

Which Animal Gave Us SARS?

Evolutionary Tree Reconstruction

Phillip Compeau and Pavel Pevzner

Bioinformatics Algorithms: An Active Learning Approach

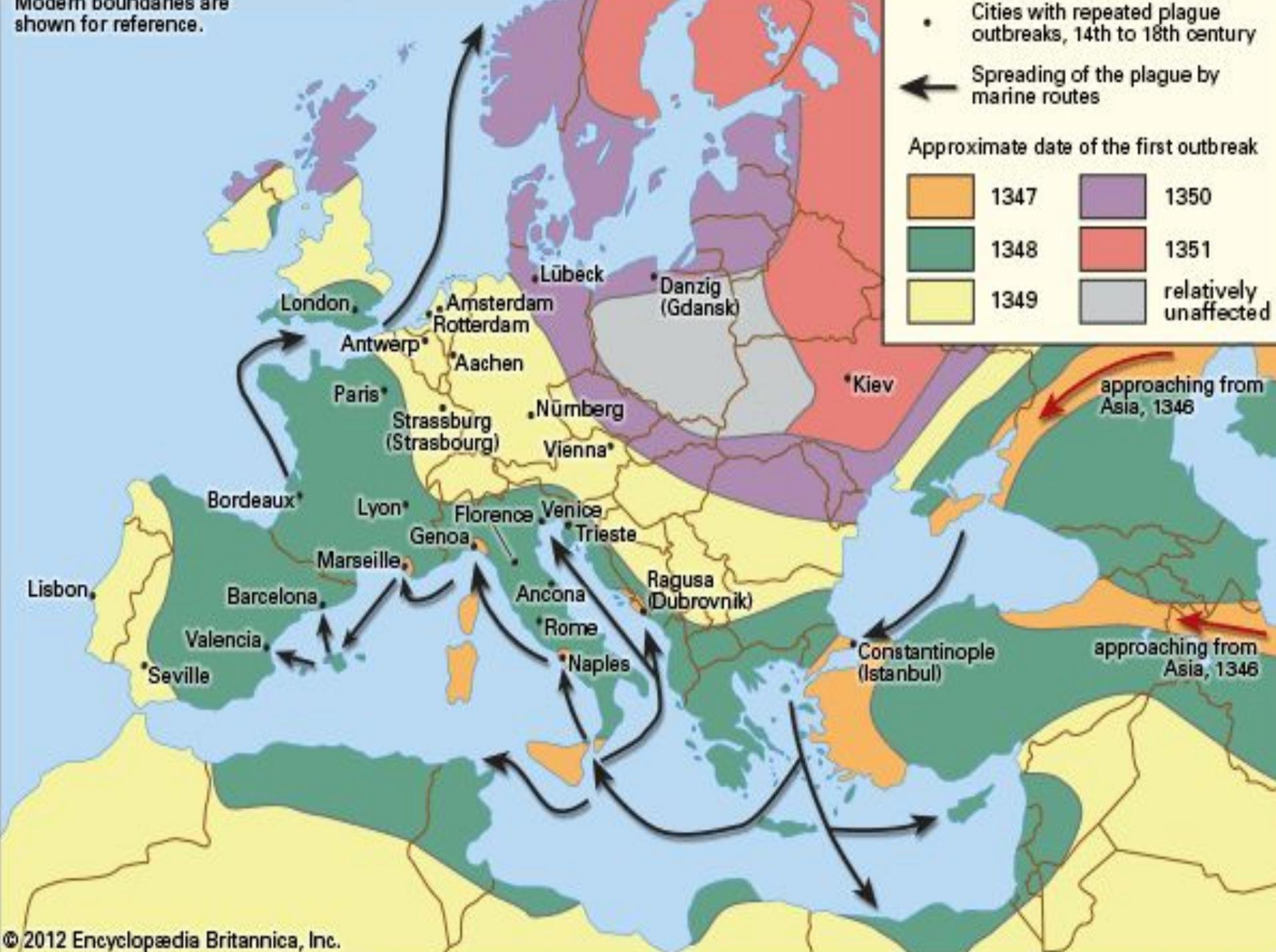
©2018 by Compeau and Pevzner. All rights reserved.

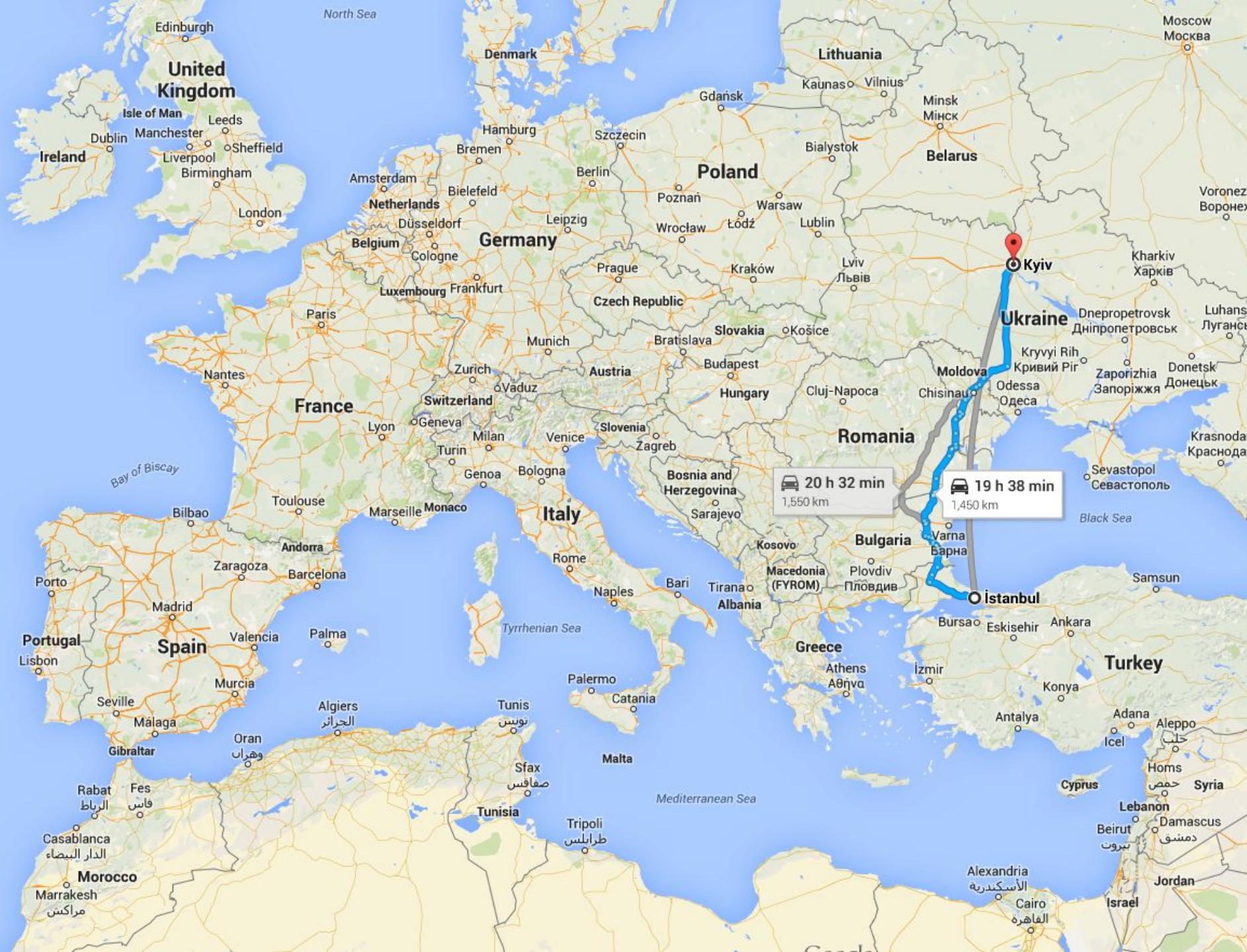
Outline

- **The Fastest Outbreak**
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era



Modern boundaries are shown for reference.

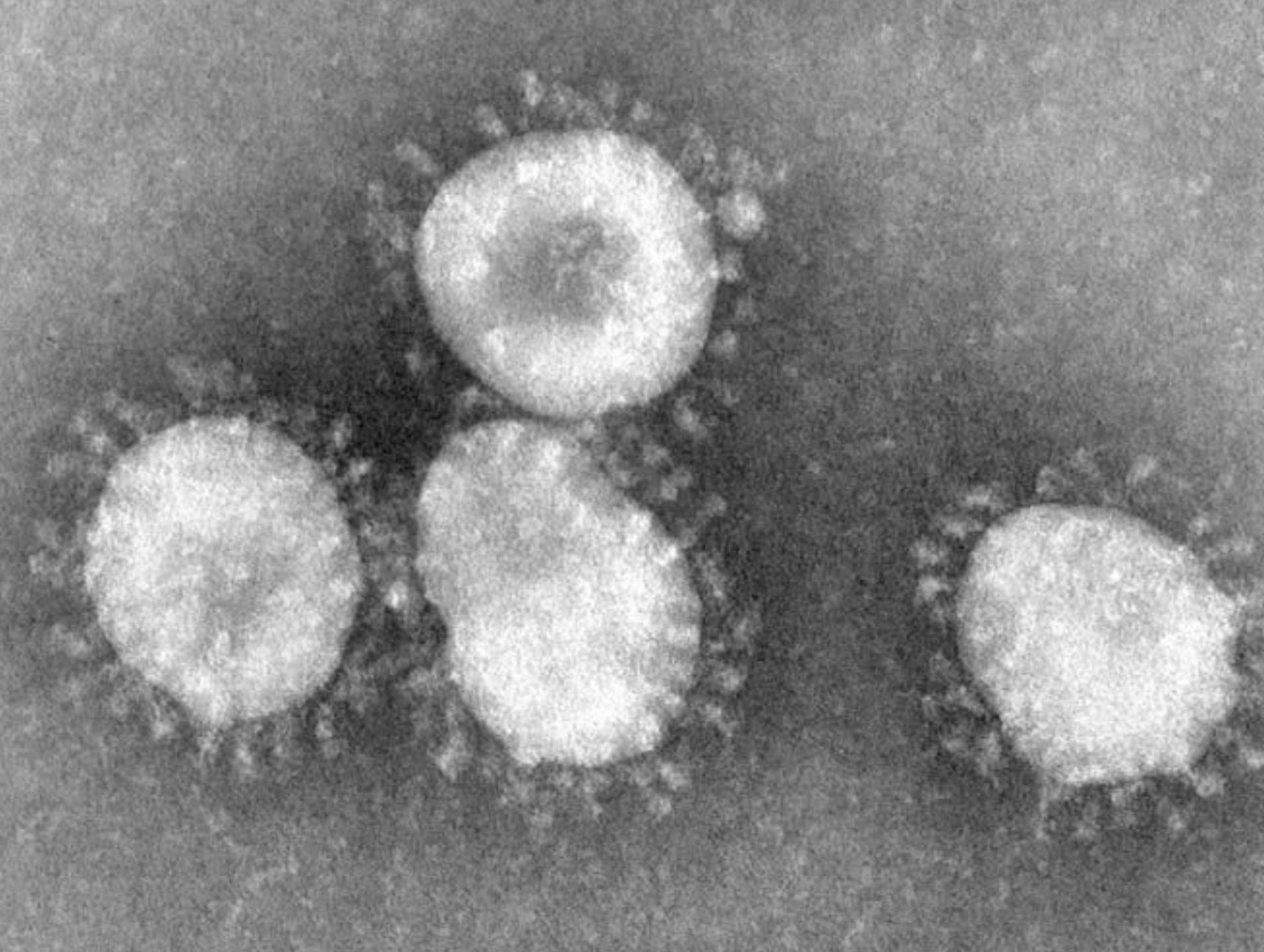


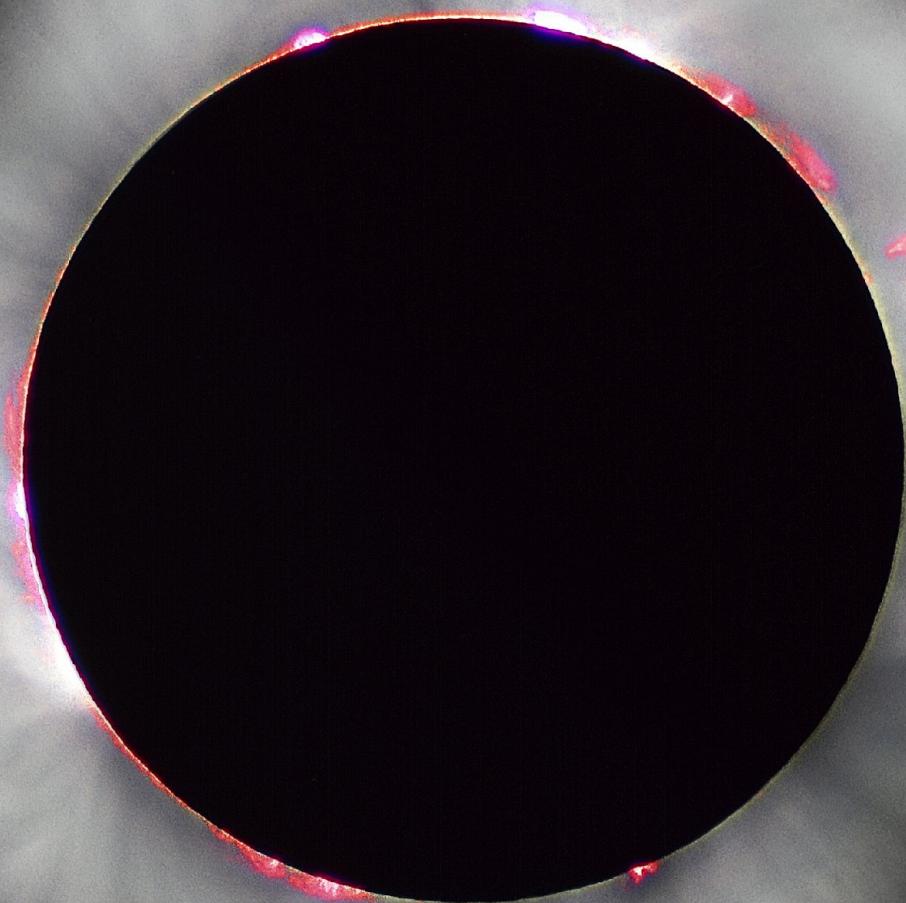






The Spread of SARS





28261 aauuaauacug cgucuugguu cacagcucuc acucagcaug gcaaggagga acuuagauuc
28321 ccucgaggcc agggcggucc aaucaacacc aauagugguc cagaugacca aauuggcuac
28381 uaccgaagag cuacccgacg aguucguggu ggugacggca aaaugaaaga gcucagcccc
28441 agaugguacu ucuaauuaccu aggaacuggc ccagaagcuu cacuuccua cggcgcuac
28501 aaagaaggca ucguauuggu ugcaacugag ggagccuuga auacacccaa agaccacaauu
28561 ggcacccgca auccuaauaa caaugcugcc accgugcuac aacuuccuca aggaacaaca
28621 uugccaaaag gcuucuacgc agagggaaagc agaggcggca gucaagccuc uucucgcucc
28681 ucaucacgua gucgcgguaa uucaagaaau ucaacuccug gcagcaguag gggaaaauuc
28741 ccugcucgaa uggcuagcgg aggugugugaa acugcccucg cgcuaauugcu gcuagacaga
28801 uugaaccagc uugagagcaa aguuucuggu aaaggccaac aacaacaagg ccaaacuguc
28861 acuaagaaau cugcugcuga ggcaucuaaa aagccucgcc aaaaacguac ugccacaaaa
28921 caguacaacg ucacucaagc auuugggaga cgugguccag aacaaaccca agaaaauuuc
28981 ggggaccaag accuaaucag acaaggaacu gauuacaaac auuggccgca aauugcaca
29041 uuugcuccaa gugccucugc auucuuugga augucacgca uuggcaugga agucacaccu
29101 ucgggaacau ggcugacuua ucauggagcc auuaaaugg augacaaaga uccacaauuc
29161 aaagacaacg ucauacugcu gaacaagcac auugacgcau acaaaacaau cccaccaaca
29221 gagccuaaaa aggacaaaaaa gaaaaagacu gaugaagcuc agccuuugcc gcagagacaa
29281 aagaagcagc ccacugugac ucuucuuccu gcggcugaca uggaugauuu cuccagacaa
29341 cuucaaaaauu ccaugagugg agcuucugcu gauucaacuc aggcauaaac acucaugaug
29401 accacacaag gcagaugggc uauguaaacg uuuucgcaau uccguuuuacg auacauaguc
29461 uacucuugug cagaaugaaau ucucguaacu aaacagcaca aguagguuua guuaacuuua
29521 aucucacaua gcaaucuuua aucaaugugu aacauuaggg aggacuugaa agagccacca
29581 cauuucauc gagggccacgc ggaguacgau cgaggguaca gugaauuaug cuagggagag
29641 cugccuauau ggaagagccc uaauguguaa aauuaauuuu aguagugcua uccccaugug
29701 auuuaauag cuucuuagga gaaugacaaa aaaaaaaaaa aaaaaaaaaa a

28261 aauuaauacug cgucuugguu cacagcucuc acucagcaug gcaaggagga acuuagauuc
28321 ccucgaggcc agggcggucc aaucaacacc aa**C**agugguc cagaugacca aauuggcuac
28381 uac**A**gaagag cuacccgacg aguucguggu ggugacggca aaaugaaaga gcucagcccc
28441 agaugguacu ucuaauuaccu aggaacuggc ccagaagcuu cacuuccua cggcgcuac
28501 aaagaaggca ucguauuggu ugcaacugag ggagccuuga auac**C**cccaa agaccacaauu
28561 ggcacccgca auccuaauaa caaugcugcc accgugcuac aacuuccuca aggaacaaca
28621 uugccaaaag gcuucuacgc agagggaaagc agaggcggca gucaagccuc uucucgcucc
28681 ucaucacgua gucgcgguaa uucaagaaau ucaacuccug gcagcaguag gggaaaauuc
28741 ccugcucgaa uggcuagcgg agguggugaa acugcccucg cgcuaauugcu gcuagacaga
28801 uugaaccagc uugagagcaa aguuucuggu aaaggccaac aacaacaagg ccaaacuguc
28861 acuaagaaau cugcugcuga ggcaucuaaa aagccucgcc aaaaacguac ugccacaaaa
28921 caguacaacg ucacucaagc auuugggaga cgugguccag aacaaaccca agaaaauuuc
28981 ggggaccaag ac**U**uaaucag acaaggaacu gauuacaaac auuggccgca aauugcaca
29041 uuugcuccaa gugccucugc auucuuugga augucacgca uuggcaugga agucacaccu
29101 ucgggaacau ggcugacuua ucauggagcc auuaaaugg augacaaaga uccacaauuc
29161 aaagacaacg ucauacugcu gaacaagcac auugacgcau aca**U**aacaau cccaccaaca
29221 ga**U**ccuaaaaa aggacaaaaaa gaaaaagacu gaugaagcuc agccuuugcc gcagagacaa
29281 aagaagcagc ccacugugac ucuucuuccu gcggcugaca uggaugauuu cuccagacaa
29341 cuucaaaaauu ccaugagugg agcuucugcu gauucaacuc aggcauaaac acucaugaug
29401 accacacaag gcagaugggc uauguaaacf uuuucgcaau uccguuuuacg auacauaguc
29461 uacucuugug cagaaugaaau ucucguaacu aaacagcaca aguagguuua guuaacuuua
29521 aucucacaua gcaaucuuua aucaa**G**ugu aacauuaggg aggacuugaa agagccacca
29581 cauu**A**ucauc gagggccacgc ggaguacgau cgaggguaca gugaaauaaug cuagggagag
29641 cugccuauau ggaagagccc uaauguguaa aauuaauuuu aguagugcua uccccaugug
29701 auuuaauag cuucuuagga gaaugacaaa aaaaaaaaaa aaaaaaaaaa a

28261 aauuaauacug cgucuugguu cacagcucuc acucagcaug gcaaggagga acuuagauuc
28321 ccucgaggcc agggcggucc aaucaacacc aa**C**agugguc cagaugacca aauuggcuac
28381 uac**A**gaagag cuacccgacg aguucguggu ggugacggca aaaugaaaga gcu**G**agcccc
28441 agaugguacu ucuaauuaccu aggaacuggc ccagaagcuu cacuuccua cggcgcuac
28501 aaagaaggca ucguau**U**ggu ugcaacugag ggagccuuga auac**C**cccaa agaccacaauu
28561 ggcacccgca auccuaauaa caaugcugcc accgugcuac aacuuccuca aggaacaaca
28621 uugccaaaag gcuuc**A**acgc agagggaaagc agaggcggca gucaagccuc uucucgcucc
28681 ucaucacgua gucgcgguaa uucaagaaau ucaacuccug gcagcaguag gggaaaauuc
28741 **cG**ugcucgaa ugcuagcgg aggugugugaa acugcccucg cgcuauugcu gcuagacaga
28801 uugaaccagc uugagagcaa aguuucuggu aaaggccaac aacaacaagg cca**G**acuguc
28861 acuaagaaau cugcugcuga ggcaucu**C**aa aagccucgcc aaaaacguac ugccacaaaa
28921 caguacaacg ucacucaagc auuugggaga cgugguccag aacaaaccca agaaaauuuc
28981 ggggaccaag ac**U**uaaucag acaaggaacu gauuacaaac auuggccgca aauugcaca
29041 uuugcuccaa gugccucugc auucuuugga augu**U**acgca uuggcaugga agucacaccu
29101 ucgggaacau ggcugacuu ucauggagcc auuaaaugg augacaaaga uccacaauuc
29161 aaagacaacg ucauacugcu gaacaagcac auugacgcau aca**U**aacaau cccaccaaca
29221 ga**U**ccuaaaaa aggacaaaaa gaaaaagacu gaugaagcuc agccuuugcc gcagagacaa
29281 aagaagcagc ccacugugac uciucuuccu gcggcugaca uggaugauuu cuccagacaa
29341 cuucaaaaauu ccaugagugg a**Ac**uuucugcu gauucaacuc ag**U**cauuaac acucaugaug
29401 acca**A**acaag gcagaugggc uauguaaacg uuuucgcaau uccguuuuacg auacauaguc
29461 uacucuugug cagaaugaaau ucucguaacu aaacagcaca aguagguuua guuaac**A**uua
29521 aucucacaua gcaaucuuua aucaa**G**ugu aacauuaggg aggacuugaa agagccacca
29581 cauu**A**ucauc gagggccacgc ggaguacgau cgaggguaca gugaaauaaug cuagggagag
29641 cugccuauau ggaagagccc uaauguguaa aauuaauuuu aguagugcua uccccaugug
29701 auuuaauag cuucuuagga gaaugacaaa aaaaaaaaaa aaaaaaaaaa a

28261 Cauaaauacug cgucuugguu caca**U**cucuc acucagcaug gcaaggagga acuuagauuc
28321 ccucgaggcc agggcguucc aaucaacacc aa**C**agugguc cagaugacca aauuggcuac
28381 uac**A**gaagag cuacccgacg aguucguggu ggugacggca aaaugaaaga gcu**G**agcccc
28441 agaugguacu ucuaauuaccu aggaacuggc ccagaagcuu cacuuccua cggcgcuac
28501 aaagaaggca ucguau**U**ggu ugcaacugag ggagccuuga auac**C**cccaa agaccacaauu
28561 ggcacccgca auccuaauaa caaugcugcc accgugcuac aacuuccuca aggaacaaca
28621 uugccaaaag gcuuc**A**acgc agagggaaagc agaggcggca gucaagccuc uucucgcucc
28681 ucaucacgua gucgcgguaa uucaag**U**aaau ucaacuccug gca**A**caguag gggaaaauuc
28741 **cG**ugcucgaa ugcuagcgg aggugugugaa acugcccucg cgcuauugcu gcuagacaga
28801 uugaaccagc uugagagcaa aguuucuggu aaaggccaac aacaacaagg cca**G**acuguc
28861 acuaagaaaau cugcugcuga ggcaucu**C**aa aagccucgcc aaaaacguac ugccacaaaa
28921 caguacaacg ucacucaagc auuugggaga cgugguccag aacaaaccca agaaaauuuc
28981 ggggaccaag ac**U**uaaucag acaaggaacu gauuacaaac auuggccgca aauugcaca
29041 uuu**A**cuccaa gugccucugc auucuuugga augu**U**acgca uuggcaugga agucacaccu
29101 ucgggaacau ggcugacuu ucauggagcc auuaaaugg augacaaaga uccacaauuc
29161 aaagacaacg ucauacugcu gaacaagcac auugacgcau aca**U**aacaau cccaccaaca
29221 ga**U**ccuaaaaa agg**G**caaaaaa gaaaaagacu gaugaagcuc agccuuugcc gcagagacaa
29281 aagaagcagc ccacugugac uciucuuccu gcggcugaca uggaugauuu cucca**U**acaa
29341 cuucaaaaauu ccaugagugg a**A**cuucugcu gauucaacuc ag**U**cauuaac acucaugaug
29401 acca**A**acaag gcagaugggc uauguaaacg uuuucgcaau uccguuuacg auacauaguc
29461 uacucuugug cagaaugaaau ucucguaacu aaacagcaca aguagguuua guuaac**A**uua
29521 aucucacaua gcaaucuuua aucaa**G**gugu acauuaggg aggacuugaa agagccacca
29581 cauu**A**ucauc gagggccacgc ggaguacgau cgaggguaca gugaauuaug cuagggagag
29641 cugccuauau ggaagagccc uaauguguaa aauuaauuuu a**U**uagugcua uccccaugug
29701 auuuaauag cuucuuagga gaaugacaaa aaaaaaaaaa aaaaaaaaaa a

28261 Cauaaauacug cgucuugguu caca**U**cucuc acucagcaug gcaaggagga acuuagauuc
28321 ccucgaggcc agggcgCcc aaucaacacc aa**C**agugguc cagaugacca aauuggcuac
28381 uac**A**gaagag cuacccgacg aguucguggu ggugacggca aaaugaaaga gcu**G**agcccc
28441 agauggguacu ucuaauuaccu aggaacuggc ccagaagcuu cacuuccua cggcgcuac
28501 aaagaaggca ucguau**U**ggu ugcaacu**U**ag ggagccuuga auac**CG**ccaa agaccacaauu
28561 ggcacccgca auccuaauaa caaugcugcc accgugcuac aacuuccuca aggaacaaca
28621 uugccaaaag gcuuc**A**acgc agagggaaagc agaggcgga gucaagccuc uucucgcucc
28681 ucaucacgua gucgcgguaa uucaag**U**aa ucaacuccug gca**A**caguag gggaaaauuc
28741 **cG**ugcucgaa ugcuagcgg aggugugugaa acugcccucg cgcuauugcu gcuagacaga
28801 uugaaccagc uugagagcaa aguuucuggu aaaggccaac aacaacaagg cca**G**acuguc
28861 acuaagaaaau cugcugcuga ggcaucu**C**aa aagccucgcc aaaaacguac ugccacaaaa
28921 caguacaacg ucacuc**C**agc auuugggaga cgugguccag aacaaaccca agaaaauuuc
28981 ggggaccaag ac**U**uaaucag acaaggaacu gauuacaaac auuggccgca aauugcaca
29041 uuu**A**cucca**C** gugccucugc auucuuugga augu**U**acgca uuggcaugga agucacaccu
29101 ucgggaacau ggcugacuu ucauggagcc auuaaaugg augacaaaga uccacaauuc
29161 aaagacaacg ucauacugcu gaacaagcac auugacgcau aca**U**aa**G**auu cccaccaaca
29221 ga**U**ccuaaaaa agg**G**caaaaaa gaaaaagacu gaugaagcuc agccuuugcc gcagagacaa
29281 aagaagcagc ccacugugac uciucuuccu gcggcugaca uggaugauuu cucca**U**acaa
29341 cuucaaaaauu ccaugagugg a**A**cuucugcu gauucaacuc ag**U**cauuaac acucaugaug
29401 acca**A**acaag gcagaugggc uauguaaacg uuuucgcaau uccguuuacg auacauaguc
29461 uacucuugug cagaaugaaau ucucguaacu aaacagcaca aguagguuua guuaac**A**uua
29521 aucucacaua gcaaucuuua aucaa**G**ugu acauuaggg aggacuugaa agagccacca
29581 cauu**A**ucauc gagggccacgc ggaguacgau cgaggguaca gugaaauaaug cua**A**ggagag
29641 cugccuauau ggaagagccc uaau**A**uguua aauuaauuuu a**U**uagugcua uccccaugug
29701 auuuaauag cuucuuagga gaaugacaaa aaaaaaaaaa aaaaaaaaaa a



2017年
勝利
勝利

Questions about SARS

- Which animal gave us SARS?

Questions about SARS

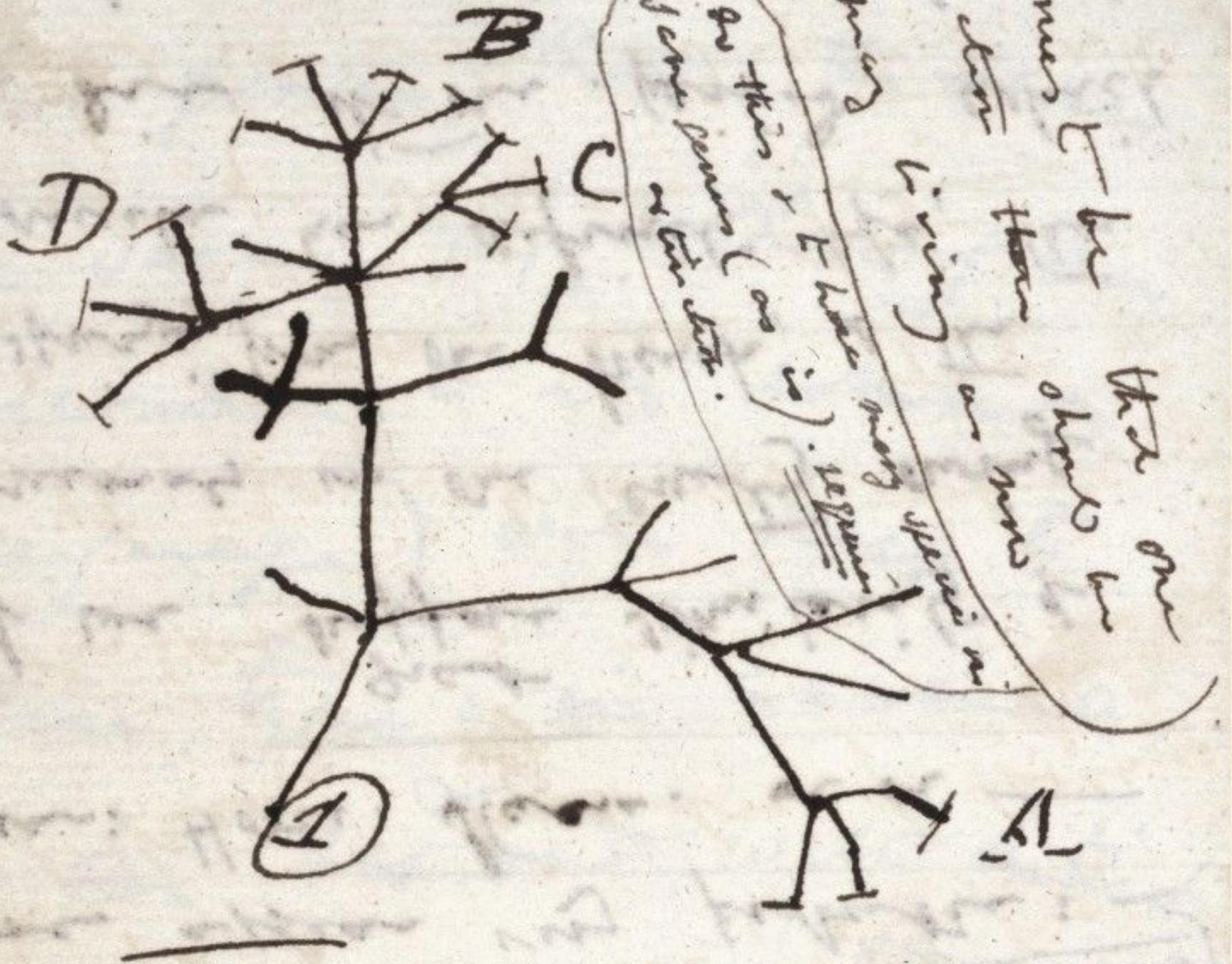
- Which animal gave us SARS?
- How were we first infected?

Questions about SARS

- Which animal gave us SARS?
- How were we first infected?
- How did SARS spread around the world?

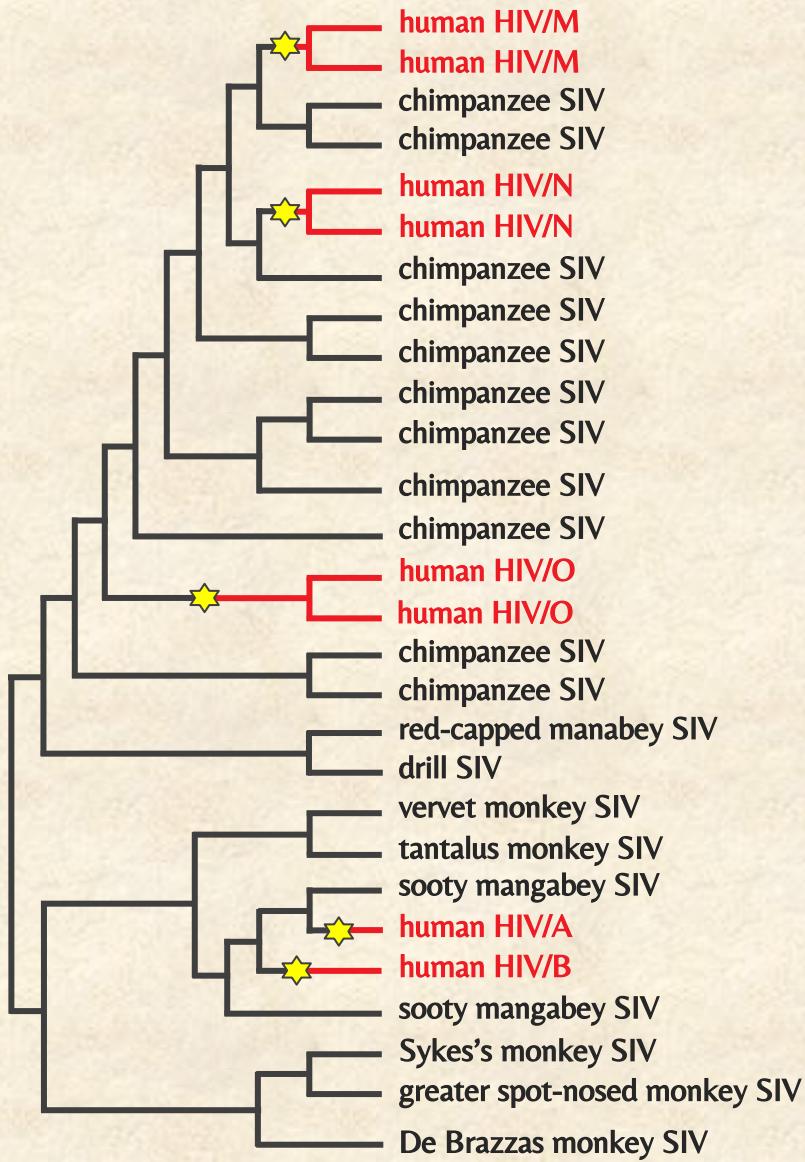
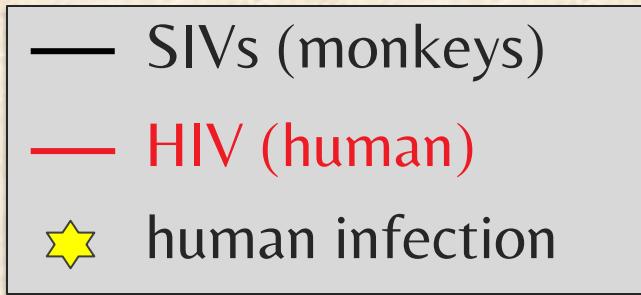
Questions about SARS

- Which animal gave us SARS?
- How were we first infected?
- How did SARS spread around the world?
- All these questions relate to constructing **evolutionary trees** (a.k.a. **phylogenies**).



(the one
that one
should be
more than
than are now
living among species in
ways or to have in). separately
as this or to have in.
so as seems (as it is).
from certain data.
and certain data.

HIV Evolutionary Tree



Outline

- The Fastest Outbreak
- **Transforming Distance Matrices into Evolutionary Trees**
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

Constructing a Distance Matrix

| SPECIES | ALIGNMENT |
|---------|------------|
| | |
| Chimp | ACGTAGGCCT |
| Human | ATGTAAGACT |
| Seal | TCGAGAGCAC |
| Whale | TCGAAAGCAT |

Constructing a 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 |

Constructing a 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 CGT A GGC CT | 0 | 3 | 6 | 4 |
| Human | A TGTA AGA CT | 3 | 0 | 7 | 5 |
| Seal | TCGAGAGCAC | 6 | 7 | 0 | 2 |
| Whale | TCGAAAGCAT | 4 | 5 | 2 | 0 |

Constructing a Distance Matrix

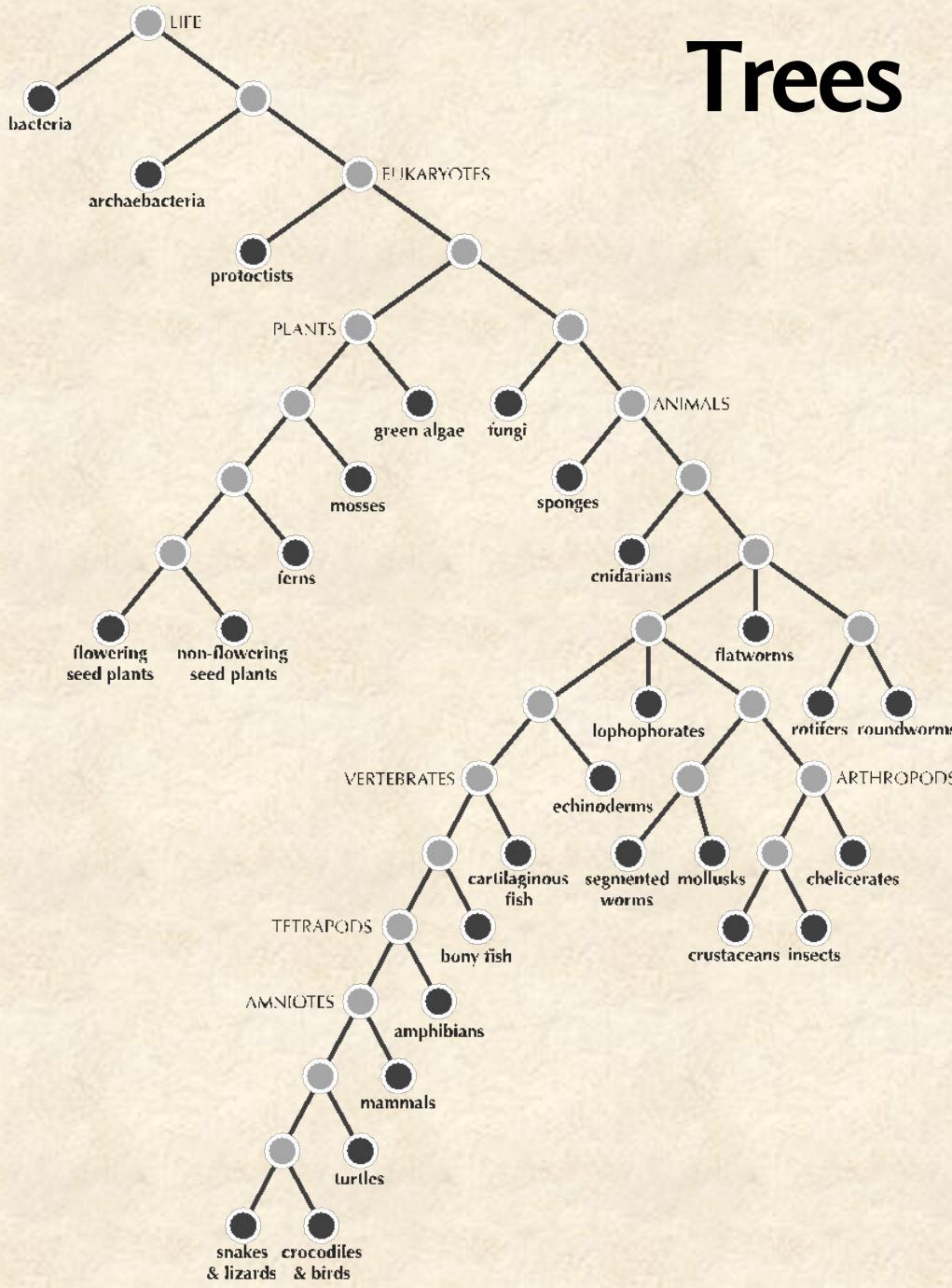
$D_{i,j}$ = number of differing symbols between i -th and j -th rows of a multiple alignment.

| 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 |

STOP and Think: How else could we form a distance matrix?

