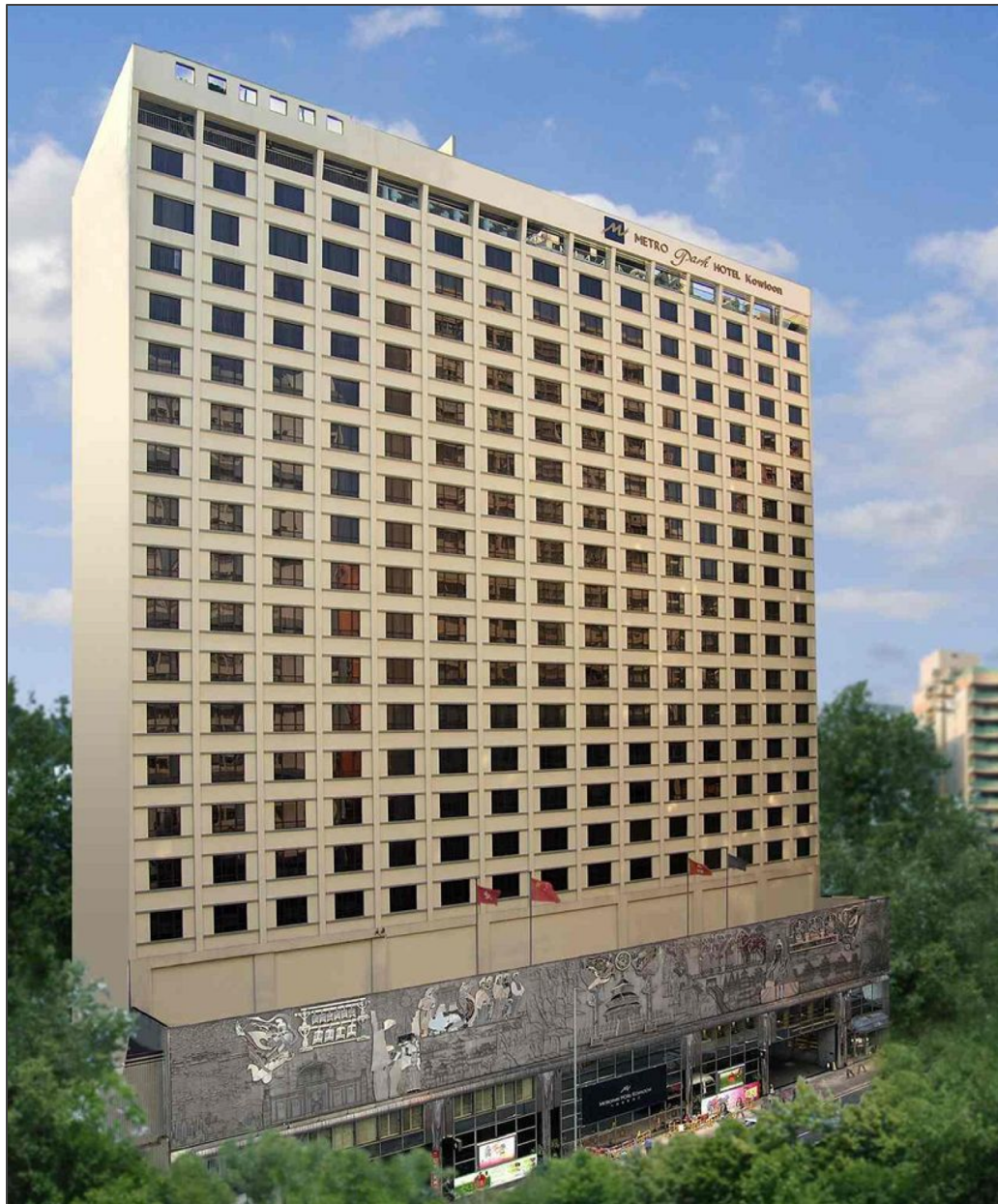# Which Animal Gave Us SARS?

