

# Statistics in BHV Tree Space

Megan Owen

Lehman College, CUNY

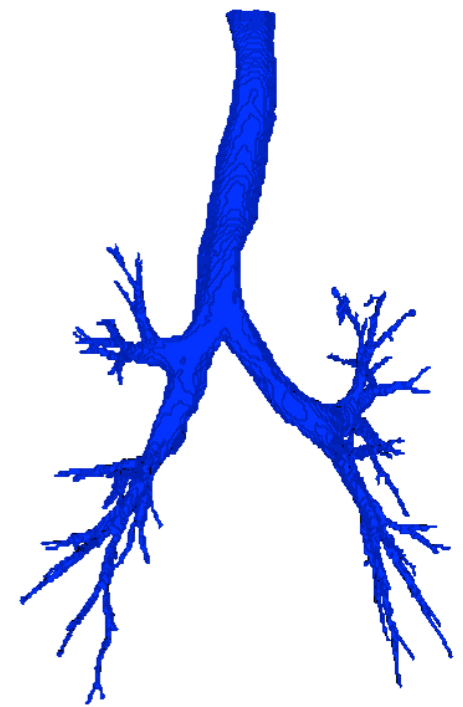
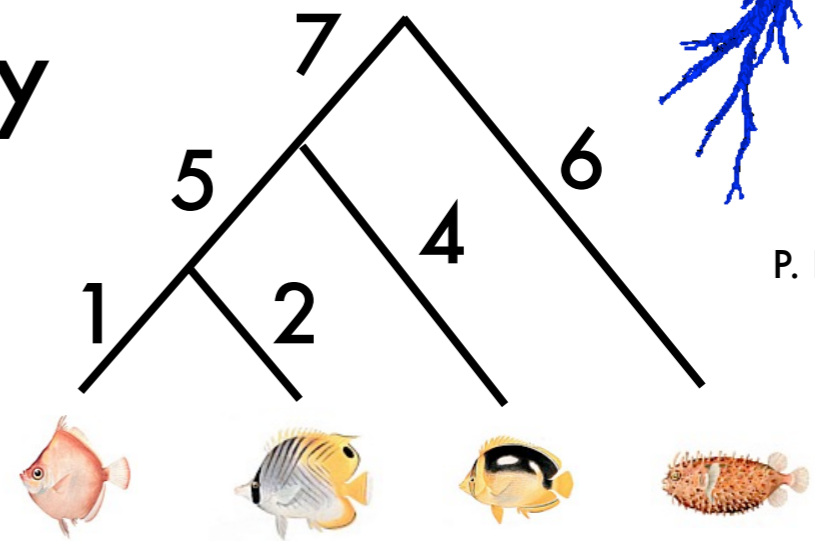
Joint with Daniel Brown (U. Waterloo)

BIRS Workshop on Mathematical Approaches to Evolutionary Trees and Networks

Feb. 13, 2017

# Overarching Goal

- many examples of tree-shaped data (phylogenies, anatomical trees, etc.)
- parameters:
  - tree shape = tree topology
  - edge lengths
- not Euclidean data!

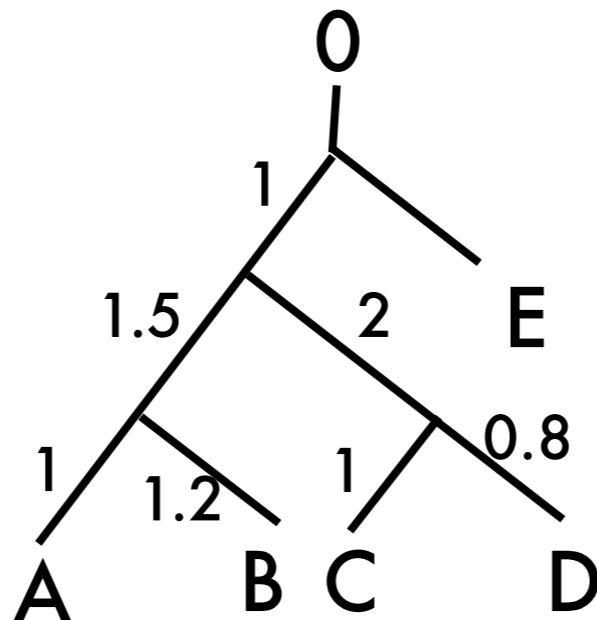


P. Lo et al. EXACT'09

**Goal:** develop methods for statistical analysis (i.e. mean, PCA) in a space of metric trees analogous to those for Euclidean space

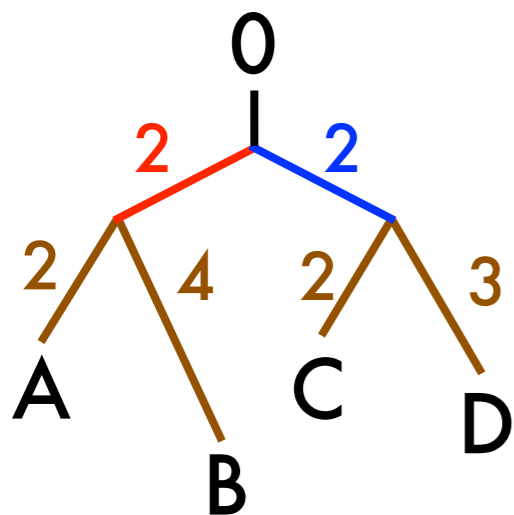
# Tree Space Framework

- constructed by Billera, Holmes, and Vogtmann (2001)
- tree space  $\mathbb{T}_n$  = set of all trees with  $n$  leaves and branch lengths
- includes degenerate trees (non-binary)

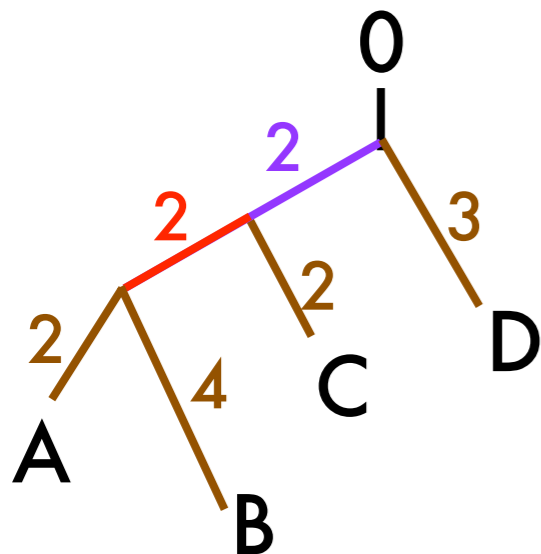


# Tree Space

- represent each tree as a vector
- coordinates = splits



$($ 
A|BCDO  
2, 
 B|ACDO  
4, 
 ...  
D|ABCO  
3, 
 A|BCDO  
2, 
 AC|BDO  
0, 
 0, 
 0, 
 ...  
ABC|DO  
0, 
 ABD|CO  
2, 
 ...
 $)$



$($ 
A|BCDO  
2, 
 B|ACDO  
4, 
 ...  
D|ABCO  
3, 
 A|BCDO  
2, 
 AC|BDO  
0, 
 0, 
 0, 
 ...  
ABC|DO  
2, 
 ABD|CO  
0, 
 ABO|CD  
...
 $)$

# Tree Space

- not all sets of splits form a tree
  - ⇒ not all vectors are possible
  - ⇒ not a Euclidean space

*A|BCDO*  
*B|ACDO* ... *D|ABCO*  
*AB|CDO*  
*AC|BDO* ... *ABC|DO*  
*ABD|CO*  
*ABO|ICD*

(2, 4, 2, 3, 2, 2, 0, 0, 0, 0, 0, ...)

# Tree Space

- not all sets of splits form a tree
  - ⇒ not all vectors are possible
  - ⇒ not a Euclidean space

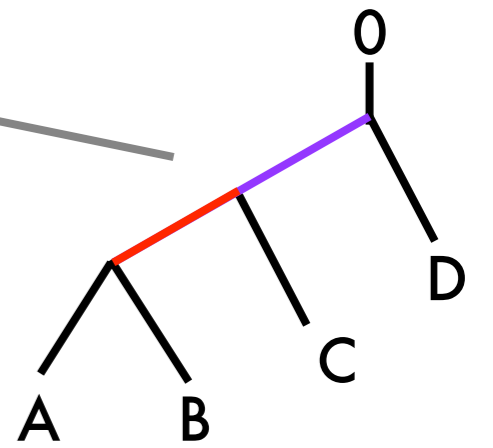
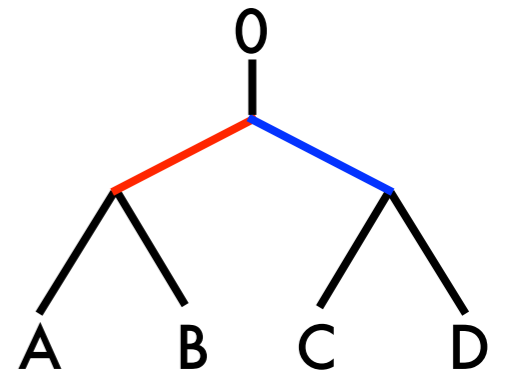
~~*A|BCDO*  
*B|ACDO*  
...  
*D|ABCO*  
*A|B|CDO*  
*AC|BDO*  
...  
*ABC|DO*  
*ABD|CO*  
*ABO|ICD*~~  
(2, 4, 2, 3, 2, 2, 0, 0, 0, 0, 0, ...)