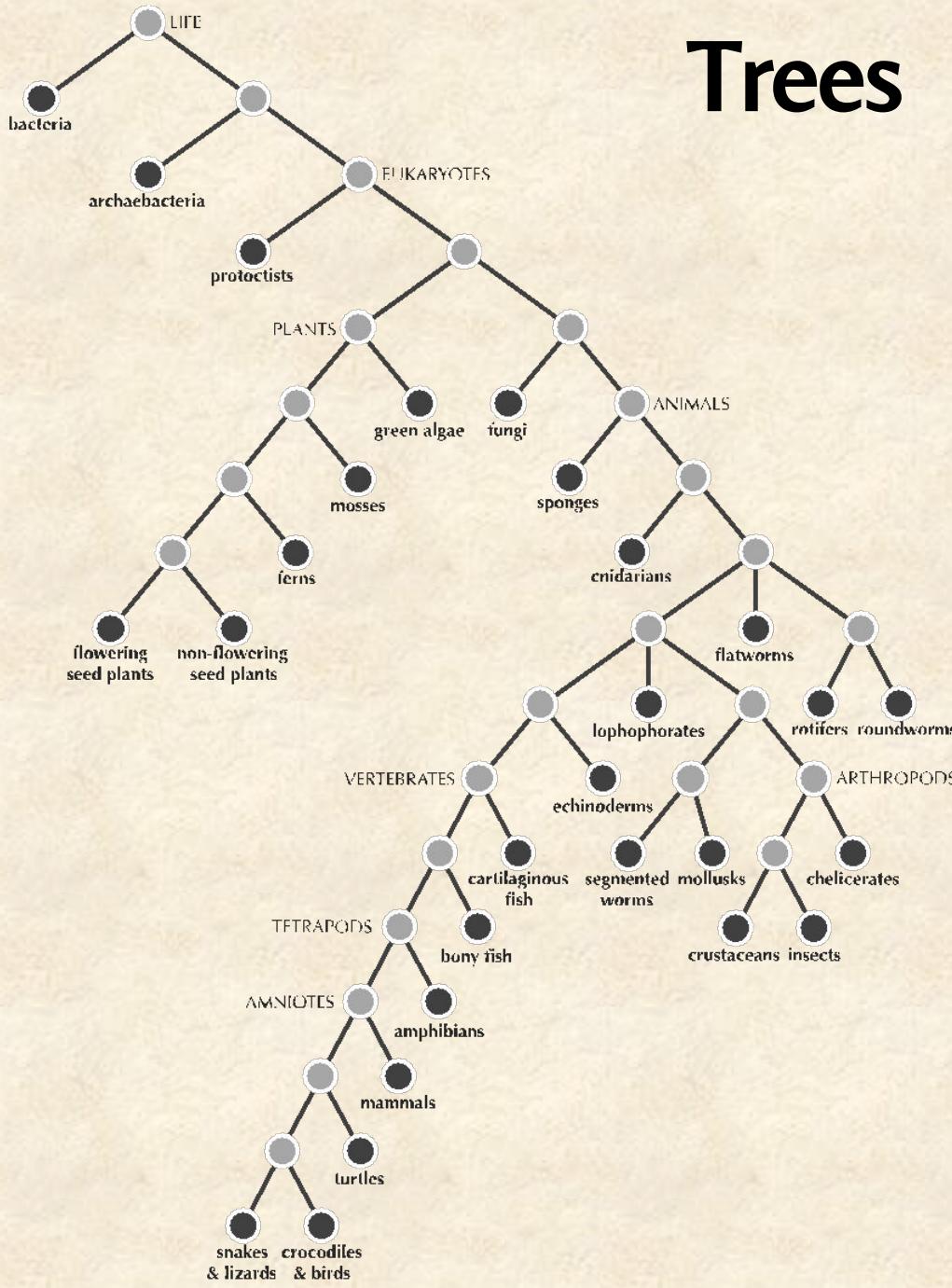
Spike Protein

1 mfifllfltl tsgsdlrct tfddvqapny tqhtssmrgv yypdeifrsd tlyltqdlf1
61 pfysnvtgfh tinhtfgnpv ipfkdgifyfa ateksnvvrg wvfgstmnnk sqsviiinns
121 tnvviracnf elcdnpffav skpmgtqtht mifdnafnct feyisdafsl dvseksgnfk
181 hlrefvfknk dgflyvykgy qpidvvrdlp sgfntlkpif klplgininitn frailtafsp
241 aqdiwgttsaa ayfvgylkpt tfmlkydeng titdavdcsg nplaelkcsv ksfeidkgiy
301 qtsnfrvvps gdvvrfpnit nlcpfgevfn atkfpsvyaw erkkisncva dysvlynstf
361 fstfkcygvs atklndlcs nvyadsfvvk gddvrqiapg qtgviadyny klpddfmvc
421 lawntrnida tstgnynyky rylrhgklrp ferdisnvpf spdgpctpp alncywplnd
481 ygfytttgig yqpyrvvvls fellnapatv cgpklstdli knqcvnfnfn gltgtgvlt
541 sskrfqpfqq fgrdvsdftd svrdpktsei ldispcafgg vsvitpgtna ssevavlyqd
601 vnctdvstai hadqltpawr iystgnvfq tqagcligae hvdtseyedi pigagicasy
661 htvsllrst qksivaytms lgadssiays nntiaiptnf sisittevmp vsmaktsvdc
721 nmyicgdste canlllqygs fctqlnrals giaaeqrnt revfaqvkqm yktptlkyfg
781 gfnfsqilpd plkptkrsfi edllfnkvtt adagfmkqyg eclgdinard licaqkfngl
841 tvlpplltdd miaaytaalv sgtatagwtf gagaalqipf amqmayrfng igvtqnvlye
901 nqkqianqfn kaisqiqesl tttstalgkl qdvvnqnaqa lntlvkqlss nfgaissvln
961 dilsrldkve aevqidrlit grlqlqltyv tqqliraaei rasanlaatk msecvlqgsk
1021 rvdfcgkgyh lmsfpqaaph gvvflhvtyv psqernftta paichegkay fpregvfvn
1081 gtswfitqrn ffspqiiittd ntfvsgncdv vigiinntvy dplqpeldsf keeldkyfkn
1141 htspdvdldg isginasvvn iqkeidrlne vaknlnesli dlqelgkyeq yikwpwyvw1
1201 gfiagliaiv mvtilccmt sccsclkac scgscckfde ddsepvlkgv klhyt



Trees

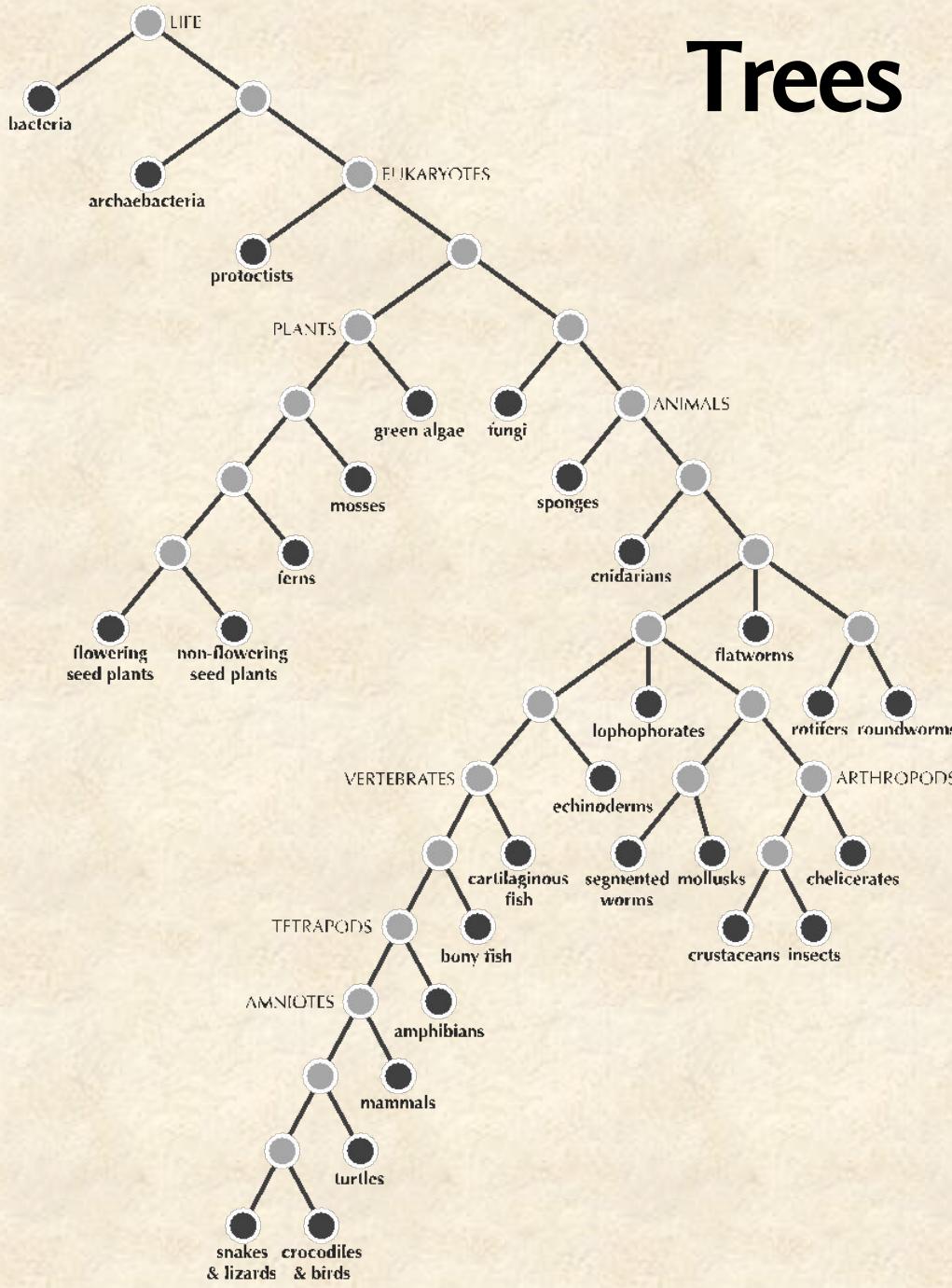
Tree: Connected graph containing no cycles.



Trees

Tree: Connected graph containing no cycles.

Leaves (degree = 1): present-day species



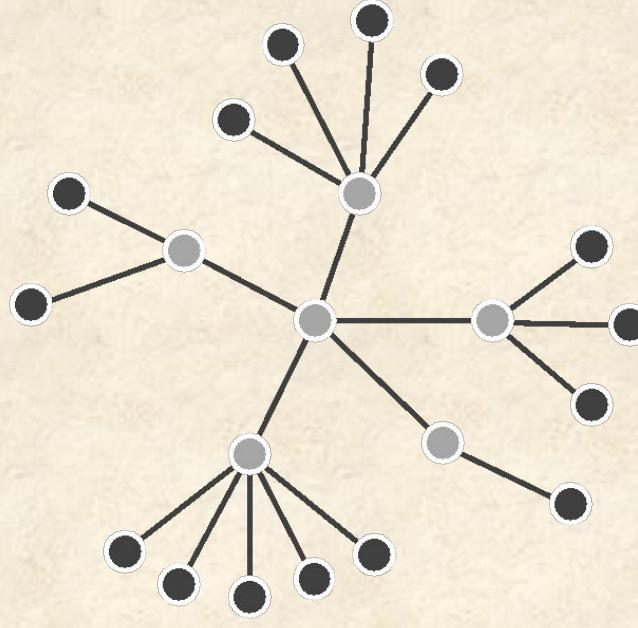
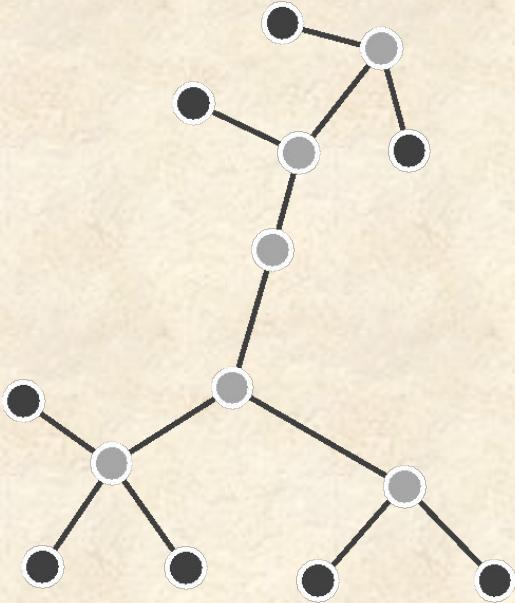
Trees

Tree: Connected graph containing no cycles.

Leaves (degree = 1): present-day species

Internal nodes (degree ≥ 1): ancestral species

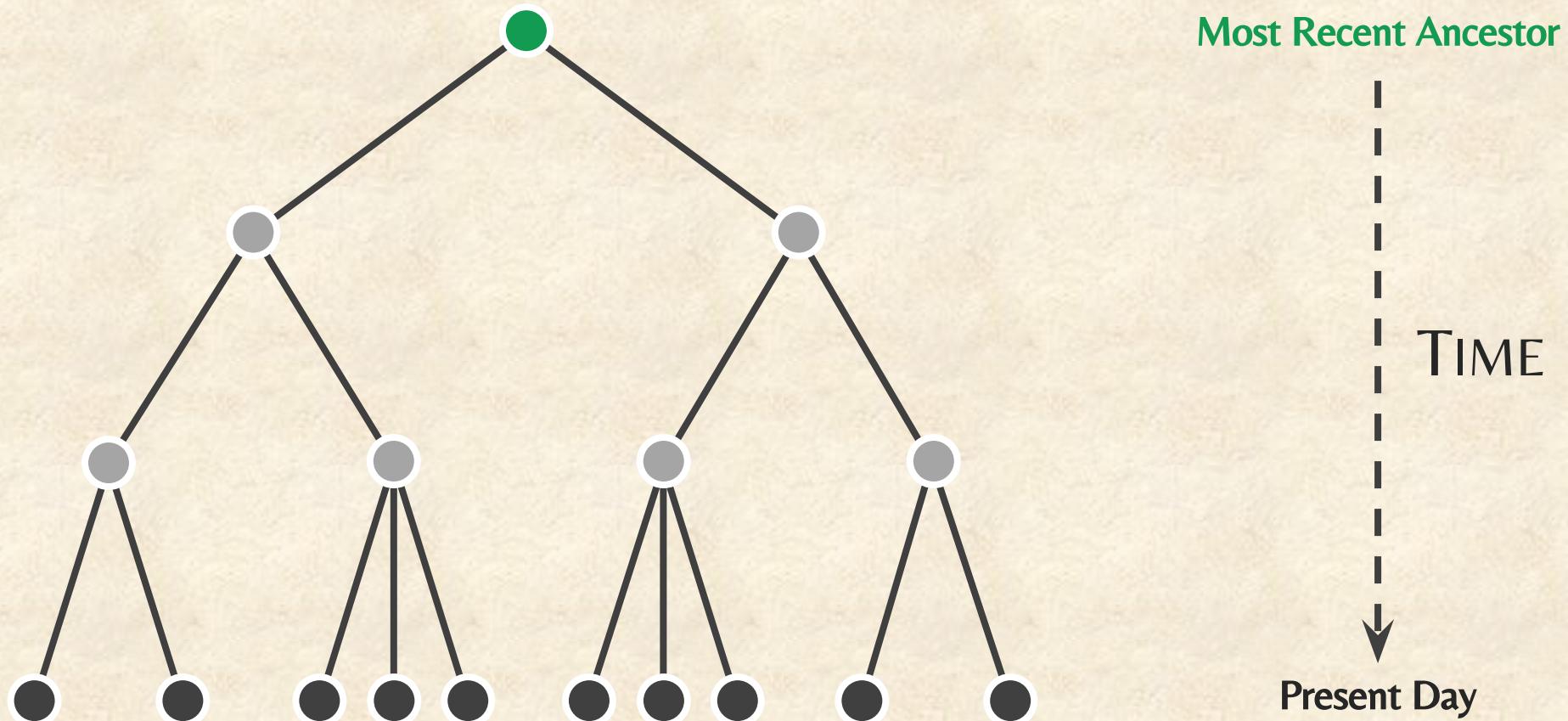
Trees



Exercise Break: Prove the following two statements.

1. Every tree with at least two nodes has at least two leaves.
2. Every tree with n nodes has exactly $n - 1$ edges.

Trees



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

Distance-Based Phylogeny

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

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

This is not a
computational
problem!

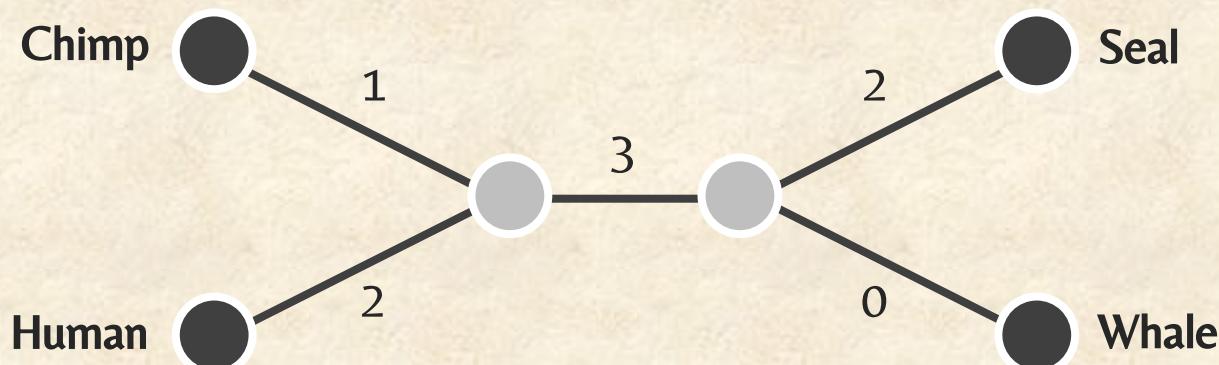


Fitting a Tree to a Matrix

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

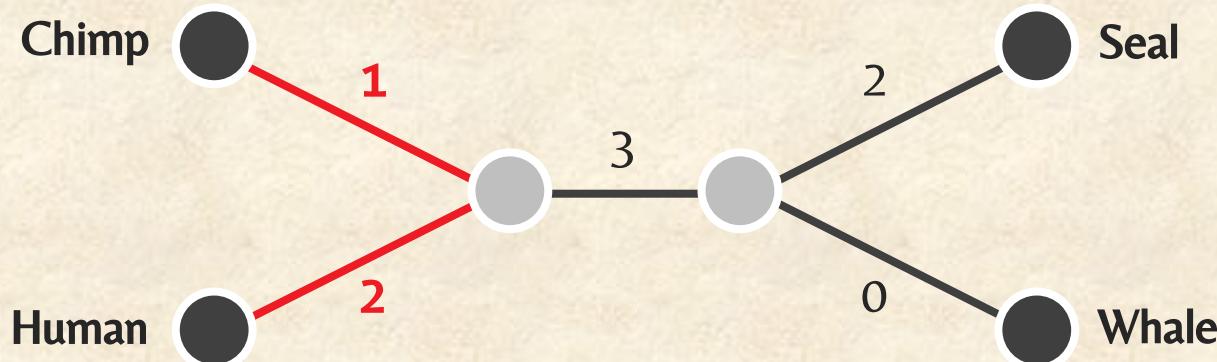
Fitting a Tree to a Matrix

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



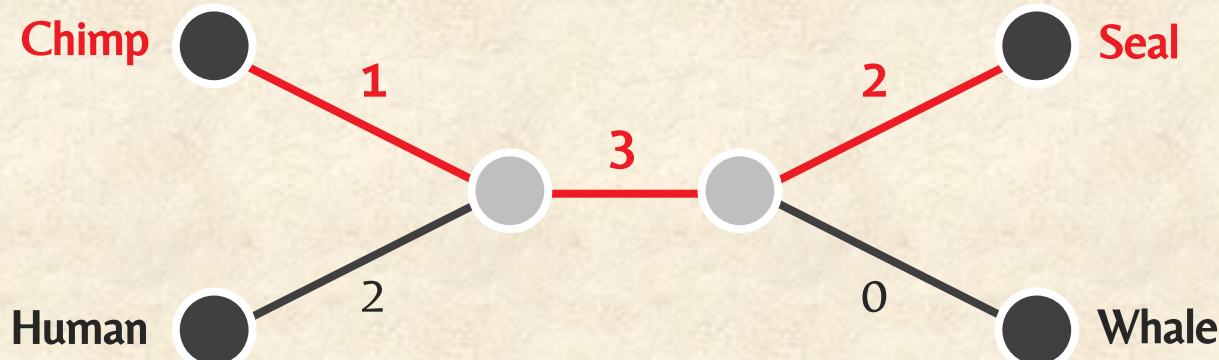
Fitting a Tree to a Matrix

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



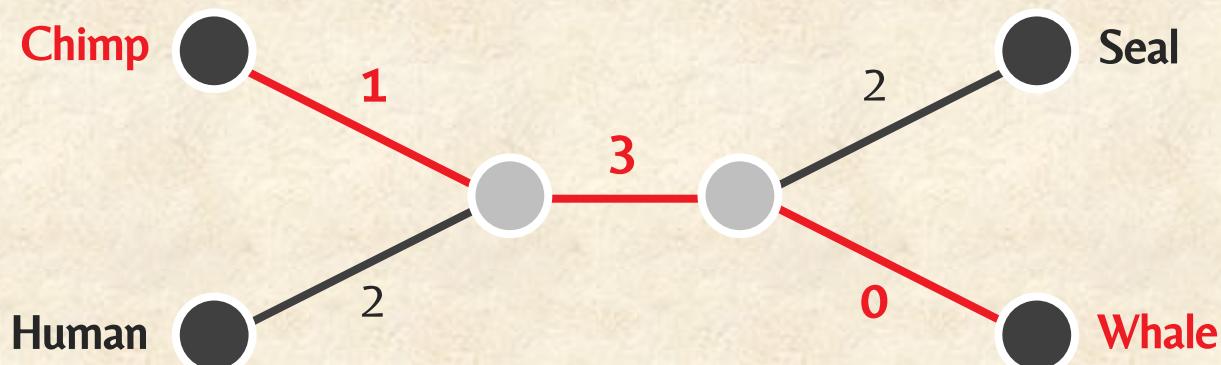
Fitting a Tree to a Matrix

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



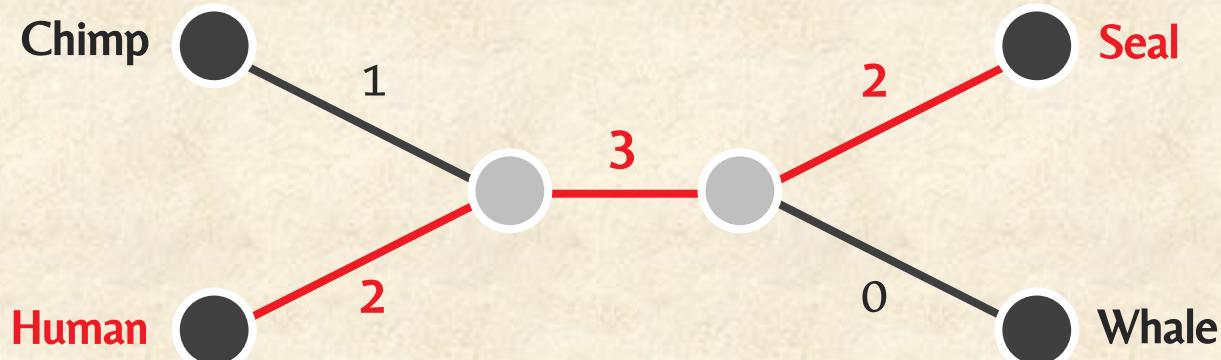
Fitting a Tree to a Matrix

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



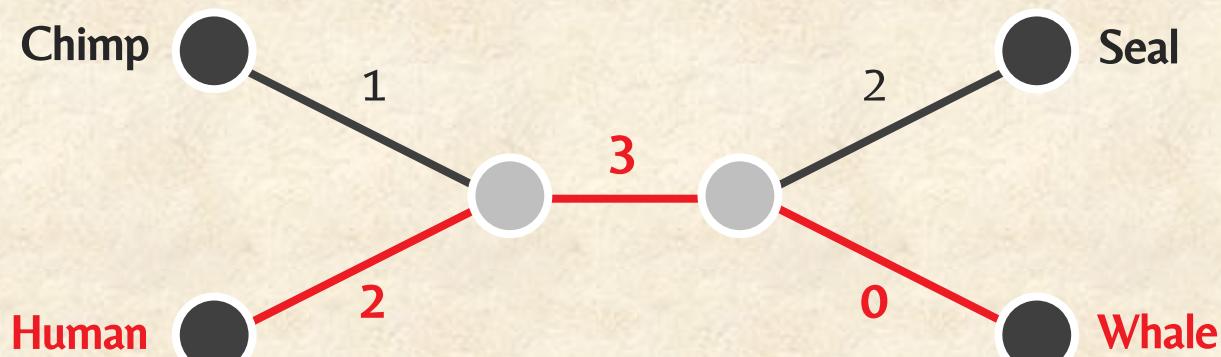
Fitting a Tree to a Matrix

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



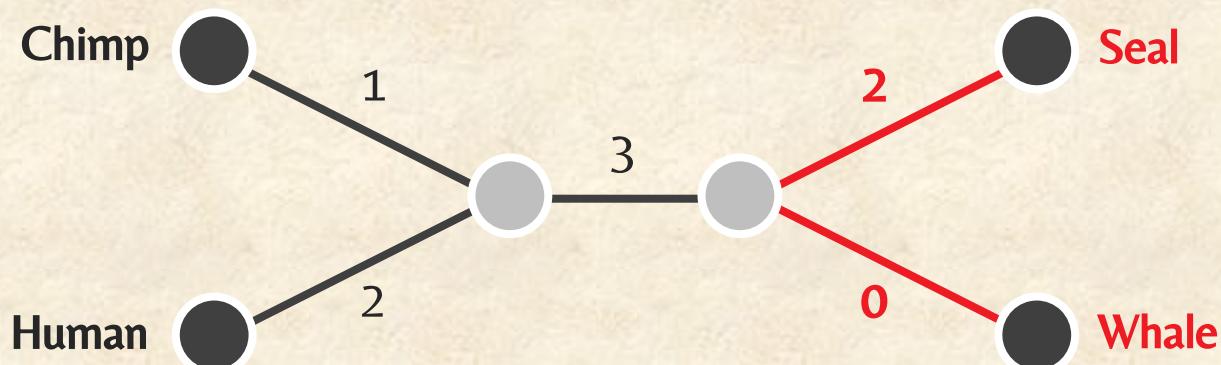
Fitting a Tree to a Matrix

| | 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 |



Distance Between Leaves

Distance Between Leaves Problem: *Compute the distances between leaves in a weighted tree.*

- **Input:** A weighted tree with n leaves.
- **Output:** An $n \times n$ matrix $(d_{i,j})$, where $d_{i,j}$ is the length of the path between leaves i and j .

Code Challenge: Solve this problem.

Return to Distance-Based Phylogeny

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

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

STOP and Think: Now is this problem well-defined?

Return to Distance-Based Phylogeny

Exercise Break: Try fitting a tree to the following matrix.

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 3 | 4 | 3 |
| <i>j</i> | 3 | 0 | 4 | 5 |
| <i>k</i> | 4 | 4 | 0 | 2 |
| <i>l</i> | 3 | 5 | 2 | 0 |

No Tree Fits a Matrix

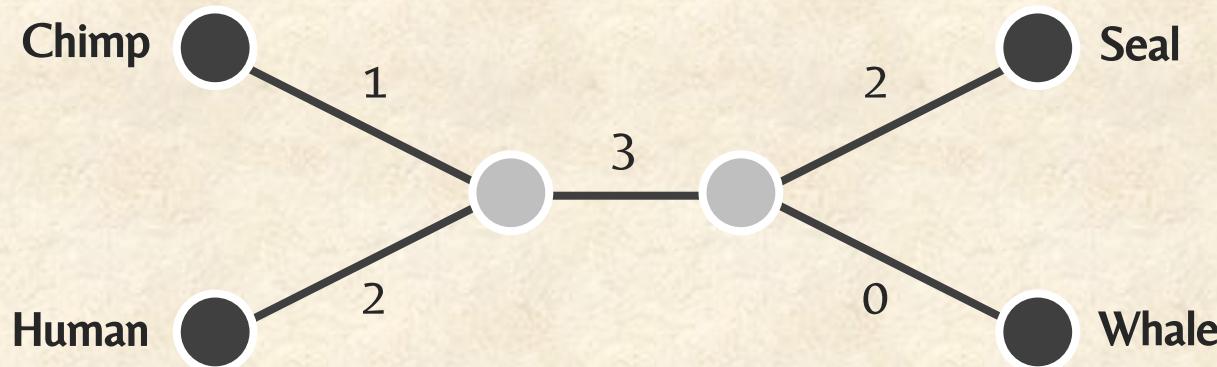
Exercise Break: Try fitting a tree to the following matrix.

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

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

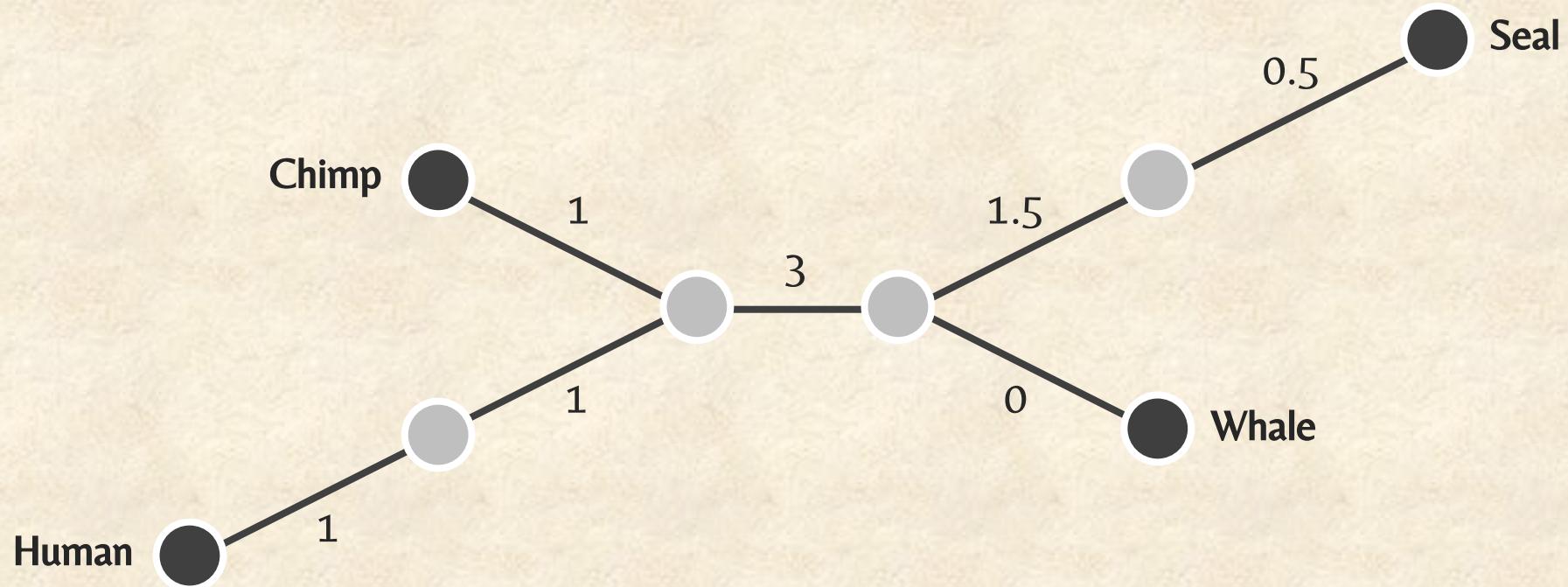
More Than One Tree Fits a Matrix

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

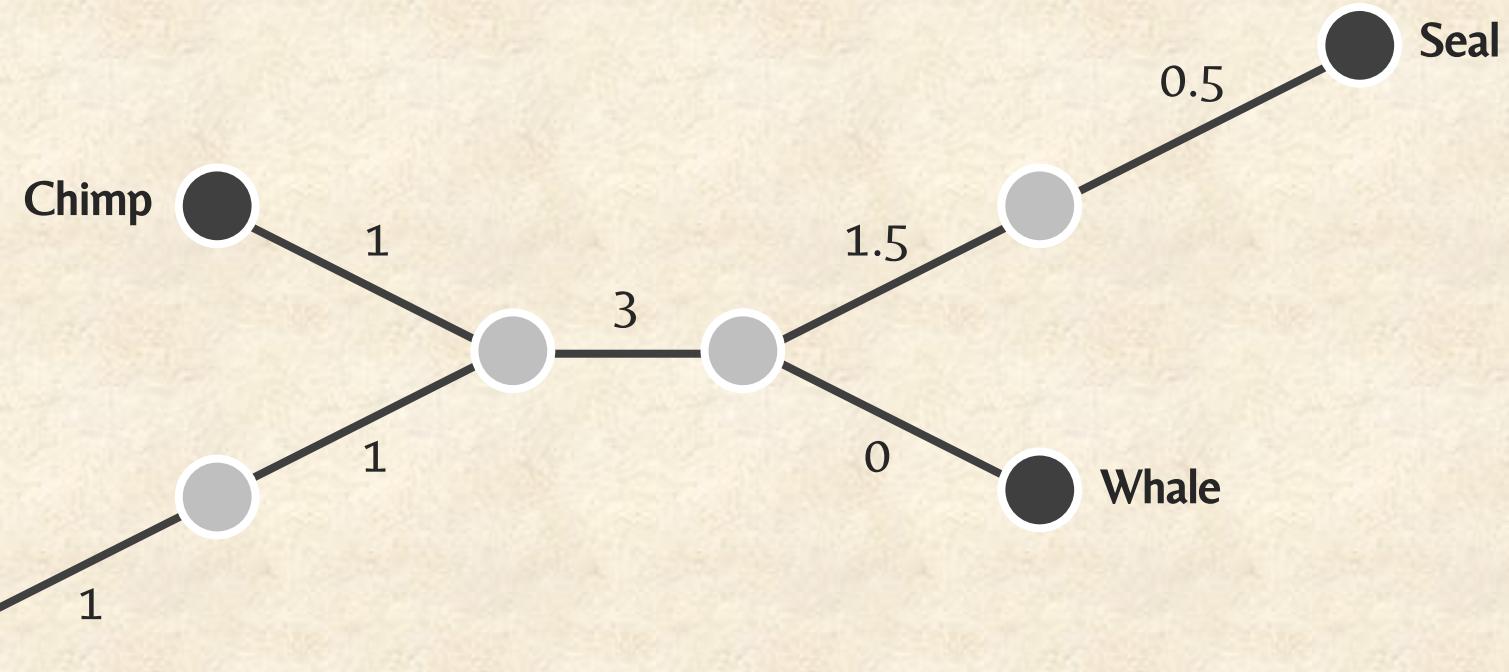
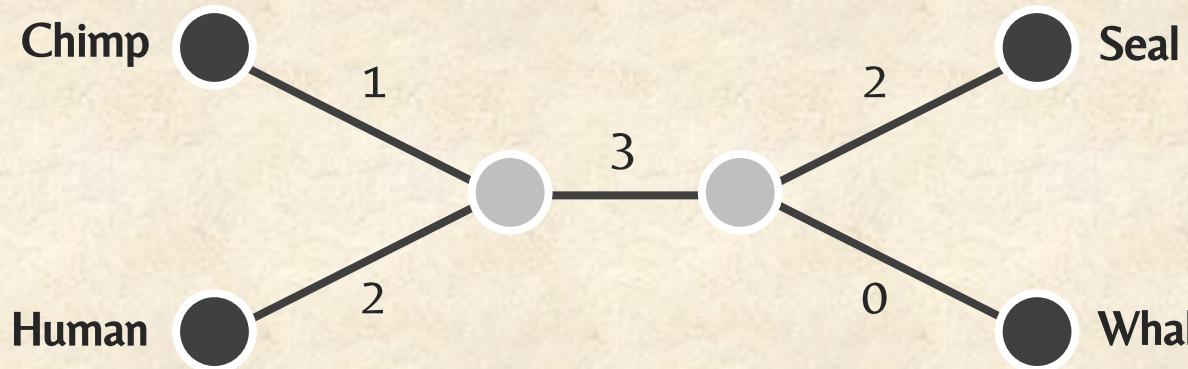


More Than One Tree Fits a Matrix

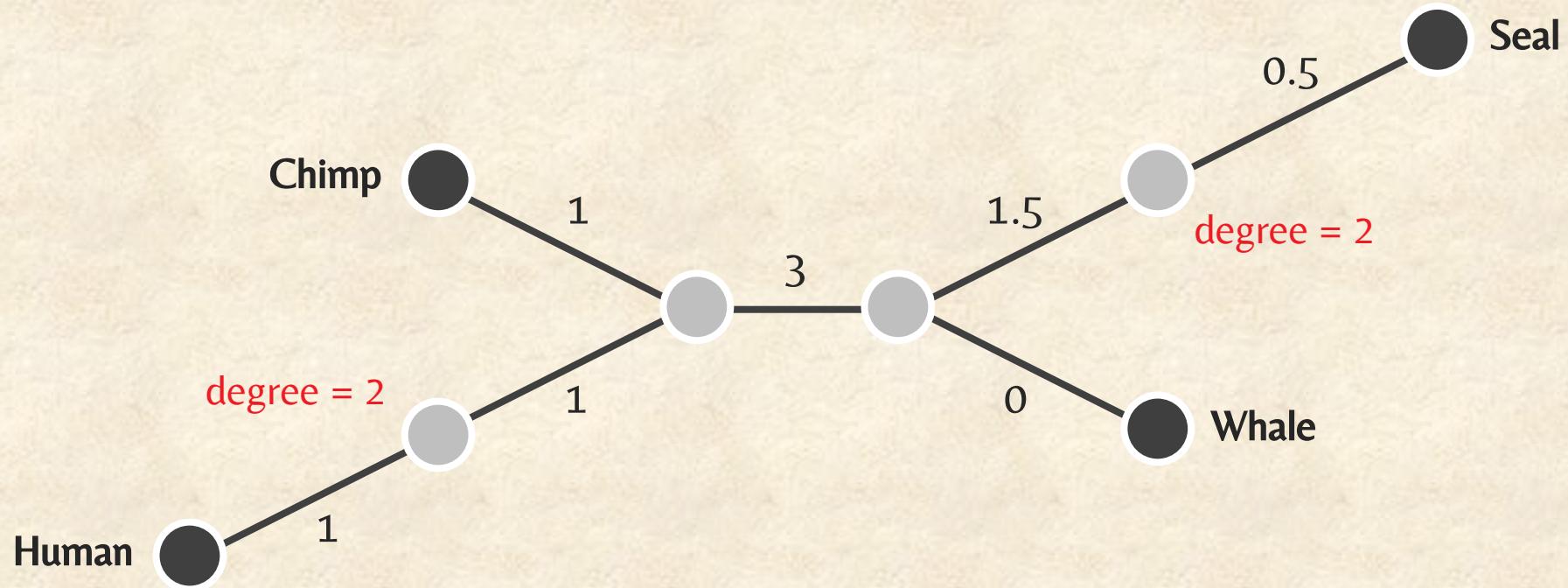
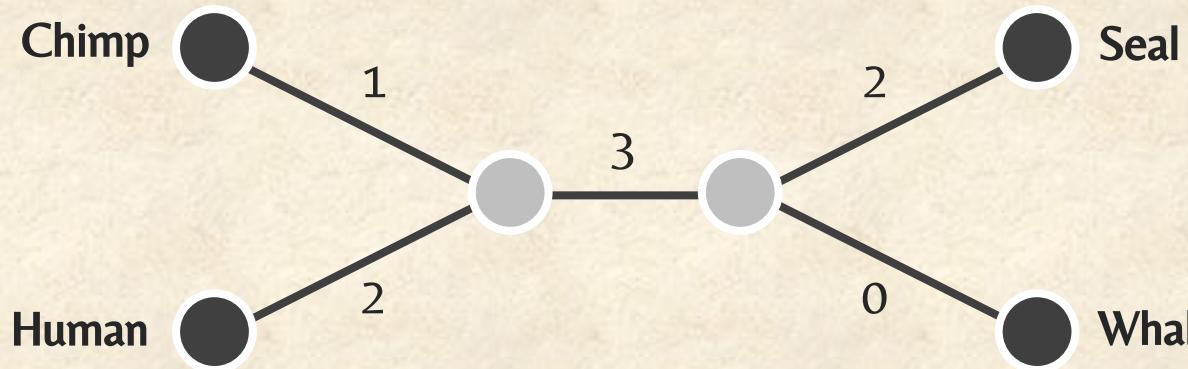
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



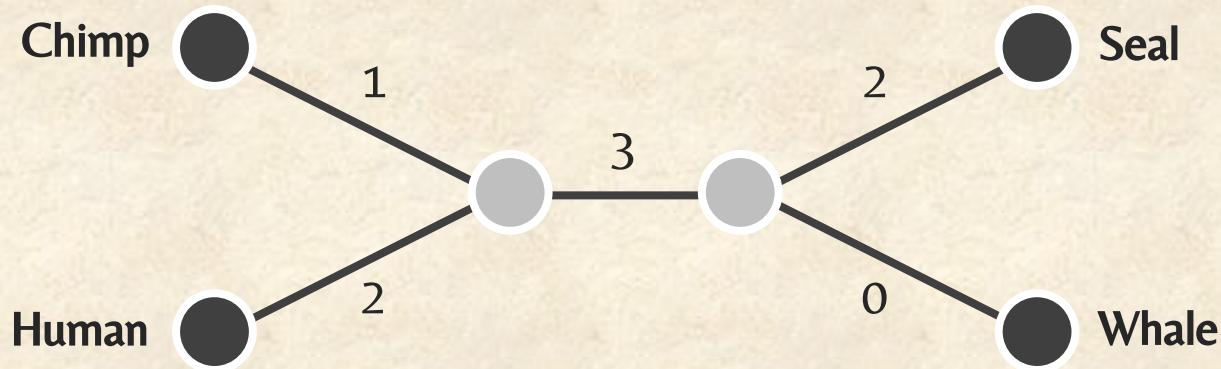
Which Tree is “Better”?



Which Tree is “Better”?

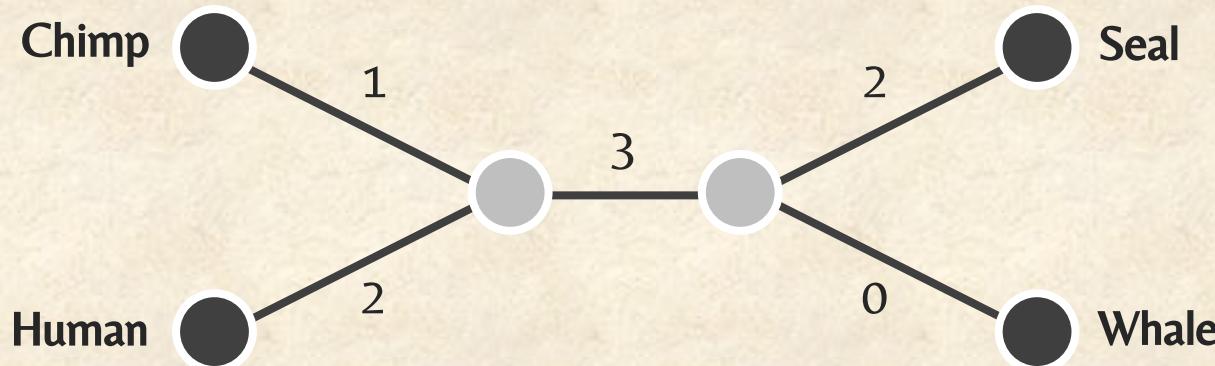


Which Tree is “Better”?