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

Phillip Compeau and Pavel Pevzner
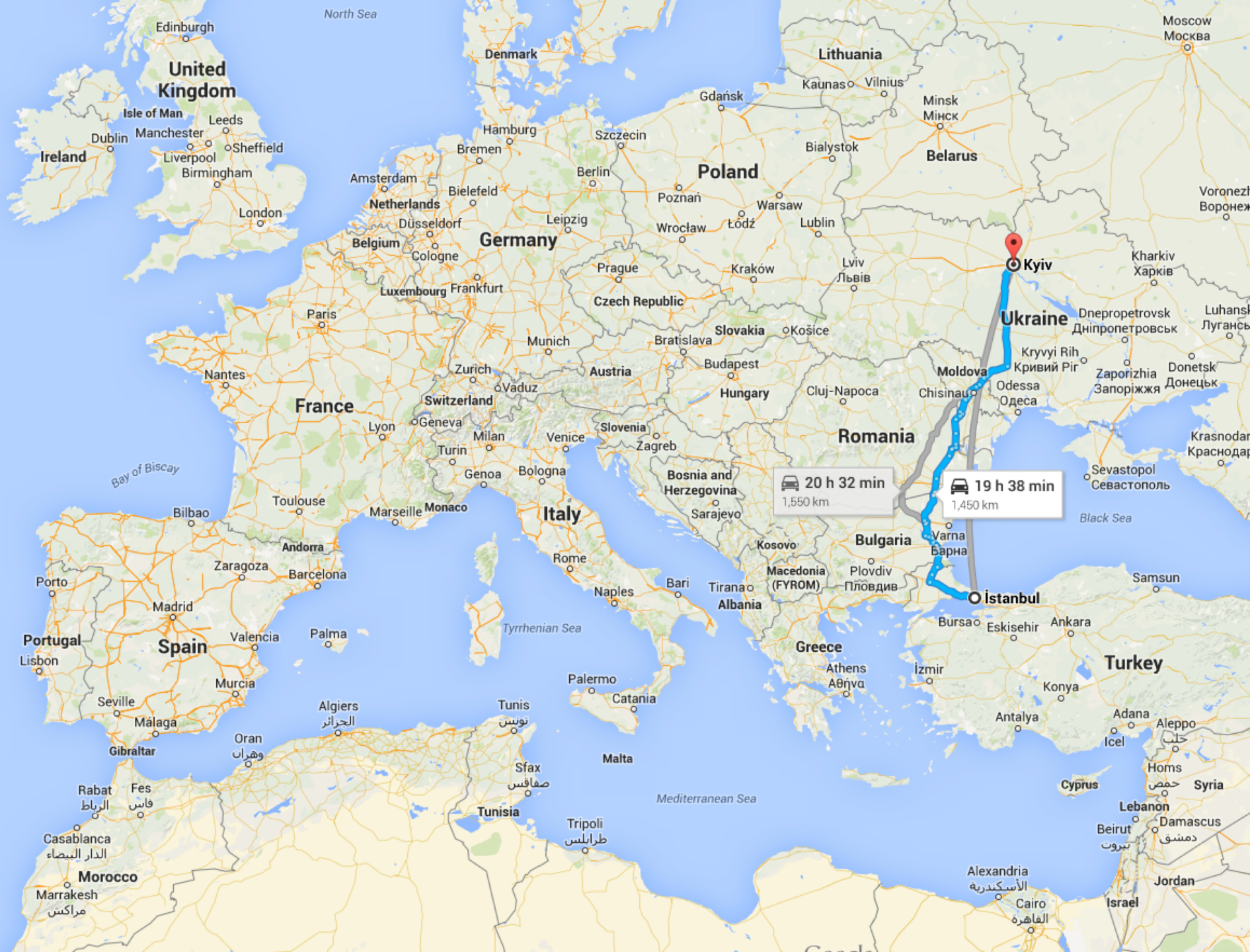
*Bioinformatics Algorithms: An Active Learning Approach*

*Bioinformatics Algorithms: An Active Learning Approach.* © 2018 Compeau and Pevzner.

Modern boundaries are shown for reference.

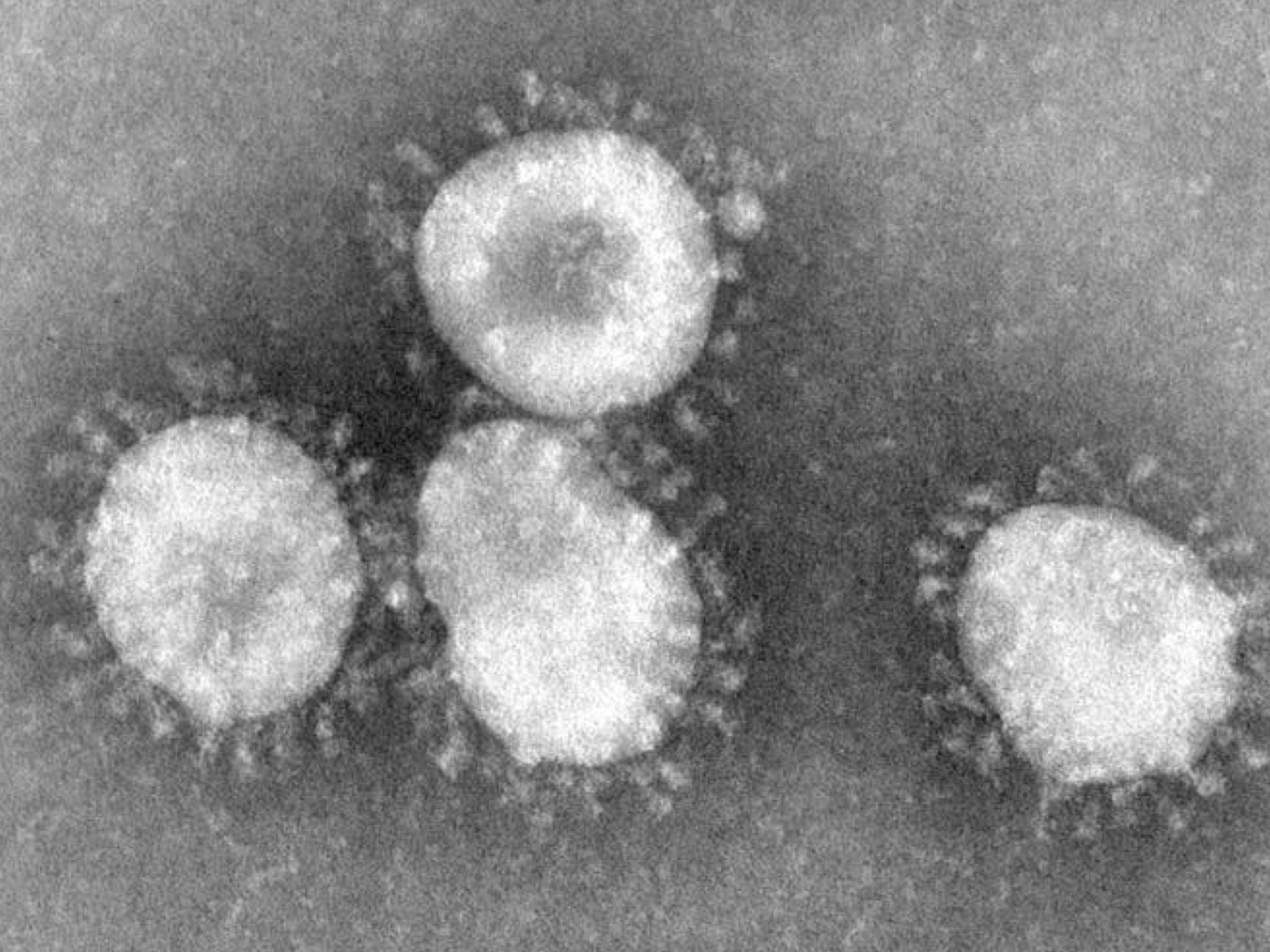Cities with repeated plague outbreaks, 14th to 18th century

