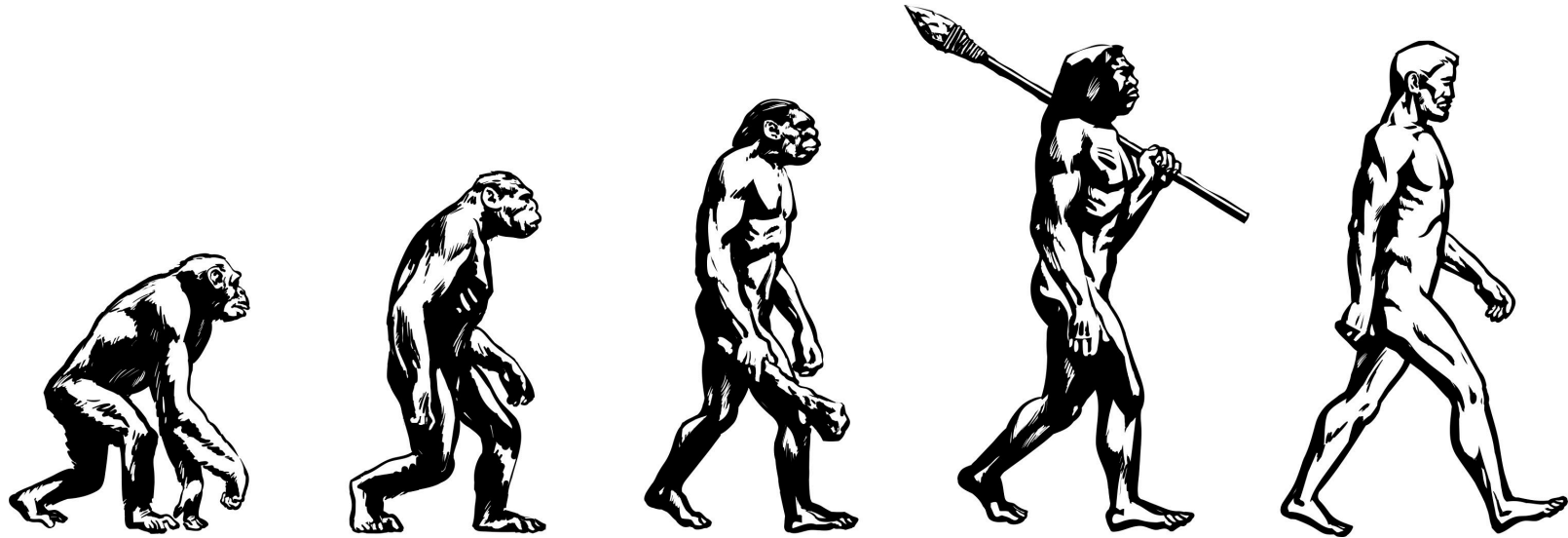


Fun with Phylogenetic Trees

Erin Molloy
University of Maryland, College Park

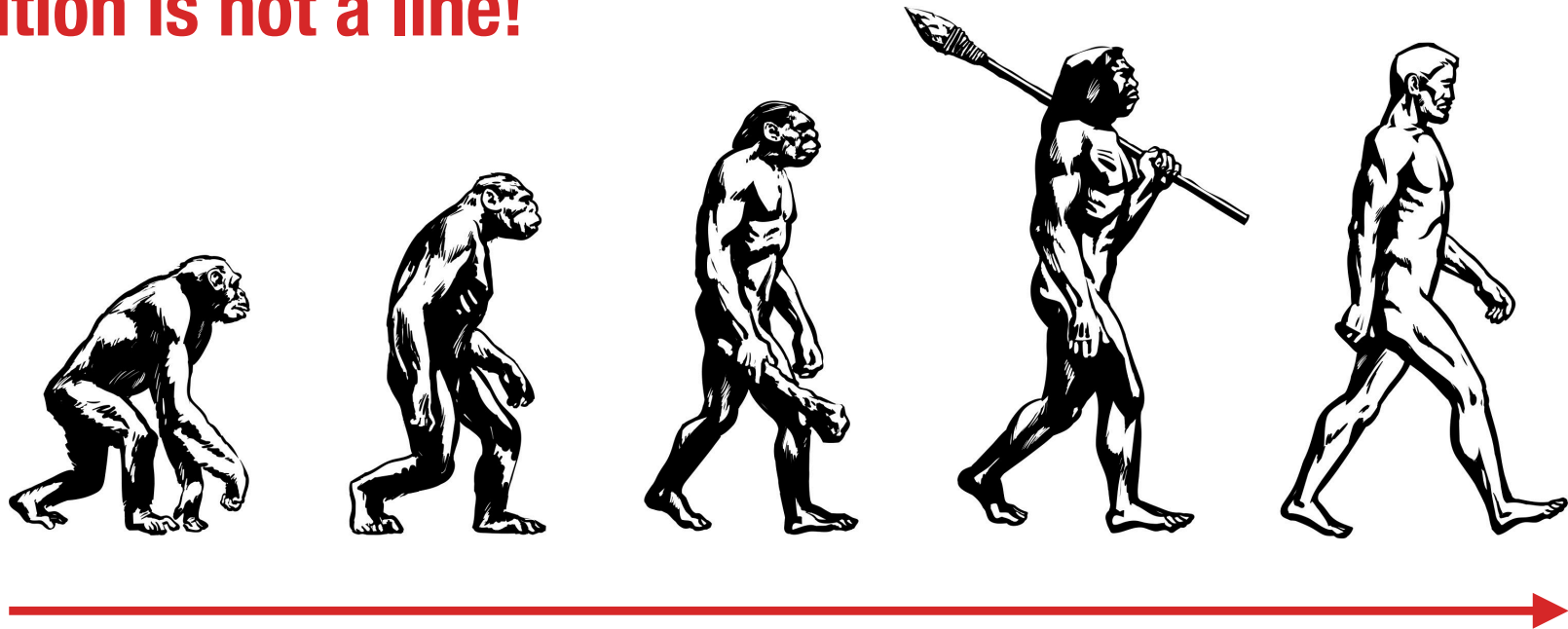
Guest Lecture in CMSC 423 Bioinformatics
April 25, 2023

Who has seen a similar picture?

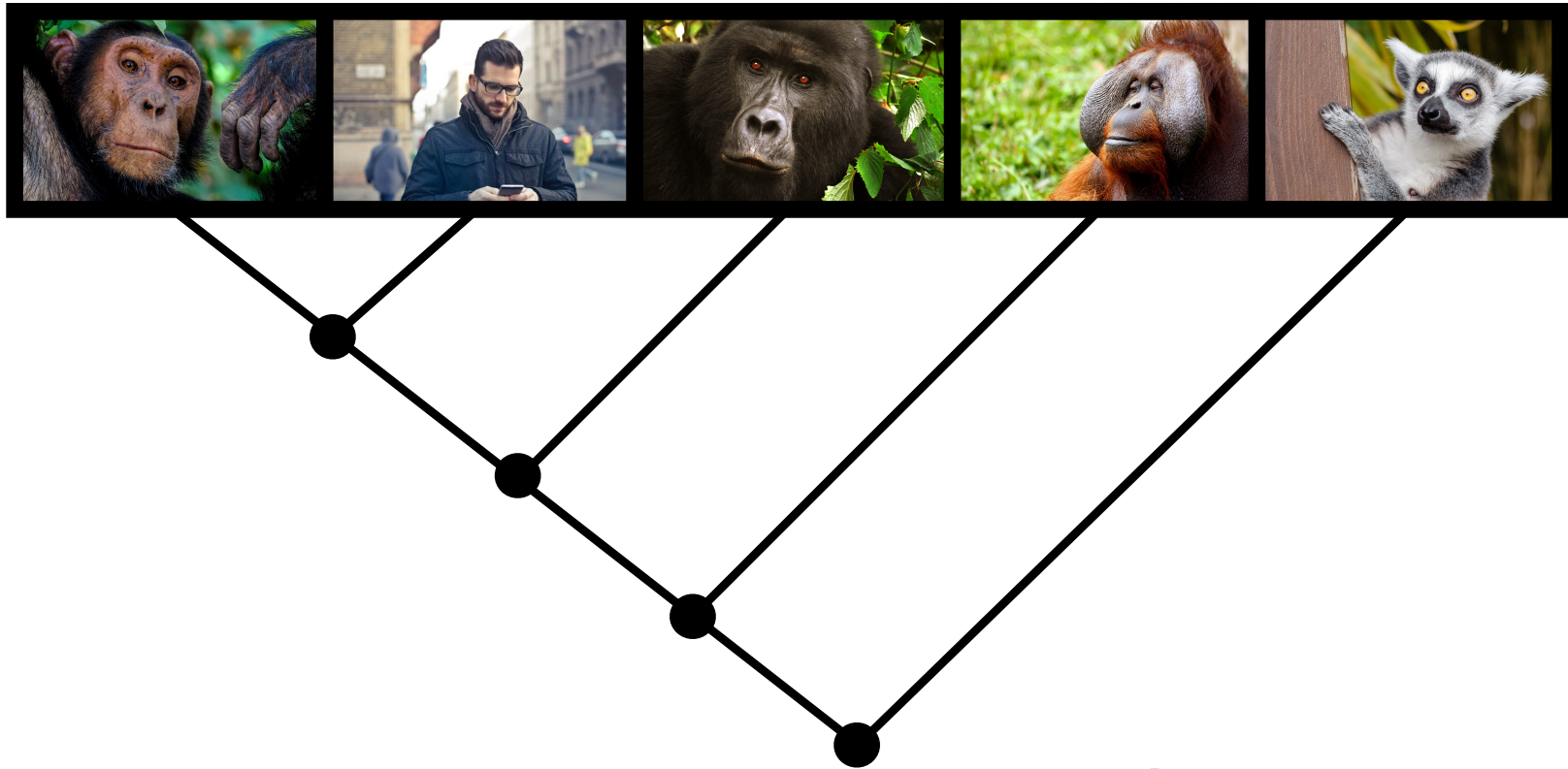


Who has seen a similar picture?

Evolution is not a line!

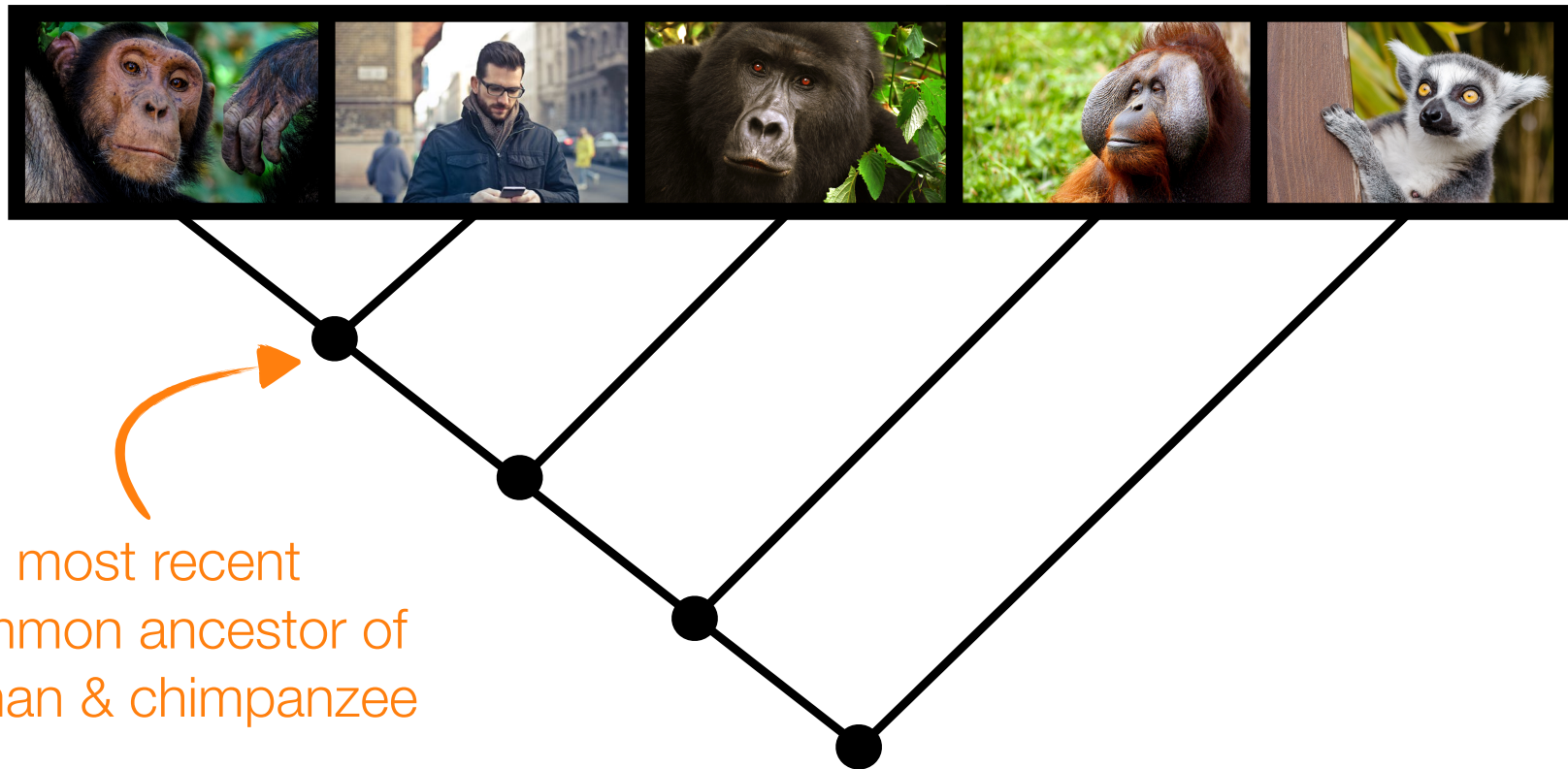


Evolution is a tree!



Evolutionary tree = phylogeny

Evolution is a tree!

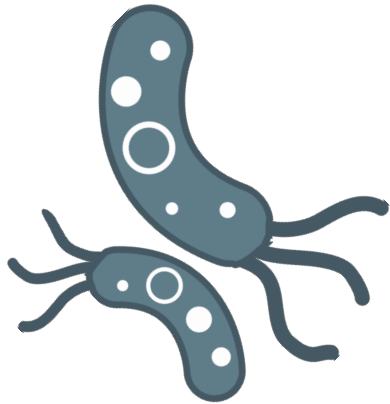


Evolutionary tree = phylogeny

Evolutionary trees are used in
bioinformatics analyses.