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.*

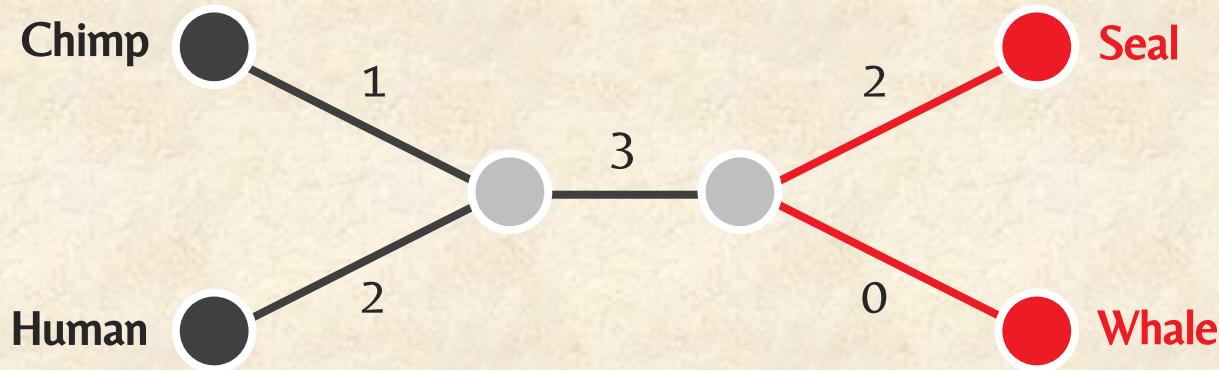
- **Input:** A distance matrix.
- **Output:** The simple tree fitting this distance matrix (if this matrix is additive).

Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- **Toward an Algorithm for Distance-Based Phylogeny Construction**
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

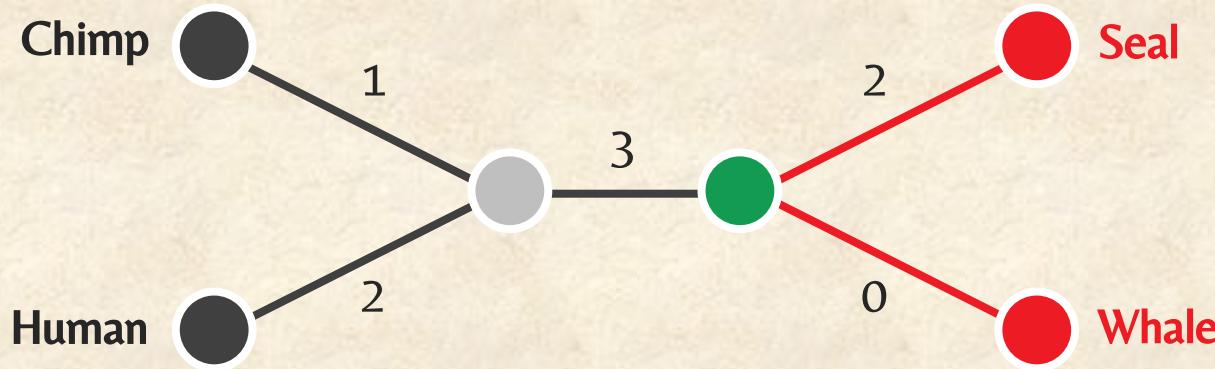
An Idea for Distance-Based Phylogeny

| | 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

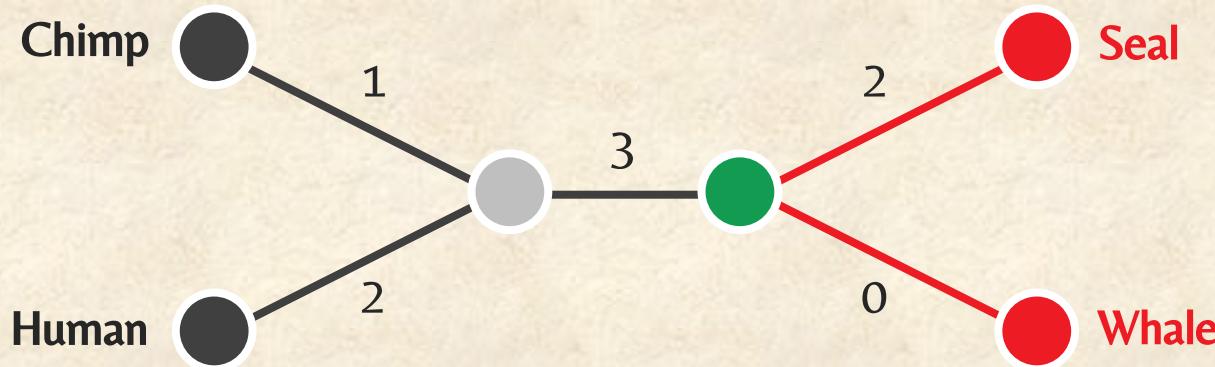
Seal and whale are **neighbors** (meaning they share the same **parent**).



An Idea for Distance-Based Phylogeny

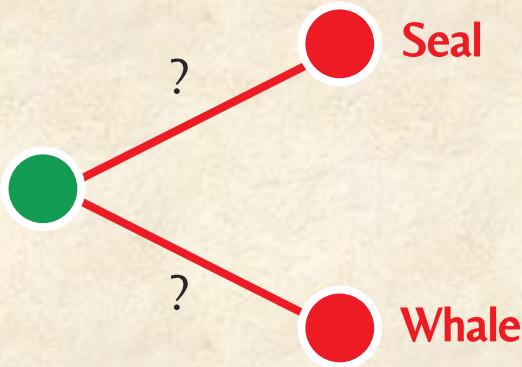
Seal and whale are **neighbors** (meaning they share the same **parent**).

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



An Idea for Distance-Based Phylogeny

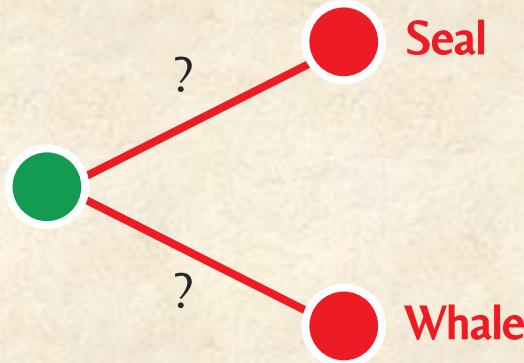
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



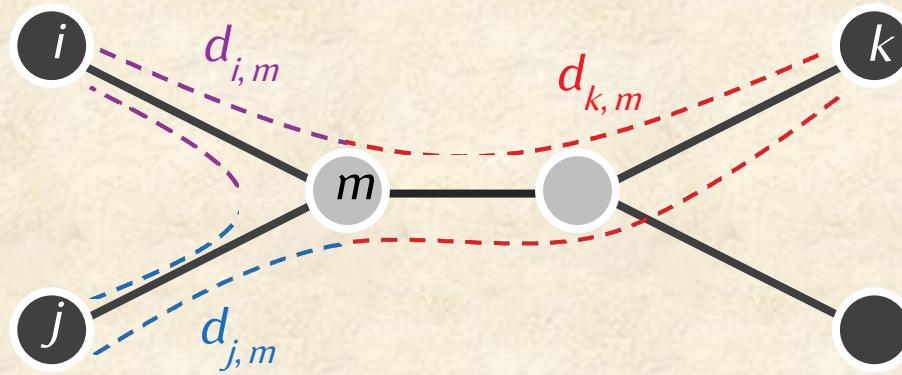
An Idea for Distance-Based Phylogeny

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

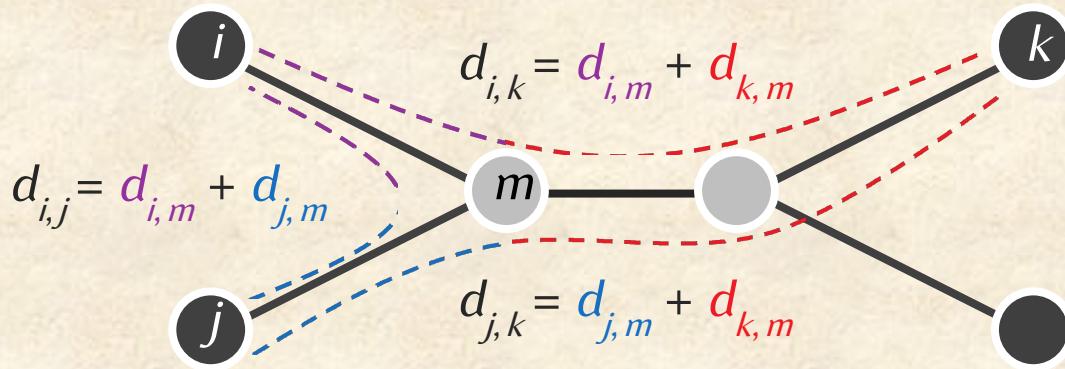
STOP and Think: How do we compute the unknown distances?



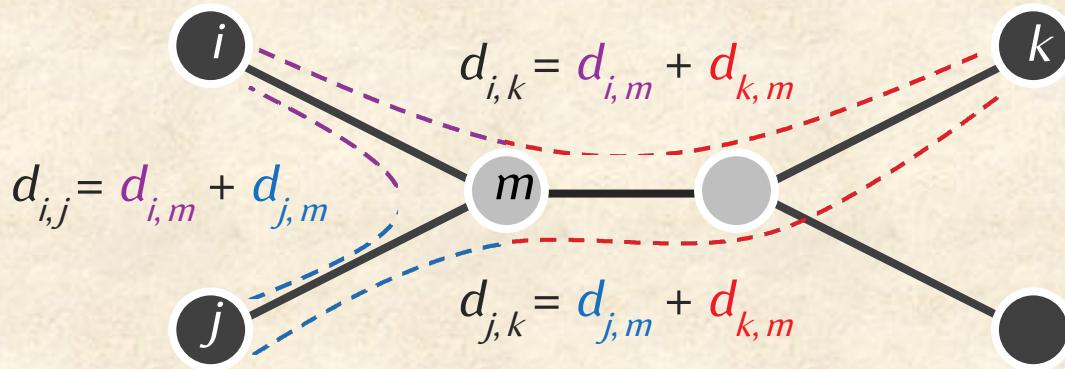
Toward a Recursive Algorithm



Toward a Recursive Algorithm

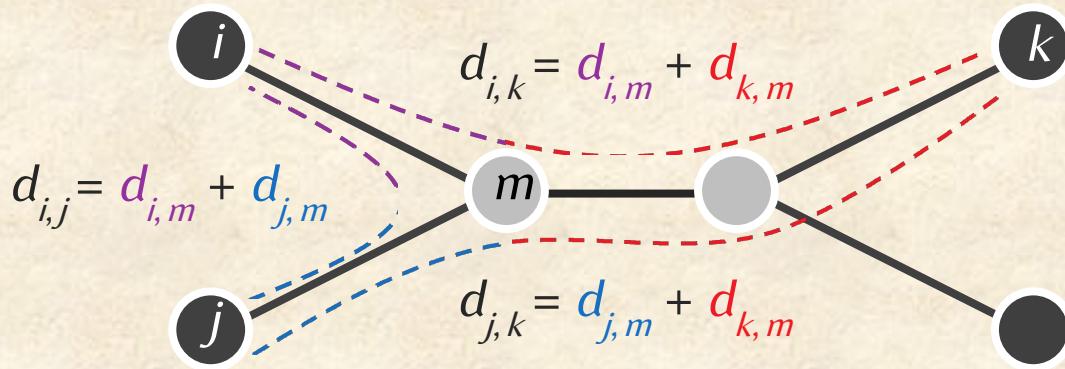


Toward a Recursive Algorithm



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

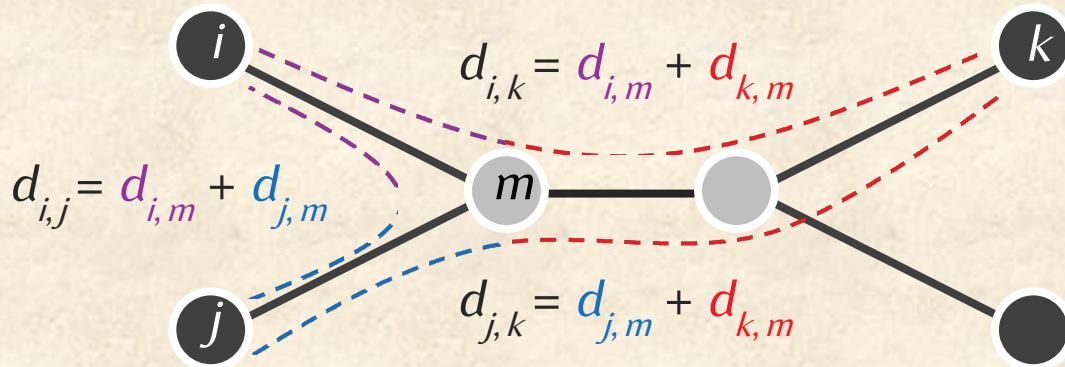
Toward a Recursive Algorithm



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

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

Toward a Recursive Algorithm

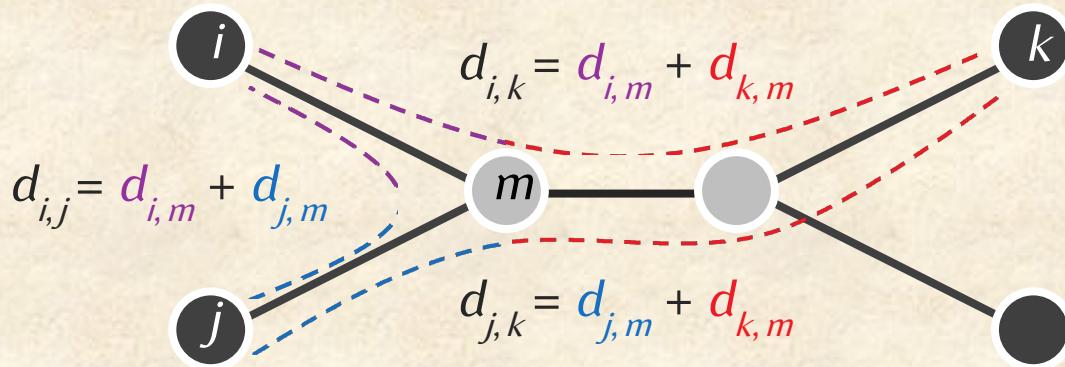


$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

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

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

Toward a Recursive Algorithm



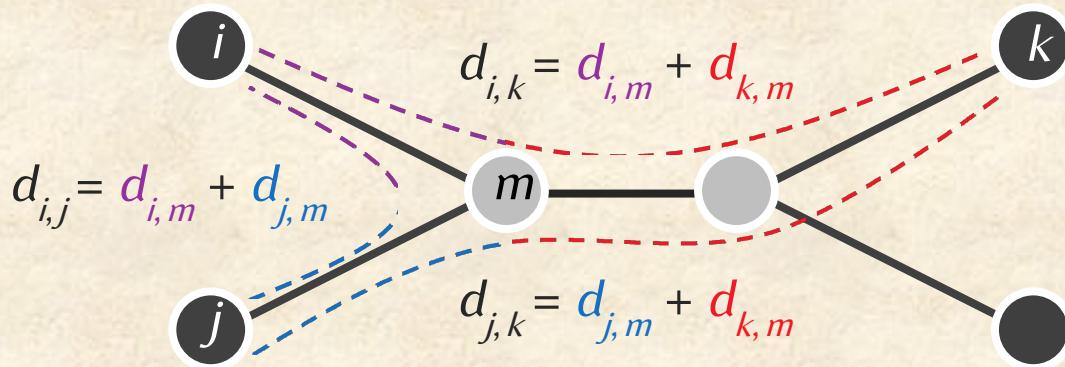
$$d_{k,m} = [(d_{i,m} + d_{k,m}) + (d_{j,m} + d_{k,m}) - (d_{i,m} + d_{j,m})] / 2$$

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

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

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

Toward a Recursive Algorithm



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

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

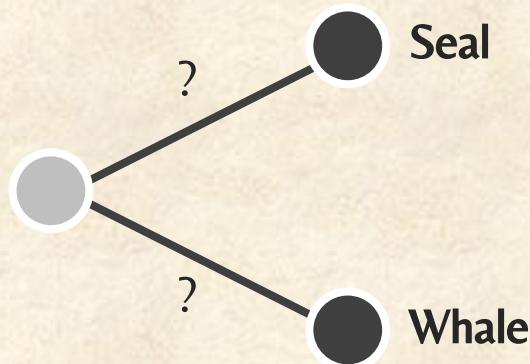
$$d_{k,m} = (D_{i,k} + D_{j,k} - D_{i,j}) / 2$$

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

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

An Idea for Distance-Based Phylogeny

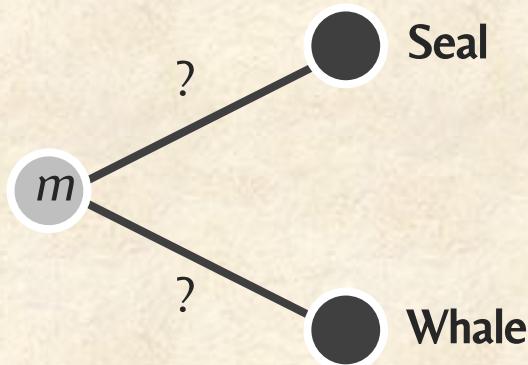
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

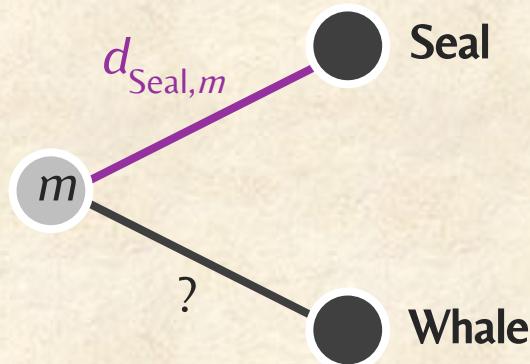
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

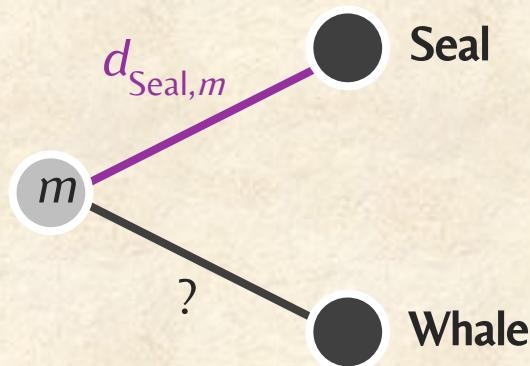
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

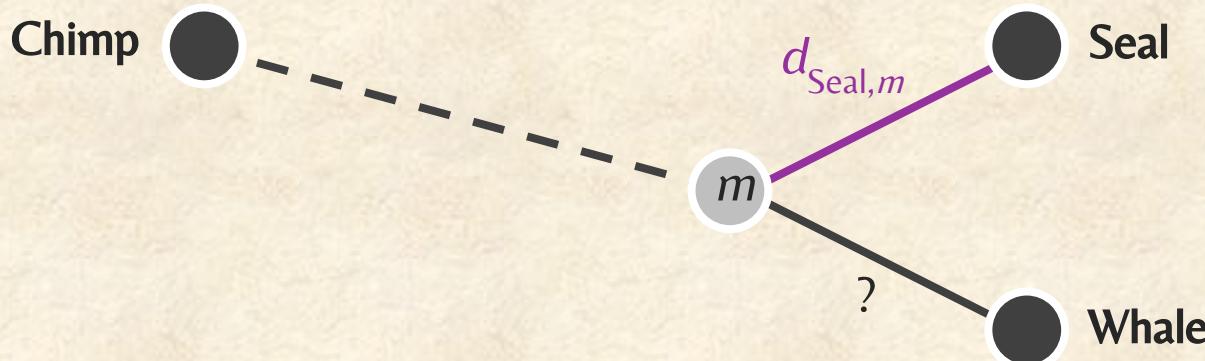
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

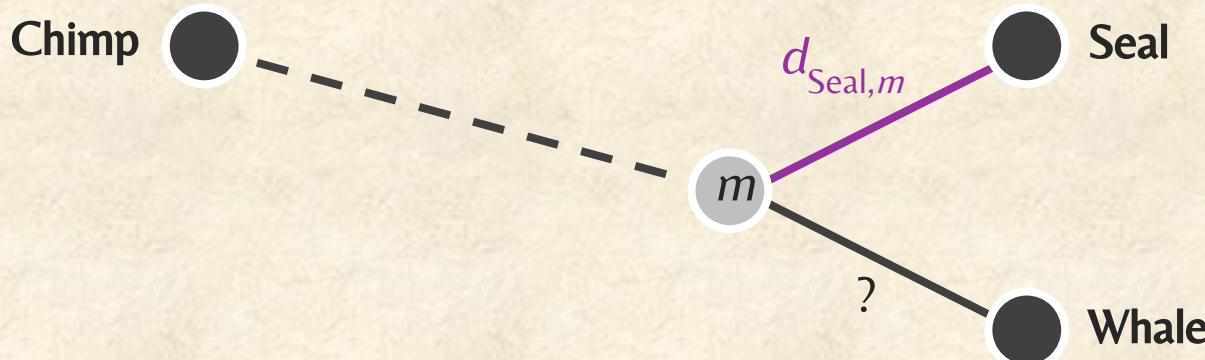
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

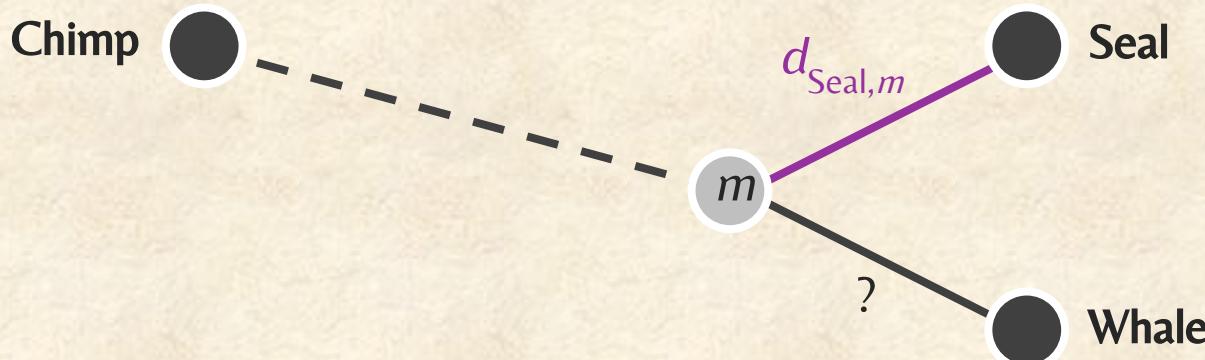
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

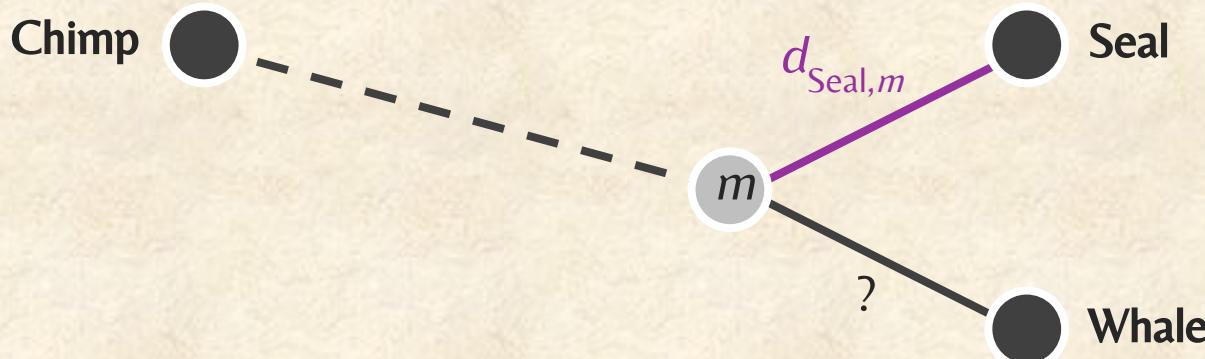
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

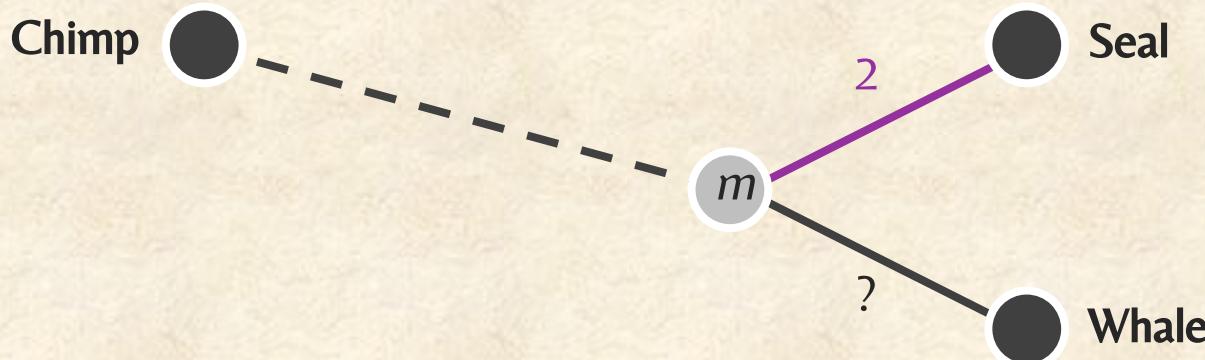
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

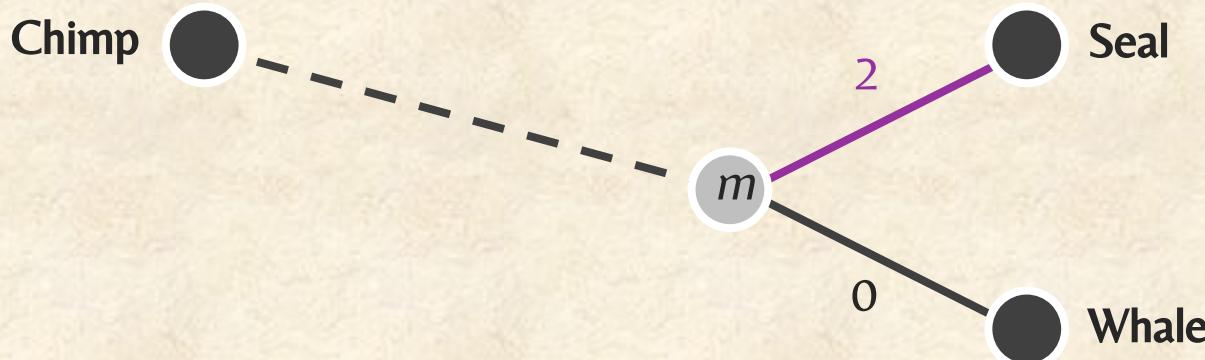
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

An Idea for Distance-Based Phylogeny

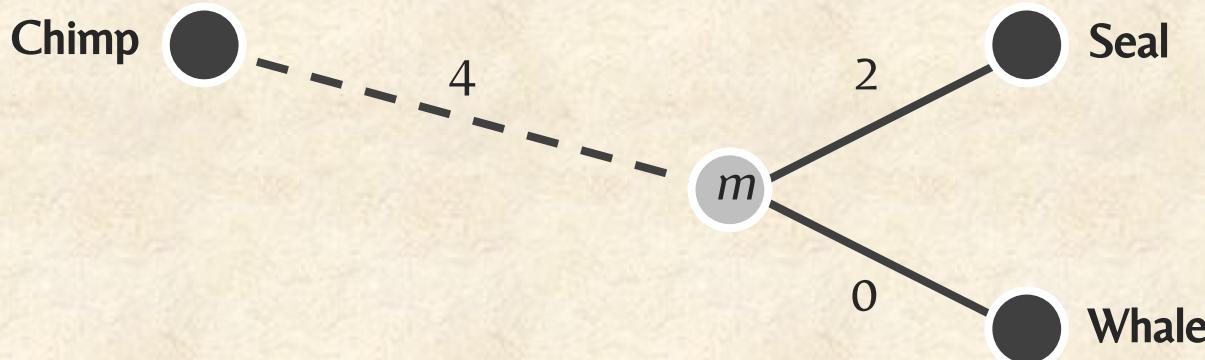
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



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

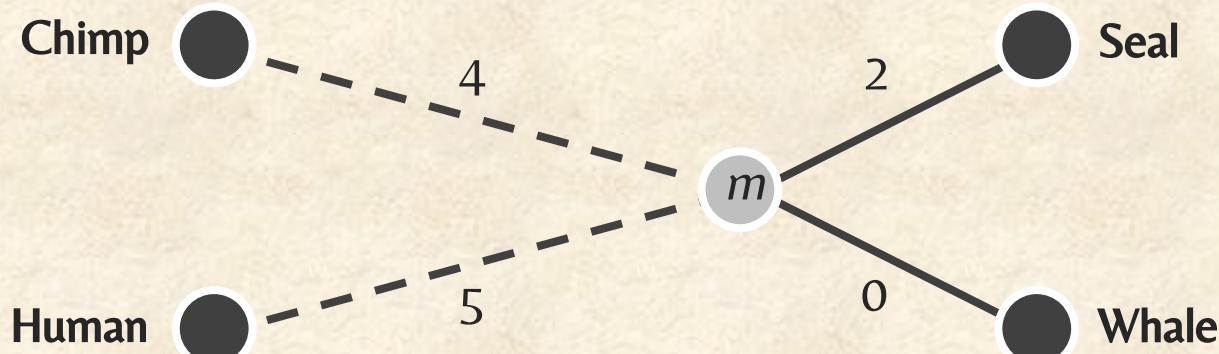
An Idea for Distance-Based Phylogeny

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



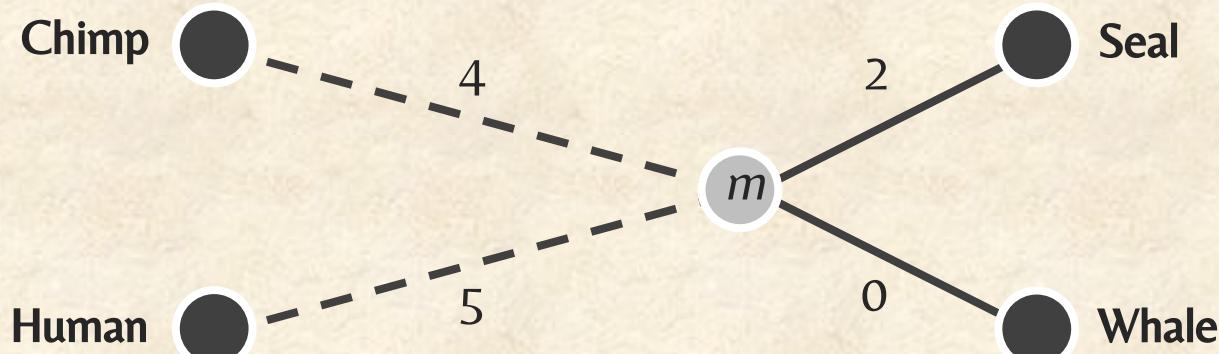
An Idea for Distance-Based Phylogeny

| | 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 |



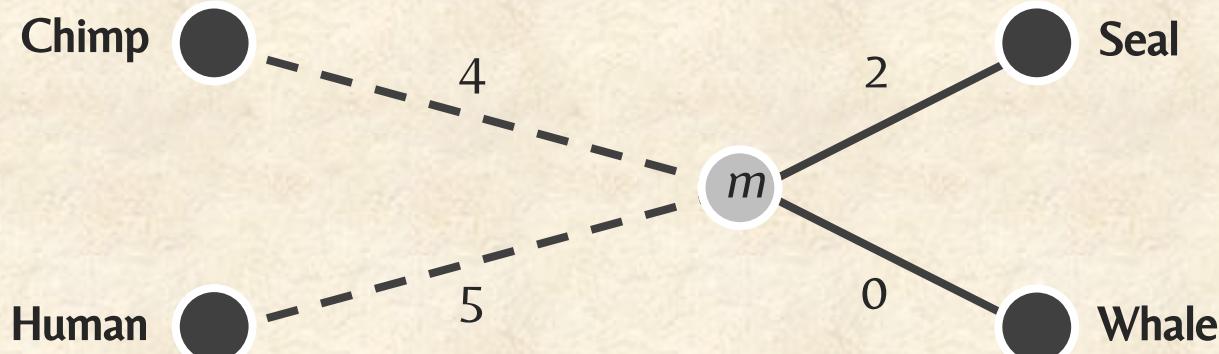
An Idea for Distance-Based Phylogeny

| | 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 |



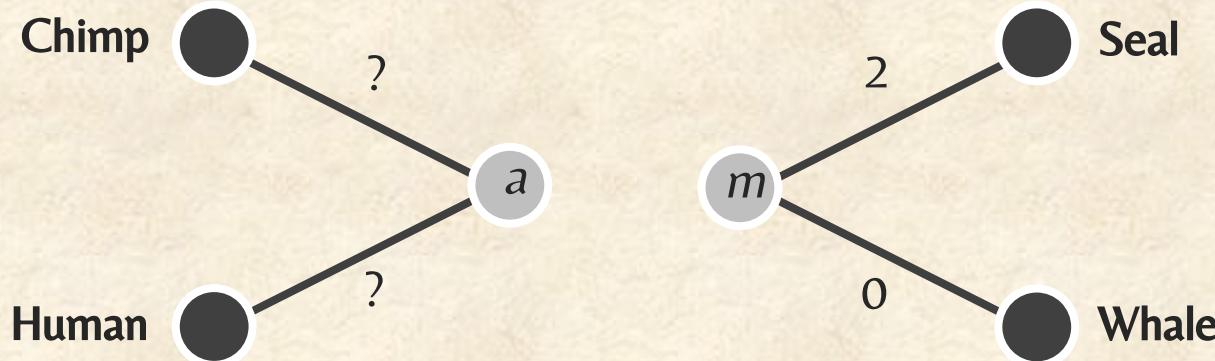
An Idea for Distance-Based Phylogeny

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



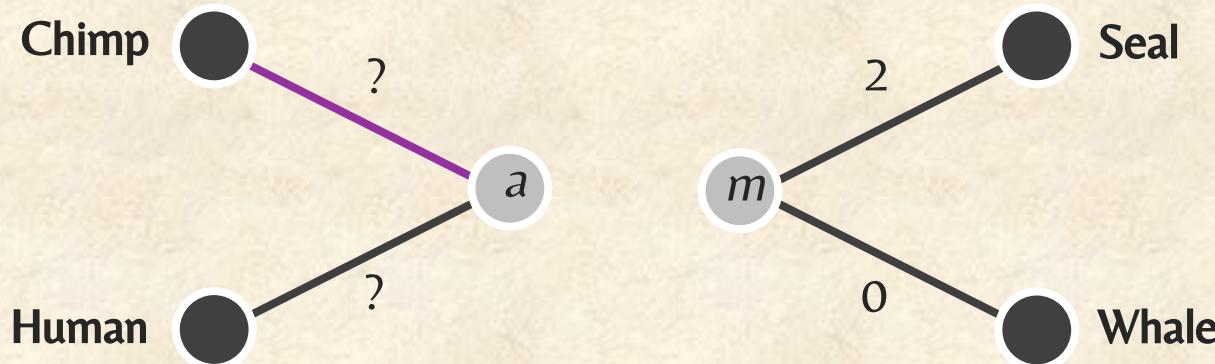
An Idea for Distance-Based Phylogeny

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



An Idea for Distance-Based Phylogeny

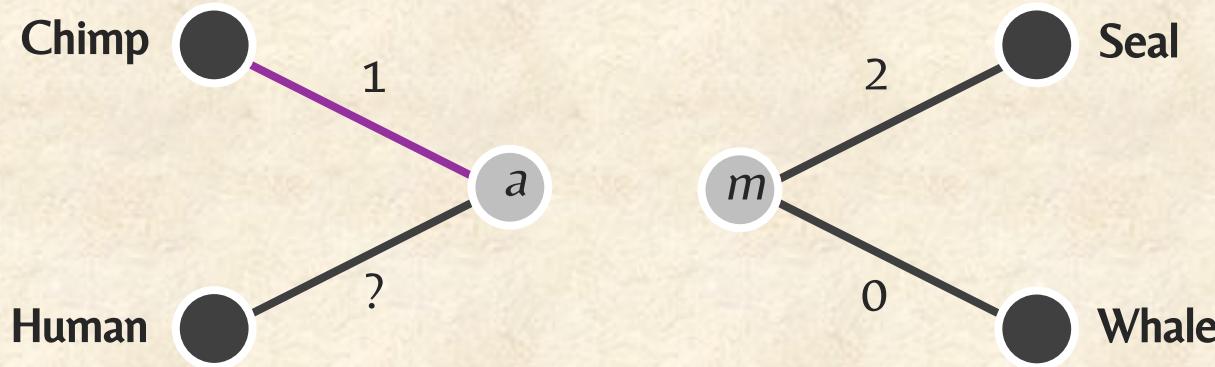
| | Chimp | Human | <i>m</i> |
|----------|-------|-------|----------|
| Chimp | 0 | 3 | 4 |
| Human | 3 | 0 | 5 |
| <i>m</i> | 4 | 5 | 0 |



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

An Idea for Distance-Based Phylogeny

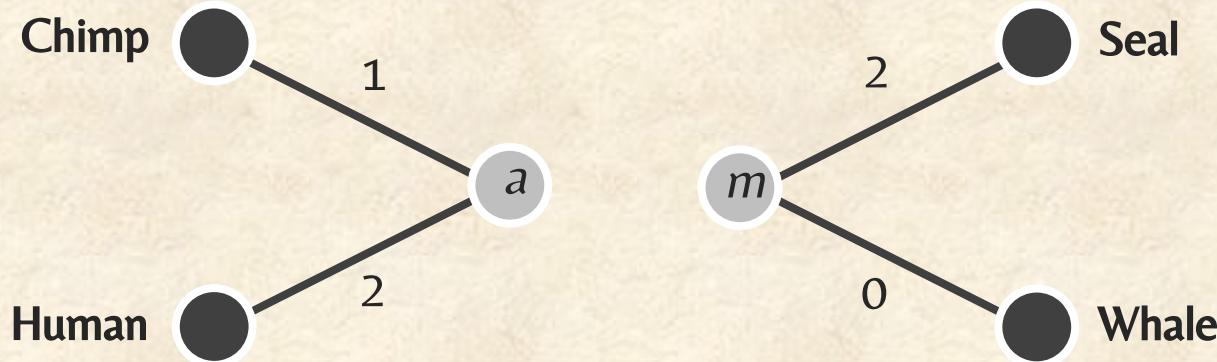
| | Chimp | Human | <i>m</i> |
|----------|-------|-------|----------|
| Chimp | 0 | 3 | 4 |
| Human | 3 | 0 | 5 |
| <i>m</i> | 4 | 5 | 0 |



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

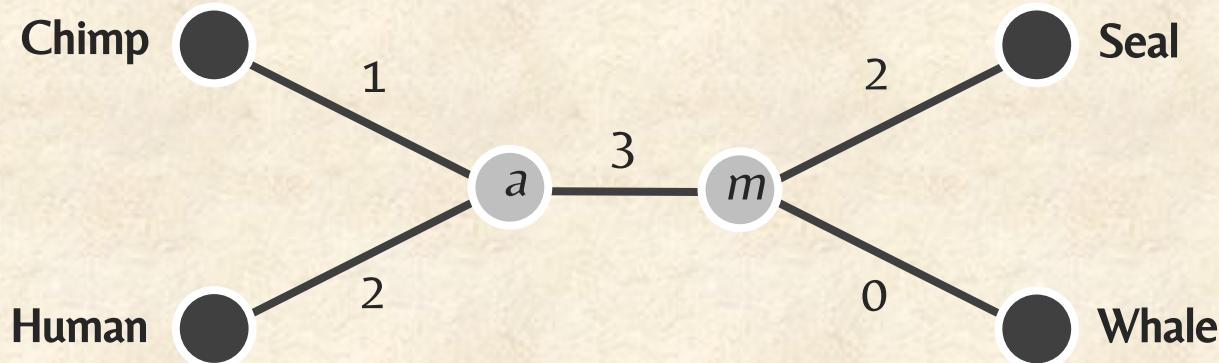
An Idea for Distance-Based Phylogeny

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



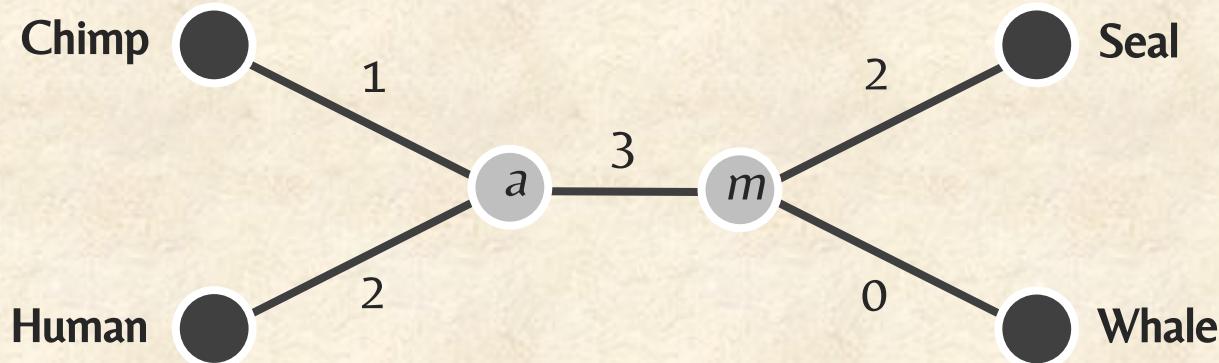
An Idea for Distance-Based Phylogeny

| | Chimp | Human | <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

Exercise Break: Apply this recursive approach to the distance matrix below.

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 13 | 21 | 22 |
| <i>j</i> | 13 | 0 | 12 | 13 |
| <i>k</i> | 21 | 12 | 0 | 13 |
| <i>l</i> | 22 | 13 | 13 | 0 |

Outline

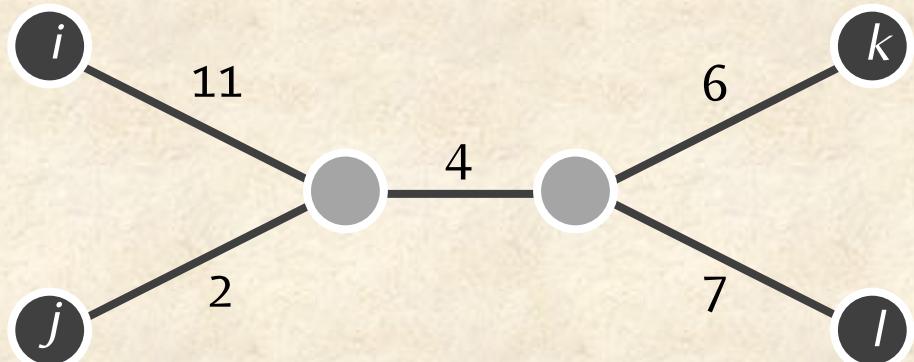
- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- **Additive Phylogeny**
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

What Was Wrong With Our Algorithm?