Spreading of the plague by marine routes

Approximate date of the first outbreak

| | | | |
|---|---|---|---|
| 1347 | | 1350 | |
| 1348 | | 1351 | |
| 1349 | | relatively unaffected | |

Lübeck
London
Amsterdam
Rotterdam
Antwerp
Aachen
Danzig (Gdansk)
Kiev
Paris
Nürnberg
approaching from Asia, 1346
Strassburg (Strasbourg)
Vienna
Bordeaux
Lyon
Florence
Venice
Trieste
Genoa
Marseille
Ancona
Ragusa (Dubrovnik)
Lisbon
Barcelona
Rome
Valencia
Naples
Seville
Constantinople (Istanbul)
approaching from Asia, 1346

© 2012 Encyclopædia Britannica, Inc.

20 h 32 min
1,550 km

19 h 38 min
1,450 km

Ireland
1 case

China

Canada
29 cases

United States
1 case

Guangdong
Province

Metropole
Hotel

Hong Kong
Special
Administrative
Region
195 cases

Vietnam
58 cases

Philippine
Sea

Arabian
Sea

Bay
of
Bengal

South
China
Sea

Singapore
71 cases

Indian

The Spread
of ?????

Ireland
1 case

Canada
29 cases

United States
1 case

China

Guangdong
Province

Metropole
Hotel

Hong Kong
Special
Administrative
Region
195 cases

Vietnam
58 cases

Philippine
Sea

Arabian
Sea

Bay
of
Bengal

South
China
Sea

Singapore
71 cases

**The Spread
of SARS**

Indian

# 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**).
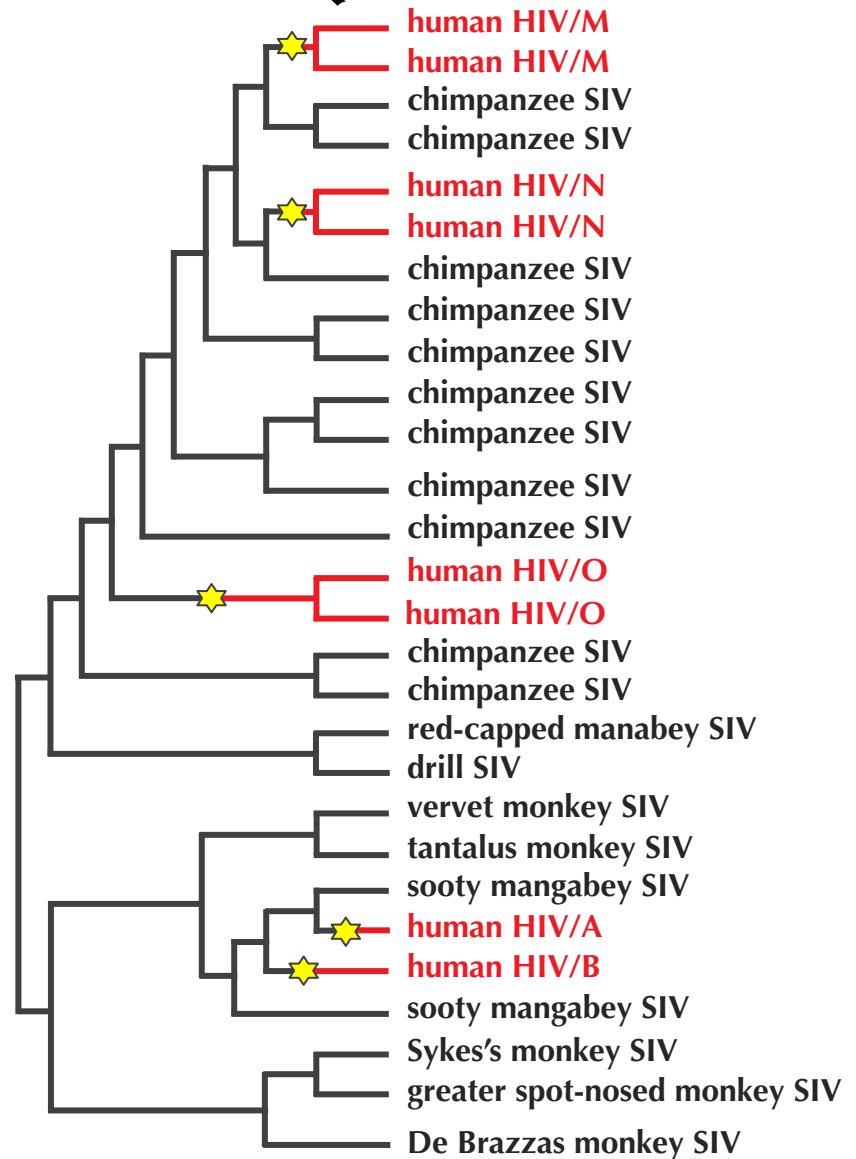
# Example: HIV Evolutionary Tree



| | |
|---|---|
| —— | SIVs (monkeys) |
| —— | HIV (human) |
| ★ | human infection |

human HIV/M
human HIV/M
chimpanzee SIV
chimpanzee SIV
human HIV/N
human HIV/N
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
human HIV/O
human HIV/O
chimpanzee SIV
chimpanzee SIV
red-capped manabey SIV
drill SIV
vervet monkey SIV
tantalus monkey SIV
sooty mangabey SIV
human HIV/A
human HIV/B
sooty mangabey SIV
Sykes's monkey SIV
greater spot-nosed monkey SIV
De Brazzas monkey SIV