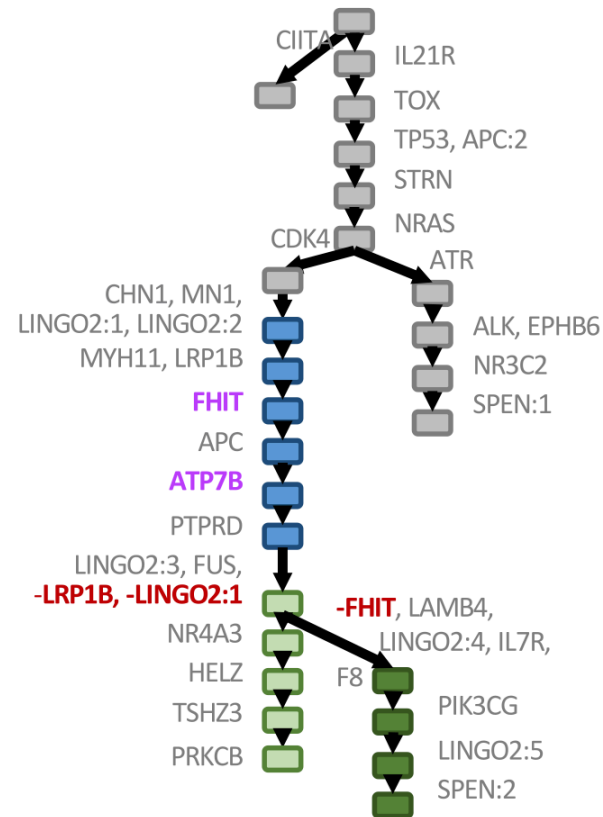
vaccine design & strain tracking

[Qiu et al., 2019; Hodcroft et al., 2021]



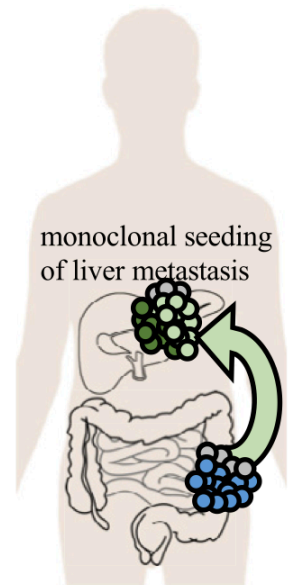
microbe identification

[Nguyen et al., 2014; Amato et al., 2019;
Meyer et al., 2019; Shah et al., 2021]



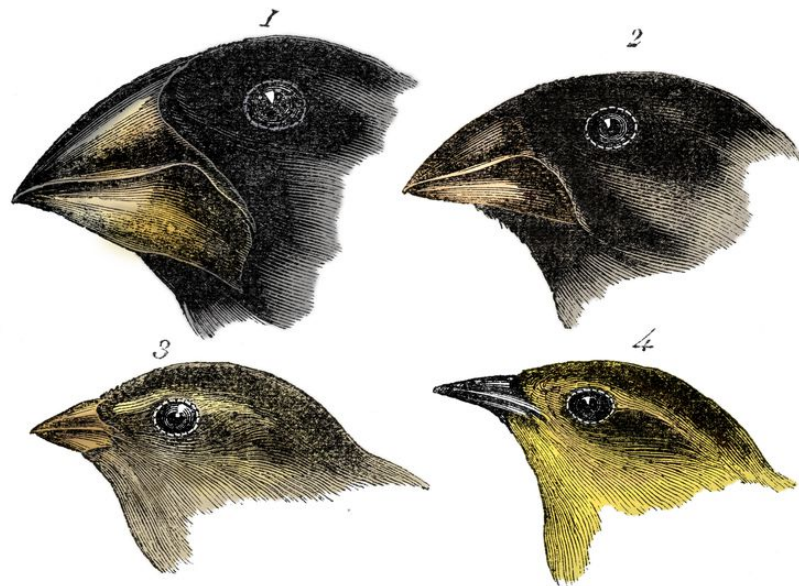
tumor evolution

[Satas et al., 2021]



Evolutionary trees are reconstructed
from **data**.

Morphology



1. *Geospiza magnirostris*.
3. *Geospiza parvula*.

2. *Geospiza fortis*.
4. *Certhidea olivacea*.

Darwin's finches

Genomics





1,000 bat genomes



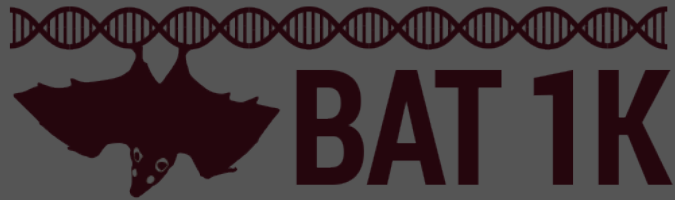
5,000 insect genomes



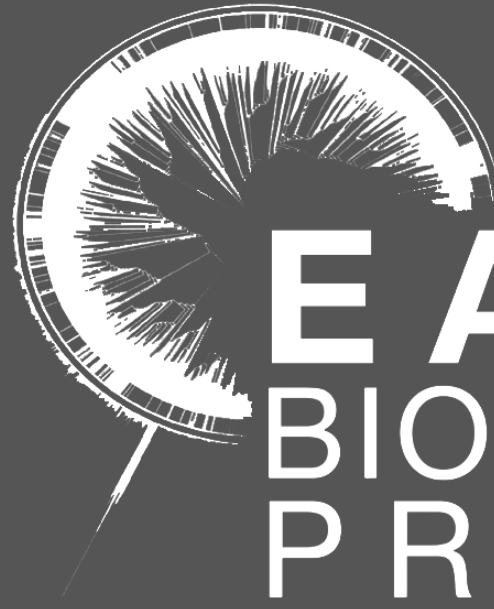
10,000 plant genomes



>60,000 vertebrate genomes



A PROJECT OF THE G10K CONSORTIUM



EARTH BIOGENOME PROJECT

~1.5 million eukaryotic genomes
in next 10 years

“Resolving the Tree of Life is unquestionably among the most complex scientific problems facing biology and presents challenges much greater than sequencing the human genome.”

From: “Assembling the Tree of Life: Harnessing Life’s History to Benefit Science and Society,” NSF, 2002.

RESEARCH ARTICLE

HUMAN GENOMICS

The complete sequence of a human genome

Sergey Nurk^{1†}, Sergey Koren^{1†}, Arang Rhie^{1†}, Mikko Rautiainen^{1†}, Andrey V. Bzikadze², Alla Mikheenko³, Mitchell R. Vollger⁴, Nicolas Altemose⁵, Lev Uralsky^{6,7}, Ariel Gershman⁸, Sergey Aganezov^{9†}, Savannah J. Hoyt¹⁰, Mark Diekhans¹¹, Glennis A. Logsdon⁴, Michael Alonge⁹, Stylianos E. Antonarakis¹², Matthew Borchers¹³, Gerard G. Bouffard¹⁴, Shelise Y. Brooks¹⁴, Gina V. Caldas¹⁵, Nae-Chyun Chen⁹, Haoyu Cheng^{16,17}, Chen-Shan Chin¹⁸, William Chow¹⁹, Leonardo G. de Lima¹³, Philip C. Dishuck⁴, Richard Durbin^{19,20}, Tatiana Dvorkina³, Ian T. Fiddes²¹, Giulio Formenti^{22,23}, Robert S. Fulton²⁴, Arkarachai Functammasan¹⁸, Erik Garrison^{11,25}, Patrick G. S. Grady¹⁰, Tina A. Graves-Lindsay²⁶, Ira M. Hall²⁷, Nancy F. Hansen²⁸, Gabrielle A. Hartley¹⁰, Marina Haukness¹¹, Kerstin Howe¹⁹, Michael W. Hunkapiller²⁹, Chirag Jain^{1,30}, Miten Jain¹¹, Erich D. Jarvis^{22,23}, Peter Kerpedjiev³¹, Melanie Kirsche⁹, Mikhail Kolmogorov³², Jonas Korf²⁹, Milinn Kremitzki²⁶, Heng Li^{16,17}, Valerie V. Maduro³³, Tobias Marschall³⁴, Ann M. McCartney¹, Jennifer McDaniel³⁵, Danny E. Miller^{4,36}, James C. Mullikin^{14,28}, Eugene W. Myers³⁷, Nathan D. Olson³⁵, Benedict Paten¹¹, Paul Peluso²⁹, Pavel A. Pevzner³², David Porubsky⁴, Tamara Potapova¹³, Evgeny I. Rogaev^{6,7,38,39}, Jeffrey A. Rosenfeld⁴⁰, Steven L. Salzberg^{9,41}, Valerie A. Schneider⁴², Fritz J. Sedlazeck⁴³, Kishwar Shafin¹¹, Colin J. Shew⁴⁴, Alaina Shumate⁴¹, Ying Sims¹⁹, Arian F. A. Smit⁴⁵, Daniela C. Soto⁴⁴, Ivan Sović^{29,46}, Jessica M. Storer⁴⁵, Aaron Streets^{5,47}, Beth A. Sullivan⁴⁸, Françoise Thibaud-Nissen⁴², James Torrance¹⁹, Justin Wagner³⁵, Brian P. Walenz¹⁴, Aaron Wenger²⁹, Jonathan M. D. Wood¹⁹, Chunlin Xiao⁴², Stephanie M. Yan⁴⁹, Alice C. Young¹⁴, Samantha Zarate⁹, Urvashi Surti⁵⁰, Rajiv C. McCoy⁴⁹, Megan Y. Dennis⁴⁴, Ivan A. Alexandrov^{3,7,51}, Jennifer L. Gerton^{13,52}, Rachel J. O'Neill¹⁰, Winston Timp^{8,41}, Justin M. Zook³⁵, Michael C. Schatz^{9,49}, Evan E. Eichler^{4,53*}, Karen H. Miga^{11,54*}, Adam M. Phillippy^{1*}

The human genome project “ended” in 2003 with ~92%...

And then it took ~20 more years to “complete” final ~8%.

We cannot analyze these forthcoming **BIG** datasets with the methods we have.

Challenge #1: Too many species

Best methods are heuristics for NP-hard optimization problems

Solution space (i.e., set of all possible phylogenetic trees) grows exponentially in number of species!



Challenge #1: Too many species

# leaves	#trees
4	3
5	15
6	105
7	945
8	10,395
9	135,135
10	2,027,025

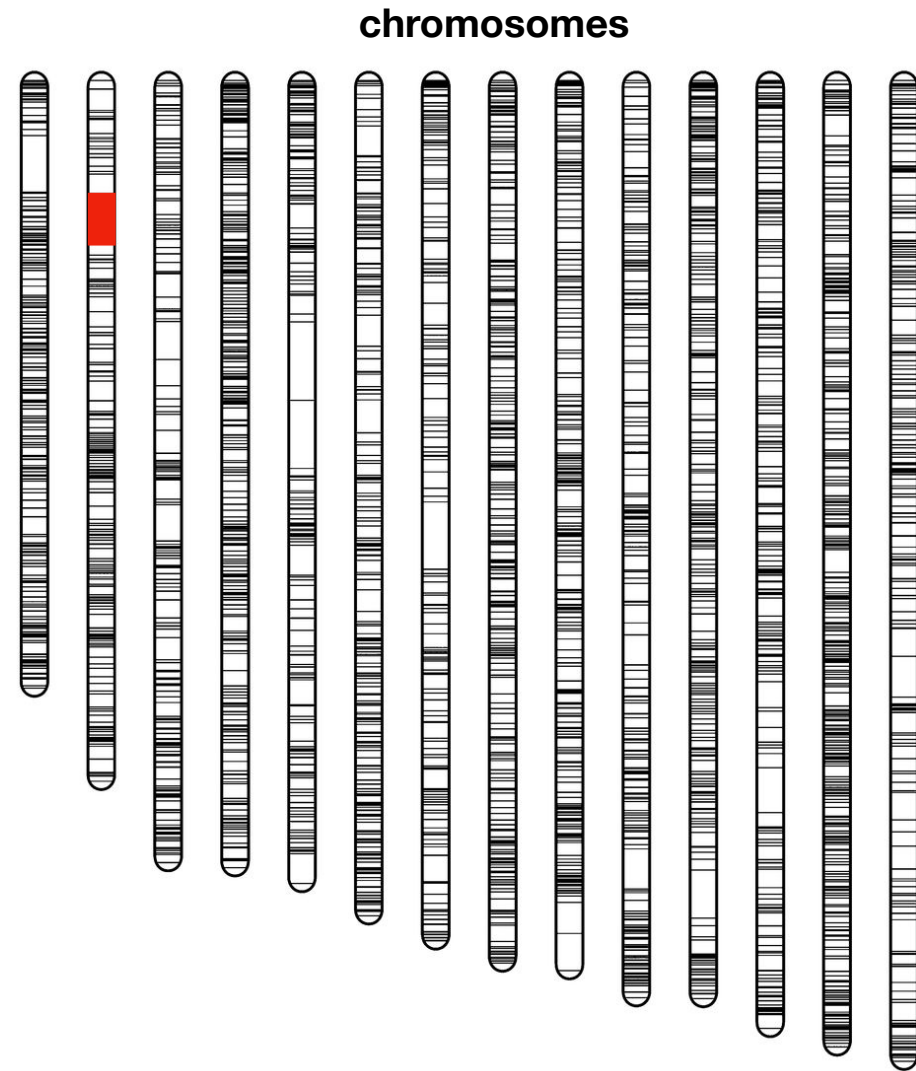


Challenge #2: Genome-scale data

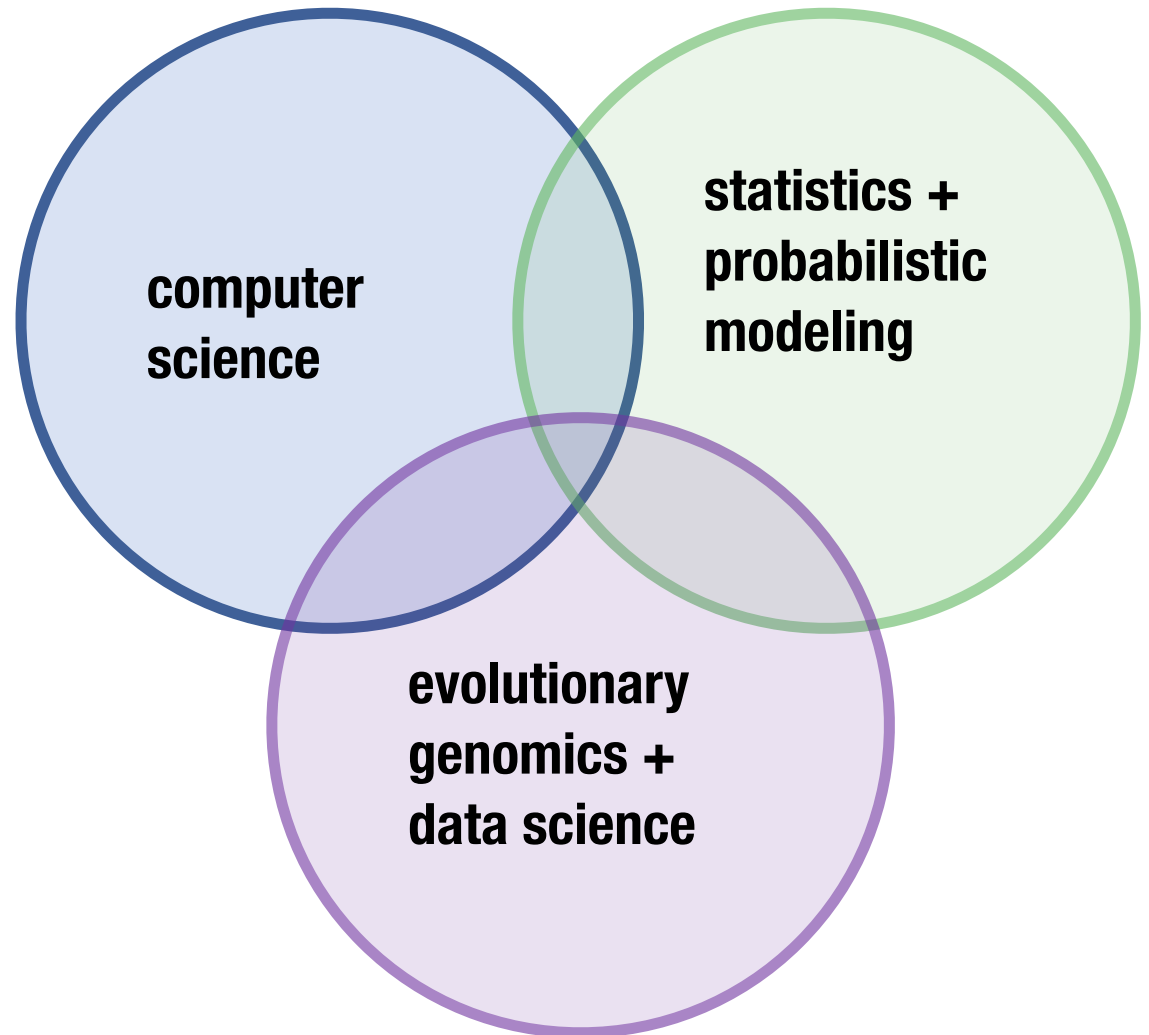
More data per species

Error, missing data, etc.

