

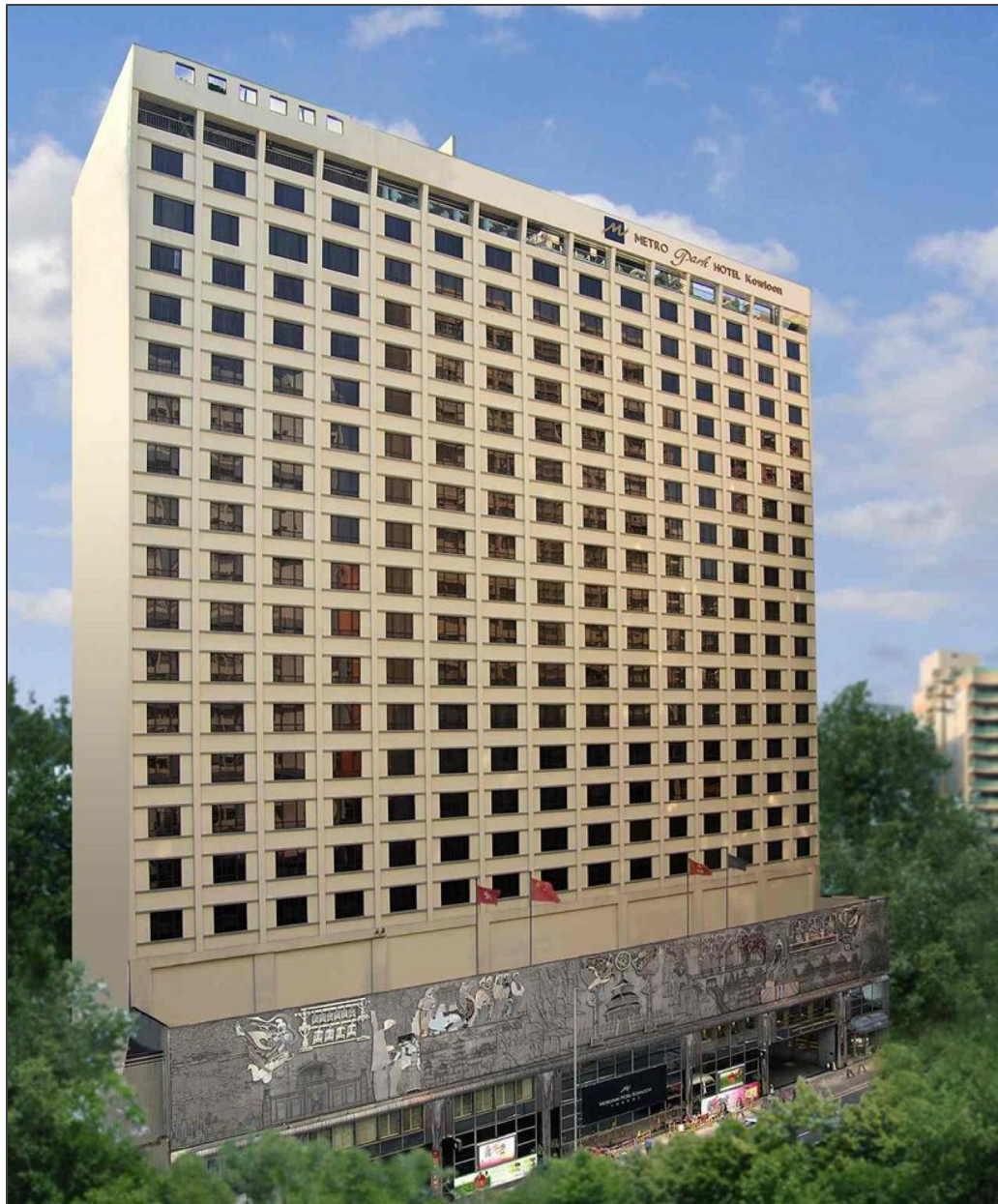
Which Animal Gave Us SARS?

*Evolutionary Trees Part 1:
The Neighbor-Joining Algorithm*

Phillip Compeau and Pavel Pevzner

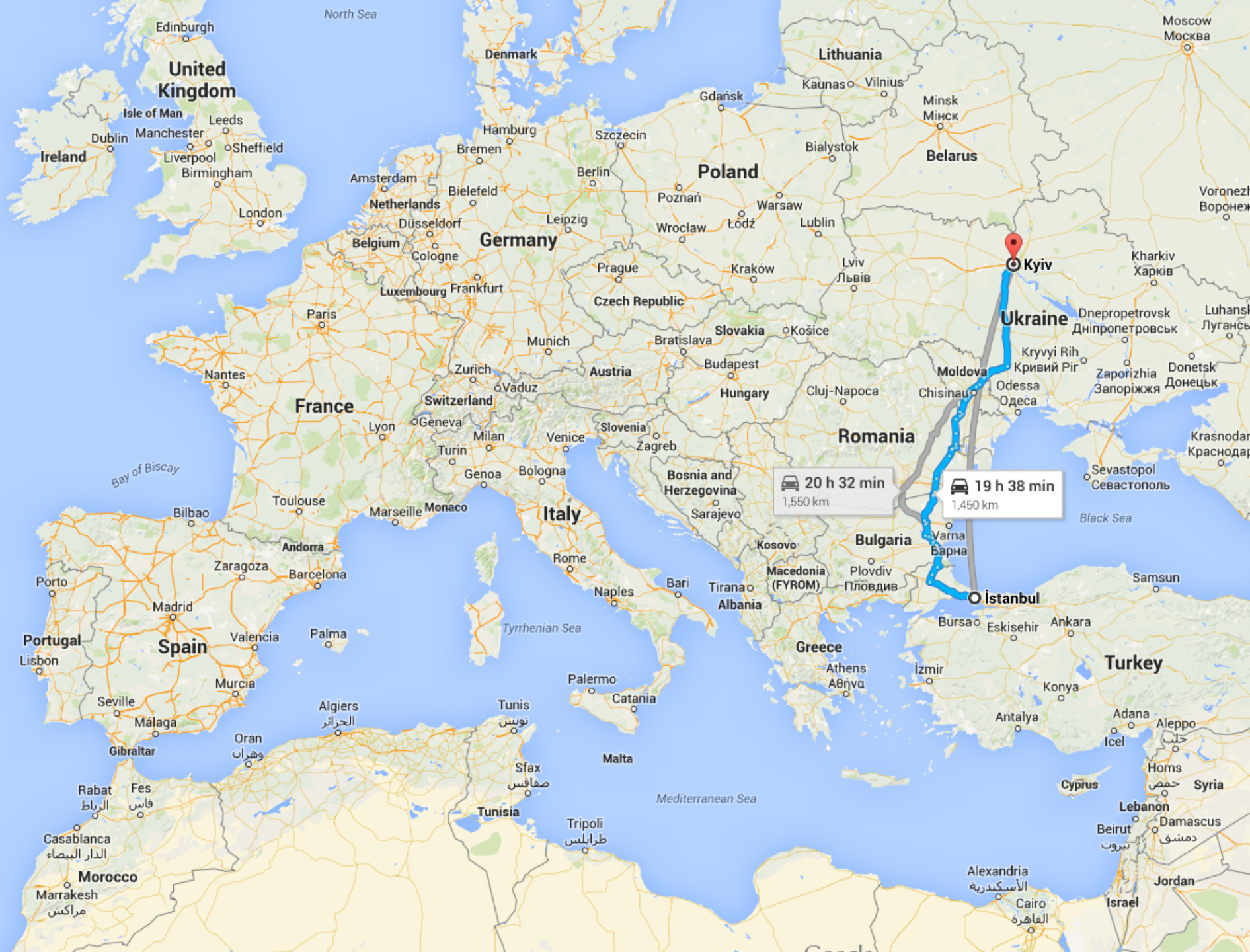
Bioinformatics Algorithms: An Active Learning Approach

©2018 by Compeau and Pevzner. All rights reserved.



Modern boundaries are shown for reference.



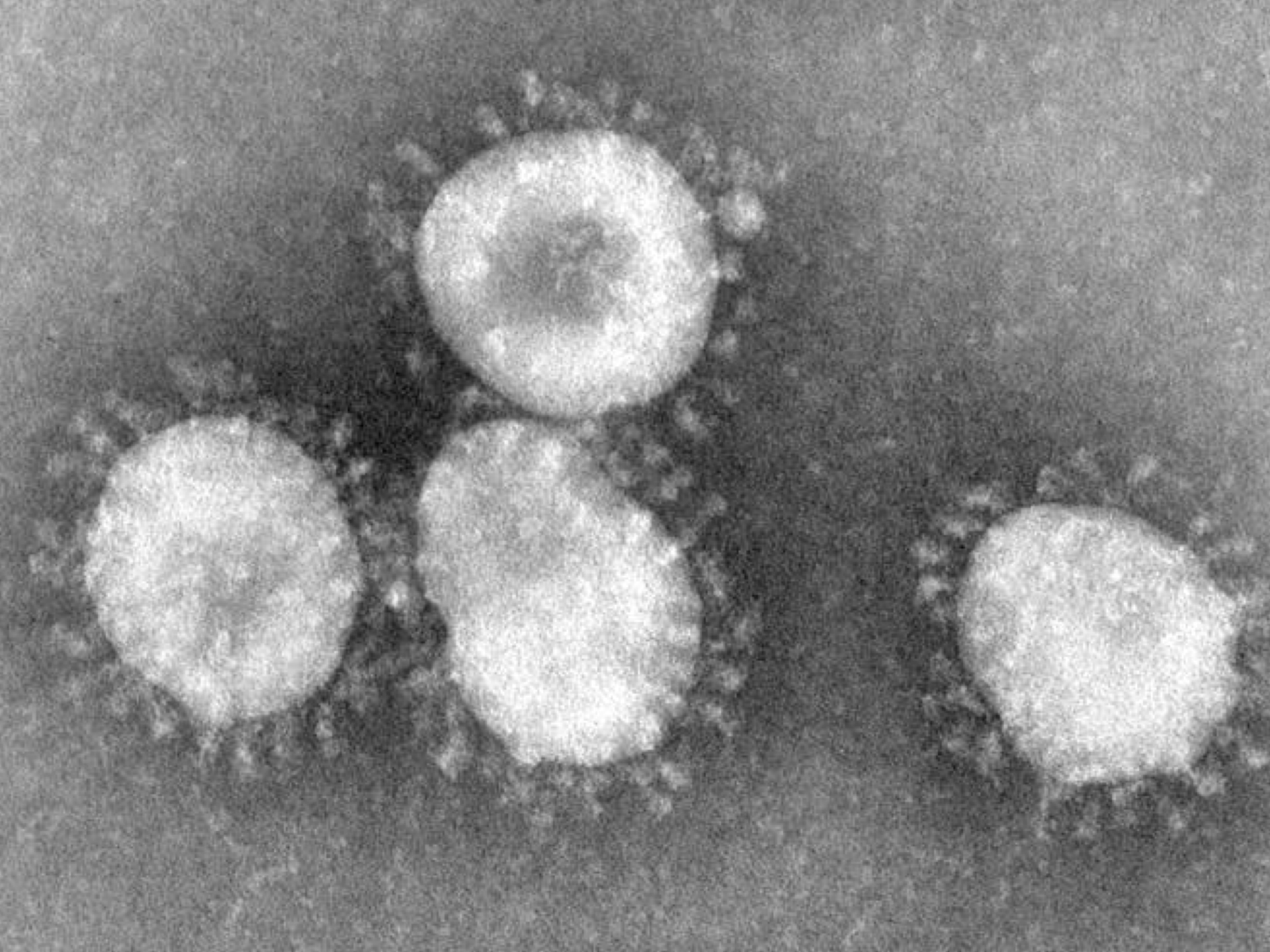




**The Spread
of ??????**



The Spread of SARS





Questions about SARS

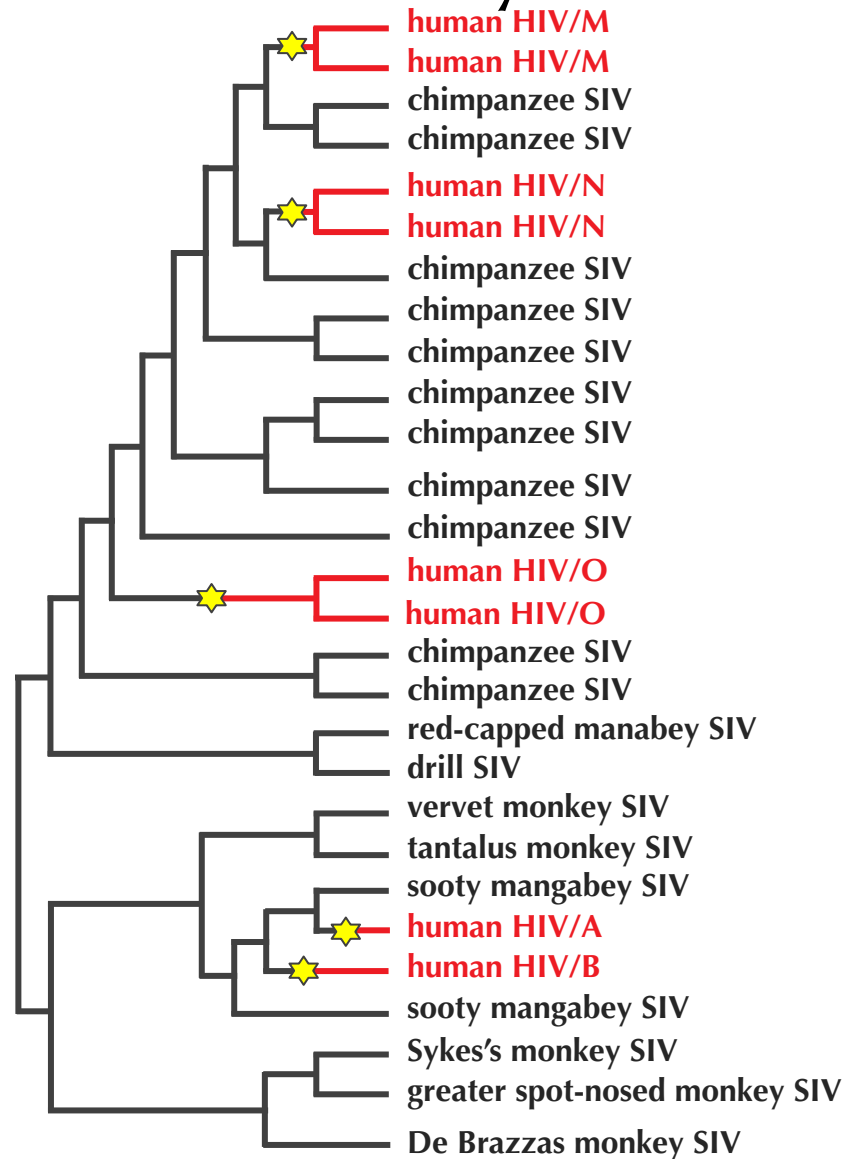
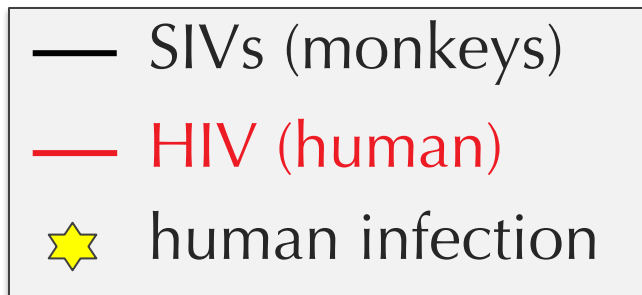
Which animal gave us SARS? How does SARS compare to other viruses and how did it mutate over time?

Questions about SARS

Which animal gave us SARS? How does SARS compare to other viruses and how did it mutate over time?

To answer these questions, we need to learn how to construct **evolutionary trees** (a.k.a. **phylogenies**).

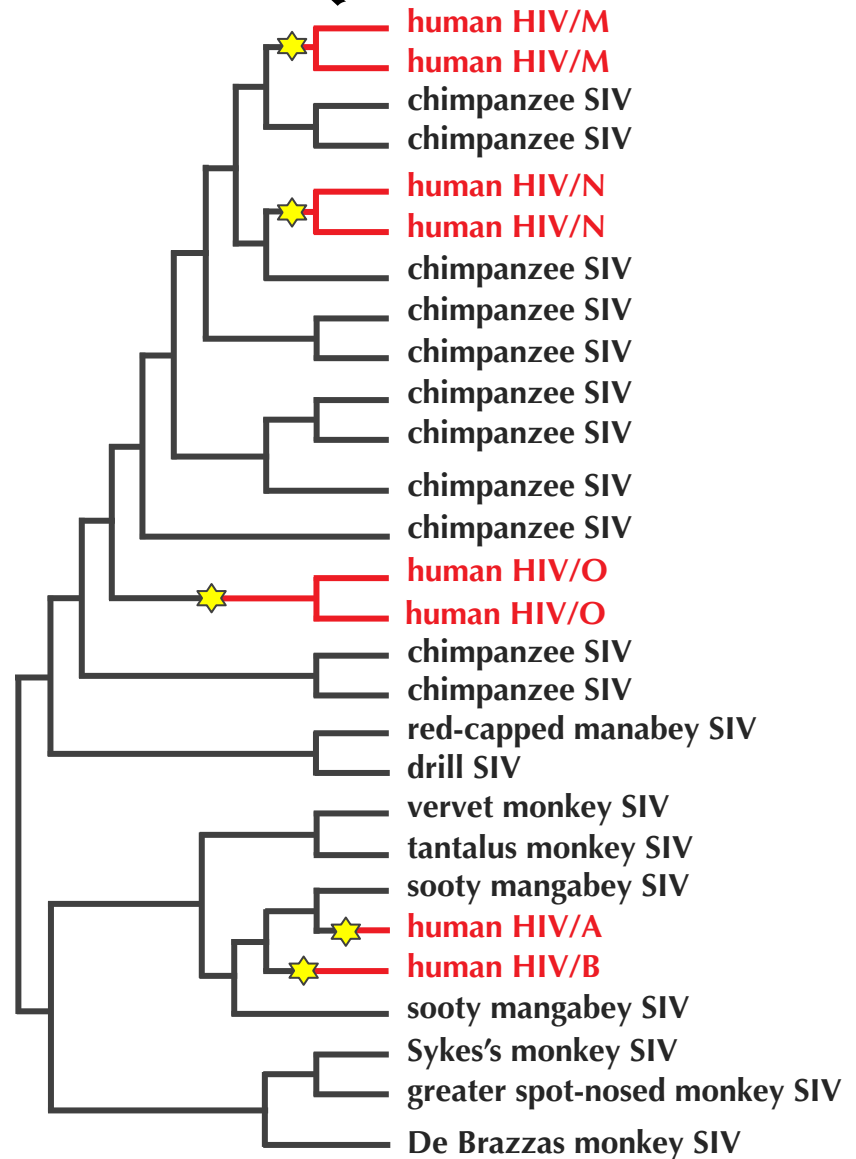
Example: HIV Evolutionary Tree



Two Computational Questions

How do we construct the tree's *structure*?

Can we infer anything about the ancestors?

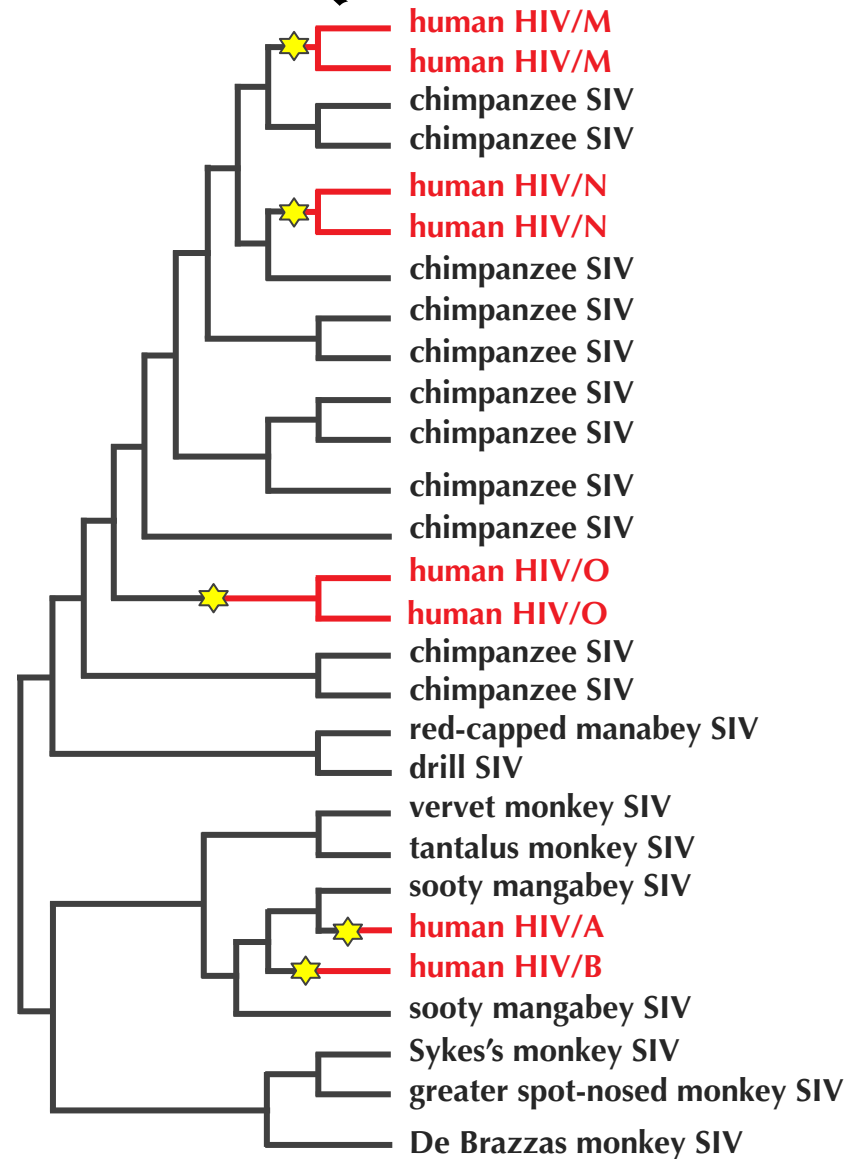


Two Computational Questions

How do we construct the tree's *structure*?

Can we infer anything about the ancestors?

Checkpoint: Any thoughts on how we could answer either question?



Trees

```
graph TD
    LIFE((LIFE)) --- Bacteria((bacteria))
    LIFE --- Node1(( ))
    Node1 --- Archaeobacteria((archaeobacteria))
    Node1 --- Eukaryotes((EUKARYOTES))
    Eukaryotes --- Protocists((protocists))
    Eukaryotes --- Node2(( ))
    Node2 --- Plants((PLANTS))
    Node2 --- Node3(( ))
    Node3 --- GreenAlgae((green algae))
    Node3 --- Fungi((fungi))
    Node3 --- Animals((ANIMALS))
    Plants --- Node4(( ))
    Node4 --- Mosses((mosses))
    Node4 --- Node5(( ))
    Node5 --- Ferns((ferns))
    Node5 --- Node6(( ))
    Node6 --- FloweringPlants((flowering seed plants))
    Node6 --- NonFloweringPlants((non-flowering seed plants))
    Animals --- Sponges((sponges))
    Animals --- Node7(( ))
    Node7 --- Cnidarians((cnidarians))
    Node7 --- Node8(( ))
    Node8 --- Flatworms((flatworms))
    Node8 --- Node9(( ))
    Node9 --- Lophophorates((lophophorates))
    Node9 --- Rotifers((rotifers))
    Node9 --- Roundworms((roundworms))
    Lophophorates --- Echinoderms((echinoderms))
    Echinoderms --- SegmentedWorms((segmented worms))
    SegmentedWorms --- Mollusks((mollusks))
    Mollusks --- Crustaceans((crustaceans))
    Mollusks --- Insects((insects))
    Lophophorates --- Chelicerates((chelicerates))
    Chelicerates --- Arthropods((ARTHROPODS))
    Arthropods --- Crustaceans
    Arthropods --- Insects
    Vertebrates((VERTEBRATES)) --- Node10(( ))
    Node10 --- CartilaginousFish((cartilaginous fish))
    Node10 --- Node11(( ))
    Node11 --- BonyFish((bony fish))
    Node11 --- Tetrapods((TETRAPODS))
    Tetrapods --- Amphibians((amphibians))
    Amphibians --- Mammals((mammals))
    Mammals --- Node12(( ))
    Node12 --- Turtles((turtles))
    Node12 --- Node13(( ))
    Node13 --- SnakesAndLizards((snakes & lizards))
    Node13 --- CrocodilesAndBirds((crocodiles & birds))
```

The diagram illustrates the hierarchical structure of life, starting from the root 'LIFE' and branching down to specific organisms. The tree is organized into several main branches: Bacteria, Archaeobacteria, Eukaryotes, Plants, and Animals. The Eukaryotes branch further into Protocists, Green algae, Fungi, and Animals. The Animals branch further into Sponges, Cnidarians, Flatworms, Lophophorates, Rotifers, Roundworms, Echinoderms, Segmented worms, Mollusks, Chelicerates, and Arthropods. The Arthropods branch further into Crustaceans and Insects. The Vertebrates branch further into Cartilaginous fish, Bony fish, Tetrapods, Amphibians, Mammals, and Turtles. The Mammals branch further into Snakes & lizards and Crocodiles & birds.

Tree: Connected graph containing no cycles.

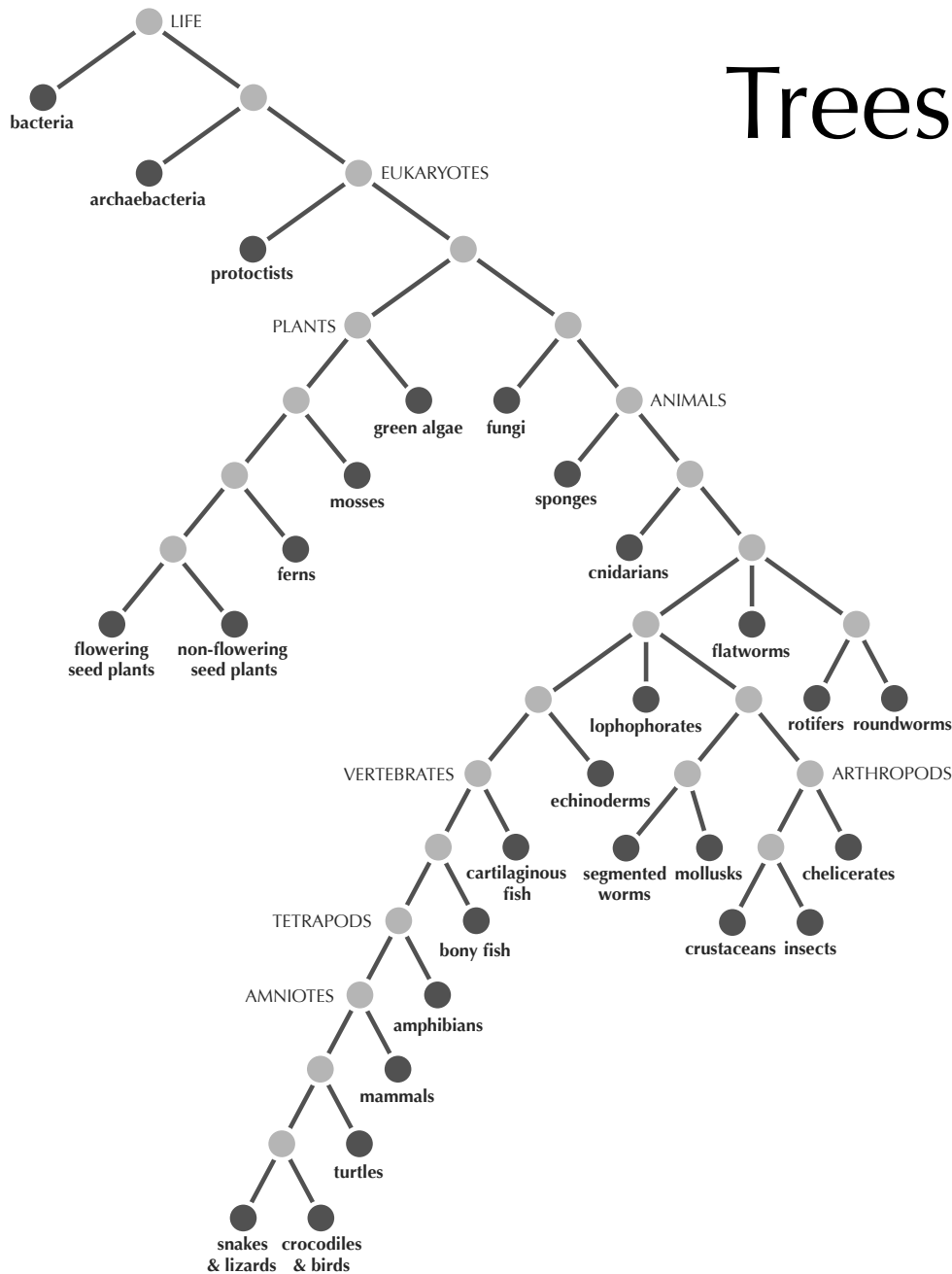
Trees

```
graph TD; LIFE --> Bacteria[bacteria]; LIFE --> Node1(( )); Node1 --> Archaeobacteria[archaeobacteria]; Node1 --> Node2(( )); Node2 --> Eukaryotes[EUKARYOTES]; Node2 --> Node3(( )); Node3 --> Protocists[protocists]; Node3 --> Node4(( )); Node4 --> Plants[PLANTS]; Node4 --> Node5(( )); Node5 --> GreenAlgae[green algae]; Node5 --> Fungi[fungi]; Node5 --> Animals[ANIMALS]; Plants --> Node6(( )); Node6 --> Mosses[mosses]; Node6 --> Node7(( )); Node7 --> Ferns[ferns]; Node7 --> Node8(( )); Node8 --> FloweringPlants[flowering seed plants]; Node8 --> NonFloweringPlants[non-flowering seed plants]; Animals --> Sponges[sponges]; Animals --> Node9(( )); Node9 --> Cnidarians[cnidarians]; Node9 --> Node10(( )); Node10 --> Flatworms[flatworms]; Node10 --> Node11(( )); Node11 --> Rotifers[rotifers]; Node11 --> Roundworms[roundworms]; Flatworms --> Lophophorates[lophophorates]; Flatworms --> Node12(( )); Node12 --> Echinoderms[echinoderms]; Node12 --> Node13(( )); Node13 --> SegmentedWorms[segmented worms]; Node13 --> Mollusks[mollusks]; Node13 --> Node14(( )); Node14 --> Crustaceans[crustaceans]; Node14 --> Insects[insects]; Node14 --> Chelicerates[chelicerates]; Lophophorates --> Vertebrates[VERTEBRATES]; Lophophorates --> Node15(( )); Node15 --> CartilaginousFish[cartilaginous fish]; Node15 --> BonyFish[bony fish]; Vertebrates --> Tetrapods[TETRAPODS]; Tetrapods --> Amphibians[amphibians]; Tetrapods --> Mammals[mammals]; Tetrapods --> Turtles[turtles]; Amphibians --> Mammals; Mammals --> SnakesAndLizards[snakes & lizards]; Mammals --> CrocodilesAndBirds[crocodiles & birds];
```

The diagram illustrates the hierarchical structure of life, starting from the root 'LIFE' and branching down to various organisms. The tree is organized into several main categories: Bacteria, Archaeobacteria, Eukaryotes, Protocists, Plants, Animals, Vertebrates, Tetrapods, and Amniotes. The organisms are represented by black dots, and the relationships are shown by lines connecting the dots. The tree is a binary tree, meaning each node has at most two children. The organisms are listed in a way that reflects their evolutionary relationships, with more closely related organisms sharing a more recent common ancestor. The tree is a good example of how a hierarchical structure can be used to represent complex information in a clear and concise way.