# Two Computational Questions



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

Can we infer anything about the ancestors?

human HIV/M
human HIV/M
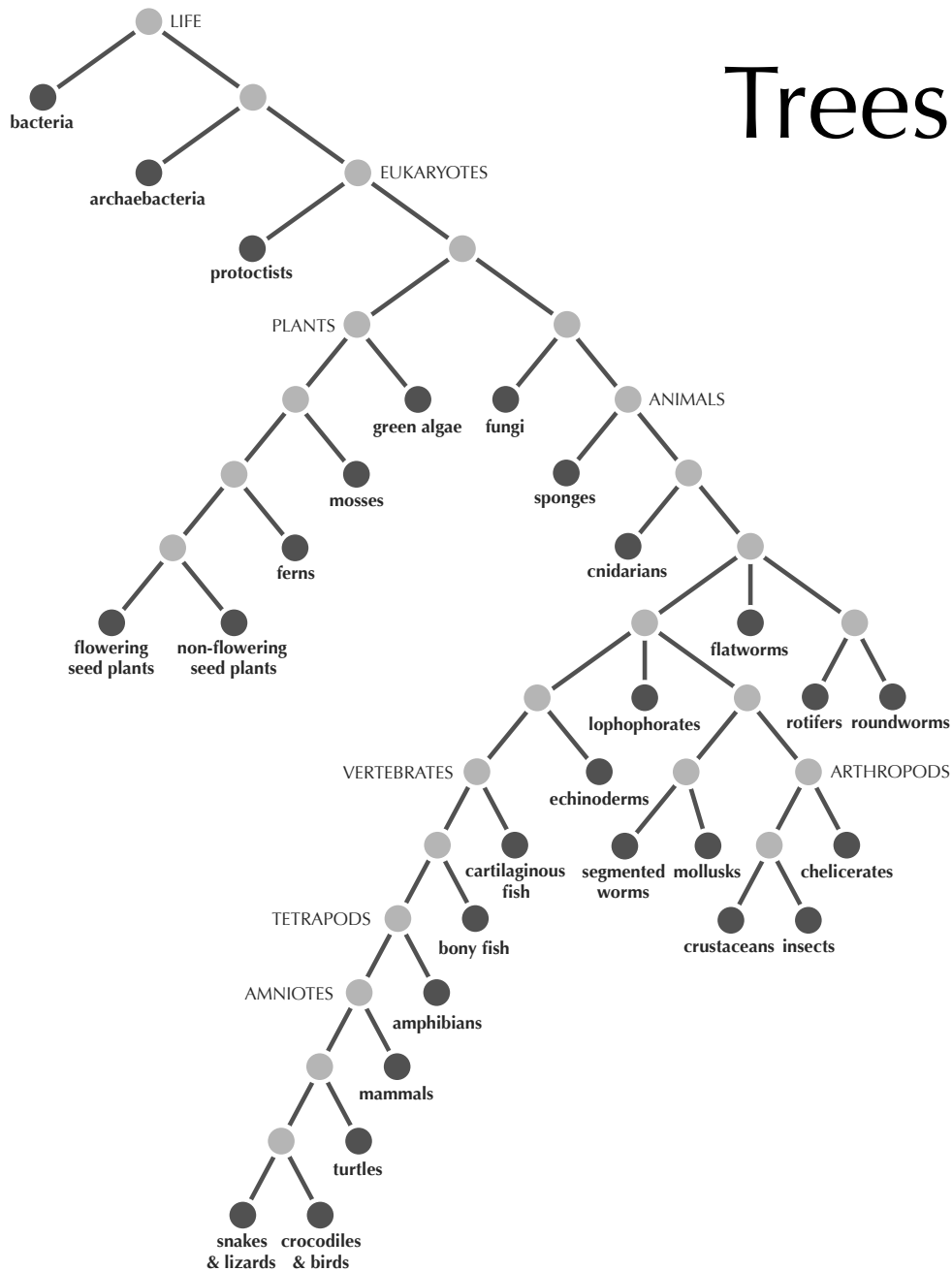chimpanzee SIV
chimpanzee SIV
human HIV/N
human HIV/N
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
human HIV/O
human HIV/O
chimpanzee SIV
chimpanzee SIV
red-capped manabey SIV
drill SIV
vervet monkey SIV
tantalus monkey SIV
sooty mangabey SIV
human HIV/A
human HIV/B
sooty mangabey SIV
Sykes's monkey SIV
greater spot-nosed monkey SIV
De Brazzas monkey SIV

# 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?

human HIV/M
human HIV/M
chimpanzee SIV
chimpanzee SIV
human HIV/N
human HIV/N
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
chimpanzee SIV
human HIV/O
human HIV/O
chimpanzee SIV
chimpanzee SIV
red-capped manabey SIV
drill SIV
vervet monkey SIV
tantalus monkey SIV
sooty mangabey SIV
human HIV/A
human HIV/B
sooty mangabey SIV
Sykes's monkey SIV
greater spot-nosed monkey SIV
De Brazzas monkey SIV

# Trees



LIFE

bacteria

archaebacteria

protoctists

EUKARYOTES

PLANTS

green algae  fungi

mosses

ferns

flowering
seed plants

non-flowering
seed plants

ANIMALS

sponges

cnidarians

flatworms

lophophorates

rotifers  roundworms

VERTEBRATES

echinoderms

ARTHROPODS

cartilaginous
fish

segmented  mollusks
worms

chelicerates

TETRAPODS

bony fish

crustaceans insects

AMNIOTES

amphibians

mammals

turtles

snakes  crocodiles
& lizards  & birds

**Tree:** Connected graph containing no cycles.

*Bioinformatics Algorithms: An Active Learning Approach.* © 2018 Compeau and Pevzner.