Modeling evolution of single **gene** vs.
genome (many genes)



Addressing
these
challenges
requires
interdisciplinary
research!



Agenda

Part 1: Perfect Phylogenies

Part 2: Small Parsimony Problem & Fitch's Algorithm

Part 3: Large Parsimony Problem

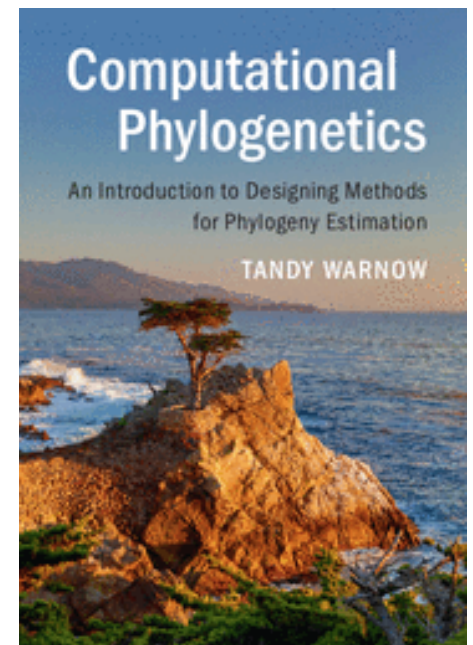
Part 4: Maximum Parsimony Methods

Acknowledgements

Material in slides based on
Computational Phylogenetics by
Prof. Tandy Warnow.



Tandy Warnow
University of Illinois at
Urbana-Champaign



Let's begin with a formal definition of **phylogenetic tree**.

Definition. A rooted phylogenetic tree T is a triplet (t, S, ϕ) , where

- t is a **rooted binary tree**
- S is a **set** of labels typically representing species, and
- ϕ is a **bijection** i.e., 1-1 function, mapping leaf vertices in t to labels in S

In this presentation

For simplicity, we do not make an explicit distinction between leaf nodes and their labels.

We will simply say that T is a phylogenetic tree on label set S , omitting the reference to ϕ .

We let

- $L(T)$ denote the **leaf set** of T
- $V(T)$ denote the **vertex set** of T
- $E(T)$ denote the **edge set** of T

Now we have a definition for phylogenetic trees — but what about the **data**?

Let's consider a simple example.



Definition. A character c is a surjection $c : S \rightarrow \{1, 2, \dots, k\}$ mapping labels (species) onto k states.



state 0 = trait is absent
state 1 = trait present



Definition. A character c is a surjection $c : S \rightarrow \{1, 2, \dots, k\}$ mapping labels (species) onto k states.

	c1
	Black Stripe
A	= 0
B	= 1
C	= 0
D	= 0



C



A

state 0 = trait is absent
state 1 = trait present



B



D

Class question — What are the characters implied by these birds?

		c1
		Black Stripe
A	=	0
B	=	1
C	=	0
D	=	0

state 0 = trait is absent
state 1 = trait present



Class question — What are the characters implied by these birds?

		c1	c2	c3	c4
		Black Stripe	Orange Wings	Orange Head	Yellow Tail
A	=	0	1	1	1
B	=	1	1	1	1
C	=	0	0	1	1
D	=	0	0	0	1



C



A



B



D

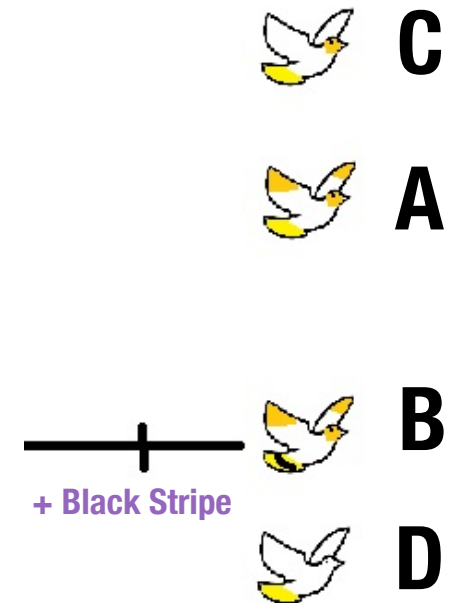
In practice, we are **given data** (characters) and **seek a phylogenetic tree** that best explains it.

Let's consider a simple example, again.



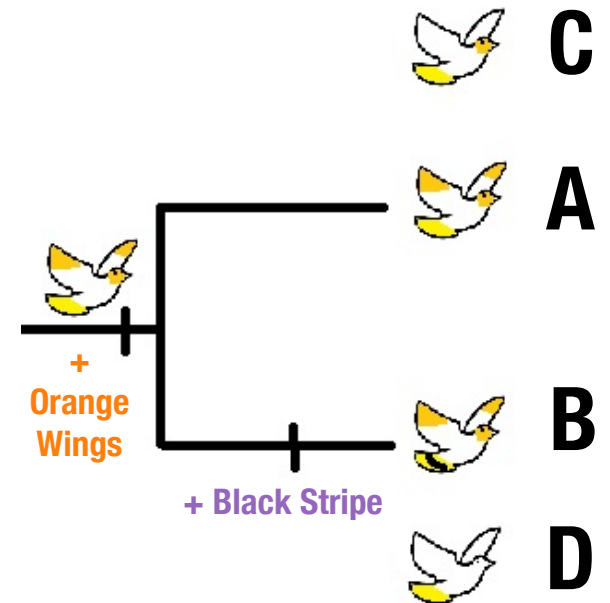
Let's consider a simple example, again.

		c1
		Black Stripe
A	=	0
B	=	1
C	=	0
D	=	0



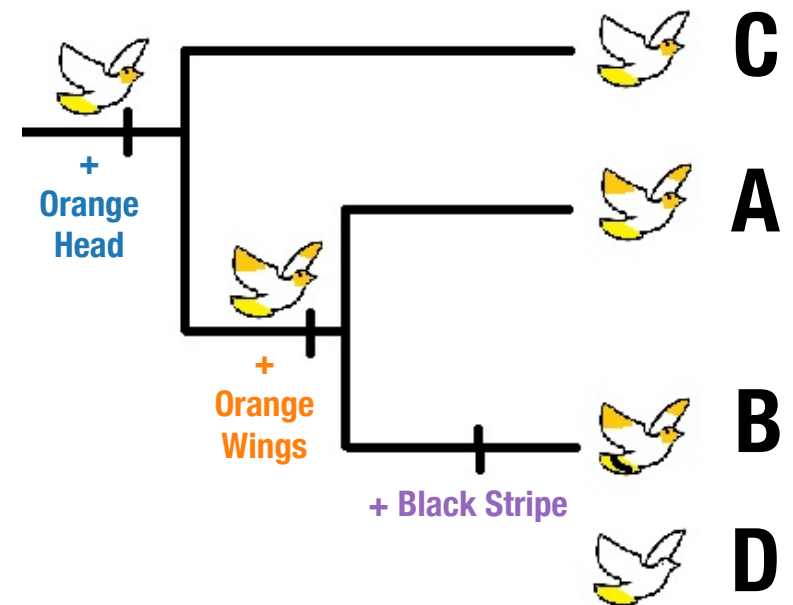
Let's consider a simple example, again.

		c1	c2
		Black Stripe	Orange Wings
A	=	0	1
B	=	1	1
C	=	0	0
D	=	0	0



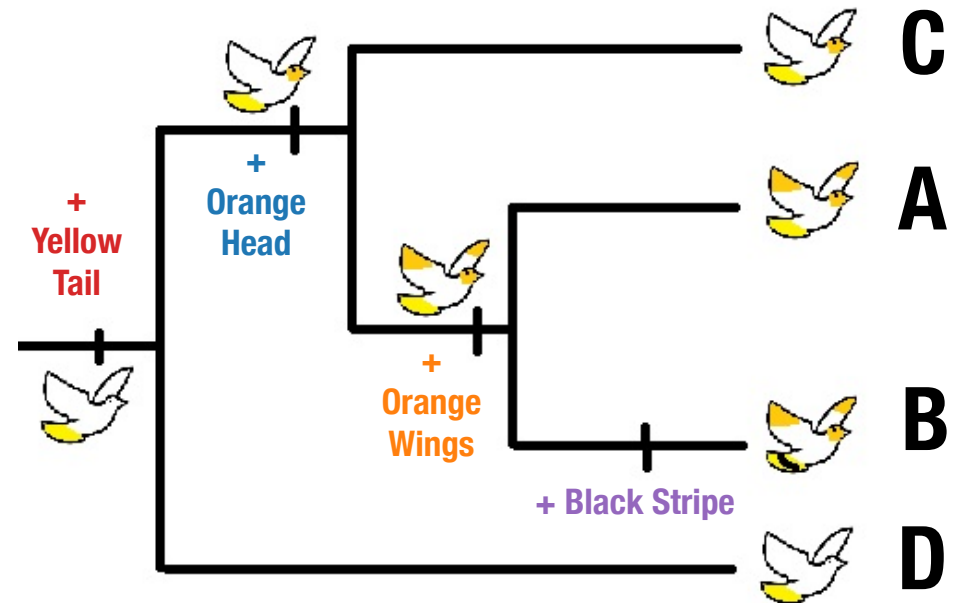
Let's consider a simple example, again.

		c1	c2	c3
		Black Stripe	Orange Wings	Orange Head
A	=	0	1	1
B	=	1	1	1
C	=	0	0	1
D	=	0	0	0



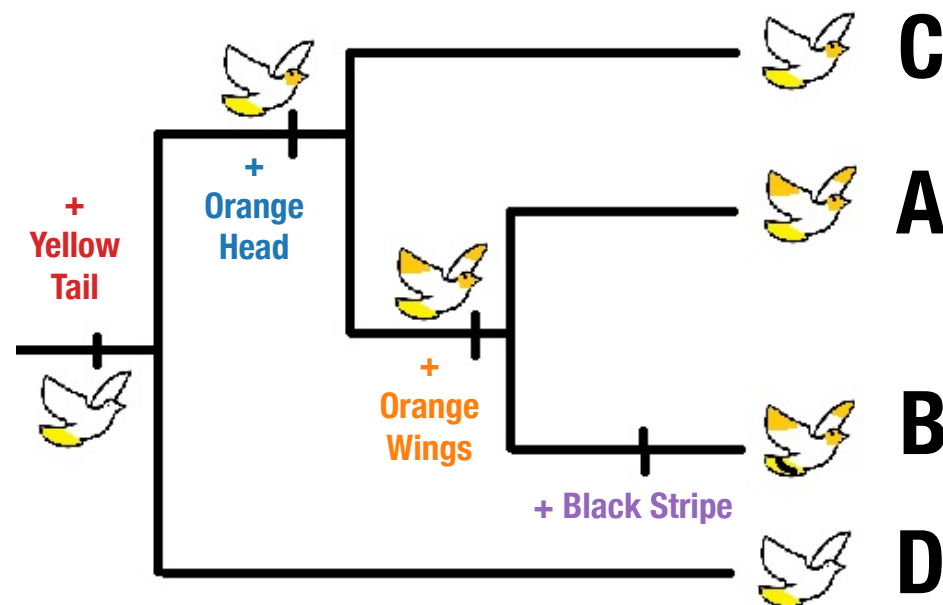
Let's consider a simple example, again.

		c1	c2	c3	c4
		Black Stripe	Orange Wings	Orange Head	Yellow Tail
A	=	0	1	1	1
B	=	1	1	1	1
C	=	0	0	1	1
D	=	0	0	0	1



Definition. A phylogenetic tree T is called a *perfect phylogeny* for a set \mathcal{C} of characters if every character can be explained by a trait arising on exactly one branch of T .

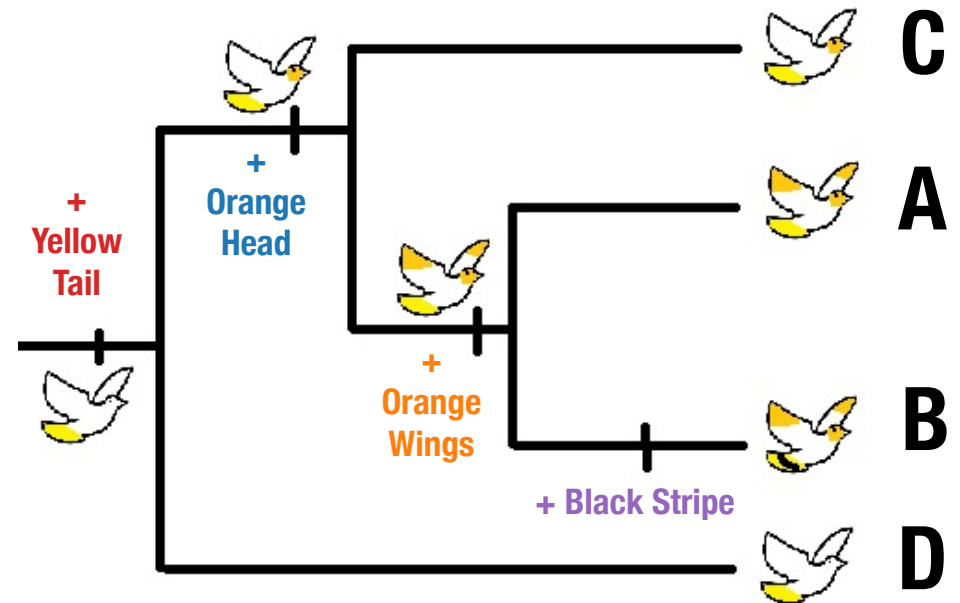
		c1	c2	c3	c4
		Black Stripe	Orange Wings	Orange Head	Yellow Tail
A	=	0	1	1	1
B	=	1	1	1	1
C	=	0	0	1	1
D	=	0	0	0	1



Now we want an **algorithm** to find a perfect phylogeny for our data if one exists.

Observation — There is a relationship between characters and clades (subsets of species).

		c1	c2	c3	c4
		Black Stripe	Orange Wings	Orange Head	Yellow Tail
A	=	0	1	1	1
B	=	1	1	1	1
C	=	0	0	1	1
D	=	0	0	0	1



Observation — There is a relationship between characters and clades (subsets of species).

Algorithm Sketch:

1. Assume 0 is ancestral state and 1 is mutated state. Write down subset of species implied by each character:

{ {B}, {A, B}, {A, B, C}, {A,B,C,D} }

	c1	c2	c3	c4
	Black Stripe	Orange Wings	Orange Head	Yellow Tail
A =	0	1	1	1
B =	1	1	1	1
C =	0	0	1	1
D =	0	0	0	1

Observation — There is a relationship between characters and clades (subsets of species).

Algorithm Sketch:

1. Assume 0 is ancestral state and 1 is mutated state. Write down subset of species implied by each character:

{ {B}, {A, B}, {A, B, C}, {A,B,C,D} }

2. Add trivial sets (set of 1 species and set of all species).

	c1	c2	c3	c4
	Black Stripe	Orange Wings	Orange Head	Yellow Tail
A =	0	1	1	1
B =	1	1	1	1
C =	0	0	1	1
D =	0	0	0	1

Observation — There is a relationship between characters and clades (subsets of species).