Leaves (degree = 1):
present-day species

Trees

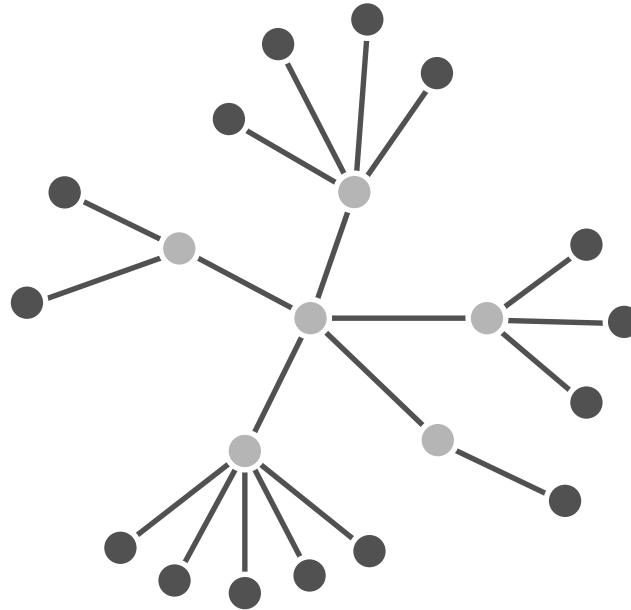
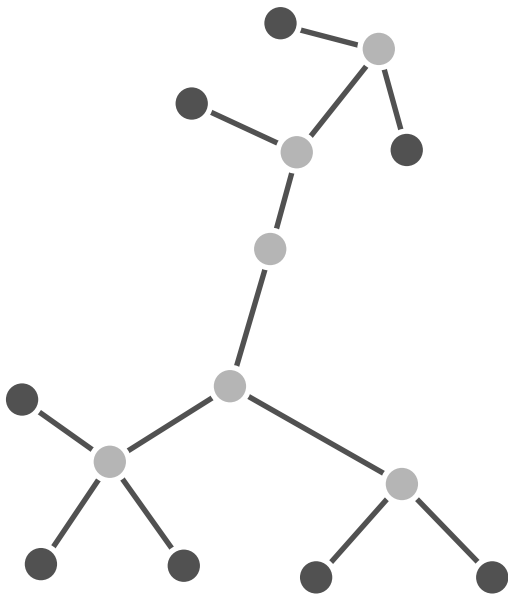


Tree: Connected graph containing no cycles.

Leaves (degree = 1):
present-day species

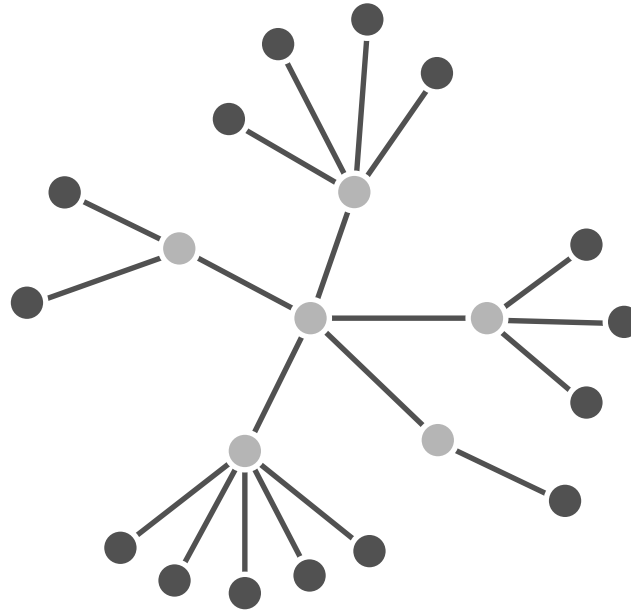
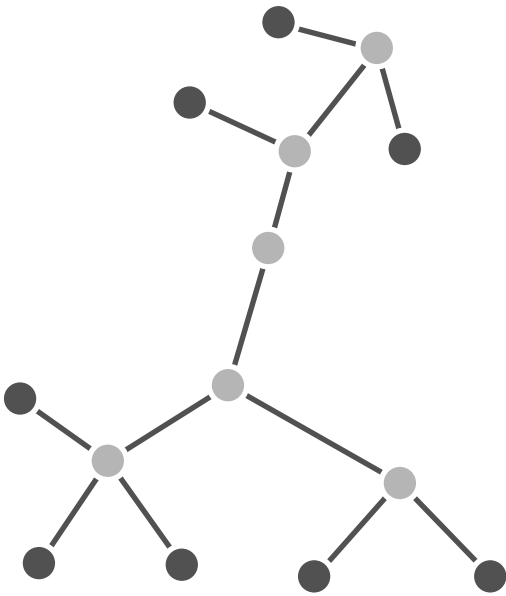
Internal nodes
(degree ≥ 2):
ancestral species

Trees



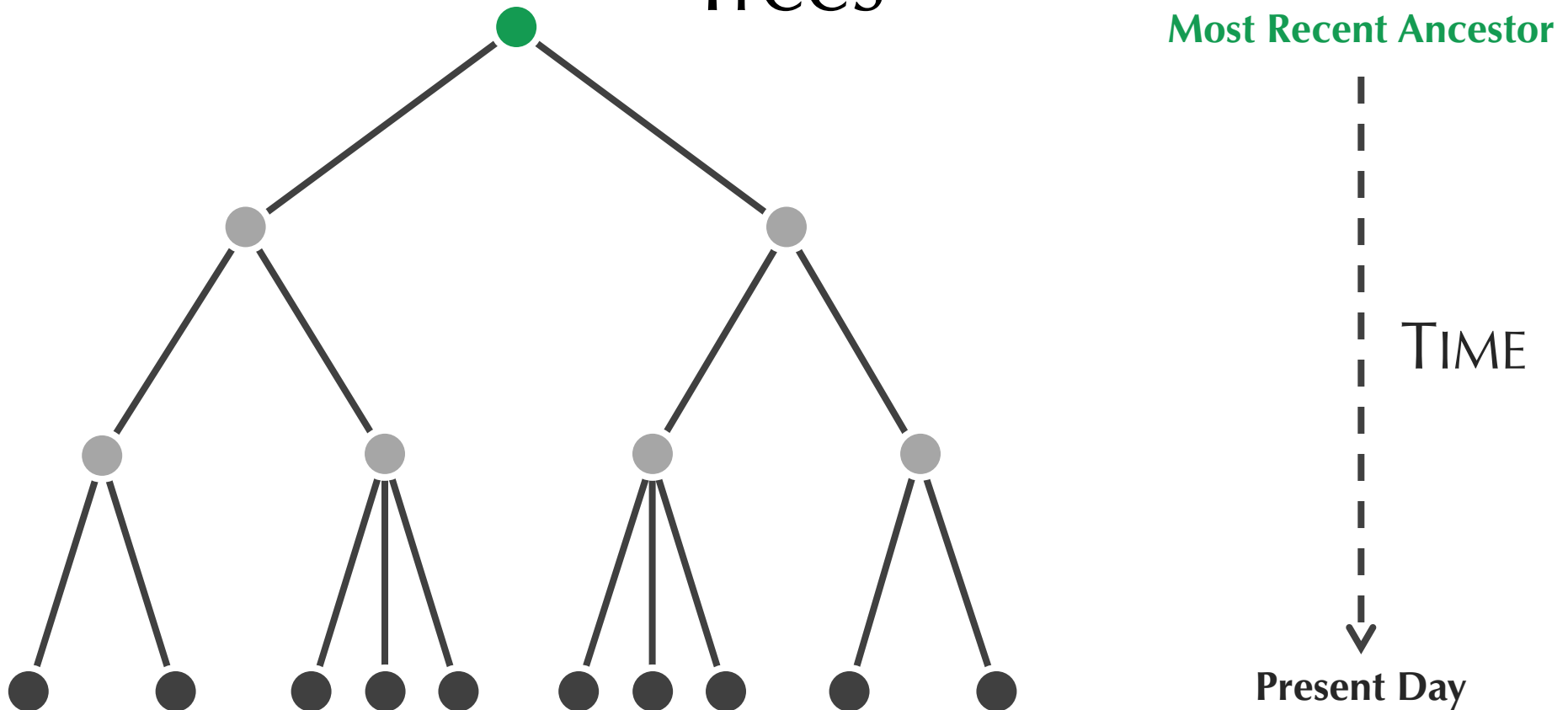
Note: We proved in a previous lecture that every tree with n nodes has exactly $n - 1$ edges.

Trees



Exercise: Prove that there is a unique path connecting any two nodes in a tree.

Trees

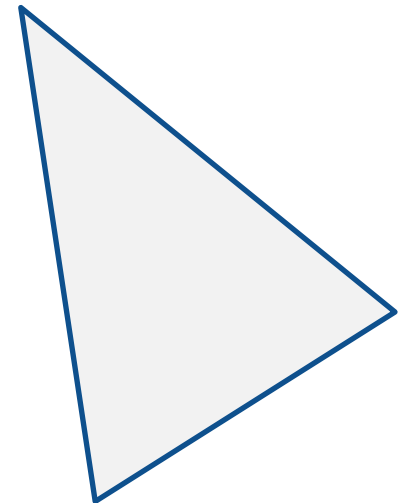


Rooted tree: one node is designated as the **root** (most recent common ancestor).

Definition of a Distance Matrix

Distance matrix: A matrix D representing distances between pairs of n organisms that satisfies three properties:

1. **Symmetry:** $D_{i,j} = D_{j,i}$ for all pairs i, j
2. **Non-negativity:** $D_{i,j} \geq 0$ for all pairs i, j
3. **Triangle inequality:** For all i, j , and k , $D_{i,j} + D_{j,k} \geq D_{i,k}$.



A Multiple Alignment Defines a Simple Distance Matrix

SPECIES	ALIGNMENT
Chimp	ACGTAGGCCT
Human	ATGTAAGACT
Seal	TCGAGAGCAC
Whale	TCGAAAGCAT

A Multiple Alignment Defines a Simple Distance Matrix

$D_{i,j}$ = number of differing symbols between i -th and j -th rows of a multiple alignment.

SPECIES	ALIGNMENT	DISTANCE MATRIX			
		Chimp	Human	Seal	Whale
Chimp	ACGTAGGCCT	0	3	6	4
Human	ATGTAAGACT	3	0	7	5
Seal	TCGAGAGCAC	6	7	0	2
Whale	TCGAAAGCAT	4	5	2	0

A Multiple Alignment Defines a Simple Distance Matrix

$D_{i,j}$ = number of differing symbols between i -th and j -th rows of a multiple alignment.

SPECIES	ALIGNMENT	DISTANCE MATRIX			
		Chimp	Human	Seal	Whale
Chimp	A C GTA G C CT	0	3	6	4
Human	A T GTA A G ACT	3	0	7	5
Seal	TCGAGAGCAC	6	7	0	2
Whale	TCGAAAGCAT	4	5	2	0

A Multiple Alignment Defines a Simple Distance Matrix

Exercise: Prove that for any multiple sequence alignment, this way of defining D produces a distance matrix.

SPECIES	ALIGNMENT	DISTANCE MATRIX			
		Chimp	Human	Seal	Whale
Chimp	A C GTA G G C CT	0	3	6	4
Human	A T GTA A G A CT	3	0	7	5
Seal	TCGAGAGCAC	6	7	0	2
Whale	TCGAAAGCAT	4	5	2	0

Distance-Based Phylogeny

Distance-Based Phylogeny Problem.

- **Input:** A distance matrix.
- **Output:** The unrooted tree “fitting” this distance matrix.

Distance-Based Phylogeny

Distance-Based Phylogeny Problem.

- **Input:** A distance matrix.
- **Output:** The unrooted tree “*fitting*” this distance matrix.

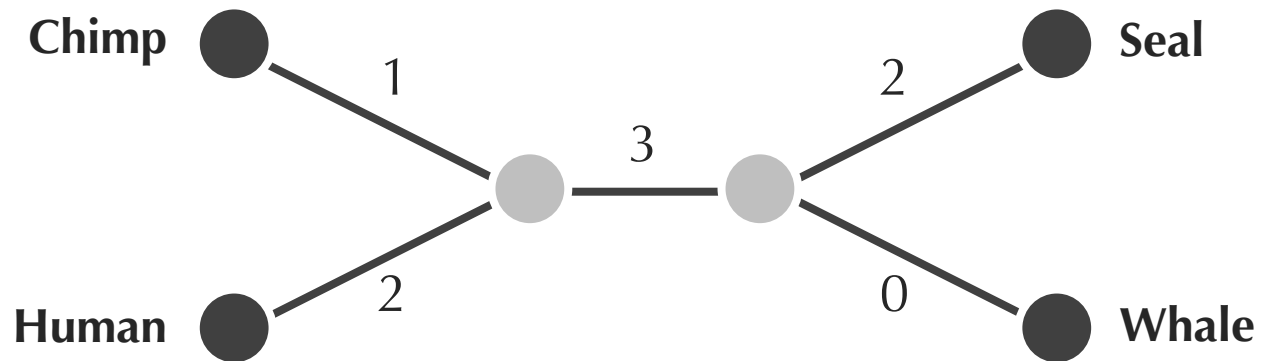
Of course, we are getting a bit ahead of ourselves – we should define what we mean by “fitting”!

“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

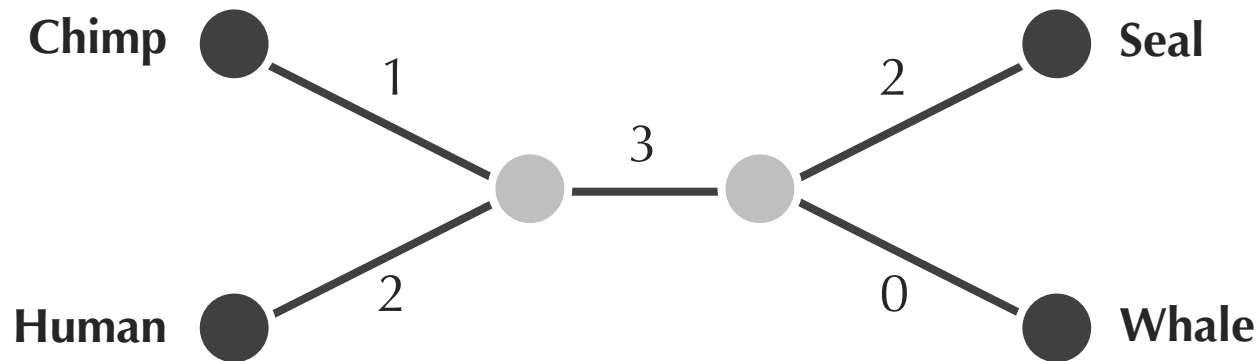
“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



“Fitting” a Tree to a Matrix

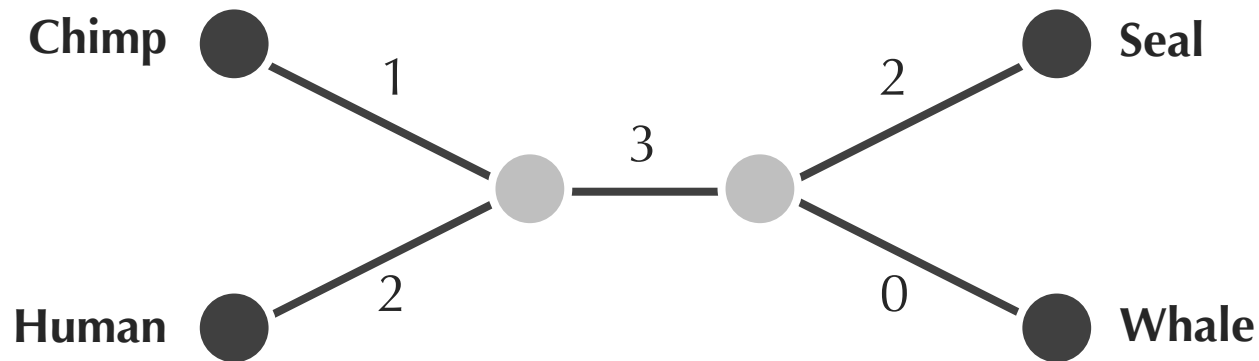
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$d_{i,j}(T)$ = distance between nodes i and j in tree T , computed by summing edge weights from i to j .

“Fitting” a Tree to a Matrix

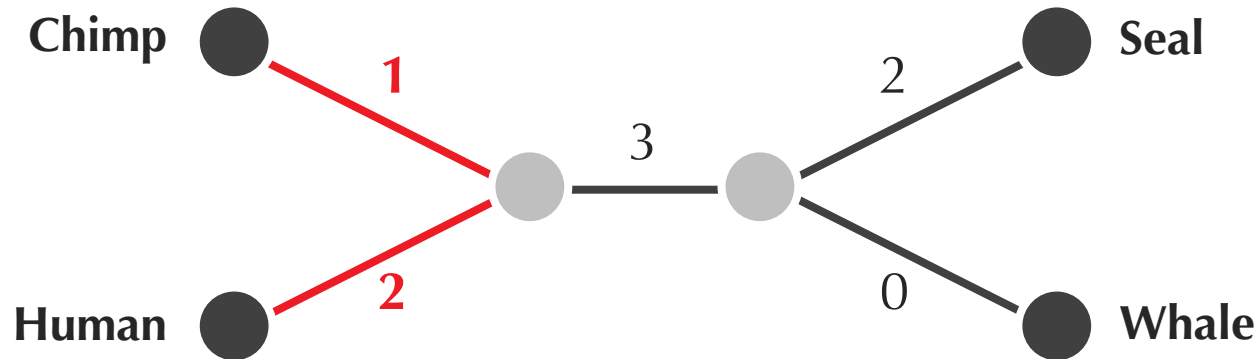
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



We say that T **fits** matrix D if for every pair i and j , $d_{i,j}(T) = D_{i,j}$.

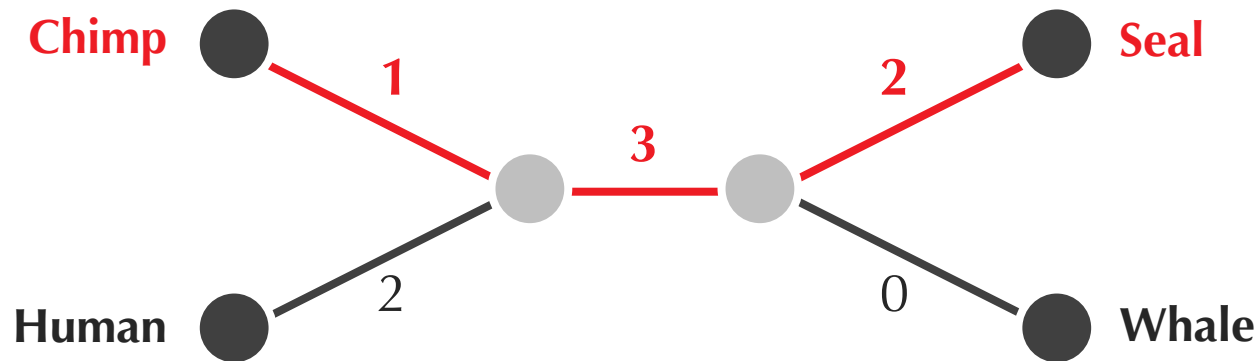
“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



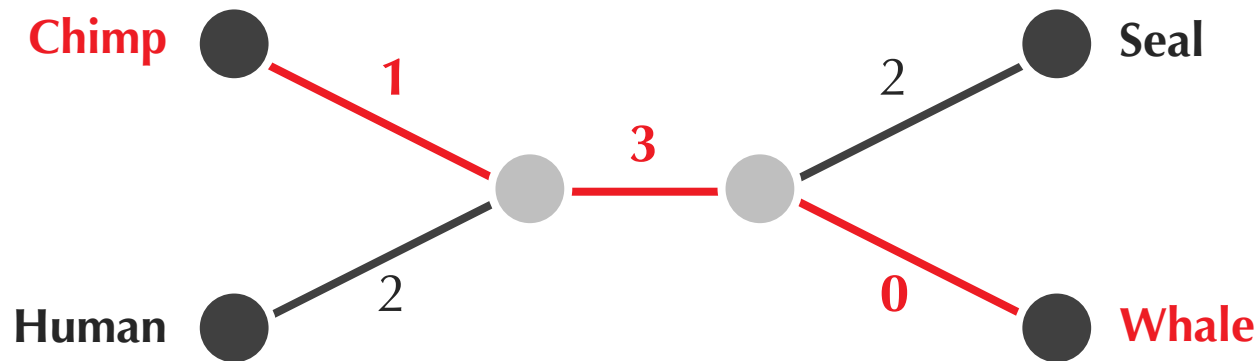
“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



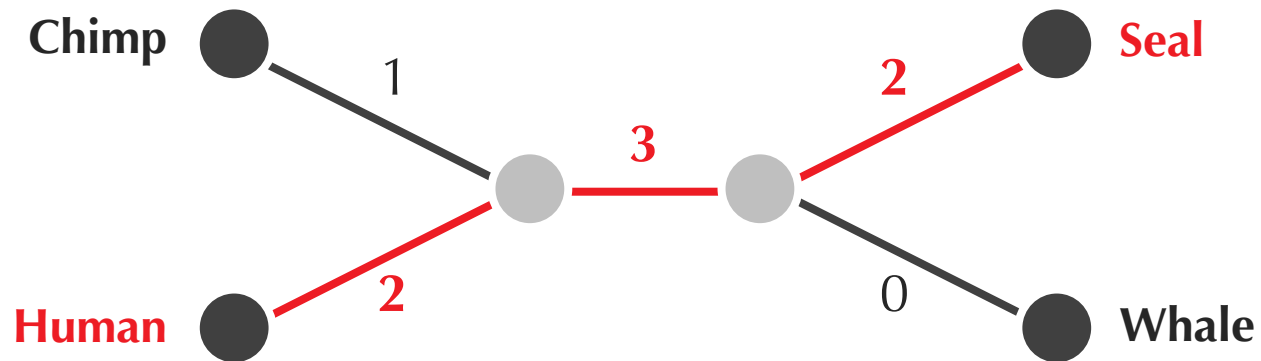
“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



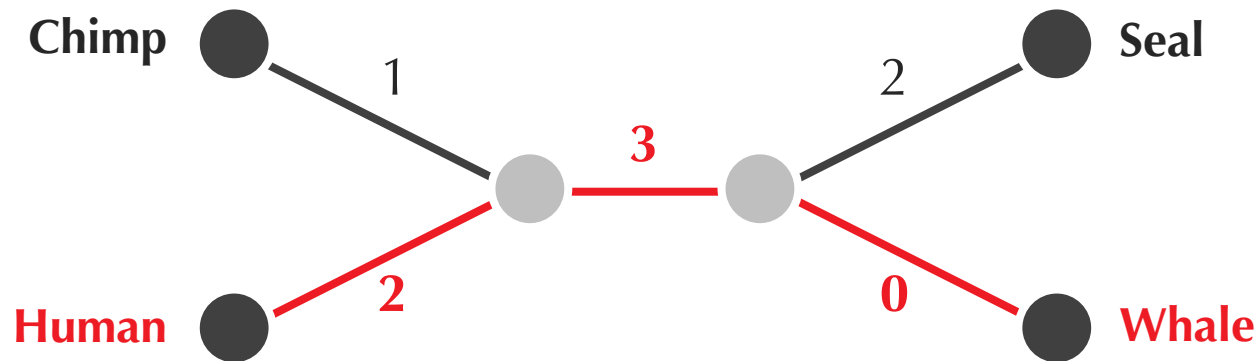
“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



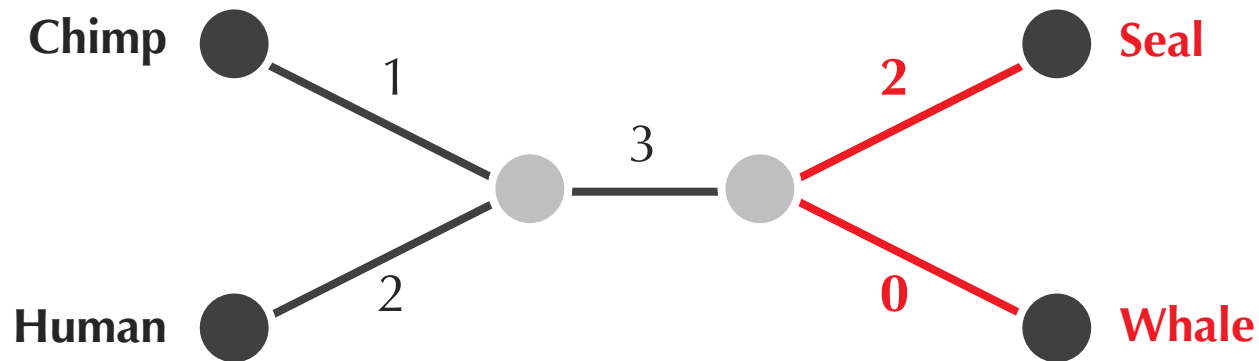
“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



“Fitting” a Tree to a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



Return to Distance-Based Phylogeny

Exercise: Find a tree fitting the following matrix.

	v_1	v_2	v_3	v_4
v_1	0	3	4	3
v_2	3	0	4	5
v_3	4	4	0	2
v_4	3	5	2	0

Sometimes, **No** Tree Fits a Matrix

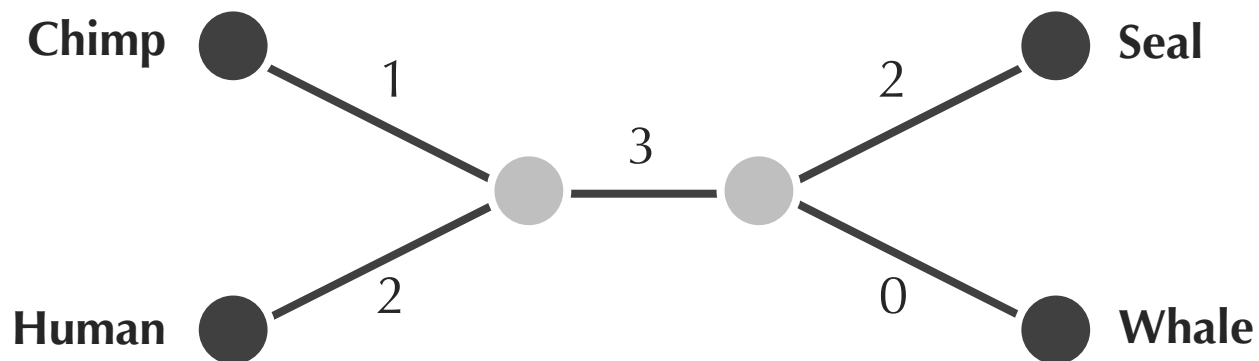
Exercise: Find a tree fitting the following matrix.

	v_1	v_2	v_3	v_4
v_1	0	3	4	3
v_2	3	0	4	5
v_3	4	4	0	2
v_4	3	5	2	0

Additive matrix: distance matrix such that there exists an unrooted tree fitting it.