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

What Was Wrong With Our Algorithm?

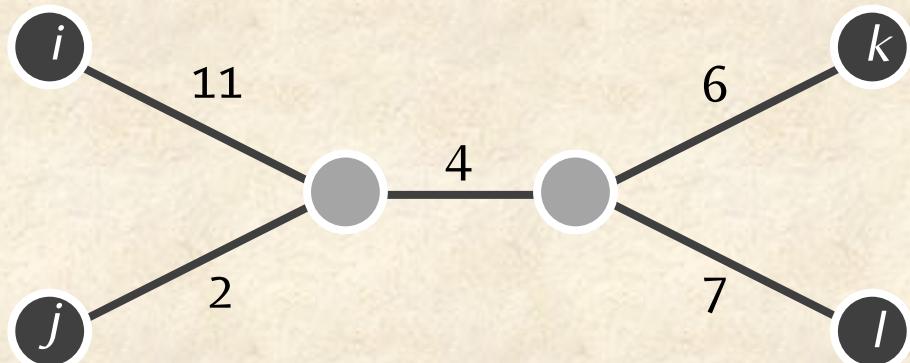
| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |



What Was Wrong With Our Algorithm?

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

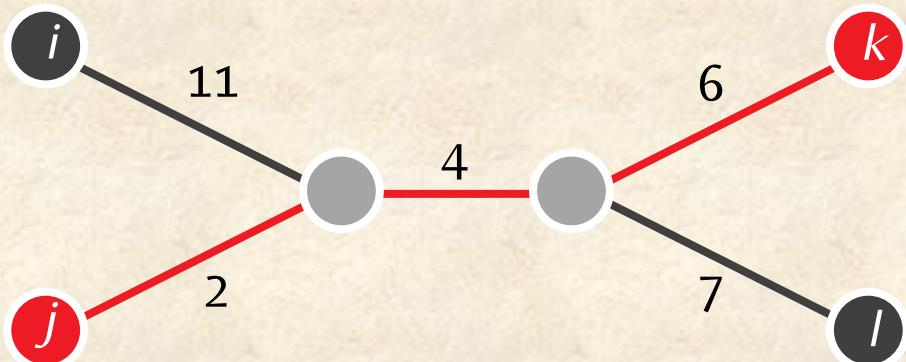
minimum element is $D_{j,k}$



What Was Wrong With Our Algorithm?

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

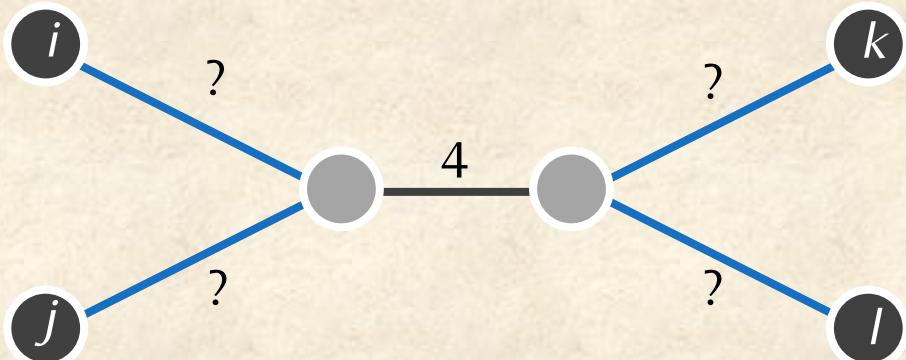
minimum element is $D_{j,k}$



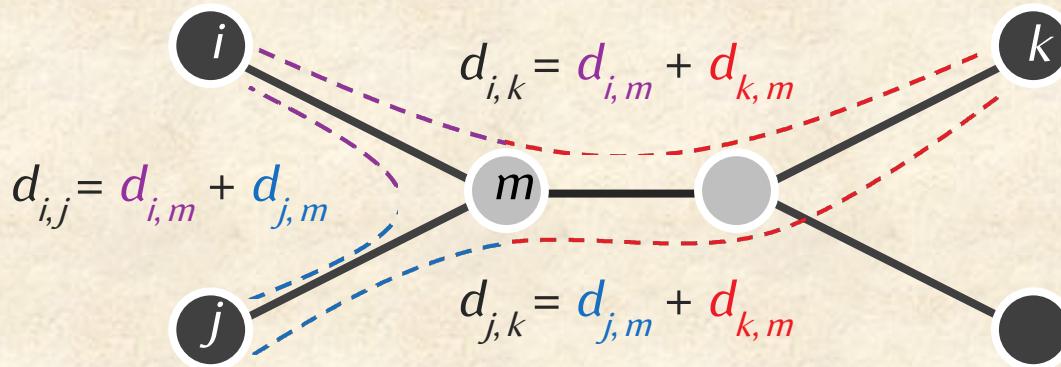
j and k are
not neighbors!

From Neighbors to Limbs

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



From Neighbors to Limbs



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

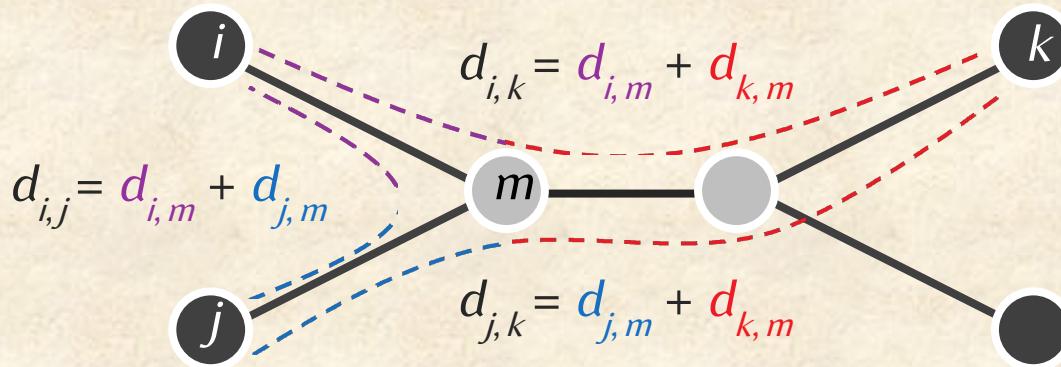
$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

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

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

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

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

Limb Length Problem: *Compute the length of a limb in the simple tree fitting an additive distance matrix.*

- **Input:** An additive distance matrix D and an integer j .
- **Output:** The length of the limb connecting leaf j to its parent, $\text{LimbLength}(j)$.

Code Challenge: Solve the Limb Length Problem.

Computing Limb Lengths

Limb Length Theorem: $\text{LimbLength}(\text{chimp})$ is equal to the minimum value of $(D_{\text{chimp},k} + D_{\text{chimp},j} - D_{j,k})/2$ 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: $\text{LimbLength}(\text{chimp})$ is equal to the minimum value of $(D_{\text{chimp},k} + D_{\text{chimp},j} - D_{j,k})/2$ 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 |

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

Computing Limb Lengths

Limb Length Theorem: $\text{LimbLength}(\text{chimp})$ is equal to the minimum value of $(D_{\text{chimp},k} + D_{\text{chimp},j} - D_{j,k})/2$ 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 |

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

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

Computing Limb Lengths

Limb Length Theorem: $\text{LimbLength}(\text{chimp})$ is equal to the minimum value of $(D_{\text{chimp},k} + D_{\text{chimp},j} - D_{j,k})/2$ 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 |

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

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

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

Computing Limb Lengths

Limb Length Theorem: $\text{LimbLength}(\text{chimp})$ is equal to the minimum value of $(D_{\text{chimp},k} + D_{\text{chimp},j} - D_{j,k})/2$ 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 |

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

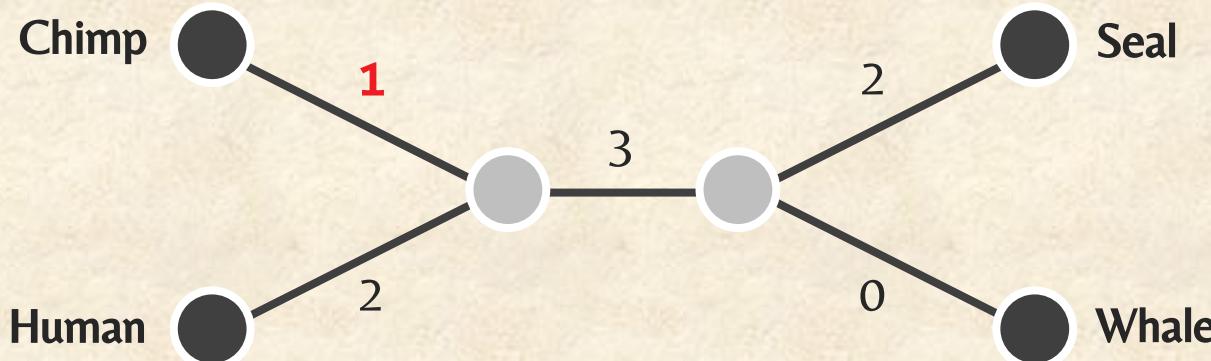
$$(D_{\text{human, chimp}} + D_{\text{chimp, whale}} - D_{\text{human, whale}}) / 2 = (3 + 4 - 5) / 2 = \mathbf{1}$$

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

Computing Limb Lengths

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

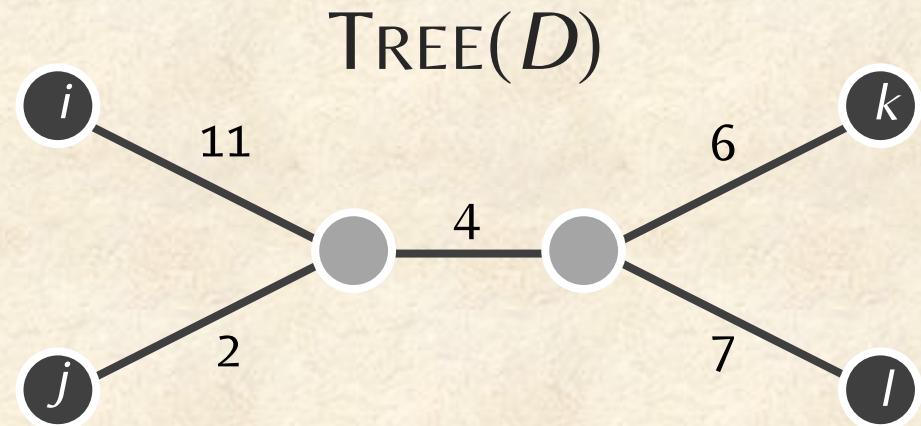
| | Chimp | Human | Seal | Whale |
|-------|-------|-------|------|-------|
| Chimp | 0 | 3 | 6 | 4 |
| Human | 3 | 0 | 7 | 5 |
| Seal | 6 | 7 | 0 | 2 |
| Whale | 4 | 5 | 2 | 0 |



Additive Phylogeny In Action

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |



Additive Phylogeny In Action

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

1. Pick an arbitrary leaf j .

Additive Phylogeny In Action

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

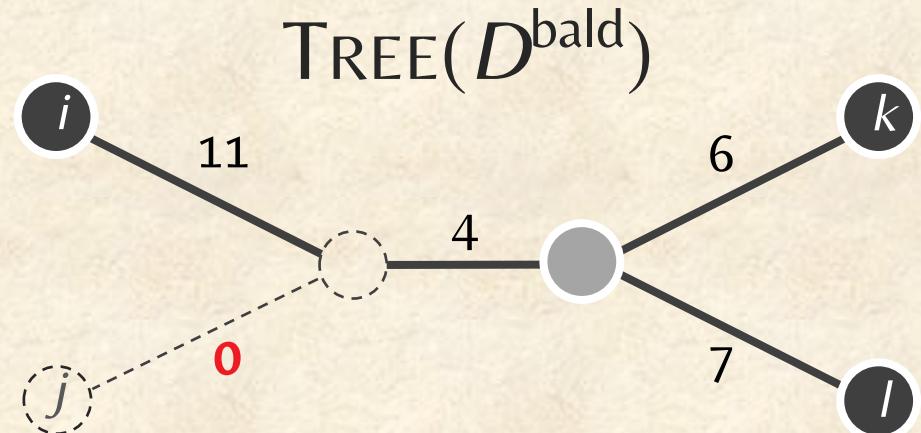
$$\text{LimbLength}(j) = 2$$

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

Additive Phylogeny In Action

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 11 | 21 | 22 |
| j | 11 | 0 | 10 | 11 |
| k | 21 | 10 | 0 | 13 |
| l | 22 | 11 | 13 | 0 |

D^{bald}



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

Additive Phylogeny In Action

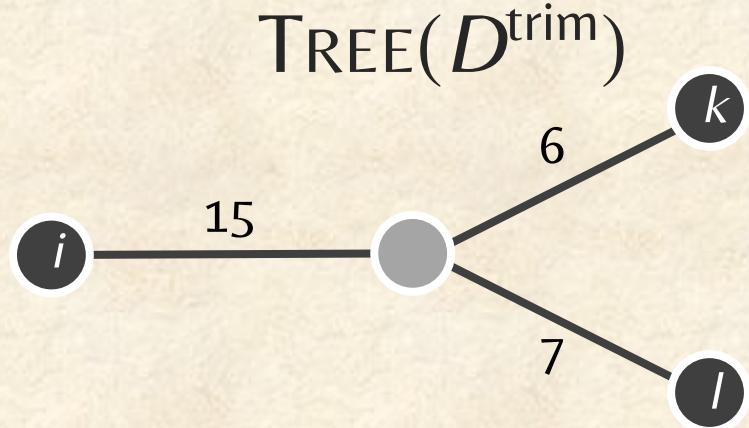
| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 11 | 21 | 22 |
| <i>j</i> | 11 | 0 | 10 | 11 |
| <i>k</i> | 21 | 10 | 0 | 13 |
| <i>l</i> | 22 | 11 | 13 | 0 |

D^{trim}

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

Additive Phylogeny In Action

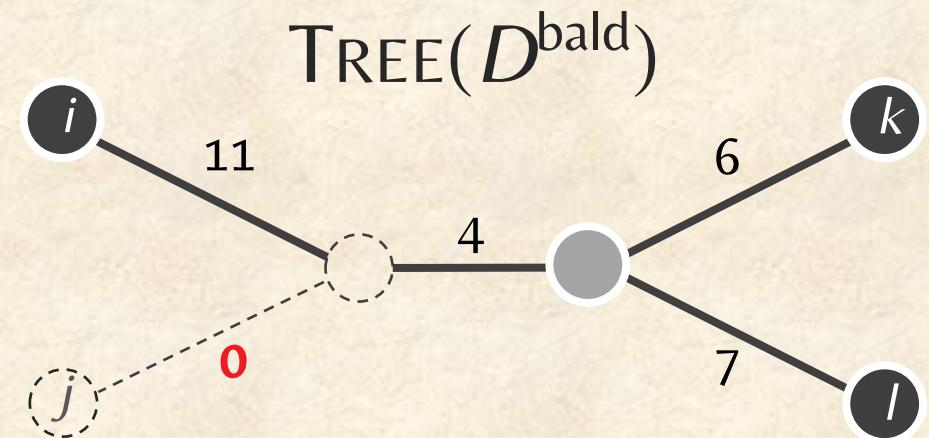
| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 11 | 21 | 22 |
| j | 11 | 0 | 10 | 11 |
| k | 21 | 10 | 0 | 13 |
| l | 22 | 11 | 13 | 0 |



5. Construct $\text{Tree}(D^{\text{trim}})$.

Additive Phylogeny In Action

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 11 | 21 | 22 |
| j | 11 | 0 | 10 | 11 |
| k | 21 | 10 | 0 | 13 |
| l | 22 | 11 | 13 | 0 |

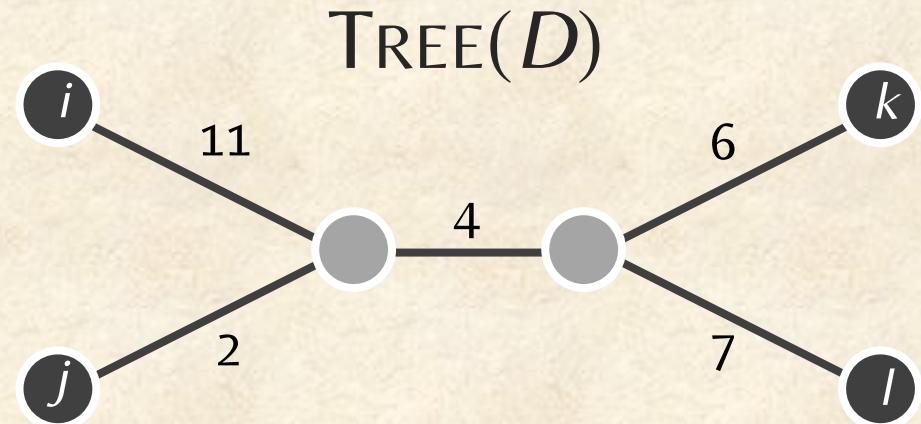


6. Identify the point in $\text{Tree}(D^{\text{trim}})$ where leaf j should be attached.

Additive Phylogeny In Action

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |



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

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

Additive Phylogeny

AdditivePhylogeny(D):

1. Pick an arbitrary leaf j .
2. Compute its limb length, $\text{LimbLength}(j)$.
3. Subtract $\text{LimbLength}(j)$ 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. Construct $\text{Tree}(D^{\text{trim}})$.
6. Identify the point in $\text{Tree}(D^{\text{trim}})$ where leaf j should be attached.
7. Attach j by an edge of length $\text{LimbLength}(j)$ in order to form $\text{Tree}(D)$.

Additive Phylogeny

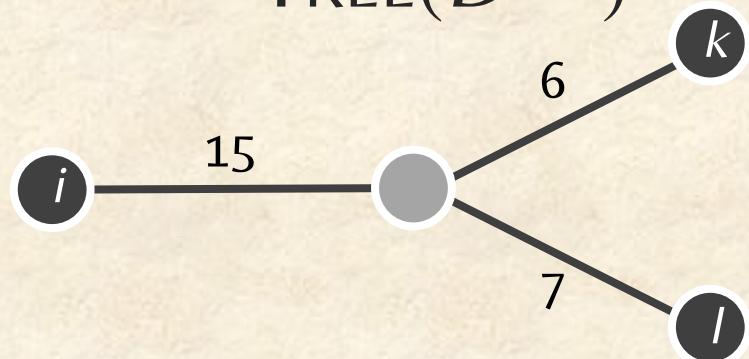
AdditivePhylogeny(D):

1. Pick an arbitrary leaf j .
2. Compute its limb length, $\text{LimbLength}(j)$.
3. Subtract $\text{LimbLength}(j)$ 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. Construct $\text{Tree}(D^{\text{trim}})$.
6. Identify the point in $\text{Tree}(D^{\text{trim}})$ where leaf j should be attached.
7. Attach j by an edge of length $\text{LimbLength}(j)$ in order to form $\text{Tree}(D)$.

Attaching a Limb

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 11 | 21 | 22 |
| j | 11 | 0 | 10 | 11 |
| k | 21 | 10 | 0 | 13 |
| l | 22 | 11 | 13 | 0 |

$\text{TREE}(D^{\text{trim}})$

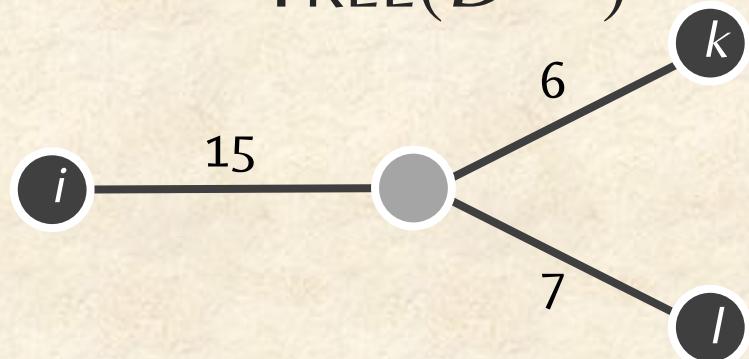


Limb Length Theorem: the length of the limb of j is equal to the minimum value of $(D^{\text{bald}}_{i,j} + D^{\text{bald}}_{j,k} - D^{\text{bald}}_{i,k})/2$ over all leaves i and k .

Attaching a Limb

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 11 | 21 | 22 |
| j | 11 | 0 | 10 | 11 |
| k | 21 | 10 | 0 | 13 |
| l | 22 | 11 | 13 | 0 |

TREE(D^{trim})

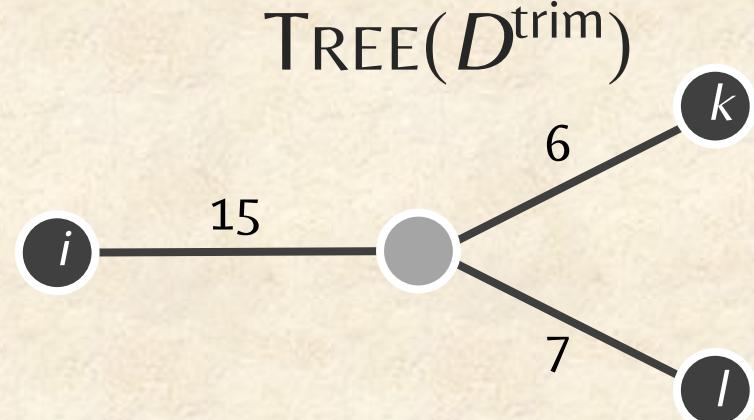


Limb Length Theorem: the length of the limb of j is equal to the minimum value of $(D^{\text{bald}}_{i,j} + D^{\text{bald}}_{j,k} - D^{\text{bald}}_{i,k})/2$ over all leaves i and k .

$$(D^{\text{bald}}_{i,j} + D^{\text{bald}}_{j,k} - D^{\text{bald}}_{i,k})/2 = 0$$

Attaching a Limb

| | i | j | k | l |
|-----|-----|-----------|-----------|-----|
| i | 0 | 11 | 21 | 22 |
| j | 11 | 0 | 10 | 11 |
| k | 21 | 10 | 0 | 13 |
| l | 22 | 11 | 13 | 0 |

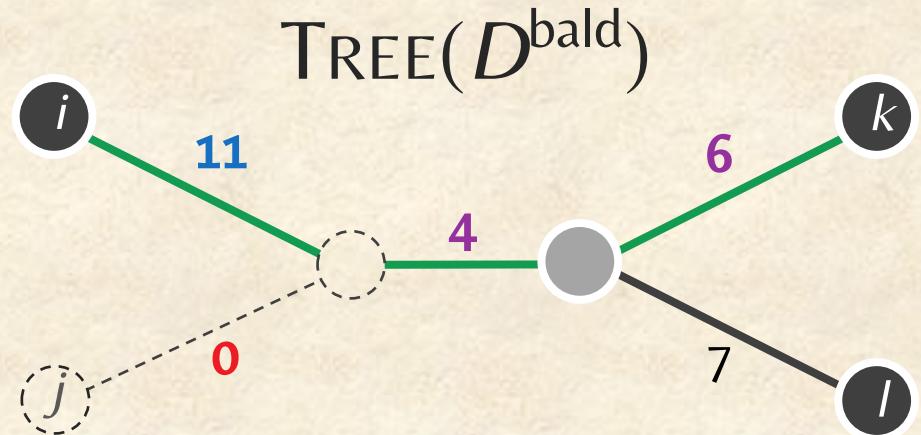


$$(D^{\text{bald}}_{i,j} + D^{\text{bald}}_{j,k} - D^{\text{bald}}_{i,k})/2 = \mathbf{0}$$

$$D^{\text{bald}}_{i,j} + D^{\text{bald}}_{j,k} = D^{\text{bald}}_{i,k}$$

Attaching a Limb

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 11 | 21 | 22 |
| j | 11 | 0 | 10 | 11 |
| k | 21 | 10 | 0 | 13 |
| l | 22 | 11 | 13 | 0 |



The attachment point for j is found on the path between leaves i and k at distance $D^{\text{bald}}_{i,j}$ from i .

$$D^{\text{bald}}_{i,j} + D^{\text{bald}}_{j,k} = D^{\text{bald}}_{i,k}$$

AdditivePhylogeny

AdditivePhylogeny(D):

1. Pick an arbitrary leaf j .
2. Compute its limb length, $\text{LimbLength}(j)$.
3. Subtract $\text{LimbLength}(j)$ 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. Construct $\text{Tree}(D^{\text{trim}})$.
6. Identify the point in $\text{Tree}(D^{\text{trim}})$ where leaf j should be attached.
7. Attach j by an edge of length $\text{LimbLength}(j)$ in order to form $\text{Tree}(D)$.

Code Challenge: Implement AdditivePhylogeny.

Spike Protein Distance Matrix

| Cow | Pig | Horse | Mouse | Dog | Cat | Turkey | Civet | Human | |
|--------|-----|-------|-------|-----|-----|--------|-------|-------|-----|
| Cow | 0 | 226 | 249 | 436 | 958 | 916 | 730 | 787 | 785 |
| Pig | 226 | 0 | 292 | 436 | 903 | 905 | 744 | 802 | 813 |
| Horse | 249 | 292 | 0 | 426 | 927 | 907 | 735 | 795 | 791 |
| Mouse | 436 | 436 | 426 | 0 | 917 | 946 | 725 | 767 | 782 |
| Dog | 958 | 903 | 927 | 917 | 0 | 706 | 730 | 844 | 846 |
| Cat | 916 | 905 | 907 | 946 | 706 | 0 | 736 | 840 | 836 |
| Turkey | 730 | 744 | 735 | 725 | 730 | 736 | 0 | 763 | 760 |
| Civet | 787 | 802 | 795 | 767 | 844 | 840 | 763 | 0 | 16 |
| Human | 785 | 813 | 791 | 782 | 846 | 836 | 760 | 16 | 0 |

This matrix isn't additive! Let's fudge it a little...

Spike Protein Distance Matrix

| | Cow | Pig | Horse | Mouse | Dog | Cat | Turkey | Civet | Human |
|--------|------|------|-------|-------|------|------|--------|-------|-------|
| Cow | 0 | 295 | 306 | 497 | 1081 | 1091 | 1003 | 956 | 954 |
| Pig | 295 | 0 | 309 | 500 | 1084 | 1094 | 1006 | 959 | 957 |
| Horse | 306 | 309 | 0 | 489 | 1073 | 1083 | 995 | 948 | 946 |
| Mouse | 497 | 500 | 489 | 0 | 1092 | 1102 | 1014 | 967 | 965 |
| Dog | 1081 | 1084 | 1073 | 1092 | 0 | 818 | 1056 | 1053 | 1051 |
| Cat | 1091 | 1094 | 1083 | 1102 | 818 | 0 | 1066 | 1063 | 1061 |
| Turkey | 1003 | 1006 | 995 | 1014 | 1056 | 1066 | 0 | 975 | 973 |
| Civet | 956 | 959 | 948 | 967 | 1053 | 1063 | 975 | 0 | 16 |
| Human | 954 | 957 | 946 | 965 | 1051 | 1061 | 973 | 16 | 0 |

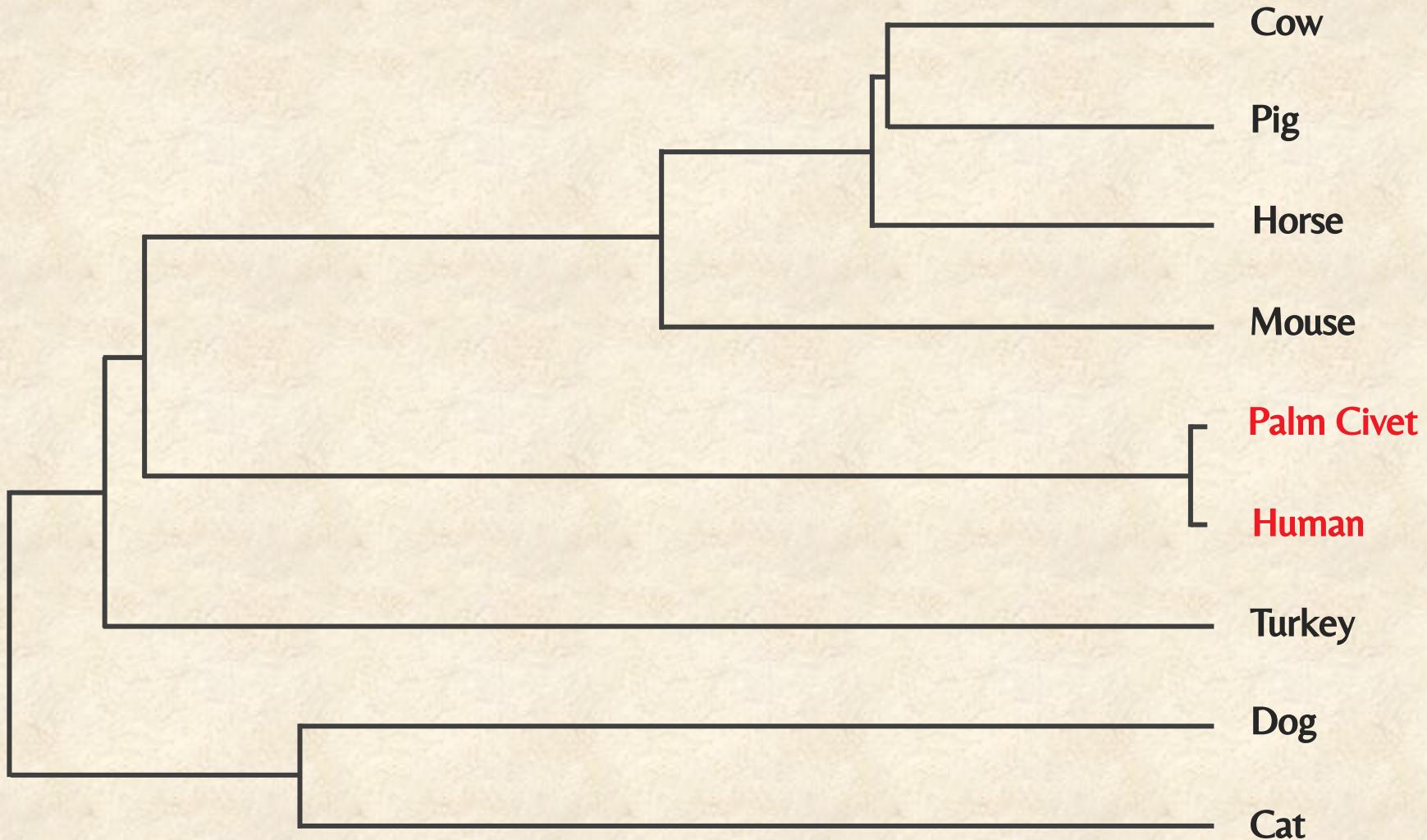
STOP and Think: Which animal gave us SARS?

Spike Protein Distance Matrix

| Cow | Pig | Horse | Mouse | Dog | Cat | Turkey | Civet | Human | |
|--------|------|-------|-------|------|------|--------|-------|-------|------|
| Cow | 0 | 295 | 306 | 497 | 1081 | 1091 | 1003 | 956 | 954 |
| Pig | 295 | 0 | 309 | 500 | 1084 | 1094 | 1006 | 959 | 957 |
| Horse | 306 | 309 | 0 | 489 | 1073 | 1083 | 995 | 948 | 946 |
| Mouse | 497 | 500 | 489 | 0 | 1092 | 1102 | 1014 | 967 | 965 |
| Dog | 1081 | 1084 | 1073 | 1092 | 0 | 818 | 1056 | 1053 | 1051 |
| Cat | 1091 | 1094 | 1083 | 1102 | 818 | 0 | 1066 | 1063 | 1061 |
| Turkey | 1003 | 1006 | 995 | 1014 | 1056 | 1066 | 0 | 975 | 973 |
| Civet | 956 | 959 | 948 | 967 | 1053 | 1063 | 975 | 0 | 16 |
| Human | 954 | 957 | 946 | 965 | 1051 | 1061 | 973 | 16 | 0 |

STOP and Think: Which animal gave us SARS?

Coronavirus Phylogeny











Spike Protein Distance Matrix

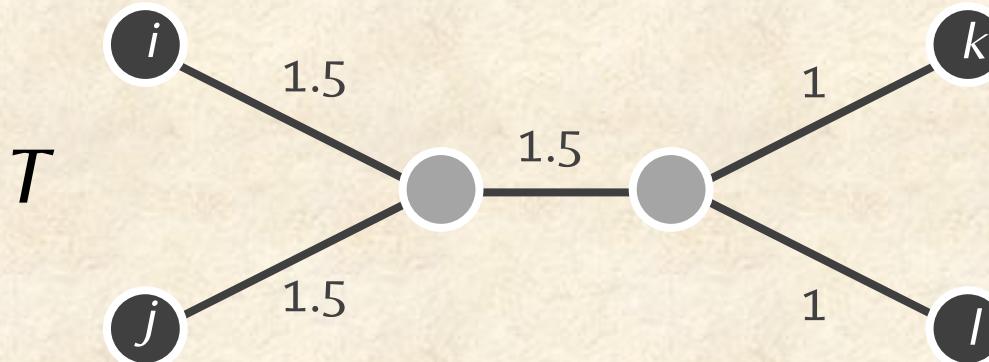
| Cow | Pig | Horse | Mouse | Dog | Cat | Turkey | Civet | Human | |
|--------|-----|-------|-------|-----|-----|--------|-------|-------|-----|
| Cow | 0 | 226 | 249 | 436 | 958 | 916 | 730 | 787 | 785 |
| Pig | 226 | 0 | 292 | 436 | 903 | 905 | 744 | 802 | 813 |
| Horse | 249 | 292 | 0 | 426 | 927 | 907 | 735 | 795 | 791 |
| Mouse | 436 | 436 | 426 | 0 | 917 | 946 | 725 | 767 | 782 |
| Dog | 958 | 903 | 927 | 917 | 0 | 706 | 730 | 844 | 846 |
| Cat | 916 | 905 | 907 | 946 | 706 | 0 | 736 | 840 | 836 |
| Turkey | 730 | 744 | 735 | 725 | 730 | 736 | 0 | 763 | 760 |
| Civet | 787 | 802 | 795 | 767 | 844 | 840 | 763 | 0 | 16 |
| Human | 785 | 813 | 791 | 782 | 846 | 836 | 760 | 16 | 0 |

But what should we do about *non-additive* matrices?

Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- **Using Least-Squares to Construct Distance-Based Phylogenies**
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

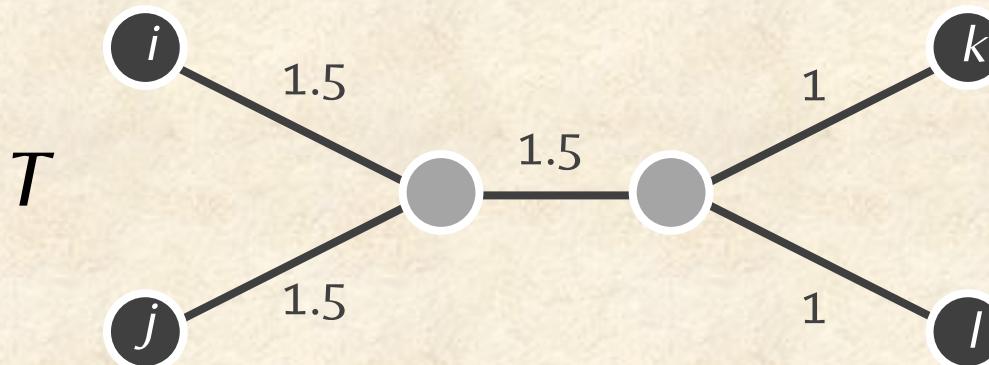
Sum of Squared Errors



D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

Sum of Squared Errors



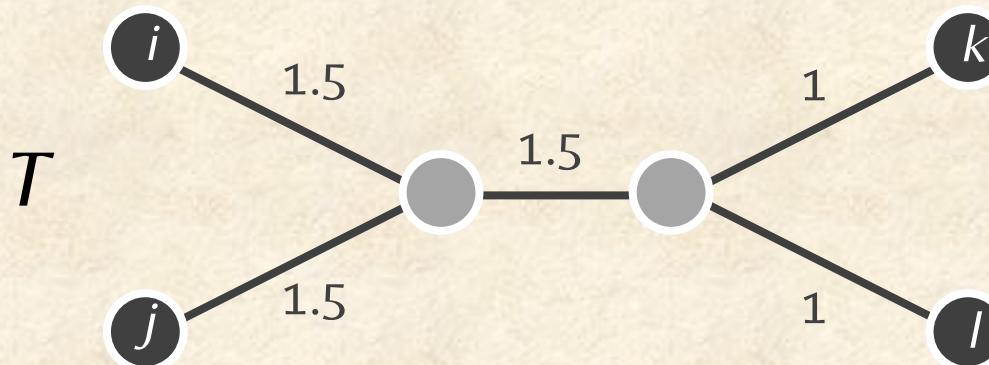
| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 4 |
| j | 3 | 0 | 4 | 4 |
| k | 4 | 4 | 0 | 2 |
| l | 4 | 4 | 2 | 0 |

d

Sum of Squared Errors



D

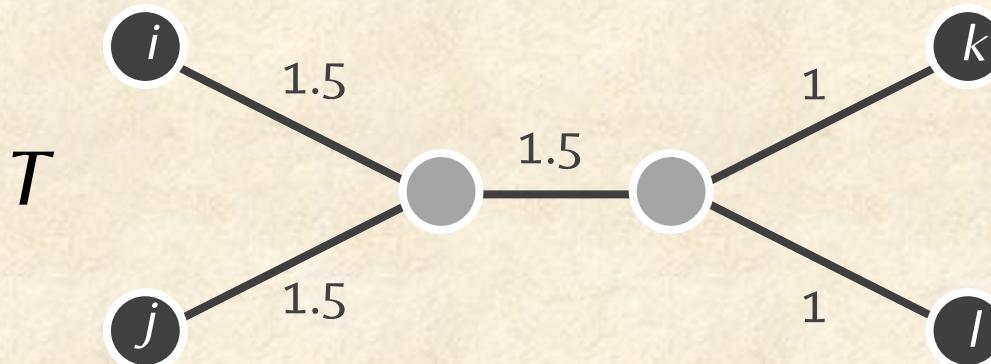
| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

d

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 4 |
| j | 3 | 0 | 4 | 4 |
| k | 4 | 4 | 0 | 2 |
| l | 4 | 4 | 2 | 0 |

Sum of Squared Errors

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



D

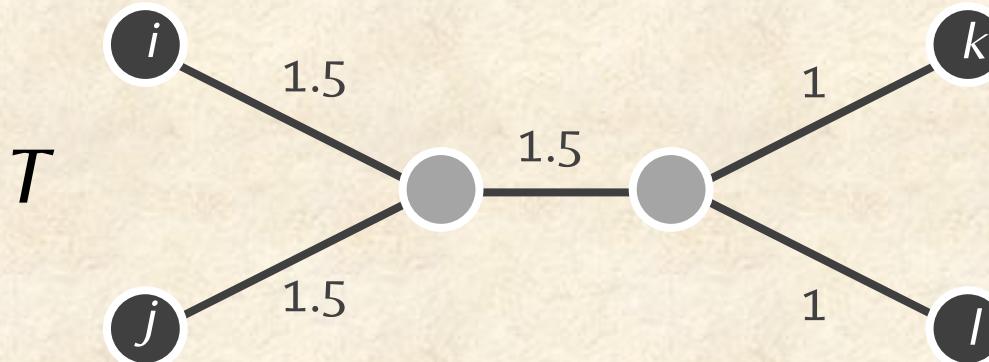
| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

d

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 4 |
| j | 3 | 0 | 4 | 4 |
| k | 4 | 4 | 0 | 2 |
| l | 4 | 4 | 2 | 0 |

Sum of Squared Errors

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



| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

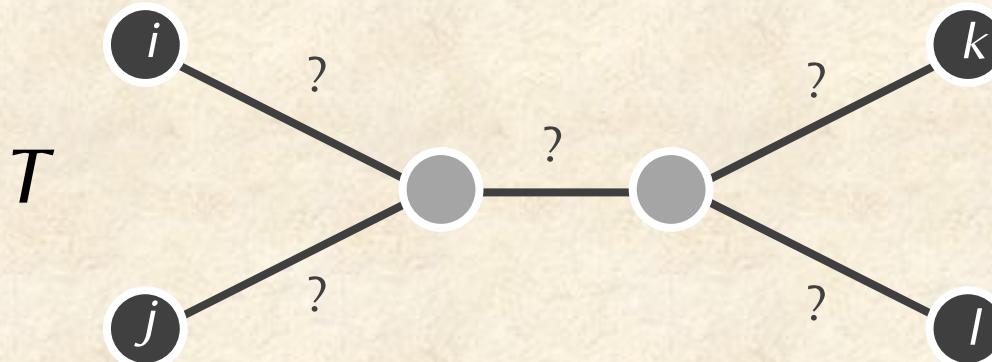
D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 4 |
| j | 3 | 0 | 4 | 4 |
| k | 4 | 4 | 0 | 2 |
| l | 4 | 4 | 2 | 0 |

d

Sum of Squared Errors

Exercise Break: Assign lengths to edges in T in order to minimize $\text{Discrepancy}(T, D)$.



| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | ? | ? | ? |
| j | ? | 0 | ? | ? |
| k | ? | ? | 0 | ? |
| l | ? | ? | ? | 0 |

d

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 $\text{Discrepancy}(T, D)$ 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 $\text{Discrepancy}(T, D)$ over all weighted trees with n leaves.

Unfortunately, this problem is NP -Complete...

Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- **Ultrametric Evolutionary Trees**
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

Modeling Speciations

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

Modeling Speciations



Squirrel
Monkey

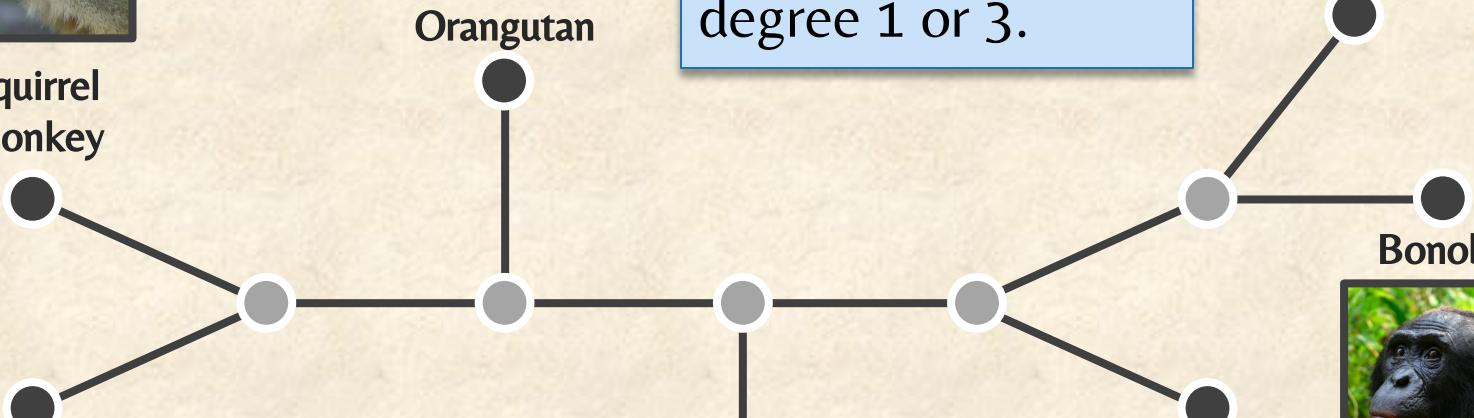


Orangutan

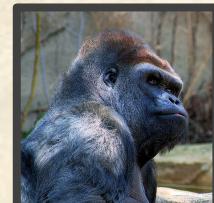
Unrooted binary tree: every node has degree 1 or 3.



Chimpanzee



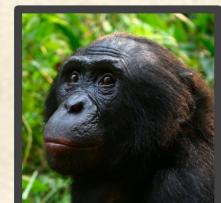
Baboon



Gorilla

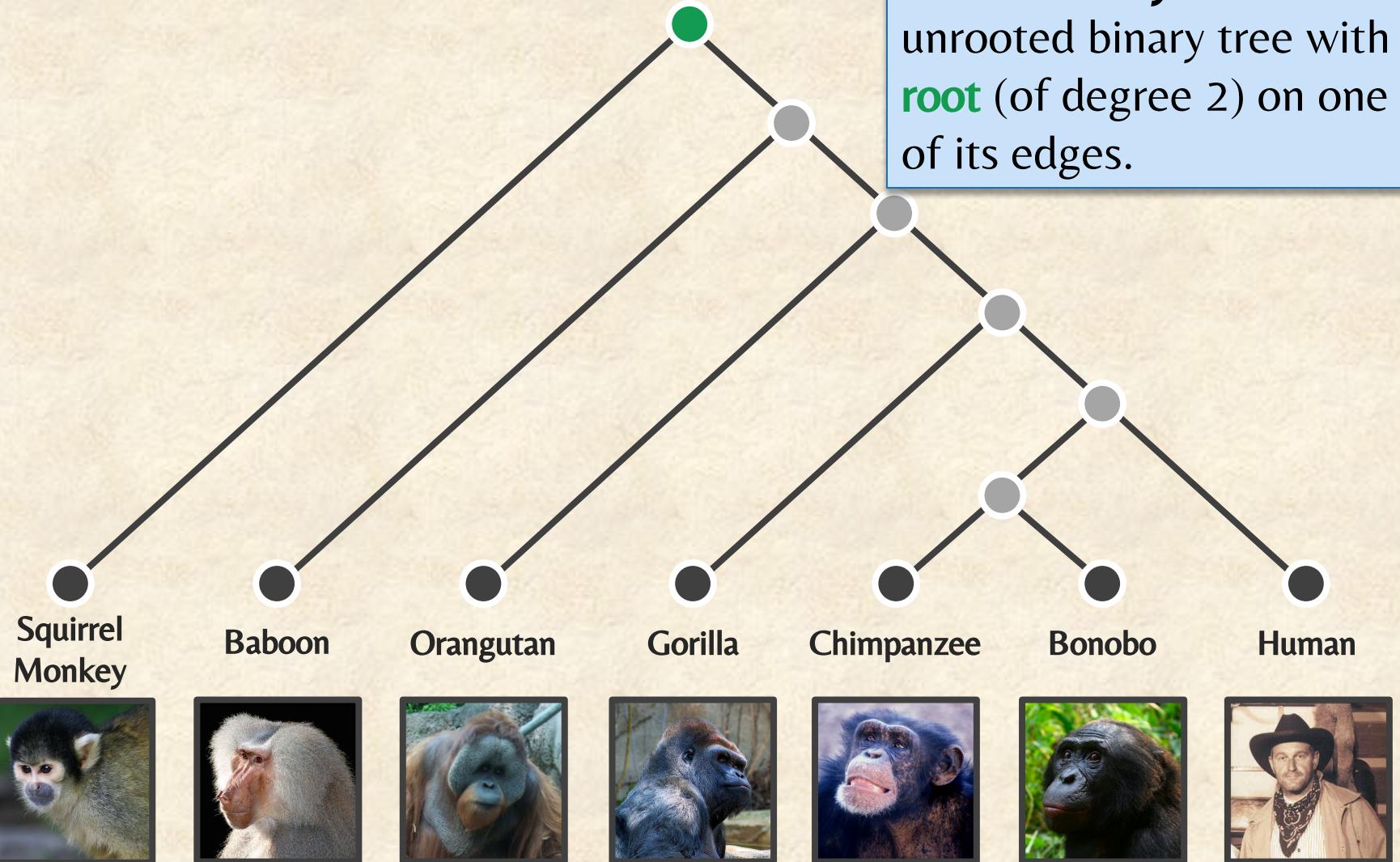


Human

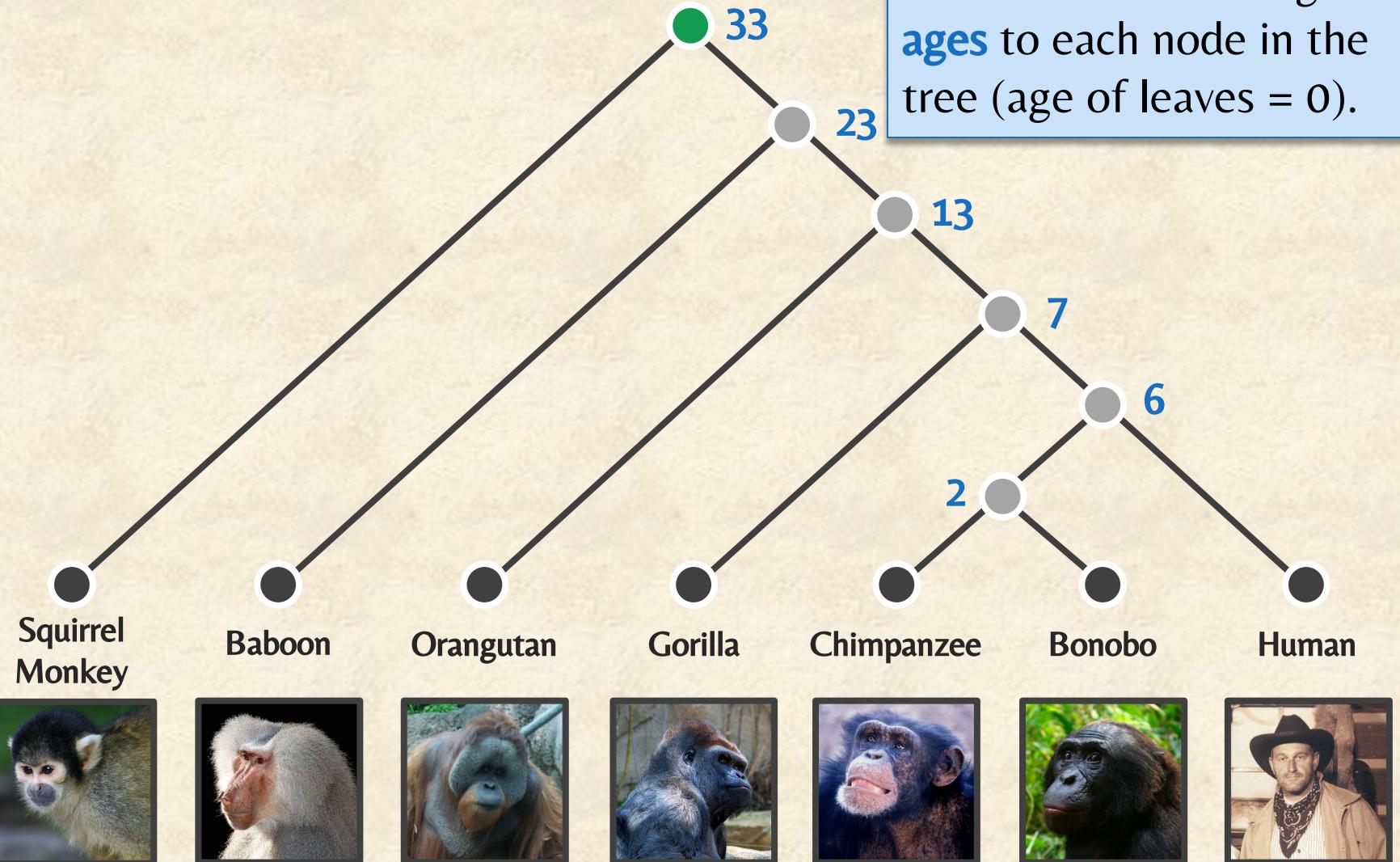


Bonobo

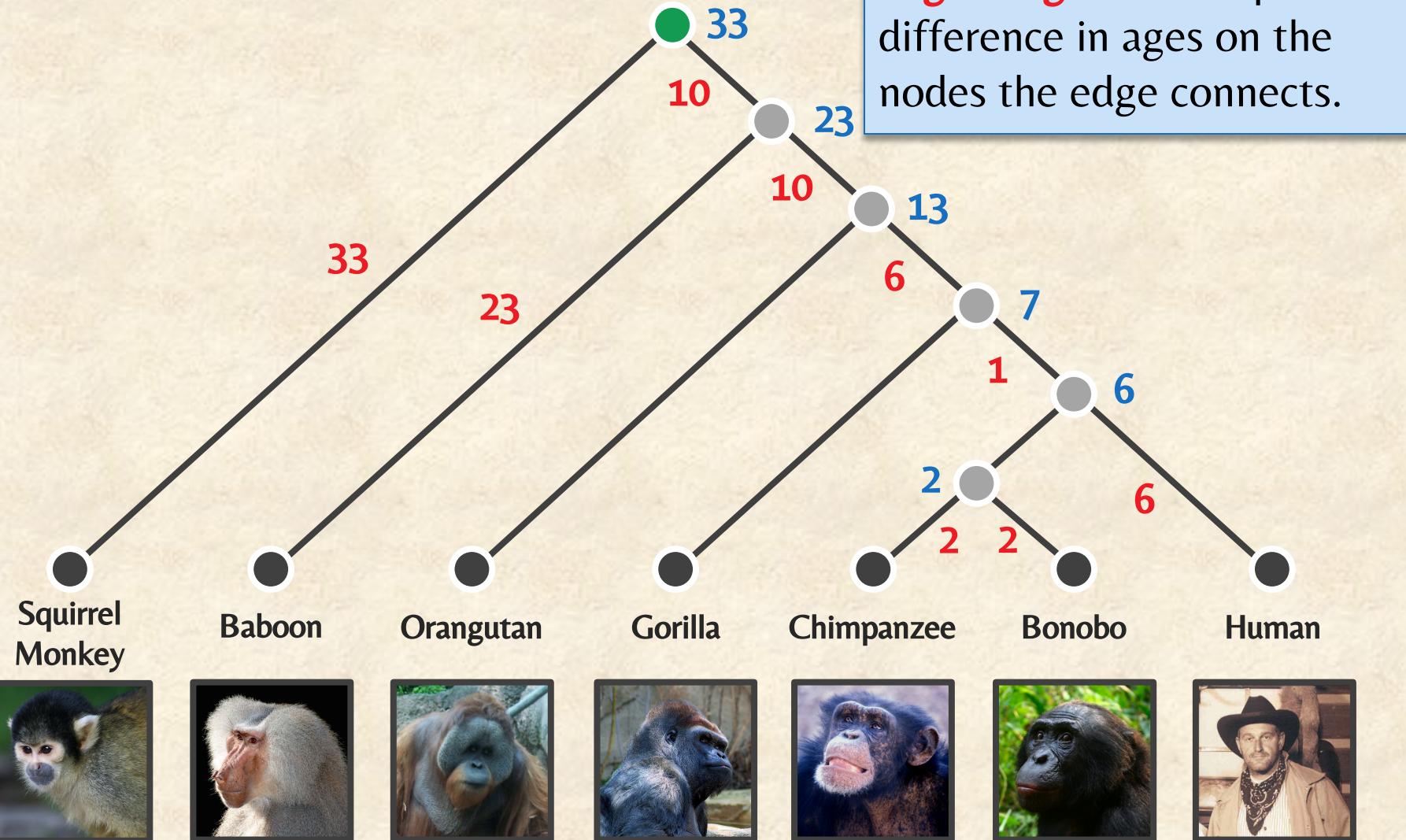
Modeling Speciations



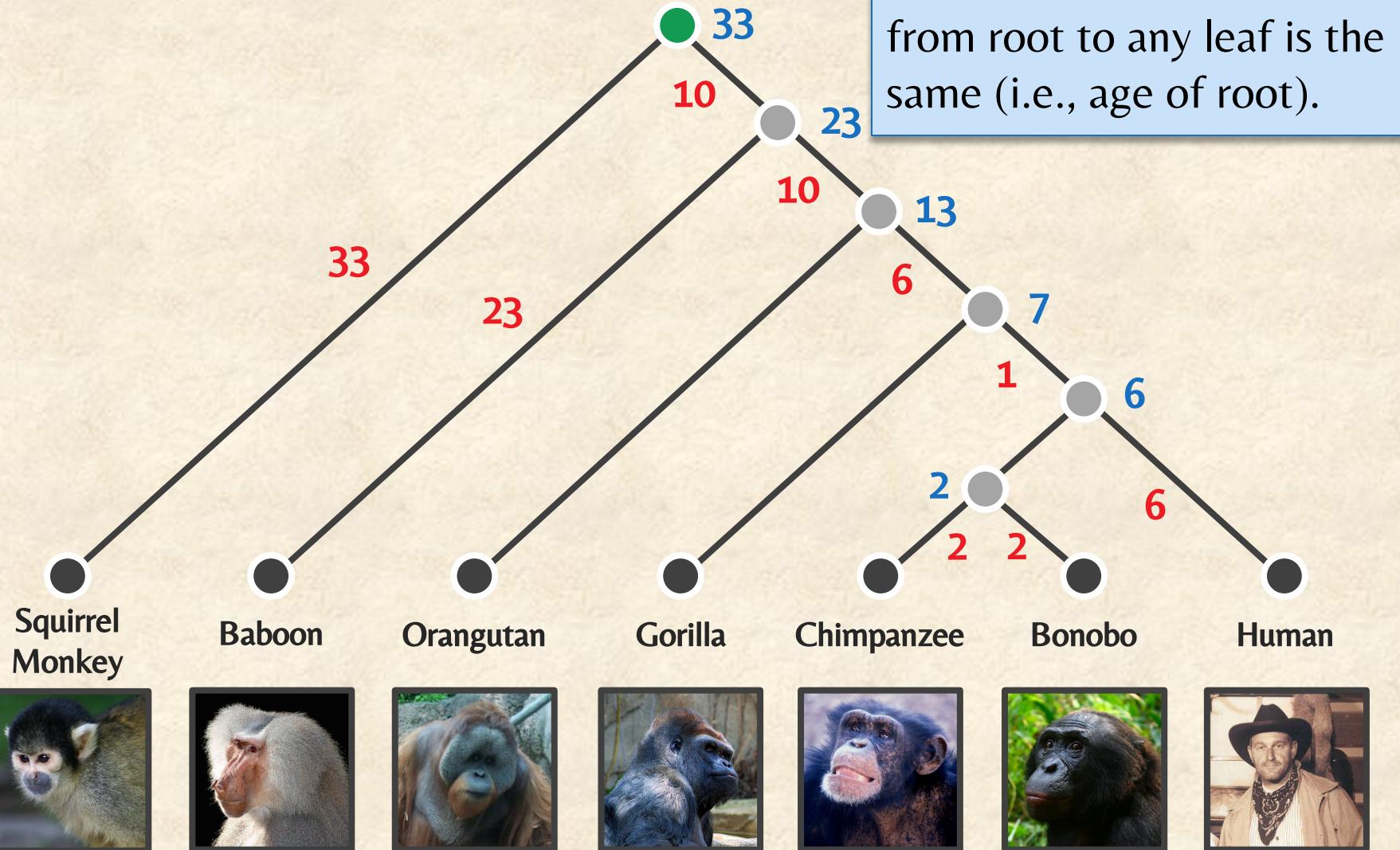
Ultrametric Trees



Ultrametric Trees



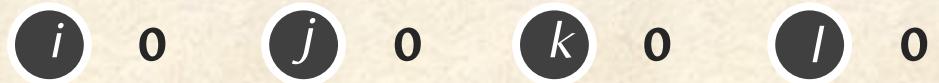
Ultrametric Trees



UPGMA: A Clustering Heuristic

1. Form a cluster for each present-day species, each containing a single leaf.

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |



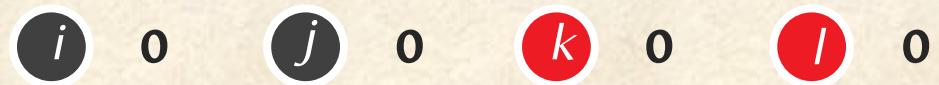
UPGMA: A Clustering Heuristic

2. Find the two closest clusters C_1 and C_2 according to the average distance

$$D_{\text{avg}}(C_1, C_2) = \sum_{i \in C_1, j \in C_2} D_{ij} / |C_1| \cdot |C_2|$$

where $|C|$ denotes the number of elements in C .

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 3 | 4 | 3 |
| <i>j</i> | 3 | 0 | 4 | 5 |
| <i>k</i> | 4 | 4 | 0 | 2 |
| <i>l</i> | 3 | 5 | 2 | 0 |



UPGMA: A Clustering Heuristic

3. Merge C_1 and C_2 into a single cluster C .

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |

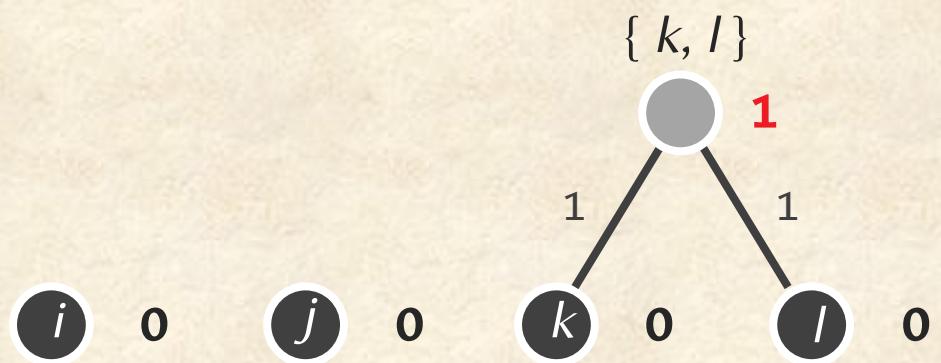
$\{ k, l \}$



UPGMA: A Clustering Heuristic

4. Form a new node for C and connect to C_1 and C_2 by an edge. Set age of C as $D_{\text{avg}}(C_1, C_2)/2$.

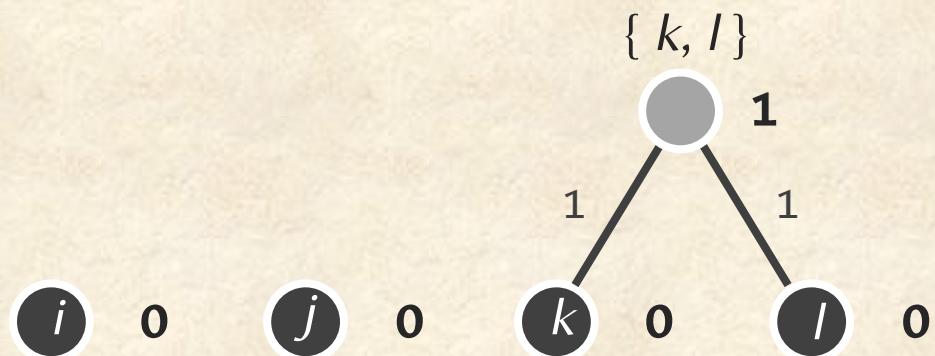
| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 3 | 4 | 3 |
| j | 3 | 0 | 4 | 5 |
| k | 4 | 4 | 0 | 2 |
| l | 3 | 5 | 2 | 0 |



UPGMA: A Clustering Heuristic

5. Update the distance matrix by computing the average distance between each pair of clusters.

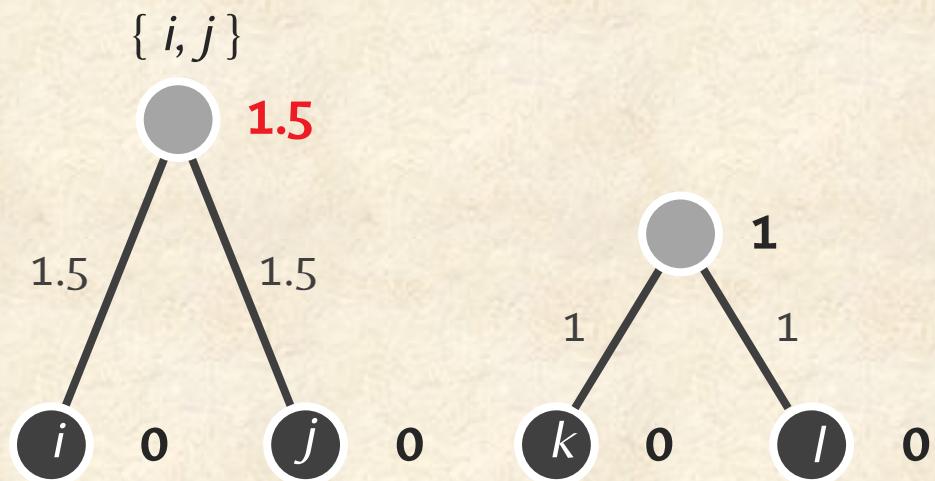
| | i | j | $\{k, l\}$ |
|------------|-----|-----|------------|
| i | 0 | 3 | 3.5 |
| j | 3 | 0 | 4.5 |
| $\{k, l\}$ | 3.5 | 4.5 | 0 |



UPGMA: A Clustering Heuristic

6. Iterate until a single cluster contains all species.

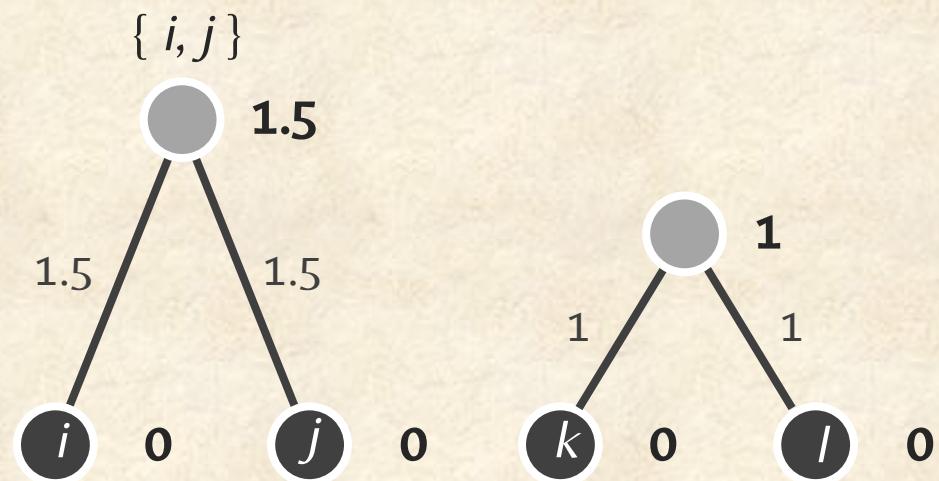
| | i | j | $\{k, l\}$ |
|------------|-----|-----|------------|
| i | 0 | 3 | 3.5 |
| j | 3 | 0 | 4.5 |
| $\{k, l\}$ | 3.5 | 4.5 | 0 |



UPGMA: A Clustering Heuristic

6. Iterate until a single cluster contains all species.

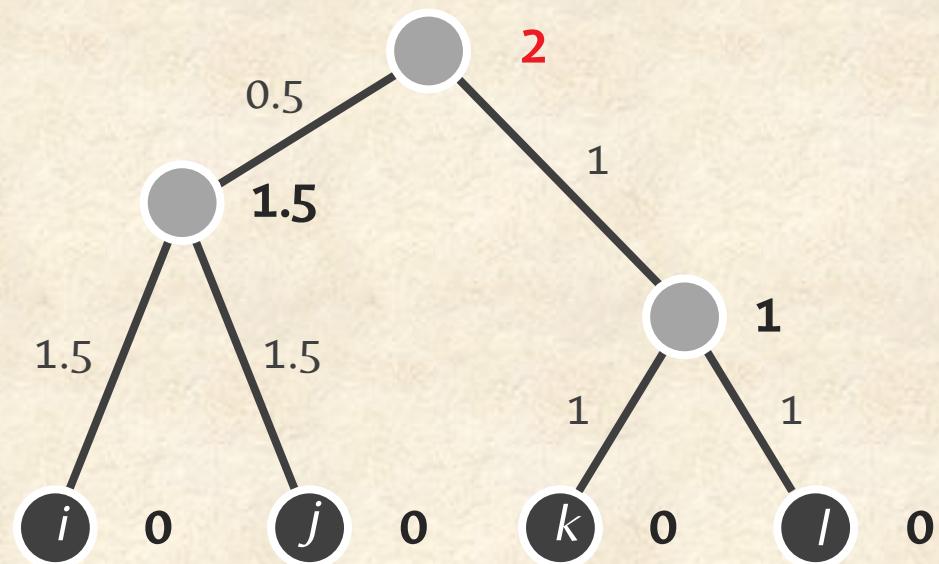
| | $\{i, j\}$ | $\{k, l\}$ |
|------------|------------|------------|
| $\{i, j\}$ | 0 | 4 |
| $\{k, l\}$ | 4 | 0 |



UPGMA: A Clustering Heuristic

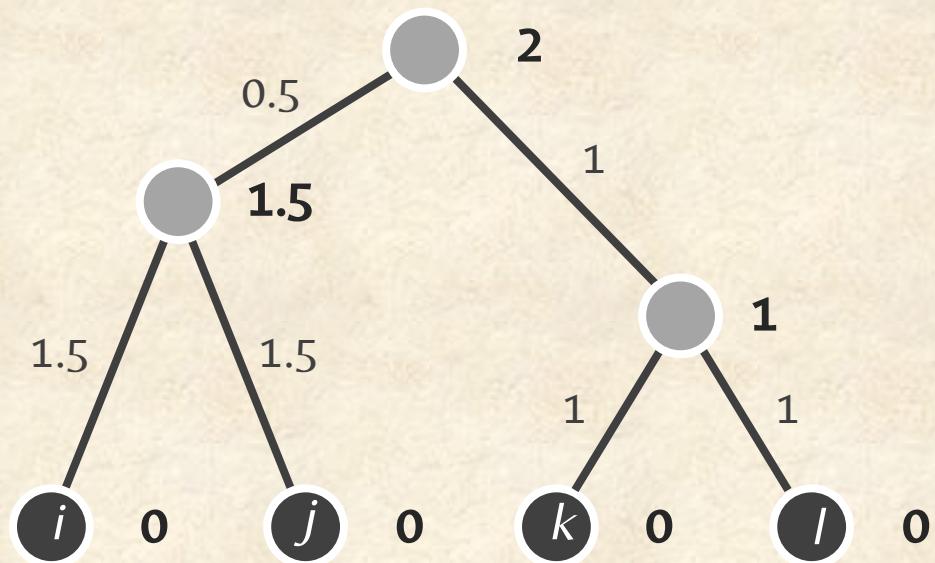
6. Iterate until a single cluster contains all species.

| | $\{i, j\}$ | $\{k, l\}$ |
|------------|------------|------------|
| $\{i, j\}$ | 0 | 4 |
| $\{k, l\}$ | 4 | 0 |



UPGMA: A Clustering Heuristic

6. Iterate until a single cluster contains all species.



UPGMA: A Clustering Heuristic

UPGMA(D):

1. Form a cluster for each present-day species, each containing a single leaf.
2. Find the two closest clusters C_1 and C_2 according to the average distance

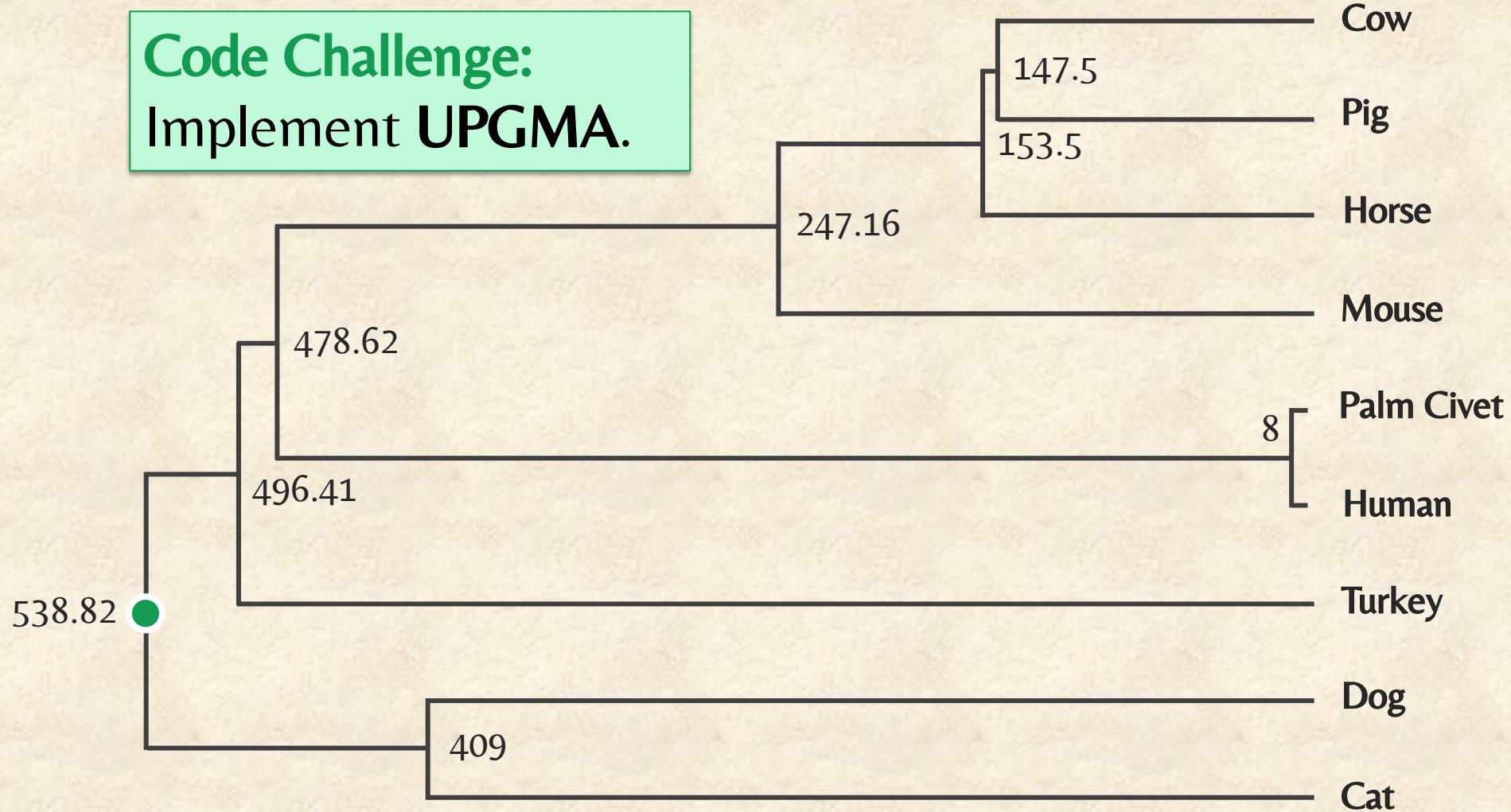
$$D_{\text{avg}}(C_1, C_2) = \sum_{i \in C_1, j \in C_2} D_{i,j} / |C_1| \cdot |C_2|$$

where $|C|$ denotes the number of elements in C

3. Merge C_1 and C_2 into a single cluster C .
4. Form a new node for C and connect to C_1 and C_2 by an edge.
Set age of C as $D_{\text{avg}}(C_1, C_2)/2$.
5. Update the distance matrix by computing the average distance between each pair of clusters.
6. Iterate steps 2-5 until a single cluster contains all species.

Applying UPGMA to Spike Proteins

Code Challenge:
Implement UPGMA.

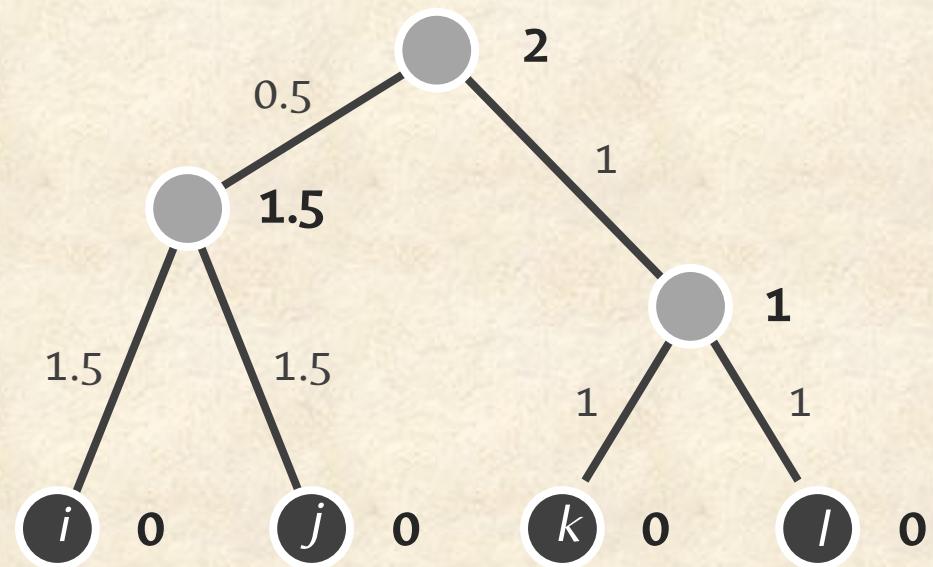




You never stated a
computational
problem!

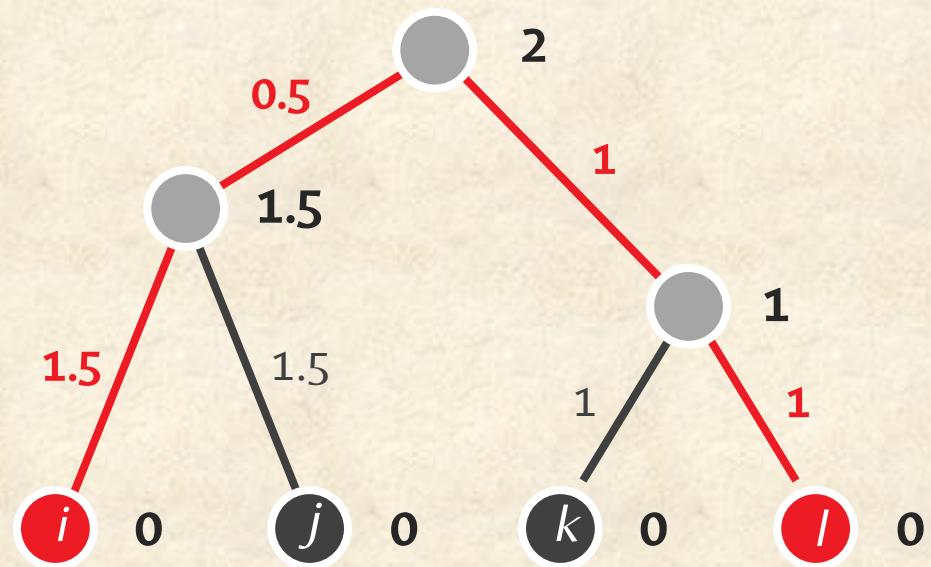
UPGMA Doesn't “Fit” a Tree to a Matrix

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 3 | 4 | 3 |
| <i>j</i> | 3 | 0 | 4 | 5 |
| <i>k</i> | 4 | 4 | 0 | 2 |
| <i>l</i> | 3 | 5 | 2 | 0 |



UPGMA Doesn't “Fit” a Tree to a Matrix

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 3 | 4 | 3 |
| <i>j</i> | 3 | 0 | 4 | 5 |
| <i>k</i> | 4 | 4 | 0 | 2 |
| <i>l</i> | 3 | 5 | 2 | 0 |



In Summary...

- Additive Phylogeny:

In Summary...

- **Additive Phylogeny:**
 - good: produces the tree fitting an *additive* matrix

In Summary...

- **Additive Phylogeny:**
 - good: produces the tree fitting an *additive* matrix
 - bad: fails completely on a *non-additive* matrix

In Summary...

- **Additive Phylogeny:**
 - good: produces the tree fitting an *additive* matrix
 - bad: fails completely on a *non-additive* matrix
- **UPGMA:**
 - good: produces a tree for any matrix

In Summary...

- **Additive Phylogeny:**
 - good: produces the tree fitting an *additive* matrix
 - bad: fails completely on a *non-additive* matrix
- **UPGMA:**
 - good: produces a tree for any matrix
 - bad: tree doesn't necessarily fit an additive matrix

In Summary...

- **Additive Phylogeny:**
 - good: produces the tree fitting an *additive* matrix
 - bad: fails completely on a *non-additive* matrix
- **UPGMA:**
 - good: produces a tree for any matrix
 - bad: tree doesn't necessarily fit an additive matrix
- **?????:**
 - good: produces the tree fitting an additive matrix
 - good: provides heuristic for a non-additive matrix

Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- **The Neighbor-Joining Algorithm**
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

Neighbor-Joining Theorem

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

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

where $\text{TotalDistance}_D(i)$ is the sum of distances from i to all other leaves.

Neighbor-Joining Theorem

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

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

where $\text{TotalDistance}_D(i)$ is the sum of distances from i to all other leaves.

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 13 | 21 | 22 |
| <i>j</i> | 13 | 0 | 12 | 13 |
| <i>k</i> | 21 | 12 | 0 | 13 |
| <i>l</i> | 22 | 13 | 13 | 0 |

Neighbor-Joining Theorem

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

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

where $\text{TotalDistance}_D(i)$ is the sum of distances from i to all other leaves.

D

| | i | j | k | l | |
|-----|-----|-----|-----|-----|--------------------------|
| i | 0 | 13 | 21 | 22 | TotalDistance_D |
| j | 13 | 0 | 12 | 13 | |
| k | 21 | 12 | 0 | 13 | |
| l | 22 | 13 | 13 | 0 | |

Neighbor-Joining Theorem

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

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

where $\text{TotalDistance}_D(i)$ is the sum of distances from i to all other leaves.