Algorithm Sketch:

1. Assume 0 is ancestral state and 1 is mutated state. Write down subset of species implied by each character:

{ {B}, {A, B}, {A, B, C}, {A,B,C,D} }

2. Add trivial sets (set of 1 species and set of all species).
3. Build tree using Hasse Diagram.

	c1	c2	c3	c4
	Black Stripe	Orange Wings	Orange Head	Yellow Tail
A =	0	1	1	1
B =	1	1	1	1
C =	0	0	1	1
D =	0	0	0	1

For step 3, need to:

Define a partial order on clades so that the Hasse Diagram produces a phylogeny.

But first some definitions...

Definition. A relation on set X is subset of the Cartesian product $X \times X$, which is the set formed by taking exactly two elements from X , in all possible ways.

But first some definitions...

Definition. A relation on set X is subset of the Cartesian product $X \times X$, which is the set formed by taking exactly two elements from X , in all possible ways.

Definition. A partial order is relation R on set X satisfying three properties:

- $\langle x, y \rangle \in R$ and $\langle y, z \rangle \in R$ implies that $\langle x, z \rangle \in R$ (TRANSITIVITY)
- $\langle x, x \rangle \in R$ for all $x \in X$
- $\langle x, y \rangle \in R$ and $\langle y, x \rangle \in R$ implies $x = y$

But first some definitions...

Definition. The Hasse Diagram for a set X with a partial order R is constructed in three steps:

- (1) create vertex for each element in X ,
- (2) add directed edge $x \rightarrow y$ if $\langle x, y \rangle \in R$ and $x \neq y$, and
- (3) remove arrows implied by transitivity.

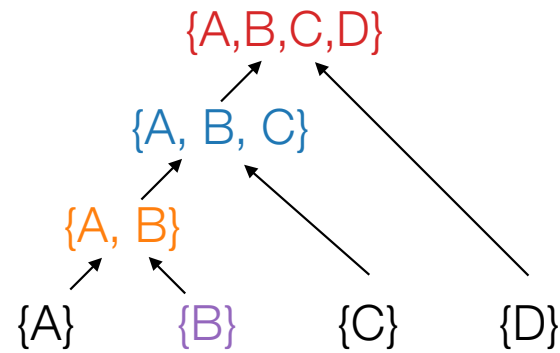
Going
back
to our
example...

	c1	c2	c3	c4
	Black Stripe	Orange Wings	Orange Head	Yellow Tail
A =	0	1	1	1
B =	1	1	1	1
C =	0	0	1	1
D =	0	0	0	1

Algorithm Sketch:

1. Assume 0 is ancestral state and 1 is mutated state. Write down subset of species implied by each character:
 $\{ \{B\}, \{A, B\}, \{A, B, C\}, \{A, B, C, D\} \}$
2. Add trivial sets (set of 1 species and set of all species).
3. Build tree using Hasse Diagram.

Going
back
to our
example...

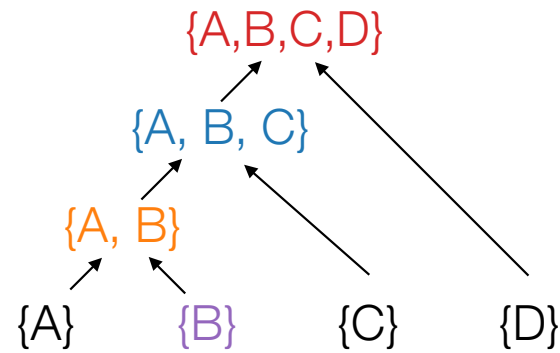


Algorithm Sketch:

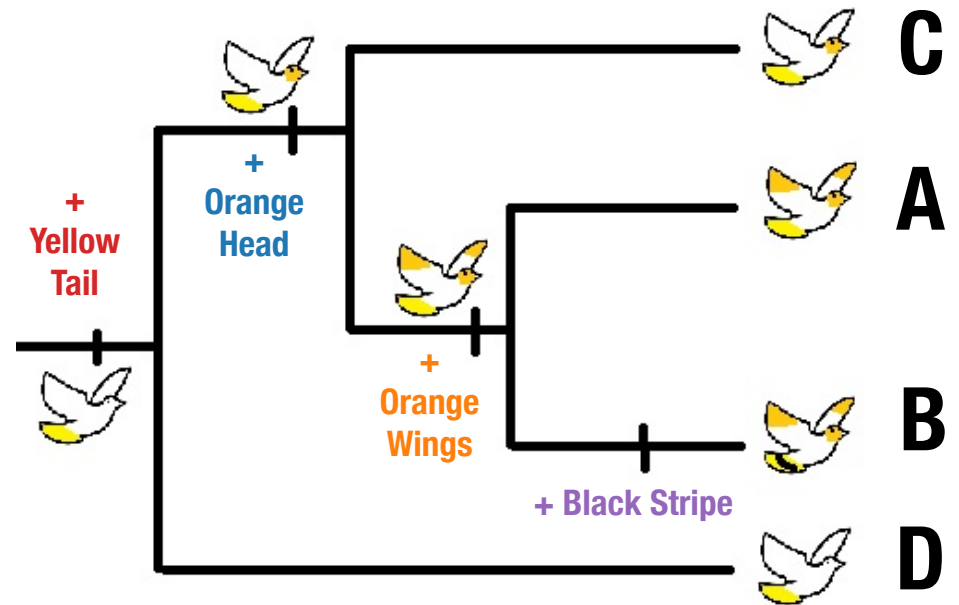
1. Assume 0 is ancestral state and 1 is mutated state. Write down subset of species implied by each character:
 $\{ \{B\}, \{A, B\}, \{A, B, C\}, \{A, B, C, D\} \}$
2. Add trivial sets (set of 1 species and set of all species).
3. Build tree using Hasse Diagram.

	c1	c2	c3	c4
	Black Stripe	Orange Wings	Orange Head	Yellow Tail
A =	0	1	1	1
B =	1	1	1	1
C =	0	0	1	1
D =	0	0	0	1

Going
back
to our
example...



		c1	c2	c3	c4
		Black Stripe	Orange Wings	Orange Head	Yellow Tail
A	=	0	1	1	1
B	=	1	1	1	1
C	=	0	0	1	1
D	=	0	0	0	1



Now we have an **algorithm** to find a perfect phylogeny for our data if one exists.

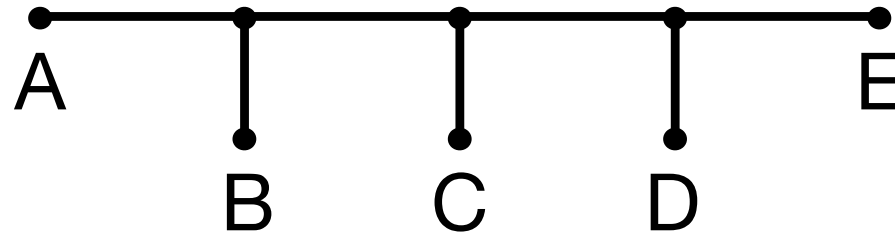
However, perfect phylogenies unlikely to exist in practice!

A perfect phylogenies will exist if characters

- **evolve without homoplasy AND**
- **are correctly called+coded** for all labels in the set S ,
 - no error
 - no missing or ambiguous states

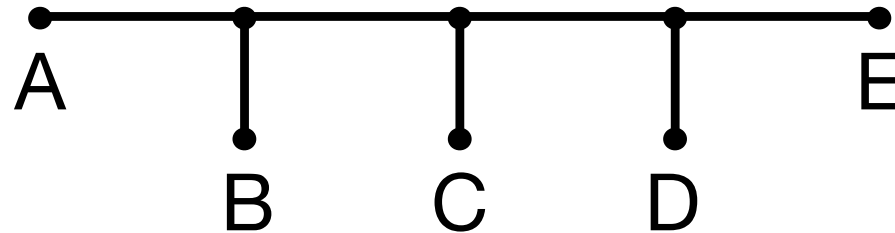
Now let's define **homoplasy** for the case where characters are undirected (i.e., we don't know which state is ancestral or mutated).

	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



Definition. Given a tree T and a k -state character c , we say that c evolves without *homoplasy* if the internal nodes can be labeled with states so that each substitution produces a new state.

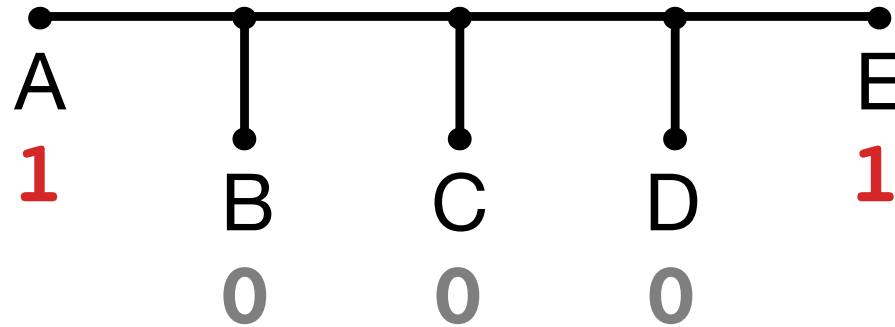
	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



Definition. Given a tree T and a k -state character c , we say that c evolves without *homoplasy* if the internal nodes can be labeled with states so that each substitution produces a new state.

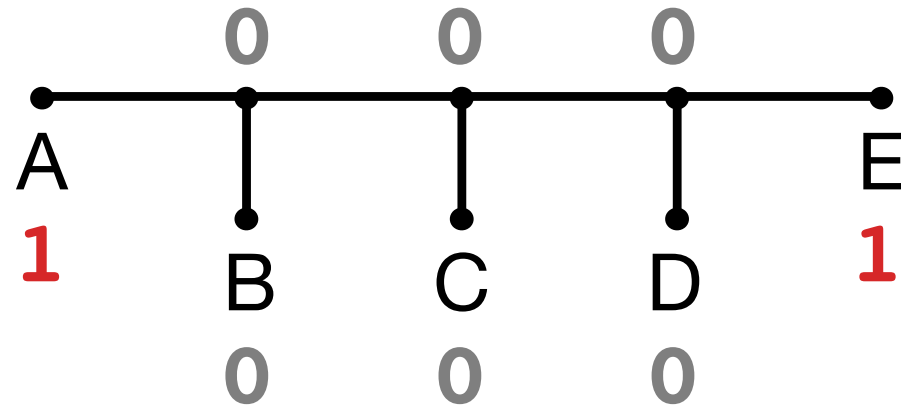
Remark. Given a tree T and a k -state character, the character evolves without *homoplasy* if it can be explained with $k - 1$ substitutions.

	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



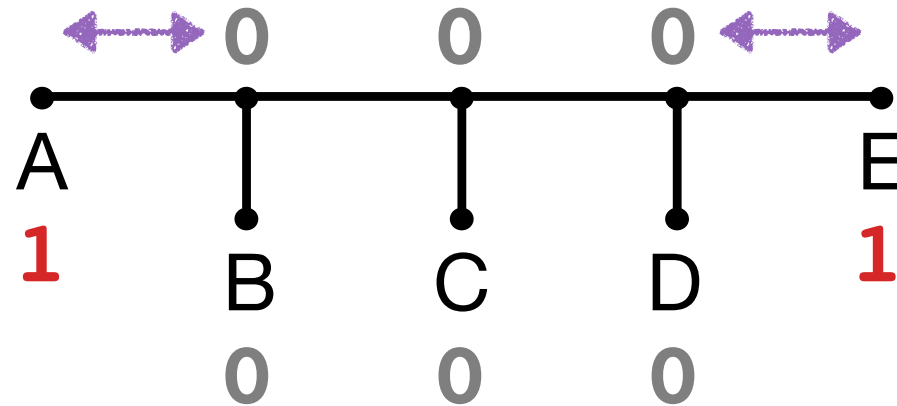
Remark. Given a tree T and a k -state character, the character evolves without homoplasy if it can be explained with $k - 1$ substitutions.

	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



Remark. Given a tree T and a k -state character, the character evolves without homoplasy if it can be explained with $k - 1$ substitutions.

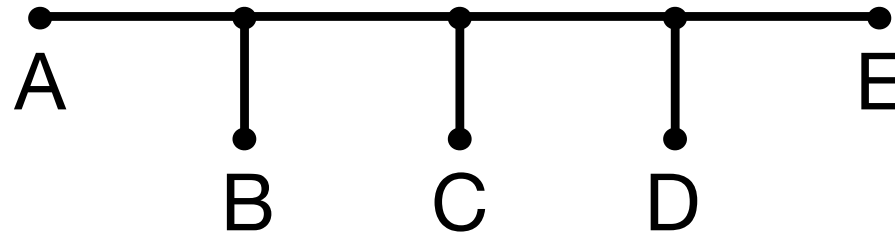
	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



At least 2 substitutions is required to explain c_1 , across all possible labelings of the internal nodes... so it evolved with homoplasy!

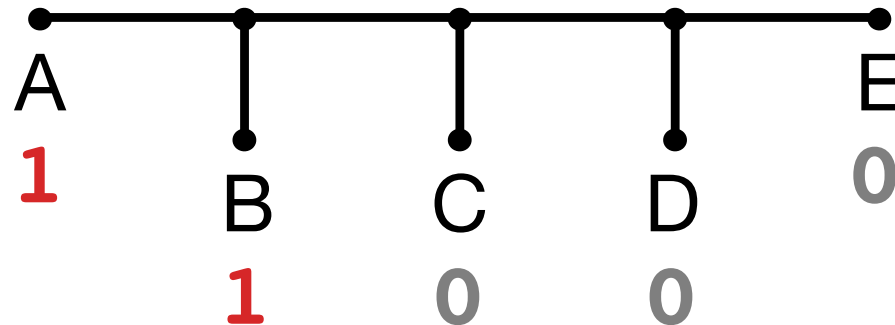
Remark. Given a tree T and a k -state character, the character evolves without homoplasy if it can be explained with $k - 1$ substitutions.

	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



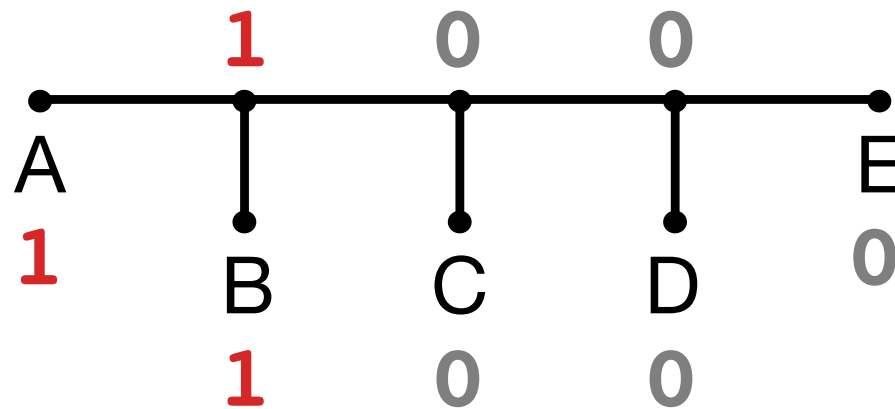
Remark. Given a tree T and a k -state character, the character evolves without homoplasy if it can be explained with $k - 1$ substitutions.

	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



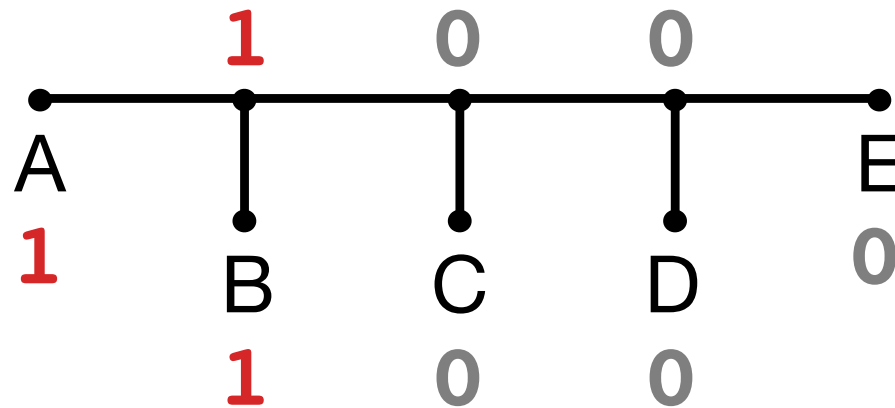
Remark. Given a tree T and a k -state character, the character evolves without homoplasy if it can be explained with $k - 1$ substitutions.