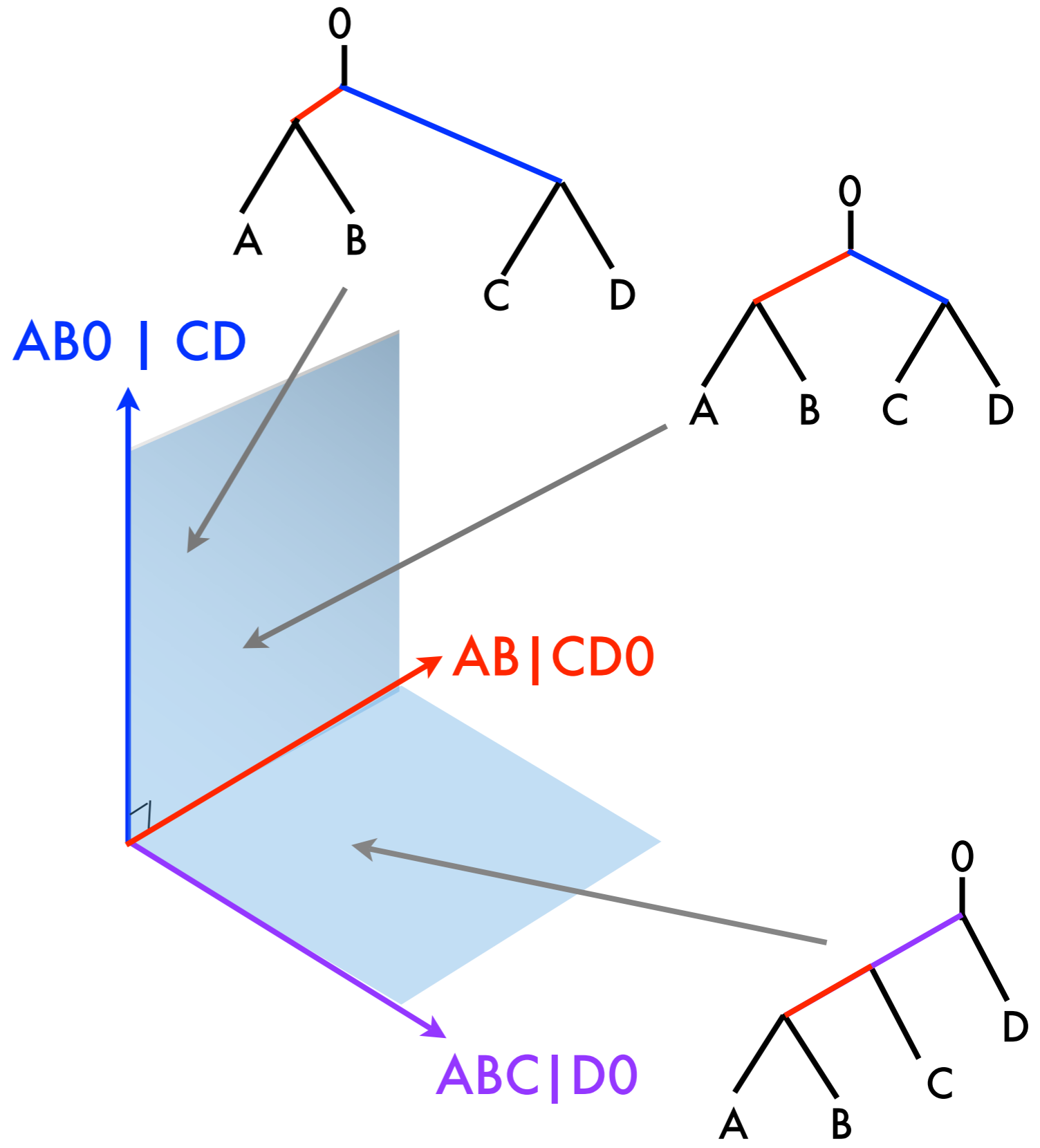
AB0 | CD



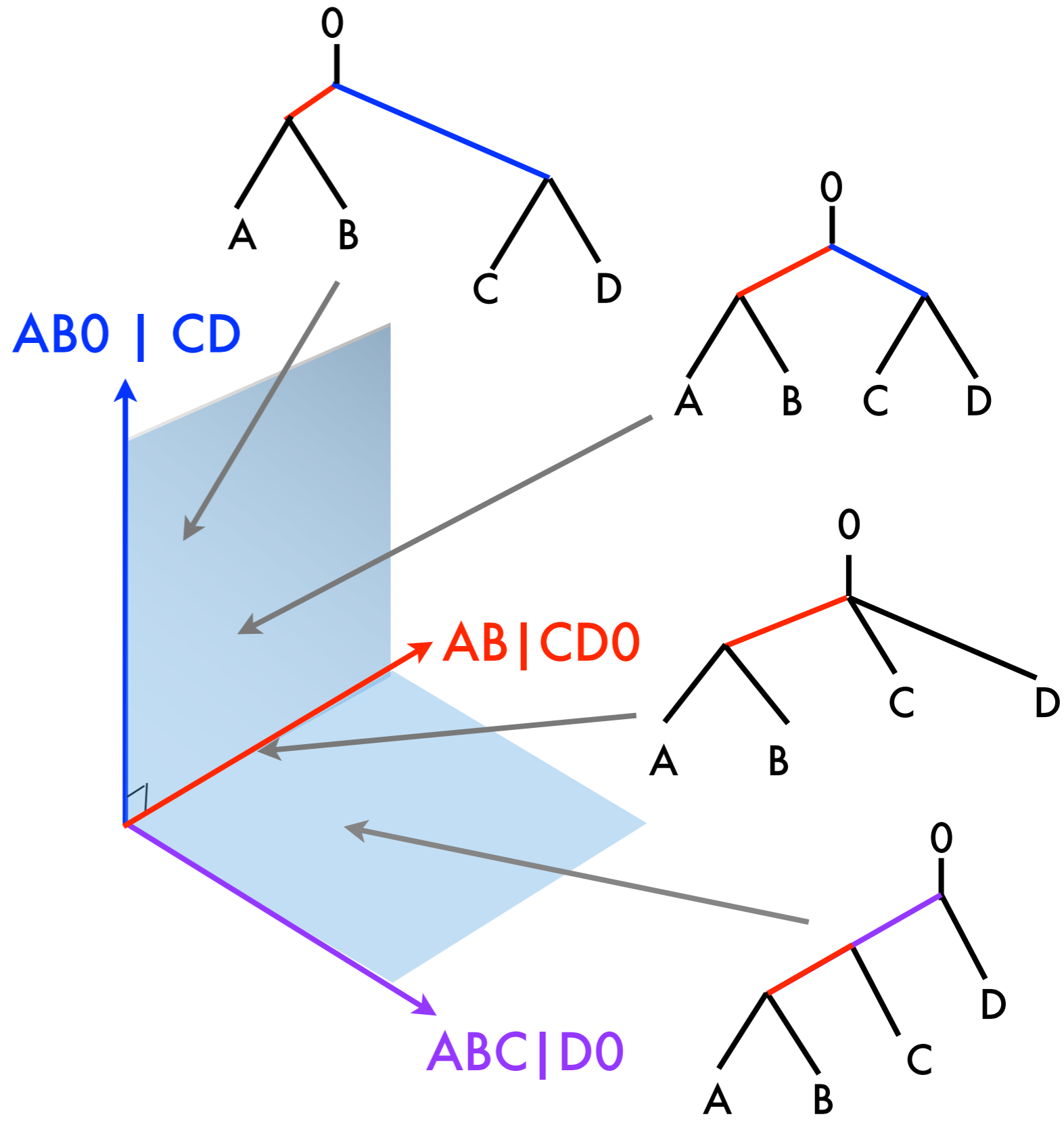
AB | CD0

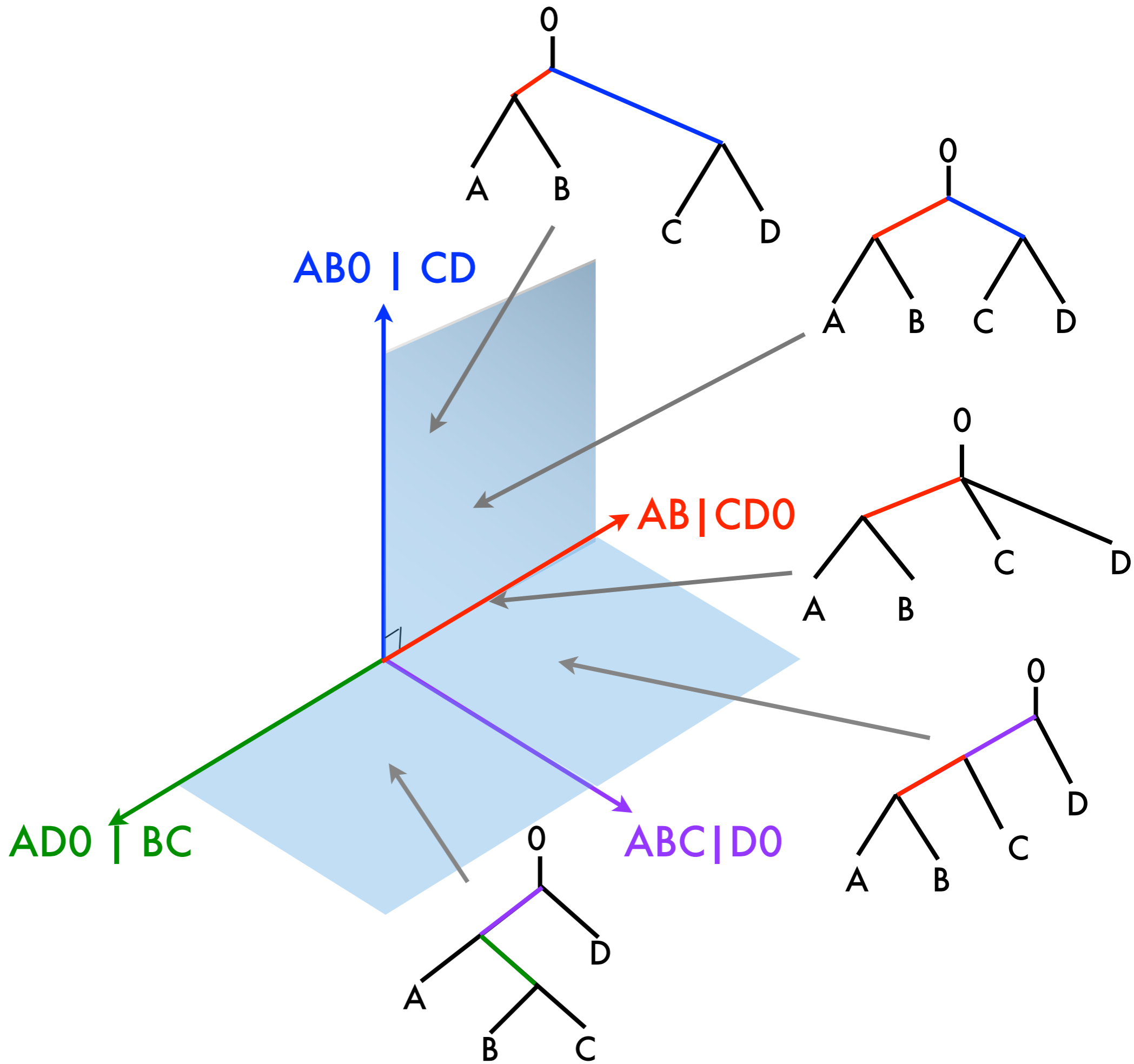
ABC | D0

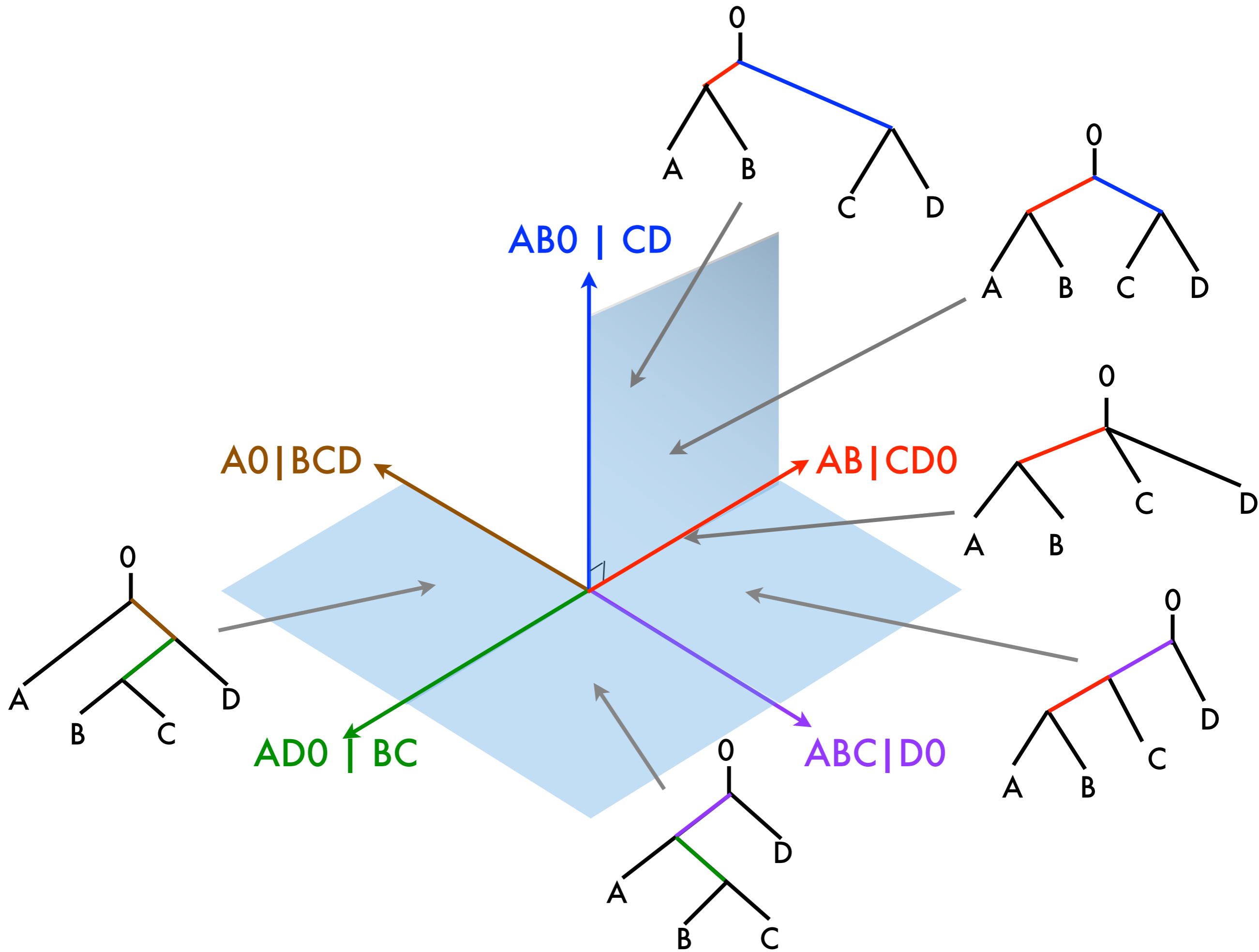


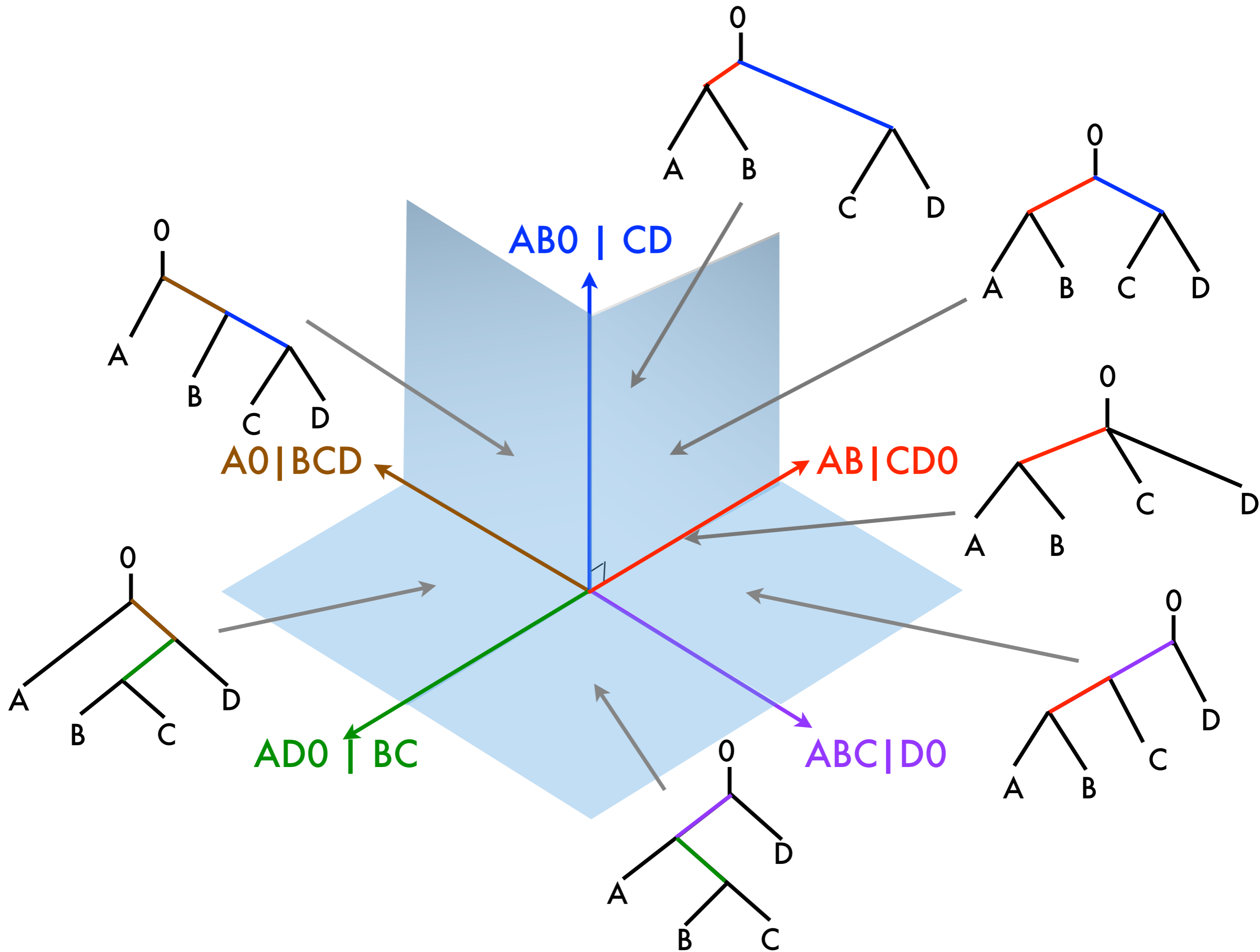


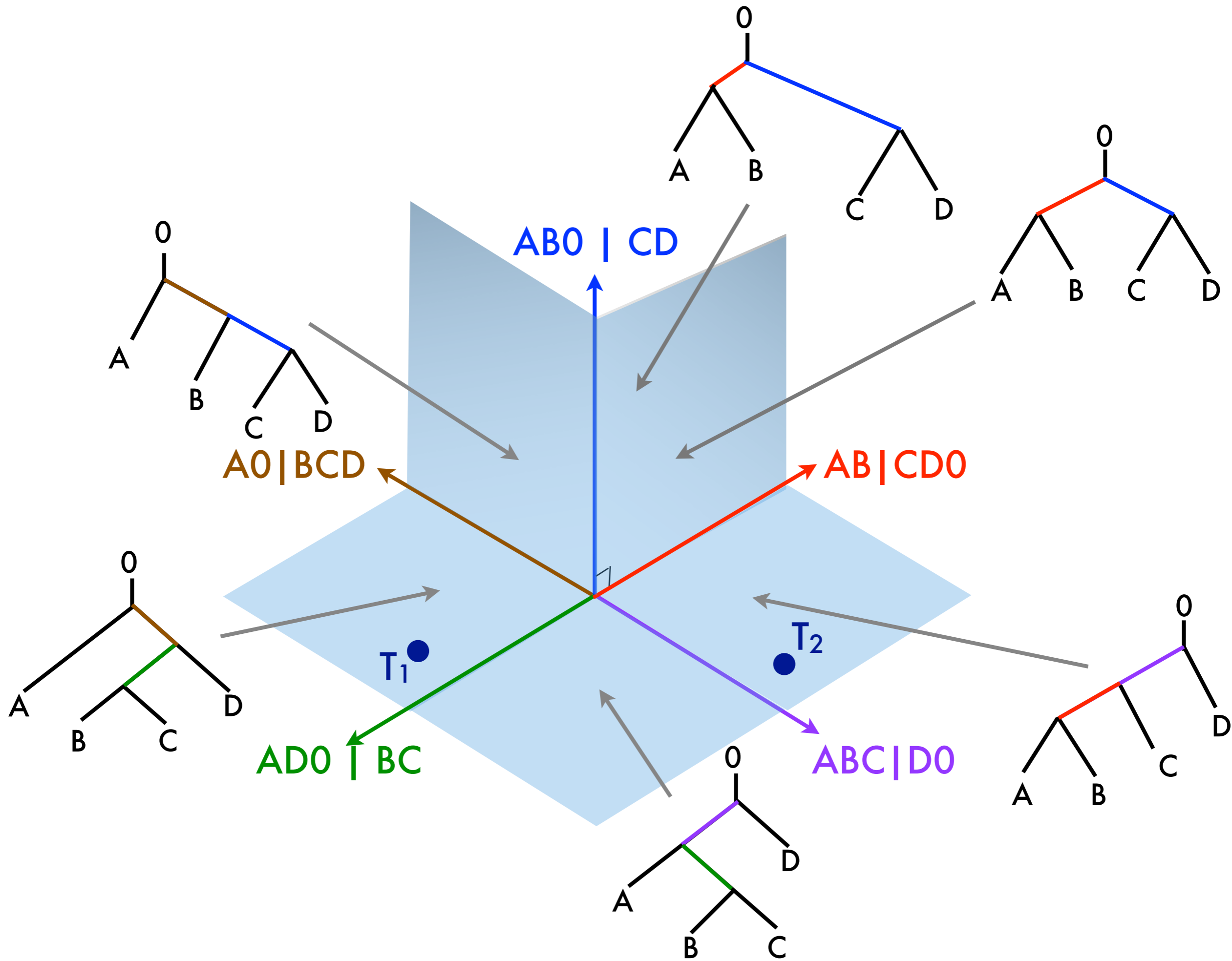


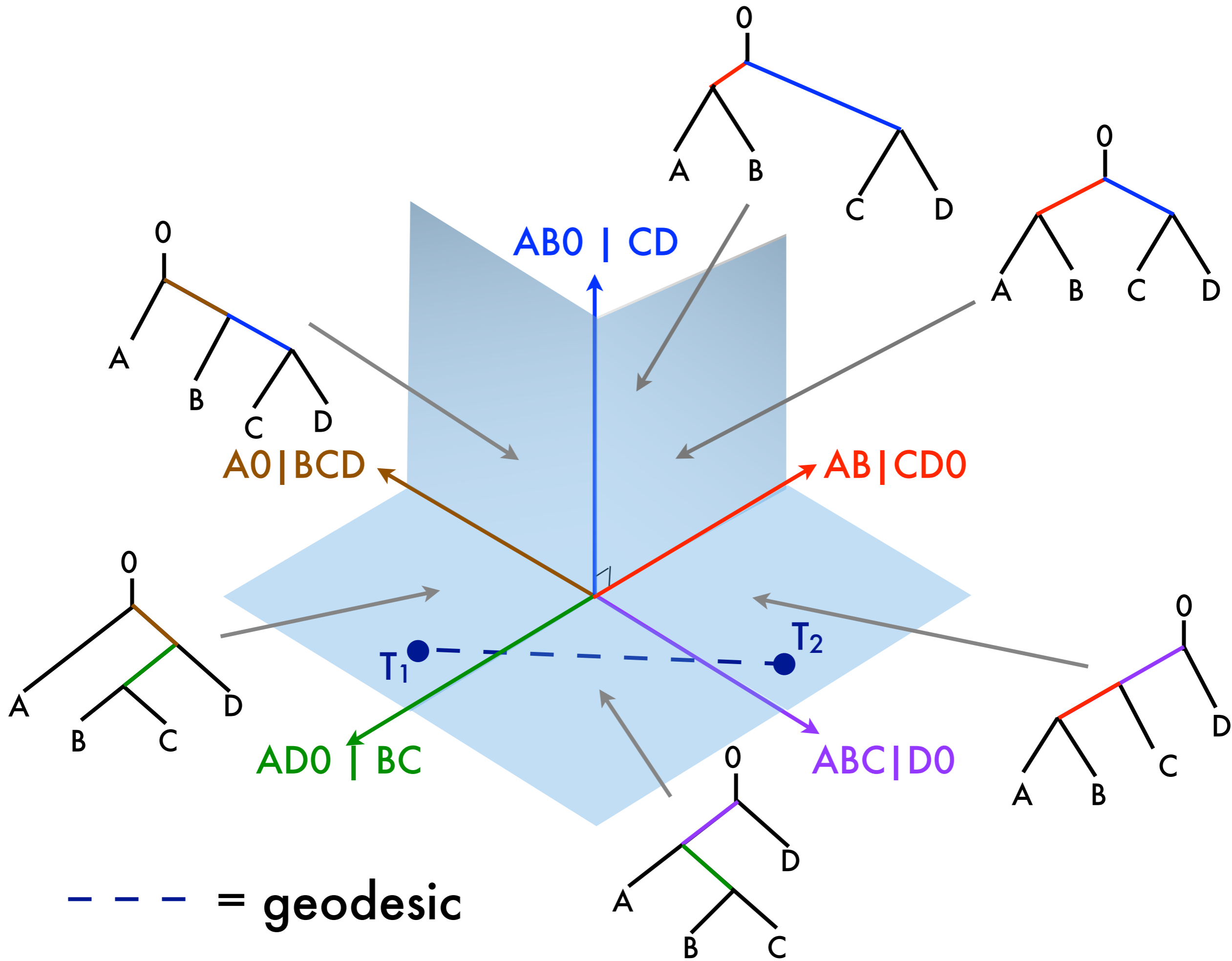


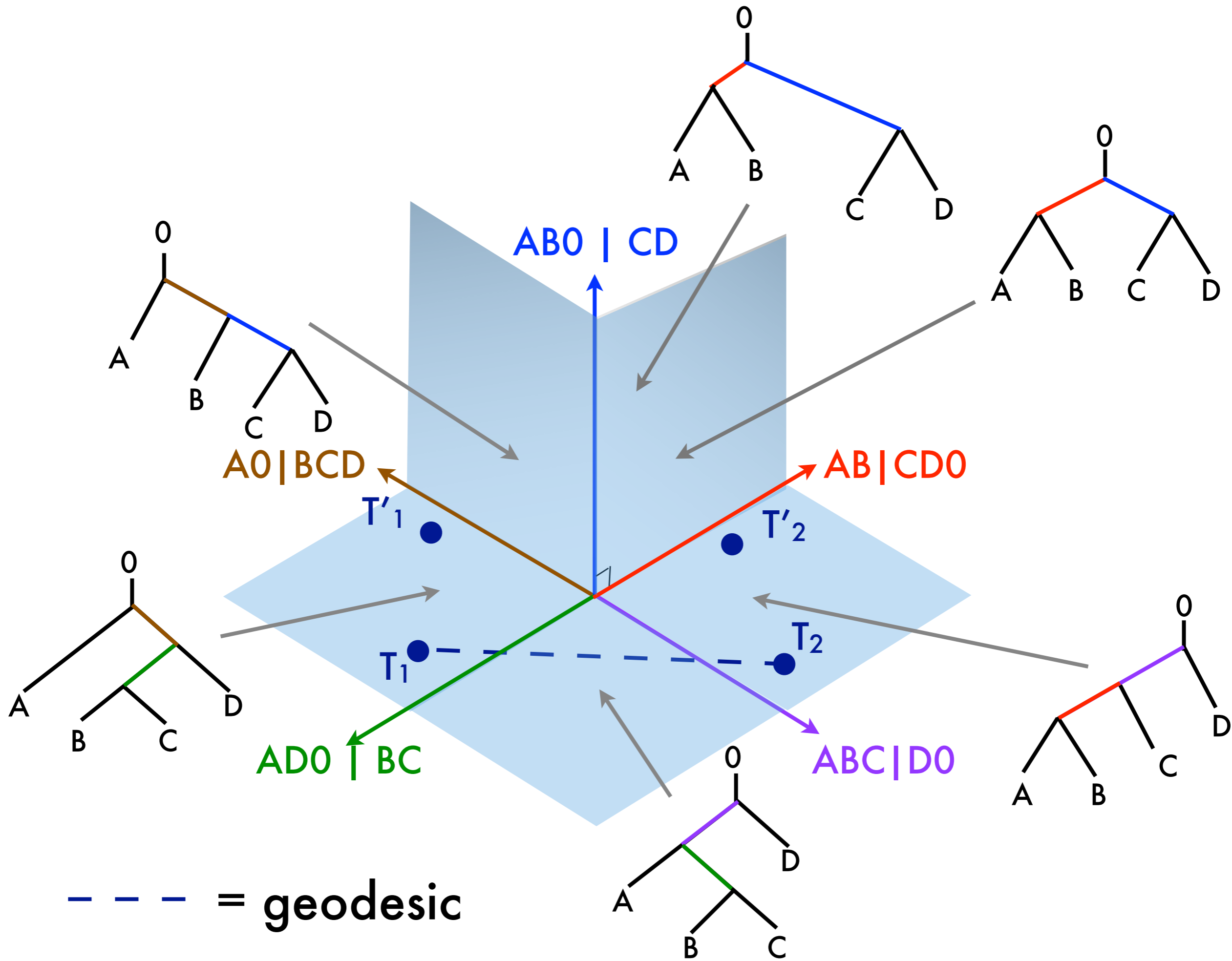


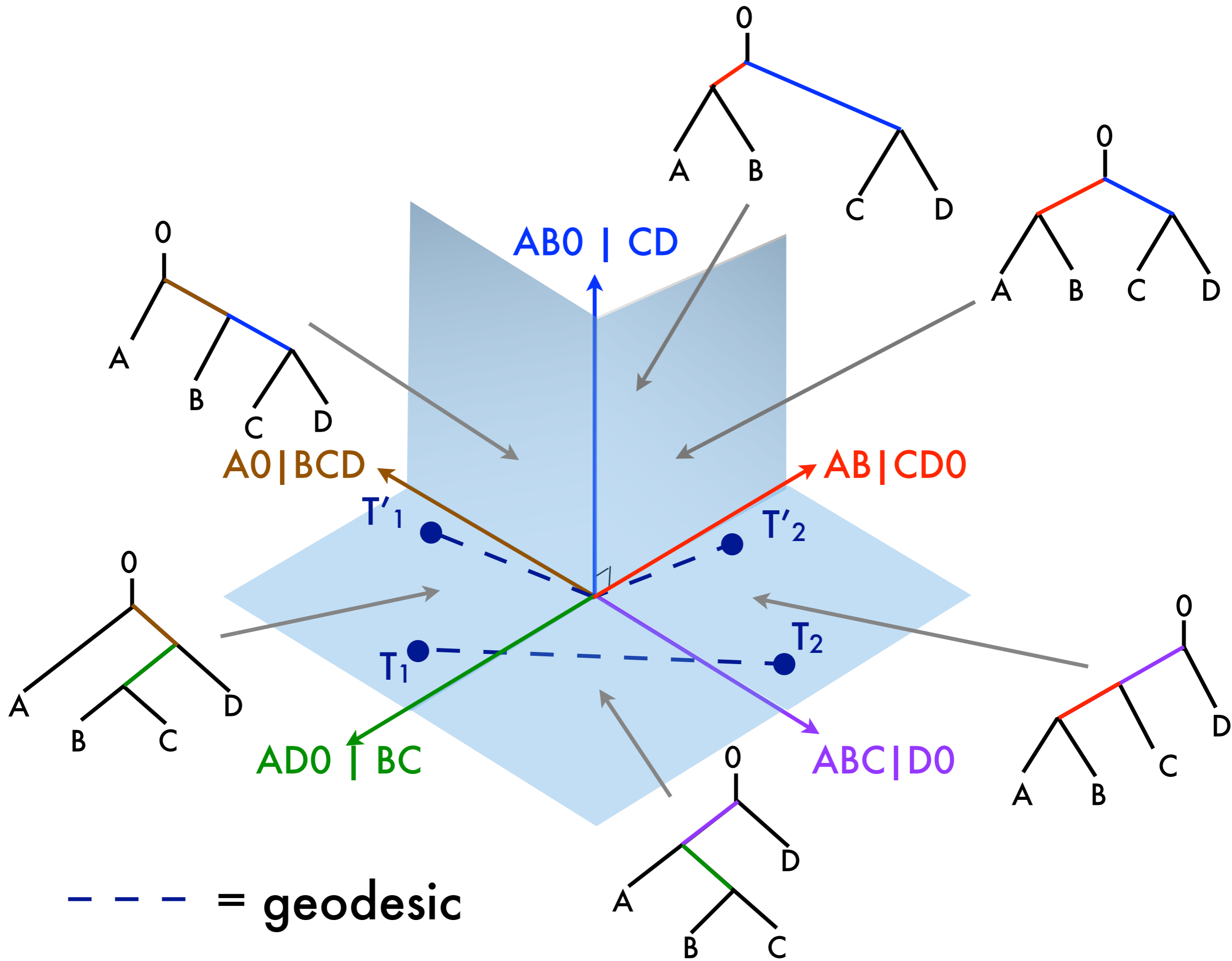






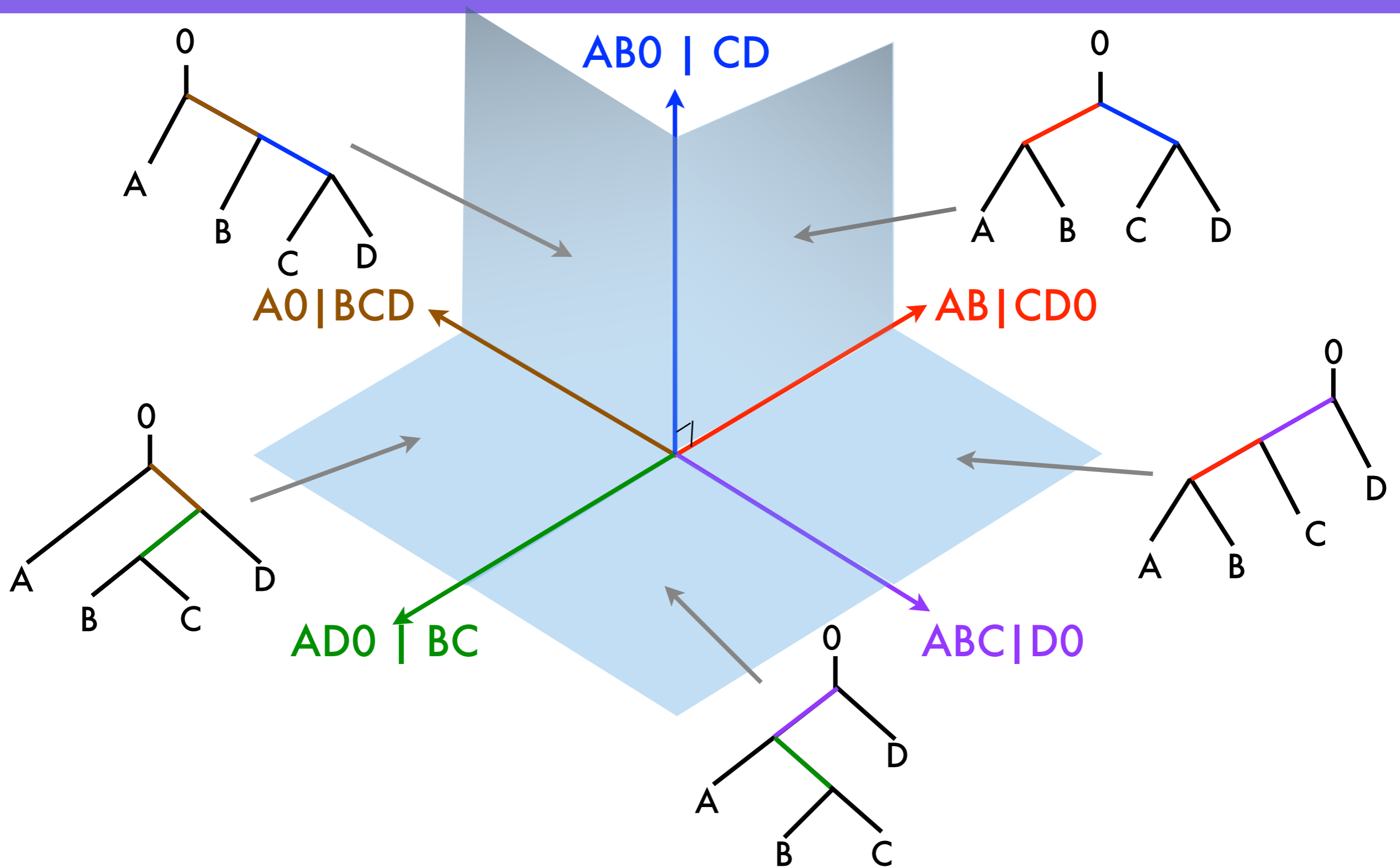




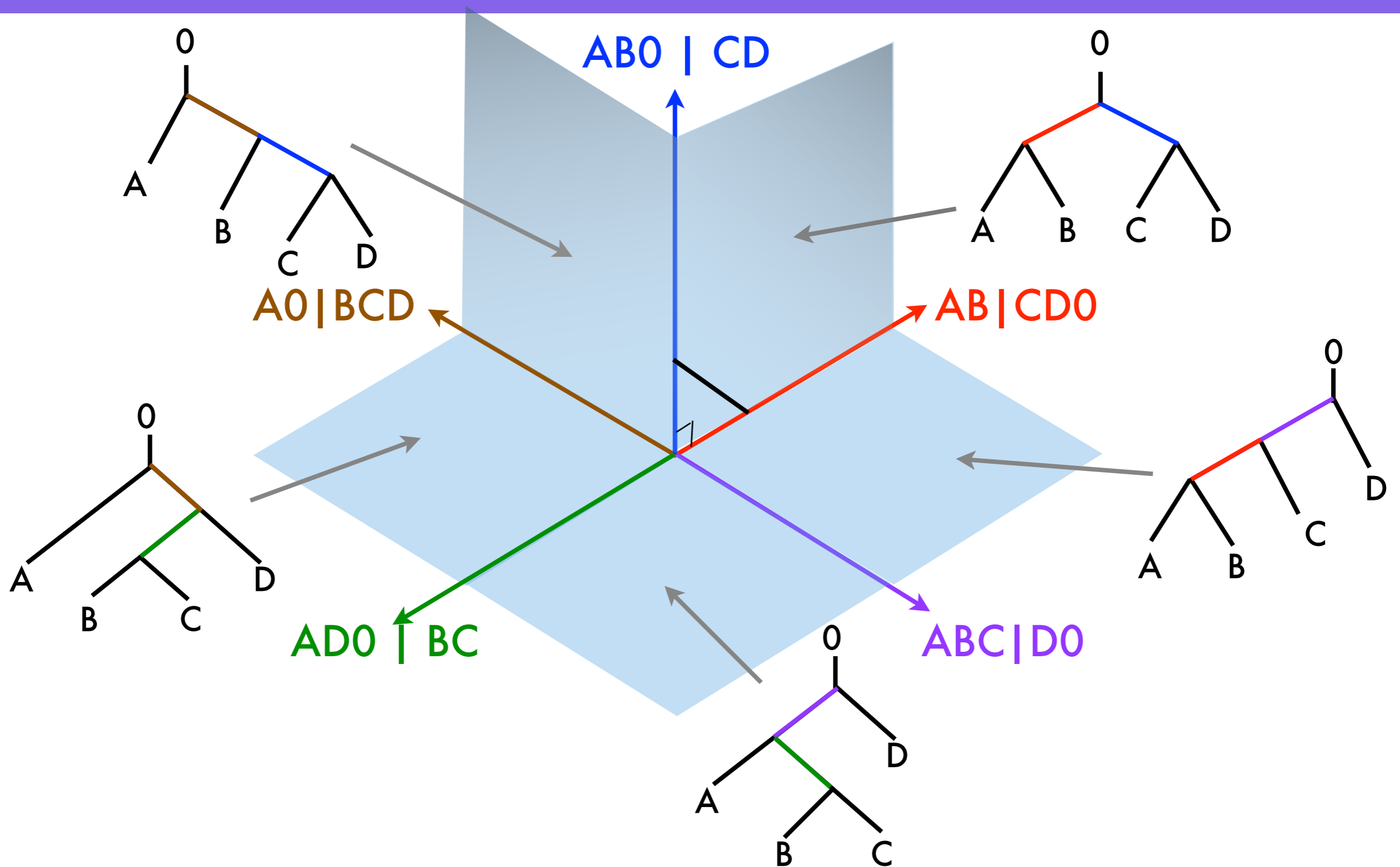




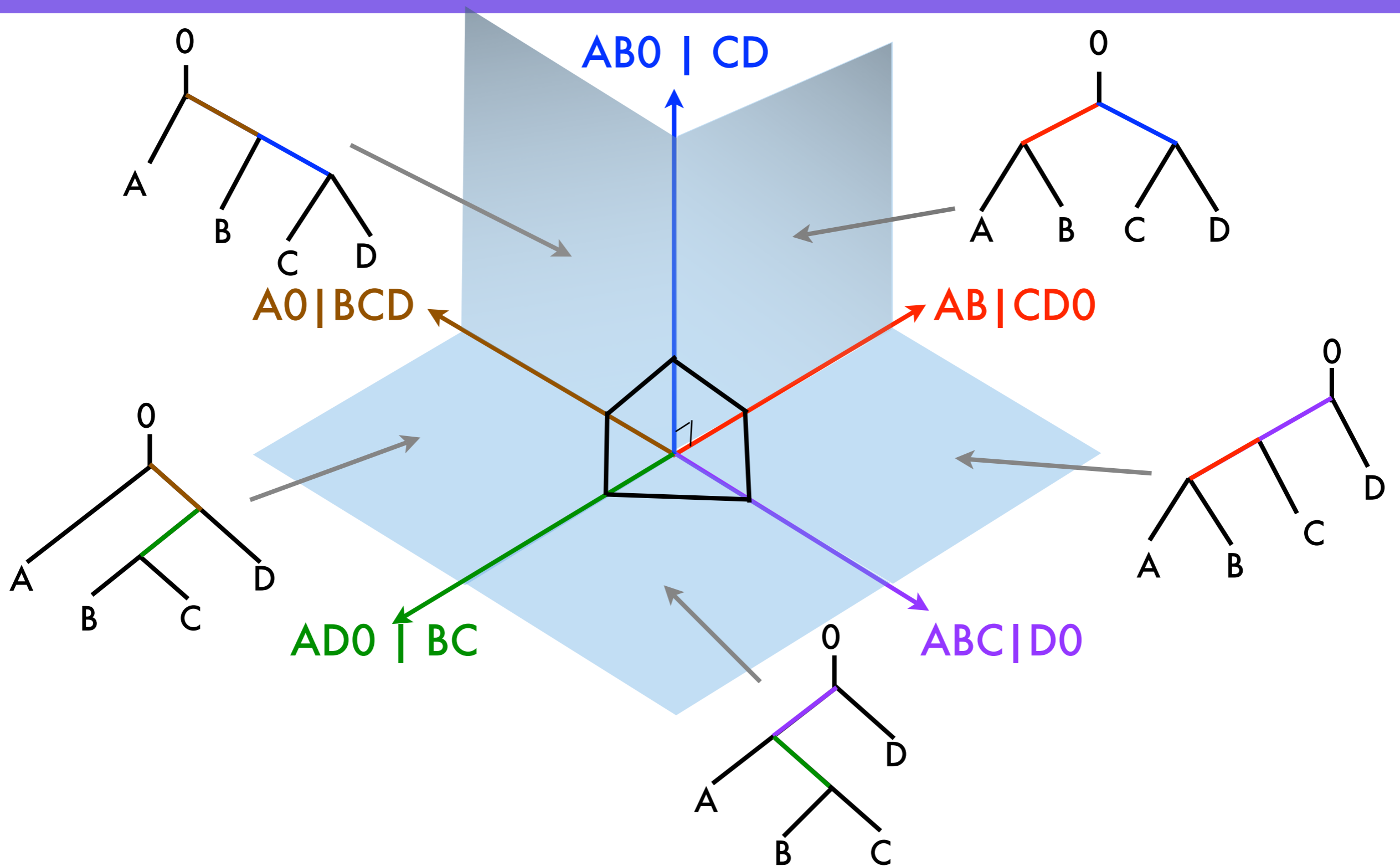
# Structure of $T_4$



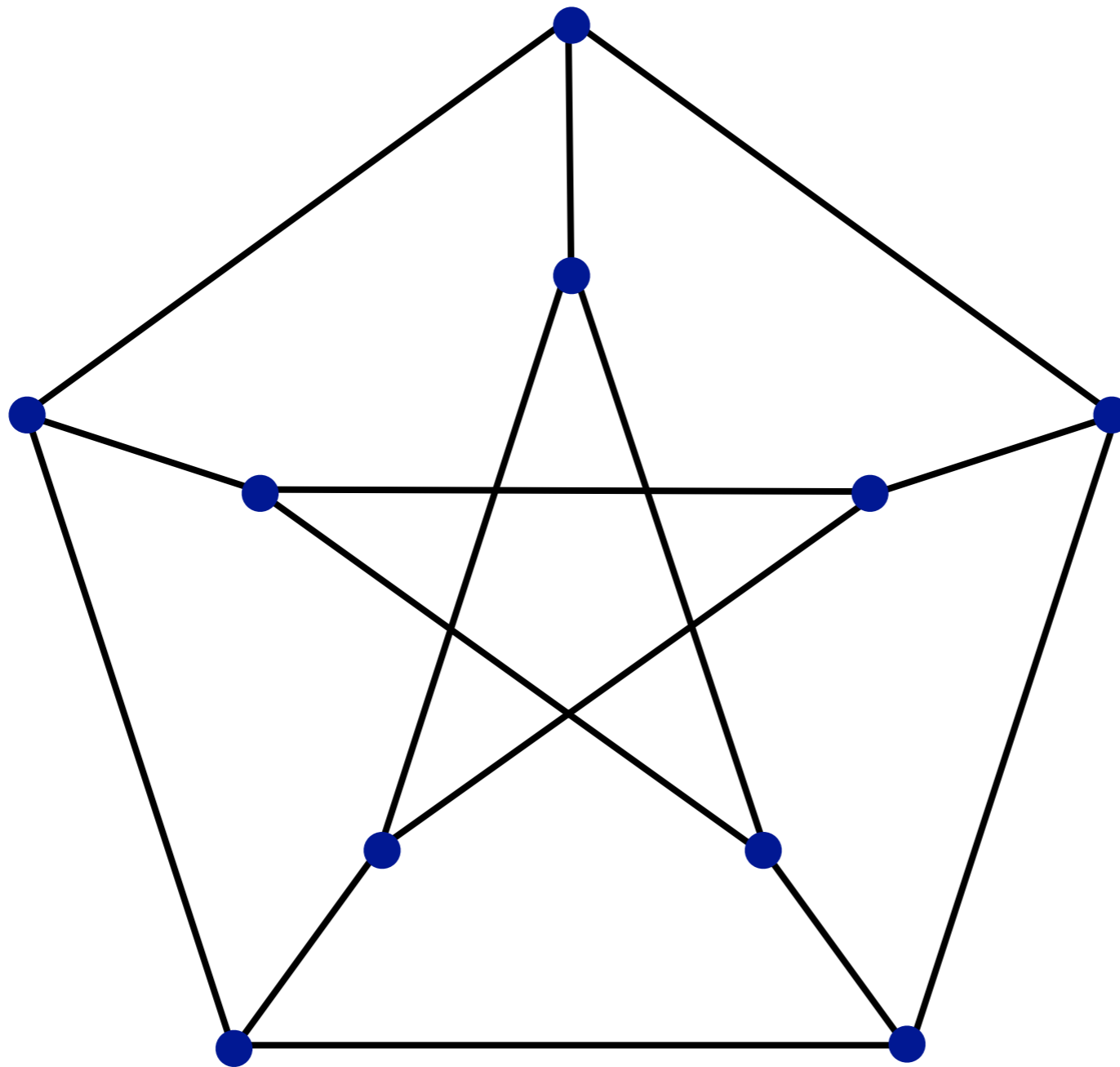
# Structure of $T_4$



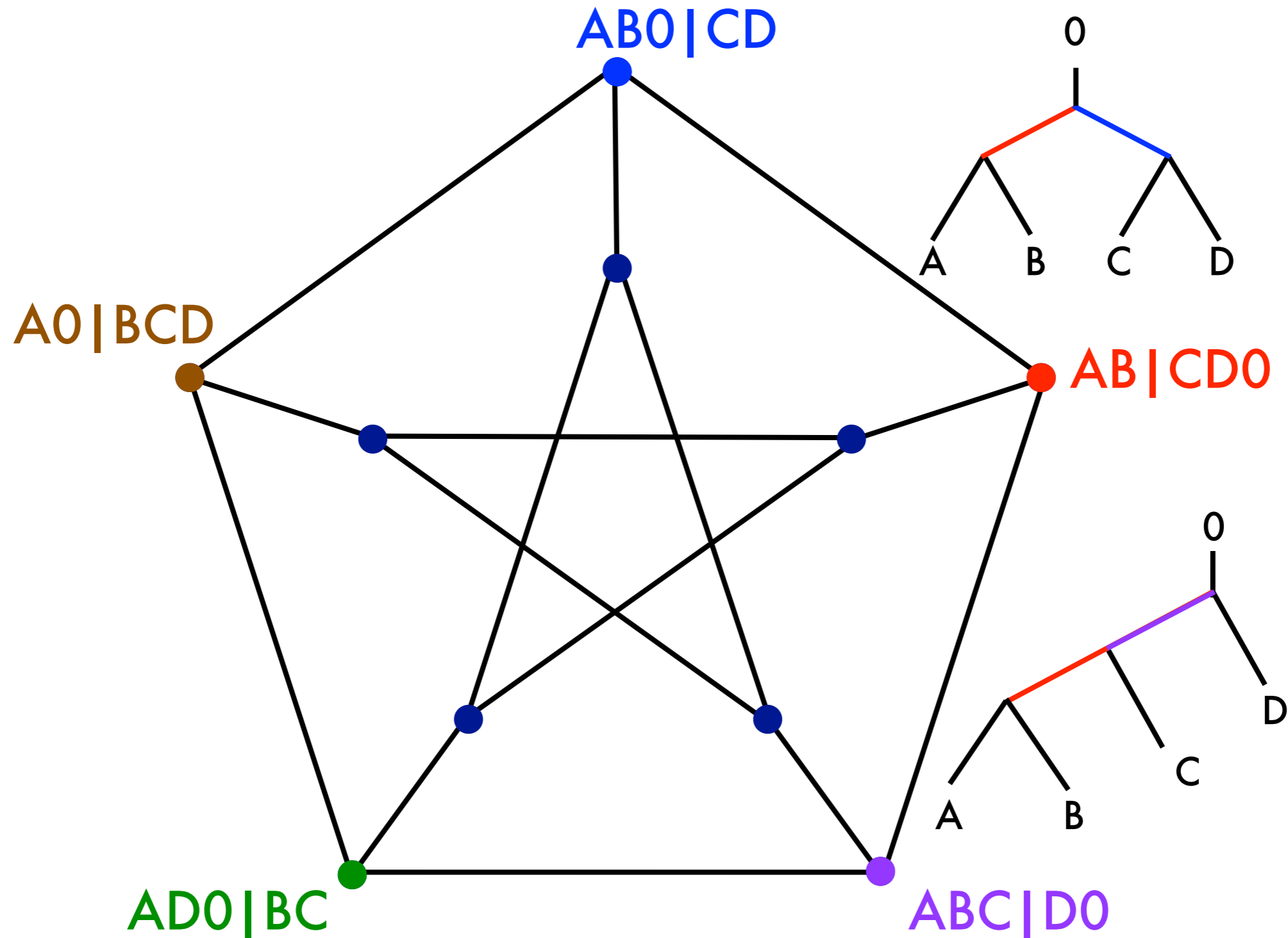
# Structure of $T_4$



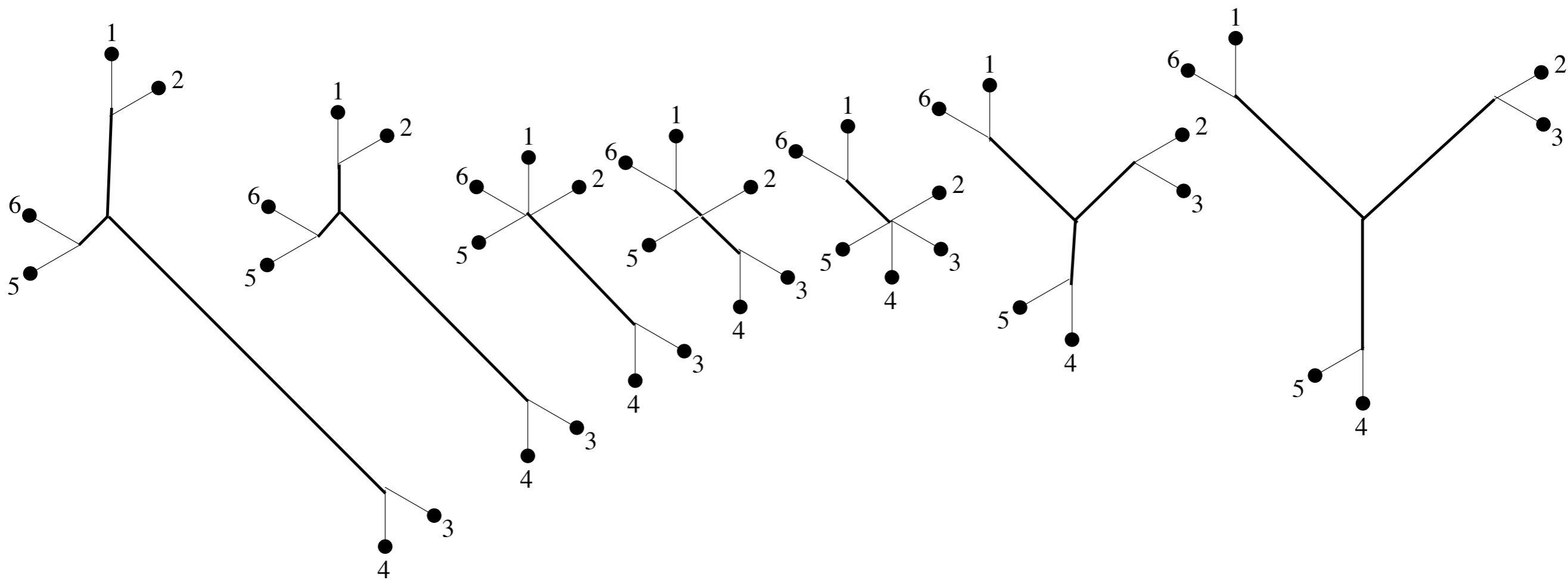
# Structure of $T_4$



# Structure of $T_4$



# Geodesics



# Tree Space Properties

**Theorem** (Billera, Holmes, Vogtmann, 2001):  
Tree space has global non-positive curvature.

⇒ unique geodesics (shortest paths)

⇒ well-defined mid-point tree