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

| | TotalDistance_D | | | |
|--|--------------------------|----|----|----|
| | 56 | 38 | 46 | 48 |

D^*

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | -68 | -60 | -60 |
| j | -68 | 0 | -60 | -60 |
| k | -60 | -60 | 0 | -68 |
| l | -60 | -60 | -68 | 0 |

Neighbor-Joining Theorem

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

D

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

| | $TotalDistance_D$ | | | |
|-----|-------------------|--|--|--|
| i | 56 | | | |
| j | 38 | | | |
| k | 46 | | | |
| l | 48 | | | |

D^*

| | i | j | k | l |
|-----|-----|-----|-----|-----|
| i | 0 | -68 | -60 | -60 |
| j | -68 | 0 | -60 | -60 |
| k | -60 | -60 | 0 | -68 |
| l | -60 | -60 | -68 | 0 |

Neighbor-Joining Theorem

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

D

| | i | j | k | l |
|-----|-----|-----------|-----------|-----|
| i | 0 | 13 | 21 | 22 |
| j | 13 | 0 | 12 | 13 |
| k | 21 | 12 | 0 | 13 |
| l | 22 | 13 | 13 | 0 |

| | $TotalDistance_D$ | | | |
|-----|-------------------|--|--|--|
| i | 56 | | | |
| j | 38 | | | |
| k | 46 | | | |
| l | 48 | | | |

D^*

| | i | j | k | l |
|-----|------------|------------|------------|------------|
| i | 0 | -68 | -60 | -60 |
| j | -68 | 0 | -60 | -60 |
| k | -60 | -60 | 0 | -68 |
| l | -60 | -60 | -68 | 0 |

Neighbor-Joining in Action

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> | <i>TotalDistance</i> _{<i>D</i>} |
|----------|----------|----------|----------|----------|--|
| <i>i</i> | 0 | -68 | -60 | -60 | 56 |
| <i>j</i> | -68 | 0 | -60 | -60 | 38 |
| <i>k</i> | -60 | -60 | 0 | -68 | 46 |
| <i>l</i> | -60 | -60 | -68 | 0 | 48 |

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

Neighbor-Joining in Action

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> | <i>TotalDistance</i> _D |
|-----------|----------|----------|----------|----------|-----------------------------------|
| <i>D*</i> | 0 | -68 | -60 | -60 | 56 |
| | -68 | 0 | -60 | -60 | 38 |
| | -60 | -60 | 0 | -68 | 46 |
| | -60 | -60 | -68 | 0 | 48 |

2. Find a minimum element D^*_{ij} of D^* .

Neighbor-Joining in Action

| | i | j | k | l | $TotalDistance_D$ |
|-----|-----|-----|-----|-----|-------------------|
| i | 0 | -68 | -60 | -60 | 56 |
| j | -68 | 0 | -60 | -60 | 38 |
| k | -60 | -60 | 0 | -68 | 46 |
| l | -60 | -60 | -68 | 0 | 48 |

2. Find a minimum element D^*_{ij} of D^* .

Neighbor-Joining in Action

| | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> | <i>TotalDistance_D</i> | |
|-----------|----------|----------|----------|----------|----------------------------------|----|
| <i>D*</i> | <i>i</i> | 0 | -68 | -60 | -60 | 56 |
| | <i>j</i> | -68 | 0 | -60 | -60 | 38 |
| | <i>k</i> | -60 | -60 | 0 | -68 | 46 |
| | <i>l</i> | -60 | -60 | -68 | 0 | 48 |

$$\Delta_{ij} = (56 - 38) / (4 - 2) = 9$$

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

Neighbor-Joining in Action

| <i>D</i> | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> | <i>TotalDistance_D</i> | |
|----------|----------|----------|----------|----------|----------------------------------|--------------------------------------|
| <i>i</i> | 0 | 13 | 21 | 22 | 56 | |
| <i>j</i> | 13 | 0 | 12 | 13 | 38 | $\Delta_{i,j} = (56 - 38) / (4 - 2)$ |
| <i>k</i> | 21 | 12 | 0 | 13 | 46 | |
| <i>l</i> | 22 | 13 | 13 | 0 | 48 | |

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

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

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

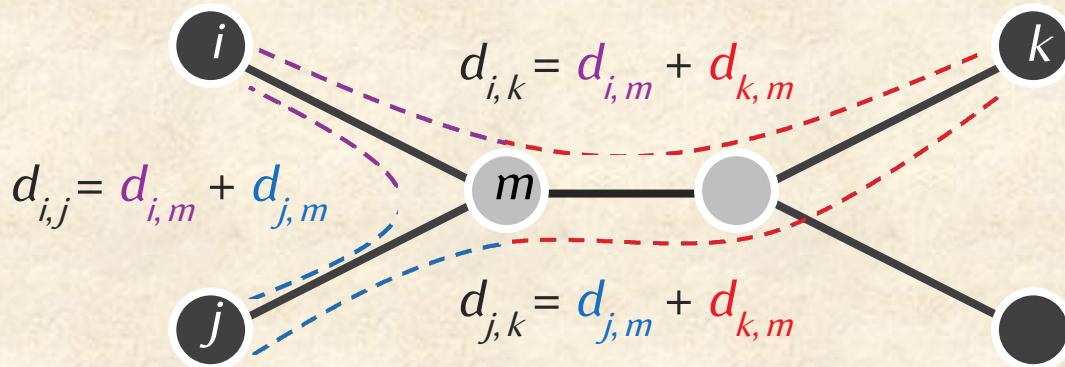
Neighbor-Joining in Action

D'

| | m | k | l | $TotalDistance_D$ |
|-----|-----|-----|-----|-------------------|
| m | 0 | 10 | 11 | 21 |
| k | 10 | 0 | 13 | 23 |
| l | 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 $d_{k,m}$



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

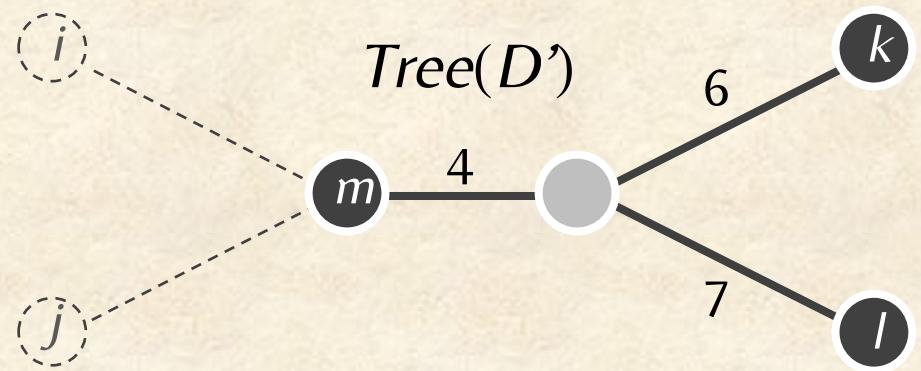
$$d_{k,m} = (d_{i,k} + d_{j,k} - d_{i,j}) / 2$$

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

Neighbor-Joining in Action

D'

| | m | k | l |
|-----|-----|-----|-----|
| m | 0 | 10 | 11 |
| k | 10 | 0 | 13 |
| l | 11 | 13 | 0 |

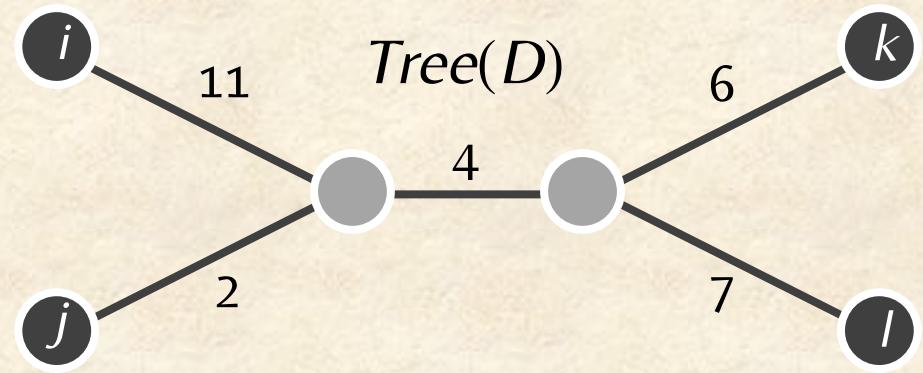


6. Apply **NeighborJoining** to D' to obtain $Tree(D')$.

Neighbor-Joining in Action

D'

| | m | k | l |
|-----|-----|-----|-----|
| m | 0 | 10 | 11 |
| k | 10 | 0 | 13 |
| l | 11 | 13 | 0 |



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

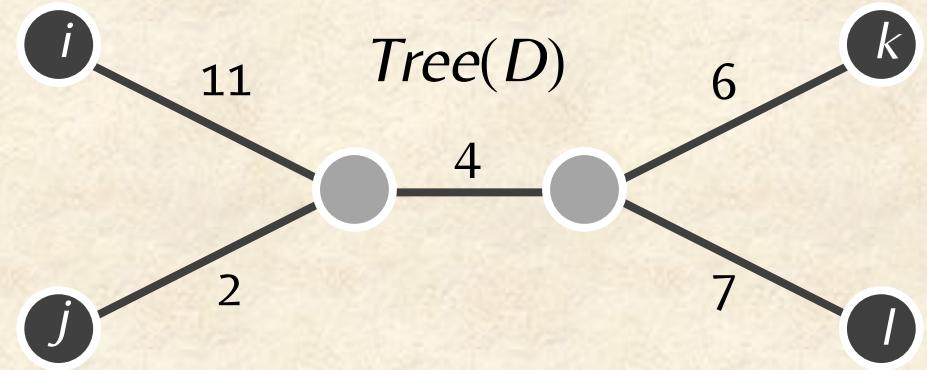
$$LimbLength(i) = \frac{1}{2}(13 - 9) = 2$$

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

Neighbor-Joining in Action

D'

| | m | k | l |
|-----|-----|-----|-----|
| m | 0 | 10 | 11 |
| k | 10 | 0 | 13 |
| l | 11 | 13 | 0 |



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

Neighbor-Joining

NeighborJoining(D):

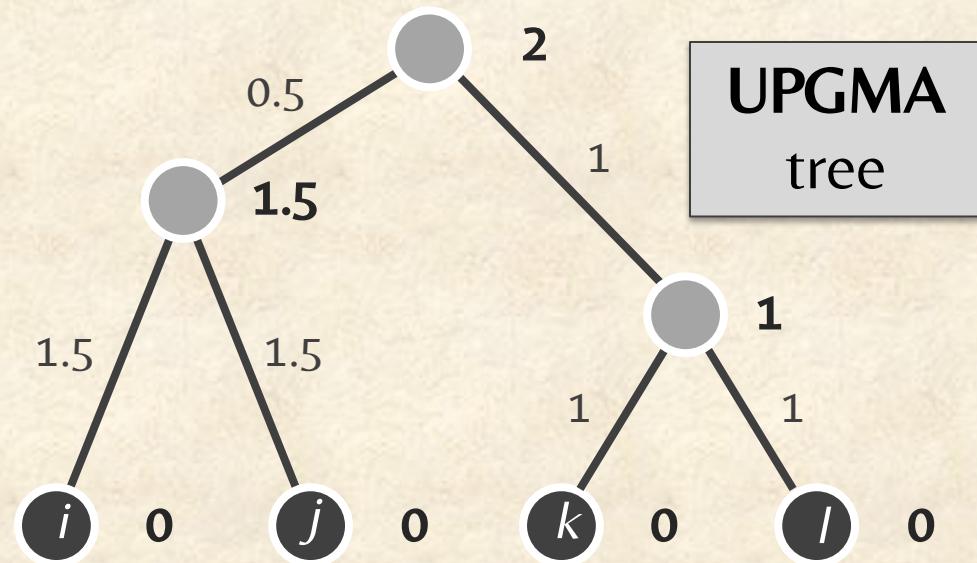
1. Construct neighbor-joining matrix D^* from D .
2. Find a minimum element $D_{i,j}^*$ of D^* .
3. Compute $\Delta_{i,j} = (\text{TotalDistance}_D(i) - \text{TotalDistance}_D(j)) / (n - 2)$.
4. Set $\text{LimbLength}(i)$ equal to $\frac{1}{2}(D_{i,j} + \Delta_{i,j})$ and $\text{LimbLength}(j)$ equal 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** to D' to obtain $\text{Tree}(D')$.
7. Reattach limbs of i and j to obtain $\text{Tree}(D)$.

Code Challenge: Implement **NeighborJoining**.

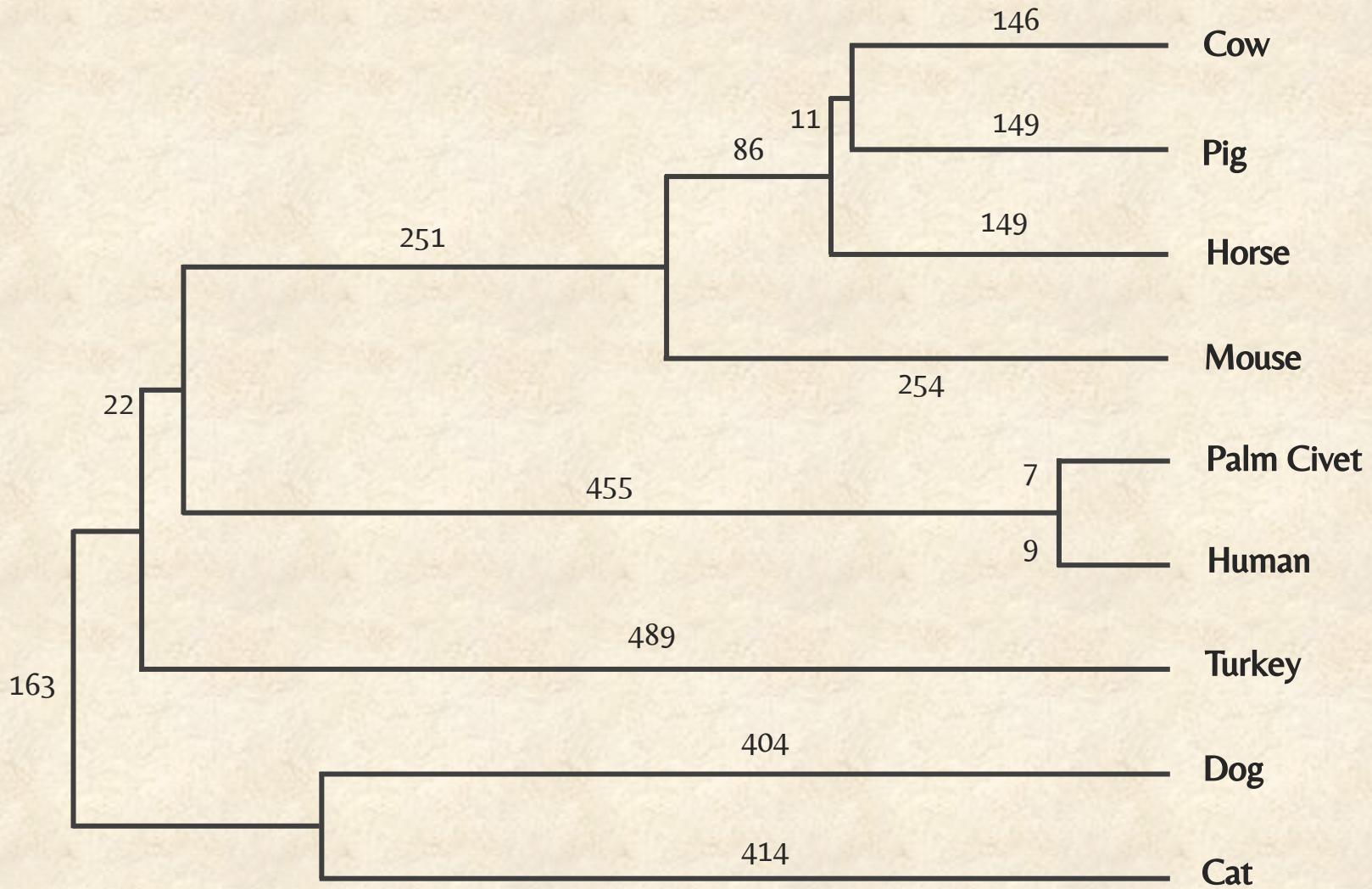
Neighbor-Joining

Exercise Break: Find the tree returned by **NeighborJoining** on the following non-additive matrix. How does the result compare with the tree produced by **UPGMA**?

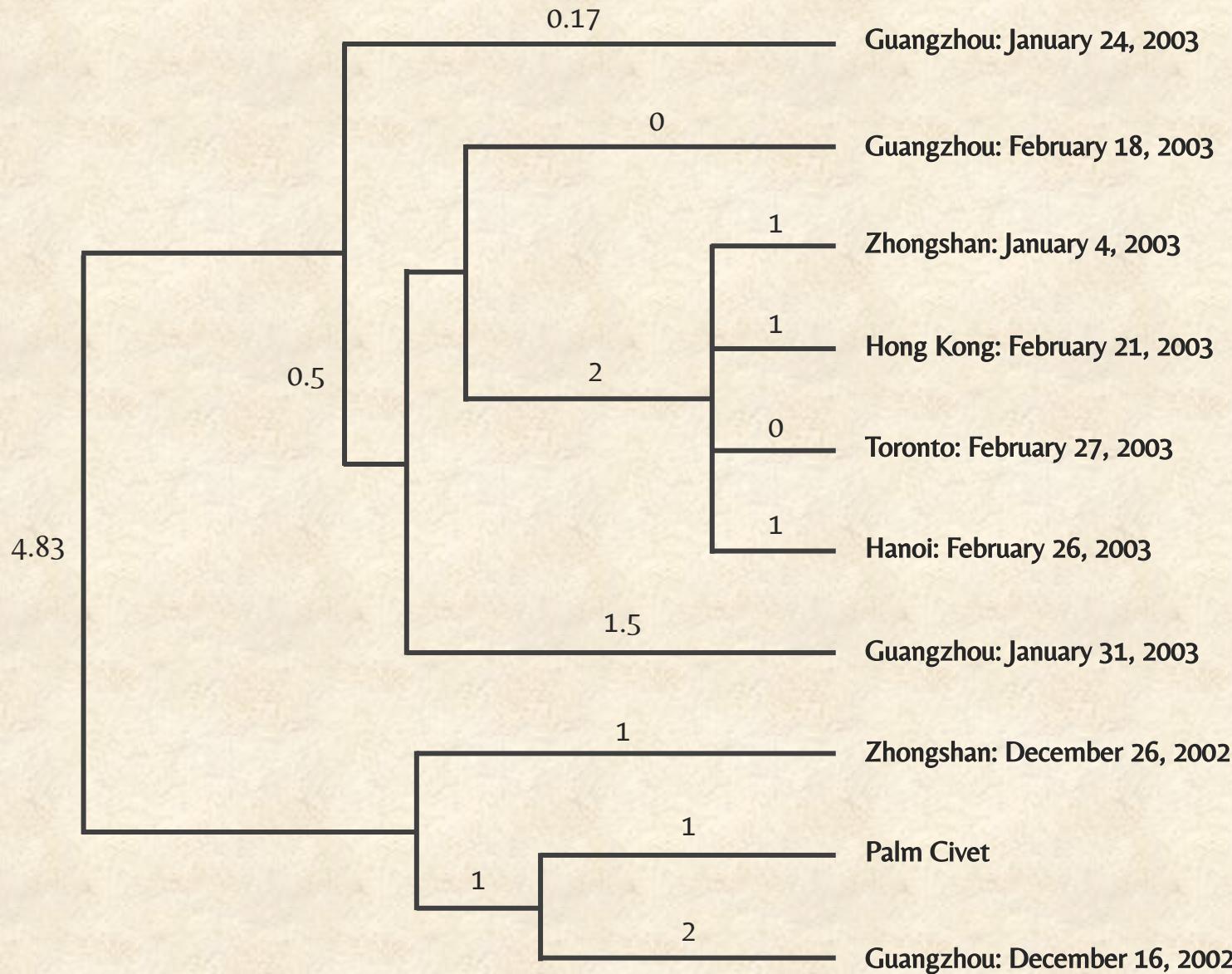
| <i>D</i> | <i>i</i> | <i>j</i> | <i>k</i> | <i>l</i> |
|----------|----------|----------|----------|----------|
| <i>i</i> | 0 | 3 | 4 | 3 |
| <i>j</i> | 3 | 0 | 4 | 5 |
| <i>k</i> | 4 | 4 | 0 | 2 |
| <i>l</i> | 3 | 5 | 2 | 0 |



Neighbor-Joining on Coronaviruses



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 |

Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- **Character-Based Tree Reconstruction**
- The Small Parsimony Problem
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

Character Tables

Fifty years ago, researchers constructed trees from anatomical/physiological properties called **characters**.



| | wings | legs |
|---------------------|-------|------|
| winged stick insect | Yes | 6 |



| | | |
|-----------------------|----|---|
| wingless stick insect | No | 6 |
|-----------------------|----|---|

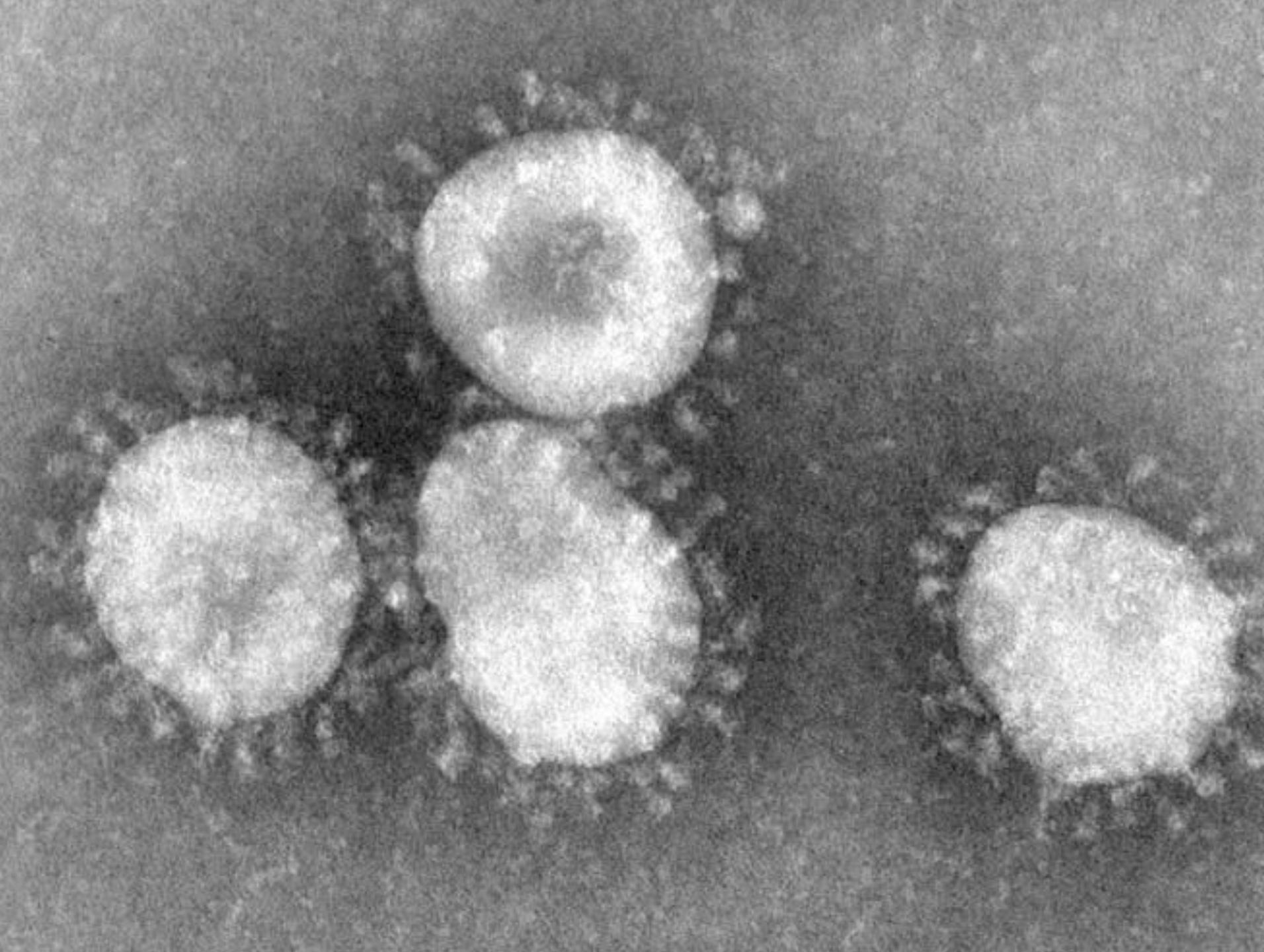


| | | |
|-----------------|----|----|
| giant centipede | No | 42 |
|-----------------|----|----|

Character-Based Phylogeny

Character-Based Phylogeny Problem: *Reconstruct a phylogeny from characters.*

- **Input:** An $n \times m$ character table for n species and m characters.
- **Output:** A tree in which species with similar character values occur near each other.

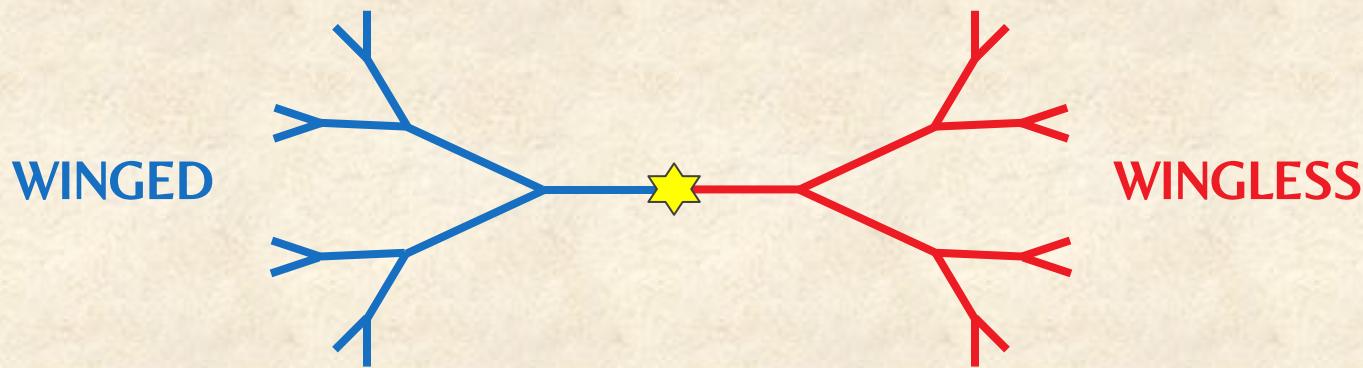


From Characters to a Phylogeny

STOP and Think: How would you construct an evolutionary tree from characters?

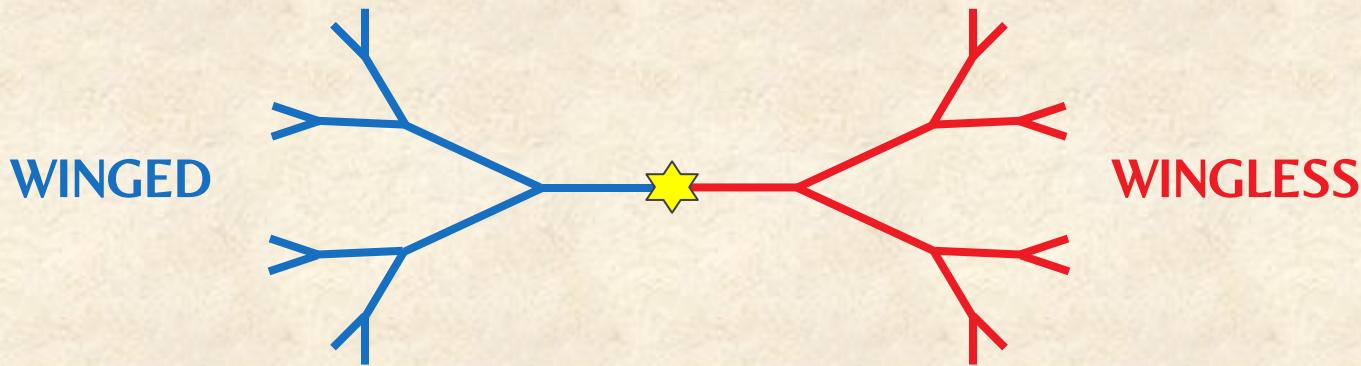
From Characters to a Phylogeny

STOP and Think: How would you construct an evolutionary tree from characters?



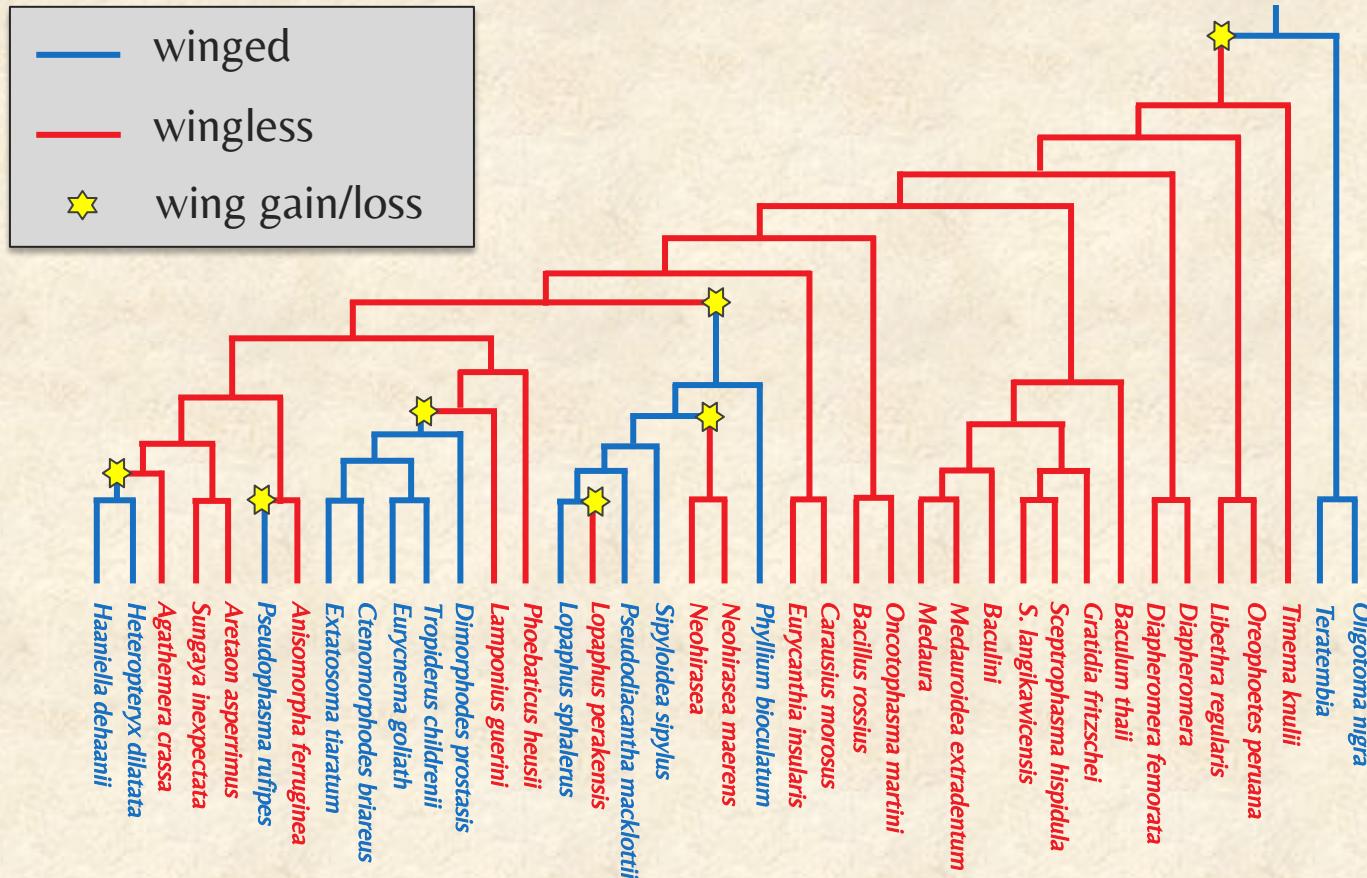
From Characters to a Phylogeny

STOP and Think: How would you construct an evolutionary tree from characters?

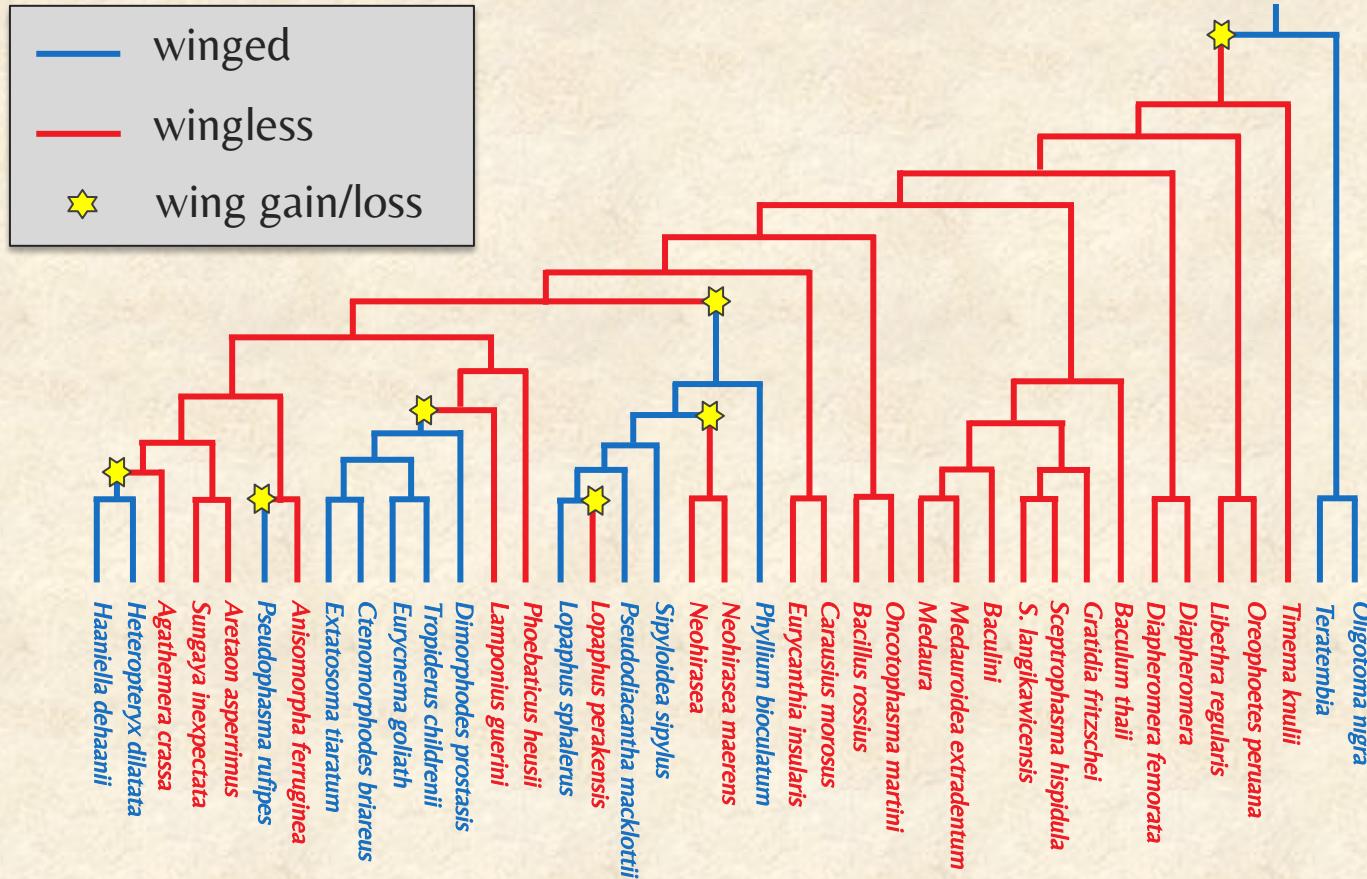


Dollo's principle of irreversibility (1893): evolution doesn't reinvent the same organ (e.g. insect wings).

Dollo's Principle Violated



Dollo's Principle Violated



STOP and Think: What do you think has happened?

An Alignment As a Character Table

| SPECIES | ALIGNMENT |
|---------|------------|
| Chimp | ACGTAGGCCT |
| Human | ATGTAAGACT |
| Seal | TCGAGAGCAC |
| Whale | TCGAAAGCAT |

An Alignment As a Character Table

| SPECIES | ALIGNMENT | |
|---------|------------|--|
| Chimp | ACGTAGGCCT | |
| Human | ATGTAAGACT | |
| Seal | TCGAGAGCAC | |
| Whale | TCGAAAGCAT | |

n species

m characters

Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- **The Small Parsimony Problem**
- The Large Parsimony Problem
- Evolutionary Tree Reconstruction in the Modern Era

Toward a Computational Problem

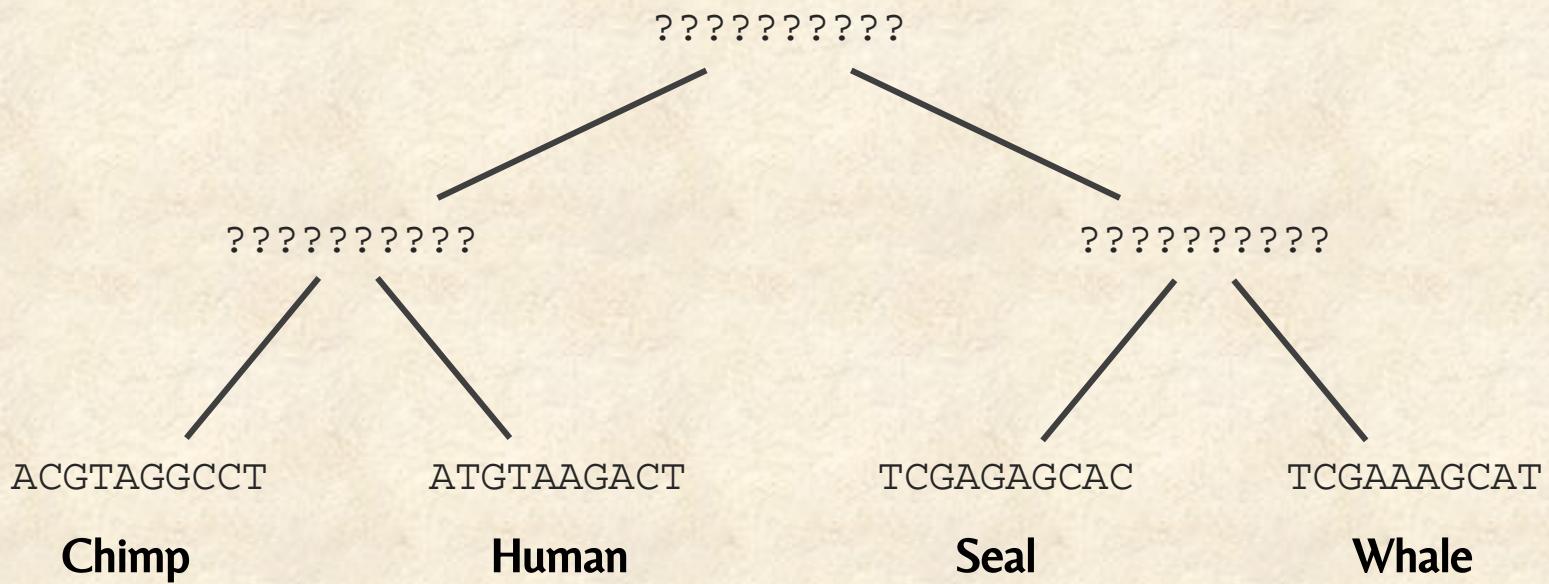
| | |
|-------|------------|
| Chimp | ACGTAGGCCT |
| Human | ATGTAAGACT |
| Seal | TCGAGAGCAC |
| Whale | TCGAAAGCAT |

n species

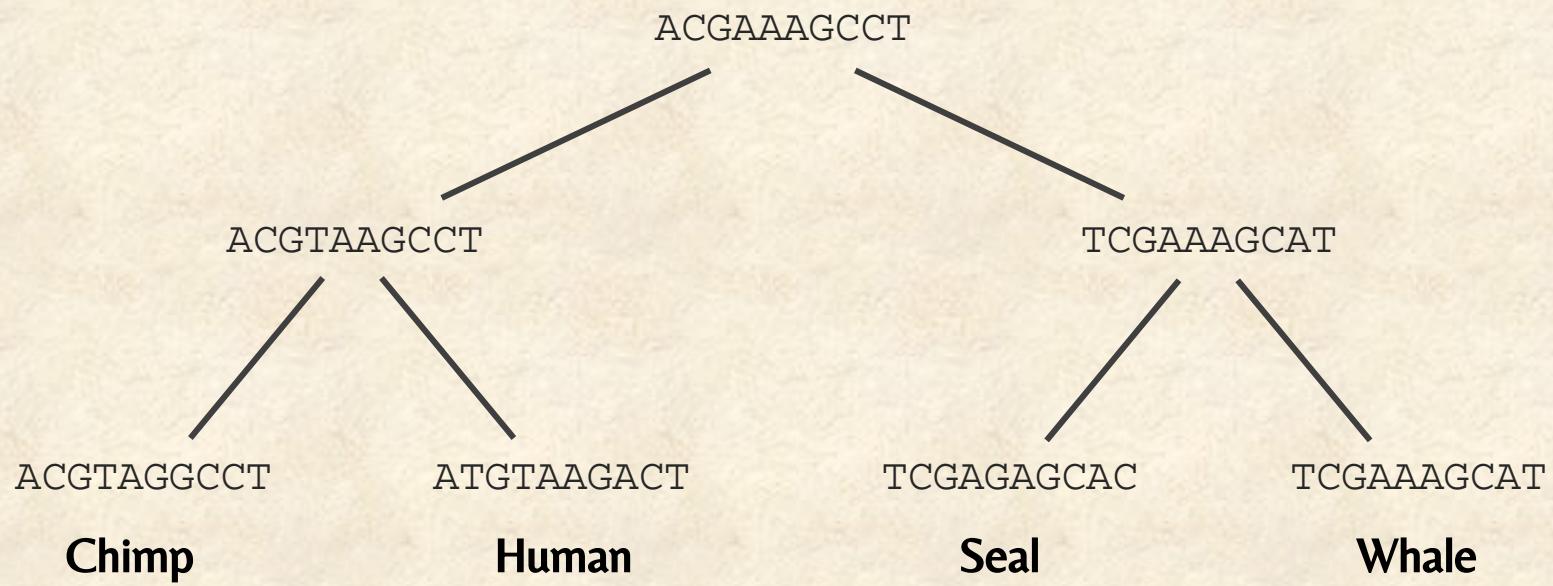
m characters

Toward a Computational Problem

| | |
|-------|------------|
| Chimp | ACGTAGGCCT |
| Human | ATGTAAGACT |
| Seal | TCGAGAGCAC |
| Whale | TCGAAAGCAT |

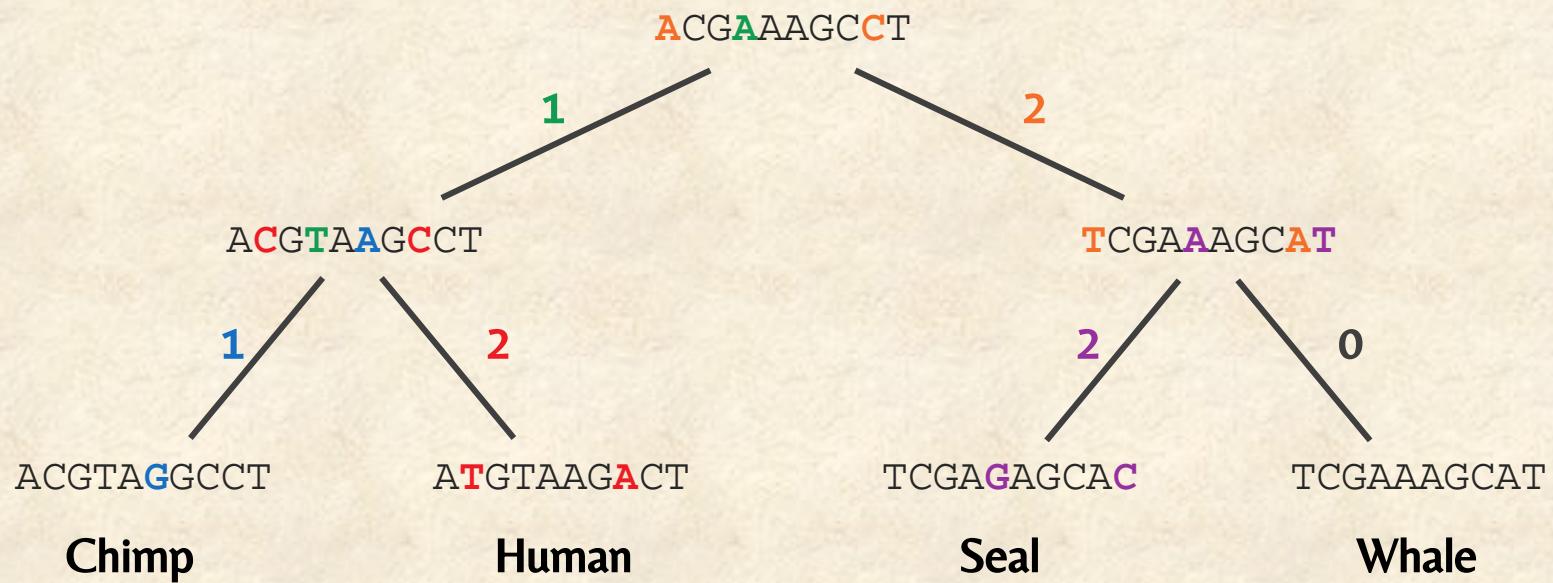


Toward a Computational Problem



Toward a Computational Problem

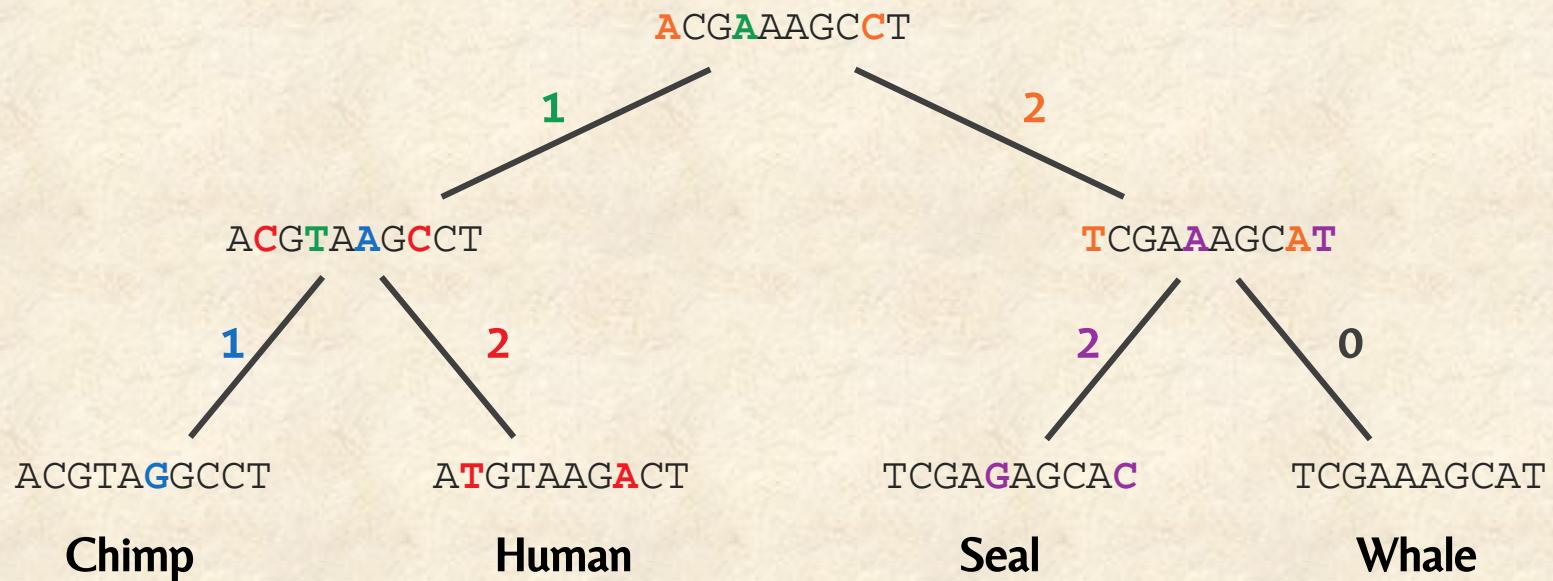
Parsimony score: sum of Hamming distances along each edge.