	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



Remark. Given a tree T and a k -state character, the character evolves without homoplasy if it can be explained with $k - 1$ substitutions.

	c_1	c_2
A =	1	1
B =	0	1
C =	0	0
D =	0	0
E =	1	0



Only 1 substitution needed to explain c_2 , across all possible labelings of the internal nodes... so it evolved **WITHOUT** homoplasy!

Remark. Given a tree T and a k -state character, the character evolves without homoplasy if it can be explained with $k - 1$ substitutions.

Agenda



Part 1: Perfect Phylogenies

Part 2: Small Parsimony Problem & Fitch's Algorithm

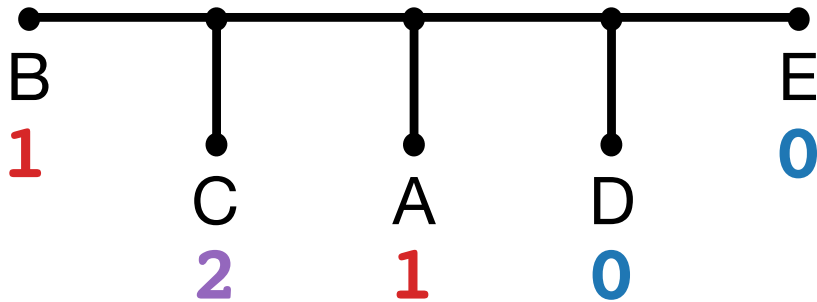
Part 3: Large Parsimony Problem

Part 4: Maximum Parsimony Methods

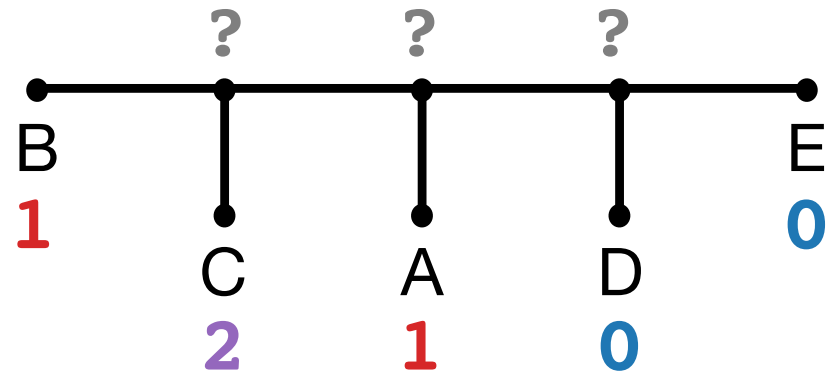
Moving on... we want to compute the unordered parsimony score to determine whether an undirected character evolves with homoplasy!

Definition 3 (parsimony score). Given a tree T and a character c , both on label set S , the *parsimony score*, denoted $length(T, c)$, is the *minimum # of substitutions* required to explain the states at the leaves.

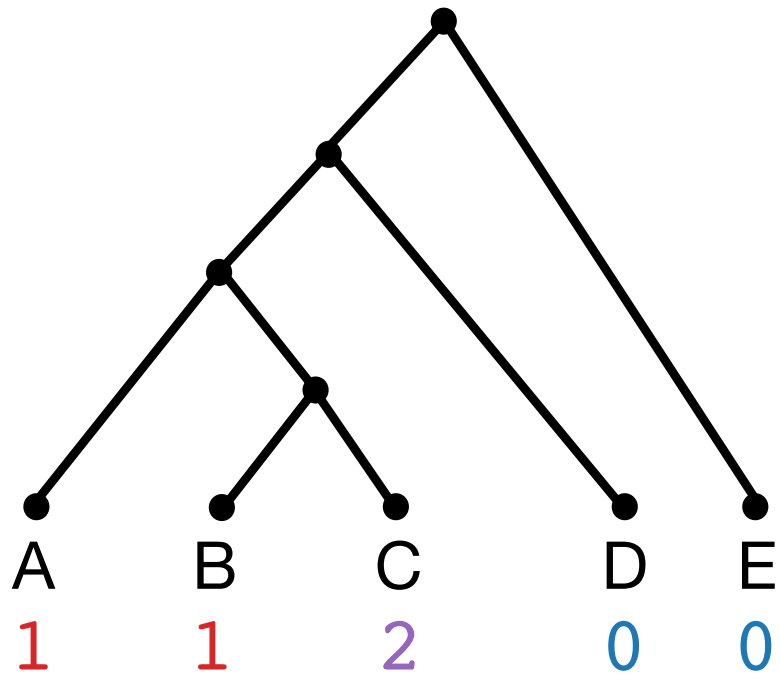
Small Parsimony Problem



Input. The pair (T, c) , where T is a an unrooted binary phylogenetic tree and c is a character, both on label set S .

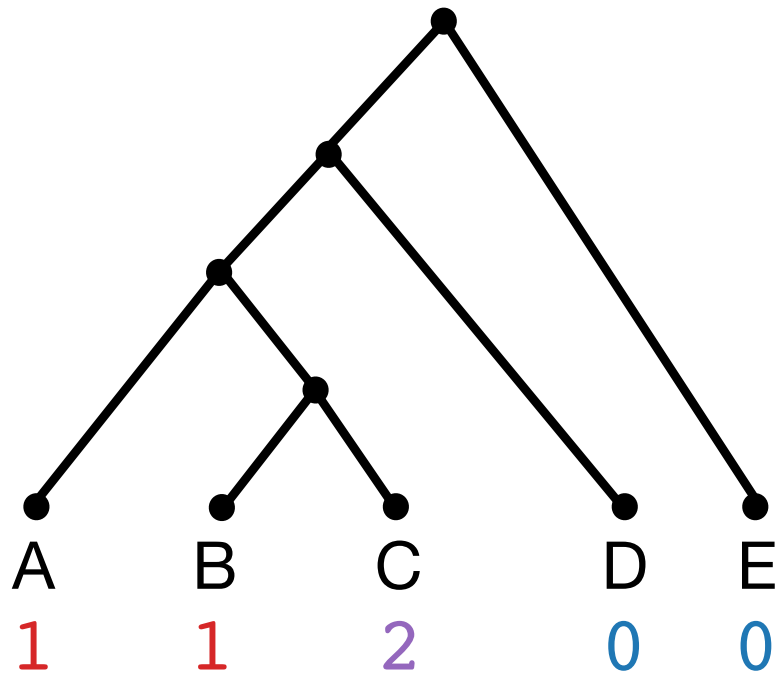


Output. An assignment of character states to the internal nodes of T to minimize the # of substitutions, i.e. the # of edges $e = (u, v)$ for which $c(u) \neq c(v)$



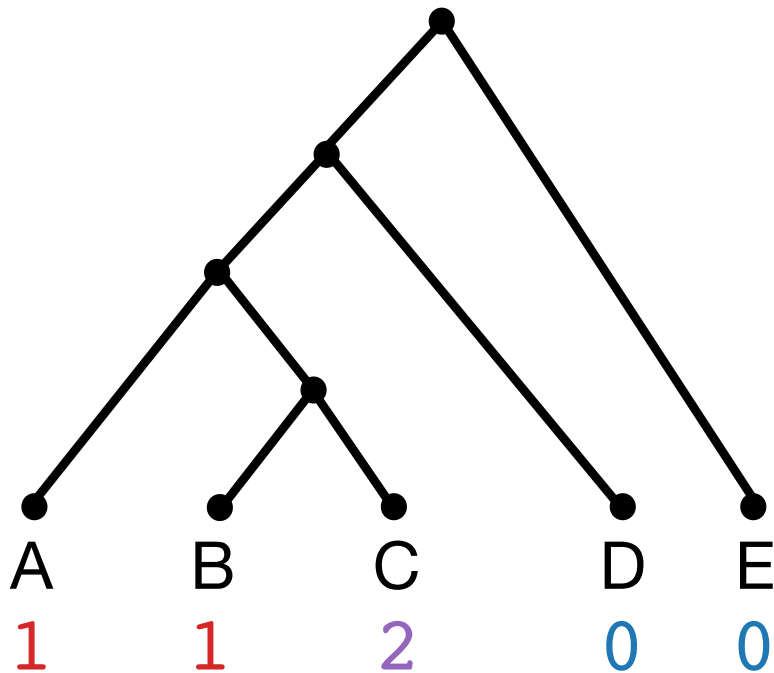
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r



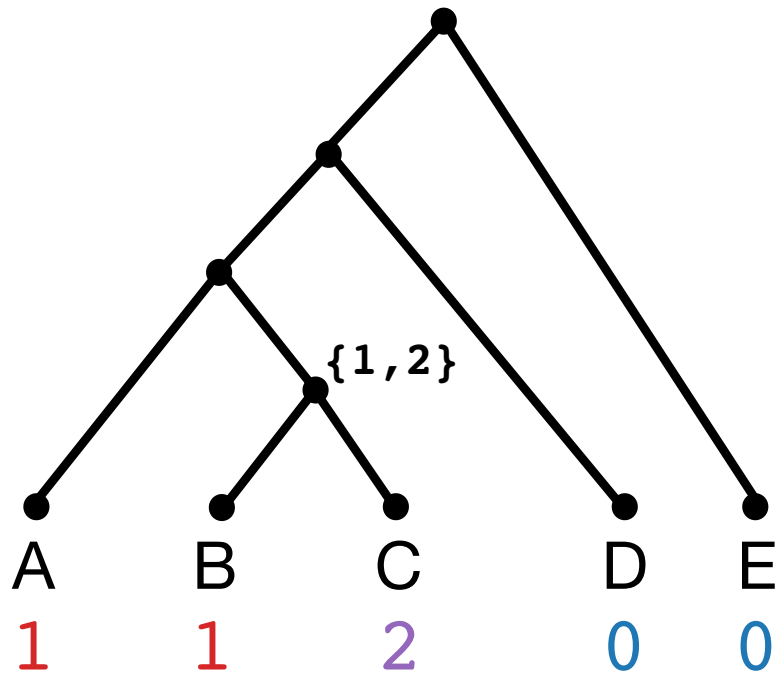
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$



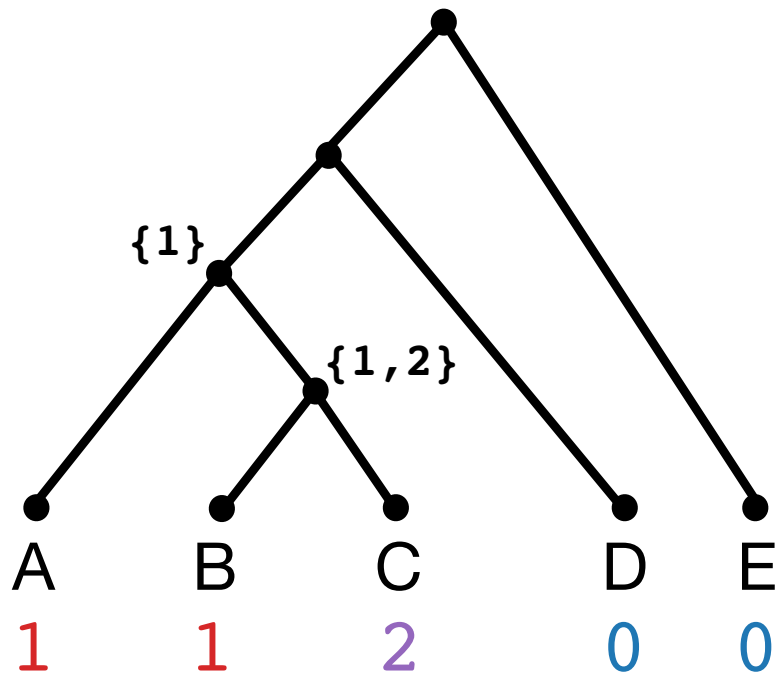
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$



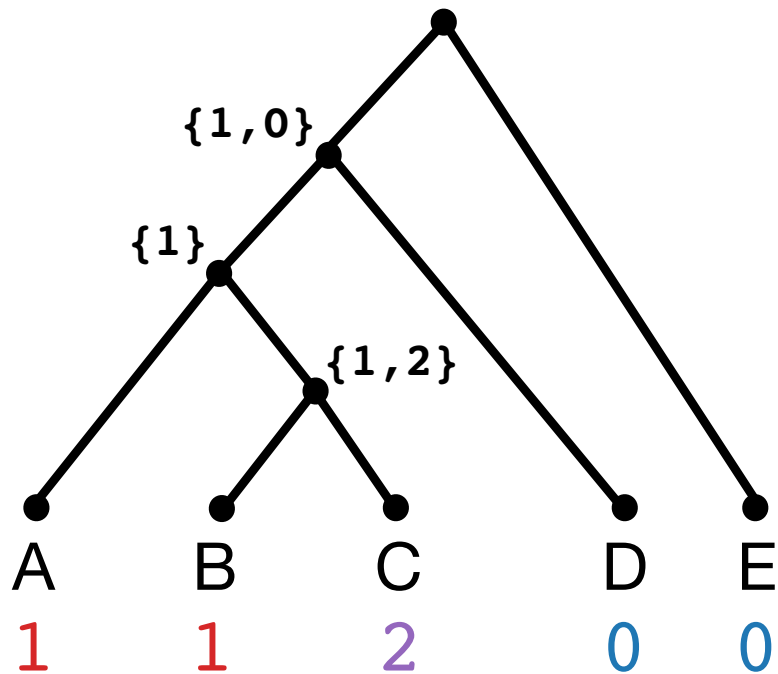
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow \text{children of } v$
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$



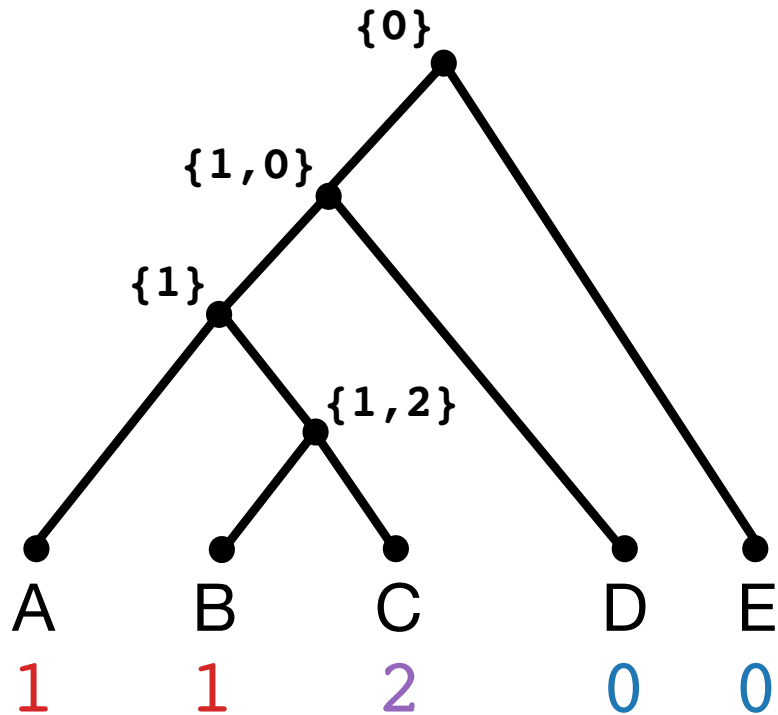
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow \text{children of } v$
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$