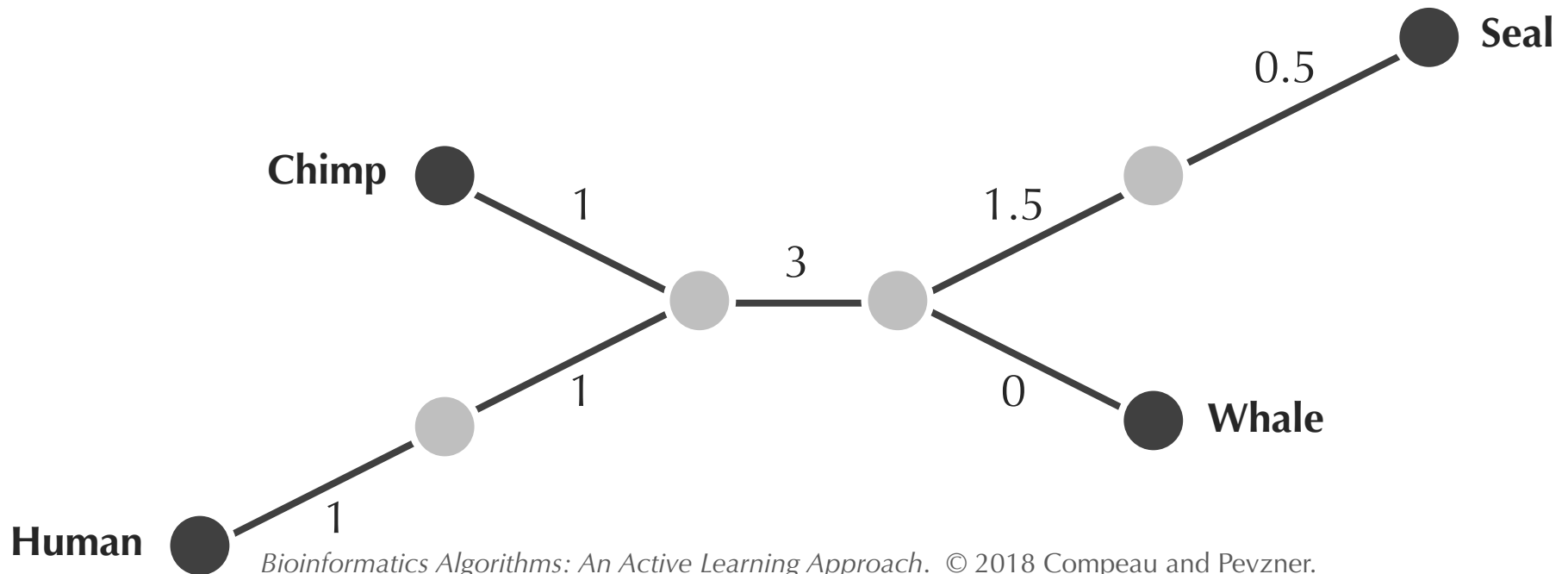
Sometimes, **More Than One** Tree Fits a Matrix

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

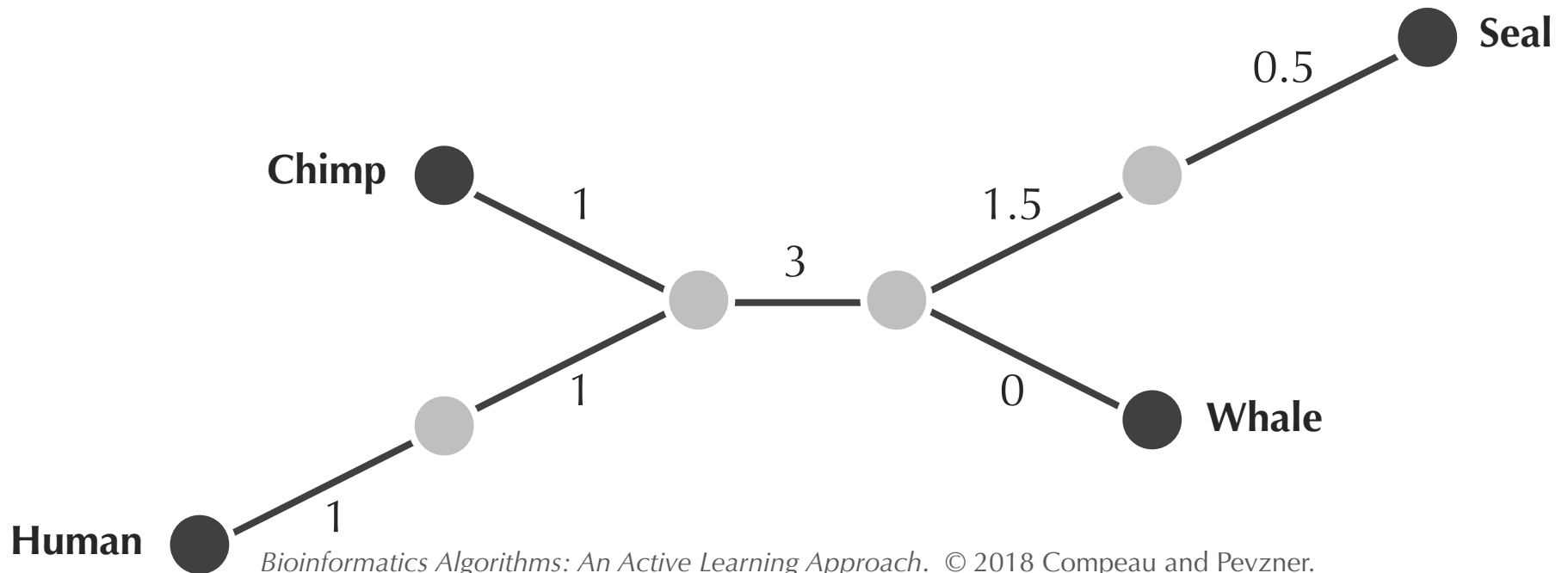
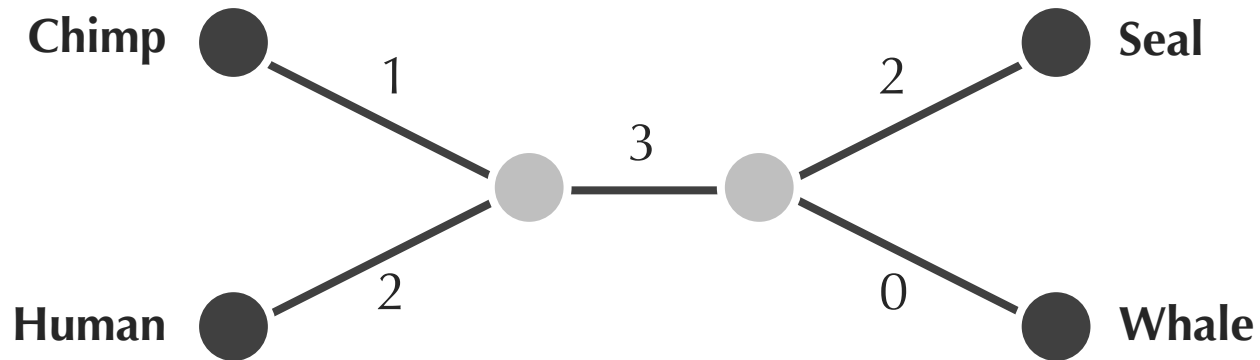


Sometimes, **More Than One** Tree Fits a Matrix

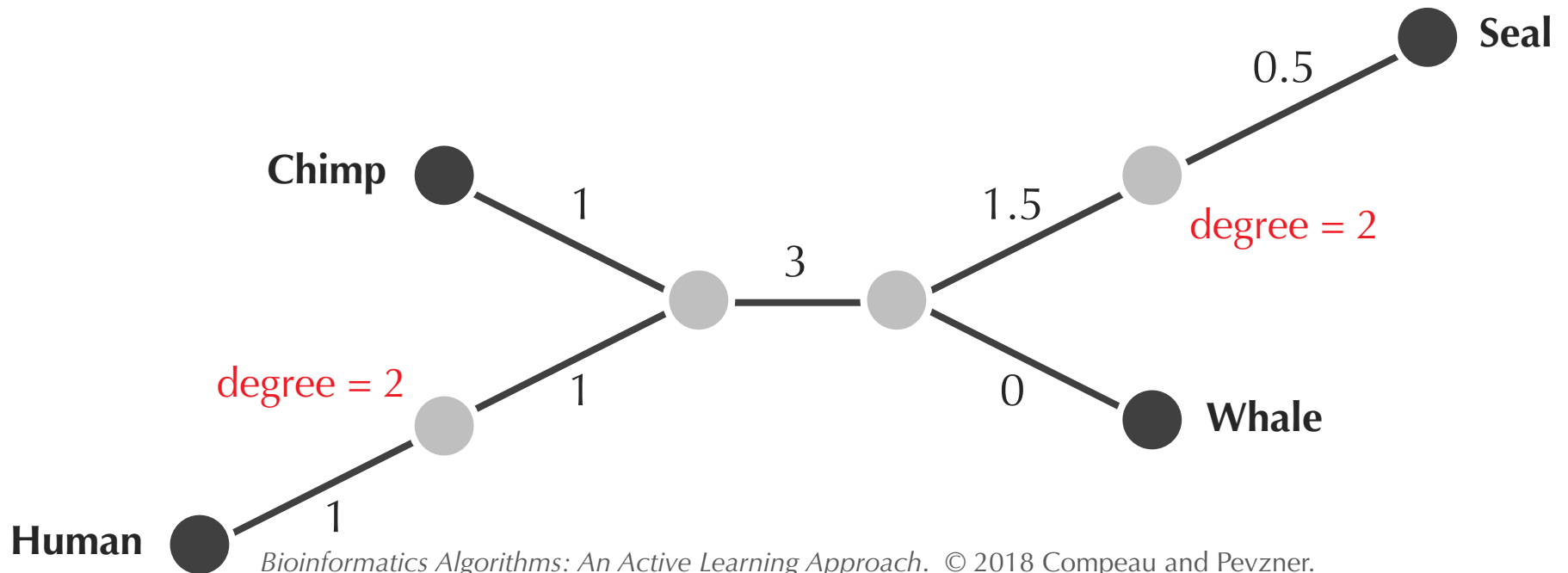
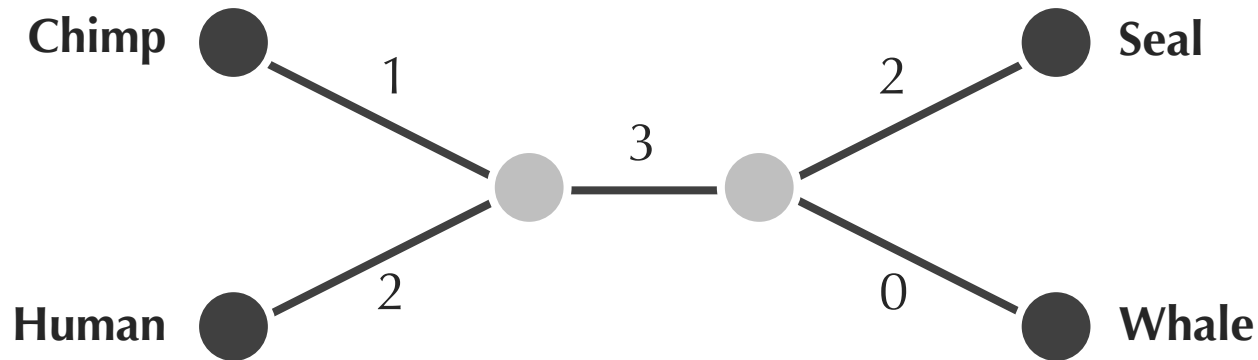
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



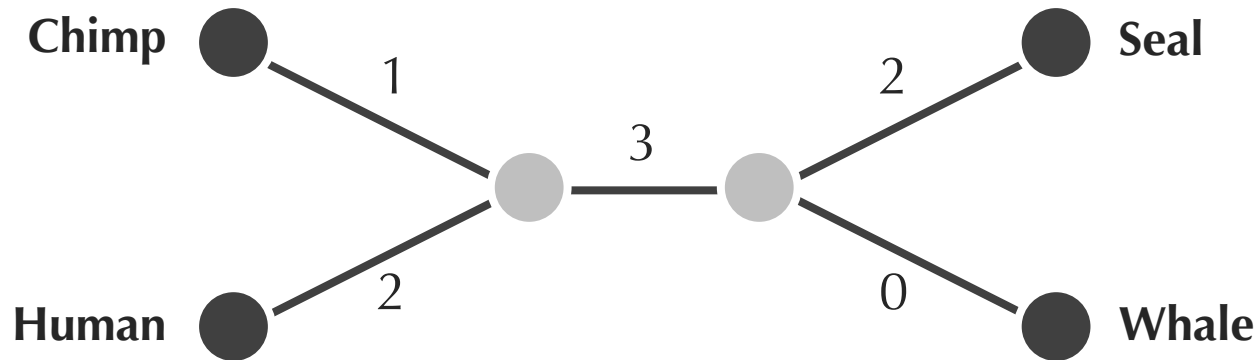
Which Tree is Better?



Which Tree is Better?

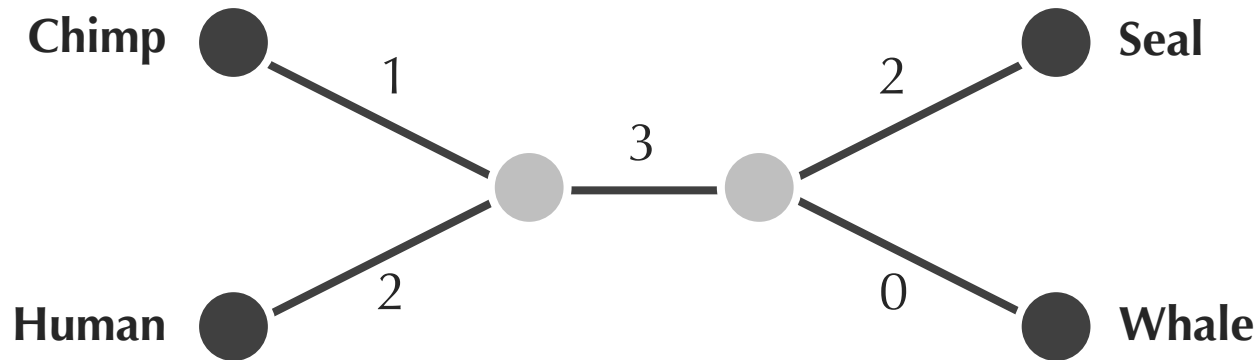


Which Tree is Better?



Simple tree: tree with no nodes of degree 2.

Which Tree is Better?



Simple tree: tree with no nodes of degree 2.

Theorem: There is a unique *simple* tree fitting an *additive* matrix.

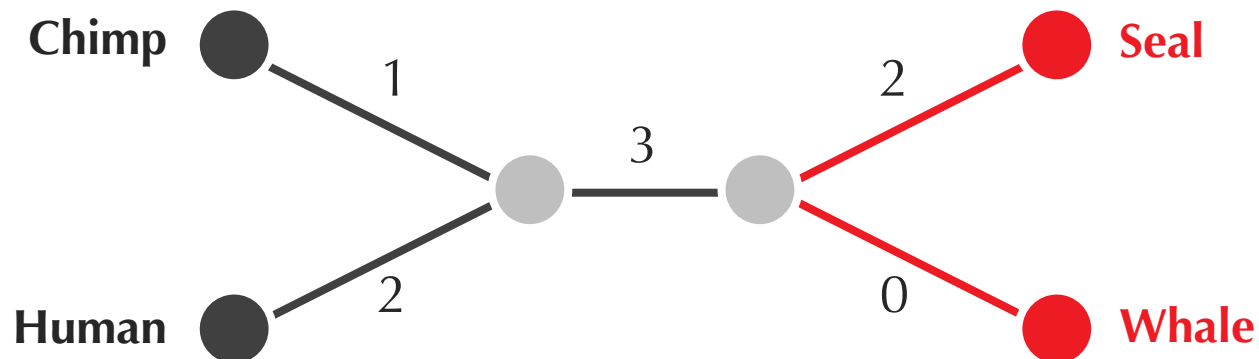
Reformulating Distance-Based Phylogeny

Distance-Based Phylogeny Problem: *Construct an evolutionary tree from a distance matrix.*

- **Input:** A distance matrix.
- **Output:** The simple tree fitting this distance matrix (if this matrix is additive).

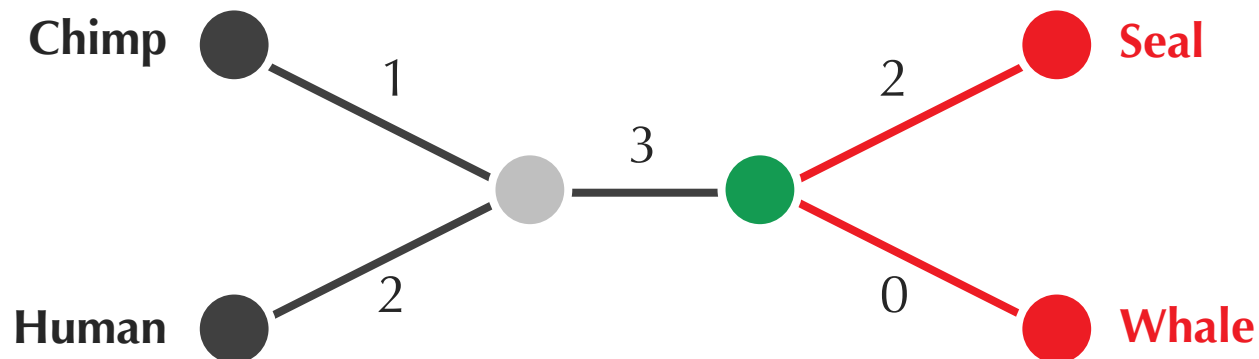
An Idea for Distance-Based Phylogeny

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



An Idea for Distance-Based Phylogeny

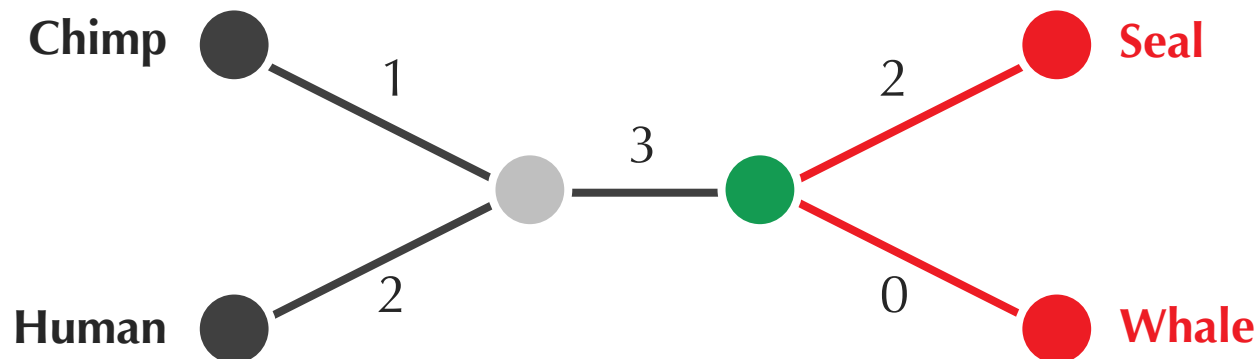
Seal and whale are **neighbors** (meaning they share the same **parent**).



An Idea for Distance-Based Phylogeny

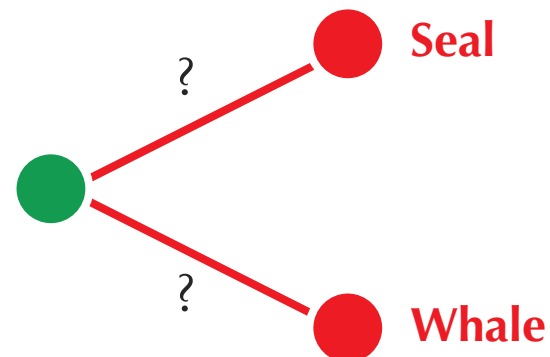
Seal and whale are **neighbors** (meaning they share the same **parent**).

Theorem: Every simple tree with at least three leaves has at least one pair of neighboring leaves.



An Idea for Distance-Based Phylogeny

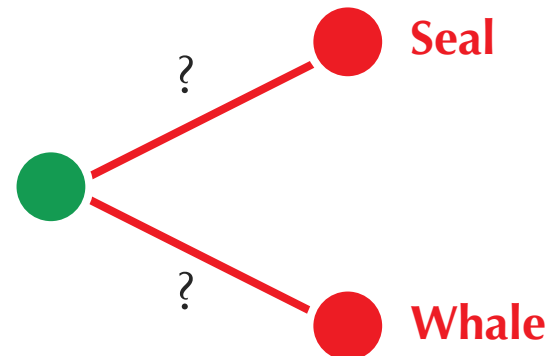
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



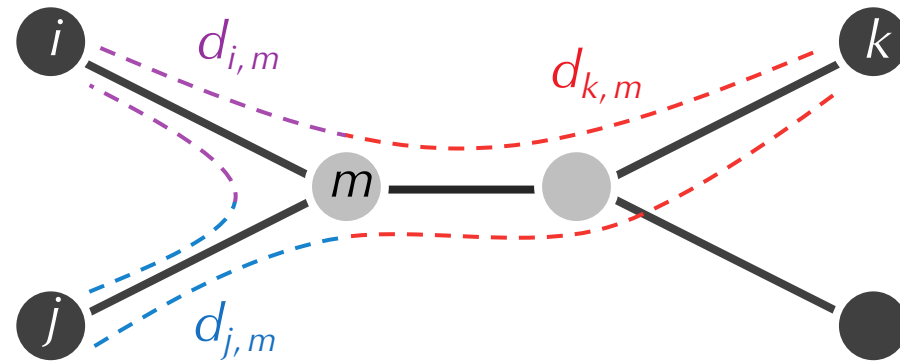
An Idea for Distance-Based Phylogeny

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

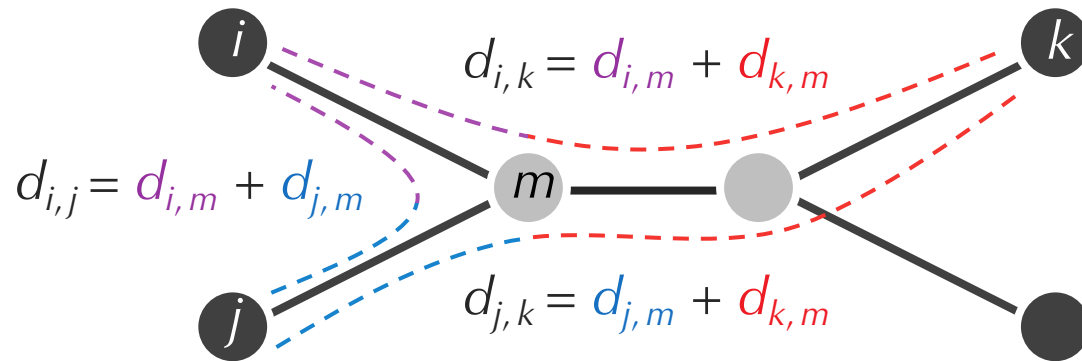
Key Point: How do we compute the unknown distances?



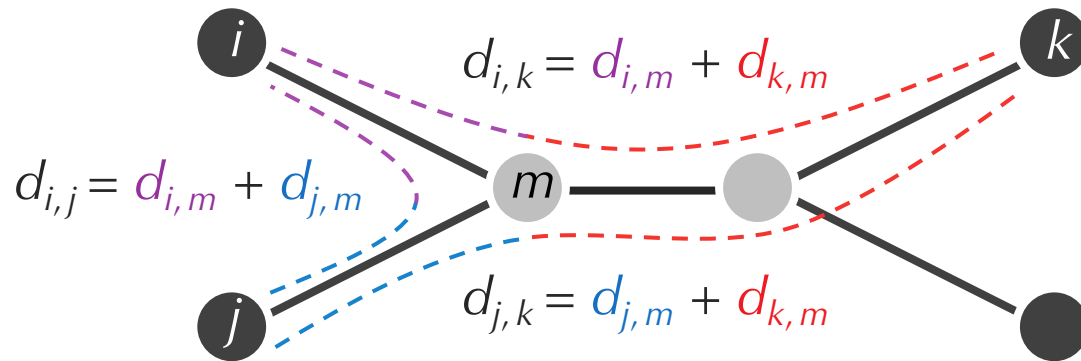
Toward a Recursive Algorithm



Toward a Recursive Algorithm

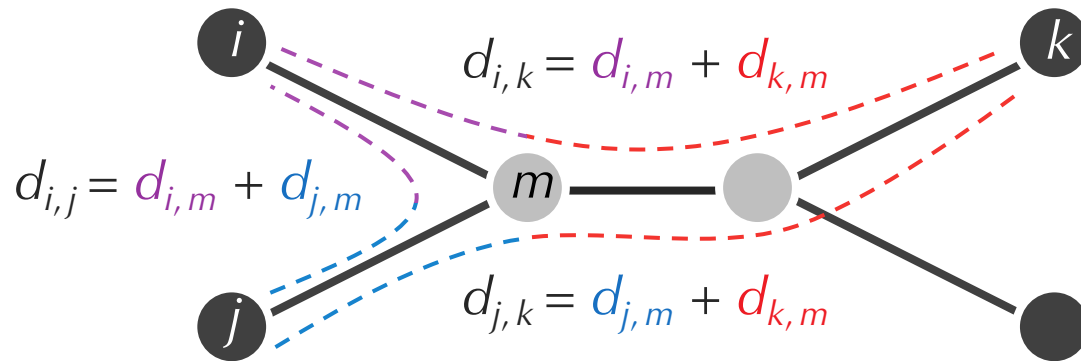


Toward a Recursive Algorithm



$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

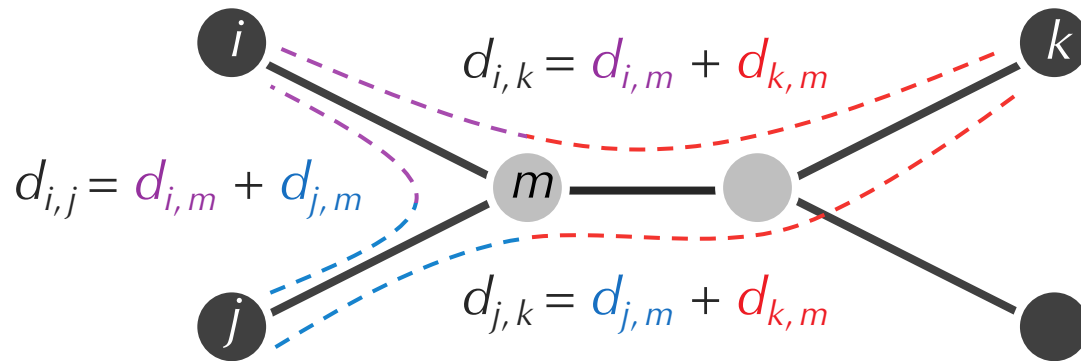
Toward a Recursive Algorithm



$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

Toward a Recursive Algorithm

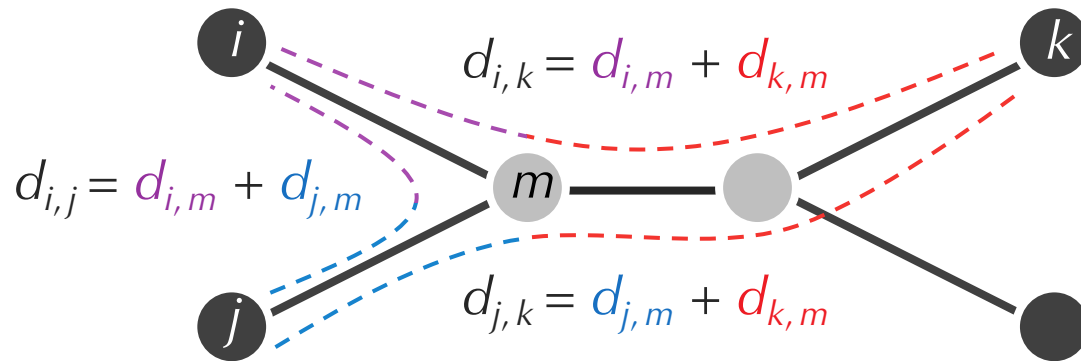


$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

$$d_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

Toward a Recursive Algorithm



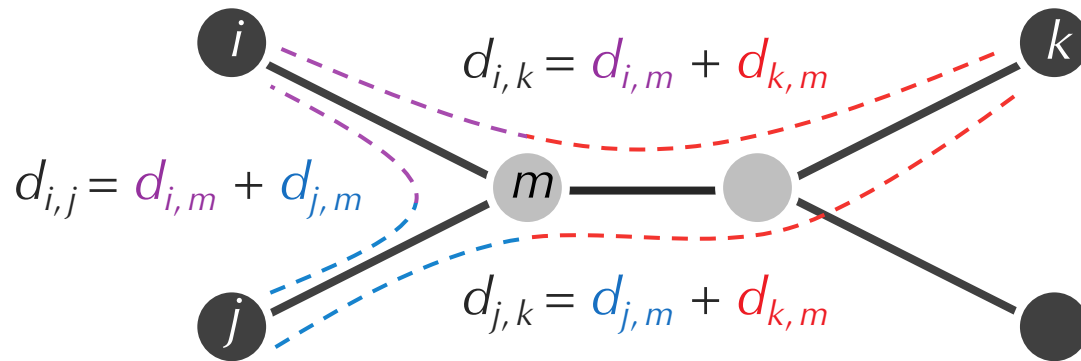
$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

$$d_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

$$\therefore d_{i,m} = D_{i,k} - (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

Toward a Recursive Algorithm



$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

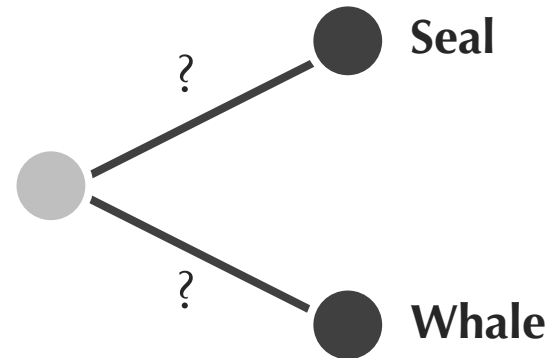
$$d_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

$$\therefore d_{i,m} = D_{i,k} - (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

$$d_{i,m} = (D_{i,k} + D_{i,j} - D_{j,k}) / 2$$

An Idea for Distance-Based Phylogeny

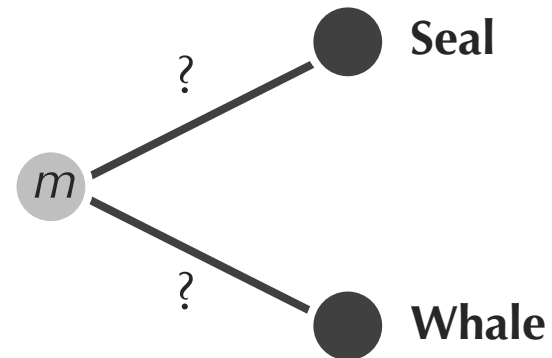
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{i,m} = (D_{i,k} + D_{i,j} - D_{j,k}) / 2$$

An Idea for Distance-Based Phylogeny

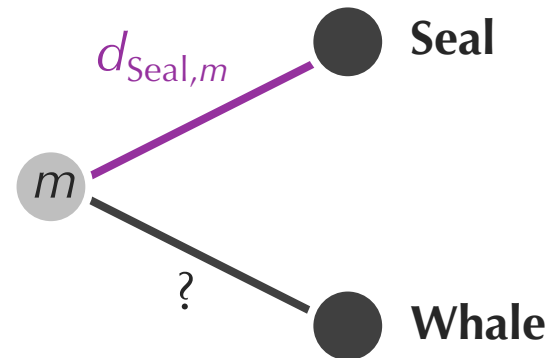
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{i,m} = (D_{i,k} + D_{i,j} - D_{j,k}) / 2$$

An Idea for Distance-Based Phylogeny

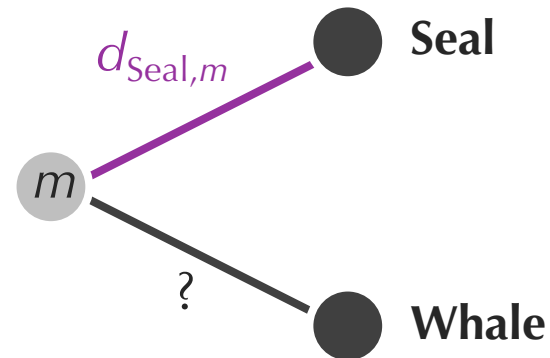
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = (D_{\text{Seal},k} + D_{\text{Seal},j} - D_{j,k}) / 2$$

An Idea for Distance-Based Phylogeny

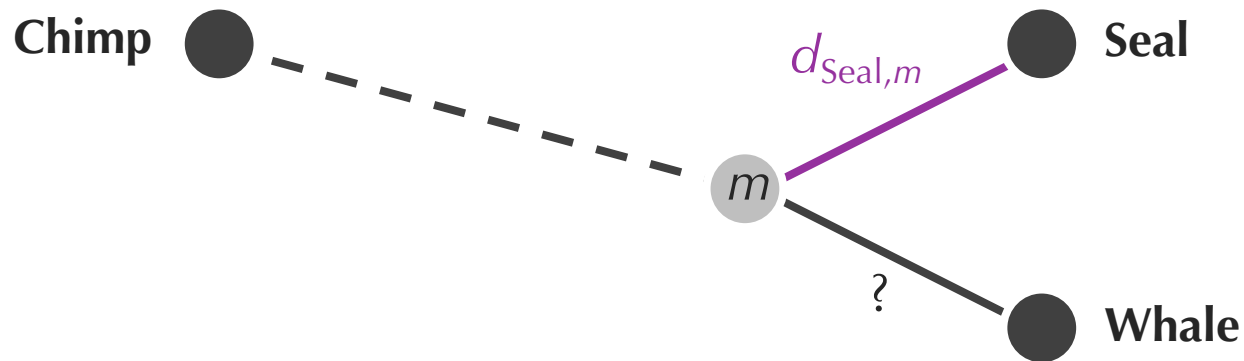
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = (D_{\text{Seal},k} + D_{\text{Seal},\text{Whale}} - D_{\text{Whale},k}) / 2$$