- **BHV or geodesic distance** = length of shortest path between two trees  $T_1$  and  $T_2$
- polynomial time algorithm to compute geodesic distance (O. and Provan, 2011)

# Mean and Variance

- weighted set  $X$  in tree space:

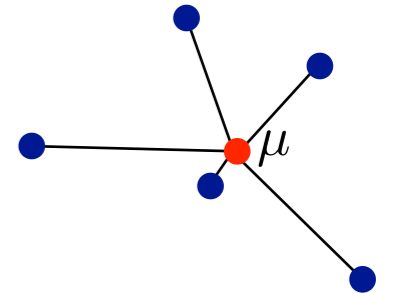
- Fréchet mean( $X$ ) = centre of mass

$$= \operatorname{argmin}_{\mu} \sum_{x \in X} p(x) d(x, \mu)^2$$

(tree minimizing sum of square BHV distances)

- variance( $X$ ) =  $\sum_{x \in X} p(x) d(x, \mu)^2$

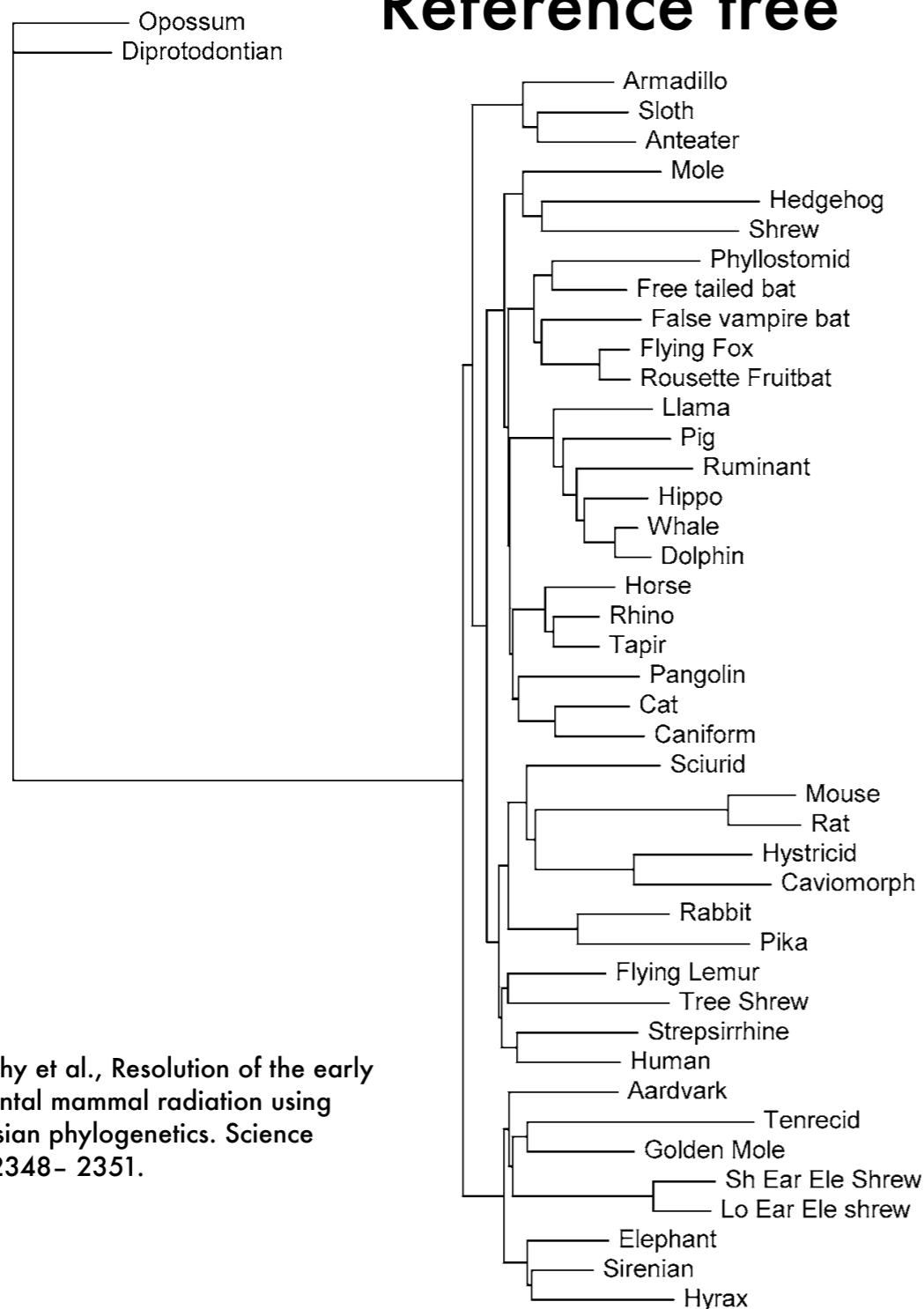
- computable by algorithm based on Law of Large Numbers (Sturm 2003; Miller, O, Provan 2015; Bačák 2014)





# Experimental Results

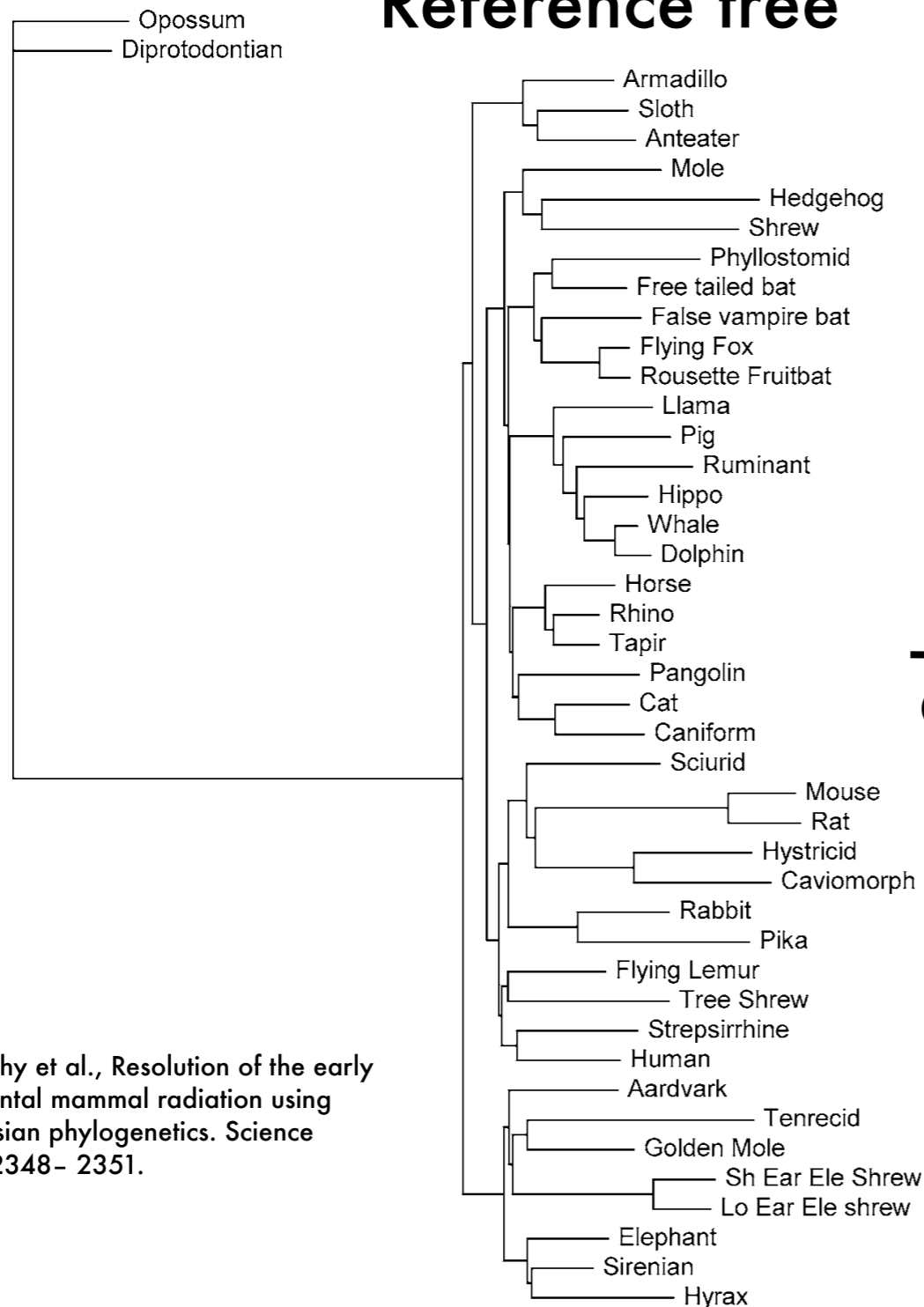
## Reference tree



Murphy et al., Resolution of the early placental mammal radiation using Bayesian phylogenetics. *Science* 294:2348- 2351.

# Experimental Results

## Reference tree



Murphy et al., Resolution of the early placental mammal radiation using Bayesian phylogenetics. Science 294:2348- 2351.