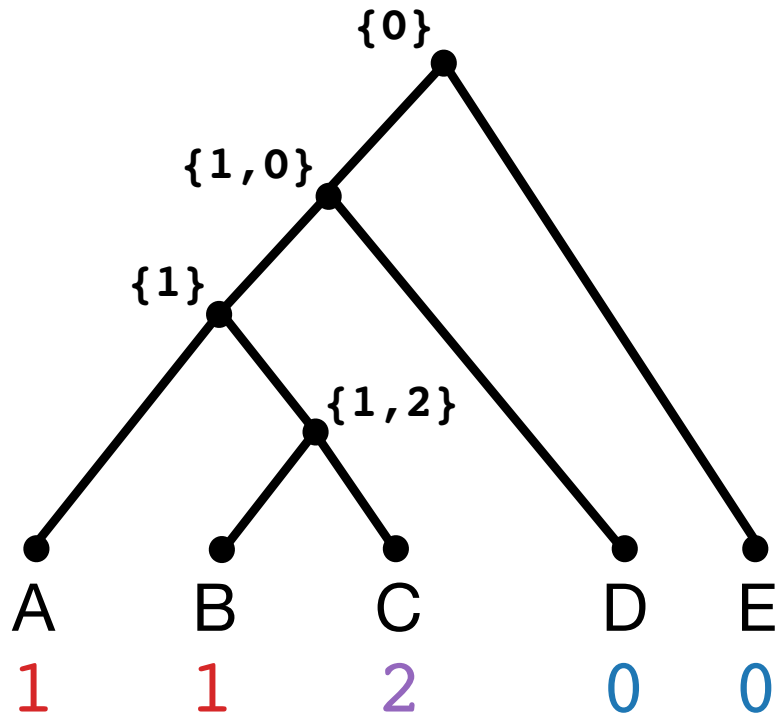
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow \text{children of } v$
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$



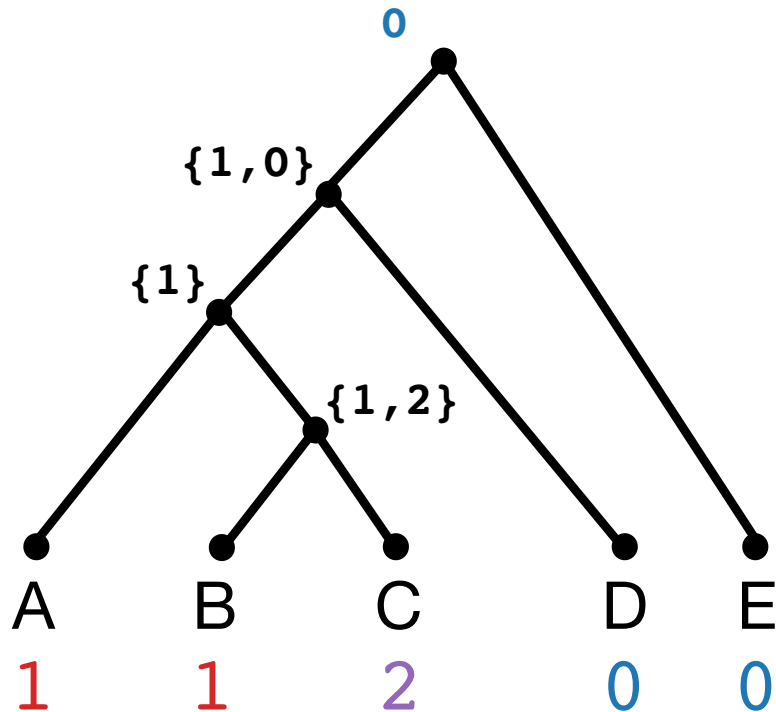
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$



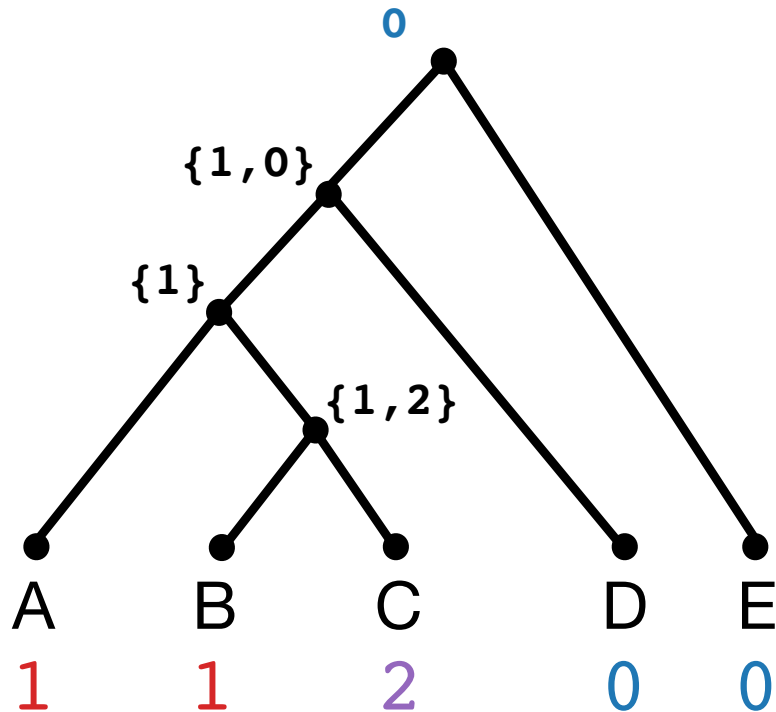
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root



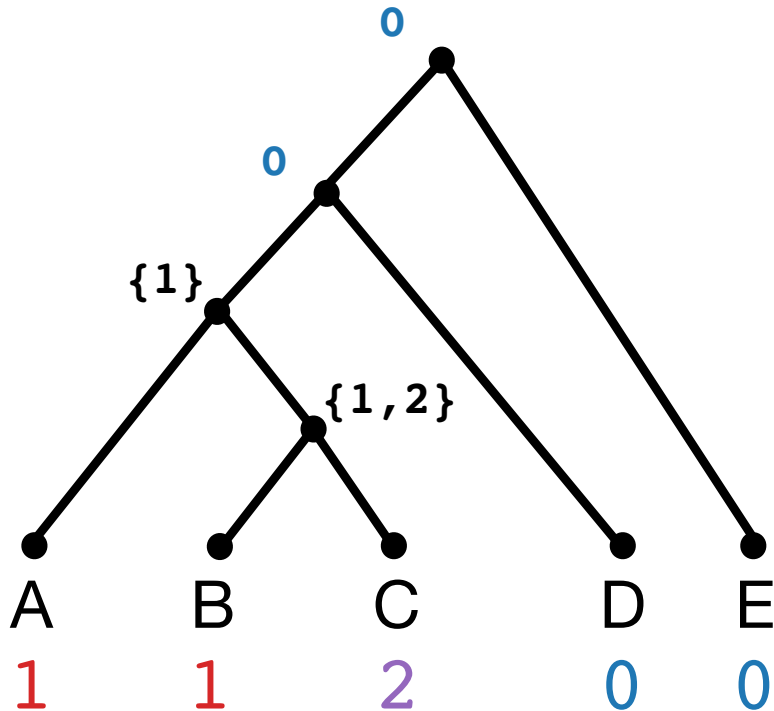
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root



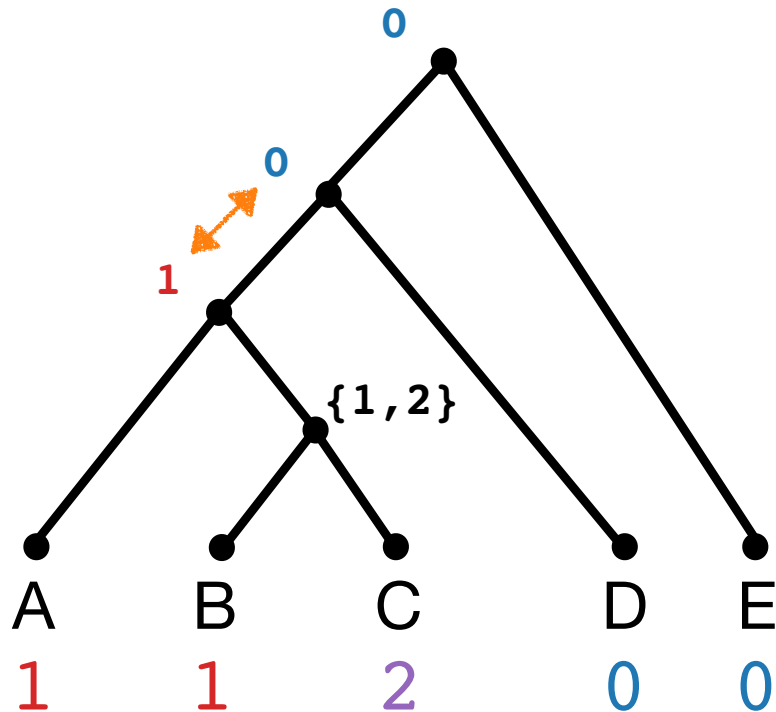
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root
5. Perform a pre-order traversal of T and for each vertex $v \in V(T) \setminus L(T) \setminus \{r\}$:
 - a. $u \leftarrow$ parent of v
 - b. If $c(u) \in A(v)$: $c(v) \leftarrow c(u)$
 - c. Else: $c(v) \leftarrow$ arbitrary state in $A(u)$



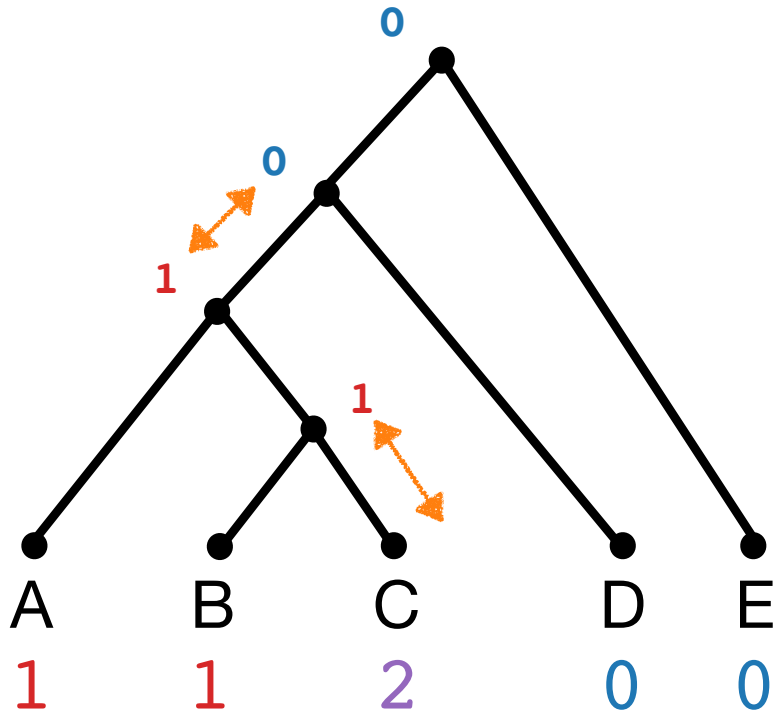
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root
5. Perform a pre-order traversal of T and for each vertex $v \in V(T) \setminus L(T) \setminus \{r\}$:
 - a. $u \leftarrow$ parent of v
 - b. If $c(u) \in A(v)$: $c(v) \leftarrow c(u)$
 - c. Else: $c(v) \leftarrow$ arbitrary state in $A(u)$



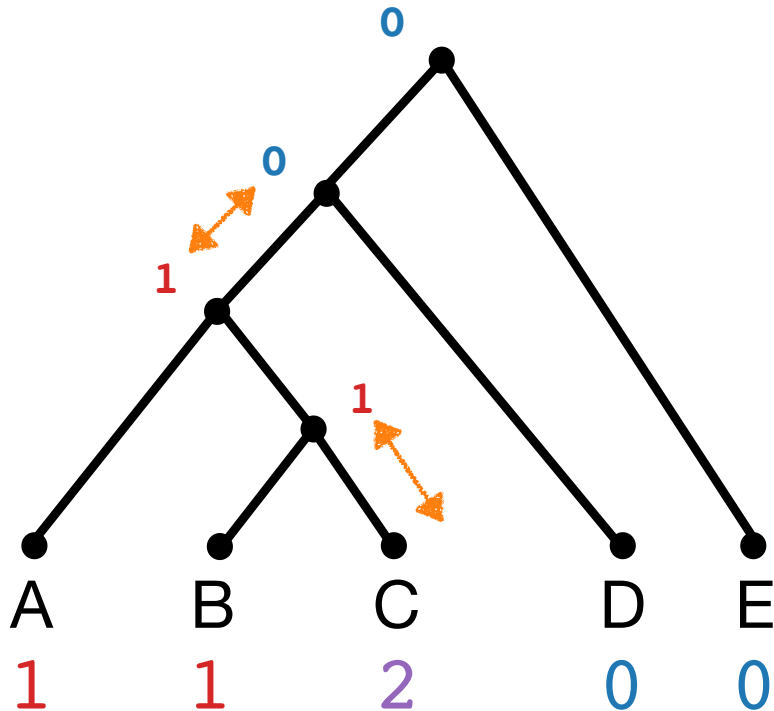
FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root
5. Perform a pre-order traversal of T and for each vertex $v \in V(T) \setminus L(T) \setminus \{r\}$:
 - a. $u \leftarrow$ parent of v
 - b. If $c(u) \in A(v)$: $c(v) \leftarrow c(u)$
 - c. Else: $c(v) \leftarrow$ arbitrary state in $A(u)$



FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root
5. Perform a pre-order traversal of T and for each vertex $v \in V(T) \setminus L(T) \setminus \{r\}$:
 - a. $u \leftarrow$ parent of v
 - b. If $c(u) \in A(v)$: $c(v) \leftarrow c(u)$
 - c. Else: $c(v) \leftarrow$ arbitrary state in $A(u)$



FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root
5. Perform a pre-order traversal of T and for each vertex $v \in V(T) \setminus L(T) \setminus \{r\}$:
 - a. $u \leftarrow$ parent of v
 - b. If $c(u) \in A(v)$: $c(v) \leftarrow c(u)$
 - c. Else: $c(v) \leftarrow$ arbitrary state in $A(u)$
6. Return c minus the root

Class Exercise — modify this algorithm to return the parsimony score?

FitchAlgorithm(T, c):

1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root
5. Perform a pre-order traversal of T and for each vertex $v \in V(T) \setminus L(T) \setminus \{r\}$:
 - a. $u \leftarrow$ parent of v
 - b. If $c(u) \in A(v)$: $c(v) \leftarrow c(u)$
 - c. Else: $c(v) \leftarrow$ arbitrary state in $A(u)$
6. Return c minus the root

Class Exercise — modify this algorithm to return the parsimony score?

Add 1 to the parsimony score every time you hit this line!

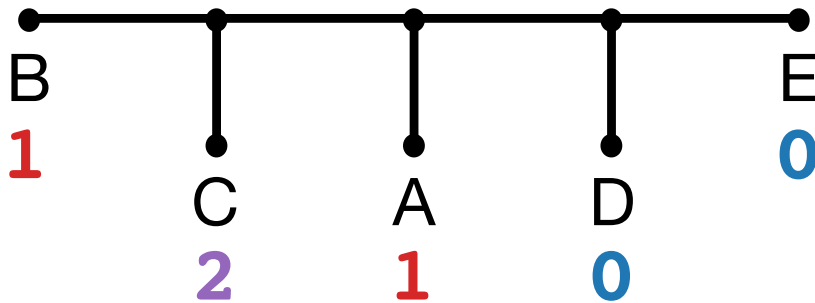
Exit at the end of step 3.