An Idea for Distance-Based Phylogeny

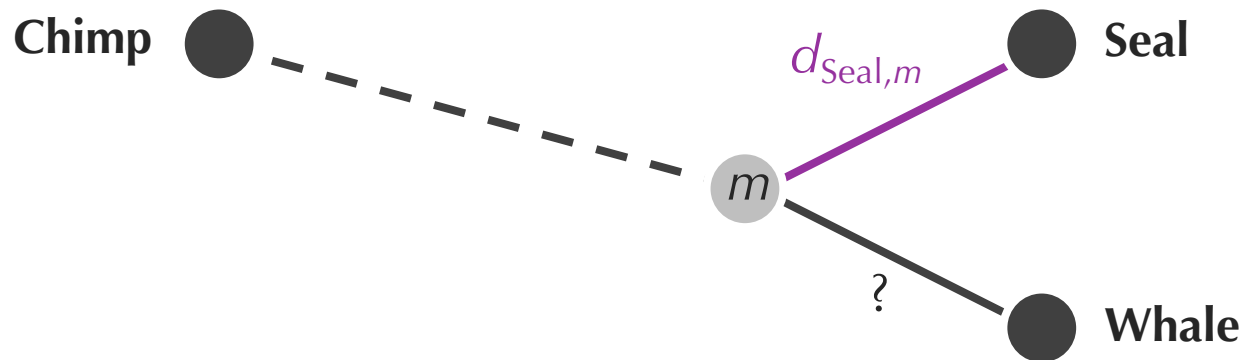
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = (D_{\text{Seal},\text{Chimp}} + D_{\text{Seal},\text{Whale}} - D_{\text{Whale},\text{Chimp}}) / 2$$

An Idea for Distance-Based Phylogeny

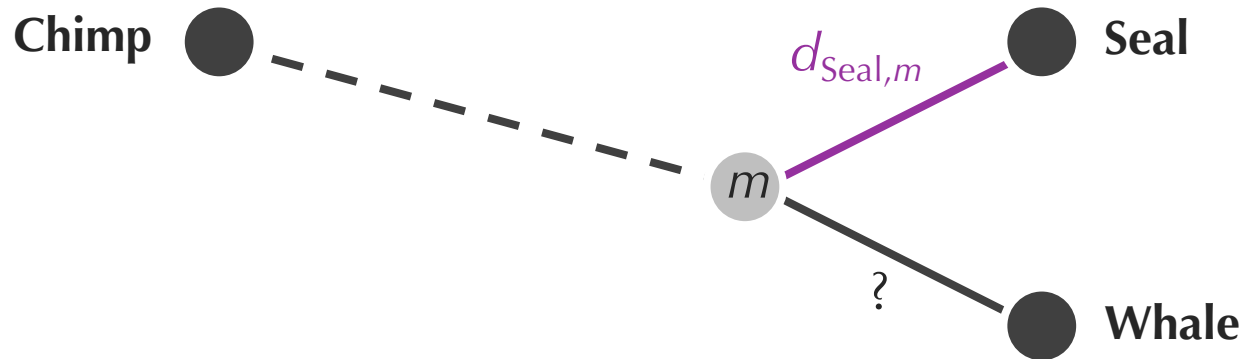
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = (\quad 6 \quad + D_{\text{Seal},\text{Whale}} - D_{\text{Whale},\text{Chimp}}) / 2$$

An Idea for Distance-Based Phylogeny

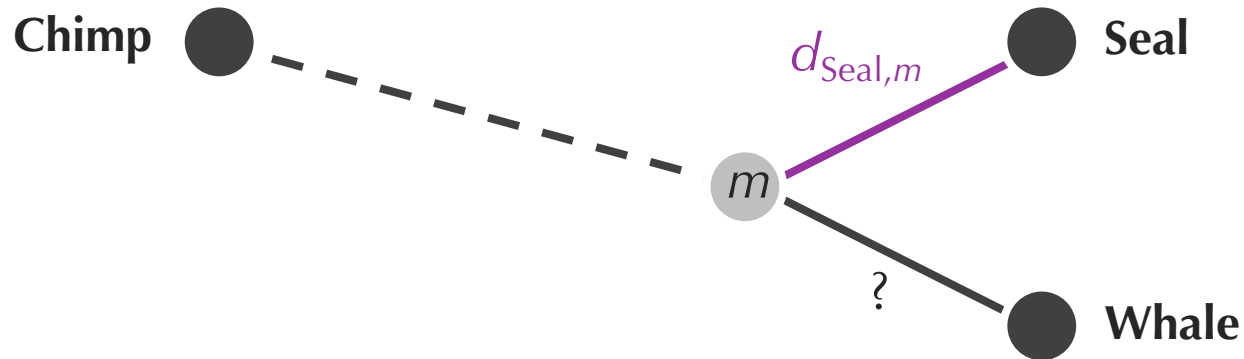
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = (6 + 2 - D_{\text{Whale,Chimp}}) / 2$$

An Idea for Distance-Based Phylogeny

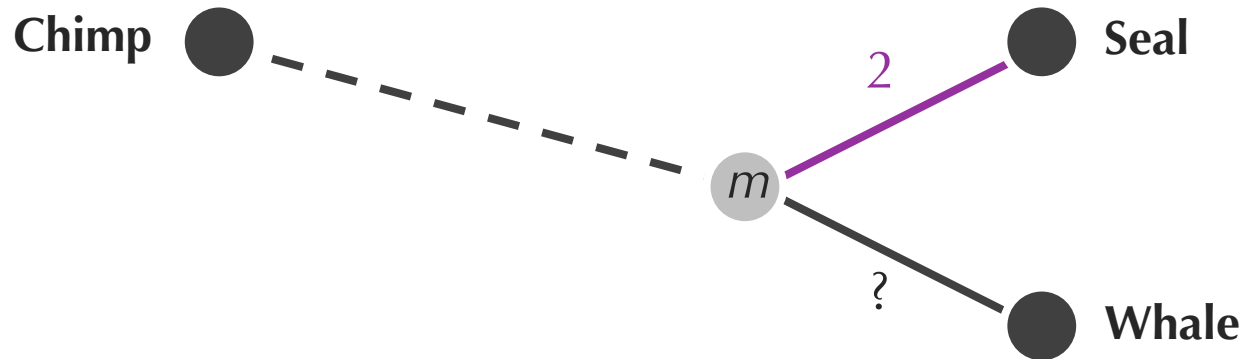
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = (6 + 2 - 4) / 2$$

An Idea for Distance-Based Phylogeny

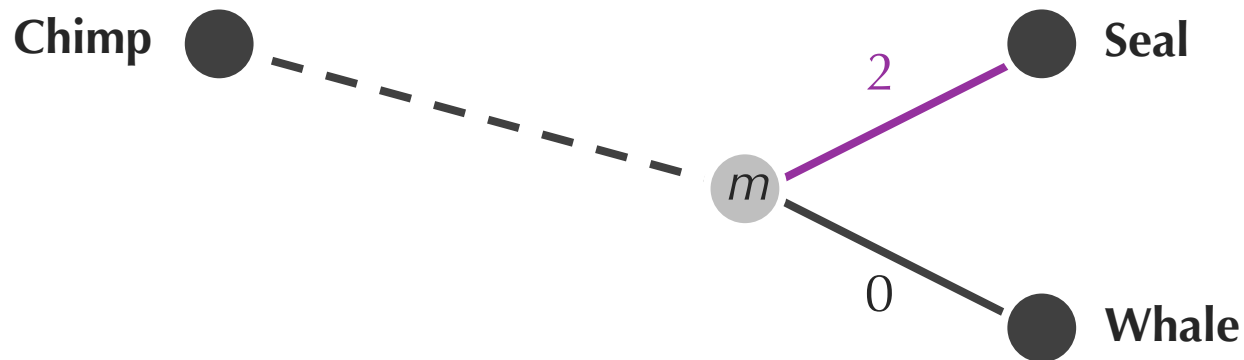
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = 2$$

An Idea for Distance-Based Phylogeny

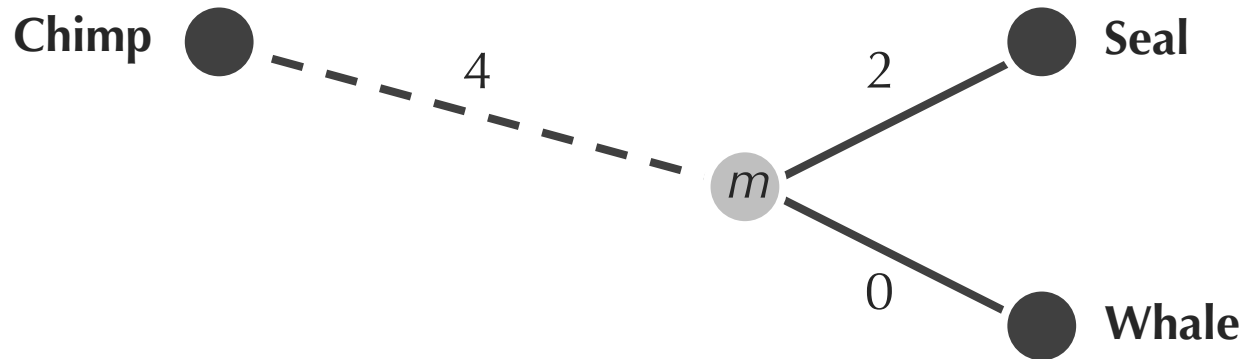
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



$$d_{\text{Seal},m} = 2$$

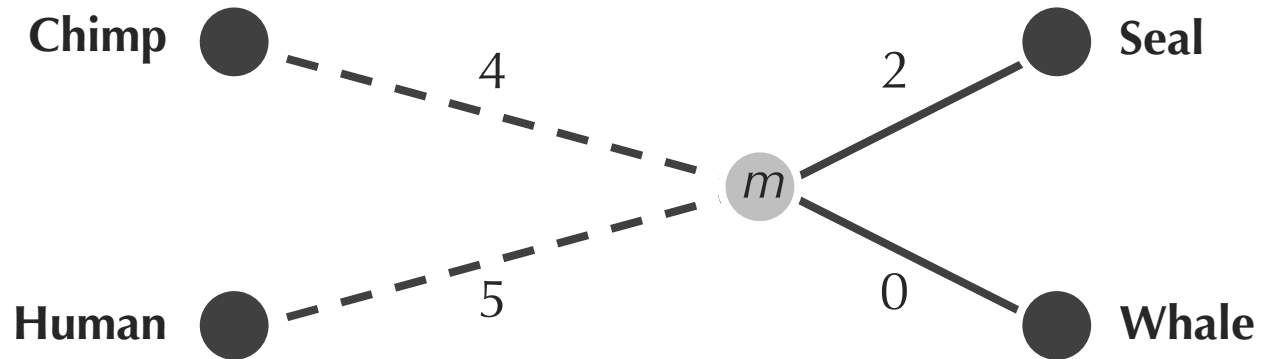
An Idea for Distance-Based Phylogeny

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



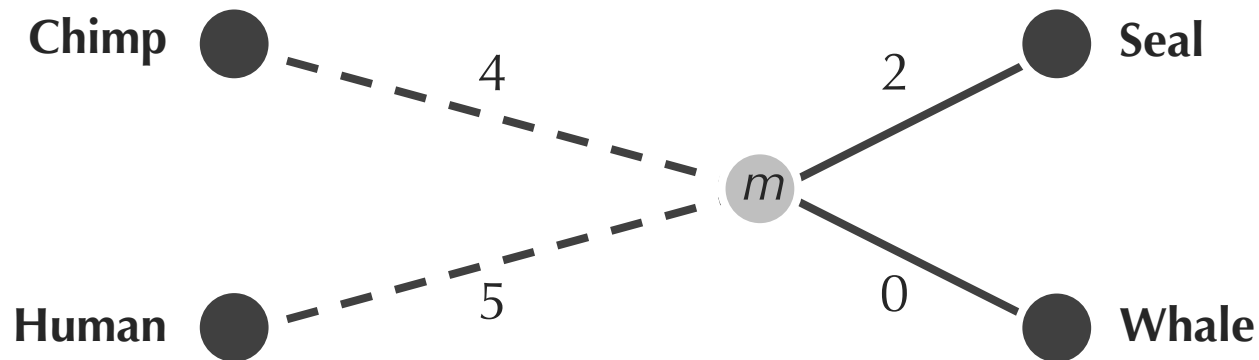
An Idea for Distance-Based Phylogeny

	Chimp	Human	Seal	Whale	<i>m</i>
Chimp	0	3	6	4	4
Human	3	0	7	5	5
Seal	6	7	0	2	2
Whale	4	5	2	0	0
<i>m</i>	4	5	2	0	0



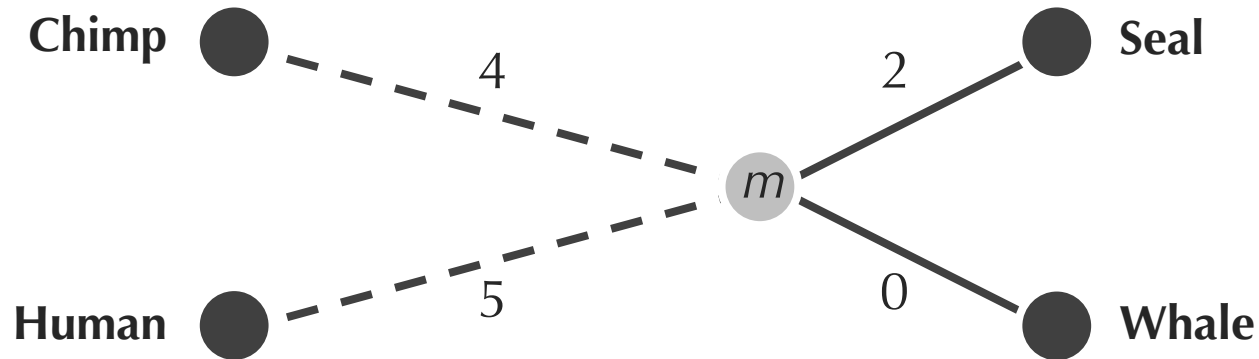
An Idea for Distance-Based Phylogeny

	Chimp	Human	Seal	Whale	<i>m</i>
Chimp	0	3	6	4	4
Human	3	0	7	5	5
Seal	6	7	0	2	2
Whale	4	5	2	0	0
<i>m</i>	4	5	2	0	0



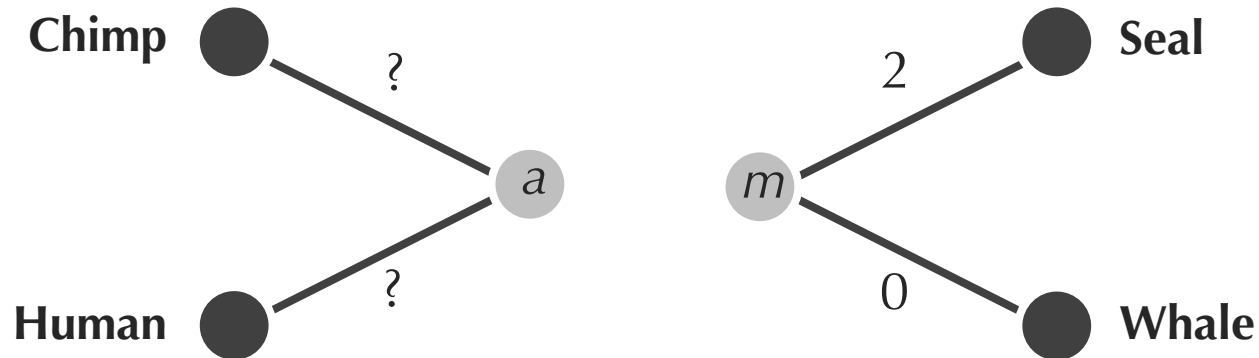
An Idea for Distance-Based Phylogeny

	Chimp	Human	m
Chimp	0	3	4
Human	3	0	5
m	4	5	0



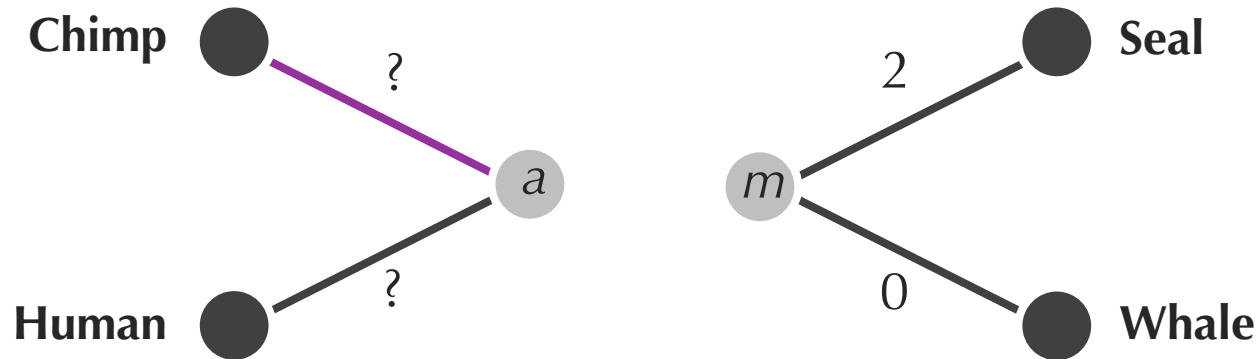
An Idea for Distance-Based Phylogeny

	Chimp	Human	m
Chimp	0	3	4
Human	3	0	5
m	4	5	0



An Idea for Distance-Based Phylogeny

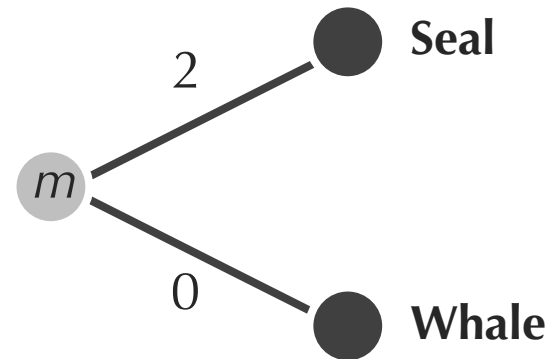
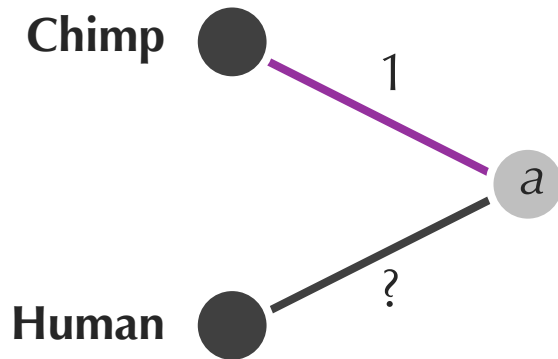
	Chimp	Human	m
Chimp	0	3	4
Human	3	0	5
m	4	5	0



$$d_{\text{Chimp},a} = (D_{\text{Chimp},m} + D_{\text{Chimp},\text{Human}} - D_{\text{Human},m}) / 2$$

An Idea for Distance-Based Phylogeny

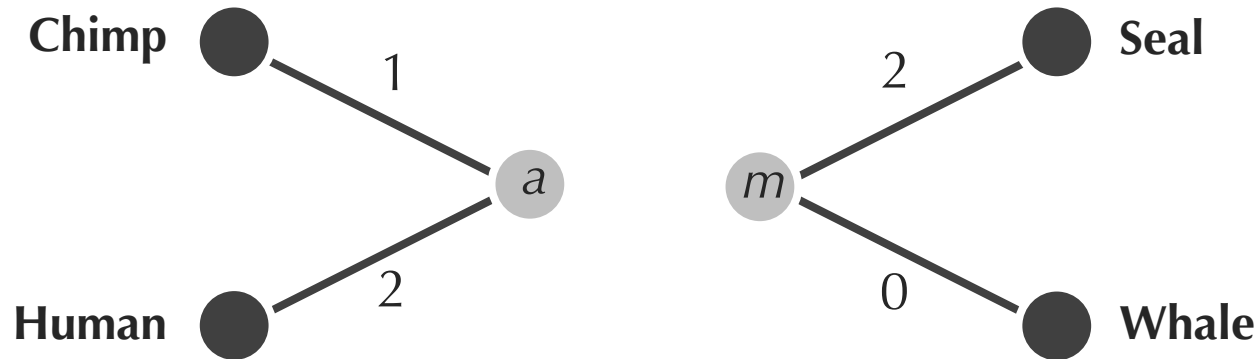
	Chimp	Human	m
Chimp	0	3	4
Human	3	0	5
m	4	5	0



$$d_{\text{Chimp},a} = 1$$

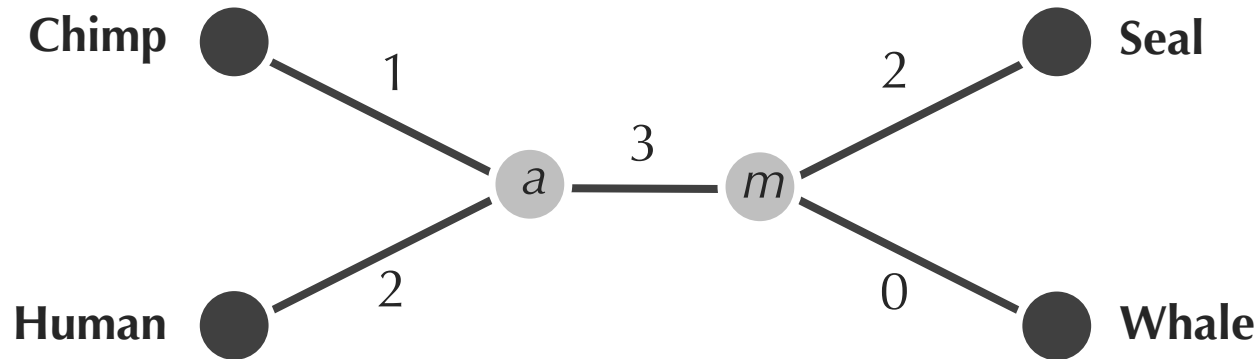
An Idea for Distance-Based Phylogeny

	Chimp	Human	m
Chimp	0	3	4
Human	3	0	5
m	4	5	0



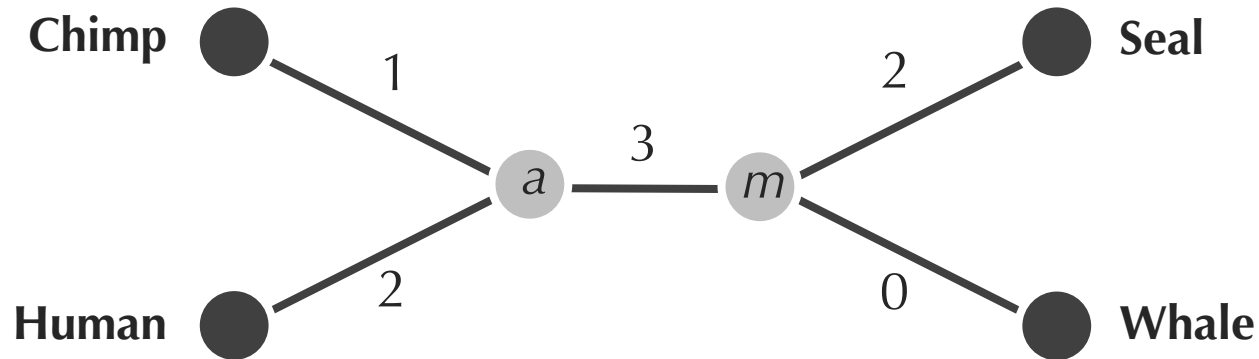
An Idea for Distance-Based Phylogeny

	Chimp	Human	<i>m</i>
Chimp	0	3	4
Human	3	0	5
<i>m</i>	4	5	0



An Idea for Distance-Based Phylogeny

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



An Idea for Distance-Based Phylogeny

	0	1	2	3
0	0	13	21	22
1	13	0	12	13
2	21	12	0	13
3	22	13	13	0

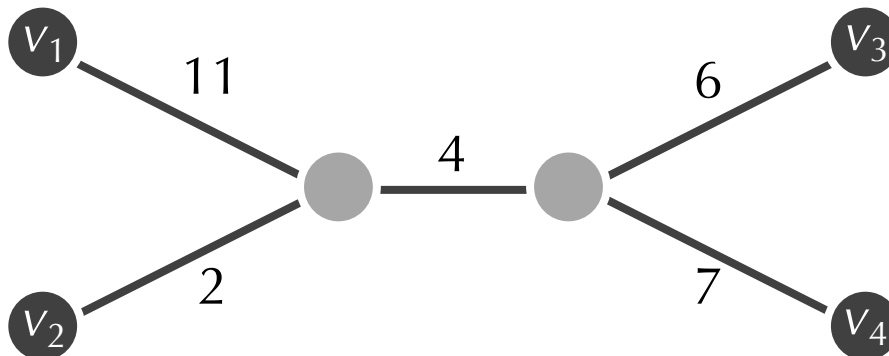
Exercise: Apply this recursive approach to this distance matrix.

What Was Wrong With Our Algorithm?

	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0

What Was Wrong With Our Algorithm?

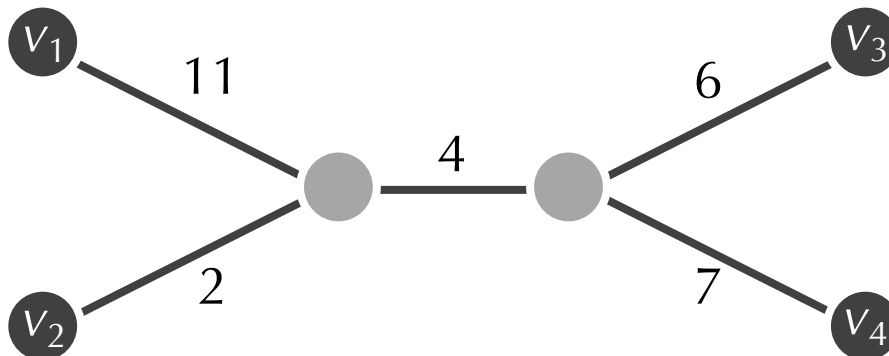
	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0



What Was Wrong With Our Algorithm?

	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0

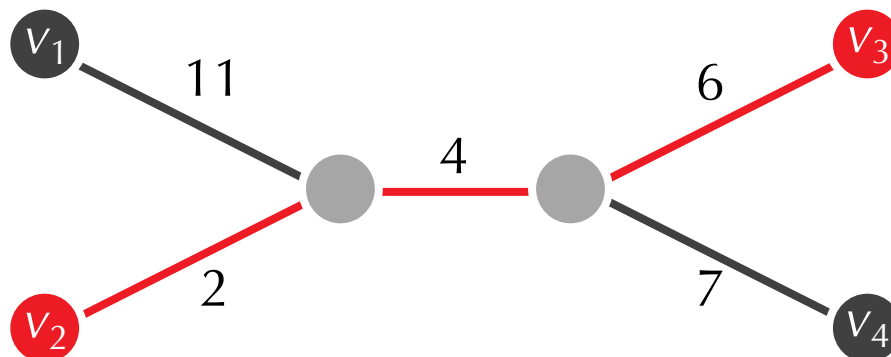
minimum
element is $D_{2,3}$



What Was Wrong With Our Algorithm?

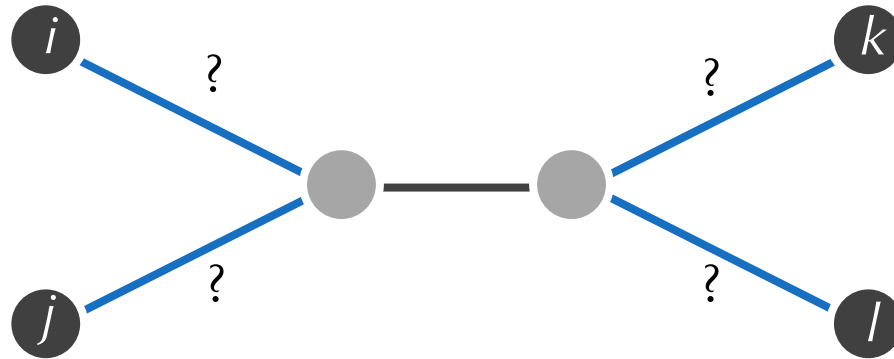
	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0

minimum
element is $D_{2,3}$



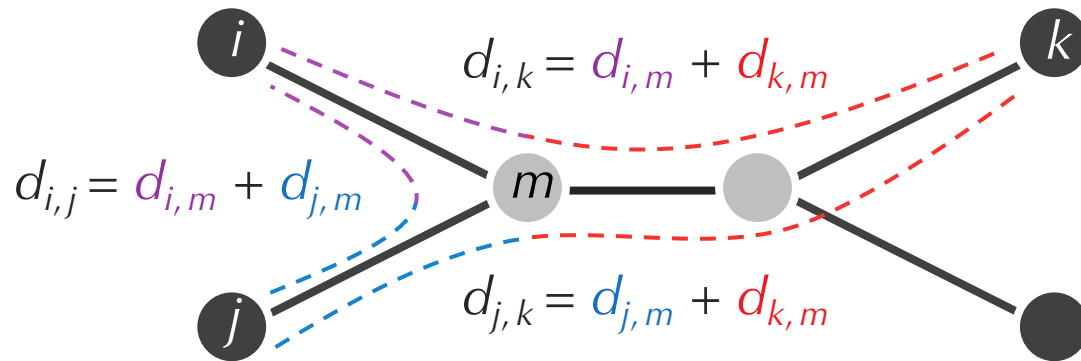
v_2 and v_3 are
not neighbors!

From Neighbors to Limbs



Rather than trying to infer **neighbors**, let's instead try to compute the length of **limbs**, the edges attached to leaves.

From Neighbors to Limbs



$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

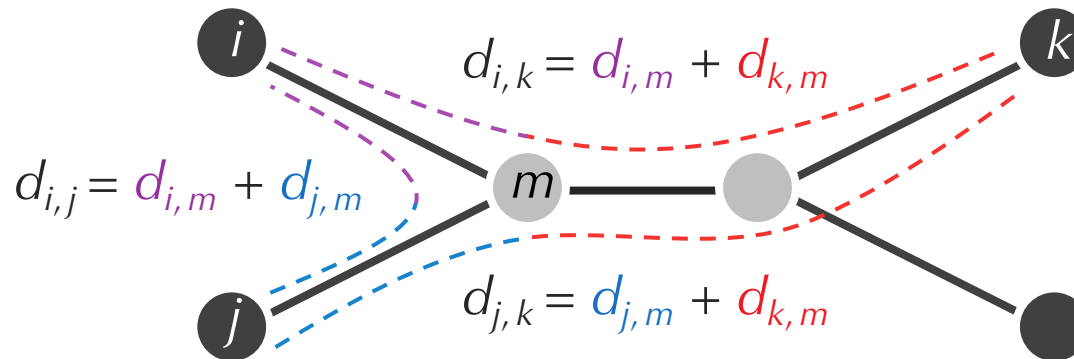
$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

$$d_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

$$\therefore d_{i,m} = D_{i,k} - (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

$$d_{i,m} = (D_{i,k} + D_{i,j} - D_{j,k}) / 2$$

From Neighbors to Limbs



$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

$$d_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

$$\therefore d_{i,m} = D_{i,k} - (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

$$d_{i,m} = (D_{i,k} + D_{i,j} - D_{j,k}) / 2$$

Assumes that i and j are neighbors...

Computing Limb Lengths

Limb Length Theorem: $LimbLength(i)$ is equal to the minimum value of $(D_{i,k} + D_{i,j} - D_{j,k})/2$ over all leaves j and k .

Computing Limb Lengths

Limb Length Theorem: $LimbLength(chimp)$ is equal to the minimum value of $(D_{chimp,k} + D_{chimp,j} - D_{j,k})/2$ over all leaves j and k .

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

Computing Limb Lengths

Limb Length Theorem: $LimbLength(chimp)$ is equal to the minimum value of $(D_{chimp,k} + D_{chimp,j} - D_{j,k})/2$ over all leaves j and k .