# Trees



**Tree:** Connected graph containing no cycles.

**Leaves** (degree = 1): present-day species

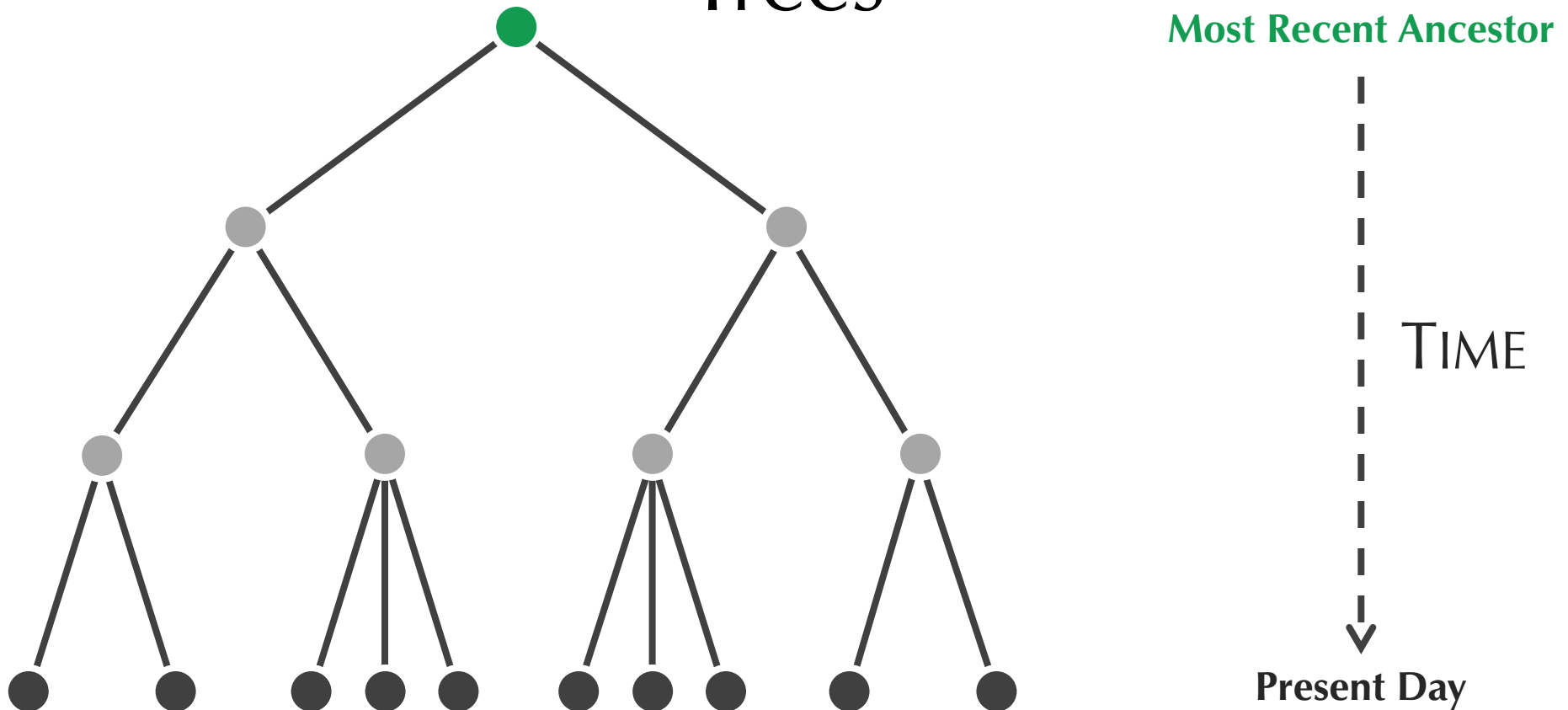*Bioinformatics Algorithms: An Active Learning Approach.* © 2018 Compeau and Pevzner.

# Trees



**Tree:** Connected graph containing no cycles.

**Leaves** (degree = 1): present-day species

**Internal nodes** (degree ≥ 2): ancestral species

*Bioinformatics Algorithms: An Active Learning Approach.* © 2018 Compeau and Pevzner.

# 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



**Most Recent Ancestor**

TIME

**Present Day**

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

# 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** | 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

**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** | ACGTAGGCCT | 0 | **3** | 6 | 4 |
| **Human** | ATGTAAGACT | **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     |

*Bioinformatics Algorithms: An Active Learning Approach.* © 2018 Compeau and Pevzner.

# "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 |



*Bioinformatics Algorithms: An Active Learning Approach.* © 2018 Compeau and Pevzner.

# "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     |



*Bioinformatics Algorithms: An Active Learning Approach.* © 2018 Compeau and Pevzner.

# 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?

**Chimp** ●

1

**Seal** ●

2

3

**Human** ●

2

0

**Whale** ●

**Seal** ●

0.5

**Chimp** ●

1

1.5

degree = 2

3

degree = 2

1

0

**Whale** ●

1

**Human** ●

# 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



$$d_{i,k} = d_{i,m} + d_{k,m}$$

$$d_{i,j} = d_{i,m} + d_{j,m}$$

$$d_{j,k} = d_{j,m} + d_{k,m}$$

# Toward a Recursive Algorithm



$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$
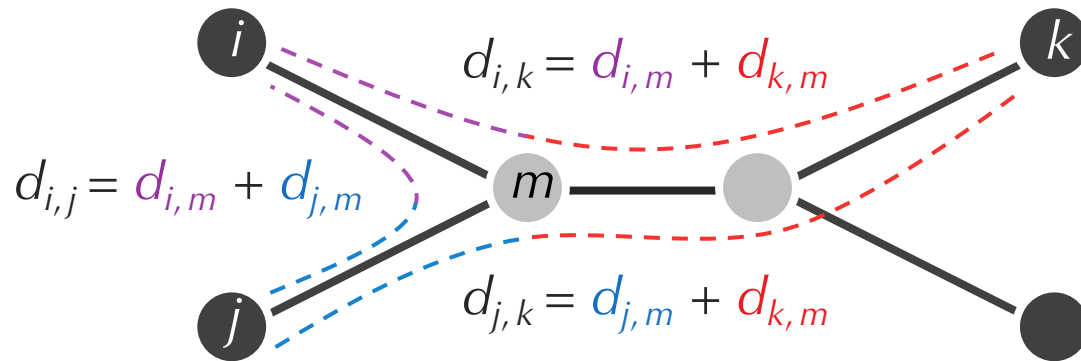
# Toward a Recursive Algorithm