FitchAlgorithm(T, c):

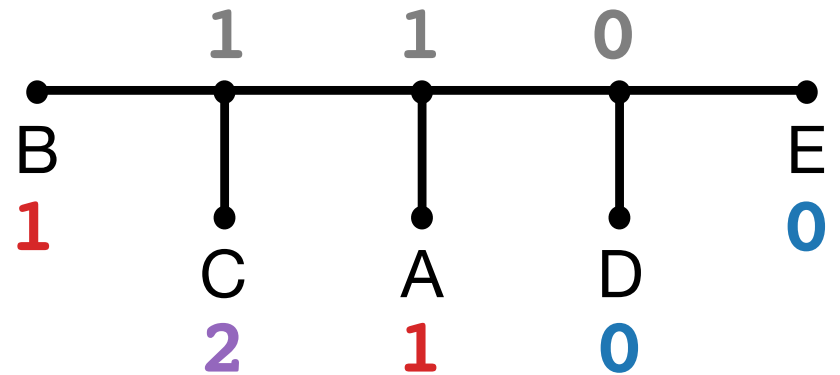
1. Root T , subdividing an arbitrary edge with root r
2. For each $l \in L(T)$: $A(l) \leftarrow \{c(l)\}$
3. Perform a post-order traversal of T and for each vertex $v \in V(T) \setminus L(T)$:
 - a. $(w, w') \leftarrow$ children of v
 - b. If $A(w) \cap A(w') \neq \emptyset$: $A(v) \leftarrow A(w) \cap A(w')$
 - c. Else: $A(v) \leftarrow A(w) \cup A(w')$
4. $c(r) \leftarrow$ arbitrary state in $A(r)$ to root
5. Perform a pre-order traversal of T and for each vertex $v \in V(T) \setminus L(T) \setminus \{r\}$:
 - a. $u \leftarrow$ parent of v
 - b. If $c(u) \in A(v)$: $c(v) \leftarrow c(u)$
 - c. Else: $c(v) \leftarrow$ arbitrary state in $A(u)$
6. Return c minus the root

Small Parsimony Problem

score = 2



Input. The pair (T, c) , where T is a an unrooted binary phylogenetic tree and c is a character, both on label set S .



Output. An assignment of character states to the internal nodes of T to minimize the # of substitutions, i.e. the # of edges $e = (u, v)$ for which $c(u) \neq c(v)$

We covered Fitch's algorithm, in which all substitutions have cost 1. Sankoff's algorithm generalizes this idea by allowing substitutions to have different costs!

ToDo: What is the time complexity of Fitch's algorithm...

To prove correctness, define subproblems $\mathbf{Cost}(v, x)$ which is the optimal parsimony score of rooted subtree T_v given the assignment $\mathbf{c}(v) = x$.

Show this holds for base case (leaves), make inductive hypothesis, and then show it holds for some vertex v .

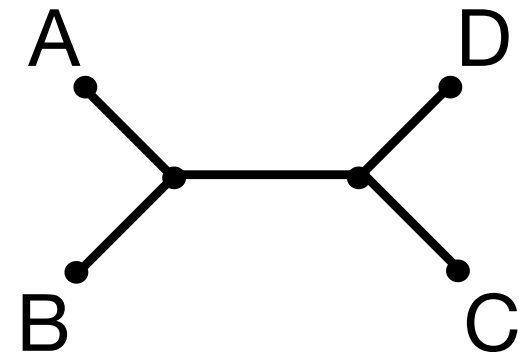
Agenda

- ✓ **Part 1:** Perfect Phylogenies
- ✓ **Part 2:** Small Parsimony Problem & Fitch's Algorithm
- Part 3:** Large Parsimony Problem
- Part 4:** Maximum Parsimony Methods

Maximum Parsimony (aka Large Parsimony Problem)

	c_1	c_2
A =	0	0
B =	0	1
C =	1	0
D =	1	0

Input. A set \mathcal{C} of k -state character , each on label set S .



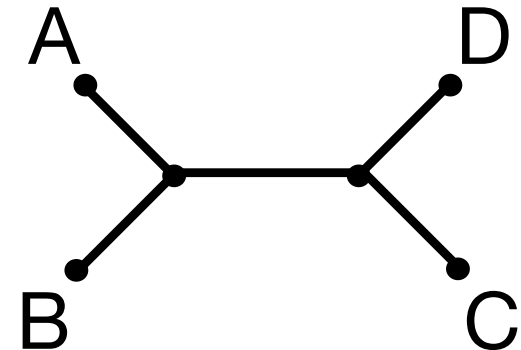
Output. A tree T on label set S that minimizes the total # of substitutions required to explain \mathcal{C} .

NP-hard (Foulds and Graham, 1982)

Maximum Parsimony (aka Large Parsimony Problem)

	c_1	c_2
A =	0	0
B =	0	1
C =	1	0
D =	1	0

Input. A set \mathcal{C} of k -state character ,
each on label set S .



Output. A tree T on label set S that minimizes the total # of substitutions required to explain \mathcal{C} .

Agenda

- ✓ **Part 1:** Perfect Phylogenies
- ✓ **Part 2:** Small Parsimony Problem & Fitch's Algorithm
- ✓ **Part 3:** Large Parsimony Problem
- Part 4:** Maximum Parsimony Methods

MP is NP-hard.
Now what?

MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

But it's challenging because **tree space** has $(2n - 5)!!$ unrooted trees on n leaves!

n	#trees
4	3
5	15
6	105
7	945
8	10,395
9	135,135
10	2,027,025

MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

Branch-and-Bound:

Like exhaustive search but better

MP is NP-hard.

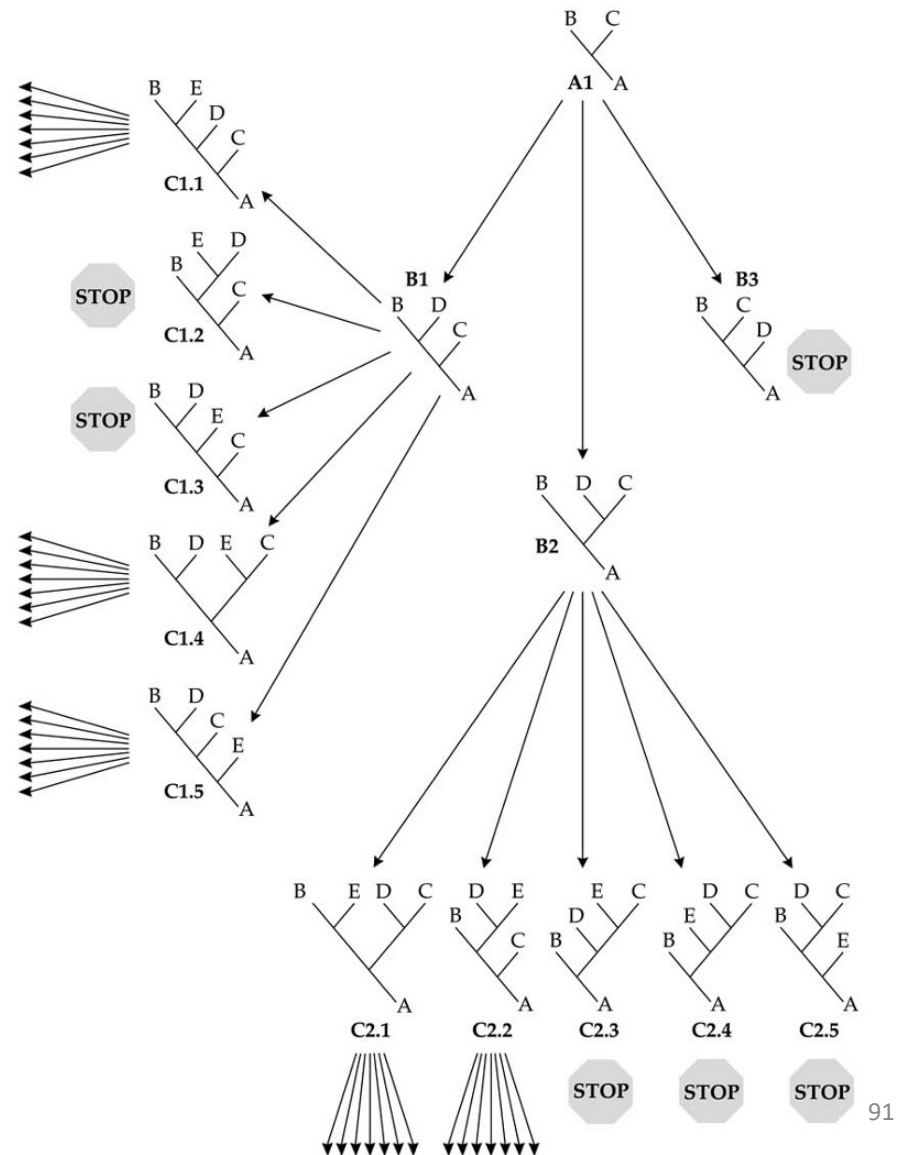
Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

Branch-and-Bound:

Like exhaustive search but better



MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

Branch-and-Bound:

Like exhaustive search but better

Heuristic (e.g. hill climbing):

Compute a starting tree. Apply some operation to **edit** the tree. Search from new tree if it's score is higher.

MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

Branch-and-Bound:

Like exhaustive search but better

Heuristic (e.g. hill climbing):

Compute a starting tree. Apply some operation to **edit** the tree. Search from new tree if it's score is higher.

But then we need methods to...

1. **Build** a starting tree.

>randomized taxon addition

2. **Edit** an existing tree.

>SPR, NNI, TBR moves

MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

Branch-and-Bound:

Like exhaustive search but better

Heuristic (e.g. hill climbing):

Compute a starting tree. Apply some operation to **edit** the tree. Search from new tree if it's score is higher.

But then we need methods to...

1. **Build** a starting tree.

>randomized taxon addition

2. **Edit** an existing tree.

>SPR, NNI, TBR moves

Input characters: $A = (1, 1, 1)$ $B = (1, 1, 1)$ $C = (0, 0, 0)$ $D = (0, 0, 1)$ $E = (0, 1, 0)$ **RandomizedTaxonAddition(\mathcal{C}, S):**

1. $L \leftarrow$ put labels in S in random order

Input characters:

$A = (1, 1, 1)$

$B = (1, 1, 1)$

$C = (0, 0, 0)$

$D = (0, 0, 1)$

$E = (0, 1, 0)$

$L = [B, C, A, D, E]$

RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order

Input characters:

$A = (1, 1, 1)$

$B = (1, 1, 1)$

$C = (0, 0, 0)$

$D = (0, 0, 1)$

$E = (0, 1, 0)$

$L = [B, C, A, D, E]$

RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L

Input characters:

A = (1, 1, 1)

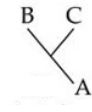
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L

Input characters:

A = (1, 1, 1)

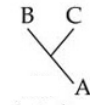
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$

Input characters:

A = (1, 1, 1)

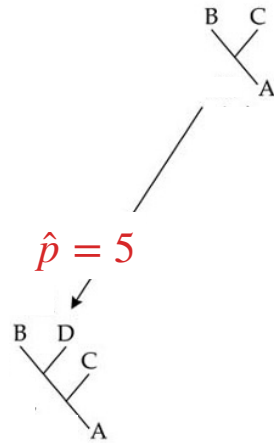
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$

Input characters:

A = (1, 1, 1)

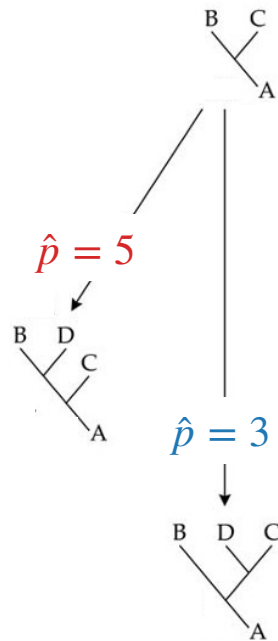
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$

Input characters:

A = (1, 1, 1)

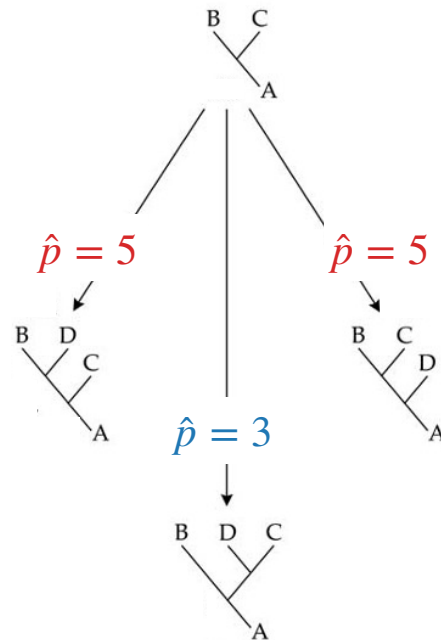
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$

Input characters:

A = (1, 1, 1)

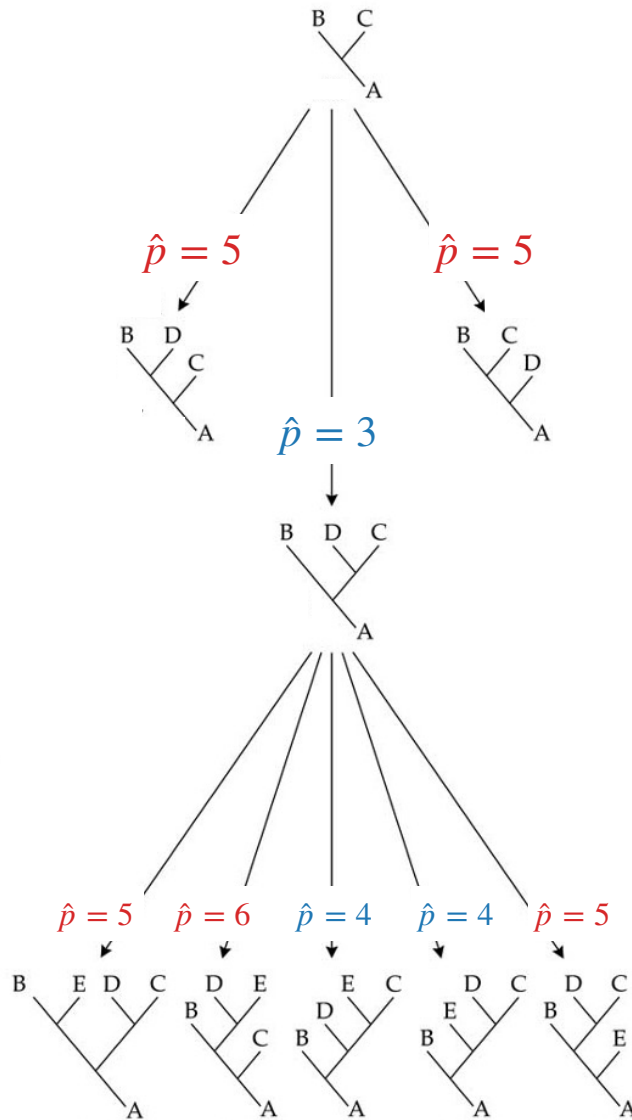
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$

Input characters:

A = (1, 1, 1)