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

$$(D_{chimp, human} + D_{chimp, seal} - D_{human, seal}) / 2 = (3 + 6 - 7) / 2 = 1$$

Computing Limb Lengths

Limb Length Theorem: $LimbLength(chimp)$ is equal to the minimum value of $(D_{chimp,k} + D_{chimp,j} - D_{j,k})/2$ over all leaves j and k .

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

$$(D_{chimp, human} + D_{chimp, seal} - D_{human, seal}) / 2 = (3 + 6 - 7) / 2 = 1$$

$$(D_{chimp, human} + D_{chimp, whale} - D_{human, whale}) / 2 = (3 + 4 - 5) / 2 = 1$$

Computing Limb Lengths

Limb Length Theorem: $LimbLength(chimp)$ is equal to the minimum value of $(D_{chimp,k} + D_{chimp,j} - D_{j,k})/2$ over all leaves j and k .

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

$$(D_{chimp, human} + D_{chimp, seal} - D_{human, seal}) / 2 = (3 + 6 - 7) / 2 = 1$$

$$(D_{chimp, human} + D_{chimp, whale} - D_{human, whale}) / 2 = (3 + 4 - 5) / 2 = 1$$

$$(D_{chimp, whale} + D_{chimp, seal} - D_{whale, seal}) / 2 = (6 + 4 - 2) / 2 = 4$$

Computing Limb Lengths

Limb Length Theorem: $LimbLength(chimp)$ is equal to the minimum value of $(D_{chimp,k} + D_{chimp,j} - D_{j,k})/2$ over all leaves j and k .

	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0

$$(D_{chimp, human} + D_{chimp, seal} - D_{human, seal}) / 2 = (3 + 6 - 7) / 2 = \mathbf{1}$$

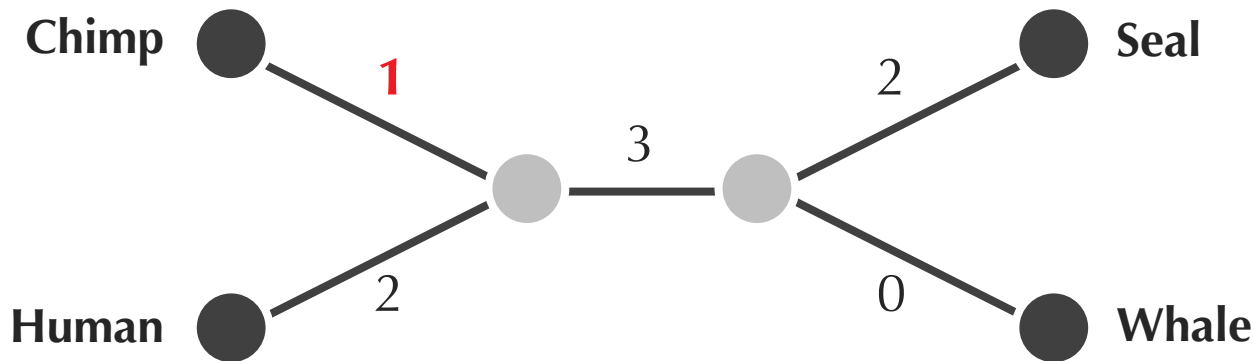
$$(D_{chimp, human} + D_{chimp, whale} - D_{human, whale}) / 2 = (3 + 4 - 5) / 2 = \mathbf{1}$$

$$(D_{chimp, whale} + D_{chimp, seal} - D_{whale, seal}) / 2 = (6 + 4 - 2) / 2 = 4$$

Computing Limb Lengths

Limb Length Theorem: $LimbLength(chimp)$ is equal to the minimum value of $(D_{chimp,k} + D_{chimp,j} - D_{j,k})/2$ over all leaves j and k .

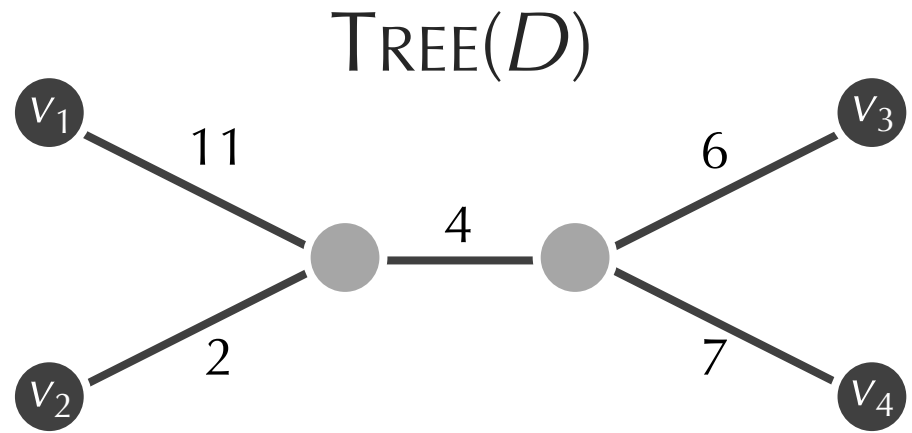
	Chimp	Human	Seal	Whale
Chimp	0	3	6	4
Human	3	0	7	5
Seal	6	7	0	2
Whale	4	5	2	0



AdditivePhylogeny In Action

D

	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0



AdditivePhylogeny In Action

		v_1	v_2	v_3	v_4
D	v_1	0	13	21	22
	v_2	13	0	12	13
	v_3	21	12	0	13
	v_4	22	13	13	0

1. Pick an arbitrary leaf j (say, $j = v_2$).

AdditivePhylogeny In Action

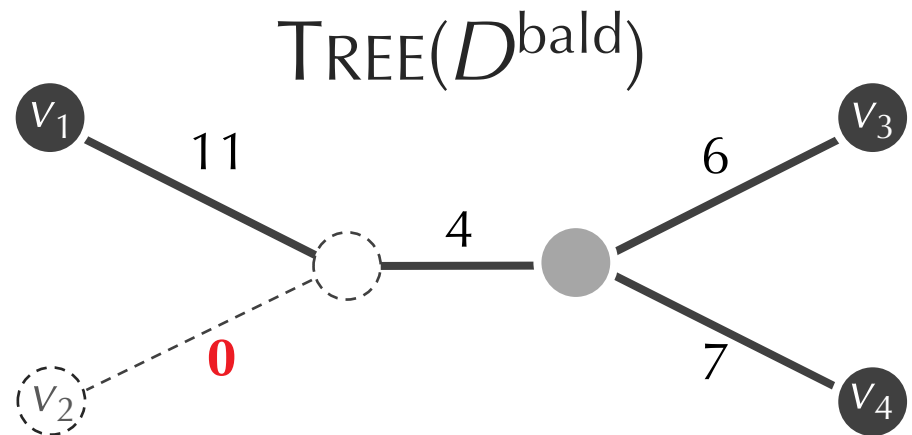
		v_1	v_2	v_3	v_4
D	v_1	0	13	21	22
	v_2	13	0	12	13
	v_3	21	12	0	13
	v_4	22	13	13	0

$$\text{LimbLength}(v_2) = 2$$

2. Compute its limb length, $\text{LimbLength}(j)$.

AdditivePhylogeny In Action

	v_1	v_2	v_3	v_4
D^{bald}				
v_1	0	11	21	22
v_2	11	0	10	11
v_3	21	10	0	13
v_4	22	11	13	0



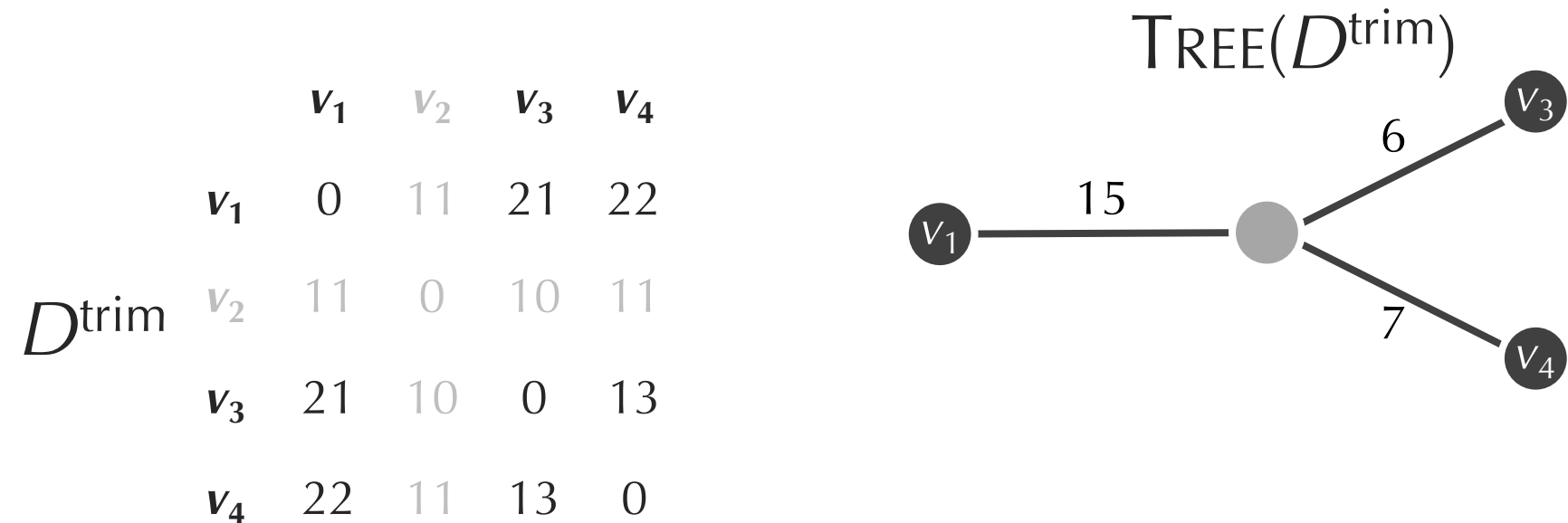
3. Subtract $\text{LimbLength}(j)$ from each row and column to produce D^{bald} in which j is a **bald limb** (length 0).

AdditivePhylogeny In Action

		v_1	v_2	v_3	v_4
D^{trim}	v_1	0	11	21	22
	v_2	11	0	10	11
	v_3	21	10	0	13
	v_4	22	11	13	0

4. Remove the j -th row and column of the matrix to form the $(n - 1) \times (n - 1)$ matrix D^{trim} .

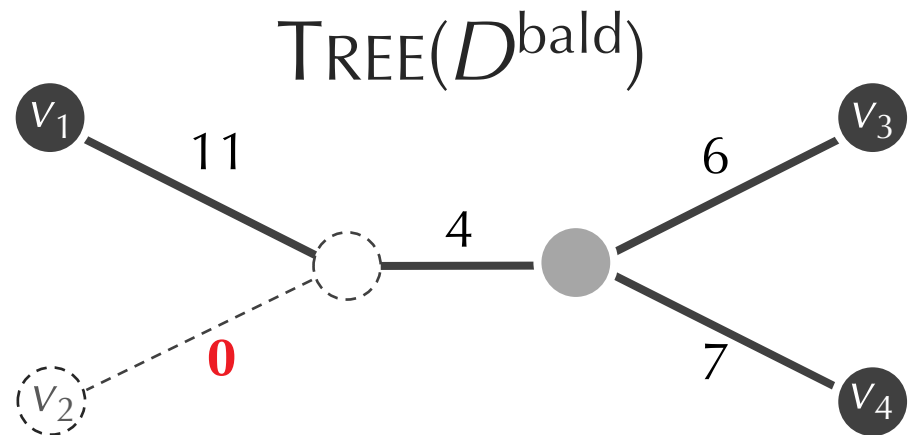
AdditivePhylogeny In Action



5. Call AdditivePhylogeny recursively to construct $\text{Tree}(D^{\text{trim}})$.

AdditivePhylogeny In Action

	v_1	v_2	v_3	v_4
v_1	0	11	21	22
v_2	11	0	10	11
v_3	21	10	0	13
v_4	22	11	13	0



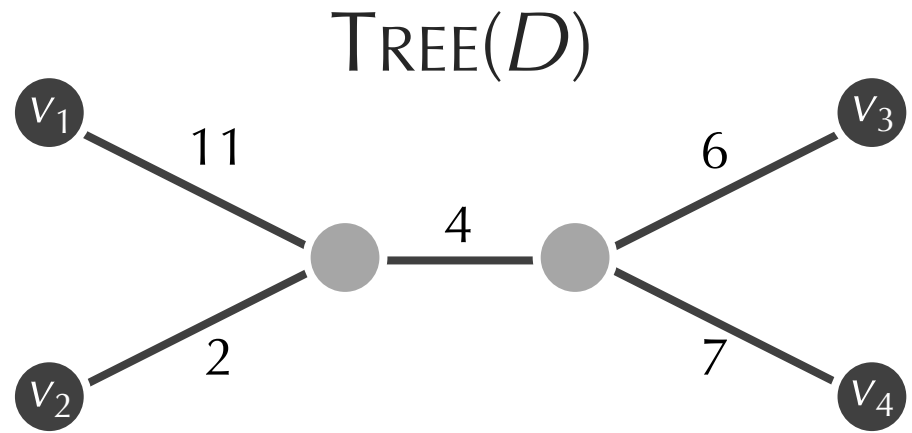
6. Identify the point in $\text{Tree}(D^{\text{trim}})$ where leaf j should be attached.

AdditivePhylogeny In Action

D

	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0

$$\text{LimbLength}(v_2) = 2$$



7. Attach j by an edge of length $\text{LimbLength}(j)$ in order to form $\text{Tree}(D)$.

AdditivePhylogeny

AdditivePhylogeny(D):

1. Pick an arbitrary leaf j .
2. Compute its limb length, $LimbLength(j)$.
3. Subtract $LimbLength(j)$ from each row and column to produce D^{bald} in which j is a bald limb (length 0).
4. Remove the j -th row and column of the matrix to form the $(n - 1) \times (n - 1)$ matrix D^{trim} .
5. Recursively call **AdditivePhylogeny**(D^{trim}) to obtain $Tree(D^{trim})$.
6. **Identify the point in $Tree(D^{trim})$ where leaf j should be attached.**
7. Attach j by an edge of length $LimbLength(j)$ in order to form $Tree(D)$.

Attaching a leaf to the tree after the recursive step is the difficult part ...

AdditivePhylogeny

AdditivePhylogeny(D):

1. Pick an arbitrary leaf j .
2. Compute its limb length, $LimbLength(j)$.
3. Subtract $LimbLength(j)$ from each row and column to produce D^{bald} in which j is a bald limb (length 0).
4. Remove the j -th row and column of the matrix to form the $(n - 1) \times (n - 1)$ matrix D^{trim} .
5. Recursively call **AdditivePhylogeny(D^{trim})** to obtain $Tree(D^{trim})$.
6. **Identify the point in $Tree(D^{trim})$ where leaf j should be attached.**
7. Attach j by an edge of length $LimbLength(j)$ in order to form $Tree(D)$.

Attaching a leaf to the tree after the recursive step is the difficult part ...

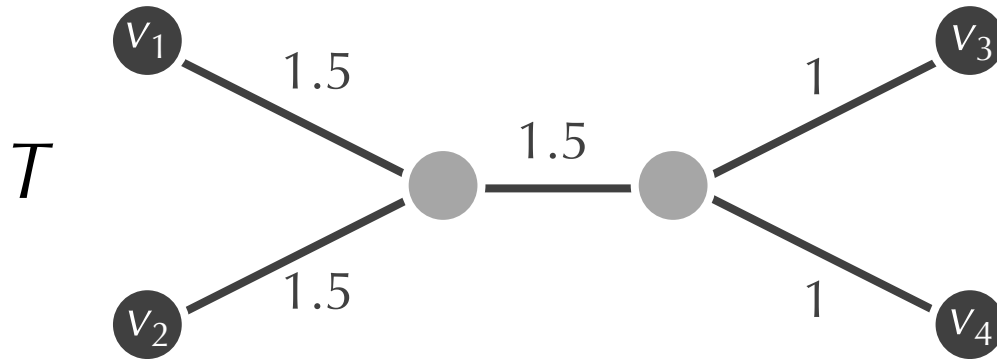
... and what do we do about *non-additive* matrices?

Least-Squares Phylogeny

Idea: Maybe we can find a tree that fits a non-additive matrix *approximately*.

Least-Squares Phylogeny

Idea: Maybe we can find a tree that fits a non-additive matrix *approximately*.

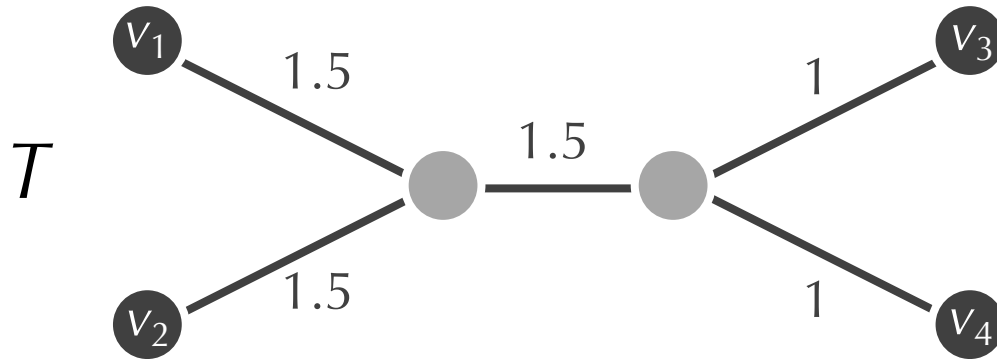


D

	v_1	v_2	v_3	v_4
v_1	0	3	4	3
v_2	3	0	4	5
v_3	4	4	0	2
v_4	3	5	2	0

Least-Squares Phylogeny

Idea: Maybe we can find a tree that fits a non-additive matrix *approximately*.

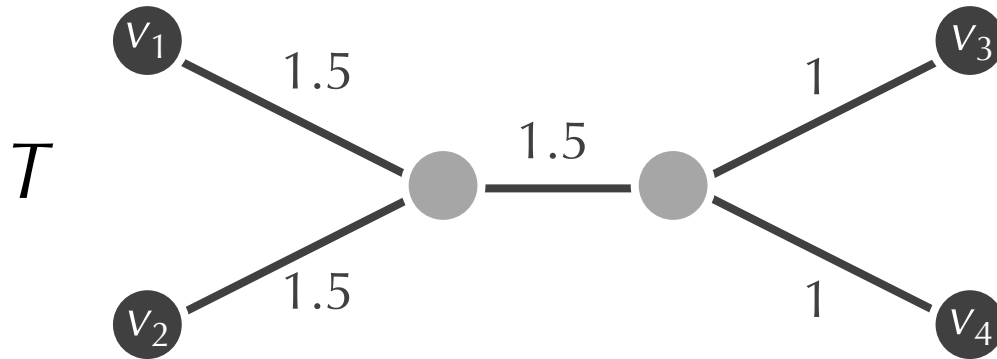


	v_1	v_2	v_3	v_4
v_1	0	3	4	3
v_2	3	0	4	5
v_3	4	4	0	2
v_4	3	5	2	0

	v_1	v_2	v_3	v_4
v_1	0	3	4	4
v_2	3	0	4	4
v_3	4	4	0	2
v_4	4	4	2	0

Least-Squares Phylogeny

Idea: Maybe we can find a tree that fits a non-additive matrix *approximately*.



D

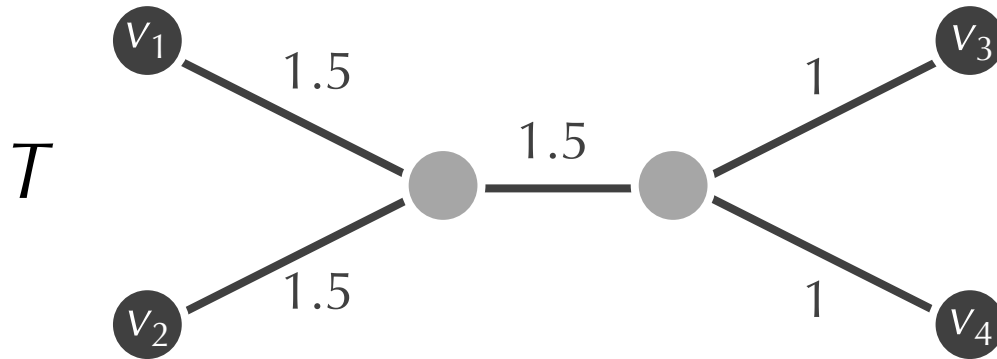
	v_1	v_2	v_3	v_4
v_1	0	3	4	3
v_2	3	0	4	5
v_3	4	4	0	2
v_4	3	5	2	0

d

	v_1	v_2	v_3	v_4
v_1	0	3	4	4
v_2	3	0	4	4
v_3	4	4	0	2
v_4	4	4	2	0

Least-Squares Phylogeny

$$\text{Discrepancy}(T, D) = \sum_{1 \leq i < j \leq n} (d_{i,j}(T) - D_{i,j})^2$$



D

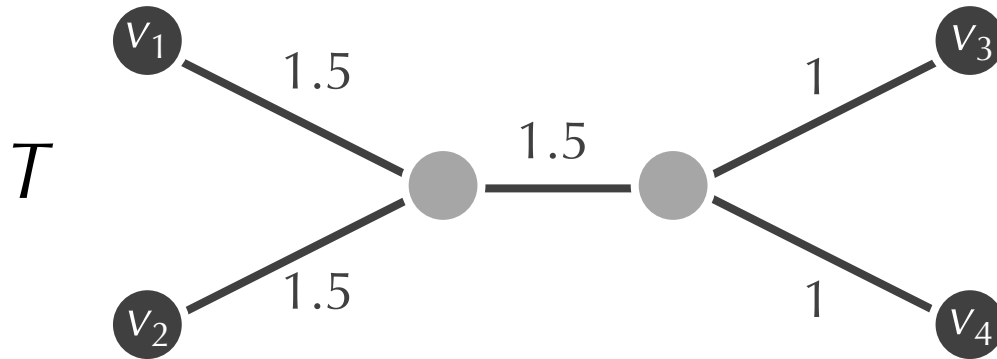
	v_1	v_2	v_3	v_4
v_1	0	3	4	3
v_2	3	0	4	5
v_3	4	4	0	2
v_4	3	5	2	0

d

	v_1	v_2	v_3	v_4
v_1	0	3	4	4
v_2	3	0	4	4
v_3	4	4	0	2
v_4	4	4	2	0

Least-Squares Phylogeny

$$\begin{aligned} \text{Discrepancy}(T, D) &= \sum_{1 \leq i < j \leq n} (d_{i,j}(T) - D_{i,j})^2 \\ &= 1^2 + 1^2 = 2 \end{aligned}$$



D

	v_1	v_2	v_3	v_4
v_1	0	3	4	3
v_2	3	0	4	5
v_3	4	4	0	2
v_4	3	5	2	0

d

	v_1	v_2	v_3	v_4
v_1	0	3	4	4
v_2	3	0	4	4
v_3	4	4	0	2
v_4	4	4	2	0

Least-Squares Phylogeny

Least-Squares Distance-Based Phylogeny Problem:

Given a distance matrix, find the tree that minimizes the sum of squared errors.

- **Input:** An $n \times n$ distance matrix D .
- **Output:** A weighted tree T with n leaves minimizing $Discrepancy(T, D)$ over all weighted trees with n leaves.

Least-Squares Phylogeny

Least-Squares Distance-Based Phylogeny Problem:

Given a distance matrix, find the tree that minimizes the sum of squared errors.

- **Input:** An $n \times n$ distance matrix D .
- **Output:** A weighted tree T with n leaves minimizing $Discrepancy(T, D)$ over all weighted trees with n leaves.

Unfortunately, this problem is *NP*-Complete...

The Critical Insight

Our first idea of using neighbors to construct the tree was a good one! The problem was not with our idea, but with the distance matrix itself!

The Neighbor-Joining Theorem

Given an $n \times n$ distance matrix D , its **neighbor-joining matrix** is the matrix D^* defined as

$$D^*_{i,j} = (n - 2) \bullet D_{i,j} - TotalDistance_D(i) - TotalDistance_D(j)$$

where $TotalDistance_D(i)$ is the sum of distances from i to all other leaves.

The Neighbor-Joining Theorem

Given an $n \times n$ distance matrix D , its **neighbor-joining matrix** is the matrix D^* defined as

$$D^*_{i,j} = (n - 2) \bullet D_{i,j} - TotalDistance_D(i) - TotalDistance_D(j)$$

where $TotalDistance_D(i)$ is the sum of distances from i to all other leaves.

	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0

The Neighbor-Joining Theorem

Given an $n \times n$ distance matrix D , its **neighbor-joining matrix** is the matrix D^* defined as

$$D^*_{i,j} = (n - 2) \bullet D_{i,j} - TotalDistance_D(i) - TotalDistance_D(j)$$

where $TotalDistance_D(i)$ is the sum of distances from i to all other leaves.

	v_1	v_2	v_3	v_4	$TotalDistance_D$
D	0	13	21	22	56
	13	0	12	13	38
	21	12	0	13	46
	22	13	13	0	48

The Neighbor-Joining Theorem

Given an $n \times n$ distance matrix D , its **neighbor-joining matrix** is the matrix D^* defined as

$$D^*_{i,j} = (n - 2) \bullet D_{i,j} - TotalDistance_D(i) - TotalDistance_D(j)$$

where $TotalDistance_D(i)$ is the sum of distances from i to all other leaves.

		v_1	v_2	v_3	v_4	$TotalDistance_D$			v_1	v_2	v_3	v_4
D	v_1	0	13	21	22	56	D^*	v_1	0	-68	-60	-60
	v_2	13	0	12	13	38		v_2	-68	0	-60	-60
	v_3	21	12	0	13	46		v_3	-60	-60	0	-68
	v_4	22	13	13	0	48		v_4	-60	-60	-68	0

The Neighbor-Joining Theorem

Given an $n \times n$ distance matrix D , its **neighbor-joining matrix** is the matrix D^* defined as

$$D^*_{i,j} = (n - 2) \bullet D_{i,j} - TotalDistance_D(i) - TotalDistance_D(j)$$

where $TotalDistance_D(i)$ is the sum of distances from i to all other leaves.

		v_1	v_2	v_3	v_4	$TotalDistance_D$			v_1	v_2	v_3	v_4
D	v_1	0	13	21	22	56	D^*	v_1	0	-68	-60	-60
	v_2	13	0	12	13	38		v_2	-68	0	-60	-60
	v_3	21	12	0	13	46		v_3	-60	-60	0	-68
	v_4	22	13	13	0	48		v_4	-60	-60	-68	0

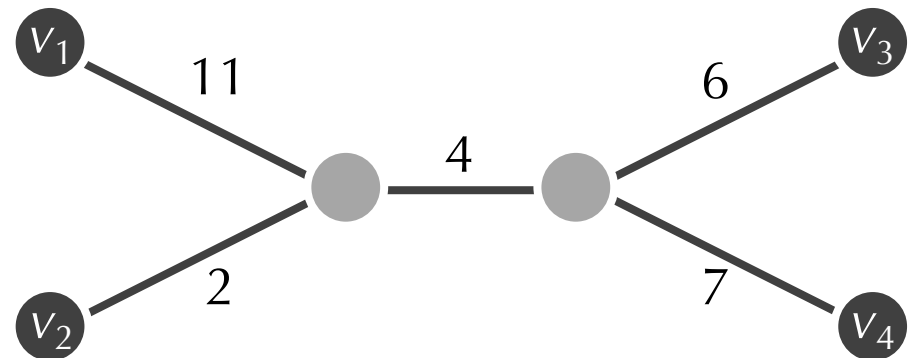
The Neighbor-Joining Theorem

Given an $n \times n$ distance matrix D , its **neighbor-joining matrix** is the matrix D^* defined as

$$D^*_{i,j} = (n - 2) \bullet D_{i,j} - \text{TotalDistance}_D(i) - \text{TotalDistance}_D(j)$$

Note: What does D^* do to outliers?

	v_1	v_2	v_3	v_4
v_1	0	13	21	22
v_2	13	0	12	13
v_3	21	12	0	13
v_4	22	13	13	0



The Neighbor-Joining Theorem

Neighbor-Joining Theorem: If D is additive, then the smallest element of D^* corresponds to neighboring leaves in $Tree(D)$!

		v_1	v_2	v_3	v_4	$TotalDistance_D$			v_1	v_2	v_3	v_4
D	v_1	0	13	21	22	56	D^*	v_1	0	-68	-60	-60
	v_2	13	0	12	13	38		v_2	-68	0	-60	-60
	v_3	21	12	0	13	46		v_3	-60	-60	0	-68
	v_4	22	13	13	0	48		v_4	-60	-60	-68	0

The Neighbor-Joining Theorem

Neighbor-Joining Theorem: If D is additive, then the smallest element of D^* corresponds to neighboring leaves in $Tree(D)$!

	v_1	v_2	v_3	v_4	$TotalDistance_D$		v_1	v_2	v_3	v_4		
D	v_1	0	13	21	22	56	D^*	v_1	0	-68	-60	-60
	v_2	13	0	12	13	38		v_2	-68	0	-60	-60
	v_3	21	12	0	13	46		v_3	-60	-60	0	-68
	v_4	22	13	13	0	48		v_4	-60	-60	-68	0

Neighbor-Joining in Action

	v_1	v_2	v_3	v_4	$TotalDistance_D$	
D^*	v_1	0	-68	-60	-60	56
	v_2	-68	0	-60	-60	38
	v_3	-60	-60	0	-68	46
	v_4	-60	-60	-68	0	48

1. Construct neighbor-joining matrix D^* from D .

Neighbor-Joining in Action

	v_1	v_2	v_3	v_4	$TotalDistance_D$	
D^*	v_1	0	-68	-60	-60	56
	v_2	-68	0	-60	-60	38
	v_3	-60	-60	0	-68	46
	v_4	-60	-60	-68	0	48

2. Find a minimum element $D^*_{i,j}$ of D^* .

Neighbor-Joining in Action

	v_1	v_2	v_3	v_4	$TotalDistance_D$	
D^*	v_1	0	-68	-60	-60	56
	v_2	-68	0	-60	-60	38
	v_3	-60	-60	0	-68	46
	v_4	-60	-60	-68	0	48

2. Find a minimum element $D^*_{i,j}$ of D^* .

Neighbor-Joining in Action

	v_1	v_2	v_3	v_4	$TotalDistance_D$	
D^*	v_1	0	-68	-60	-60	56
	v_2	-68	0	-60	-60	38
	v_3	-60	-60	0	-68	46
	v_4	-60	-60	-68	0	48

$$\Delta_{i,j} = (56 - 38) / (4 - 2) = 9$$

3. Compute $\Delta_{i,j} = (TotalDistance_D(i) - TotalDistance_D(j)) / (n - 2)$.

Neighbor-Joining in Action

	v_1	v_2	v_3	v_4	<i>TotalDistance_D</i>	
<i>D</i>	v_1	0	13	21	22	56
	v_2	13	0	12	13	38
	v_3	21	12	0	13	46
	v_4	22	13	13	0	48

$$\Delta_{i,j} = (56 - 38) / (4 - 2) = 9$$

$$LimbLength(i) = \frac{1}{2}(13 + 9) = 11$$

$$LimbLength(j) = \frac{1}{2}(13 - 9) = 2$$

4. Set $LimbLength(i)$ equal to $\frac{1}{2}(D_{i,j} + \Delta_{i,j})$ and $LimbLength(j)$ equal to $\frac{1}{2}(D_{i,j} - \Delta_{i,j})$.