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

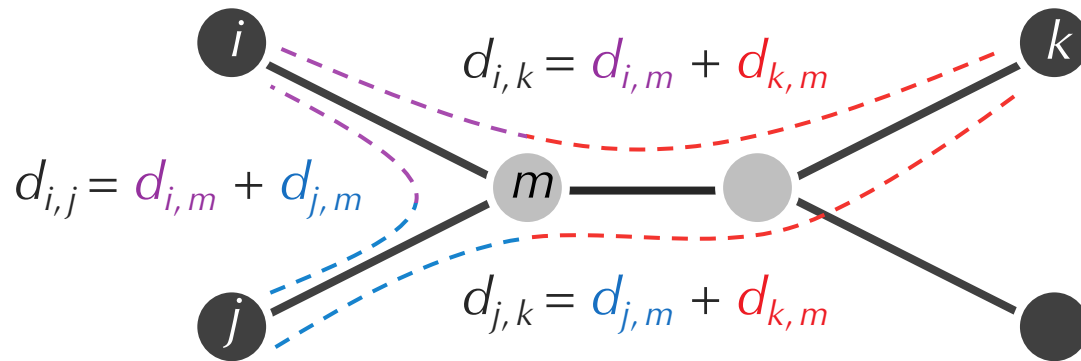# Toward a Recursive Algorithm



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

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

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

# Toward a Recursive Algorithm



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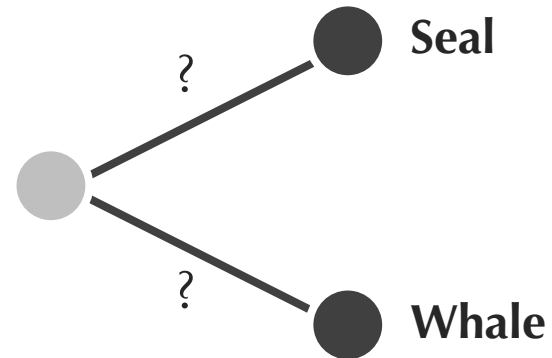
$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

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

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

# Toward a Recursive Algorithm



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

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

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

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

$$d_{i,m} = (D_{i,k} + D_{i,j} - D_{j,k}) / 2$$
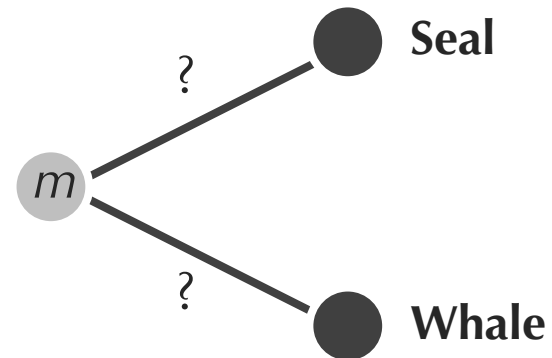
# An Idea for Distance-Based Phylogeny

|       | 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

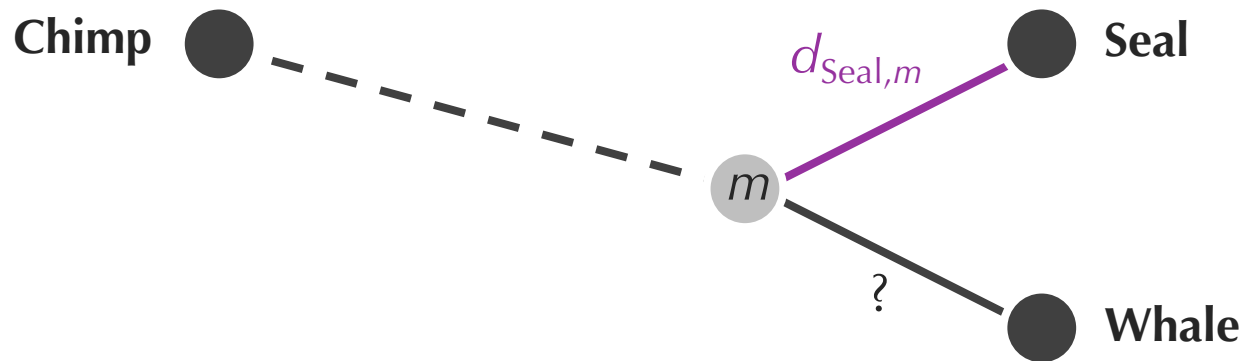|  | Chimp | Human | Seal | Whale |
|---|---|---|---|---|
| **Chimp** | 0 | 3 | 6 | 4 |
| **Human** | 3 | 0 | 7 | 5 |
| **Seal** | 6 | 7 | 0 | 2 |
| **Whale** | 4 | 5 | 2 | 0 |