Seq-gen  
GTR model

Simulated DNA sequences

500 bp x 10

1000 bp x 10

⋮

4000 bp x 10

# Experimental Results

**Simulated DNA  
sequences**

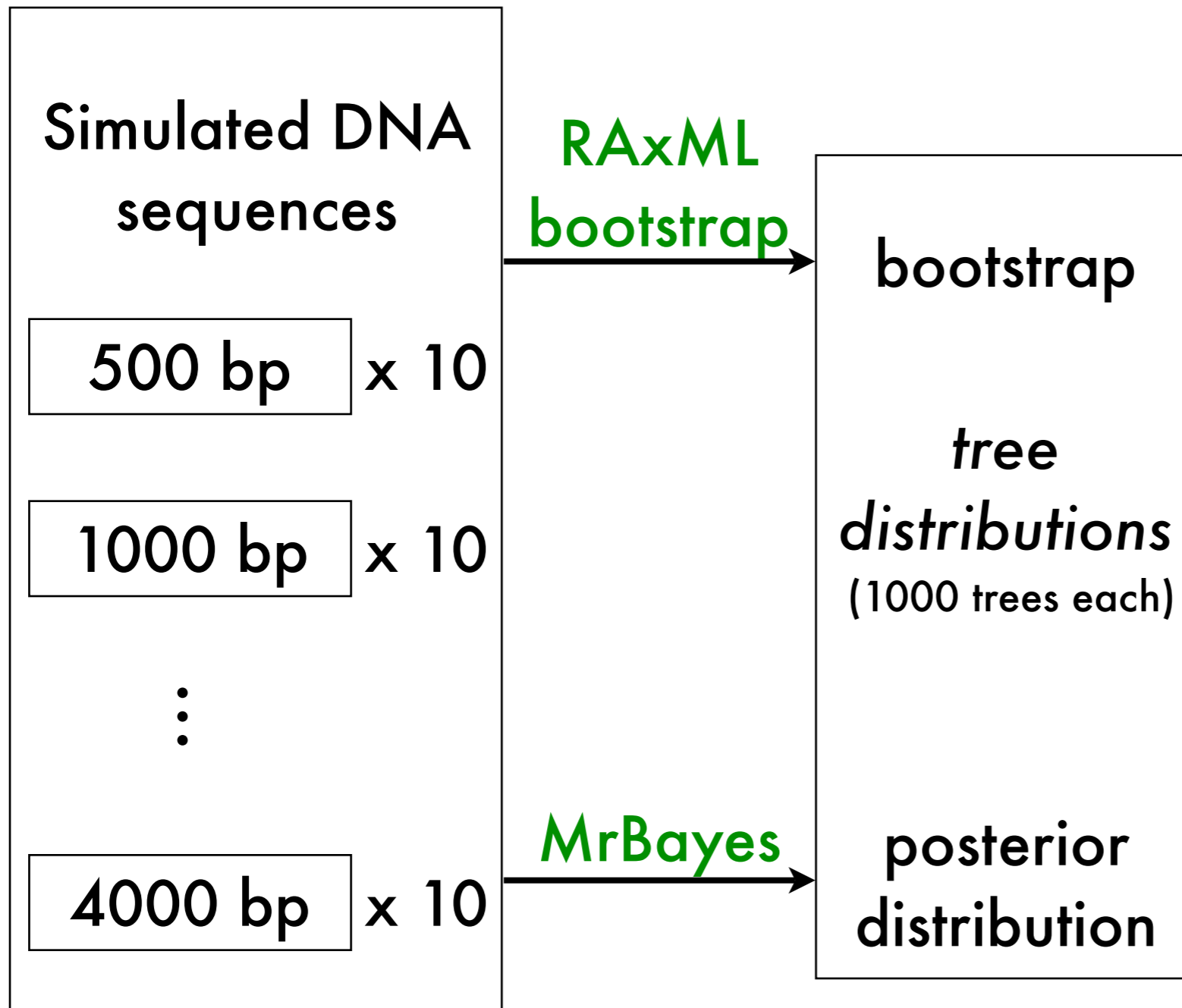
**500 bp x 10**

**1000 bp x 10**

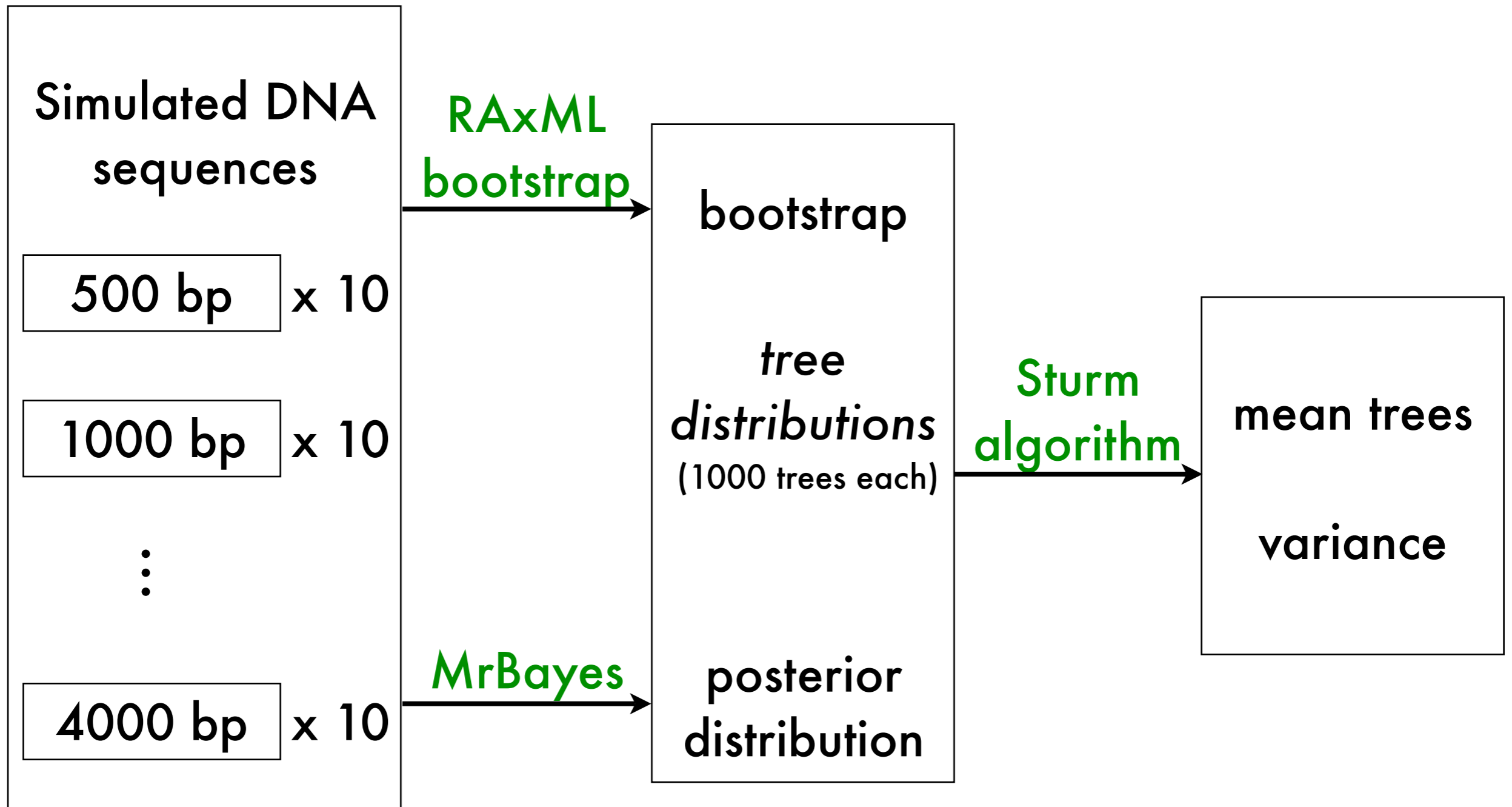
**⋮**

**4000 bp x 10**

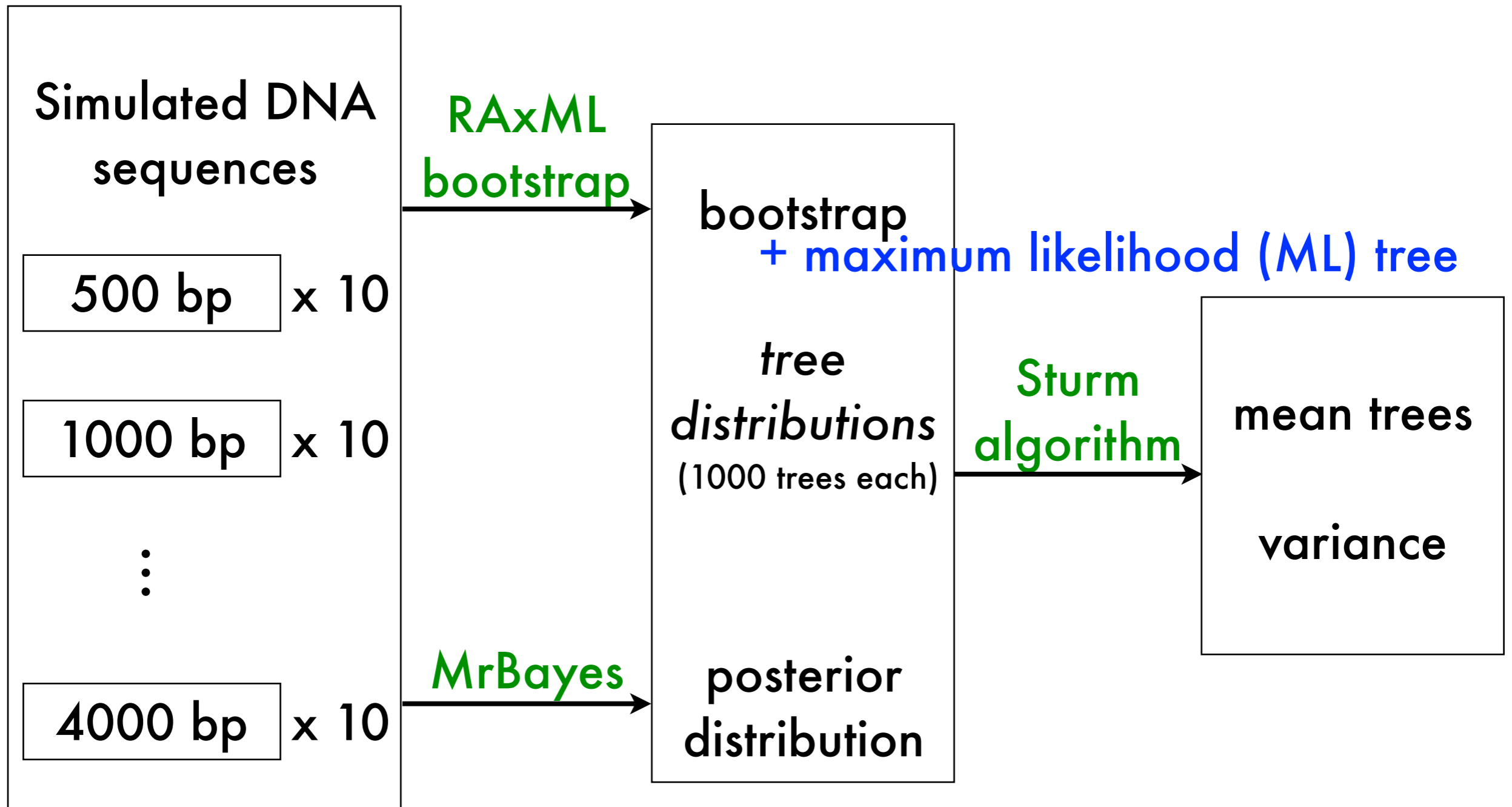
# Experimental Results



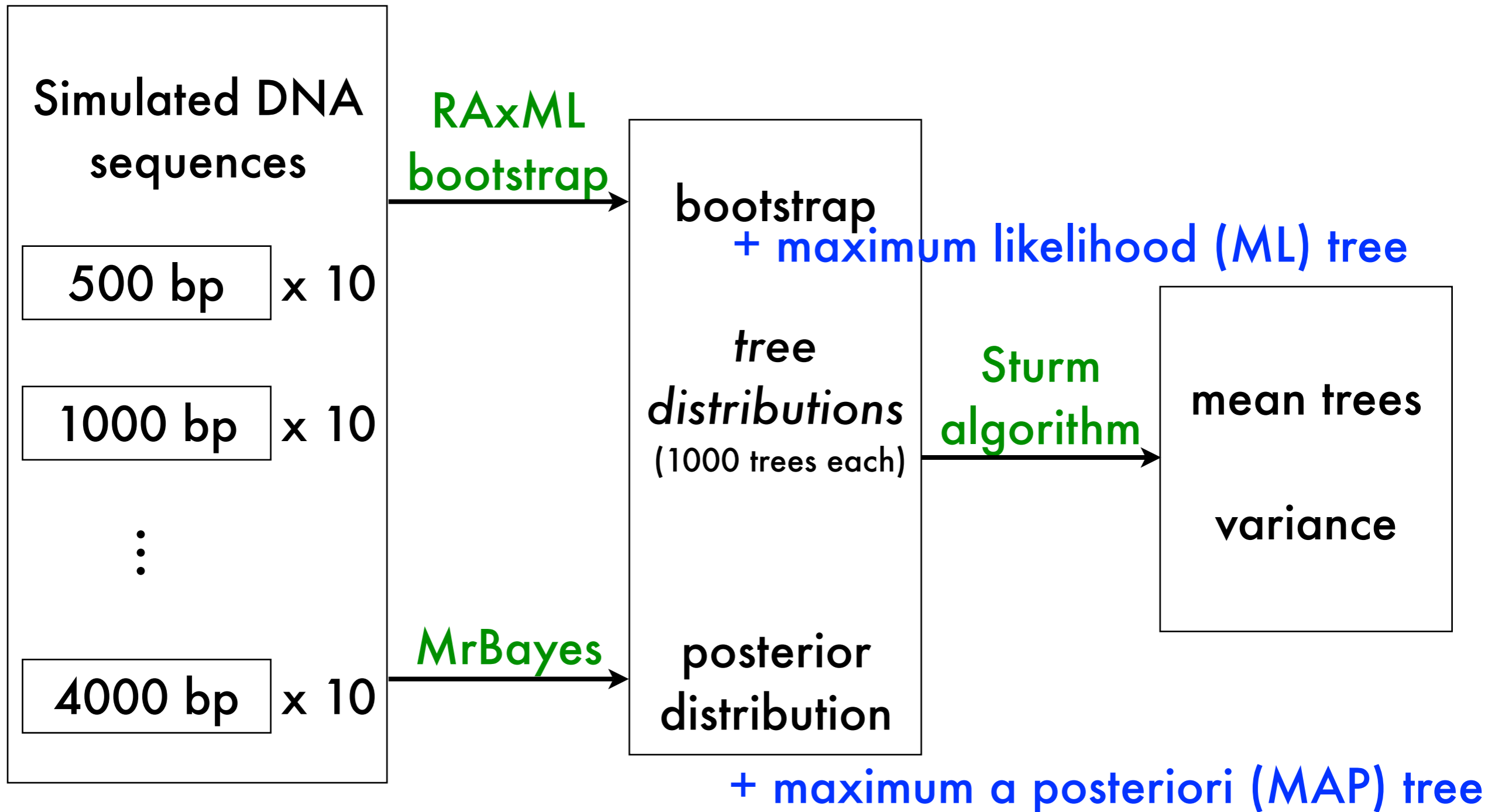
# Experimental Results



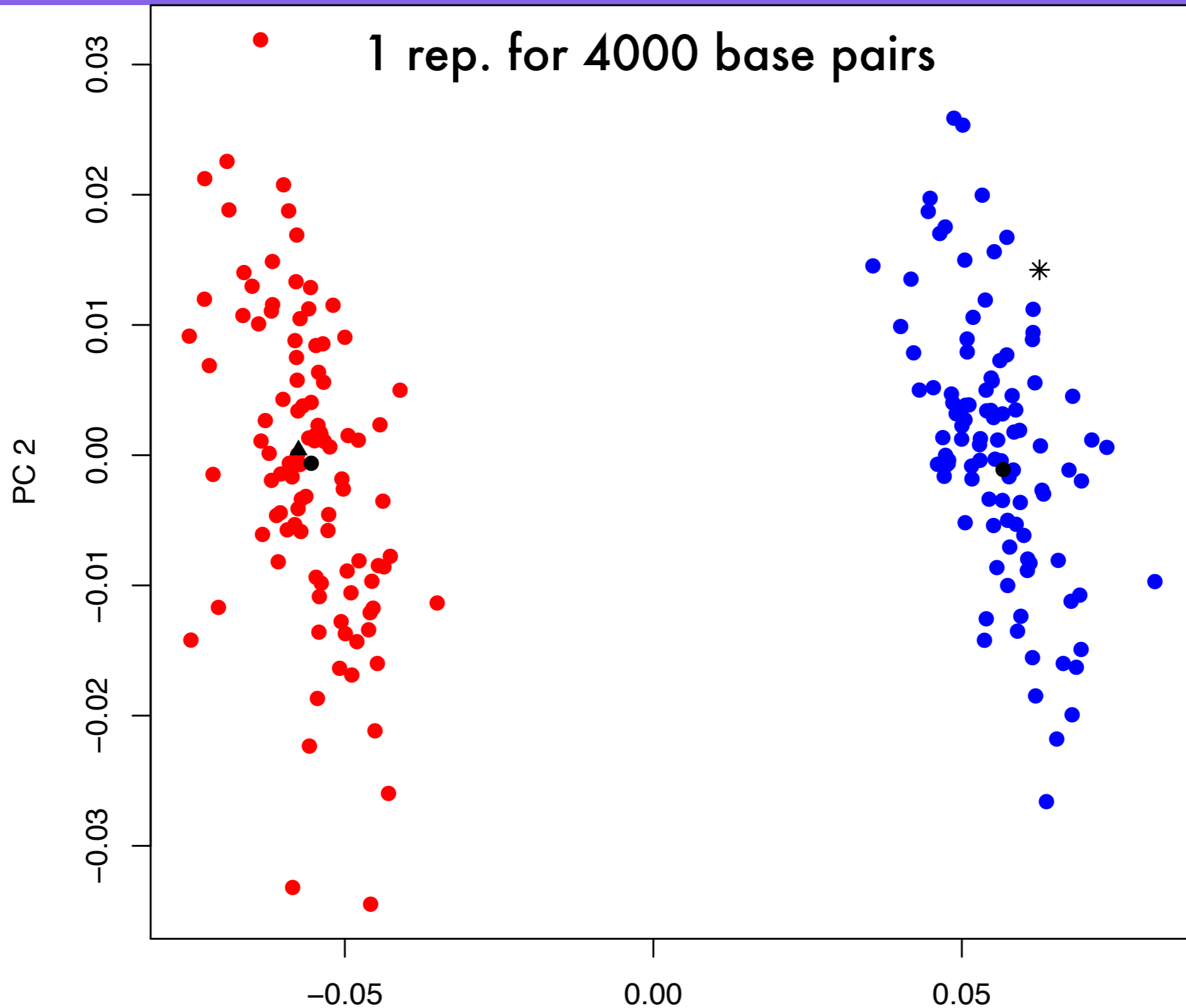
# Experimental Results



# Experimental Results



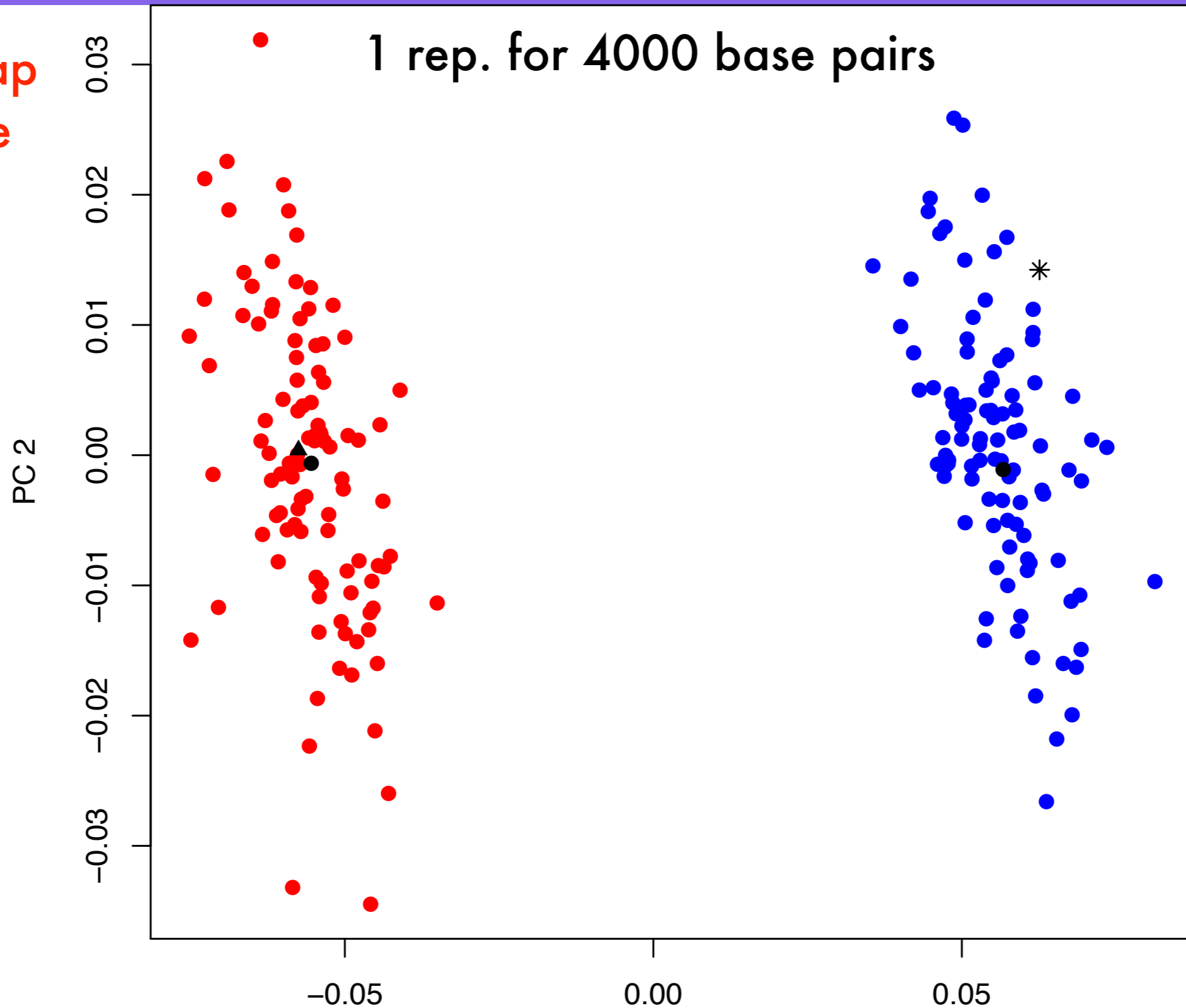