Neighbor-Joining in Action

	v_1	v_2	v_3	v_4	$TotalDistance_D$	
D	v_1	0	13	21	22	56
	v_2	13	0	12	13	38
	v_3	21	12	0	13	46
	v_4	22	13	13	0	48

$$\Delta_{i,j} = (56 - 38) / (4 - 2) = 9$$

$$LimbLength(i) = \frac{1}{2}(13 + 9) = 11$$

$$LimbLength(j) = \frac{1}{2}(13 - 9) = 2$$

4. Set $LimbLength(i)$ equal to $\frac{1}{2}(D_{i,j} + \Delta_{i,j})$ and $LimbLength(j)$ equal to $\frac{1}{2}(D_{i,j} - \Delta_{i,j})$.

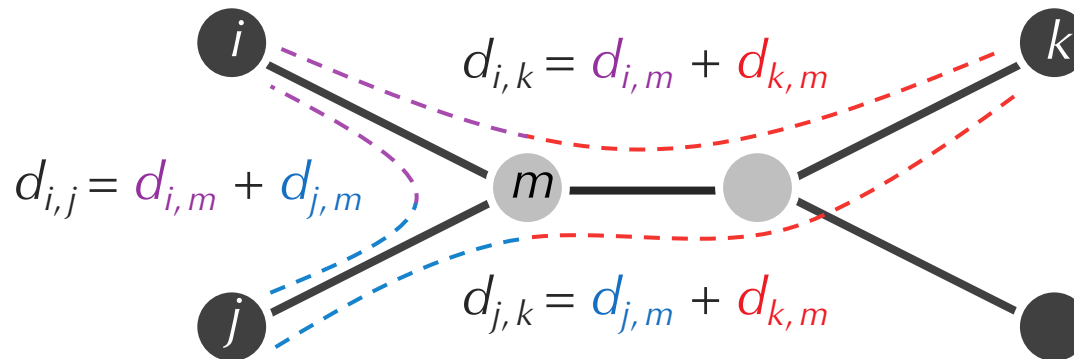
Wait ... where do these formulas come from?

Neighbor-Joining in Action

	m	v_3	v_4	$TotalDistance_D$	
D'	m	0	10	11	21
	v_3	10	0	13	23
	v_4	11	13	0	24

5. Form a matrix D' by removing i -th and j -th row/column from D and adding an m -th row/column such that for any k , $D_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$.

Flashback: Computation of $d_{k,m}$

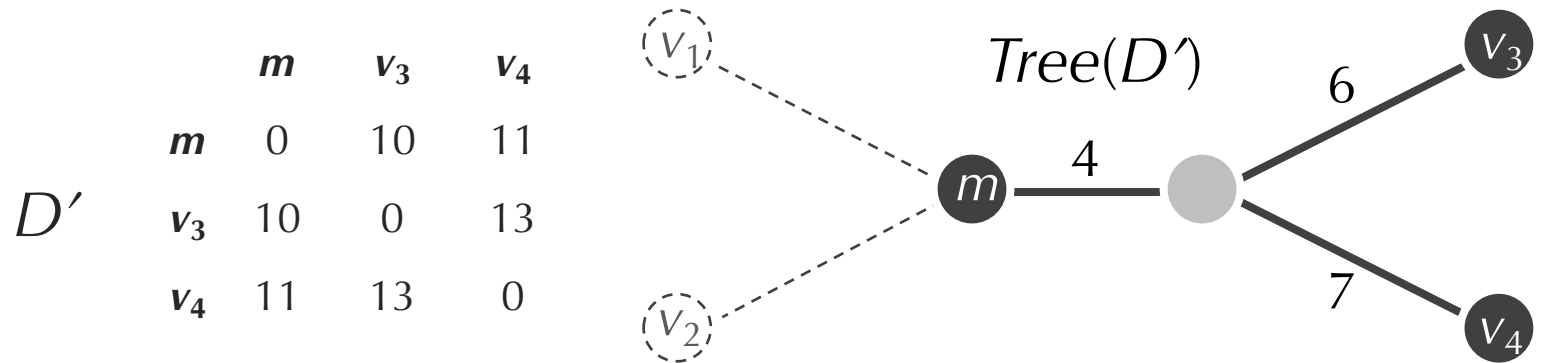


$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

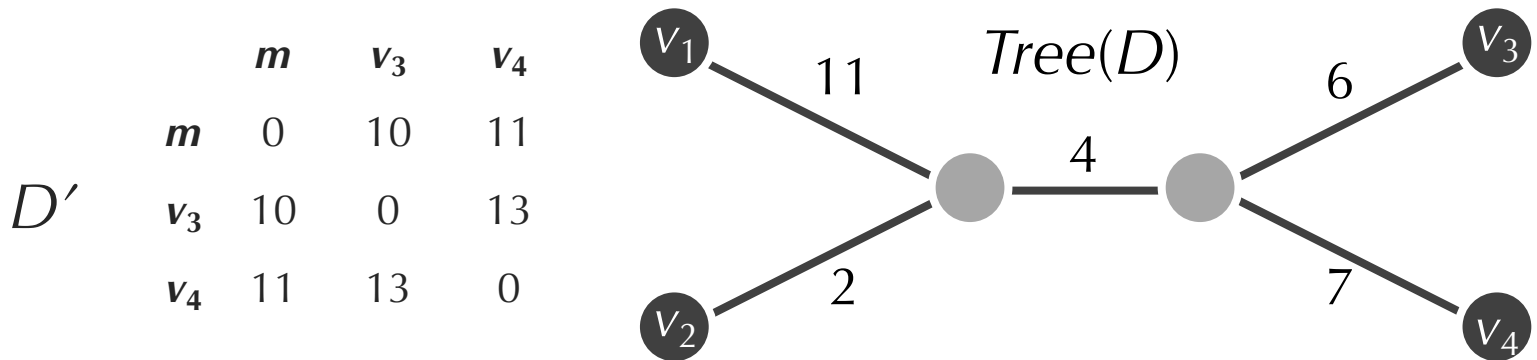
$$d_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

Neighbor-Joining in Action



6. Apply **NeighborJoining** recursively to D' to obtain $Tree(D')$.

Neighbor-Joining in Action



$$\text{LimbLength}(v_1) = \frac{1}{2}(13 + 9) = 11$$

$$\text{LimbLength}(v_2) = \frac{1}{2}(13 - 9) = 2$$

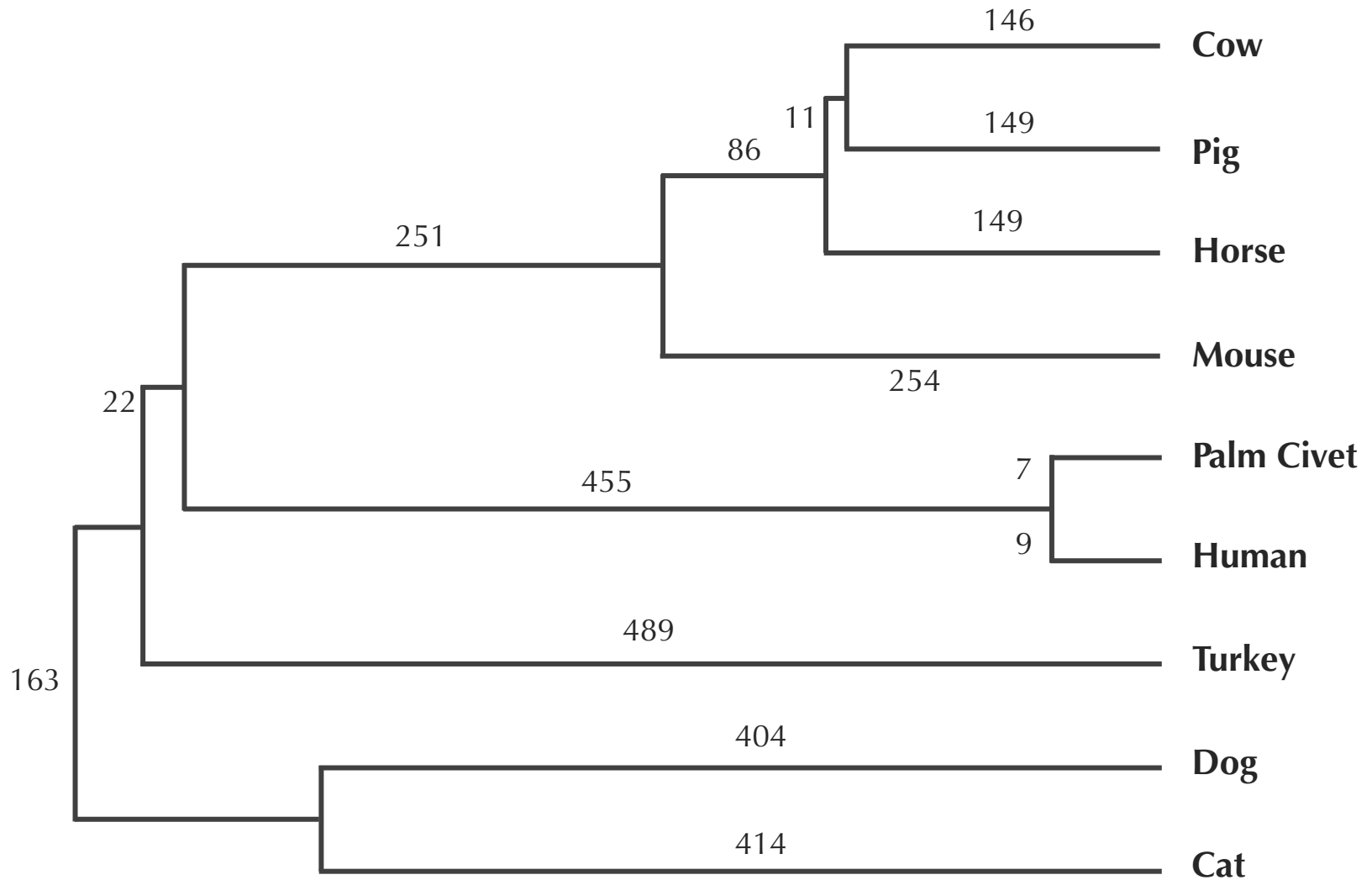
7. Reattach limbs of *i* and *j* to obtain *Tree(D)*.

Neighbor-Joining in Action

NeighborJoining(D):

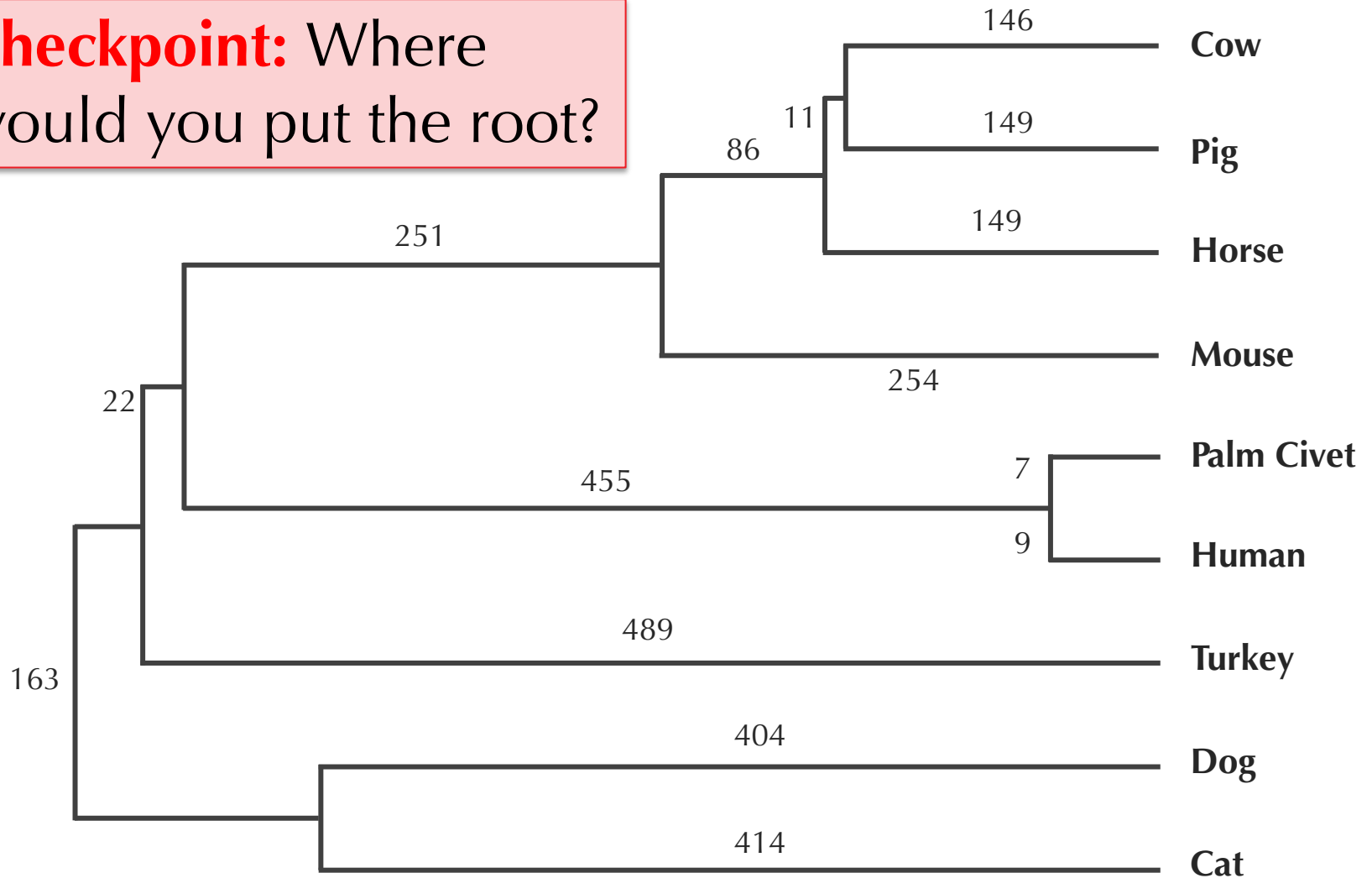
1. Construct neighbor-joining matrix D^* from D .
2. Find a minimum element $D^*_{i,j}$ of D^* .
3. Compute $\Delta_{i,j} = (TotalDistance_D(i) - TotalDistance_D(j)) / (n - 2)$.
4. Set $LimbLength(i)$ equal to $\frac{1}{2}(D_{i,j} + \Delta_{i,j})$ and $LimbLength(j)$ equal to $\frac{1}{2}(D_{i,j} - \Delta_{i,j})$.
5. Form a matrix D' by removing i -th and j -th row/column from D and adding an m -th row/column such that for any k , $D_{k,m} = (D_{k,i} + D_{k,j} - D_{i,j}) / 2$.
6. Apply **NeighborJoining** recursively to D' to obtain $Tree(D')$.
7. Reattach limbs of i and j to obtain $Tree(D)$.

Neighbor-Joining on Coronavirus “Spike Proteins”



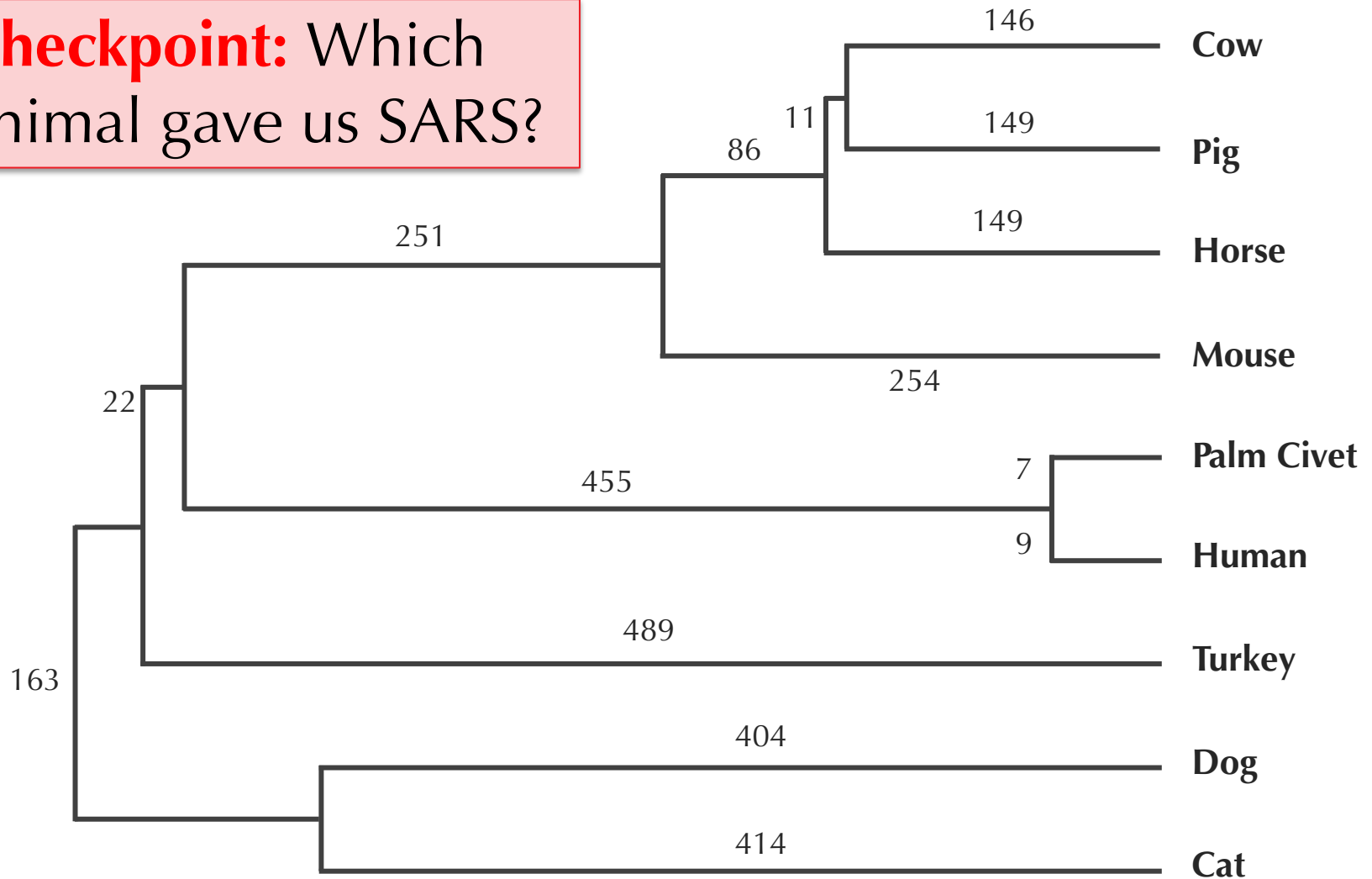
Neighbor-Joining on Coronavirus “Spike Proteins”

Checkpoint: Where would you put the root?



Neighbor-Joining on Coronavirus “Spike Proteins”

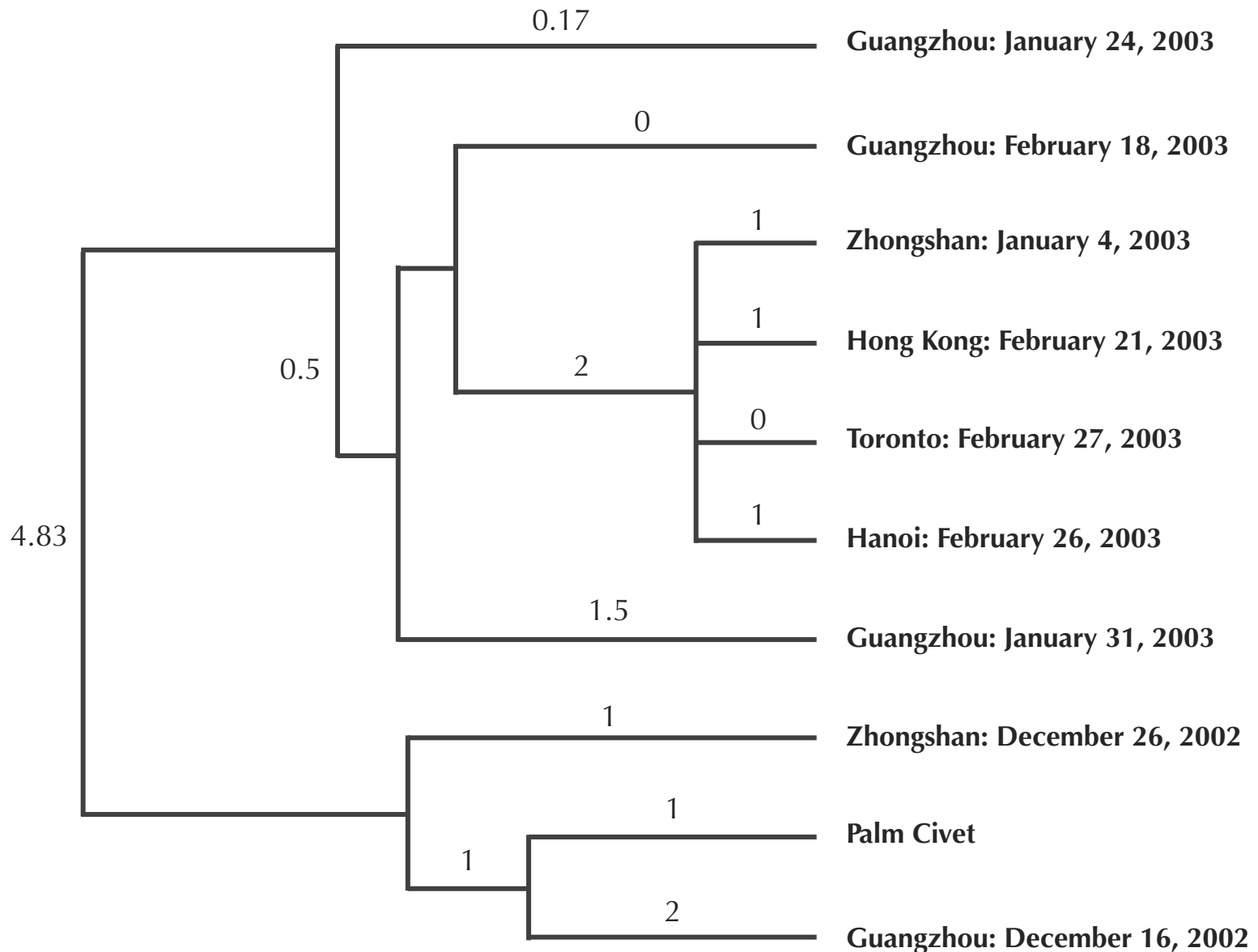
Checkpoint: Which animal gave us SARS?







Neighbor-Joining on Coronaviruses



Weakness of Distance-Based Methods

Distance-based algorithms for evolutionary tree reconstruction say nothing about ancestral states at internal nodes.

Weakness of Distance-Based Methods

Distance-based algorithms for evolutionary tree reconstruction say nothing about ancestral states at internal nodes.

We *lost* information when we converted a multiple alignment to a distance matrix...

SPECIES	ALIGNMENT	DISTANCE MATRIX			
		Chimp	Human	Seal	Whale
Chimp	ACGTAGGCCT	0	3	6	4
Human	ATGTAAGACT	3	0	7	5
Seal	TCGAGAGCAC	6	7	0	2
Whale	TCGAAAGCAT	4	5	2	0

Which Animal Gave Us SARS?

Evolutionary Trees Part 2:

The Small Parsimony Algorithm for Inferring Ancestral States

Phillip Compeau and Pavel Pevzner

Bioinformatics Algorithms: An Active Learning Approach

©2018 by Compeau and Pevzner. All rights reserved.

What is the main evolutionary division of dinosaurs?



Character-Based Construction

Biologists used to form trees based on anatomical or physiological properties called **characters**.

Character-Based Construction

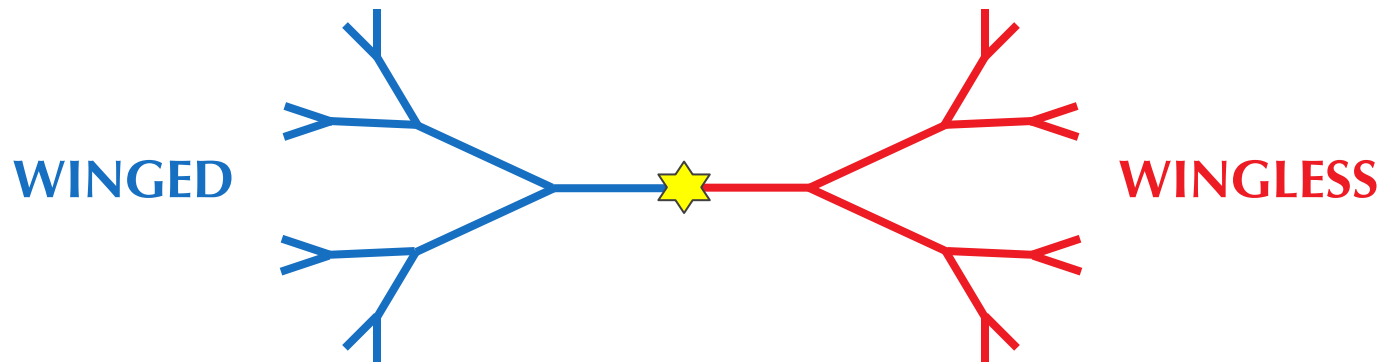
Biologists used to form trees based on anatomical or physiological properties called **characters**.

Checkpoint: Say you wanted to construct an evolutionary tree of all insects. What might be the first thing you would do?

Character-Based Construction

Biologists used to form trees based on anatomical or physiological properties called **characters**.

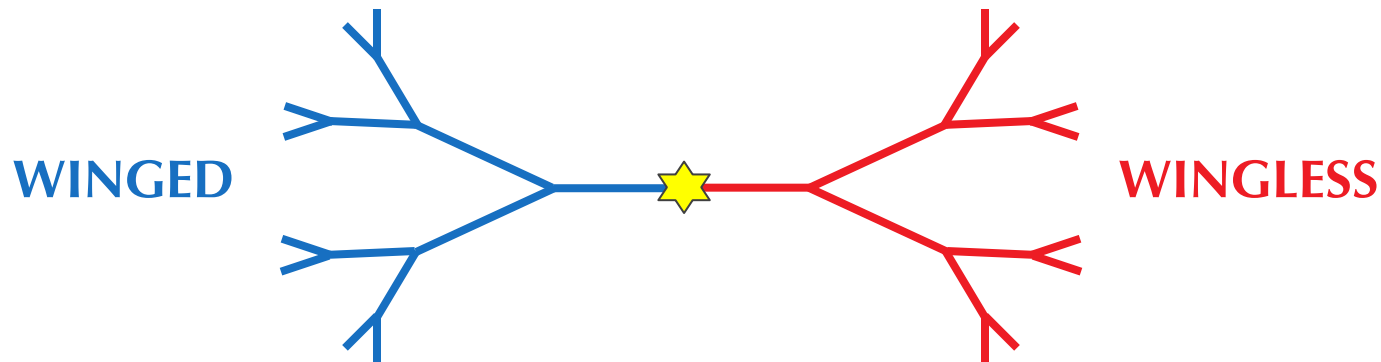
Checkpoint: Say you wanted to construct an evolutionary tree of all insects. What might be the first thing you would do?



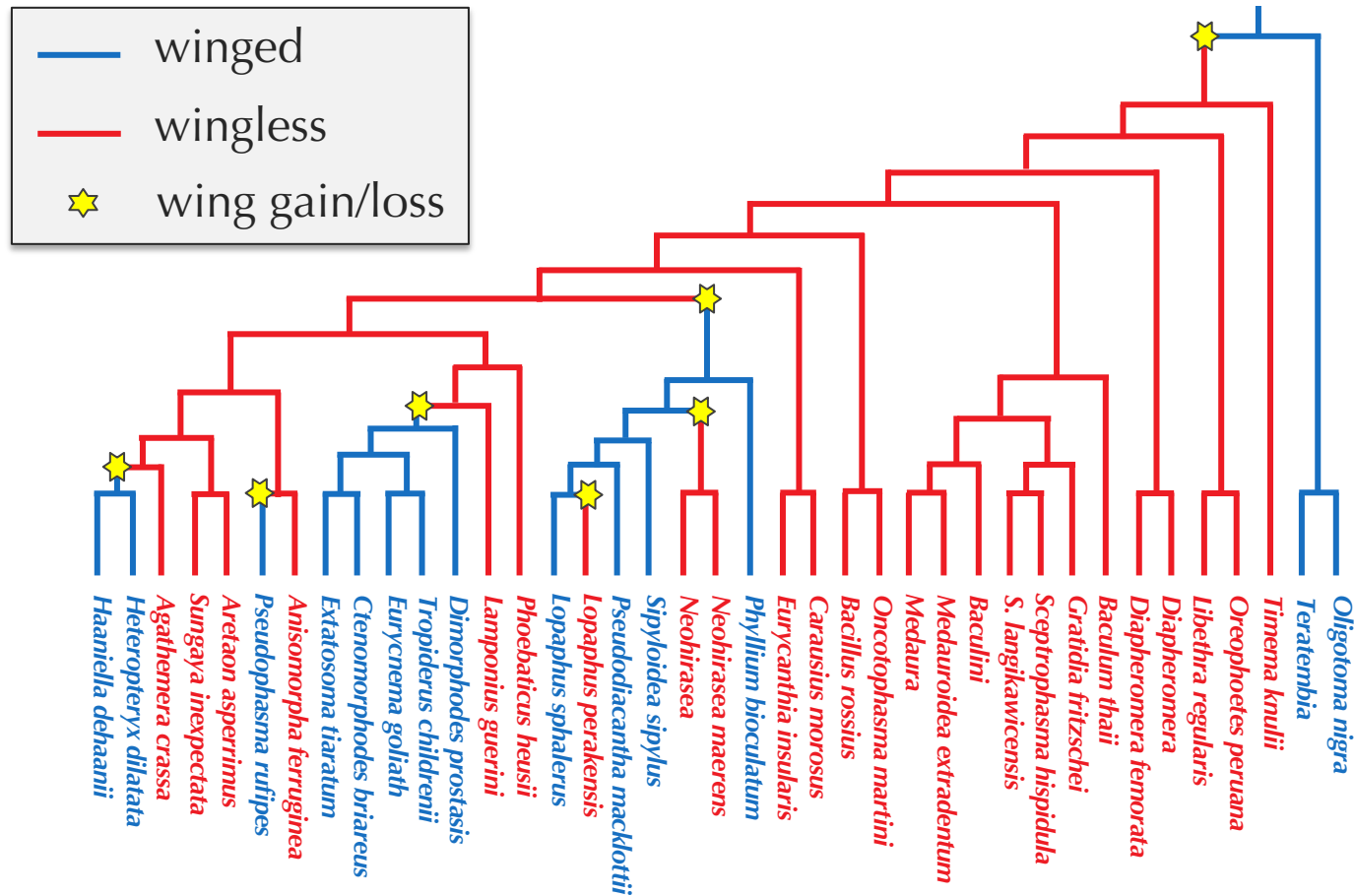
Character-Based Construction

Dollo's principle of irreversibility (1893): evolution doesn't reinvent the same organ (e.g. insect wings).