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

# An Idea for Distance-Based Phylogeny
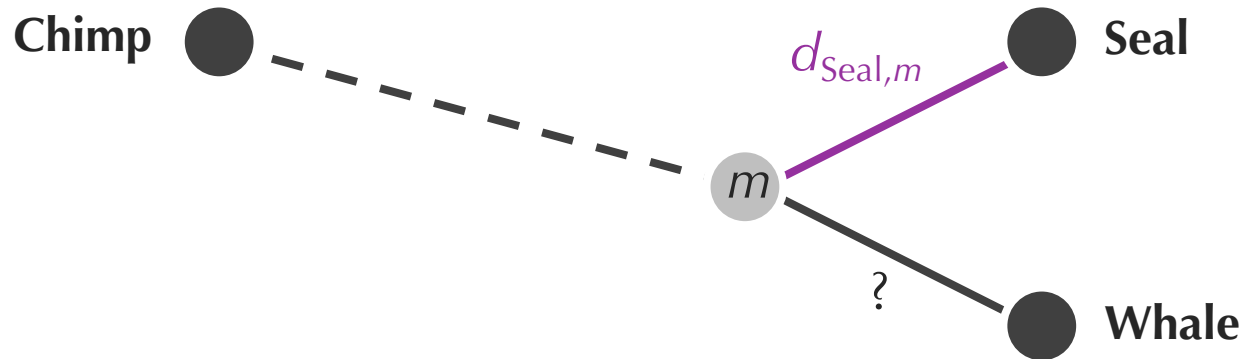
|        | Chimp | Human | Seal | Whale |
|--------|-------|-------|------|-------|
| Chimp  | 0     | 3     | 6    | 4     |
| Human  | 3     | 0     | 7    | 5     |
| Seal   | 6     | 7     | 0    | 2     |
| Whale  | 4     | 5     | 2    | 0     |



$$d_{\text{Seal},m} = (D_{\text{Seal,Chimp}} + D_{\text{Seal,Whale}} - D_{\text{Whale,Chimp}}) / 2$$
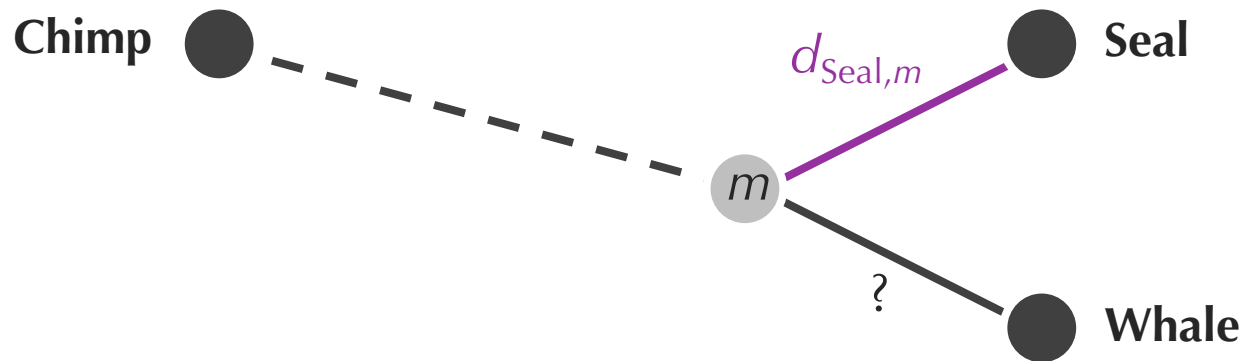
# An Idea for Distance-Based Phylogeny

|        | Chimp | Human | Seal | Whale |
|--------|-------|-------|------|-------|
| **Chimp** | 0 | 3 | 6 | 4 |
| **Human** | 3 | 0 | 7 | 5 |
| **Seal**  | 6 | 7 | 0 | 2 |
| **Whale** | 4 | 5 | 2 | 0 |

**Chimp** ● - - - - - - - - - - $m$

$d_{\text{Seal},m}$ — ● **Seal**

? ● **Whale**

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

# An Idea for Distance-Based Phylogeny

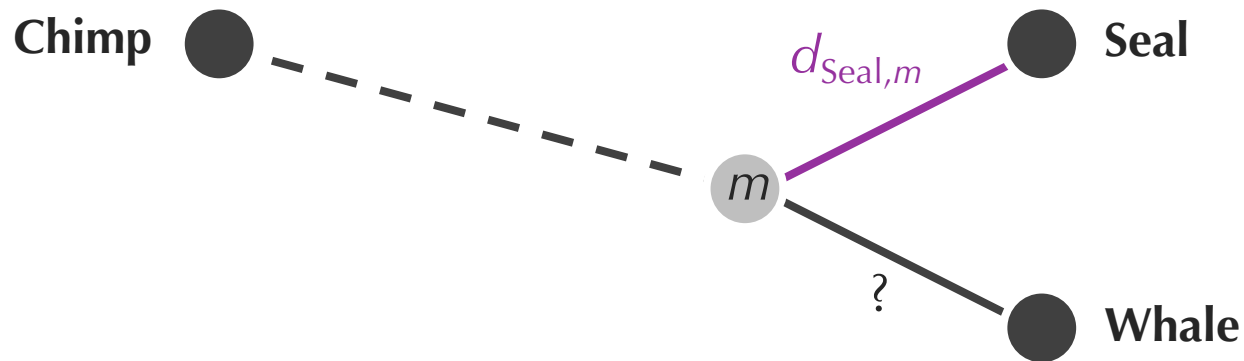|  | Chimp | Human | Seal | Whale |
|---|---|---|---|---|
| **Chimp** | 0 | 3 | 6 | 4 |
| **Human** | 3 | 0 | 7 | 5 |
| **Seal** | 6 | 7 | 0 | 2 |
| **Whale** | 4 | 5 | 2 | 0 |



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

# 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 |

**Chimp** ● - - - - - - - - - $m$ —— ● **Seal**    $d_{\text{Seal},m}$

? ●  **Whale**

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

# 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 | *m* |
|---|---|---|---|---|---|
| **Chimp** | 0 | 3 | 6 | 4 | 4 |
| **Human** | 3 | 0 | 7 | 5 | 5 |
| **Seal** | 6 | 7 | 0 | 2 | 2 |
| **Whale** | 4 | 5 | 2 | 0 | 0 |
| ***m*** | 4 | 5 | 2 | 0 | 0 |