Toward a Computational Problem

Parsimony score: sum of Hamming distances along each edge.

Parsimony Score: 8



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.

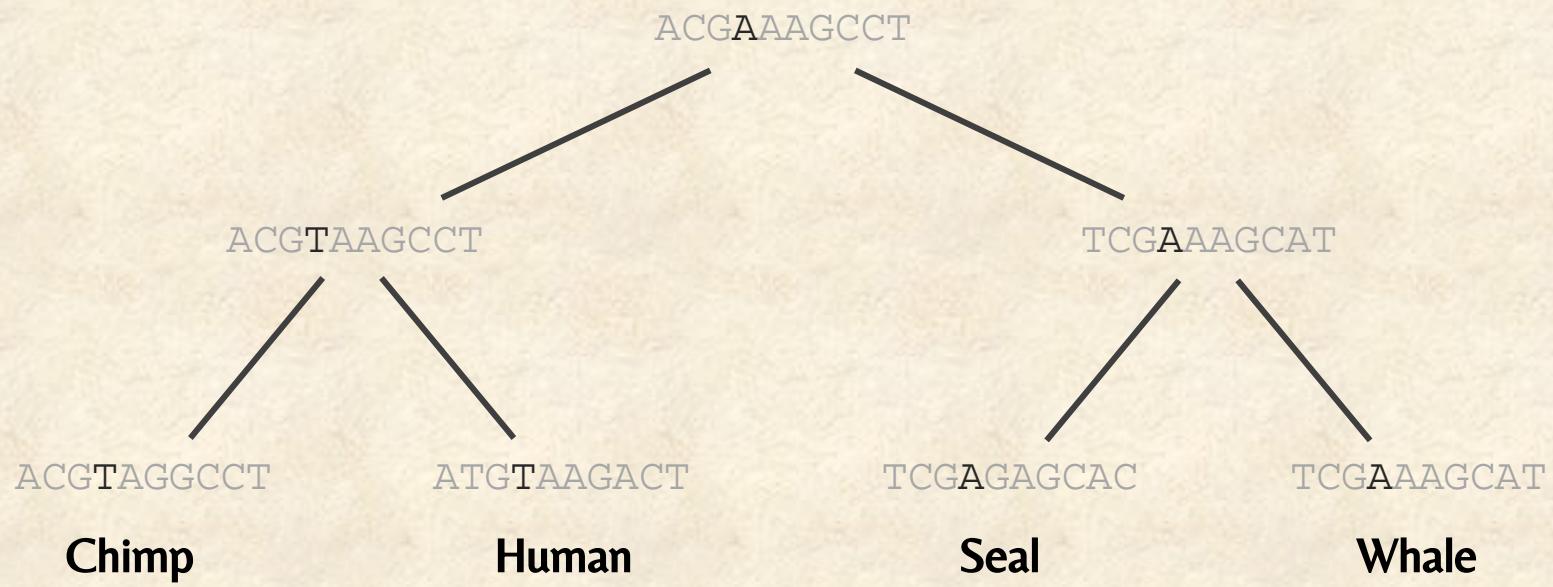
STOP and Think: Is there any way we can simplify this problem statement?

Toward a Computational Problem

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

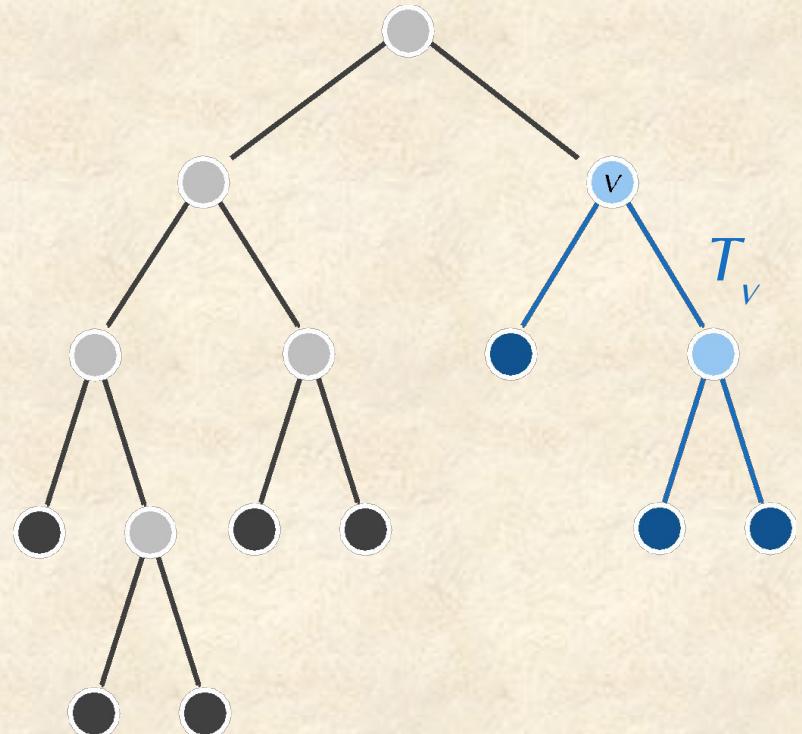
- **Input:** A rooted binary tree with each leaf labeled by a **single symbol**.
- **Output:** A labeling of all other nodes of the tree by **single symbols** that minimizes the tree's parsimony score.

Toward a Computational Problem



A Dynamic Programming Algorithm

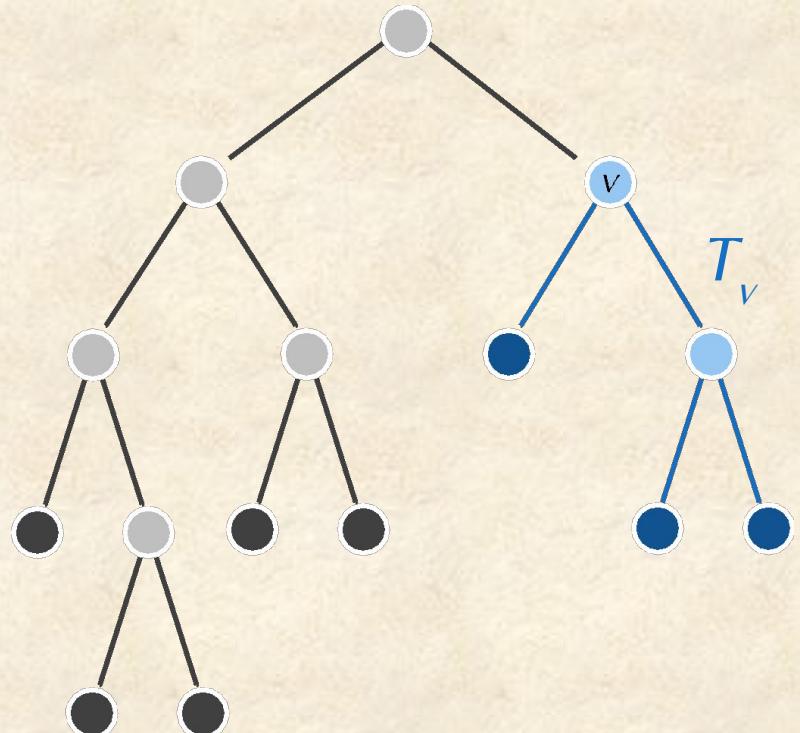
Let T_v denote the subtree of T whose root is v .



A Dynamic Programming Algorithm

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

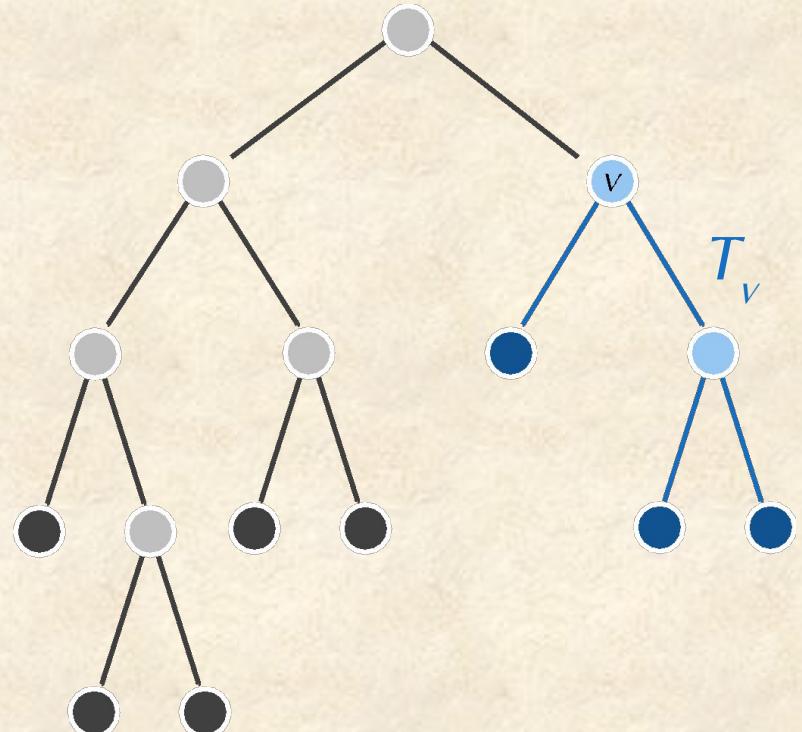
Define $s_k(v)$ as the minimum parsimony score of T_v over all labelings of T_v , assuming that v is labeled by k .



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 .

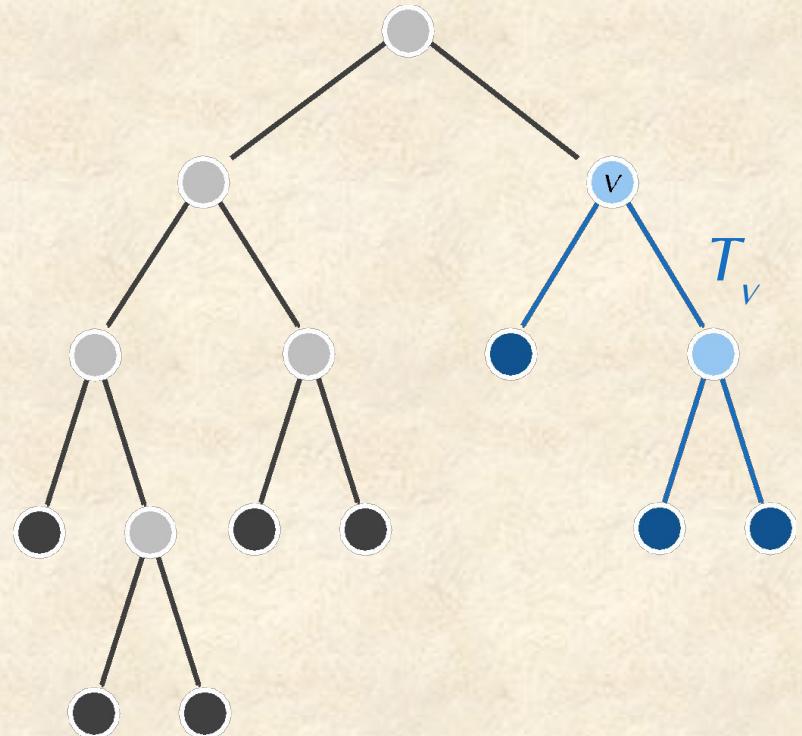


The minimum parsimony score for the tree is equal to the minimum value of $s_k(\text{root})$ over all symbols k .

A Dynamic Programming Algorithm

For symbols i and j , define

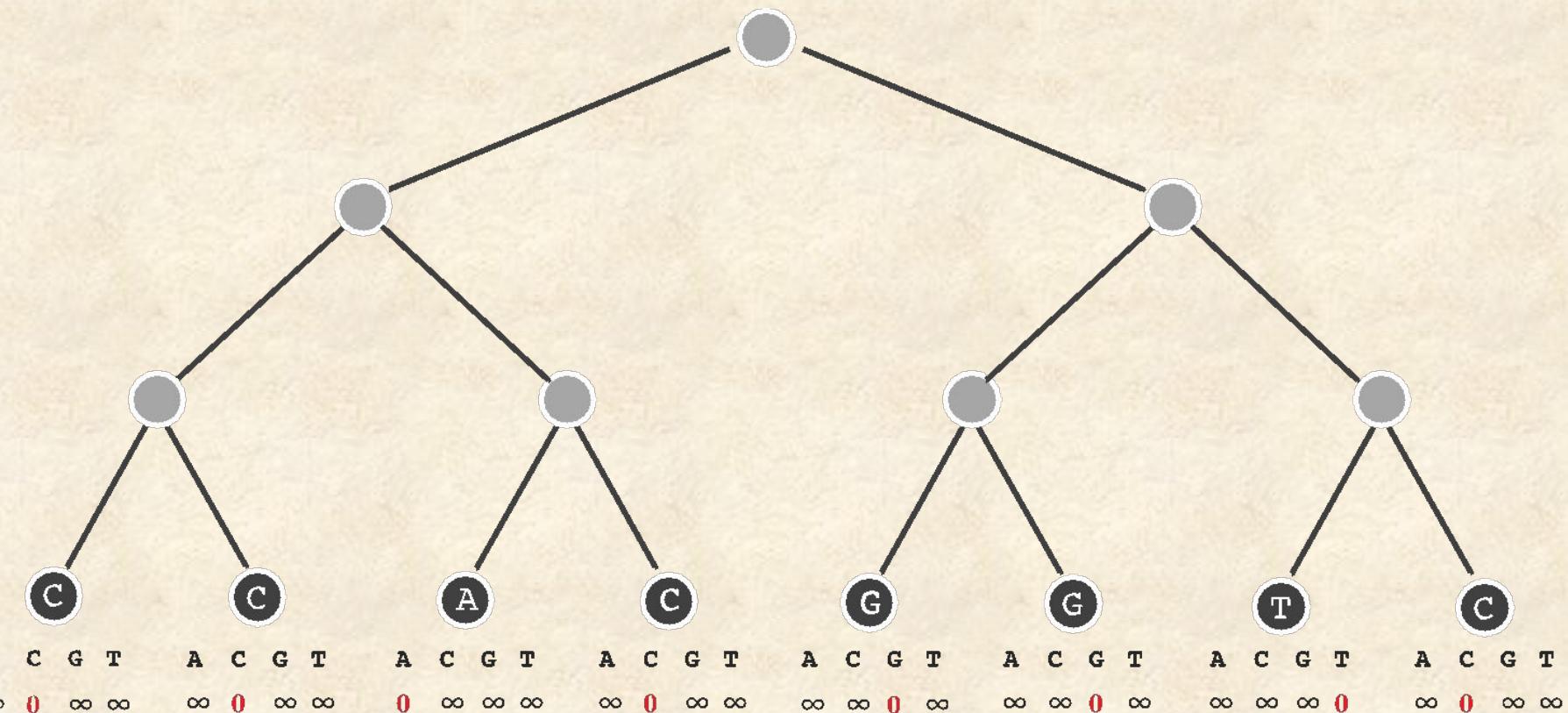
- $\Delta_{i,j} = 0$ if $i = j$
- $\Delta_{i,j} = 1$ otherwise.



Exercise Break: Prove the following recurrence relation:

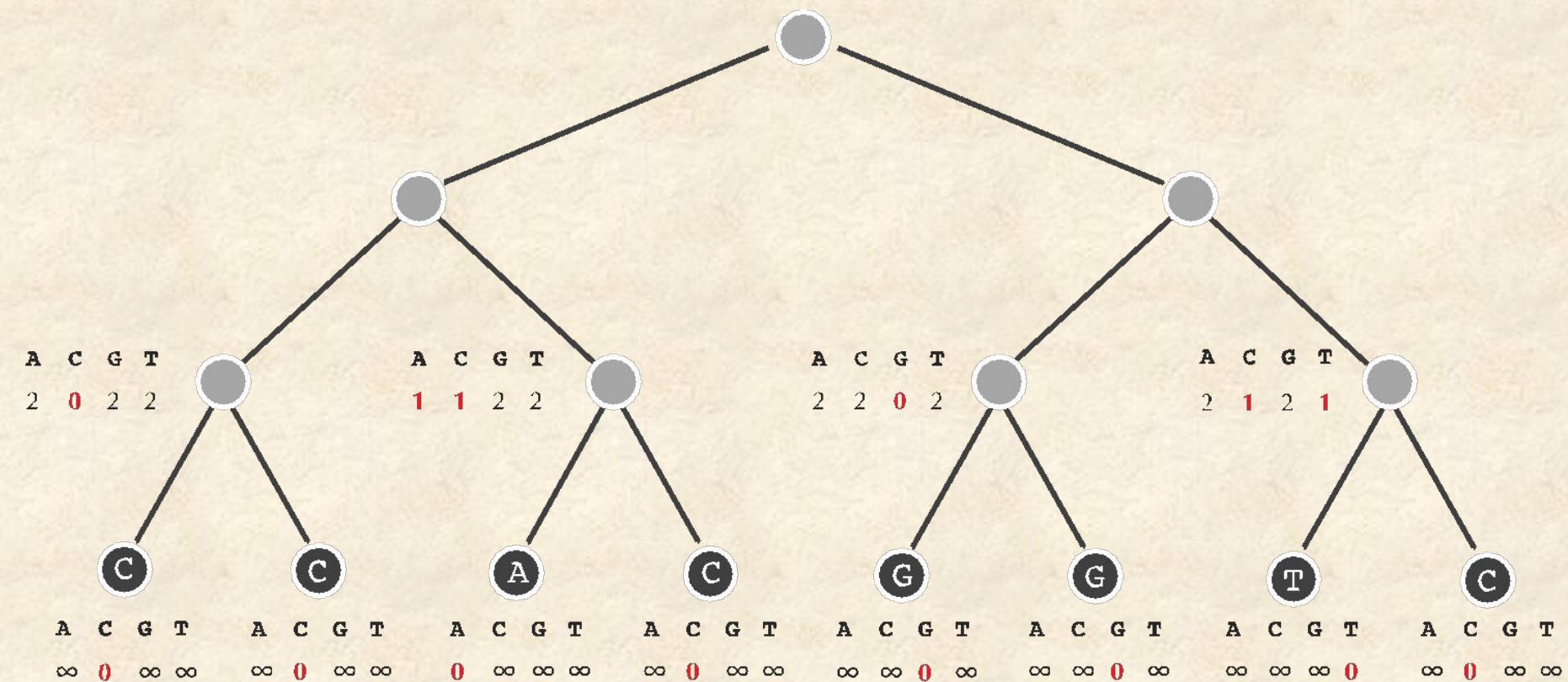
$$s_k(v) = \min_{\text{all symbols } i} \{s_i(\text{Daughter}(v)) + \Delta_{i,k}\} + \min_{\text{all symbols } i} \{s_i(\text{Son}(v)) + \Delta_{j,k}\}$$

A Dynamic Programming Algorithm



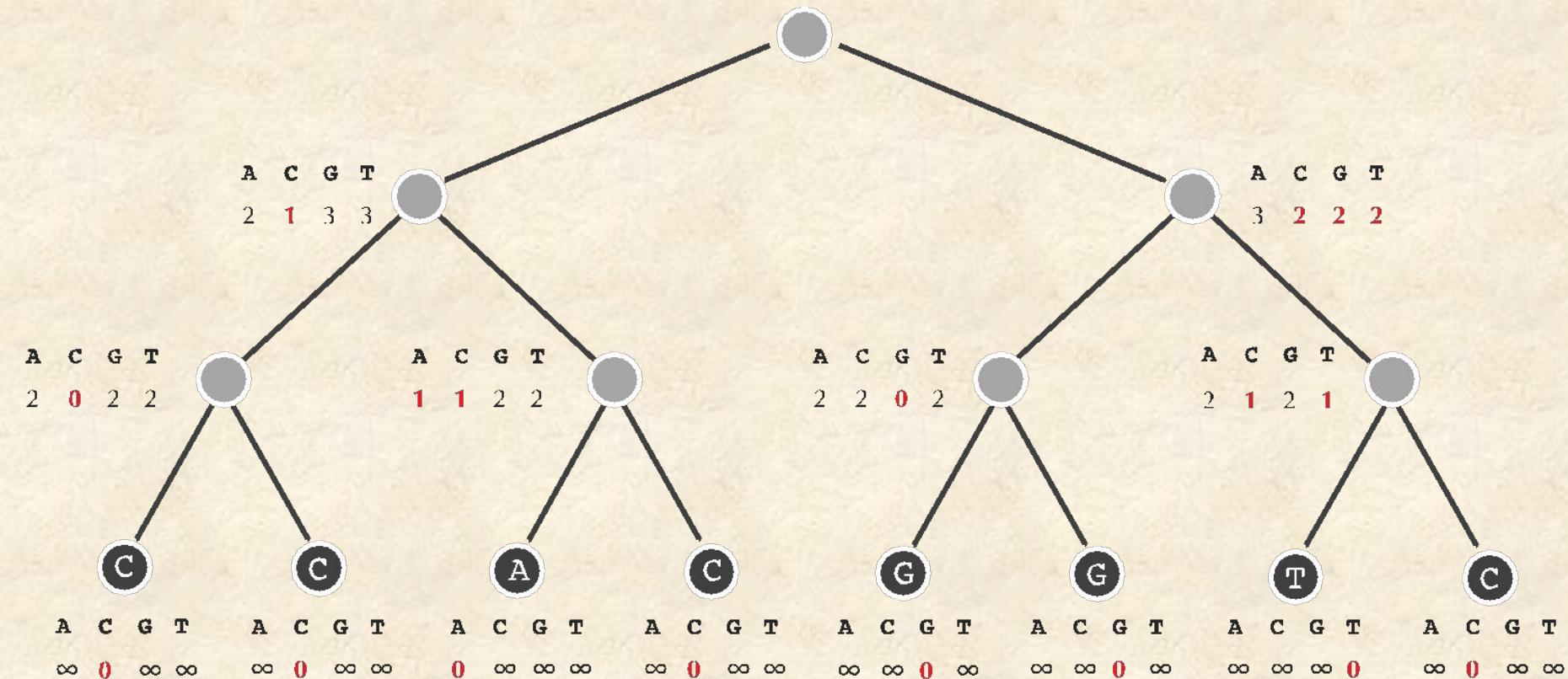
$$s_k(v) = \min_{\text{all symbols } i} \{ s_i(\text{Daughter}(v)) + \Delta_{i,k} \} + \min_{\text{all symbols } i} \{ s_i(\text{Son}(v)) + \Delta_{j,k} \}$$

A Dynamic Programming Algorithm



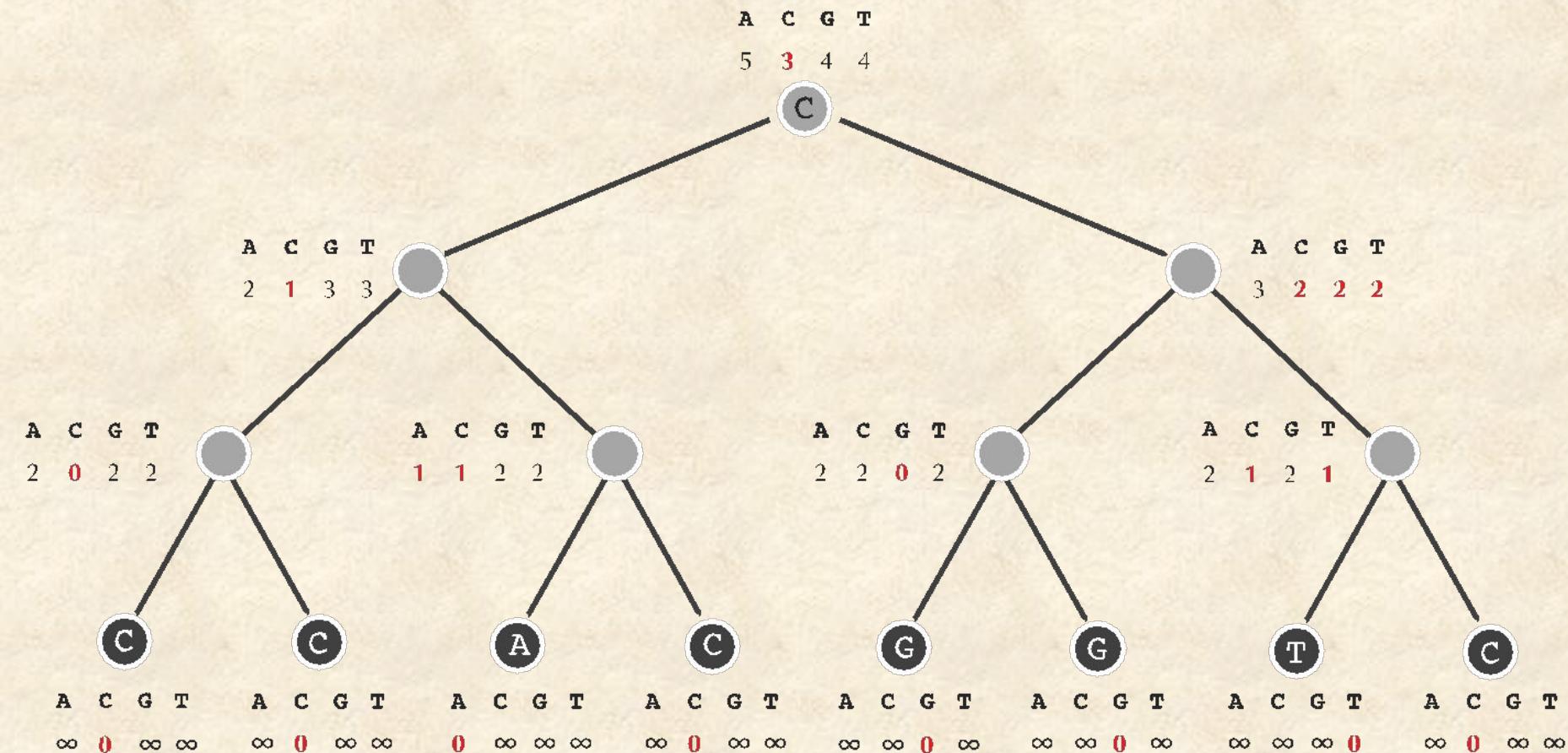
$$s_k(v) = \min_{\text{all symbols } i} \{ s_i(\text{Daughter}(v)) + \Delta_{i,k} \} + \min_{\text{all symbols } i} \{ s_i(\text{Son}(v)) + \Delta_{j,k} \}$$

A Dynamic Programming Algorithm



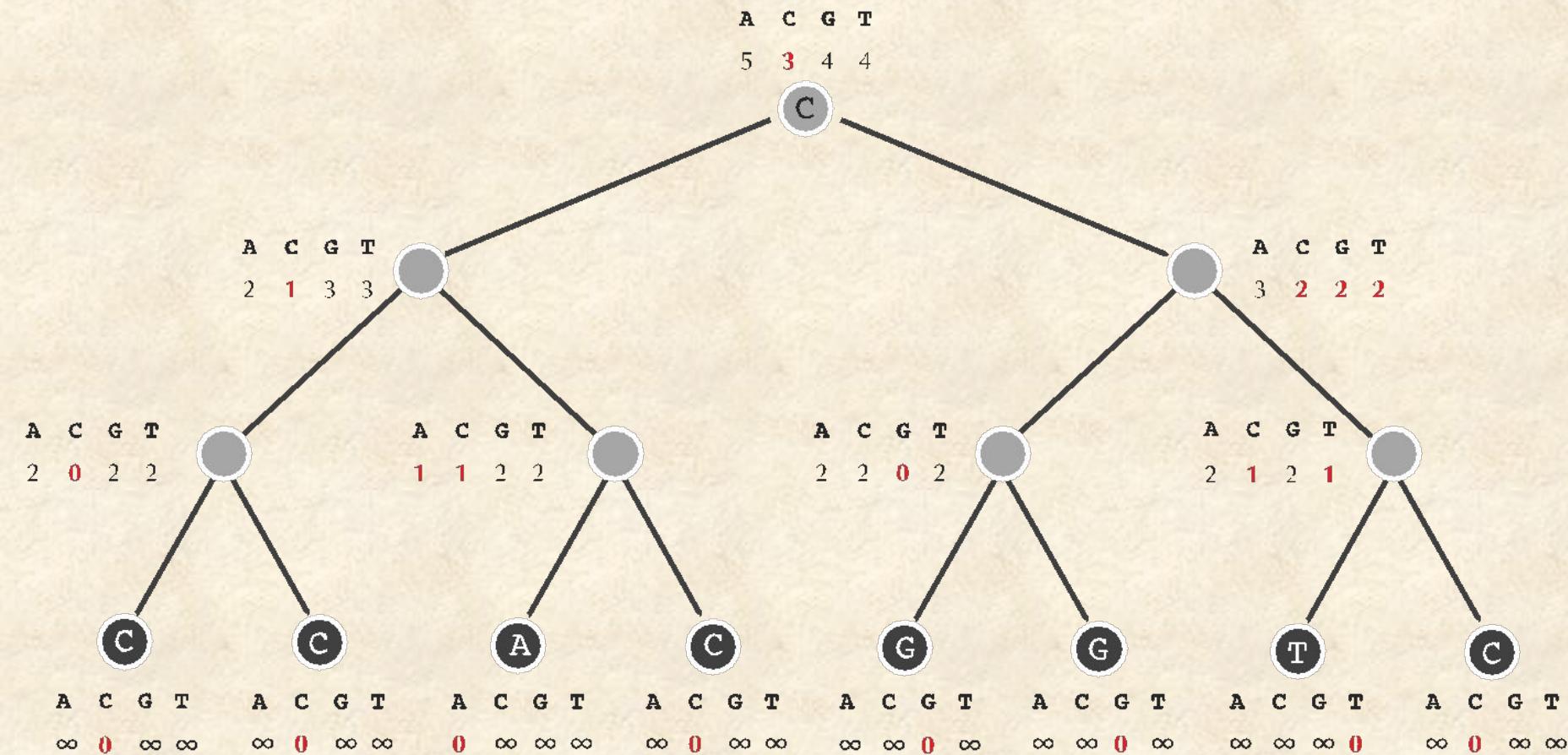
$$s_k(v) = \min_{\text{all symbols } i} \{ s_i(\text{Daughter}(v)) + \Delta_{i,k} \} + \min_{\text{all symbols } i} \{ s_i(\text{Son}(v)) + \Delta_{j,k} \}$$

A Dynamic Programming Algorithm



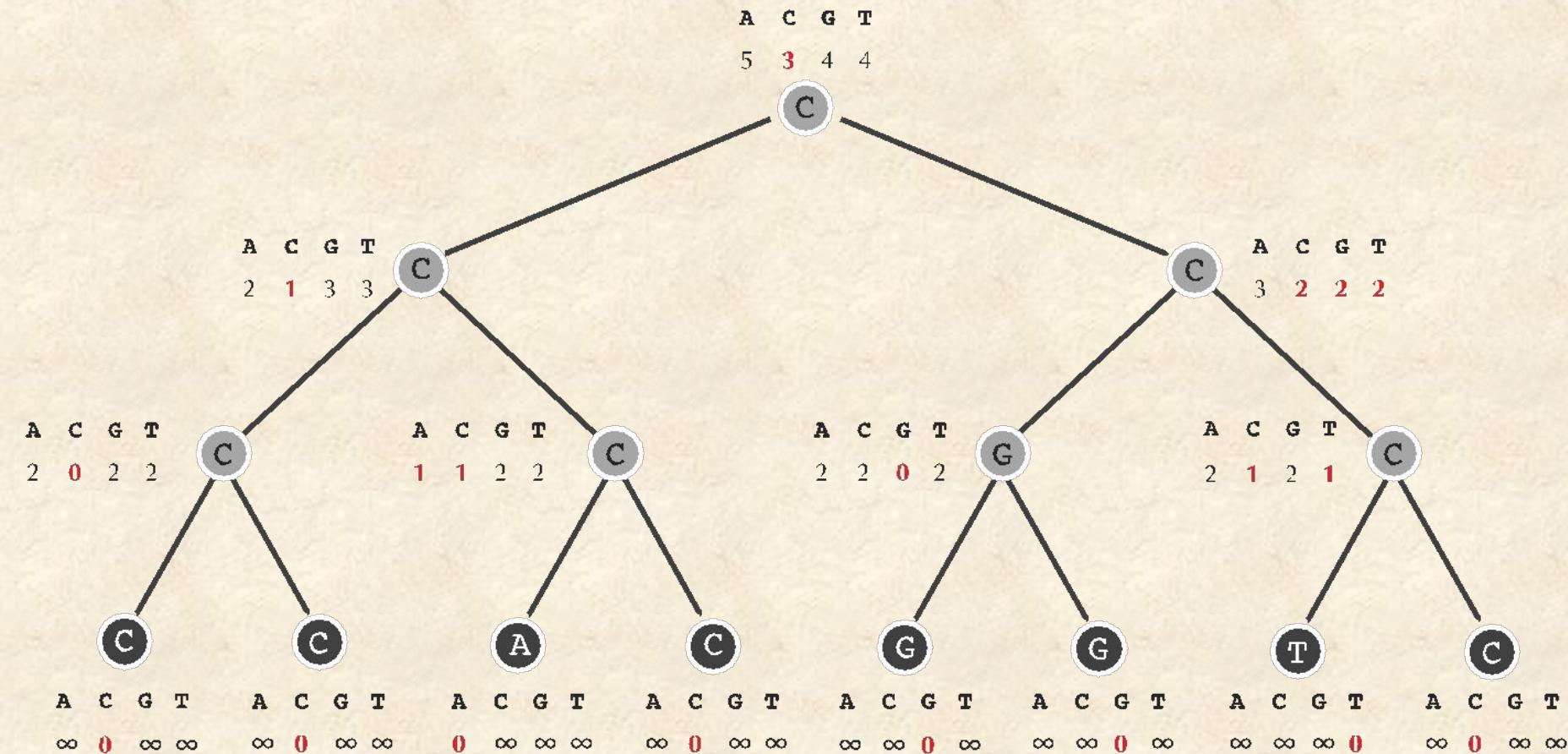
$$s_k(v) = \min_{\text{all symbols } i} \{ s_i(\text{Daughter}(v)) + \Delta_{i,k} \} + \min_{\text{all symbols } i} \{ s_i(\text{Son}(v)) + \Delta_{j,k} \}$$

A Dynamic Programming Algorithm

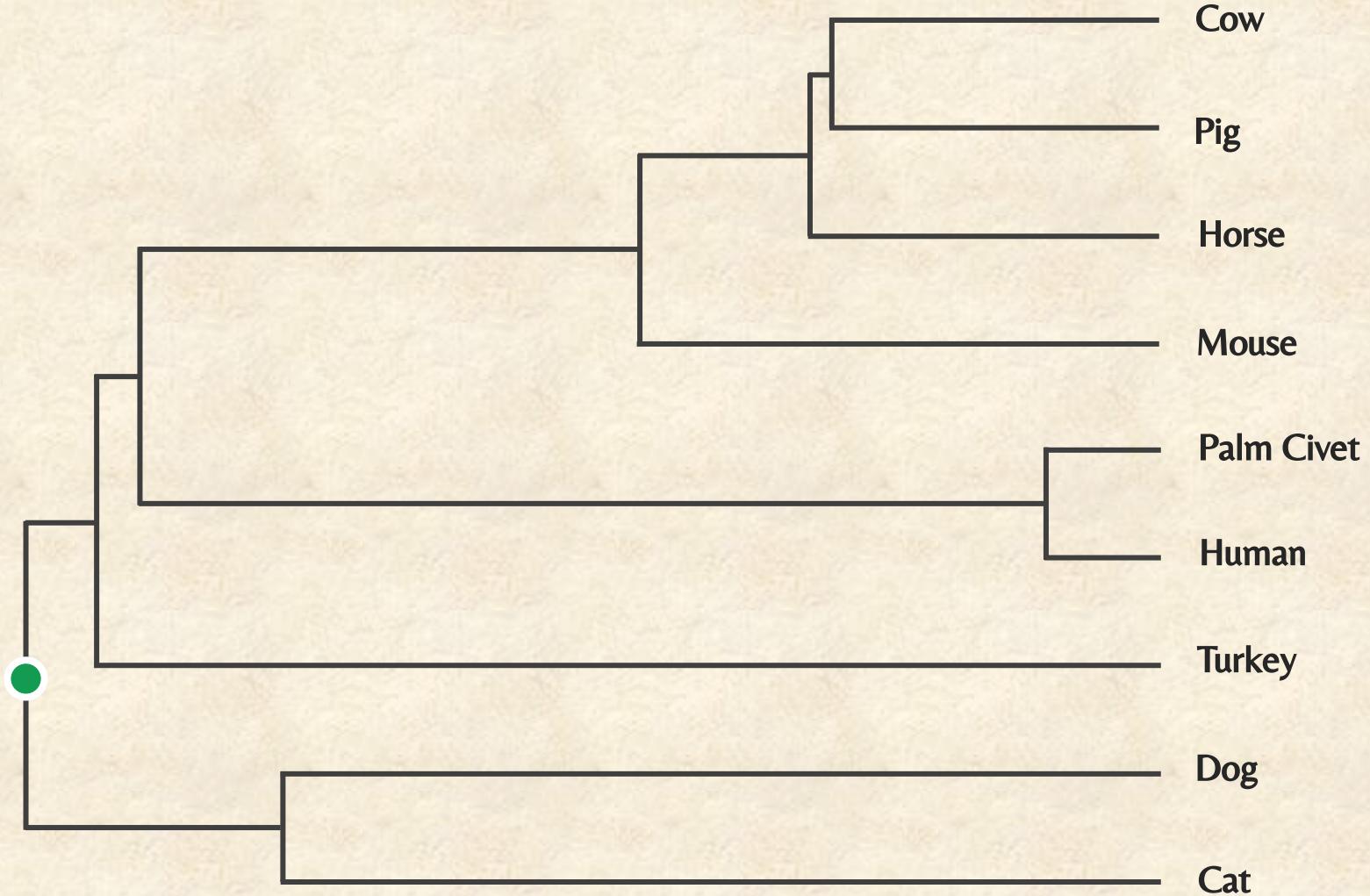


Exercise Break: “Backtrack” to fill in the remaining nodes of the tree.

A Dynamic Programming Algorithm



Code Challenge: Solve the Small Parsimony Problem.



Exercise Break: Apply **SmallParsimony** to this tree to reconstruct ancestral coronavirus sequences.