# Visualization vis MDS





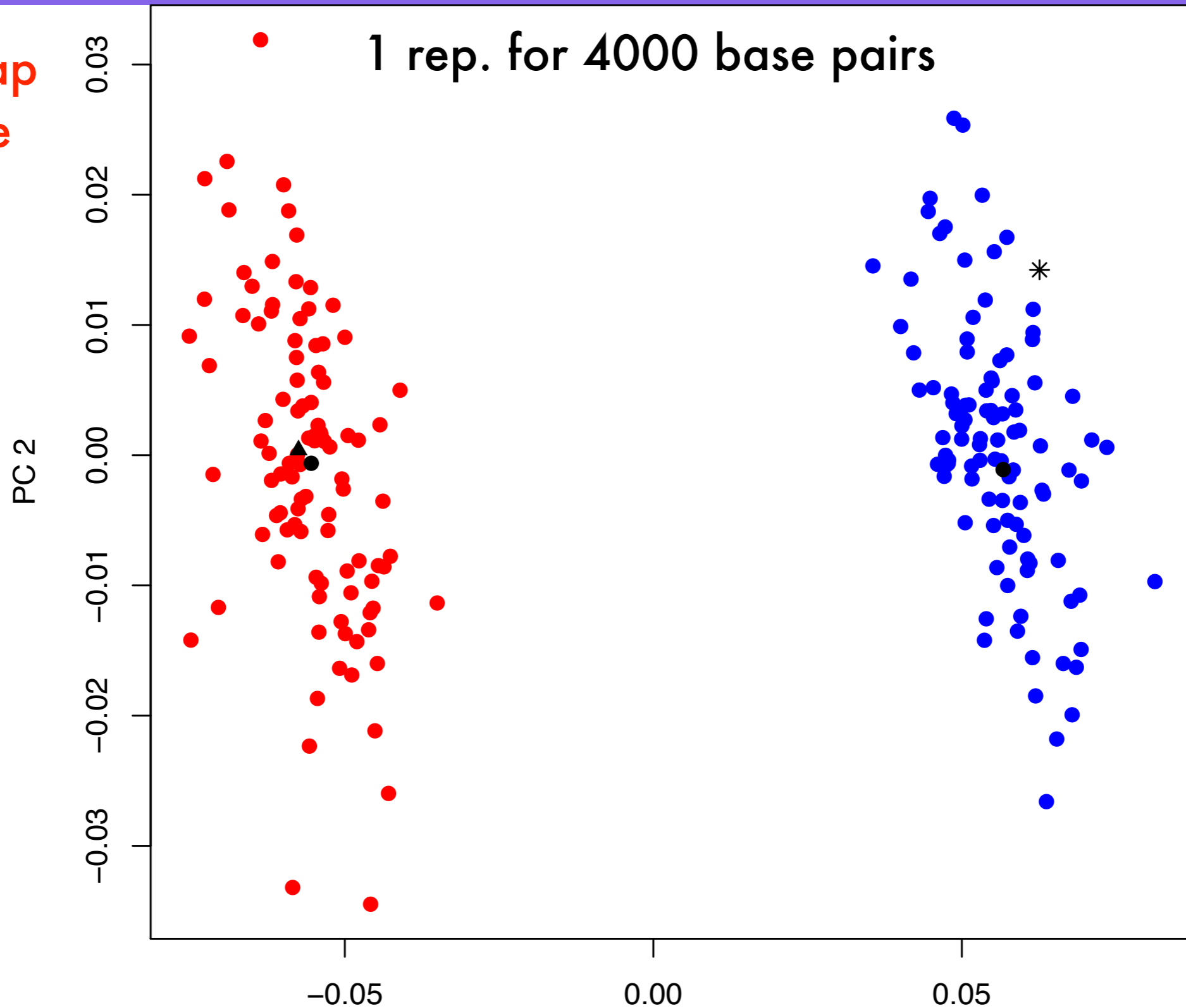
# Visualization vis MDS

bootstrap  
sample



# Visualization vis MDS

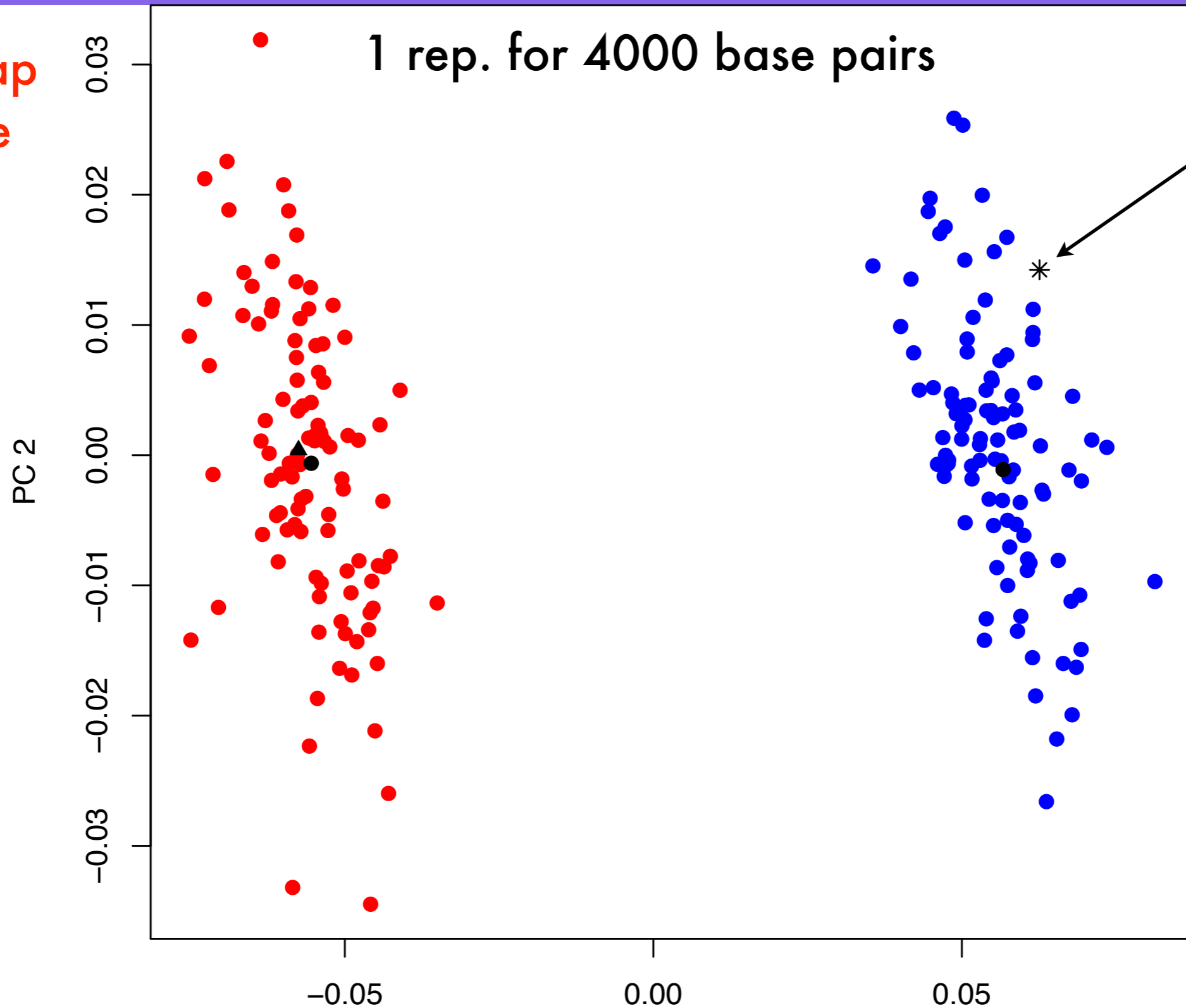
bootstrap  
sample



posterior  
sample

# Visualization vis MDS

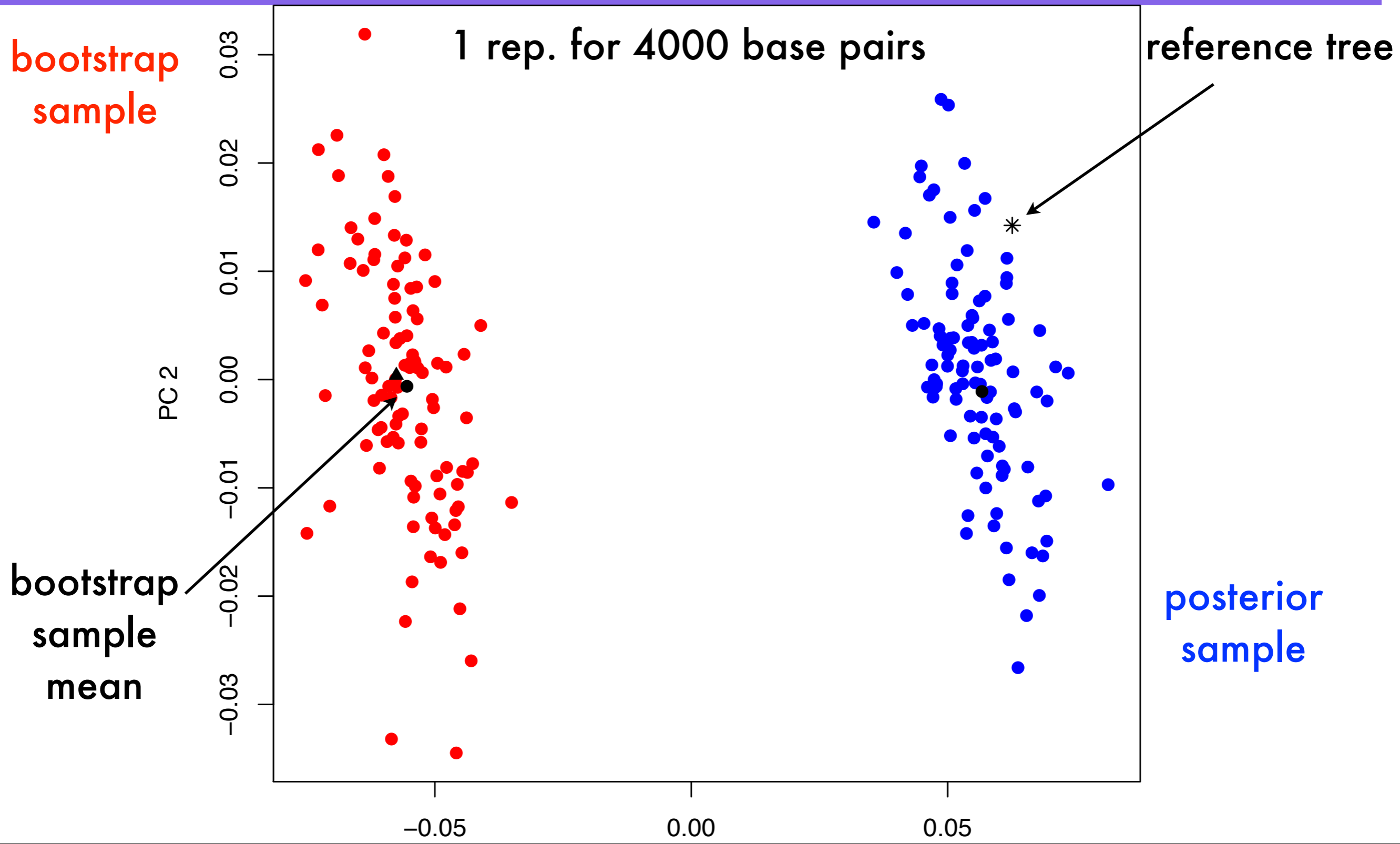
bootstrap  
sample



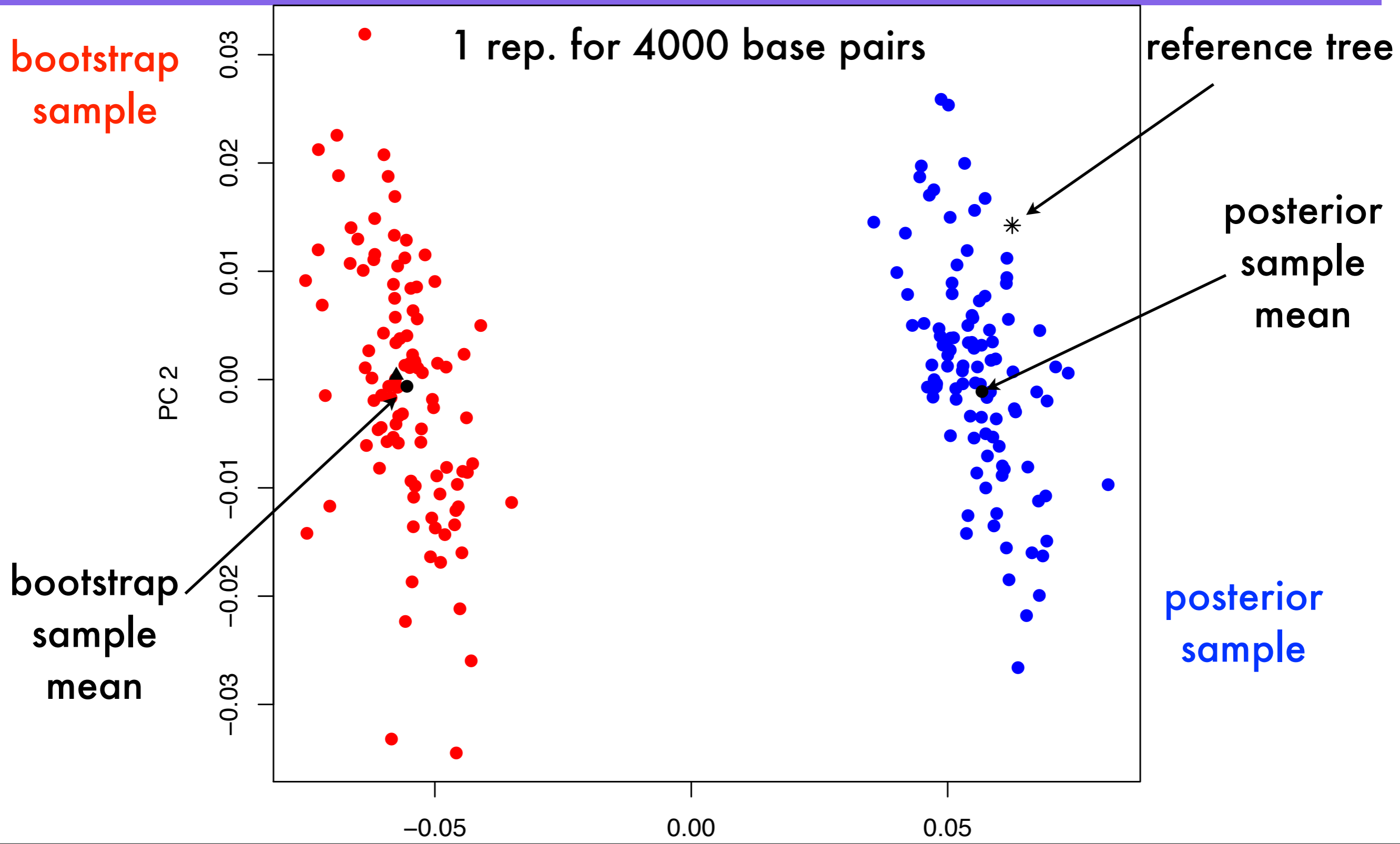
reference tree

posterior  
sample

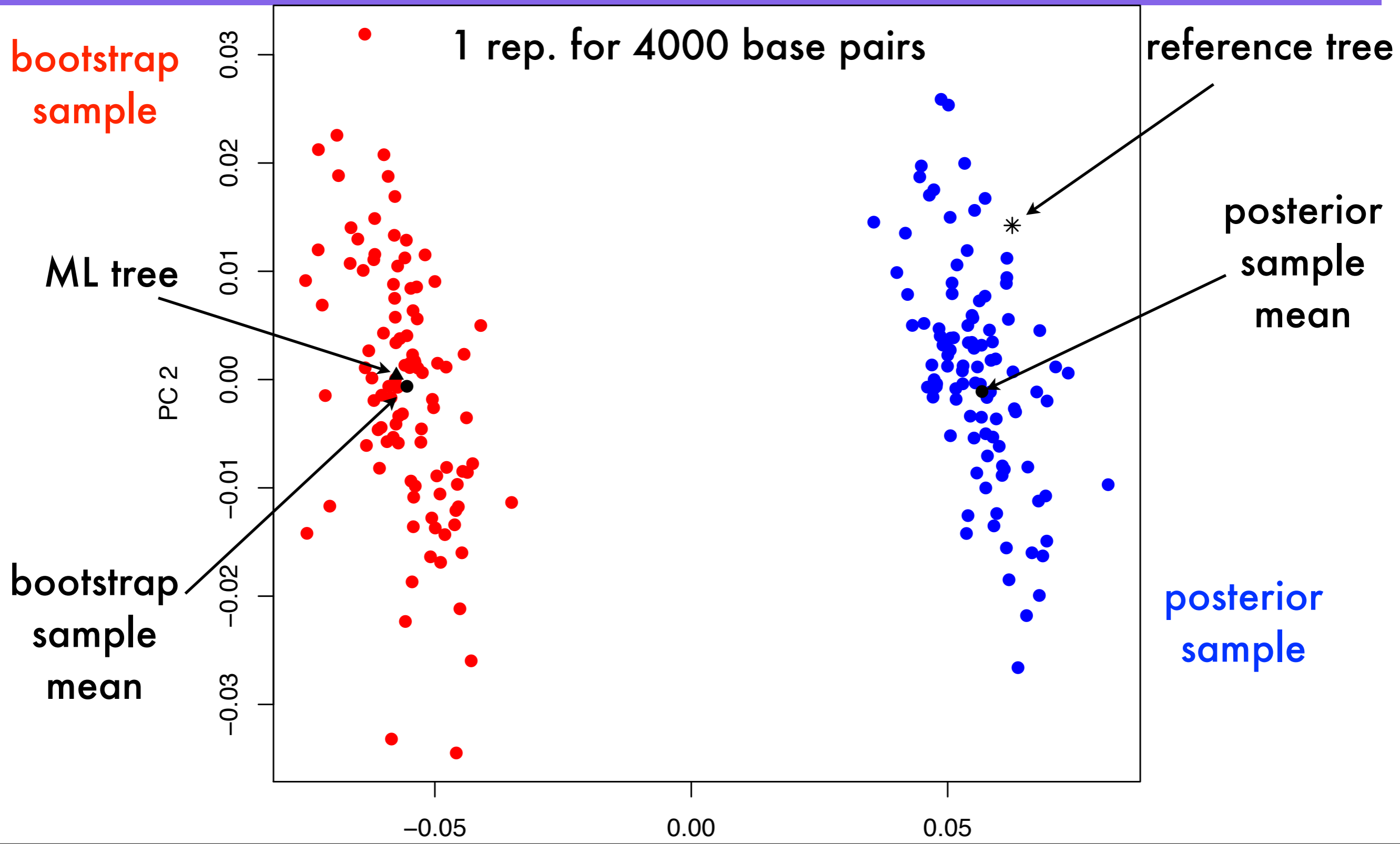
# Visualization vis MDS



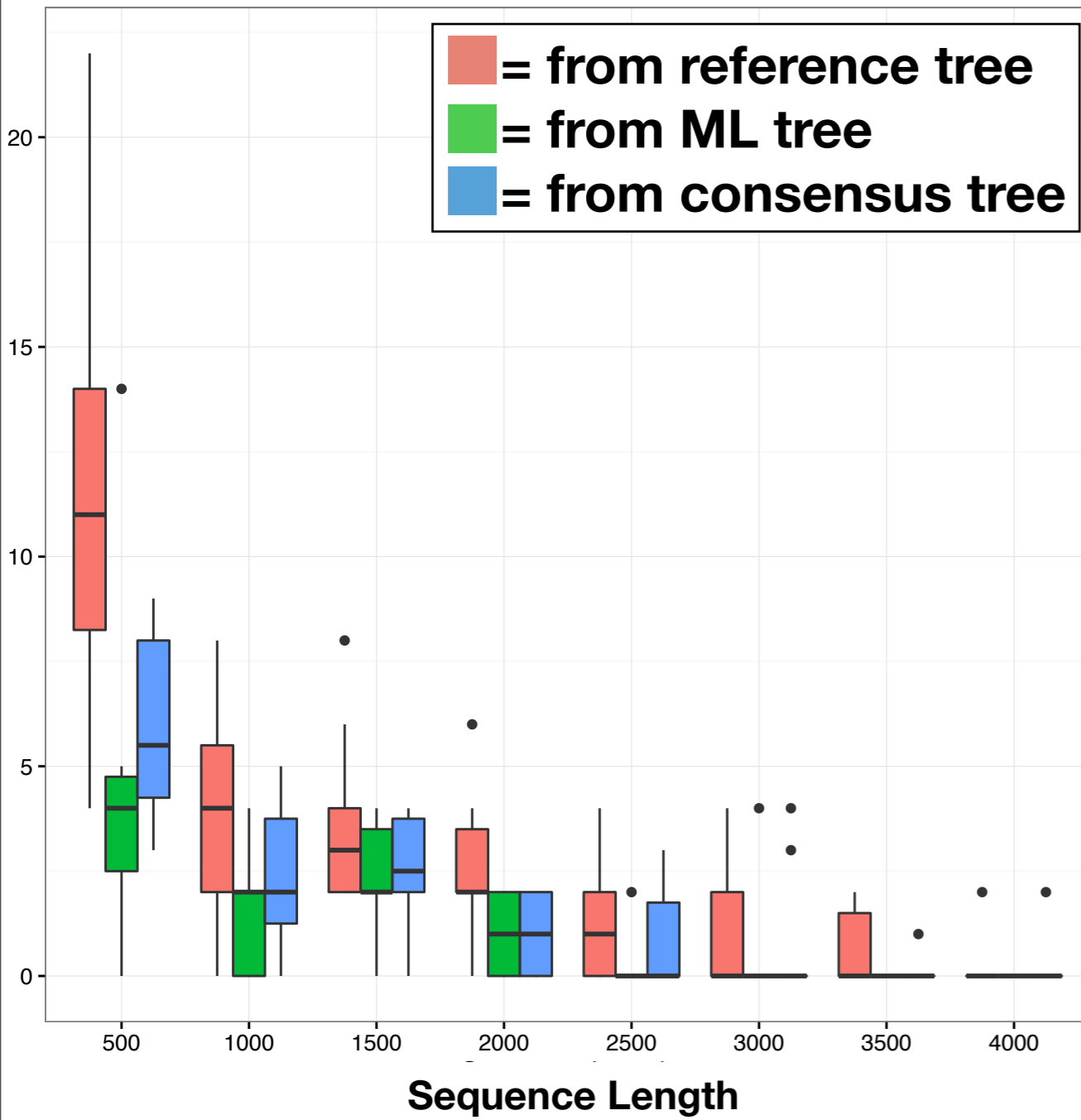
