more reticulations than necessary to explain the data

## Hybridization Number

**Input:** set of trees  $T$ , int  $k$

**Question:** Is there a network with  $\leq k$  reticulations displaying all trees in  $T$ ?

$\leadsto$  NP-hard for 2 trees [Bordewich & Semple, '07]

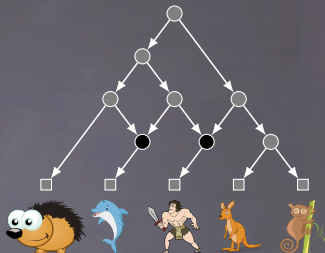
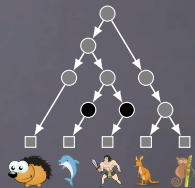
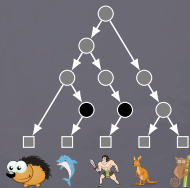
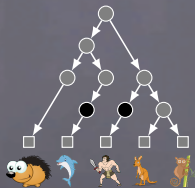
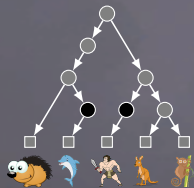
**Note:**  $HN(T_1, T_2) = \max \text{acyclic agreement forest} - 1$  [Baroni et al., '05]



# Networks Display Trees

## Observation

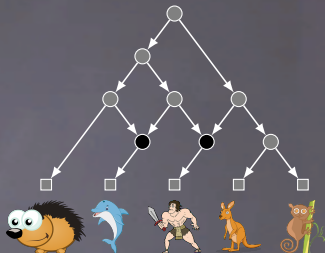
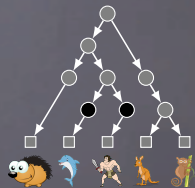
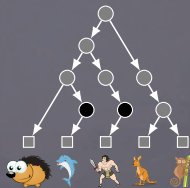
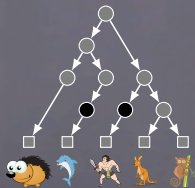
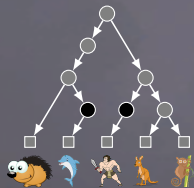
A network may display up to  $2^{|R|}$  trees.



# Networks Display Trees

## Observation

A network may display up to  $2^{|R|}$  trees.

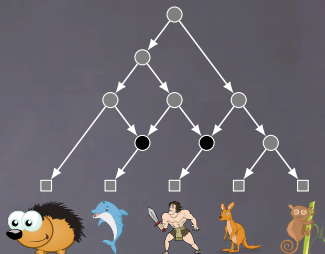
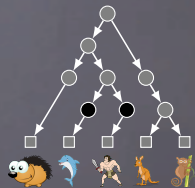
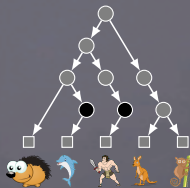
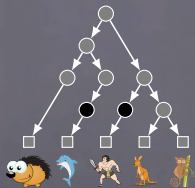
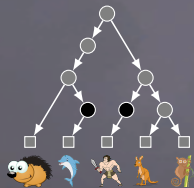


# Networks Display Trees

## Observation

A network may display up to  $2^{|R|}$  trees.

But: how to decide if a given tree is displayed?



# Networks Display Trees

## Tree Containment

Input: a network  $N$ , a tree  $T$

Question: Does  $N$  display  $T$ ?

# Networks Display Trees

## Tree Containment

Input: a network  $N$ , a tree  $T$

Question: Does  $N$  display  $T$ ?

$\rightsquigarrow$  NP-hard (from Disjoint Paths) [Kanj et al.'08]

# Networks Display Trees

## Tree Containment

Input: a network  $N$ , a tree  $T$

Question: Does  $N$  display  $T$ ?

$\rightsquigarrow$  NP-hard (from Disjoint Paths) [Kanj et al.'08]

Note: linear time on reticulation visible  $N$  [Gunawan,'18][Weller,'18]

# Networks Display Trees

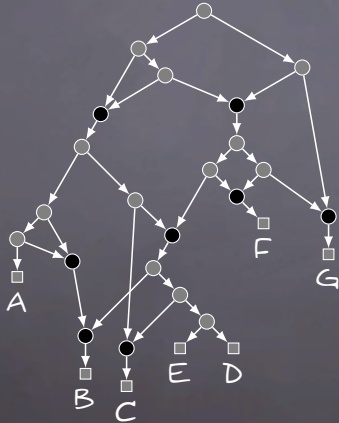
## Tree Containment

Input: a network  $N$ , a tree  $T$

Question: Does  $N$  display  $T$ ?

~ NP-hard (from Disjoint Paths) [Kanj et al.'08]

Note: linear time on reticulation visible  $N$  [Gunawan,'18][Weller,'18]



Can you see it?



# Small Taxonomy of Network Classes