# An Idea for Distance-Based Phylogeny

|  | Chimp | Human | Seal | Whale | *m* |
|---|---|---|---|---|---|
| **Chimp** | 0 | 3 | 6 | 4 | 4 |
| **Human** | 3 | 0 | 7 | 5 | 5 |
| **Seal** | 6 | 7 | 0 | 2 | 2 |
| **Whale** | 4 | 5 | 2 | 0 | 0 |
| ***m*** | 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   |

Chimp ● —— ? —— *a* —— ? —— ● Human

*m* —— 2 —— ● Seal
*m* —— 0 —— ● Whale

# An Idea for Distance-Based Phylogeny

|        | Chimp | Human | $m$ |
|--------|-------|-------|-----|
| **Chimp** | 0 | 3 | 4 |
| **Human** | 3 | 0 | 5 |
| $m$    | 4 | 5 | 0 |



$$d_{\text{Chimp},a} = (D_{\text{Chimp},m} + D_{\text{Chimp,Human}} - D_{\text{Human},m}) / 2$$
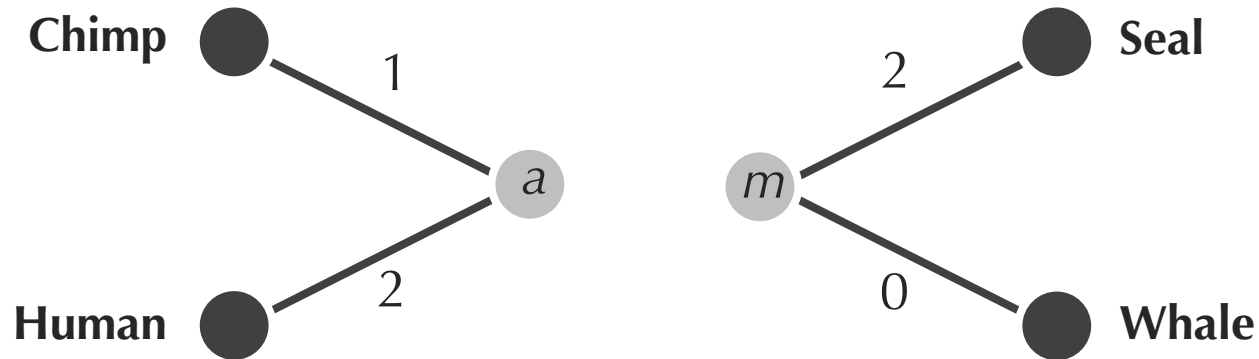
# An Idea for Distance-Based Phylogeny

|  | 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 | *m* |
|--------|-------|-------|-----|
| **Chimp** | 0     | 3     | 4   |
| **Human** | 3     | 0     | 5   |
| ***m***   | 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?

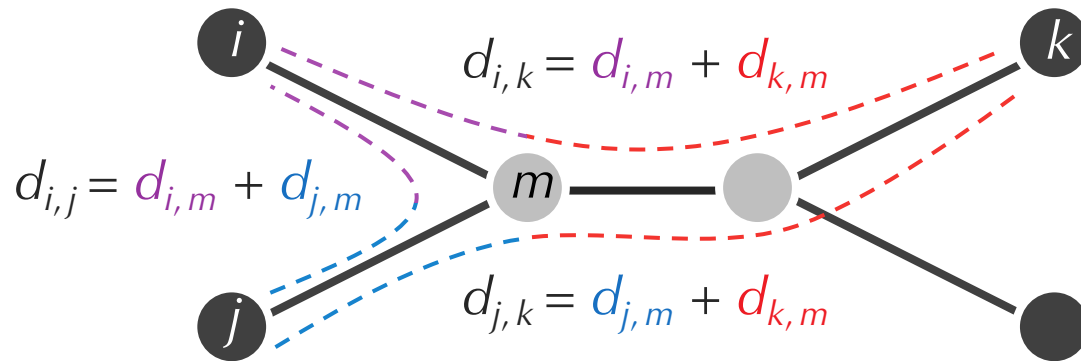|       | $v_1$ | $v_2$ | $v_3$ | $v_4$ |
|-------|-------|-------|-------|-------|
| $v_1$ | 0     | 13    | 21    | 22    |
| $v_2$ | 13    | 0     | **12** | 13   |
| $v_3$ | 21    | 12    | 0     | 13    |
| $v_4$ | 22    | 13    | 13    | 0     |

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



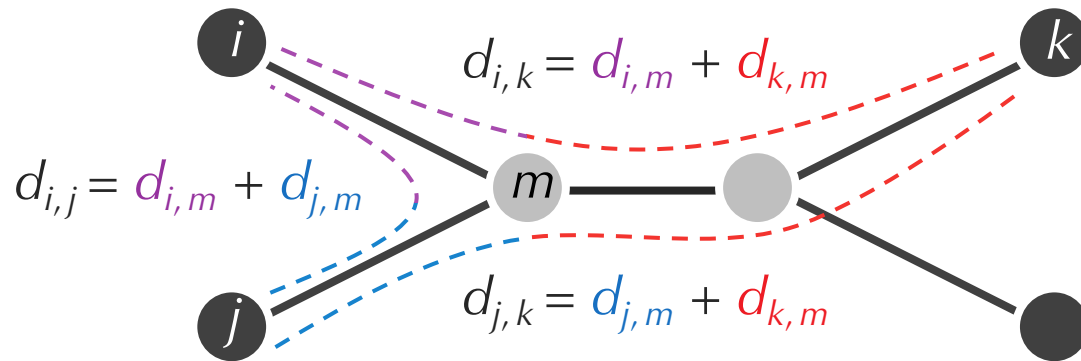$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

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

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

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

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

# From Neighbors to Limbs



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

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

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

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

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

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*.