# Visualization vis MDS



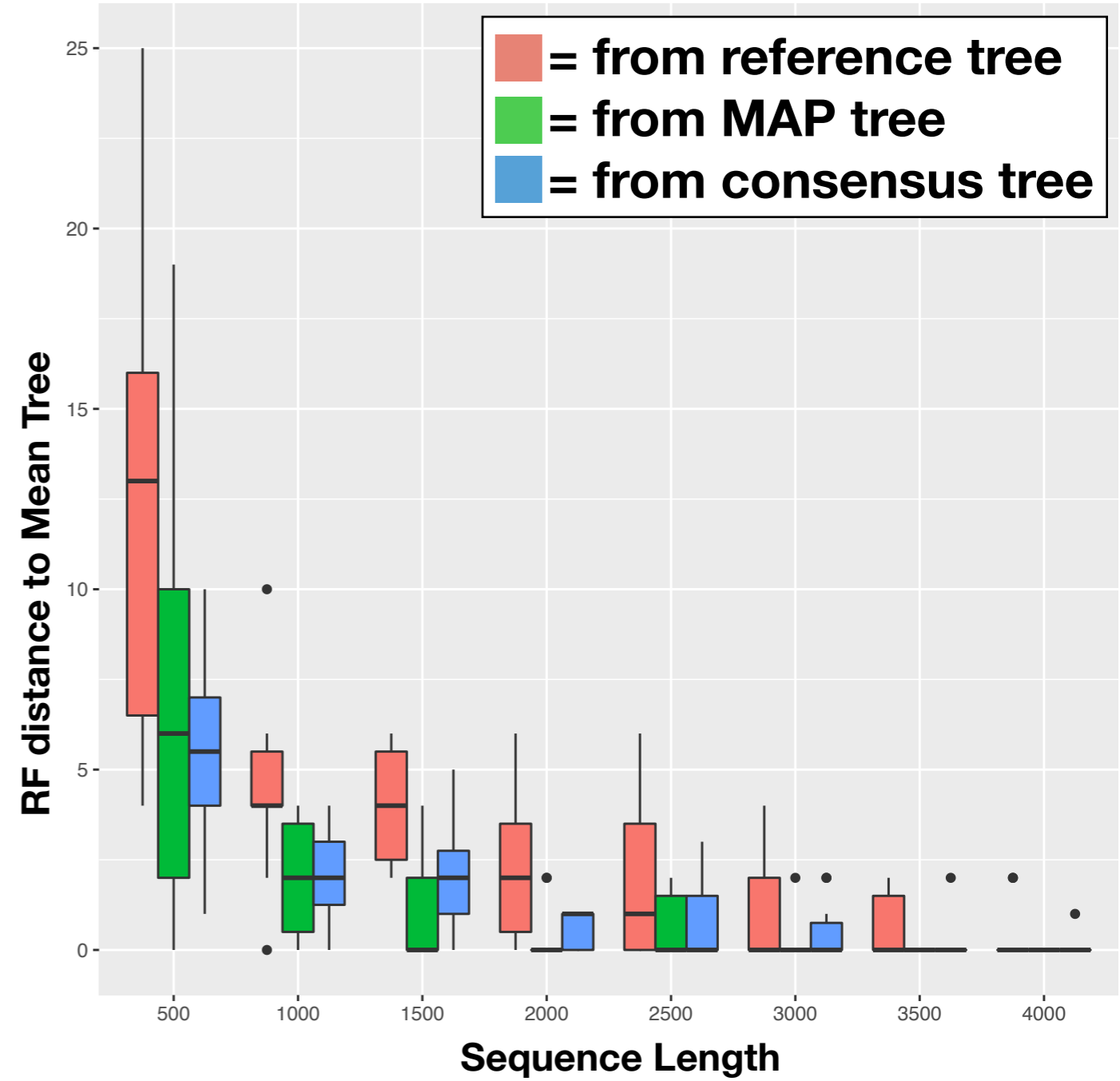
# Visualization vis MDS



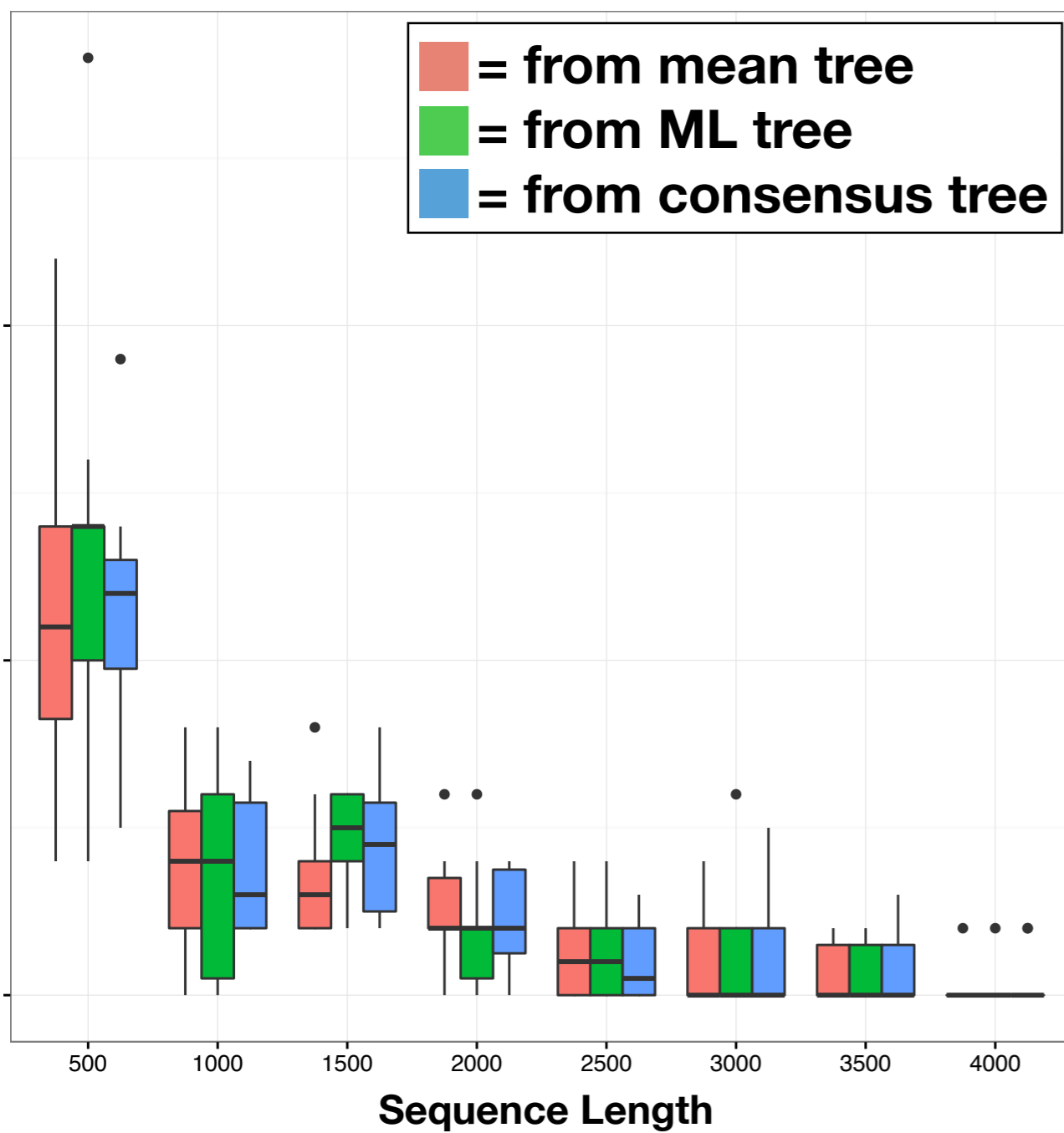
## Bootstrap Samples



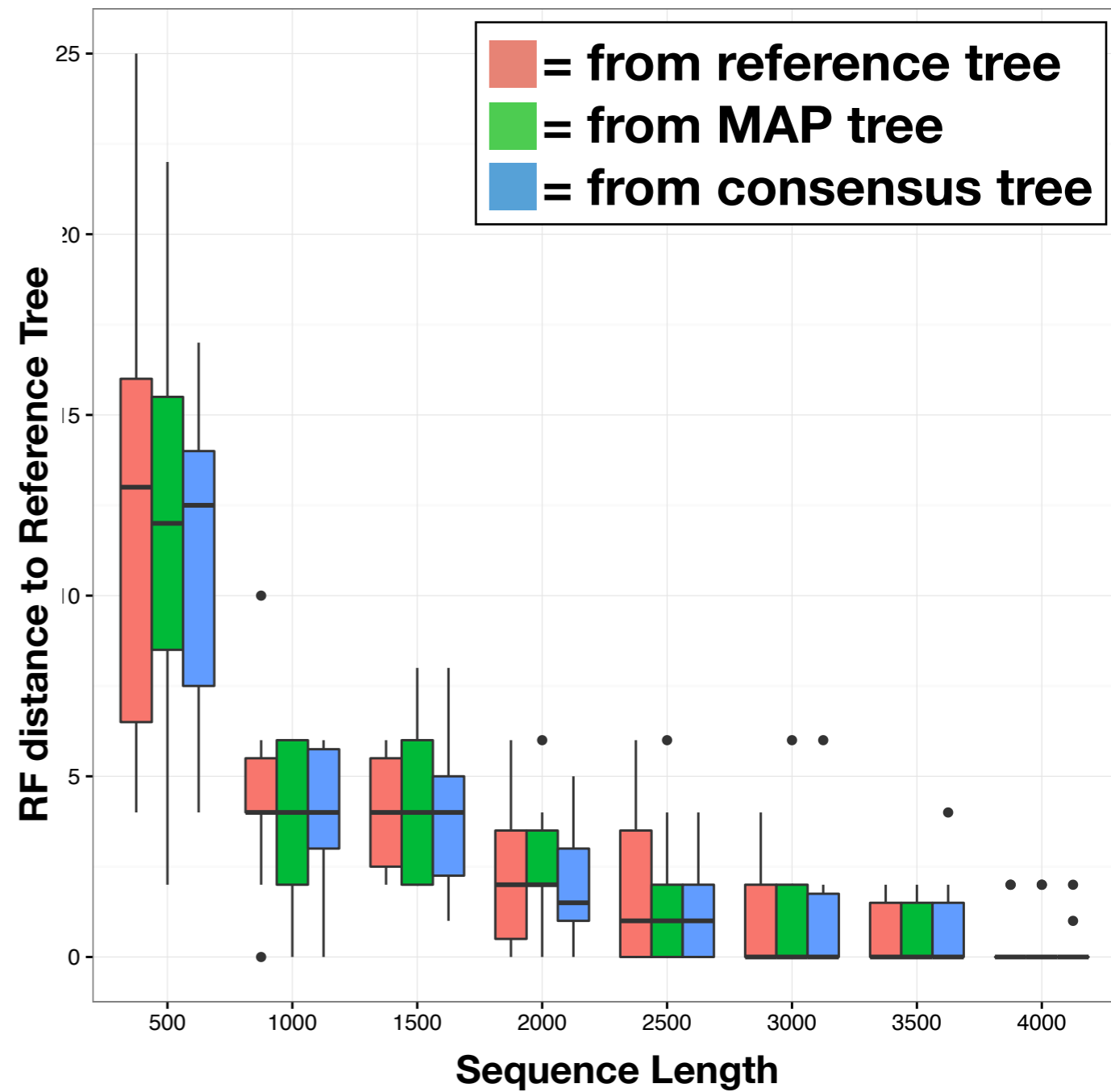
## Posterior Samples



## Bootstrap Samples



## Posterior Samples

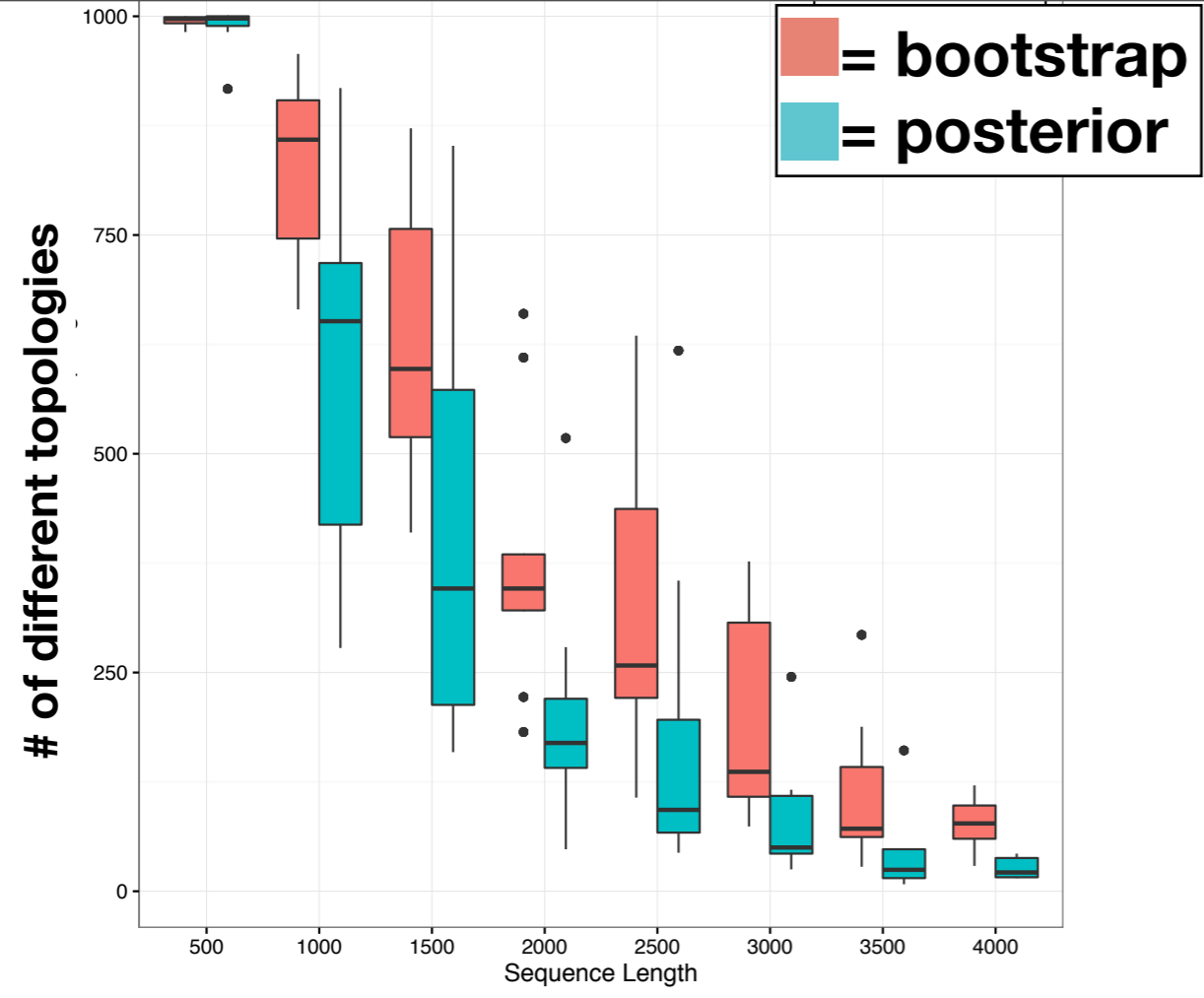


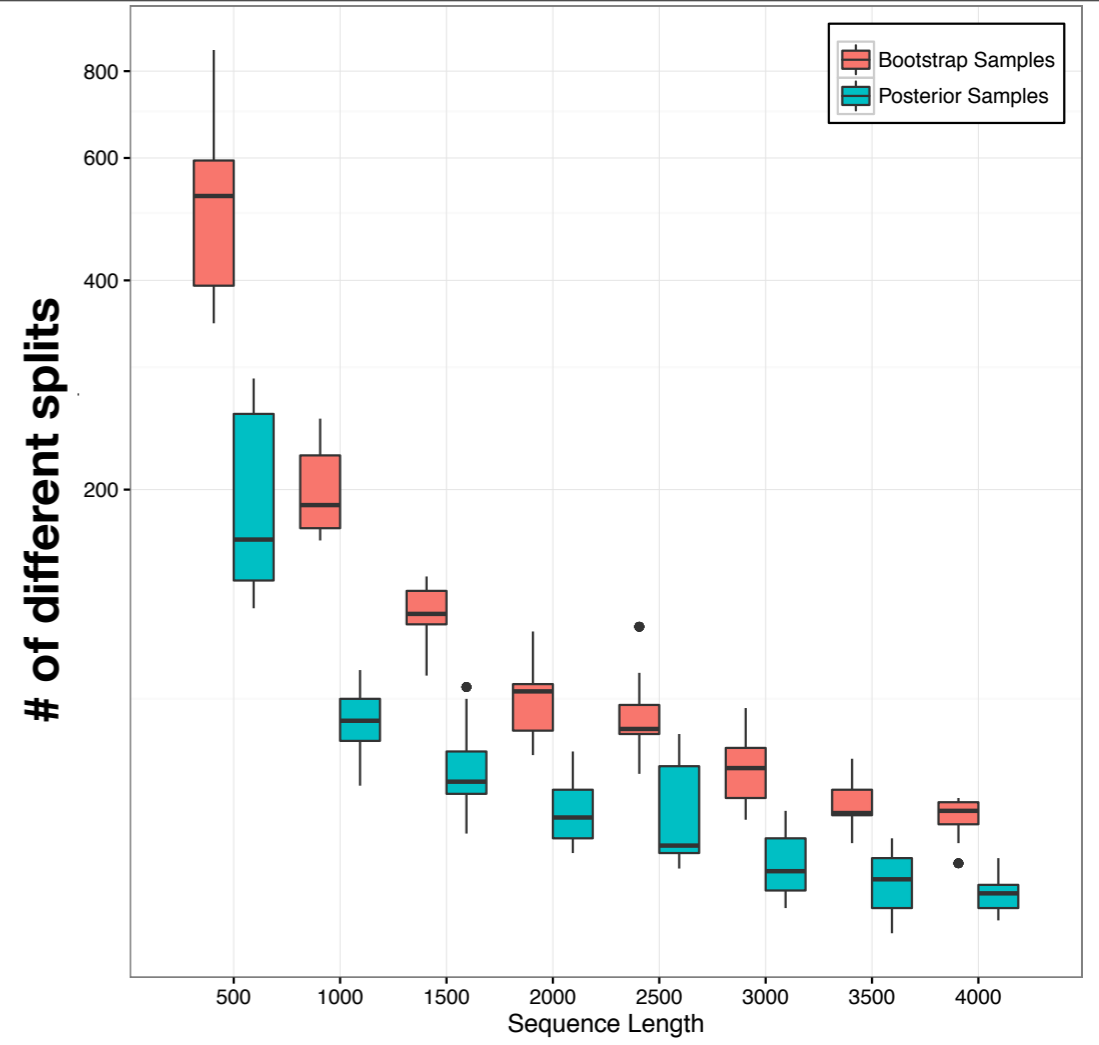
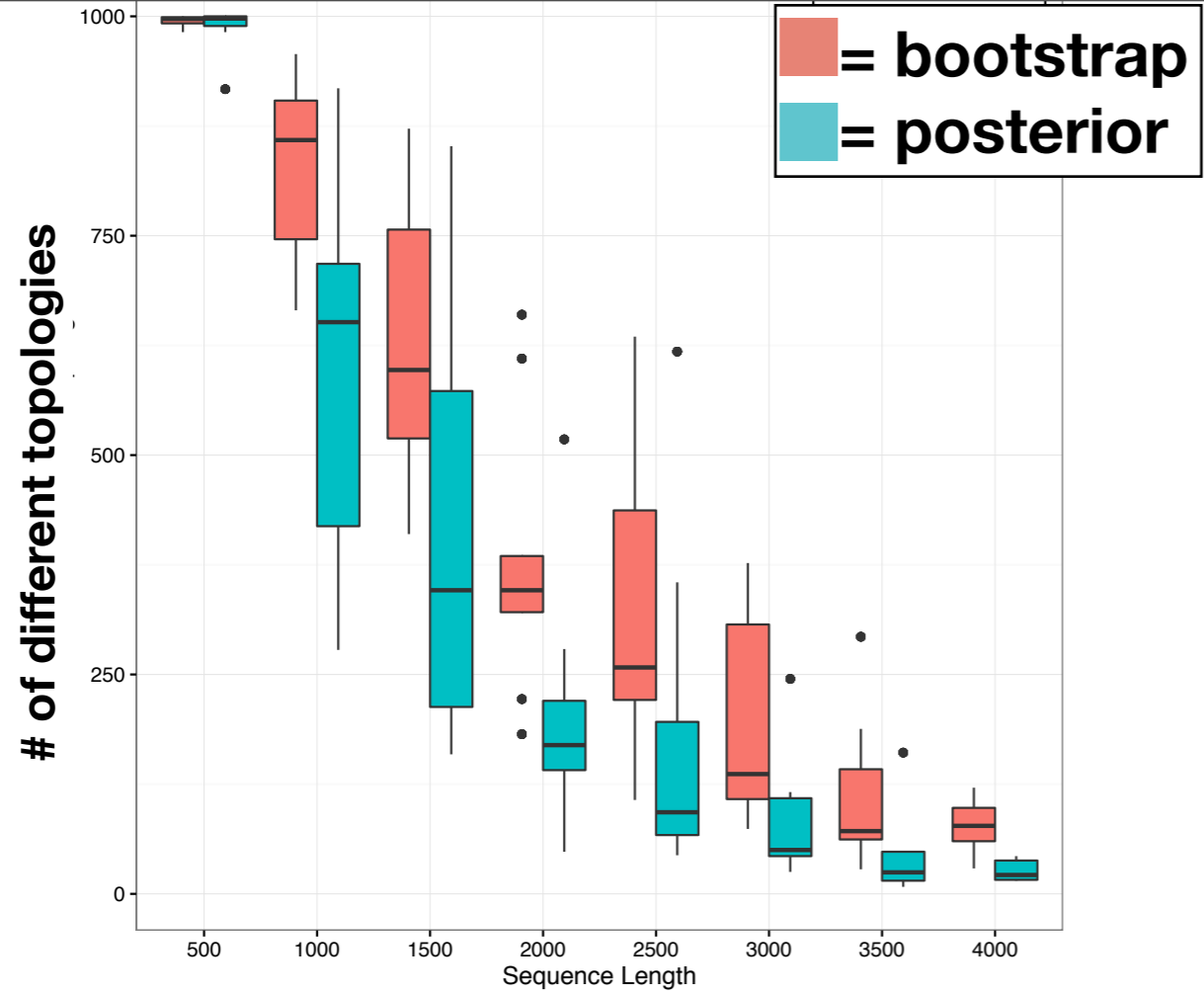


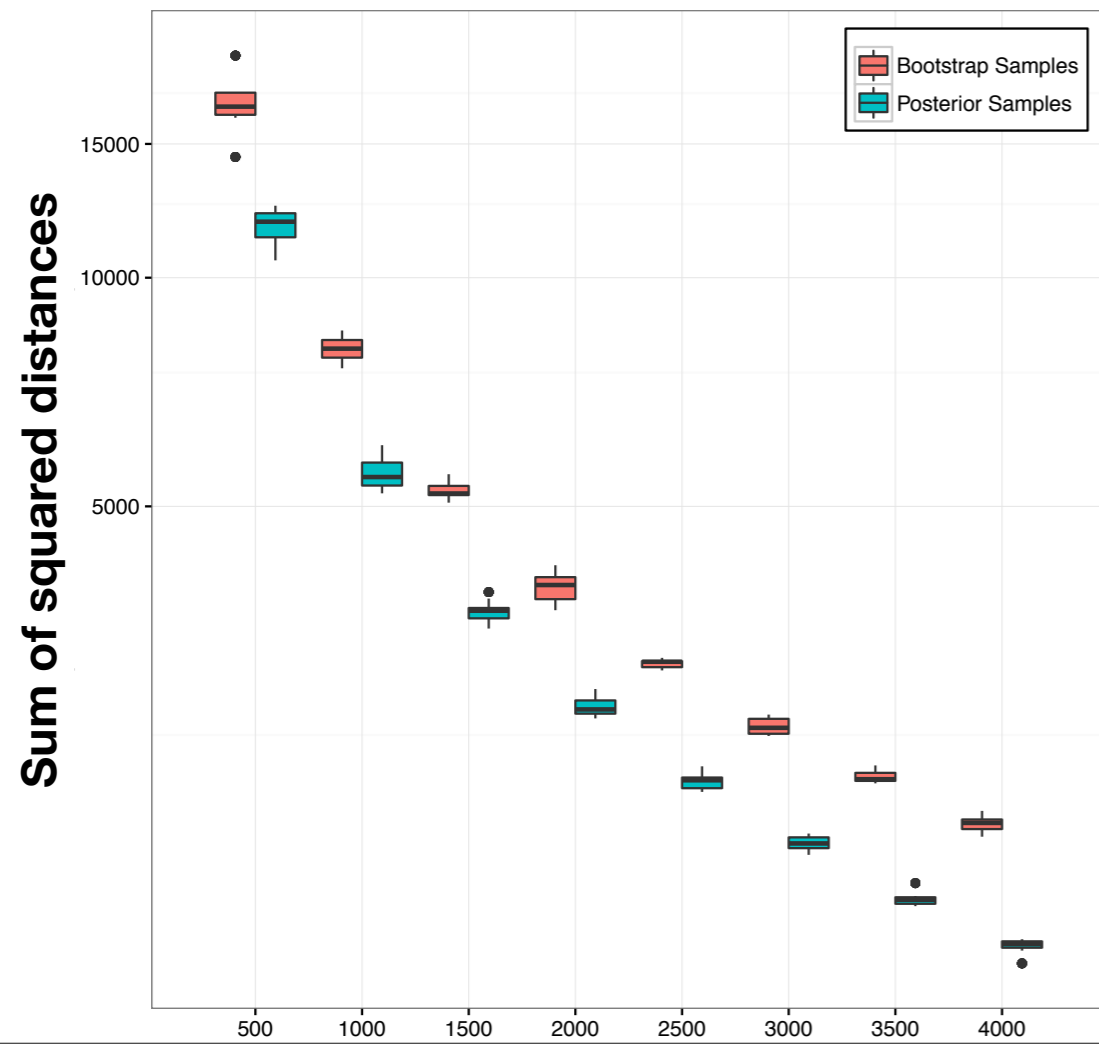
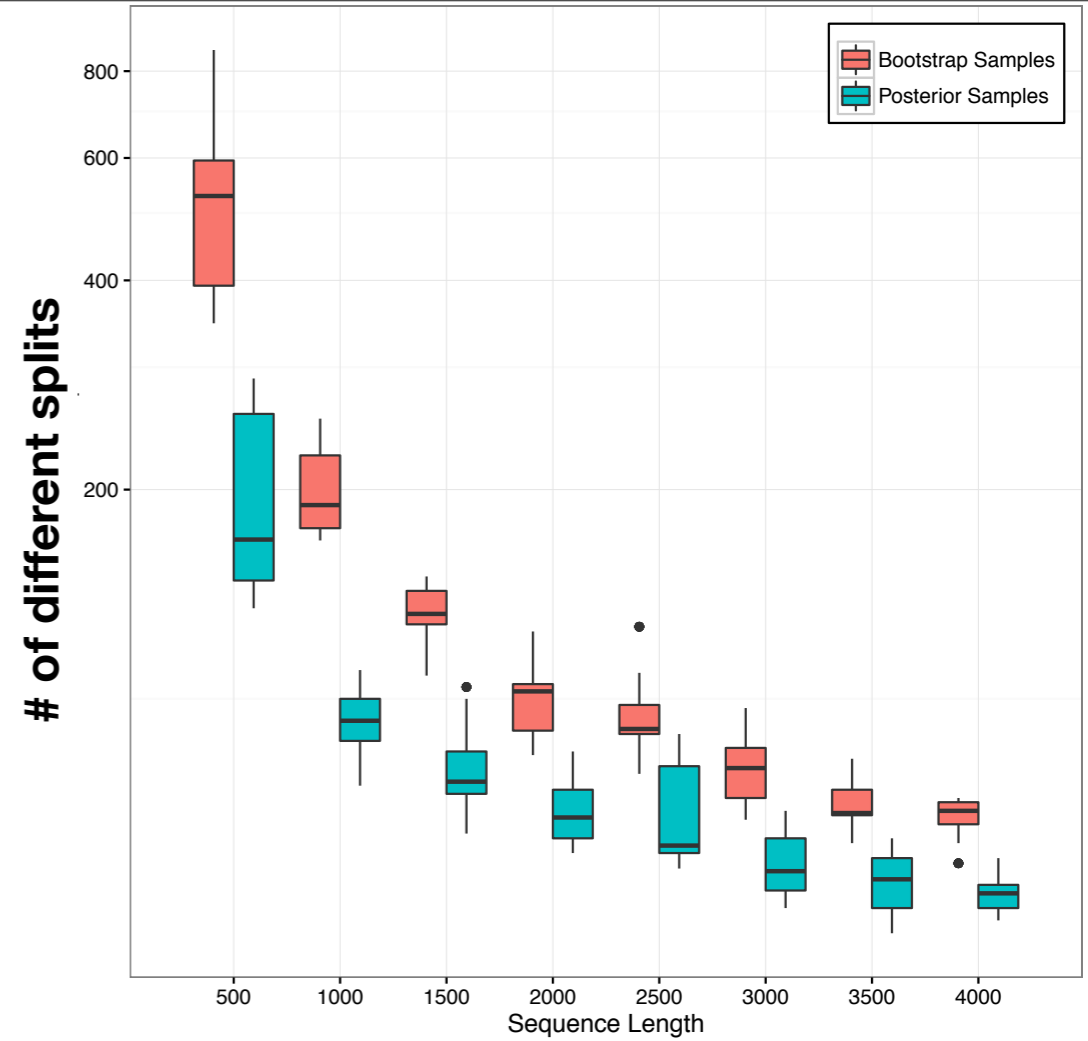
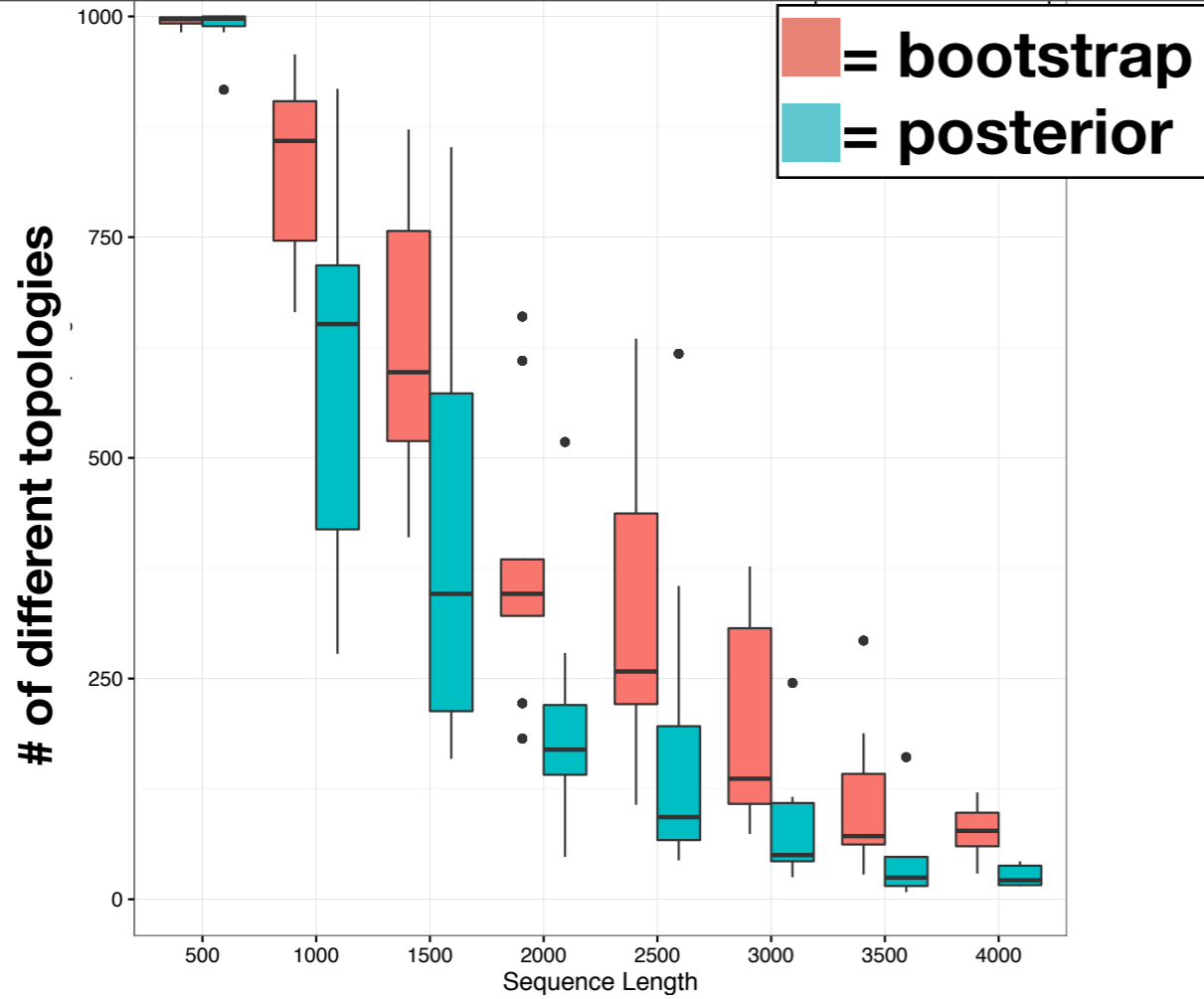
# Measures of Variance