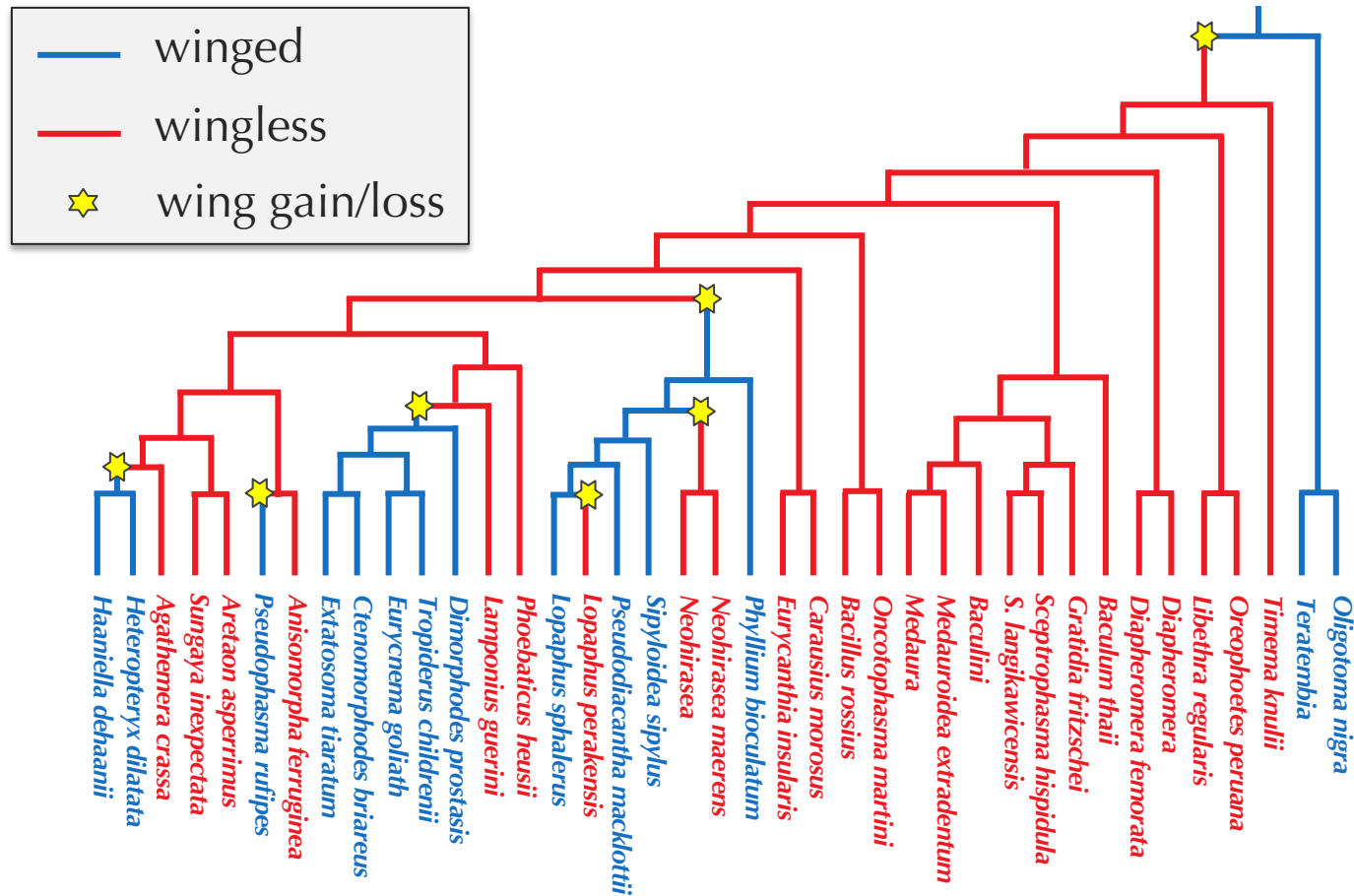
Checkpoint: Say you wanted to construct an evolutionary tree of all insects. What might be the first thing you would do?



Dollo's Principle Violated in Stick Insect Phylogeny

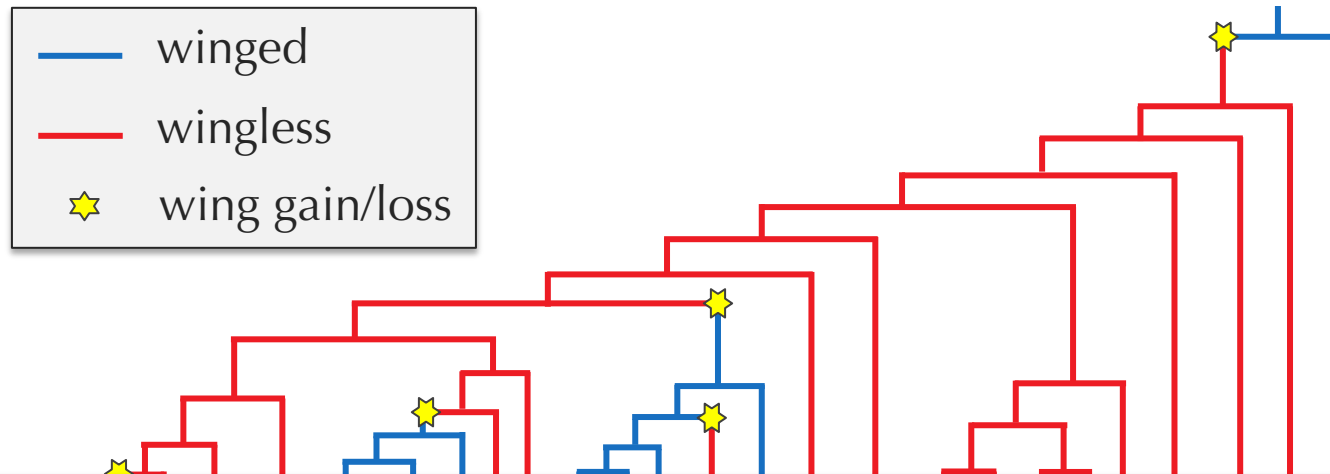


Dollo's Principle Violated in Stick Insect Phylogeny



Checkpoint: What do you think happened?

Dollo's Principle Violated in Stick Insect Phylogeny




Key Point: This approach won't help reconstruct the tree *structure*, but perhaps we can infer *ancestral* characters after constructing a tree using another method (e.g., neighbor-joining).

Checkpoint: What do you think happened?

An Alignment Is a Character Table

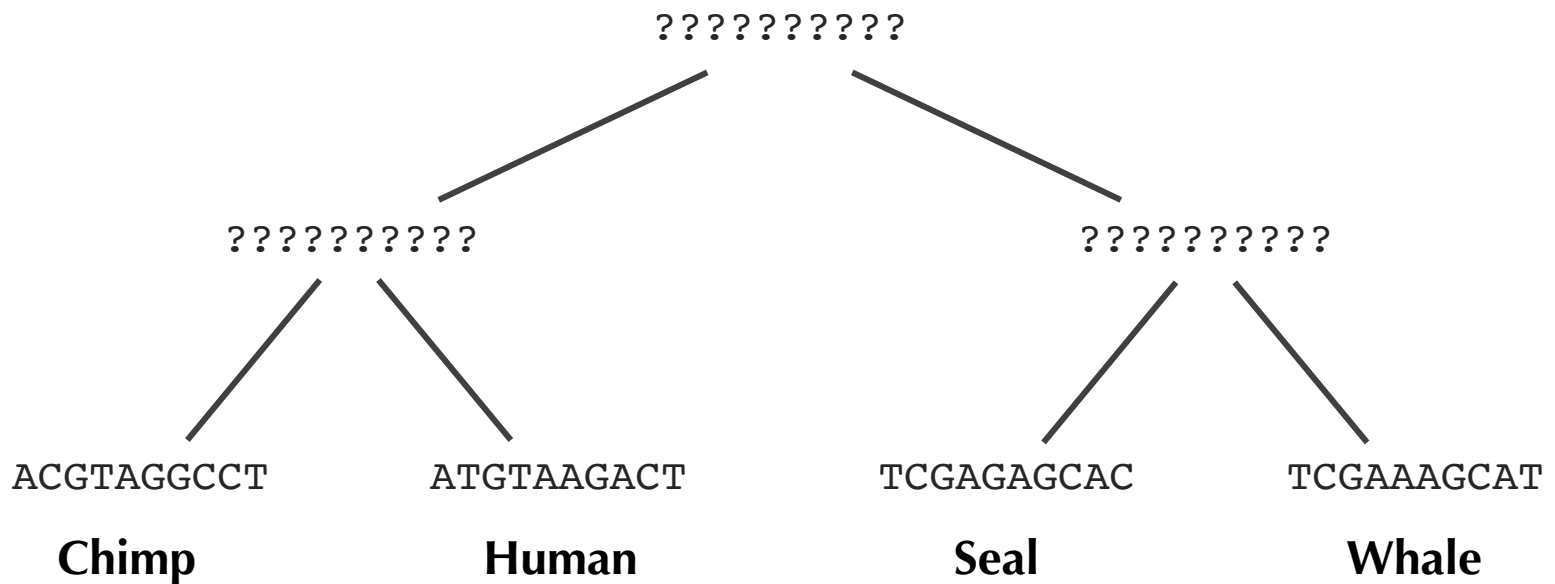
Chimp	ACGTAGGCCT
Human	ATGTAAGACT
Seal	TCGAGAGCAC
Whale	TCGAAAGCAT

An Alignment Is a Character Table

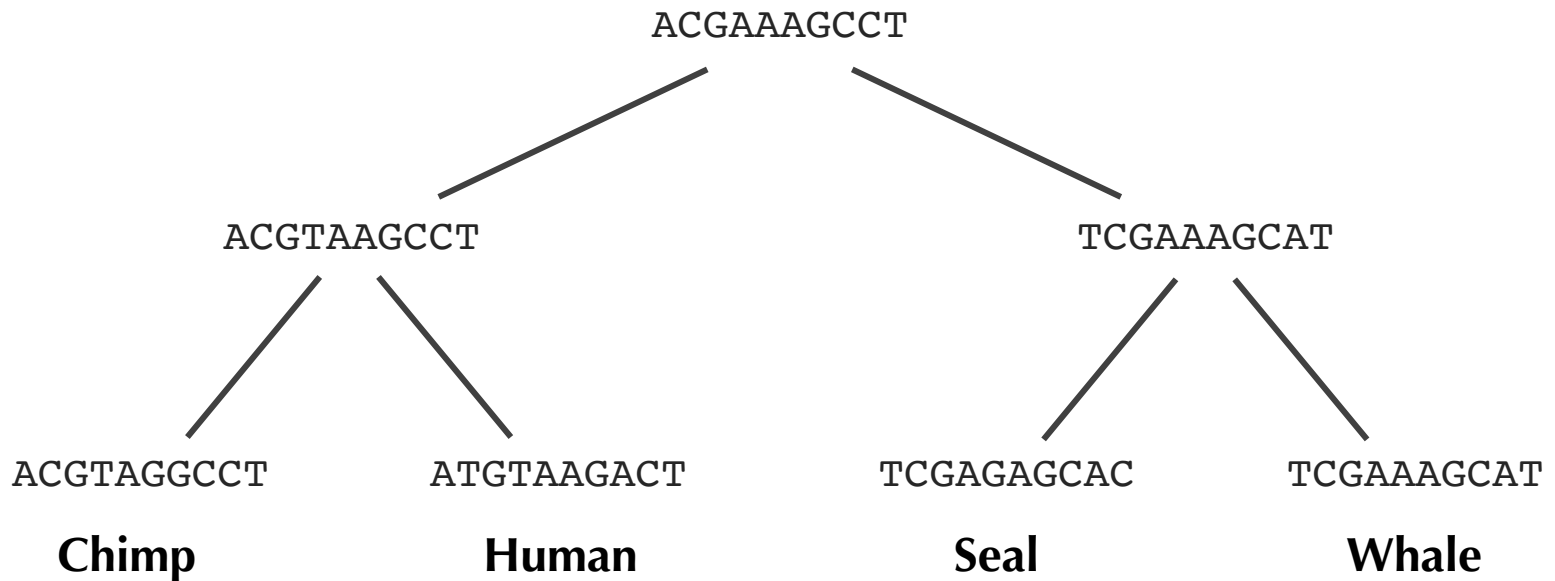
Chimp	ACGTAGGCCT	} n species
Human	ATGTAAGACT	
Seal	TCGAGAGCAC	
Whale	TCGAAAGCAT	
		
	m characters	

Toward a Computational Problem

Chimp	ACGTAGGCCT
Human	ATGTAAGACT
Seal	TCGAGAGCAC
Whale	TCGAAAGCAT

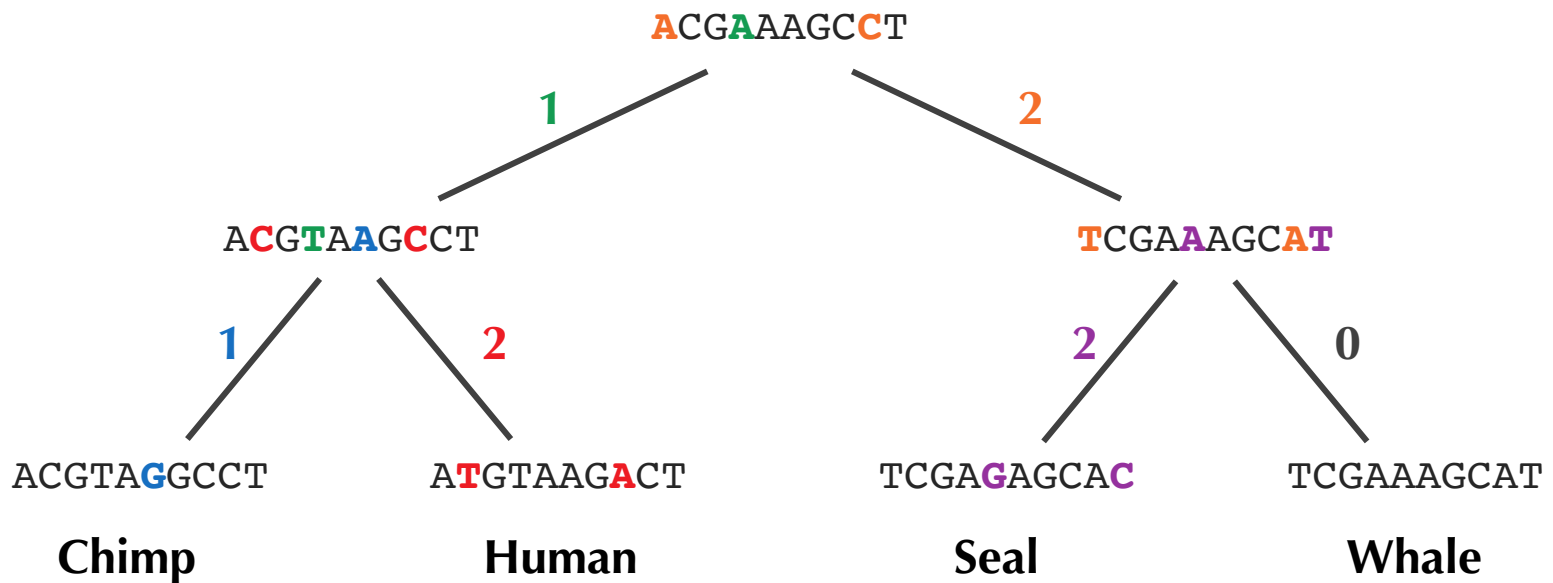


Toward a Computational Problem



Toward a Computational Problem

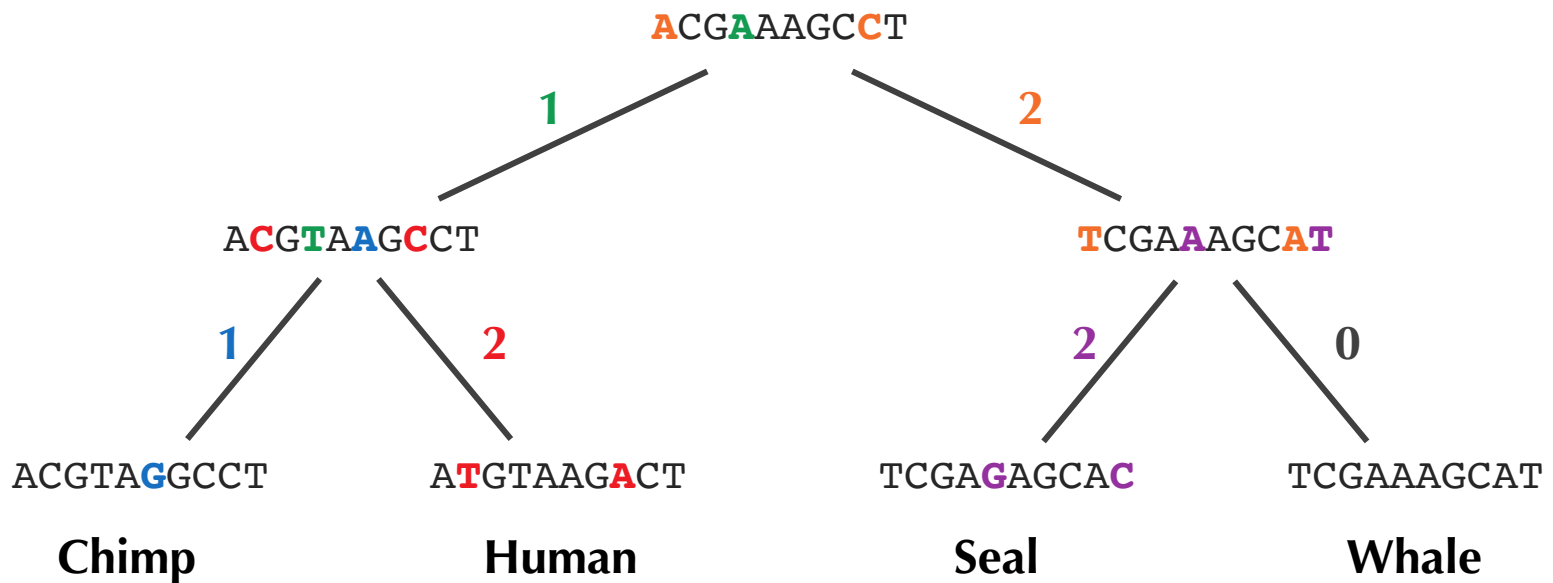
Parsimony score: sum of Hamming distances (total mismatches) along each edge.



Toward a Computational Problem

Parsimony score: sum of Hamming distances (total mismatches) along each edge.

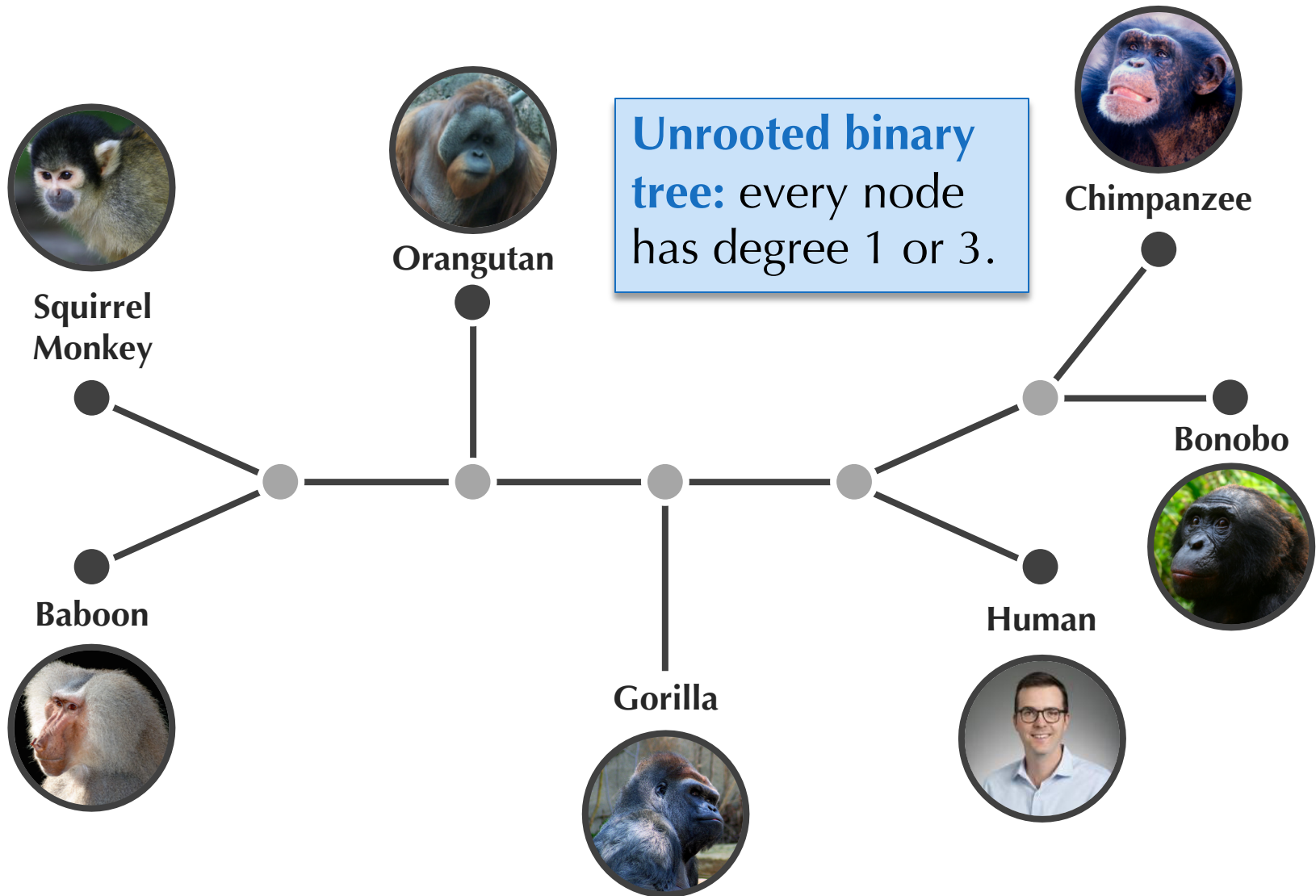
Parsimony Score: 8



Modeling Speciations

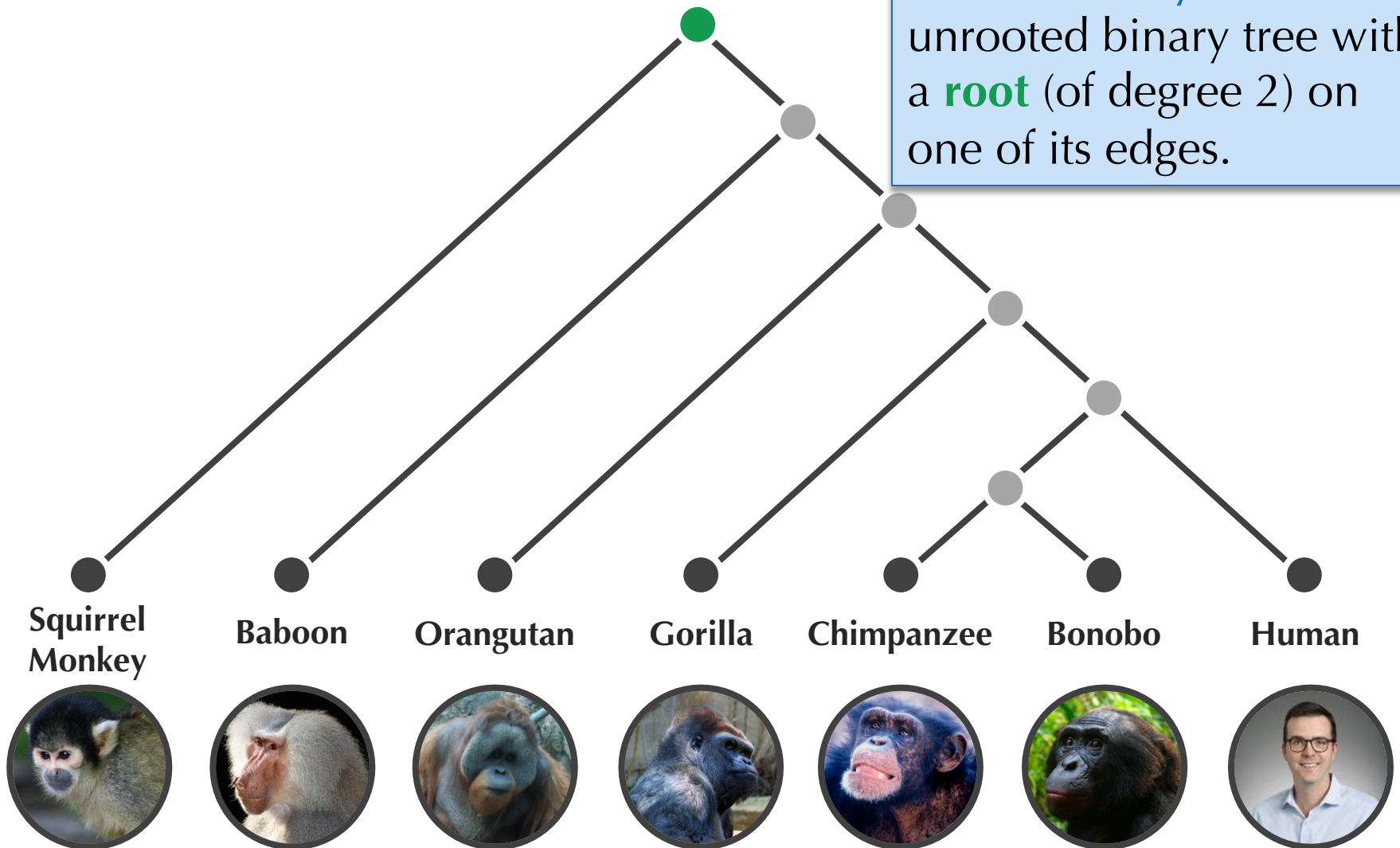
Researchers often assume that all internal nodes correspond to **speciations**, where one species splits into two.

Modeling Speciations



Modeling Speciations

Rooted binary tree: an unrooted binary tree with a **root** (of degree 2) on one of its edges.



Toward a Computational Problem

Small Parsimony Problem: *Find the most parsimonious labeling of the internal nodes of a rooted tree.*

- **Input:** A rooted binary tree with each leaf labeled by a string of length m .
- **Output:** A labeling of all other nodes of the tree by strings of length m that minimizes the tree's parsimony score.

Toward a Computational Problem

Small Parsimony Problem: *Find the most parsimonious labeling of the internal nodes of a rooted tree.*

- **Input:** A rooted binary tree with each leaf labeled by a string of length m .
- **Output:** A labeling of all other nodes of the tree by strings of length m that minimizes the tree's parsimony score.

Checkpoint: Is there any way we can simplify this problem statement?

Toward a Computational Problem

Small Parsimony Problem: *Find the most parsimonious labeling of the internal nodes of a rooted tree.*

- **Input:** A rooted binary tree with each leaf labeled by a **single symbol**.
- **Output:** A labeling of all other nodes of the tree by **single symbols** that minimizes the tree's parsimony score.

Toward a Computational Problem

Small Parsimony Problem: *Find the most parsimonious labeling of the internal nodes of a rooted tree.*

- **Input:** A rooted binary tree with each leaf labeled by a **single symbol**.
- **Output:** A labeling of all other nodes of the tree by **single symbols** that minimizes the tree's parsimony score.

Checkpoint: Why is this an acceptable simplification?

Toward a Computational Problem

Small Parsimony Problem: *Find the most parsimonious labeling of the internal nodes of a rooted tree.*

- **Input:** A rooted binary tree with each leaf labeled by a **single symbol**.
- **Output:** A labeling of all other nodes of the tree by **single symbols** that minimizes the tree's parsimony score.

Answer: We may choose to assume that the characters are *independent*.

Toward a Computational Problem

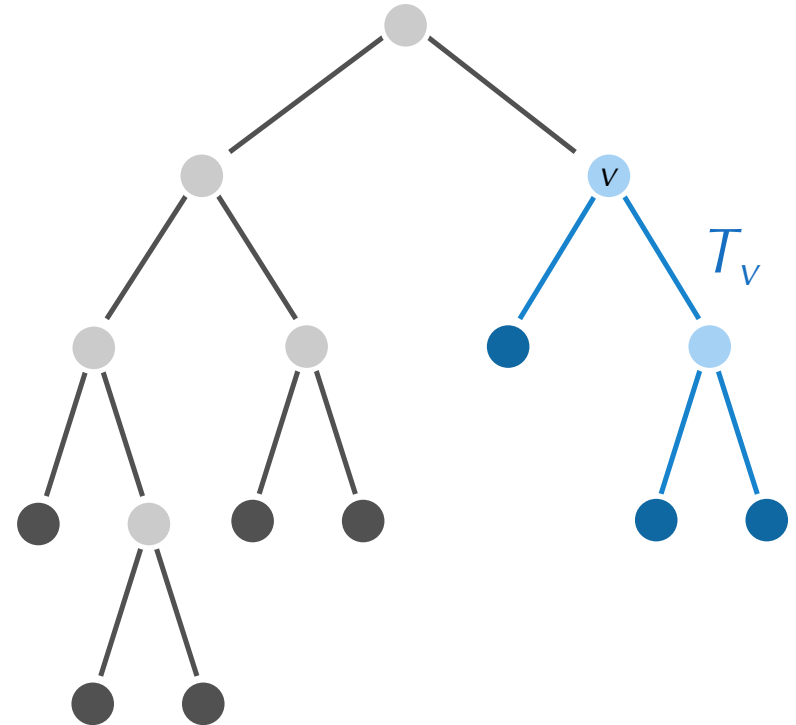
Small Parsimony Problem: *Find the most parsimonious labeling of the internal nodes of a rooted tree.*

- **Input:** A rooted binary tree with each leaf labeled by a **single symbol**.
- **Output:** A labeling of all other nodes of the tree by **single symbols** that minimizes the tree's parsimony score.

Checkpoint: Any thoughts on what approach we might use to solve this problem?

A Dynamic Programming Algorithm

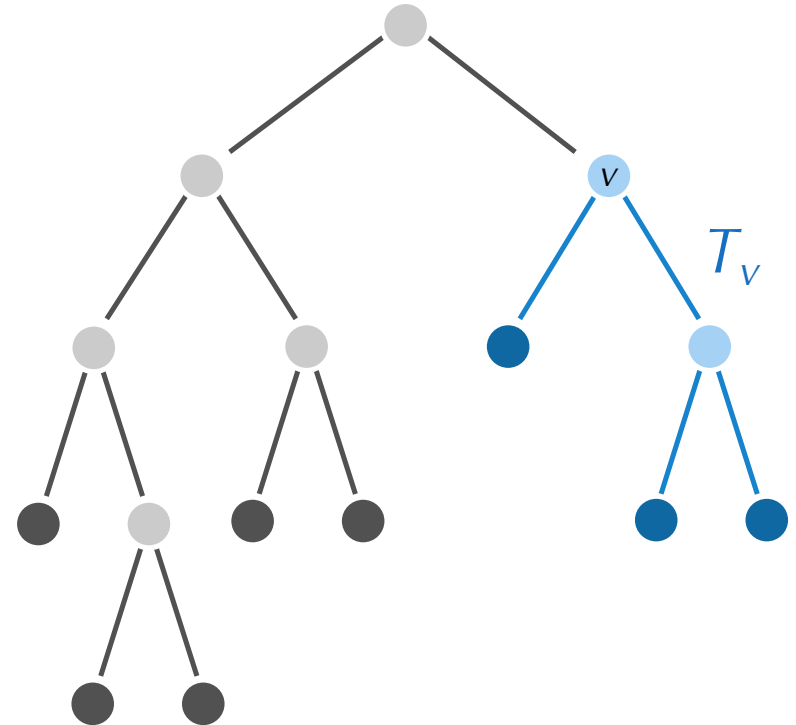
Let T_v denote the subtree of T whose root is v .