Small Parsimony for Unrooted Trees

Small Parsimony in an Unrooted Tree Problem: *Find the most parsimonious labeling of the internal nodes of an unrooted tree.*

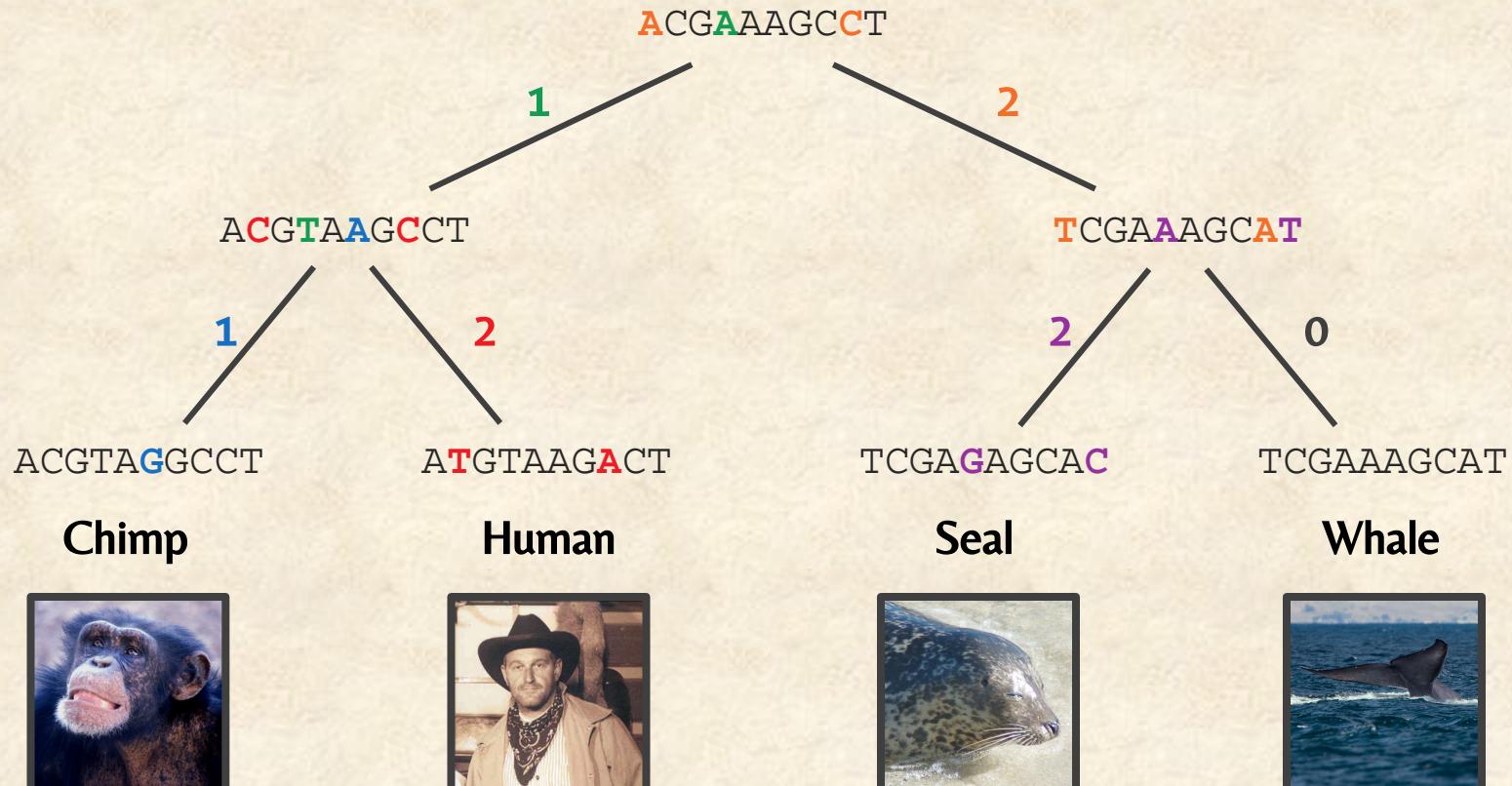
- **Input:** An unrooted binary tree with each leaf labeled by a string of length m .
- **Output:** A position of the root and a labeling of all other nodes of the tree by strings of length m that minimizes the tree's parsimony score.

Code Challenge: Solve this problem.

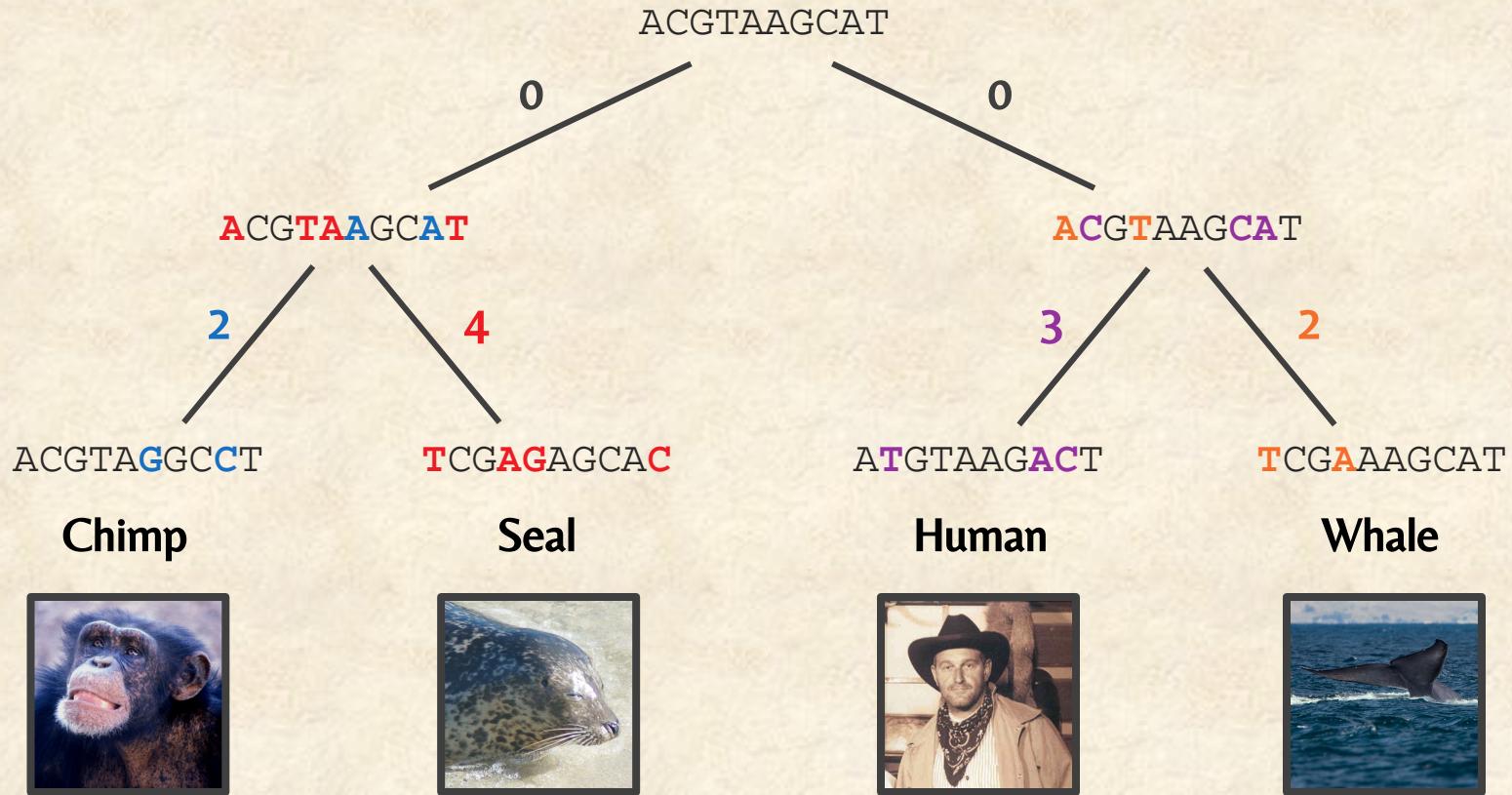
Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- **The Large Parsimony Problem**
- Evolutionary Tree Reconstruction in the Modern Era

Finding the Most Parsimonious Tree

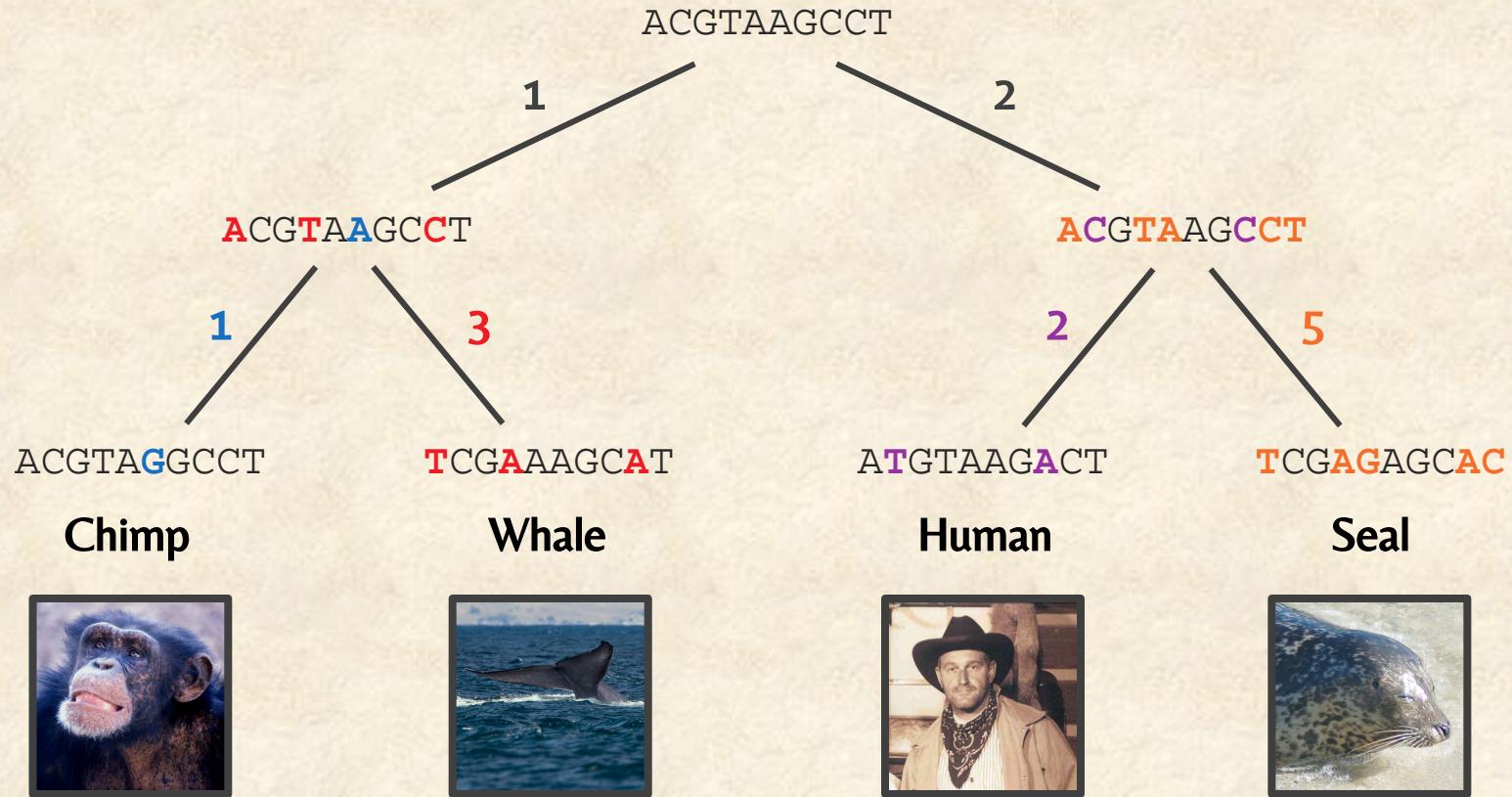


Finding the Most Parsimonious Tree



Parsimony Score: 11

Finding the Most Parsimonious Tree



Parsimony Score: 14

Finding the Most Parsimonious Tree

Large Parsimony Problem: *Given a set of strings, find a tree (with leaves labeled by all these strings) having minimum parsimony score.*

- **Input:** A collection of strings of equal length.
- **Output:** A rooted binary tree T that minimizes the parsimony score among all possible rooted binary trees with leaves labeled by these strings.

Finding the Most Parsimonious Tree

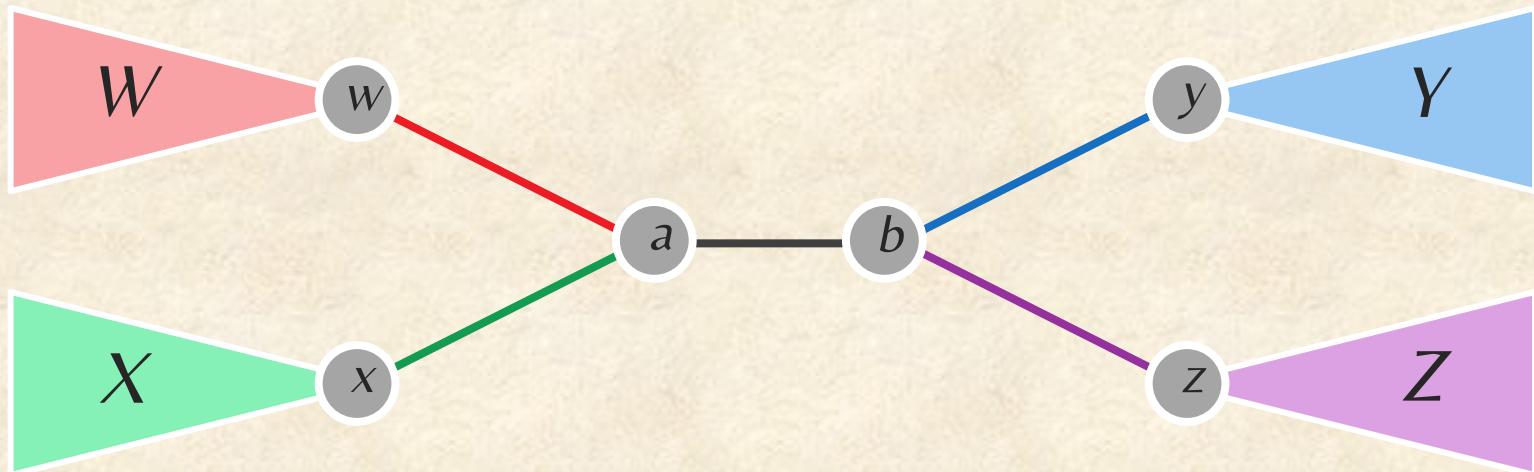
Large Parsimony Problem: *Given a set of strings, find a tree (with leaves labeled by all these strings) having minimum parsimony score.*

- **Input:** A collection of strings of equal length.
- **Output:** A rooted binary tree T that minimizes the parsimony score among all possible rooted binary trees with leaves labeled by these strings.

Unfortunately, this problem is NP -Complete...

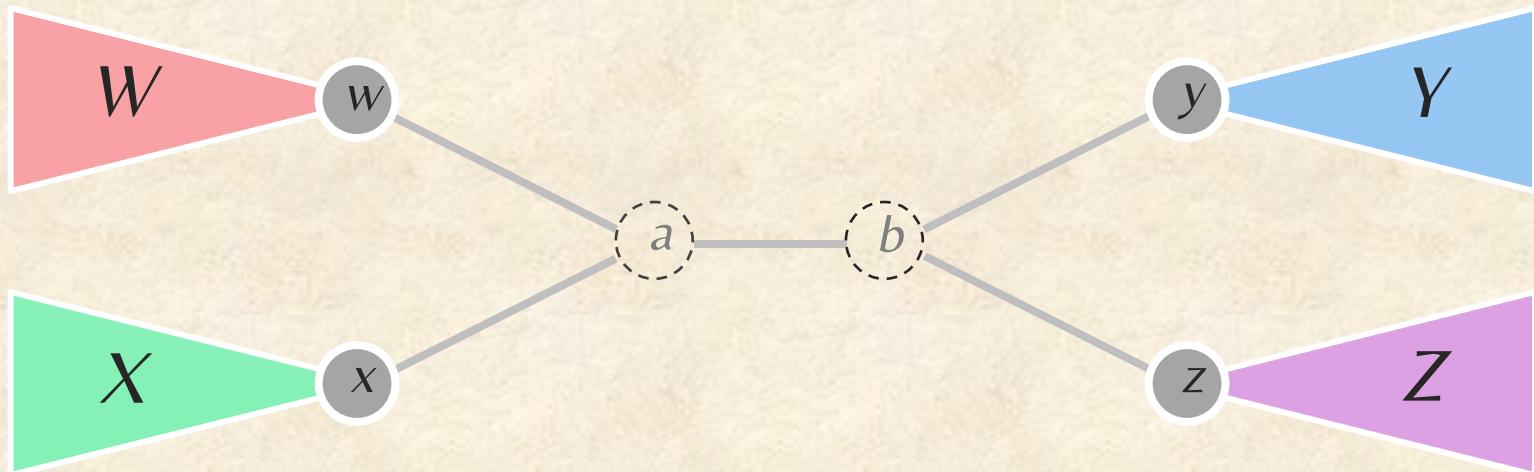
A Greedy Heuristic for Large Parsimony

Note that removing an **internal edge**, an edge connecting two internal nodes (along with the nodes), produces four subtrees (W, X, Y, Z).



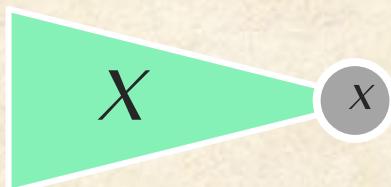
A Greedy Heuristic for Large Parsimony

Note that removing an **internal edge**, an edge connecting two internal nodes (along with the nodes), produces four subtrees (W, X, Y, Z).



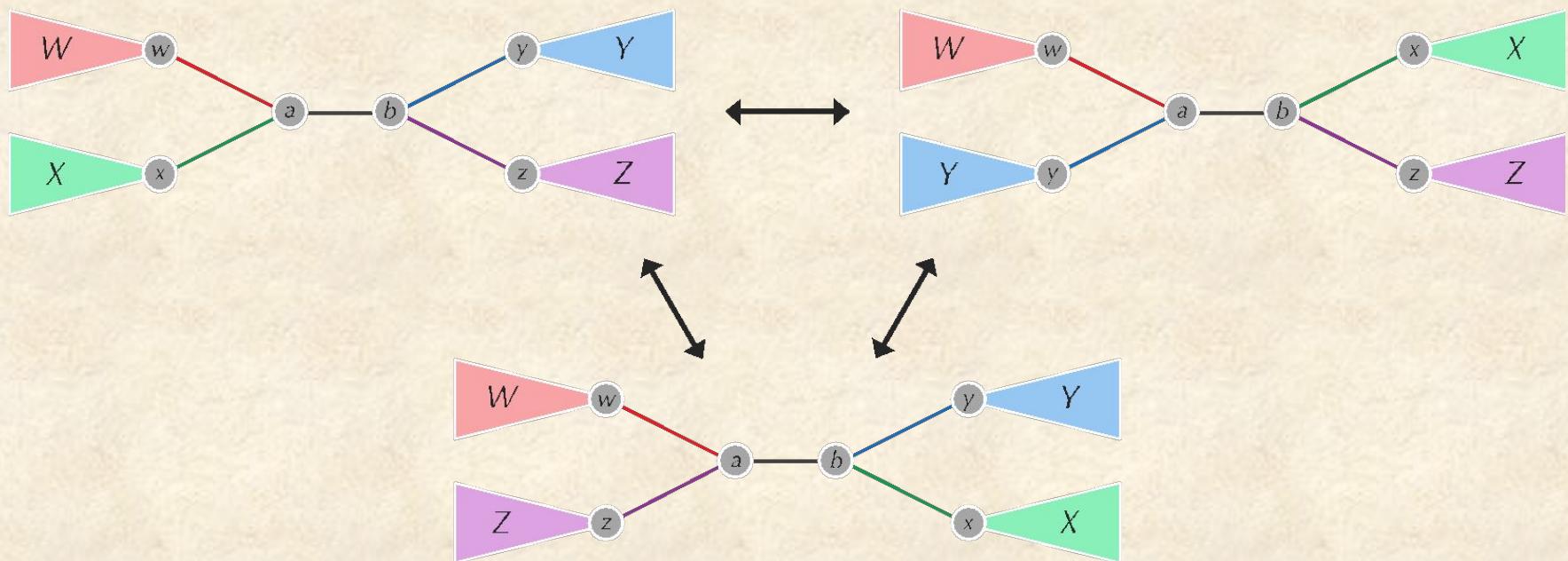
A Greedy Heuristic for Large Parsimony

Note that removing an **internal edge**, an edge connecting two internal nodes (along with the nodes), produces four subtrees (W, X, Y, Z).



A Greedy Heuristic for Large Parsimony

Rearranging these subtrees is called a **nearest neighbor interchange**.



A Greedy Heuristic for Large Parsimony

Nearest Neighbors of a Tree Problem: *Given an edge in a binary tree, generate the two neighbors of this tree.*

- **Input:** An internal edge in a binary tree.
- **Output:** The two nearest neighbors of this tree (for the given internal edge).

Code Challenge: Solve this problem.

A Greedy Heuristic for Large Parsimony

Nearest Neighbor Interchange Heuristic:

1. Set current tree equal to arbitrary binary rooted tree structure.

A Greedy Heuristic for Large Parsimony

Nearest Neighbor Interchange Heuristic:

1. Set current tree equal to arbitrary binary rooted tree structure.
2. Go through all internal edges and perform all possible nearest neighbor interchanges.

A Greedy Heuristic for Large Parsimony

Nearest Neighbor Interchange Heuristic:

1. Set current tree equal to arbitrary binary rooted tree structure.
2. Go through all internal edges and perform all possible nearest neighbor interchanges.
3. Solve Small Parsimony Problem on each tree.

A Greedy Heuristic for Large Parsimony

Nearest Neighbor Interchange Heuristic:

1. Set current tree equal to arbitrary binary rooted tree structure.
2. Go through all internal edges and perform all possible nearest neighbor interchanges.
3. Solve Small Parsimony Problem on each tree.
4. If any tree has parsimony score improving over optimal tree, set it equal to the current tree.
Otherwise, return current tree.

A Greedy Heuristic for Large Parsimony

Nearest Neighbor Interchange Heuristic:

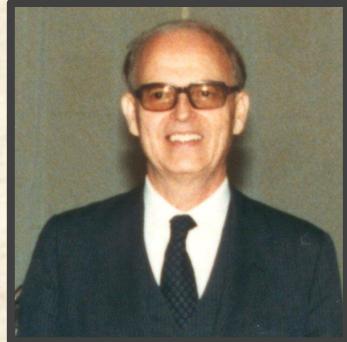
1. Set current tree equal to arbitrary binary rooted tree structure.
2. Go through all internal edges and perform all possible nearest neighbor interchanges.
3. Solve Small Parsimony Problem on each tree.
4. If any tree has parsimony score improving over optimal tree, set it equal to the current tree.
Otherwise, return current tree.

Code Challenge: Implement the nearest-neighbor interchange heuristic.

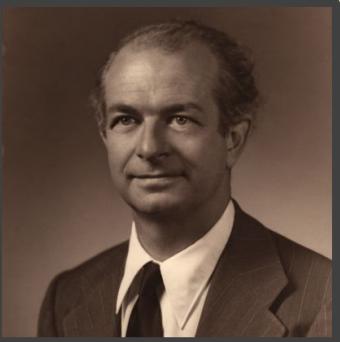
Outline

- The Fastest Outbreak
- Transforming Distance Matrices into Evolutionary Trees
- Toward an Algorithm for Distance-Based Phylogeny Construction
- Additive Phylogeny
- Using Least-Squares to Construct Distance-Based Phylogenies
- Ultrametric Evolutionary Trees
- The Neighbor-Joining Algorithm
- Character-Based Tree Reconstruction
- The Small Parsimony Problem
- The Large Parsimony Problem
- **Evolutionary Tree Reconstruction in the Modern Era**

1963: Paradigm Shift in Genetic Analysis



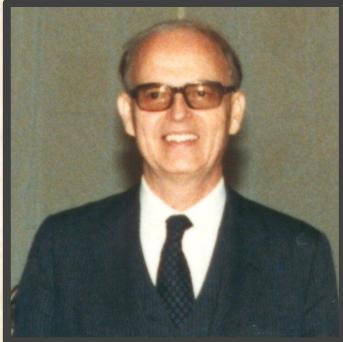
Emile
Zuckerkandl



Linus
Pauling

From the point of view of hemoglobin structure, it appears that gorilla is just an abnormal human.

1963: Paradigm Shift in Genetic Analysis

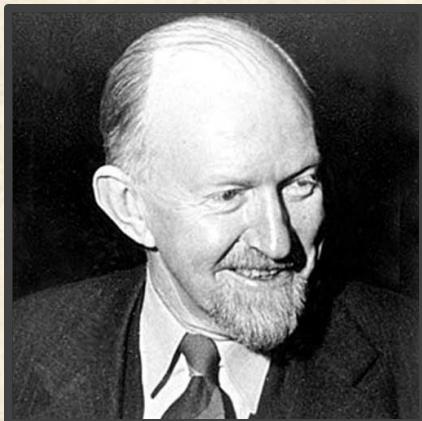


Emile
Zuckerkandl



Linus
Pauling

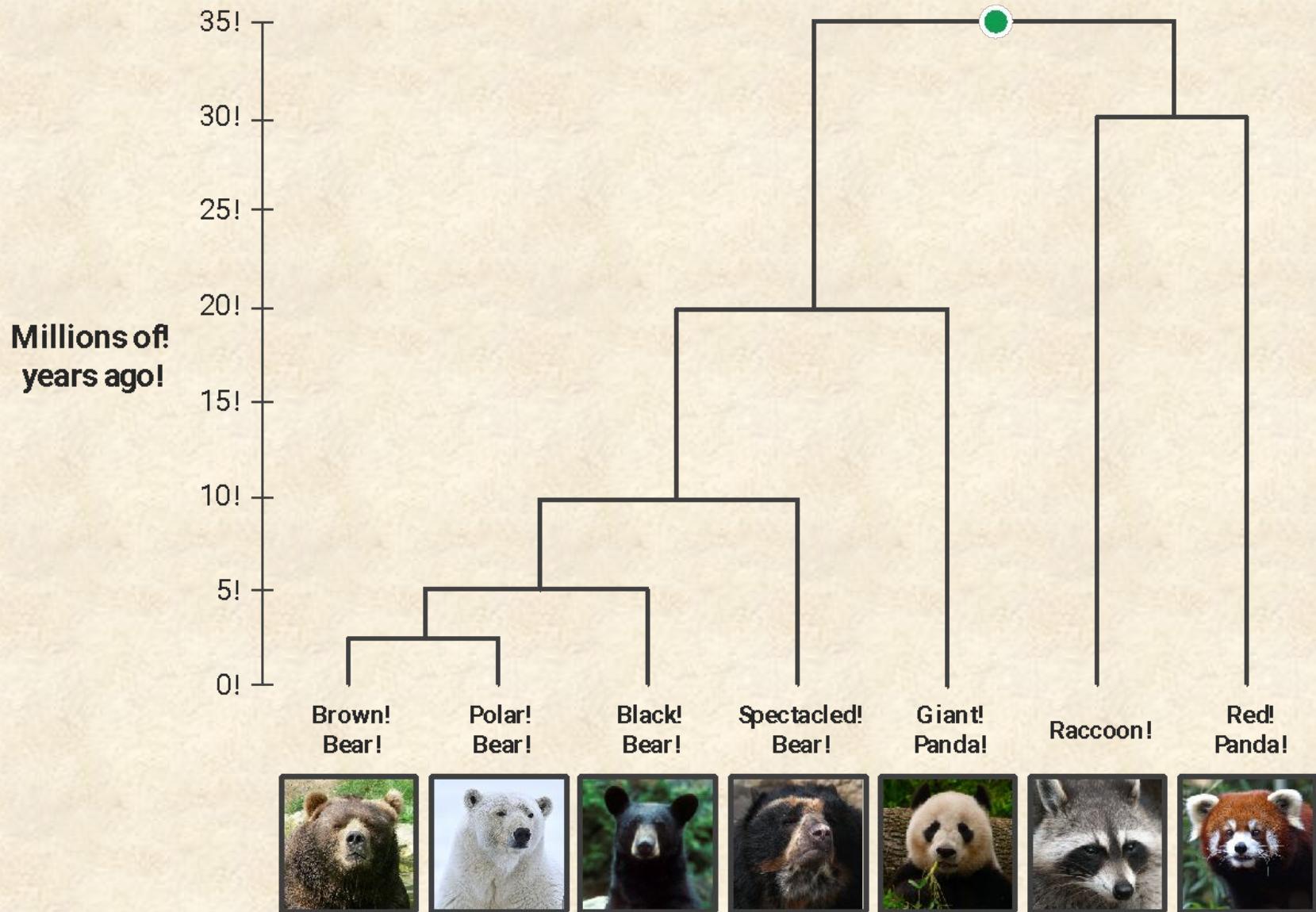
From the point of view of hemoglobin structure, it appears that gorilla is just an abnormal human.



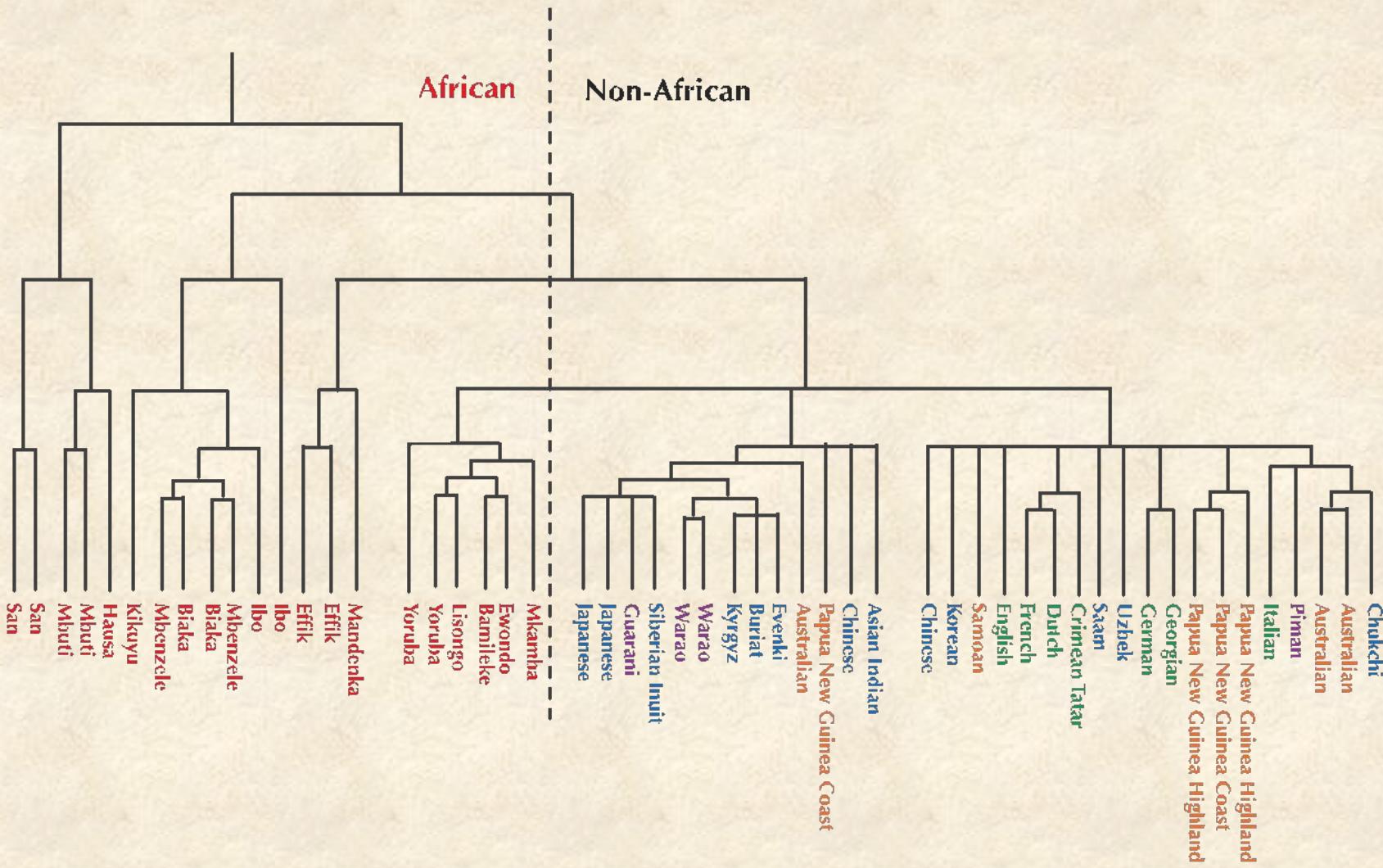
Gaylord Simpson

...that is of course nonsense. What the comparison really indicates is that hemoglobin is a bad choice and has nothing to tell us about attributes, or indeed tells us a lie.

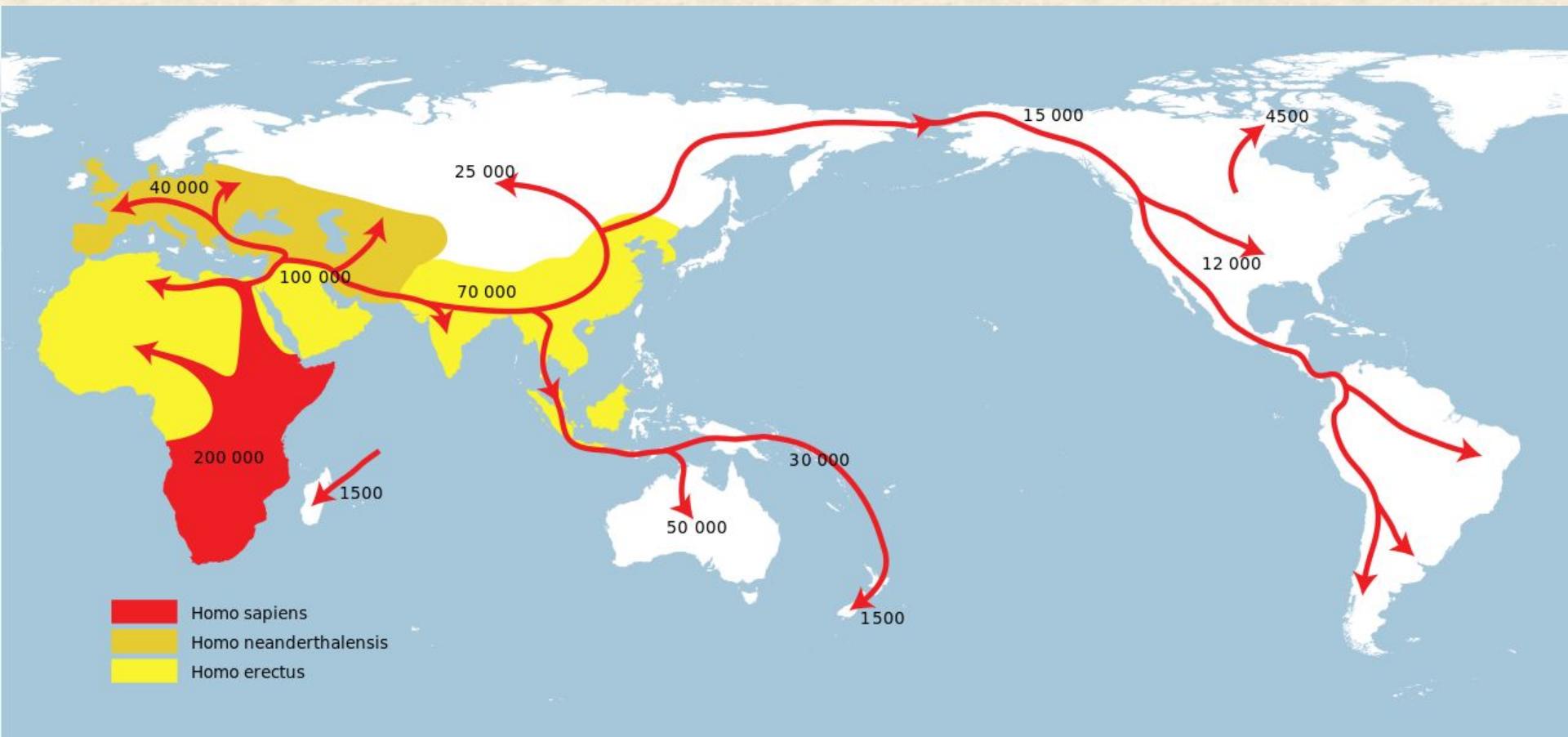
1985: Panda's Pedigree Decoded



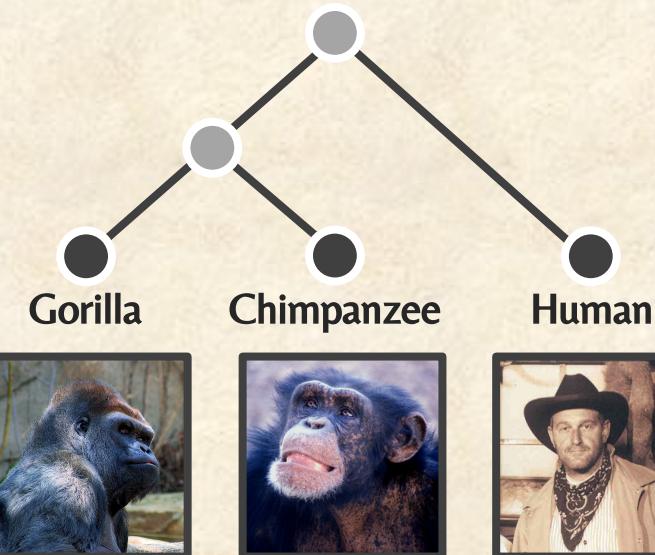
1987: Tracing Human Origins



1987: Tracing Human Origins

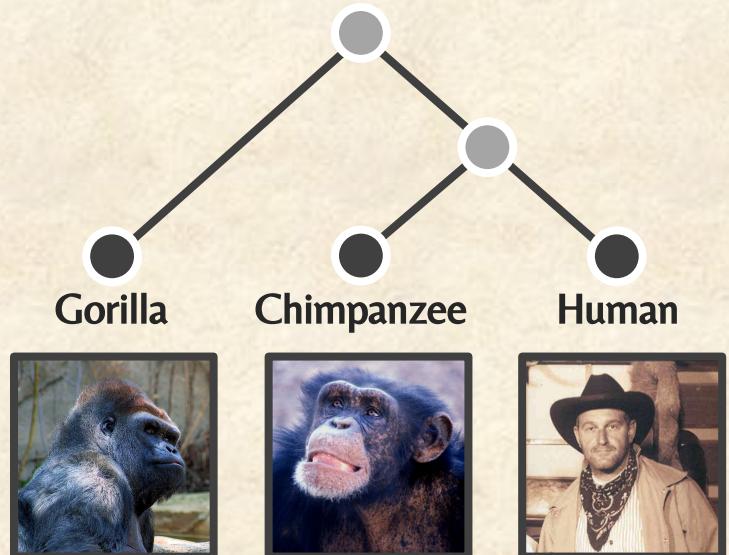


1996: Human-Chimp-Gorilla Ancestry



dopamine D4 receptor

vs.

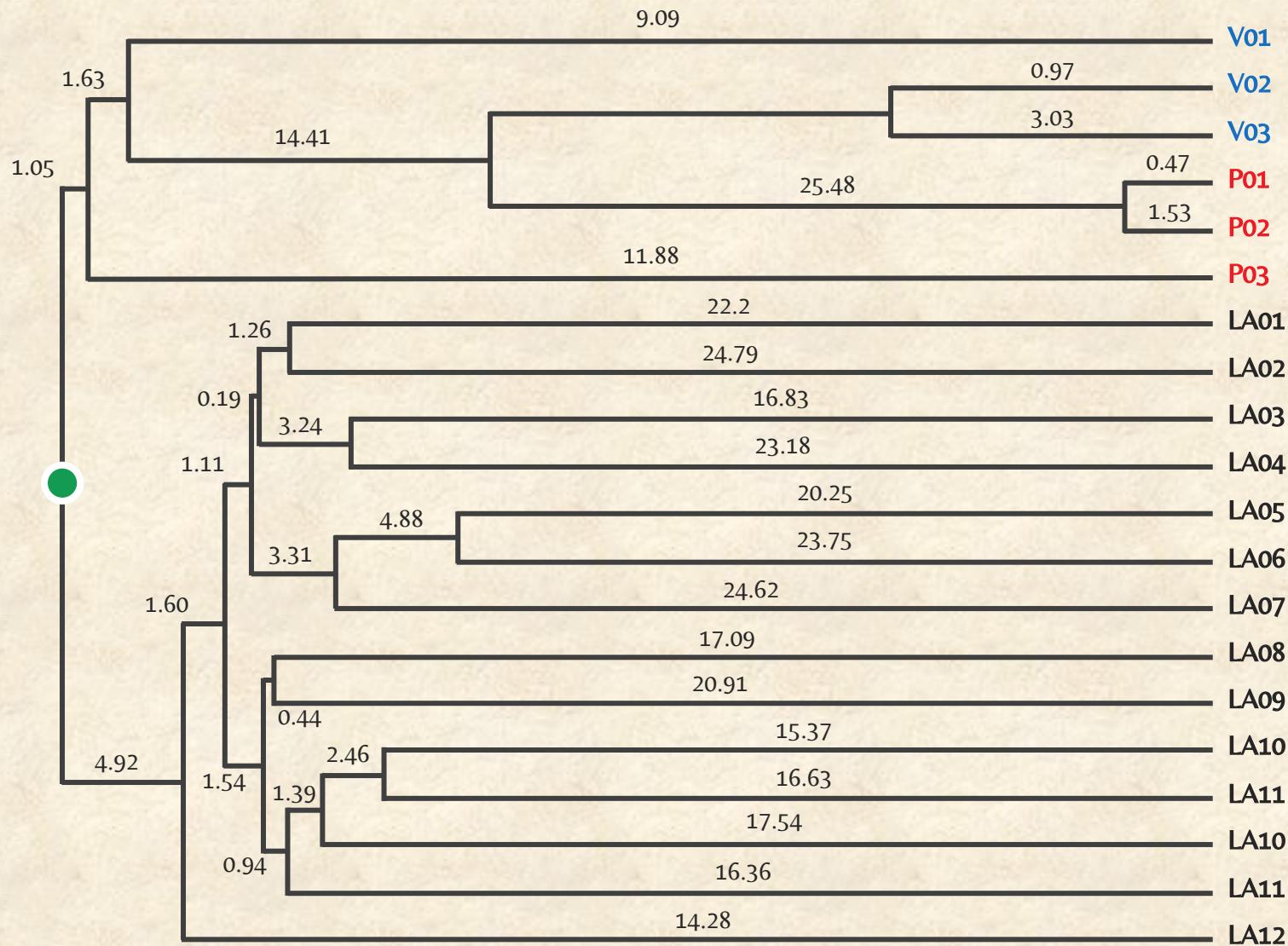


beta-globin

2000s: Humans and Mice Are Relatives



1998: Evolutionary Trees Fight Crime



1988: Evolutionary Trees Fight Crime

Challenge Problem: Given HIV sequences from AIDS patients in Lafayette, construct the evolutionary tree for each of nine HIV proteins. Does each tree support conviction? Reconstruct the ancestral HIV sequences at the internal nodes of the resulting trees.