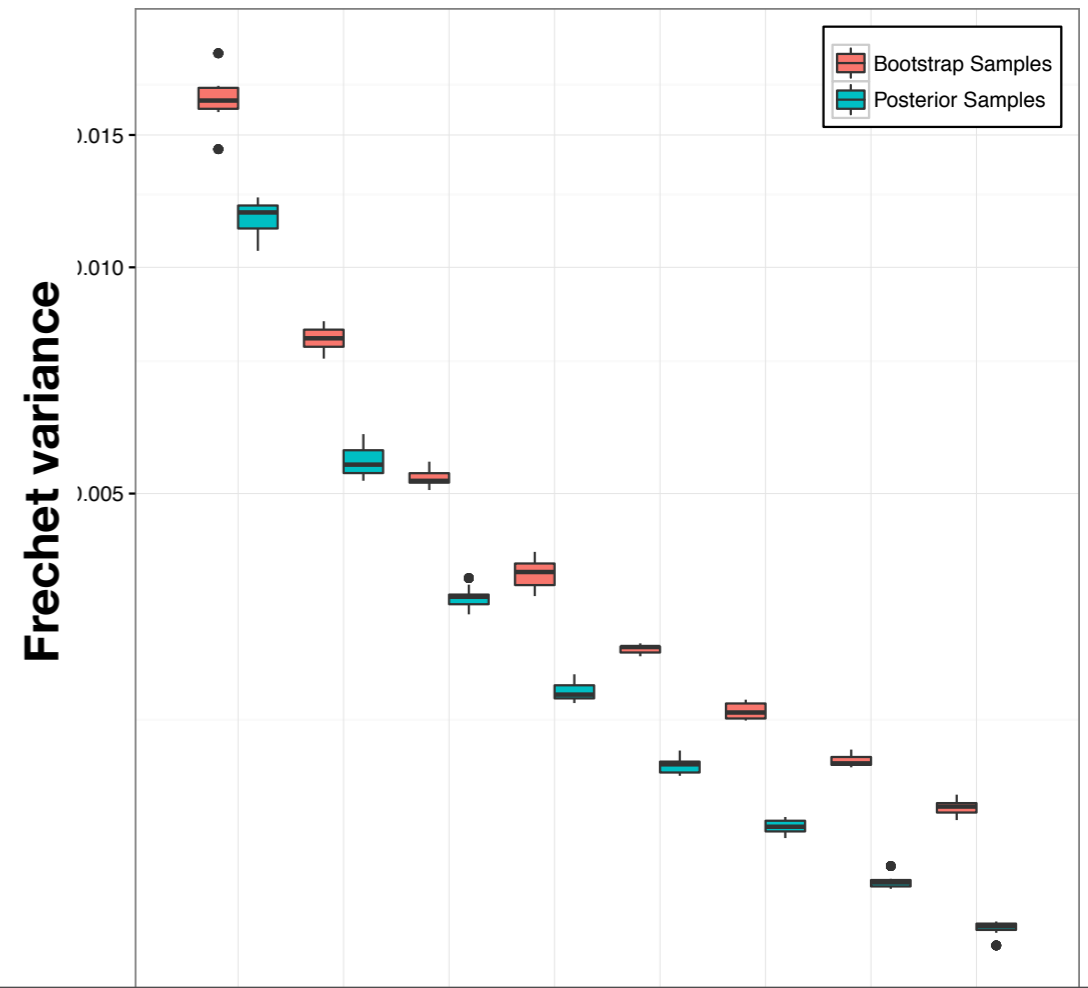
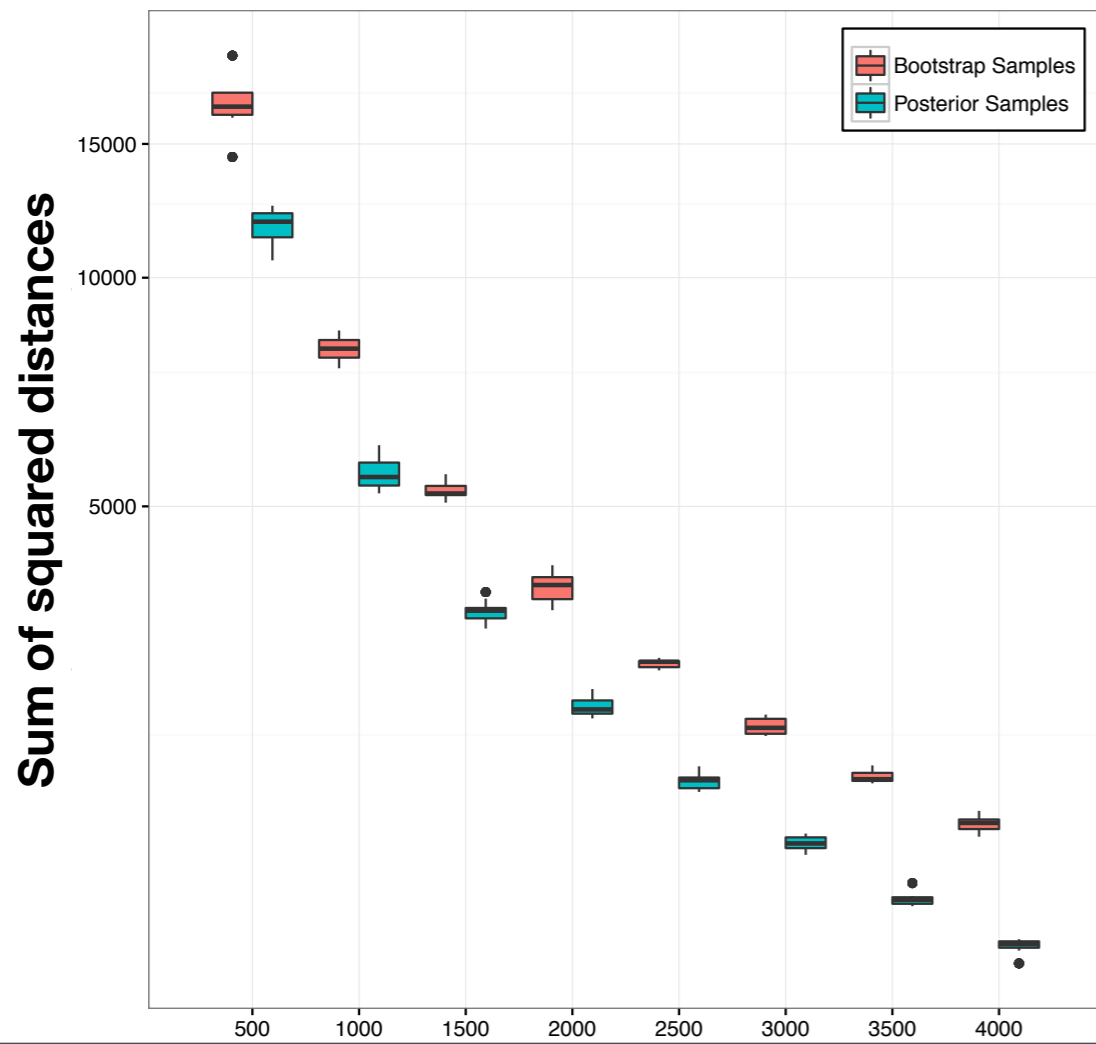
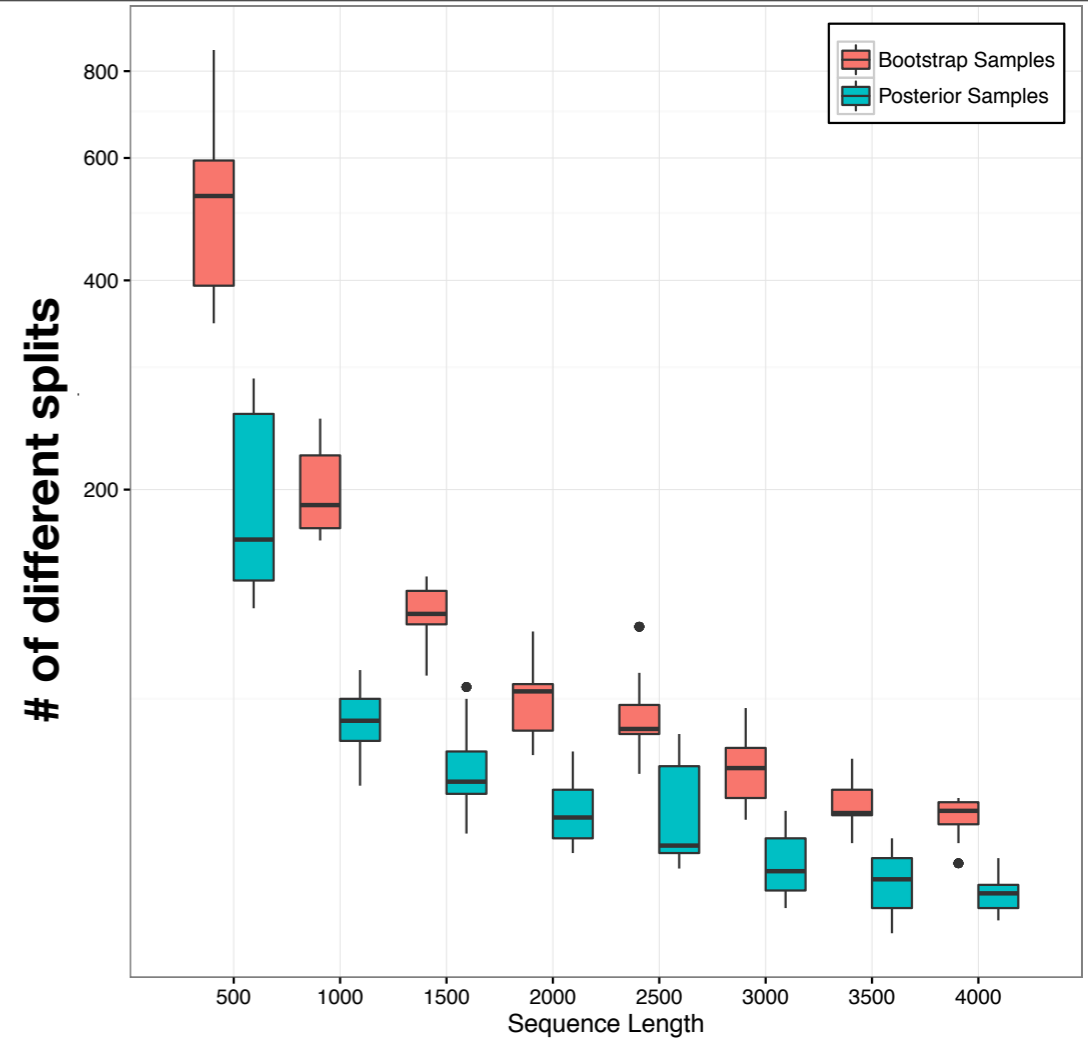
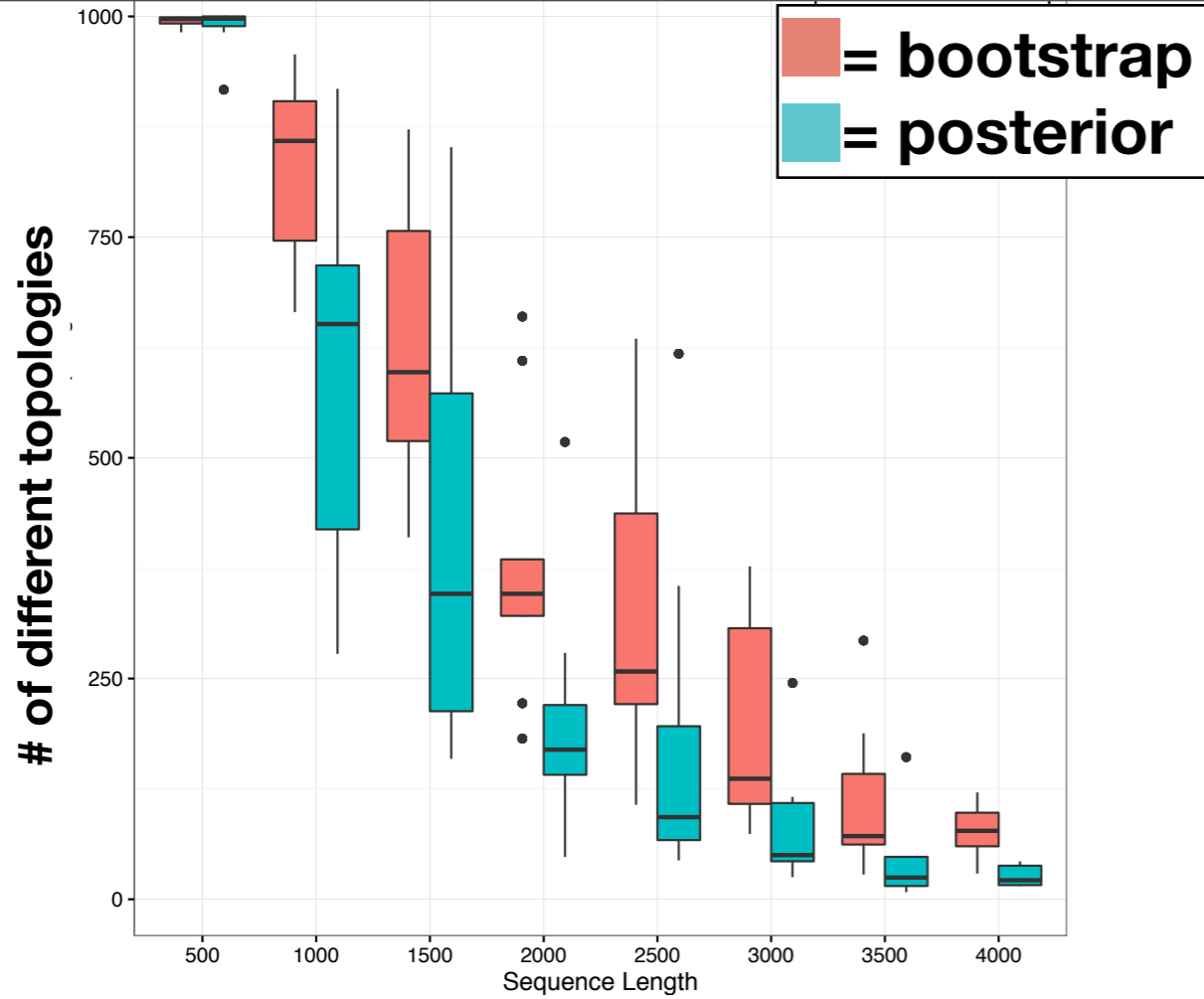
- # of different topologies in sample
- # of different splits in sample
- sum of squared distances between trees