A Dynamic Programming Algorithm

Let T_v denote the subtree of T whose root is v .

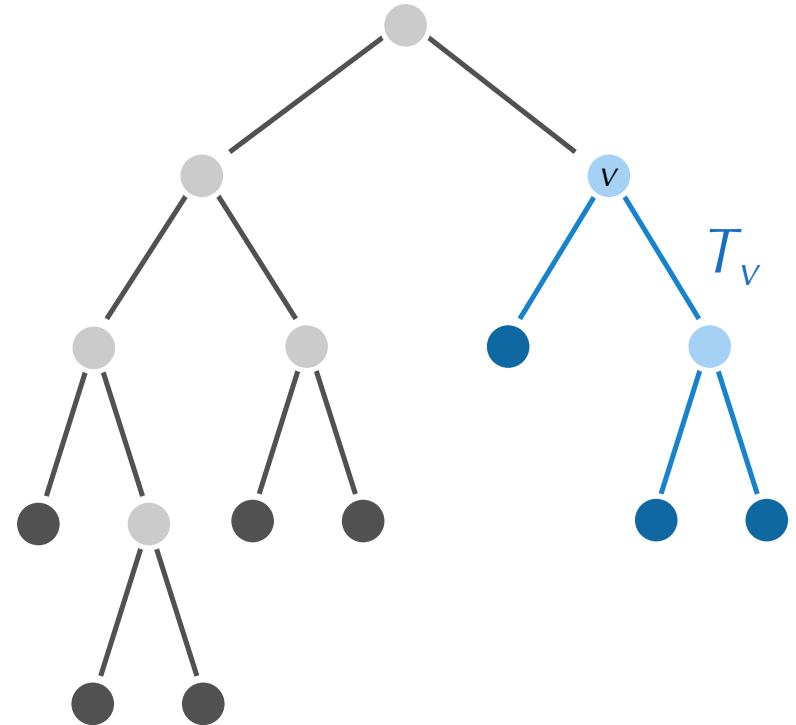
Define $s_k(v)$ as the minimum parsimony score of T_v over all labelings of T_v , assuming that v is labeled by k .



A Dynamic Programming Algorithm

Let T_v denote the subtree of T whose root is v .

Define $s_k(v)$ as the minimum parsimony score of T_v over all labelings of T_v , assuming that v is labeled by k .

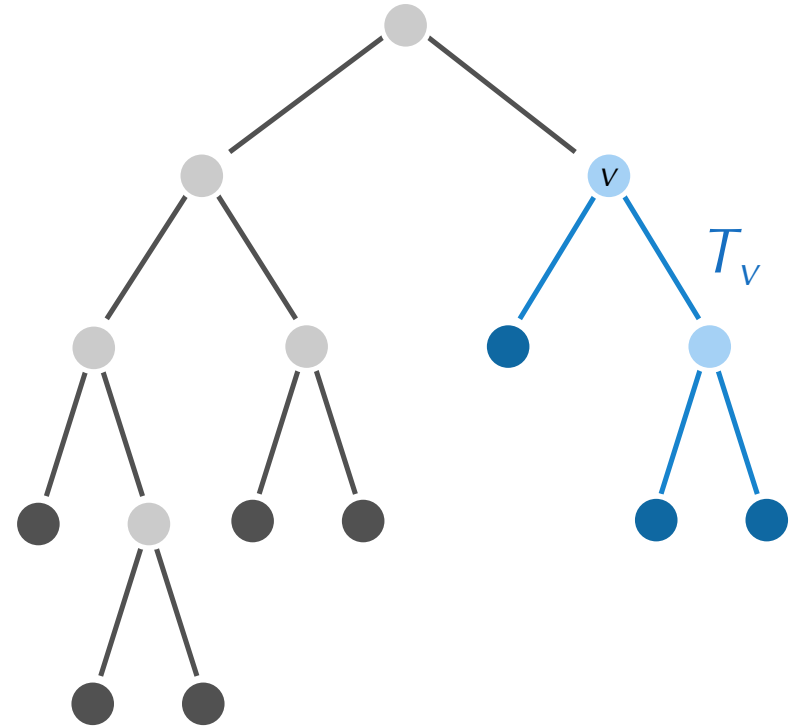


The minimum parsimony score for the tree is equal to the minimum value of $s_k(\text{root})$ over all symbols k .

A Dynamic Programming Algorithm

For symbols i and j , define

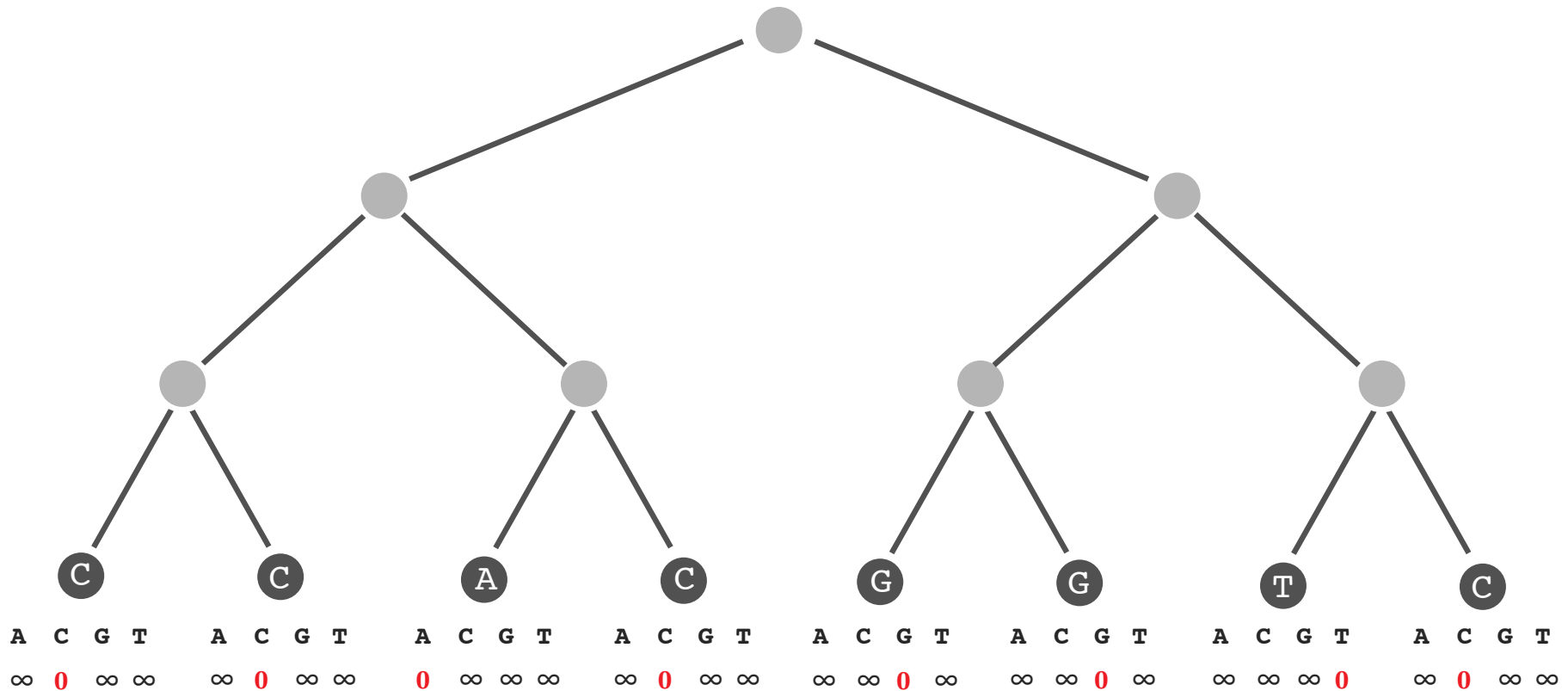
- $\alpha_{i,j} = 0$ if $i = j$
- $\alpha_{i,j} = 1$ otherwise.



Theorem: The following recurrence relation holds:

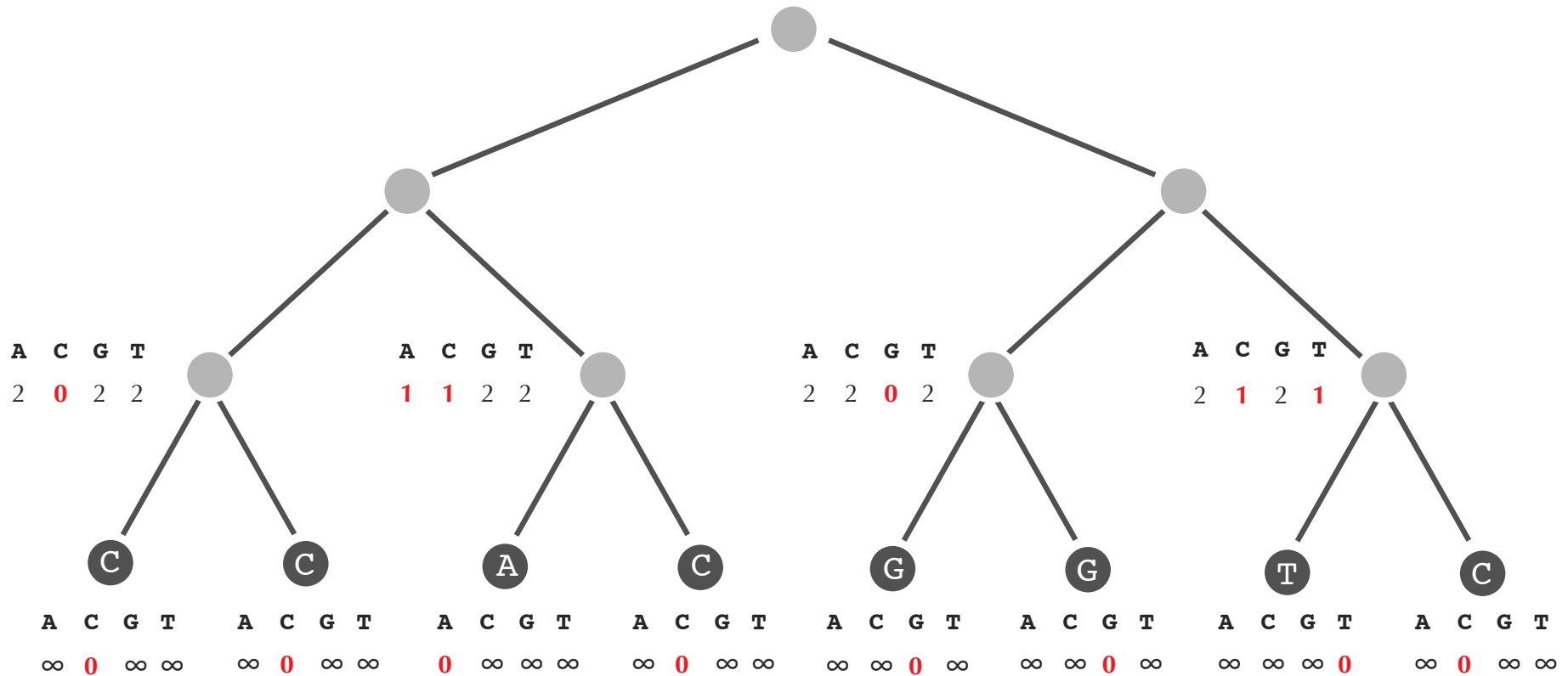
$$s_k(v) = \min_{\text{all symbols } i} \{s_i(\text{Daughter}(v)) + \delta_{i,k}\} + \min_{\text{all symbols } i} \{s_i(\text{Son}(v)) + \delta_{j,k}\}$$

A Dynamic Programming Algorithm



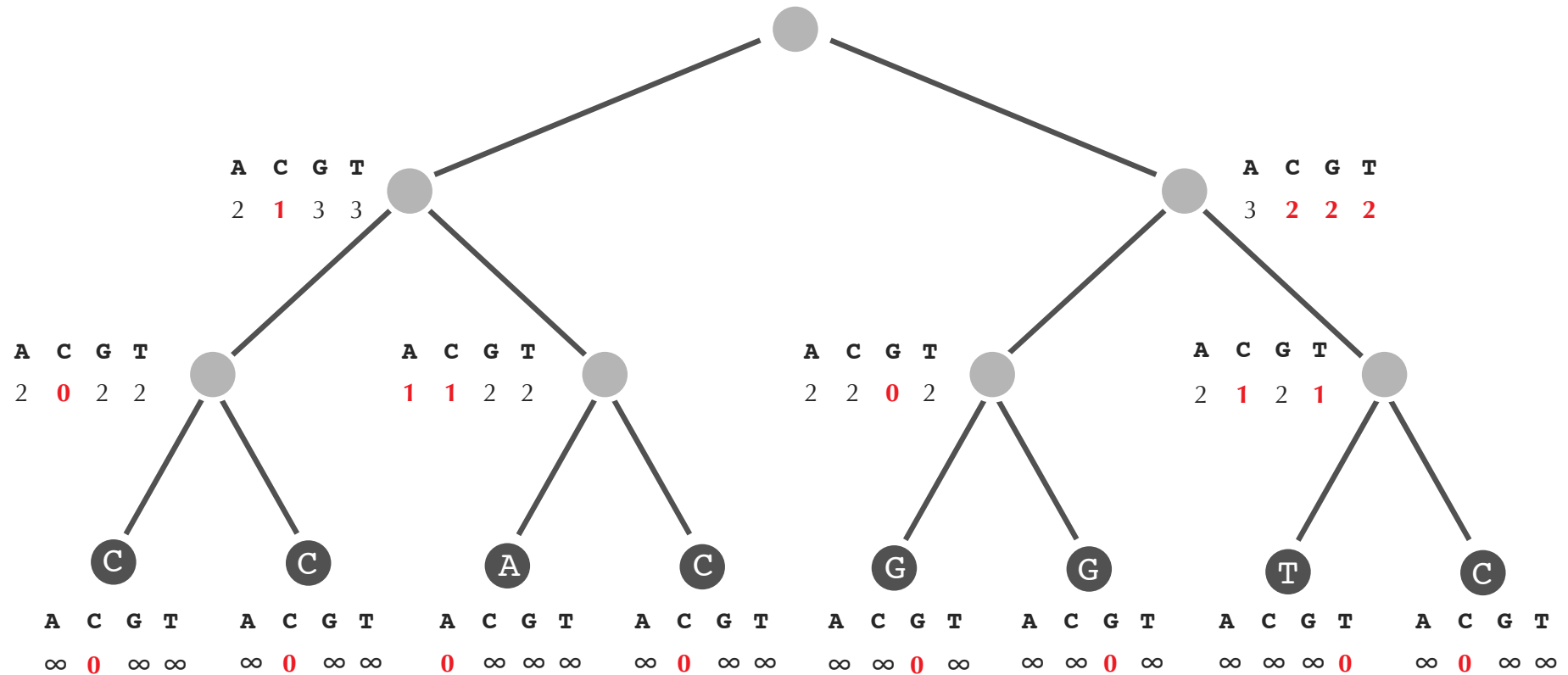
$$s_k(v) = \min_{\text{all symbols } i} \{s_i(\text{Daughter}(v)) + \delta_{i,k}\} + \min_{\text{all symbols } i} \{s_i(\text{Son}(v)) + \delta_{i,k}\}$$

A Dynamic Programming Algorithm



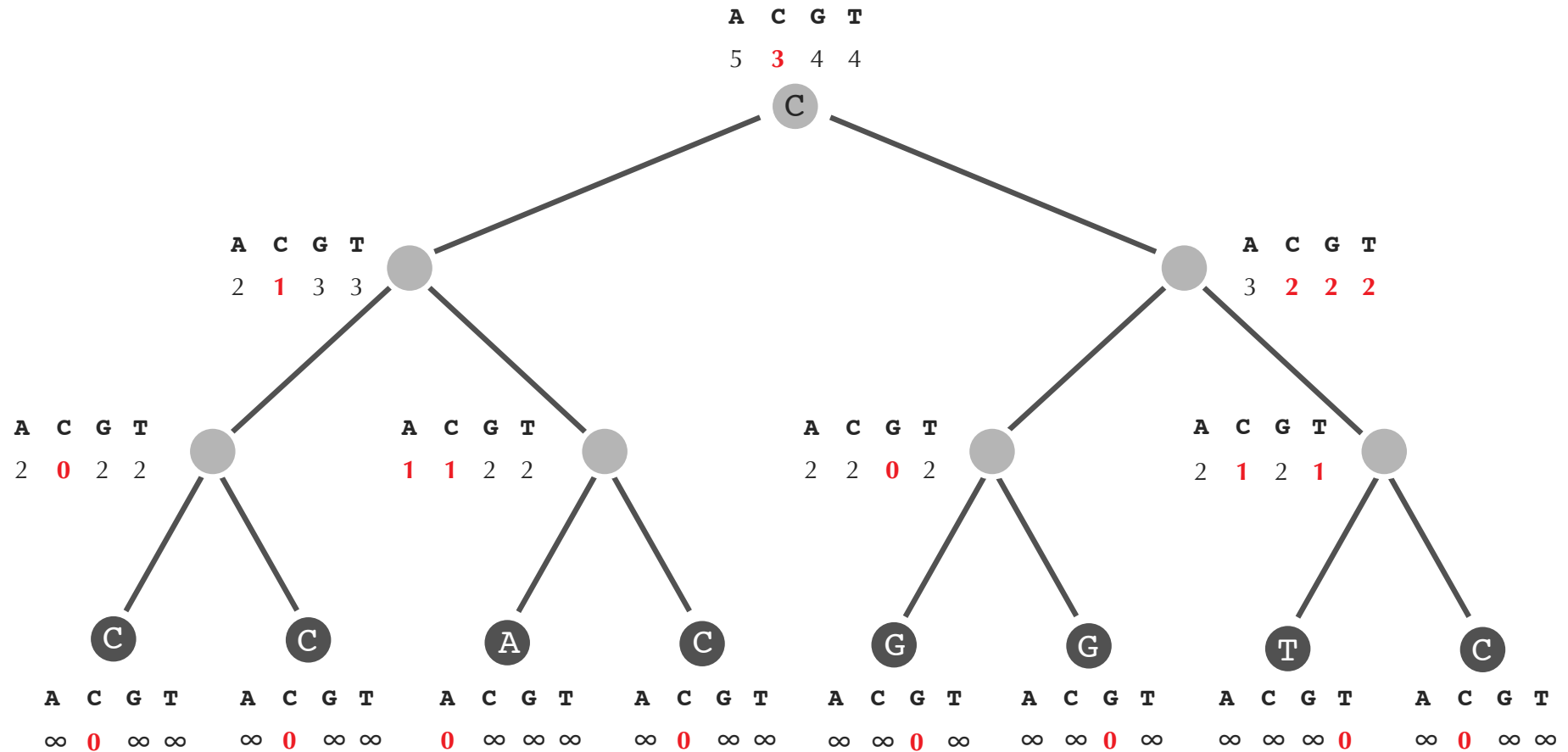
$$s_k(v) = \min_{\text{all symbols } i} \{s_i(\text{Daughter}(v)) + \delta_{i,k}\} + \min_{\text{all symbols } i} \{s_i(\text{Son}(v)) + \delta_{i,k}\}$$

A Dynamic Programming Algorithm



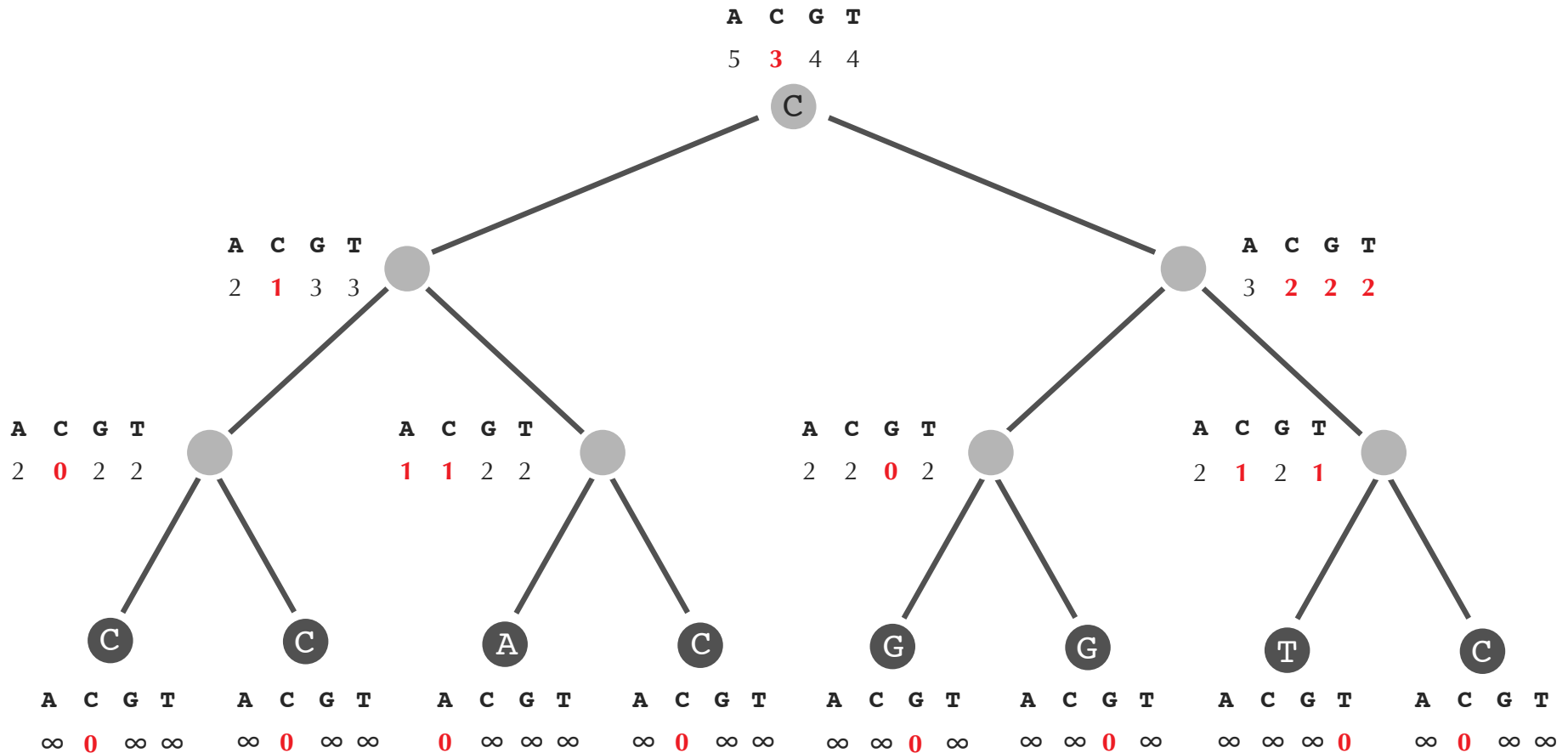
$$s_k(v) = \min_{\text{all symbols } i} \{s_i(\text{Daughter}(v)) + \delta_{i,k}\} + \min_{\text{all symbols } i} \{s_i(\text{Son}(v)) + \delta_{i,k}\}$$

A Dynamic Programming Algorithm



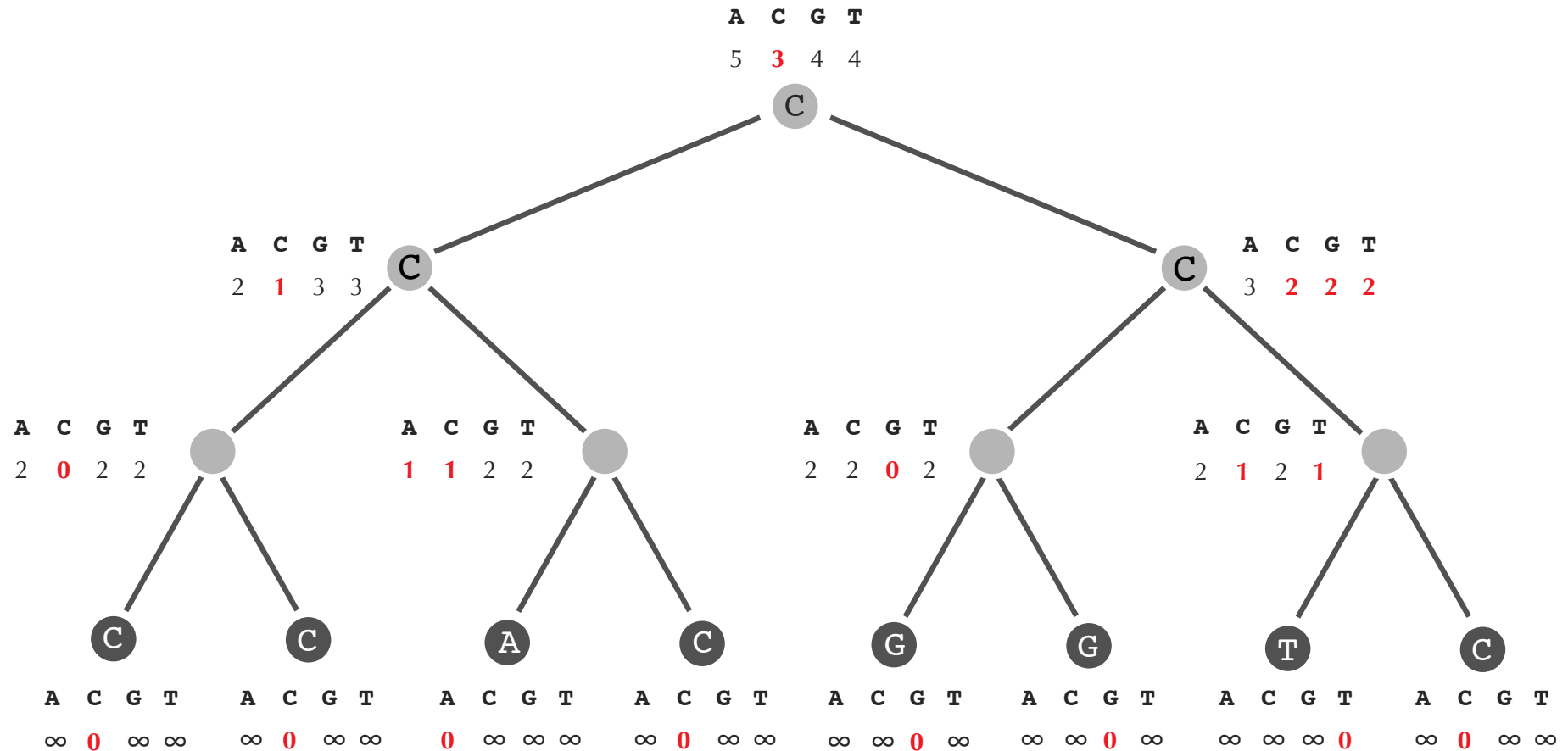
$$s_k(v) = \min_{\text{all symbols } i} \{s_i(\text{Daughter}(v)) + \delta_{i,k}\} + \min_{\text{all symbols } i} \{s_i(\text{Son}(v)) + \delta_{i,k}\}$$

A Dynamic Programming Algorithm



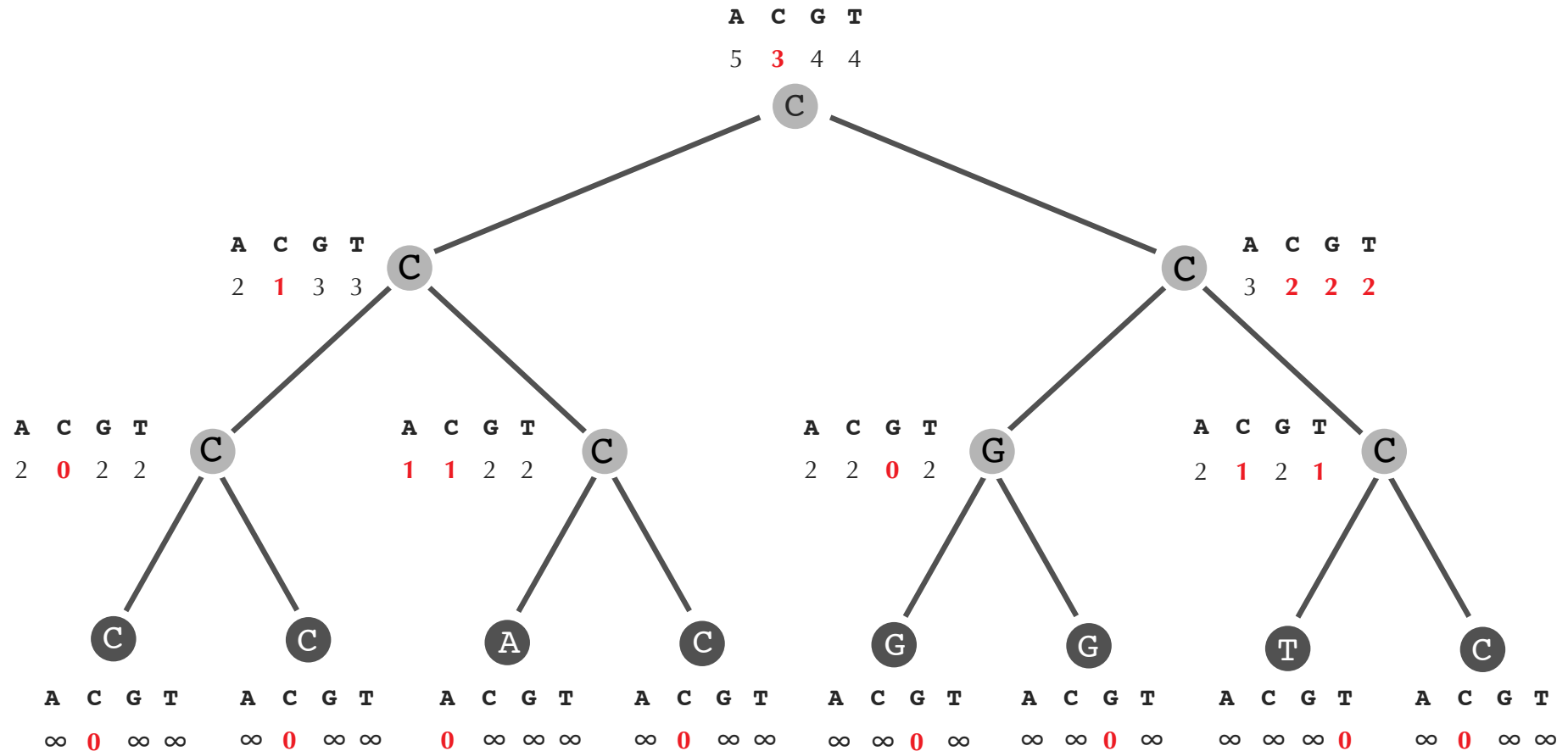
Checkpoint: How should we “backtrack” to fill in the remaining nodes of the tree?

A Dynamic Programming Algorithm



Checkpoint: How should we “backtrack” to fill in the remaining nodes of the tree?

A Dynamic Programming Algorithm



Checkpoint: How should we “backtrack” to fill in the remaining nodes of the tree?