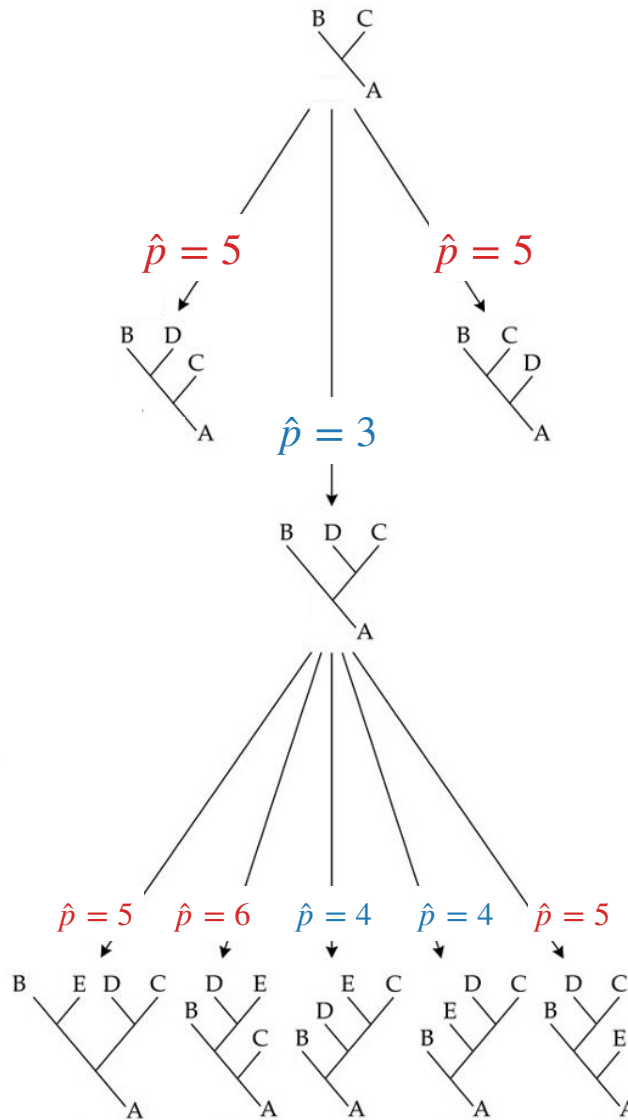
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$

4. Return T

Input characters:

A = (1, 1, 1)

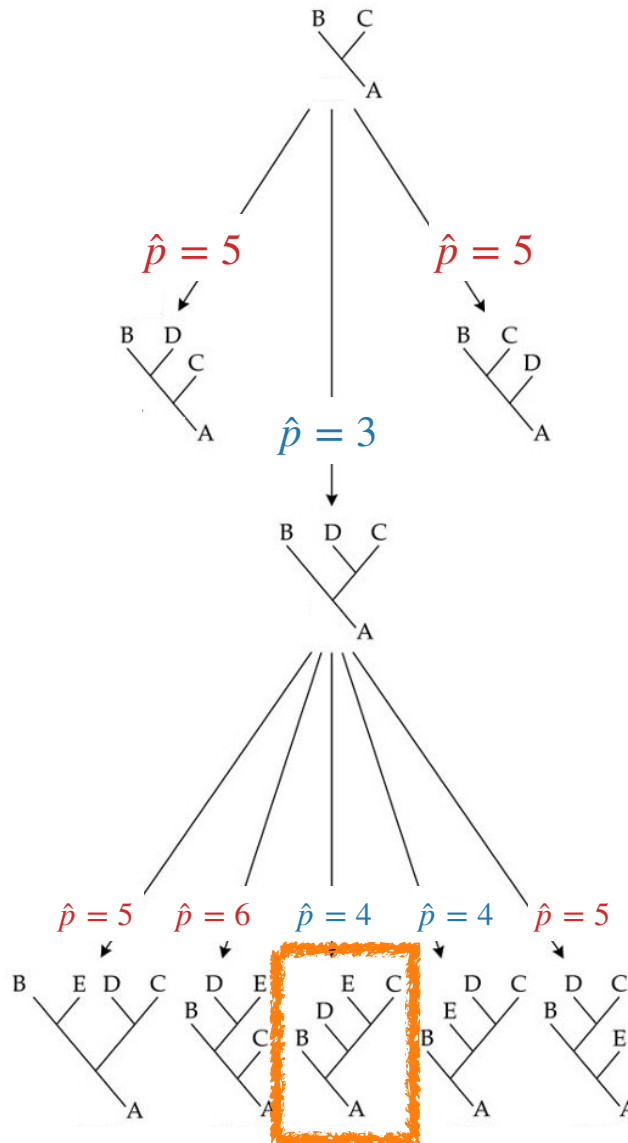
B = (1, 1, 1)

C = (0, 0, 0)

D = (0, 0, 1)

E = (0, 1, 0)

L = [B, C, A, D, E]



RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$

4. Return T

Think about — how many calls will you make to the length function?

Give your answer in Big Oh, where

- n is the # of labels
- m is the #r of characters

RandomizedTaxonAddition(\mathcal{C}, S):

1. $L \leftarrow$ put labels in S in random order
2. $T \leftarrow$ star tree on first 3 elements of L
3. For each remaining element s in L :
 - a. $p \leftarrow \infty$
 - b. For each $e \in E(T)$:
 - i. $\hat{T} \leftarrow$ Add s to T by subdividing e with new vertex v and creating edge (v, s)
 - ii. $\hat{p} \leftarrow \text{length}(\hat{T}, \mathcal{C})$
 - iii. If $\hat{p} < p$: $T_{\text{save}} \leftarrow \hat{T}$ and $p \leftarrow \hat{p}$
 - c. $T \leftarrow T_{\text{save}}$
4. Return T

MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

Branch-and-Bound:

Like exhaustive search but better

Heuristic (e.g. hill climbing):

Compute a starting tree. Apply some operation to **edit** the tree. Search from new tree if it's score is higher.

But then we need methods to...

1. **Build** a starting tree.

>randomized taxon addition

2. **Edit** an existing tree.

>SPR, NNI, TBR moves

MP is NP-hard.

Now what?

Exhaustive Search:

Evaluate the parsimony score of all trees.

Branch-and-Bound:

Like exhaustive search but better

Heuristic (e.g. hill climbing):

Compute a starting tree. Apply some operation to **edit** the tree. Search from new tree if it's score is higher.

But then we need methods to...

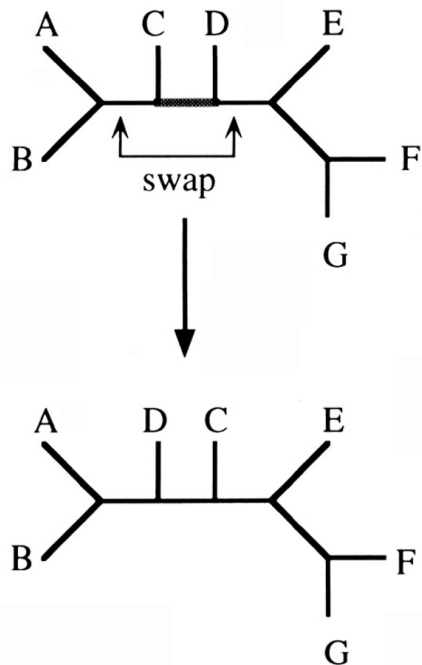
1. **Build** a starting tree.

>randomized taxon addition

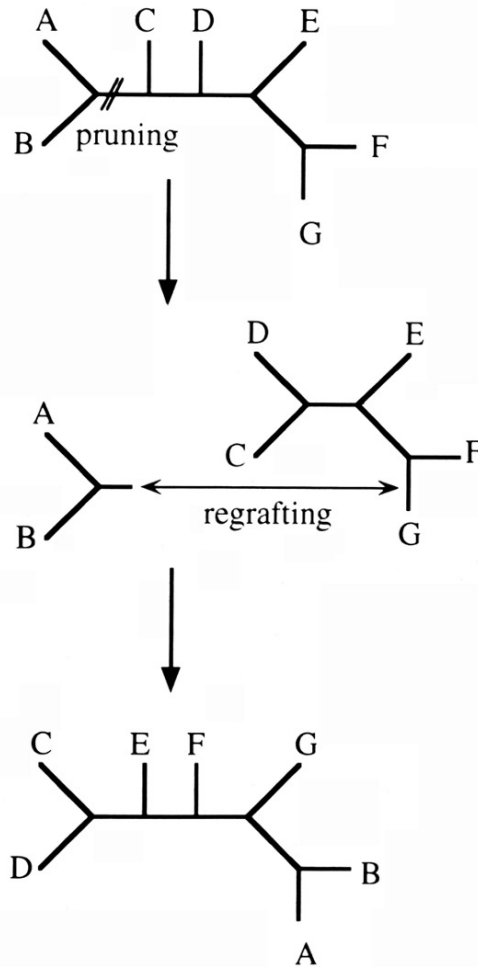
2. **Edit** an existing tree.

>SPR, NNI, TBR moves

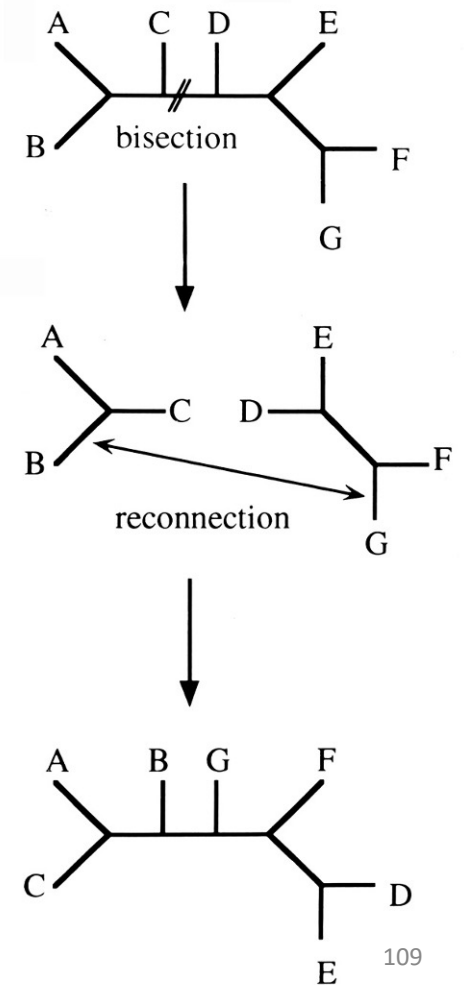
Nearest Neighbor Interchange (NNI) move



Subtree Prune and Regraft (SPR) move



Tree Bisection and Reconnection (TBR) move



Lastly, let's take a closer look at **branch-and-bound**.

Branch-and-Bound Ideas:

1. Take a tree T (e.g. compute a tree with randomized taxon addition) and compute its length L .

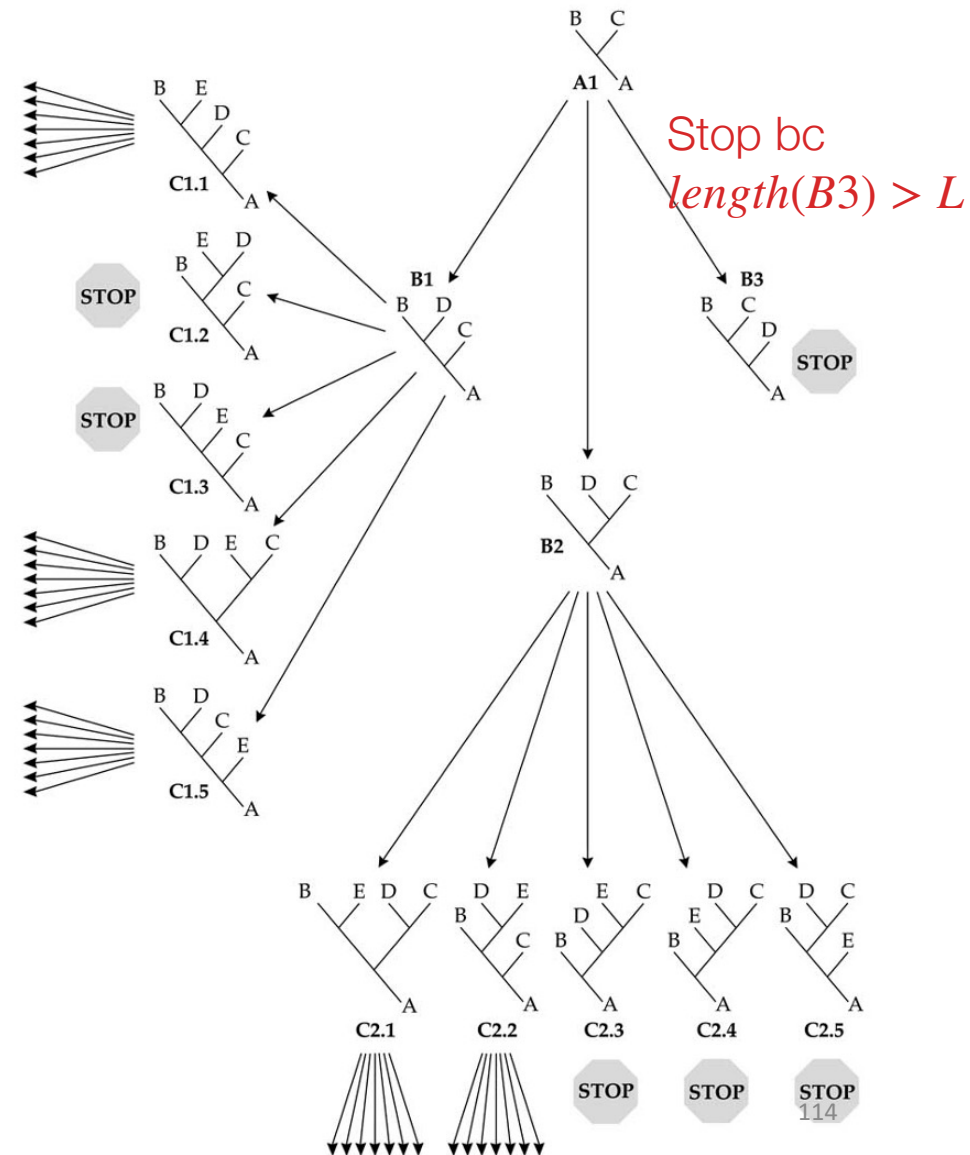
Branch-and-Bound Ideas:

1. Take a tree T (e.g. compute a tree with randomized taxon addition) and compute its length L .
2. An optimal solution to maximum parsimony tree must have length $\leq L$.

Branch-and-Bound Ideas:

1. Take a tree T (e.g. compute a tree with randomized taxon addition) and compute its length L .
2. An optimal solution to maximum parsimony tree must have length $\leq L$.
3. Now suppose you add a taxon x to a tree t , The length of the resulting tree t_x must be $\geq \text{length}(t)$.

1. Take a tree T (e.g. compute a tree with randomized taxon addition) and compute its length L .
2. An optimal solution to maximum parsimony tree must have length $\leq L$.
3. Now suppose you add a taxon x to a tree t , The length of the resulting tree t_x must be $\geq \text{length}(t)$.
4. Therefore, you can enumerate all trees via taxon addition and stop enumerating from a given tree t if $\text{length}(t) > L$.



Agenda

- ✓ **Part 1:** Perfect Phylogenies
- ✓ **Part 2:** Small Parsimony Problem & Fitch's Algorithm
- ✓ **Part 3:** Large Parsimony Problem
- ✓ **Part 4:** Maximum Parsimony Methods

Take-Aways

- Evolutionary trees are reconstructed from data (e.g., characters).
- A perfect phylogeny does not always exist for the data.
- The goal of parsimony is to find the tree that offers the simplest explanation of our data (i.e., minimum substitutions).
- The parsimony score of a given tree can be computed in polynomial time but finding a tree so that the score is minimized is NP-hard.
- Whether it makes sense to reconstruct a tree using maximum parsimony depends on the **model of evolution** — take my undergrad (498Y) or grad class (829A) to learn more!
- You covered a basic model of evolution (substitutions only) in the textbook and considered how to compute likelihood under this model (similar ideas apply for scoring, hardness, and heuristics).