$$\sum_{T, T' \in \mathcal{T}} d(T, T')^2$$



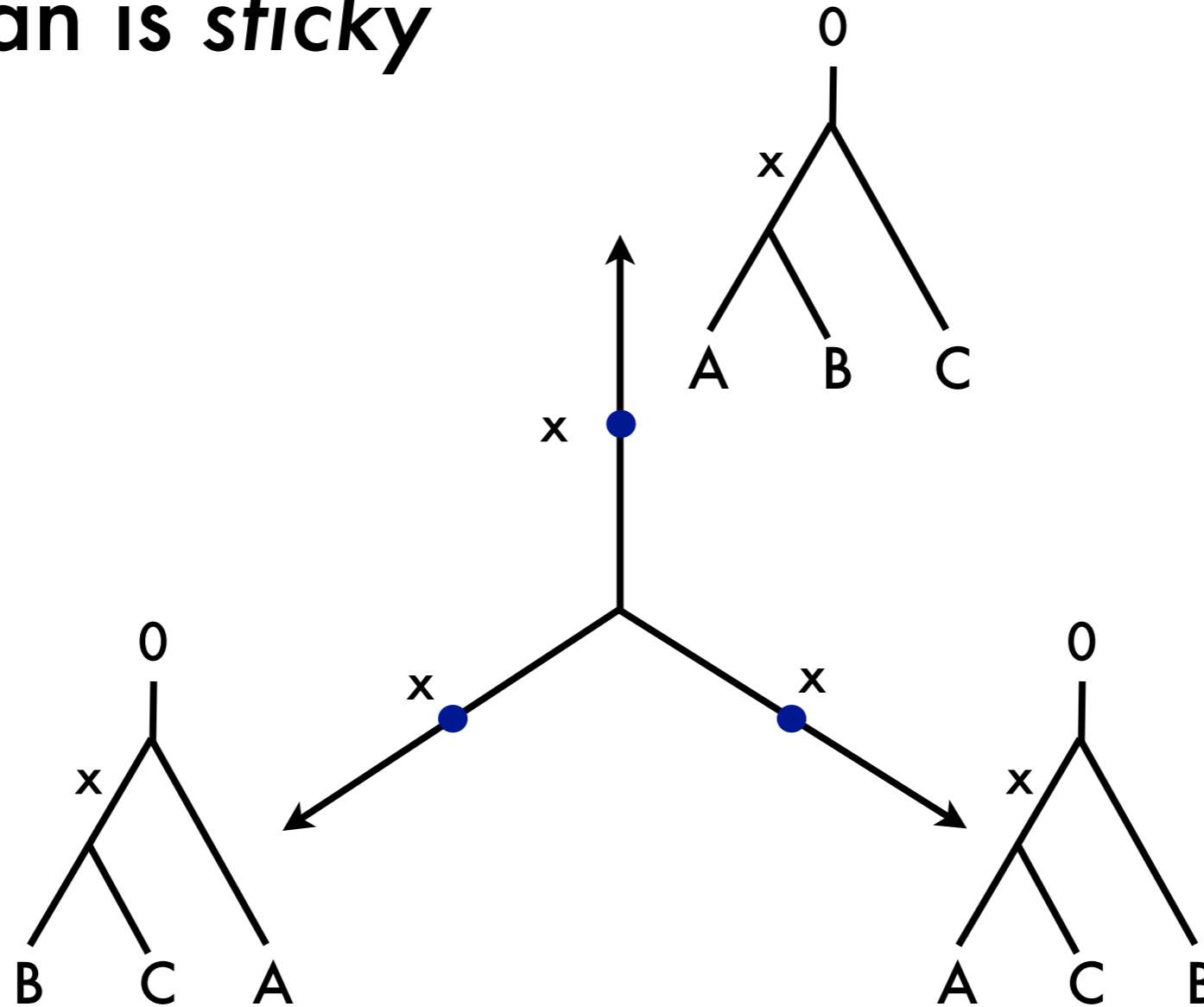






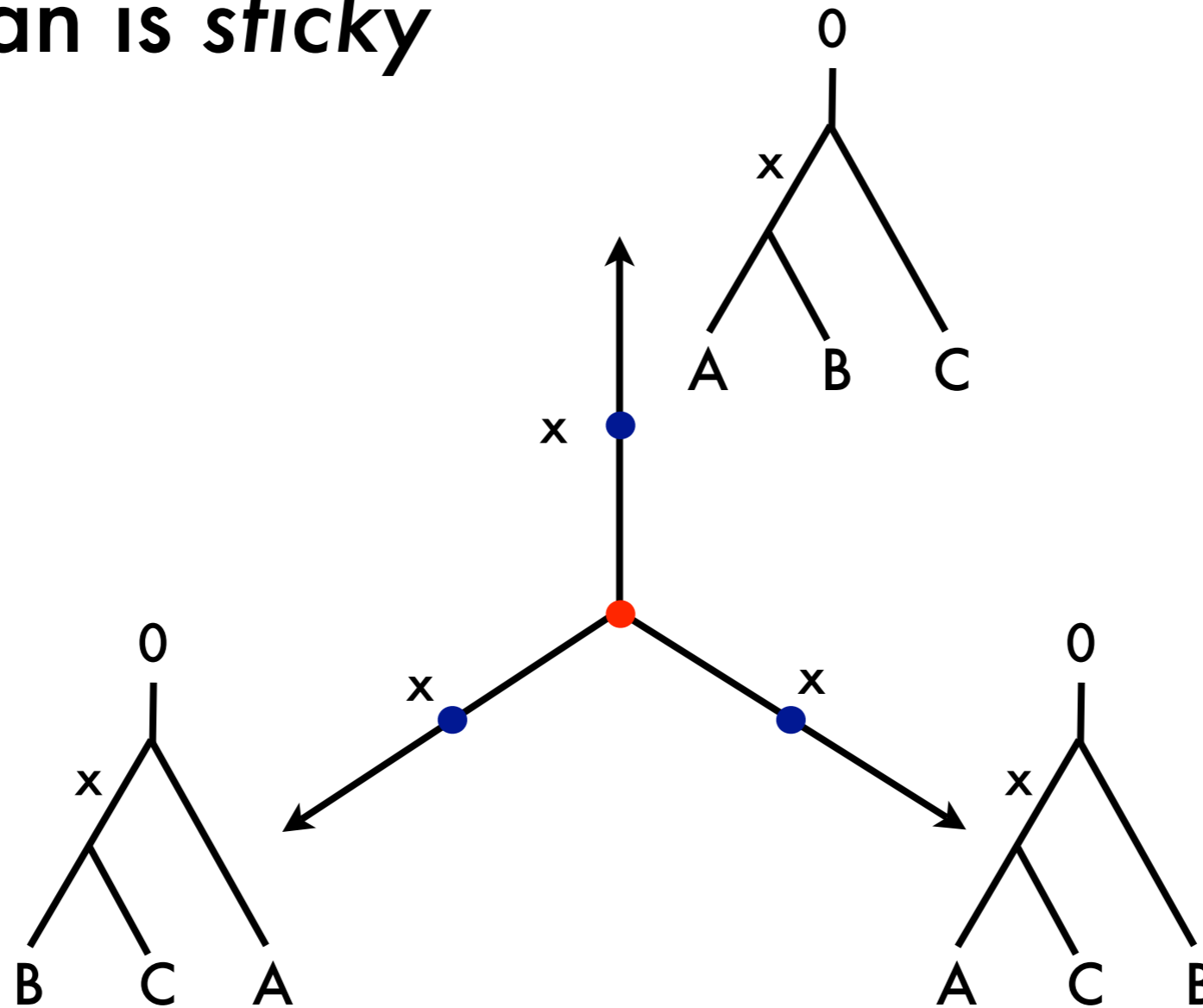
# Caveat

- Mean is *sticky*



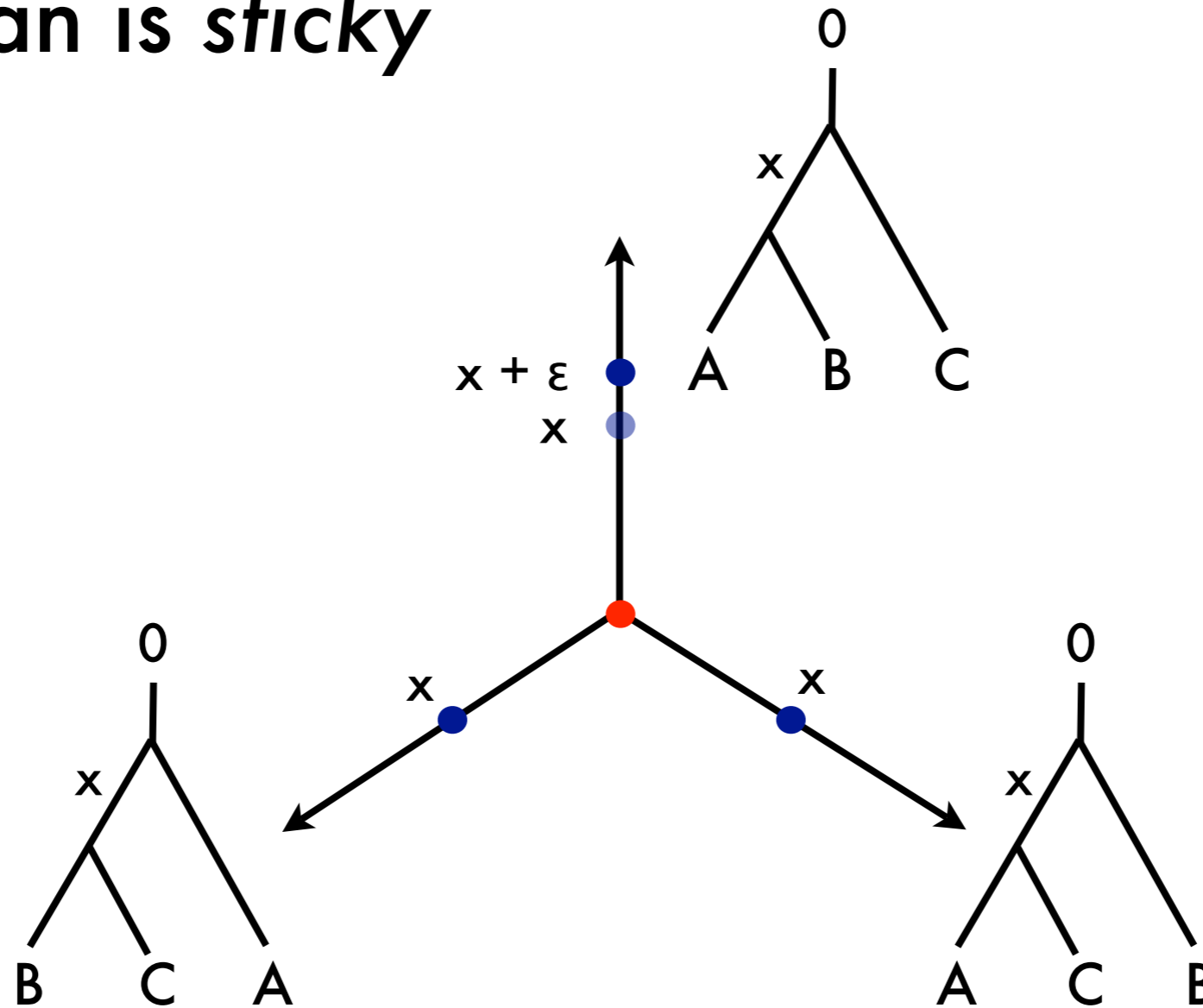
# Caveat

- Mean is *sticky*



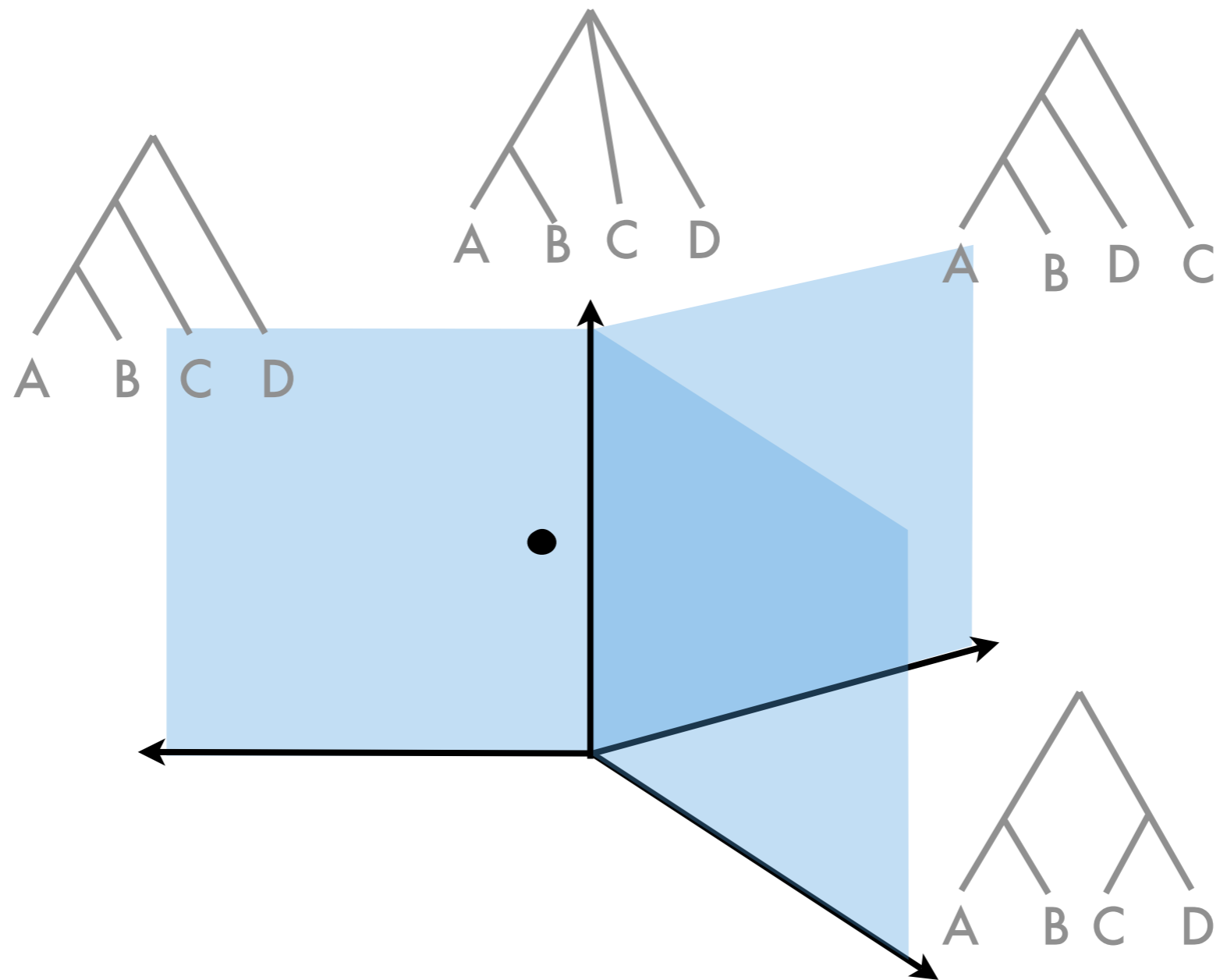
# Caveat

- Mean is *sticky*

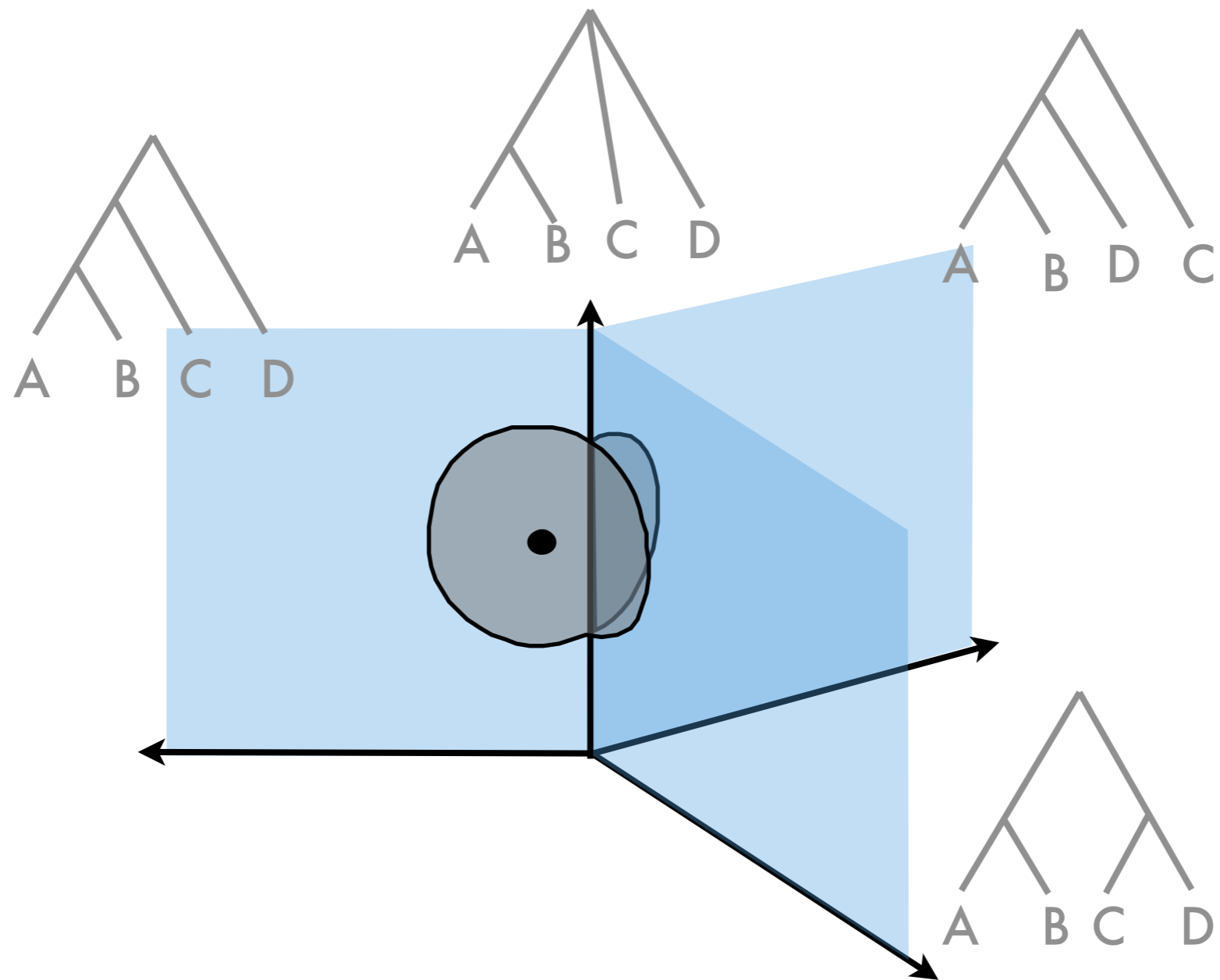




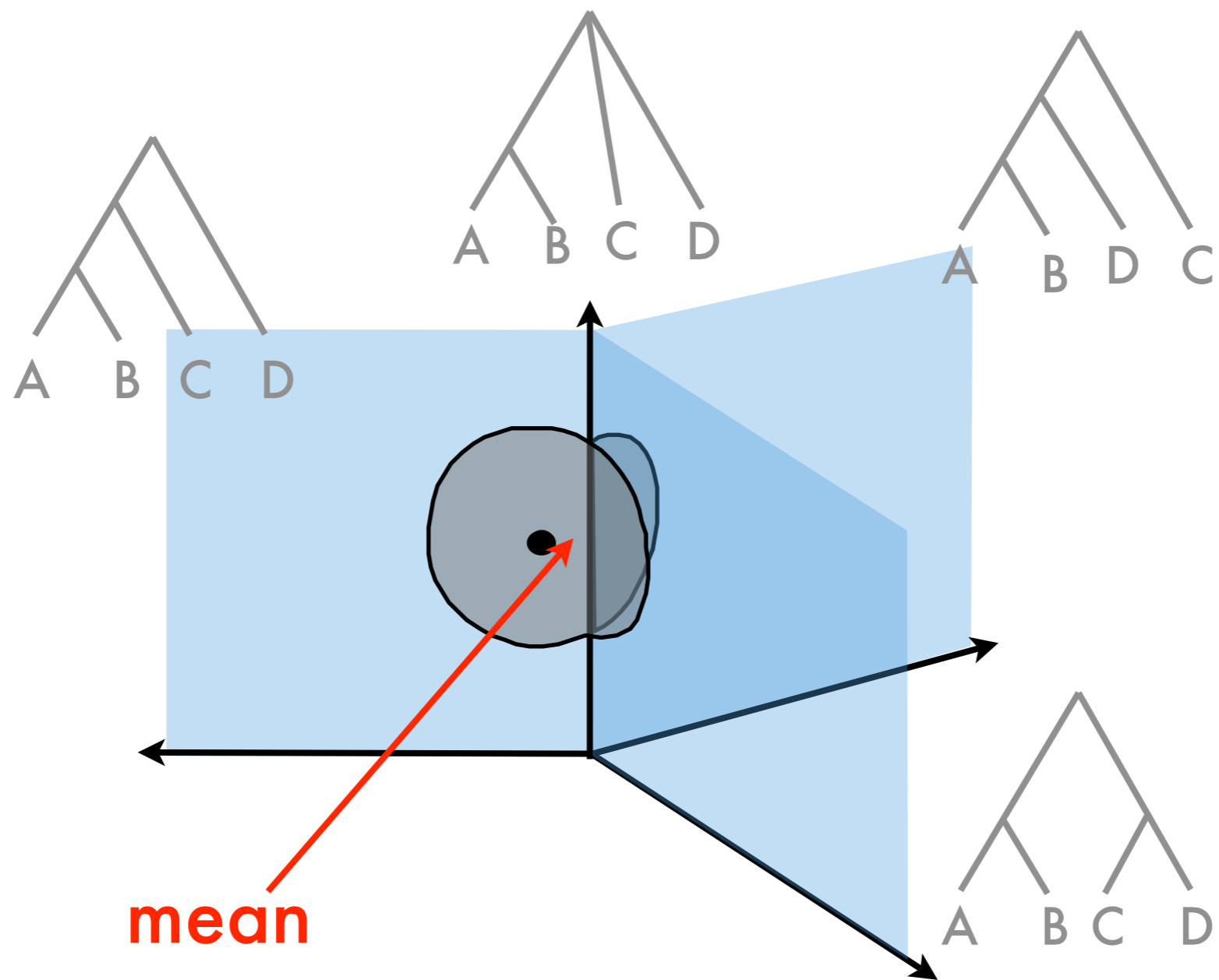
# Caveat



# Caveat



# Caveat



# Other Statistics

- Central Limit Theorem on BHV tree space:
  - special cases: Hotz, O., et al. 2012; Barden, Le, O., 2013, 2014; Huckemann et al. 2015
- Principal Components Analysis (PCA): (Nye 2011, 2014; Feragen, O. et al. 2013; Nye et al. 2016)
- confidence regions: Willis 2016
- multiple techniques: Chakerian and Holmes 2012, Zairis et al. 2016
- and more...

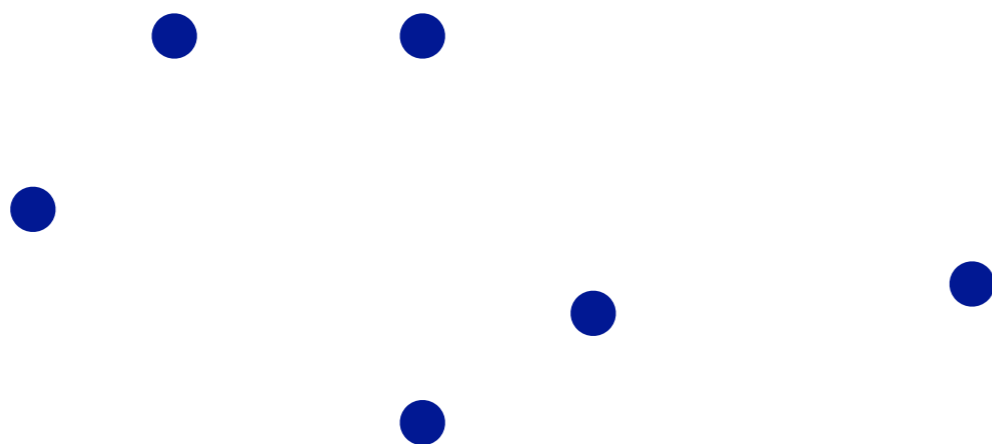
# Thank You

- **funding:** SIMONS FOUNDATION
- **webpage:** <http://comet.lehman.cuny.edu/owen>

# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )

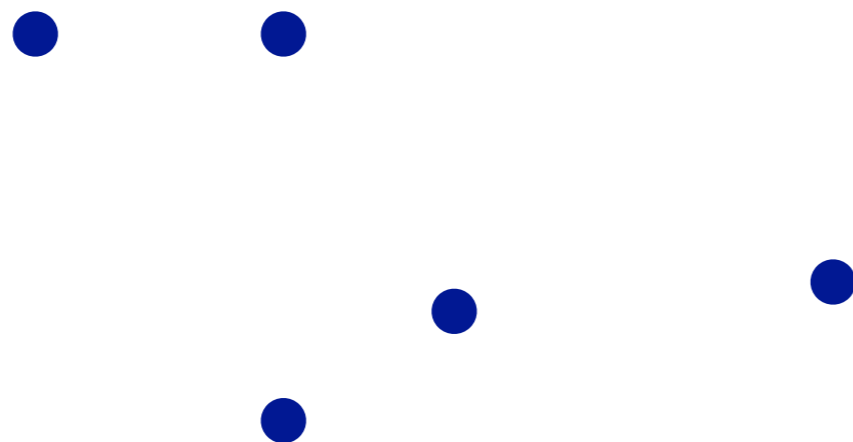


# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )

$m_1 = T_1$  ●



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )

$m_1 = T_1$  ●

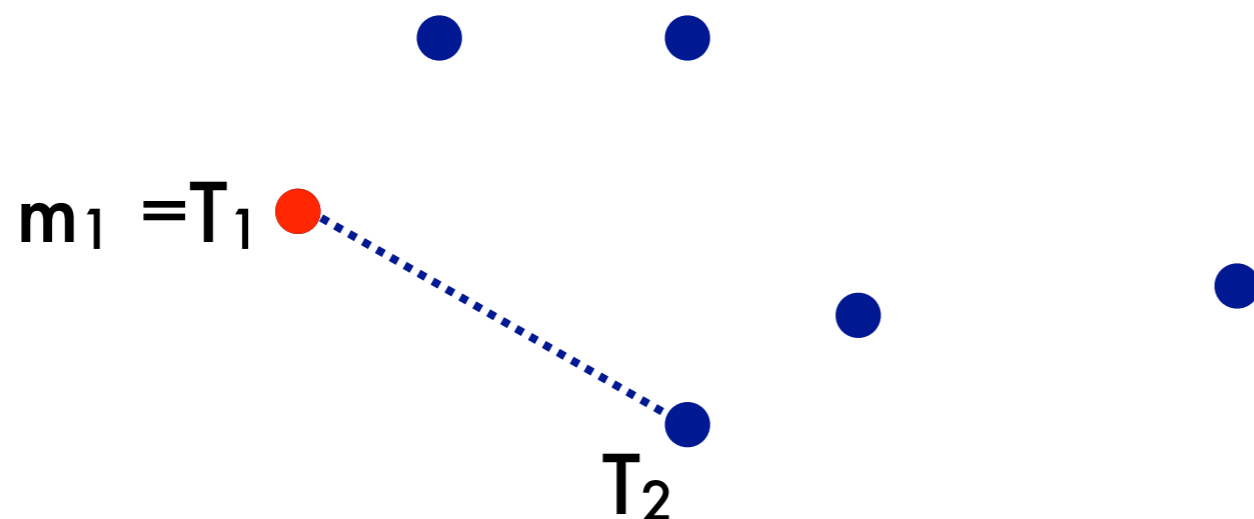
$T_2$  ●



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

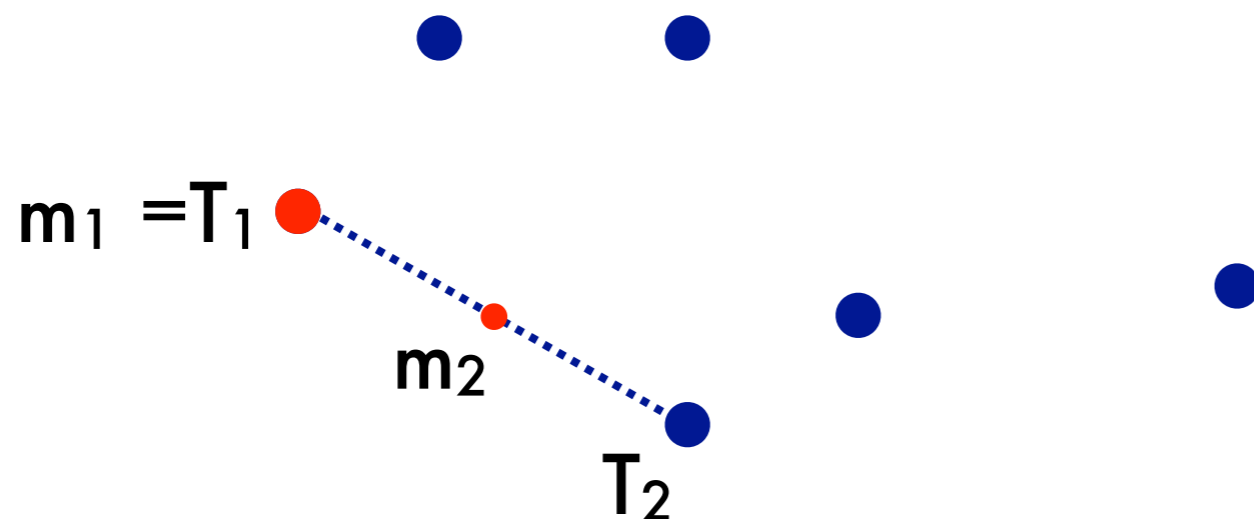
- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )

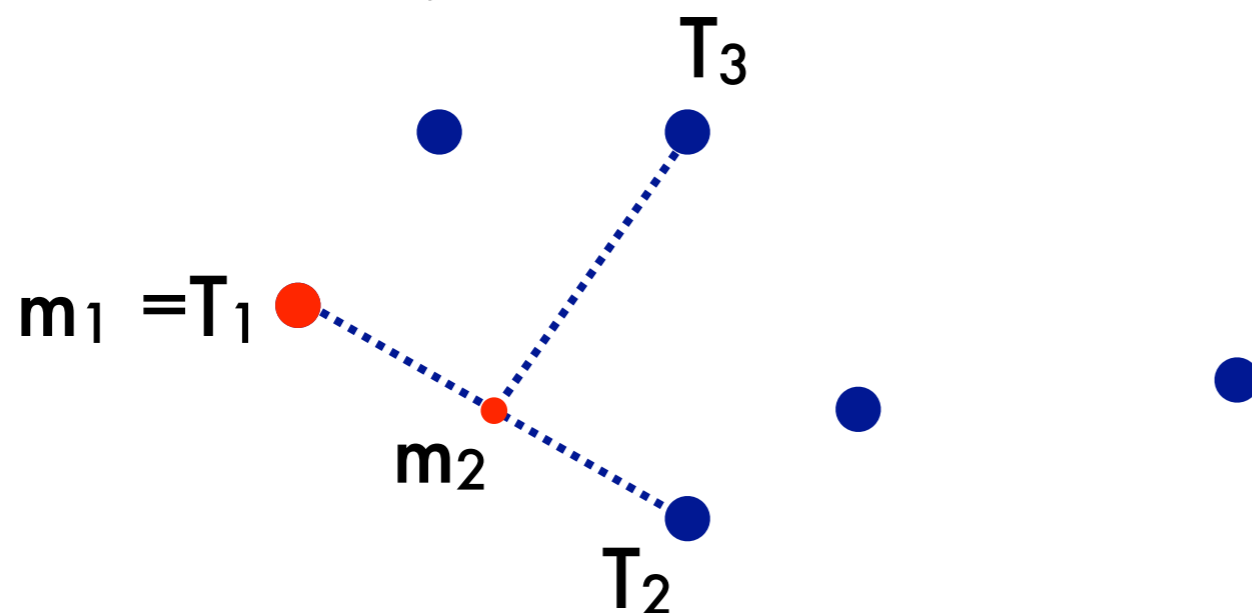




# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

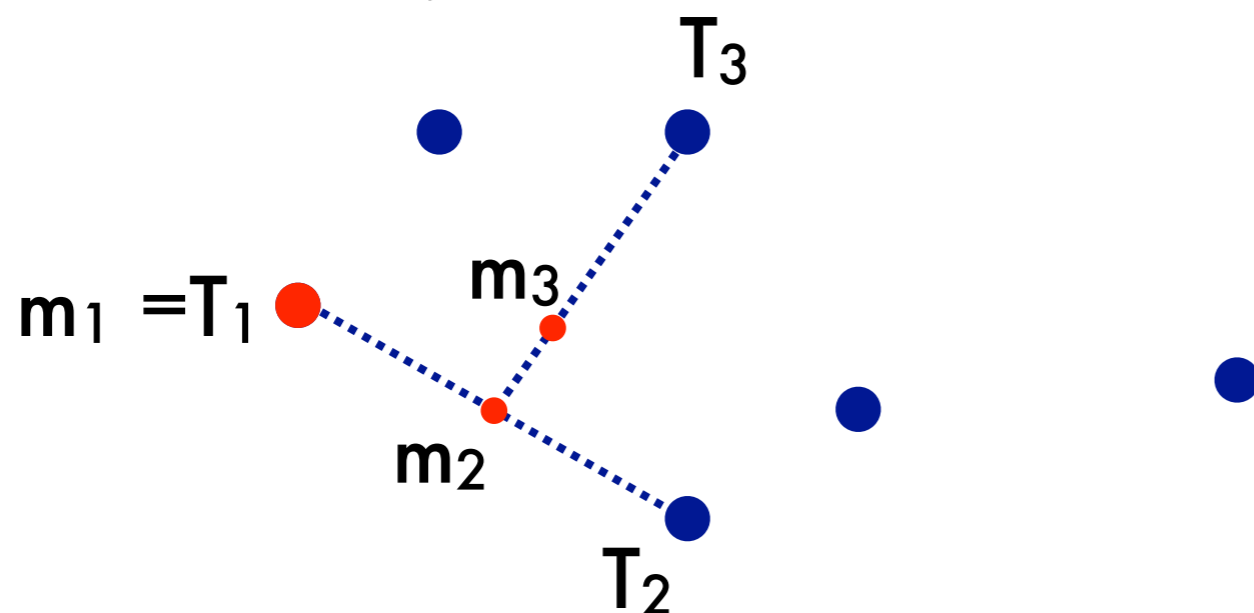
- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

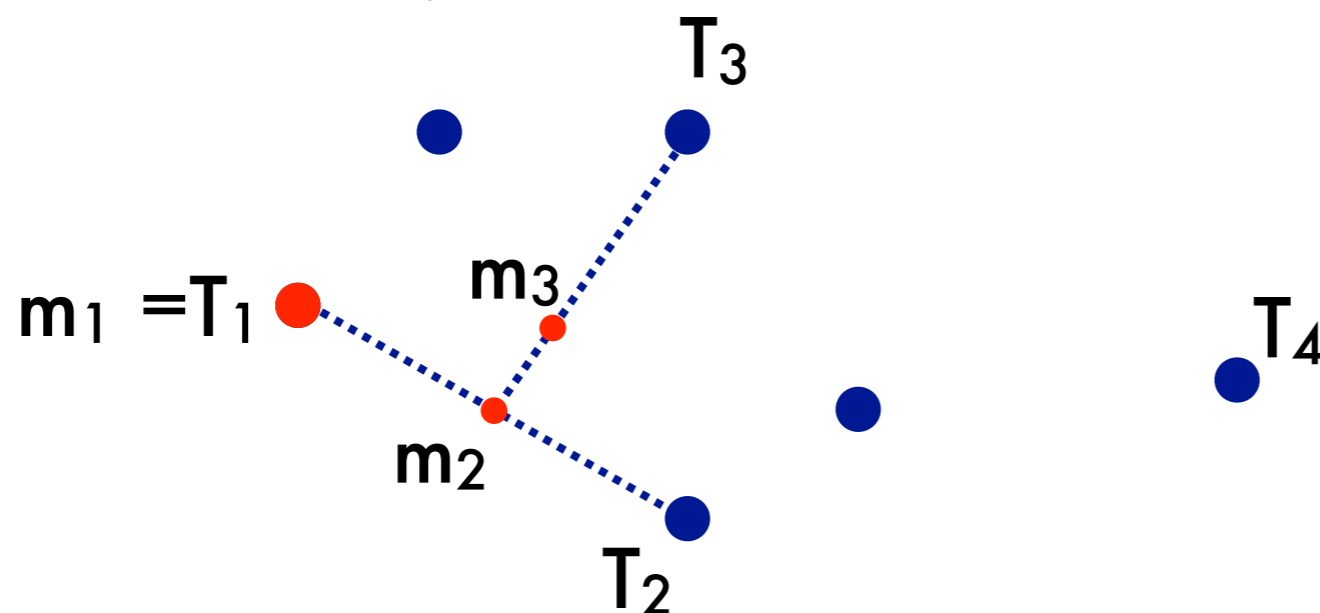
- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

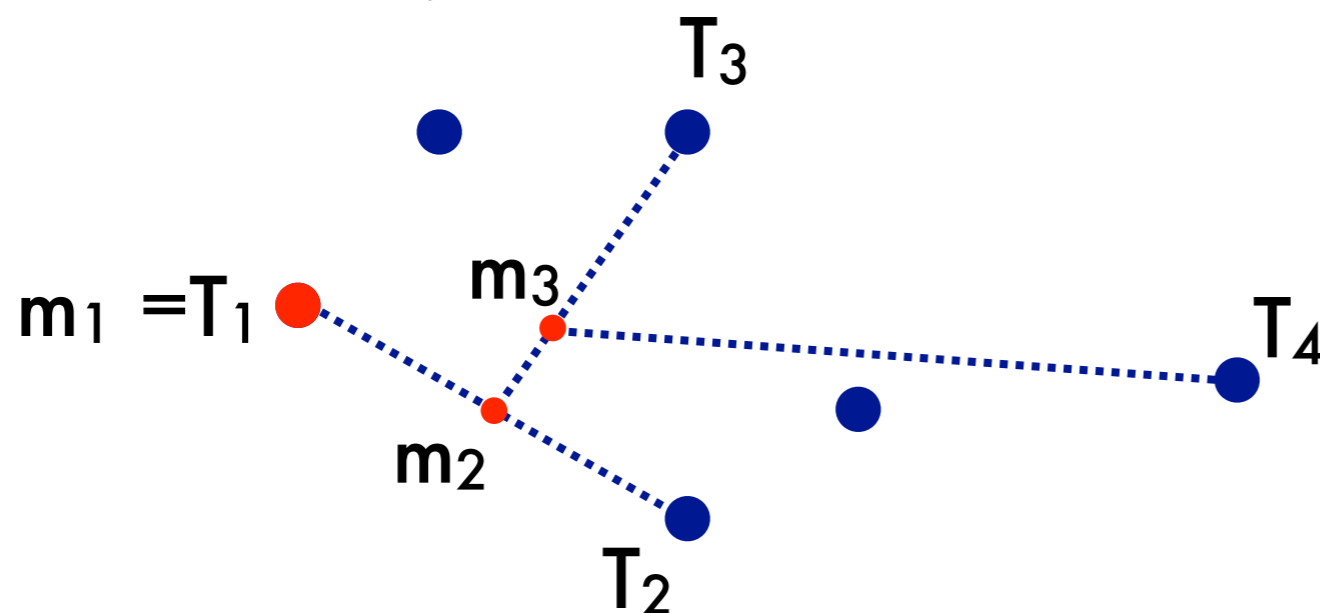
- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

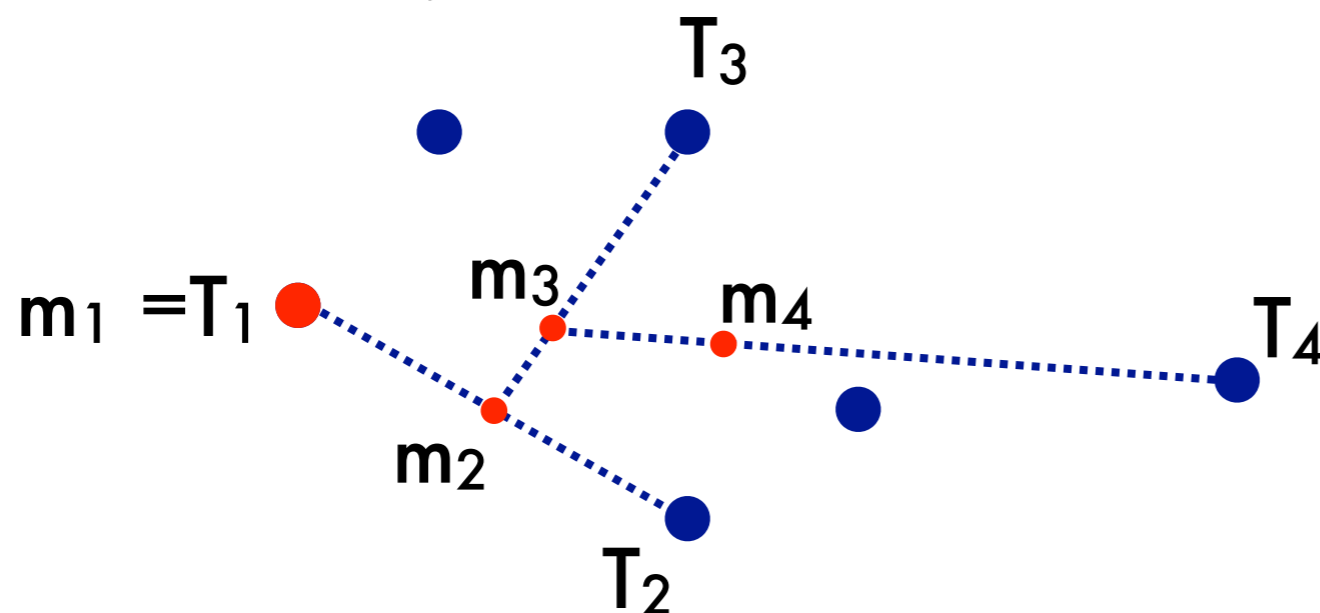
- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )

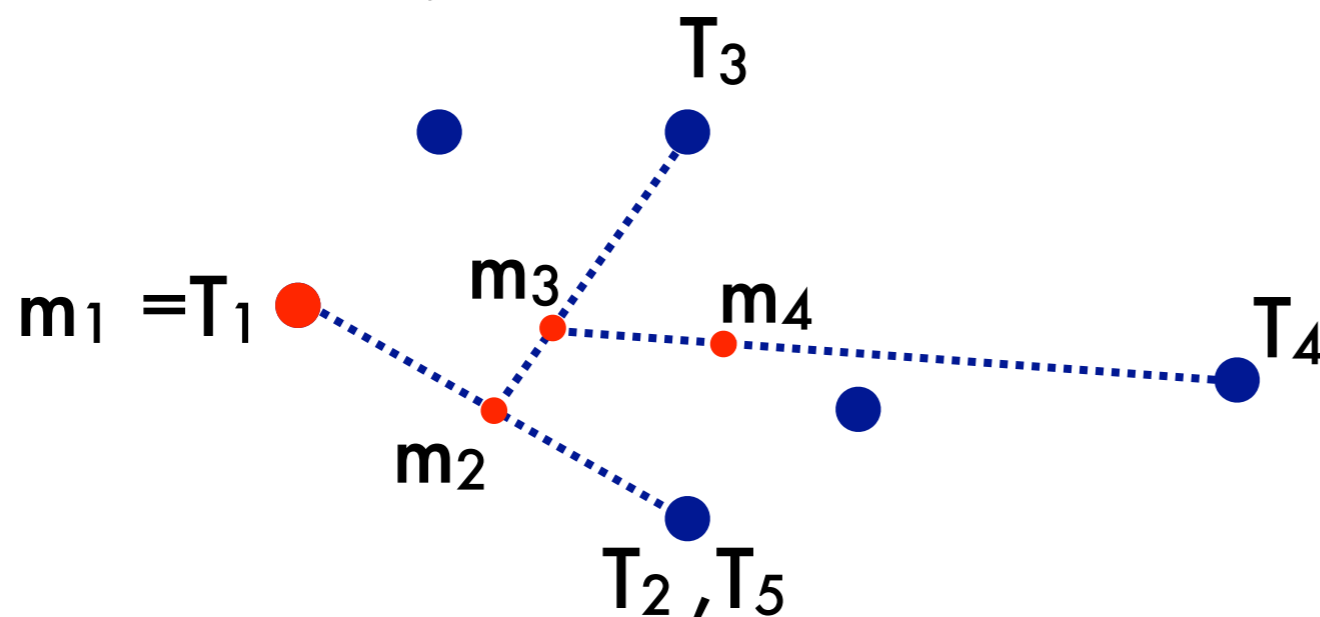




# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

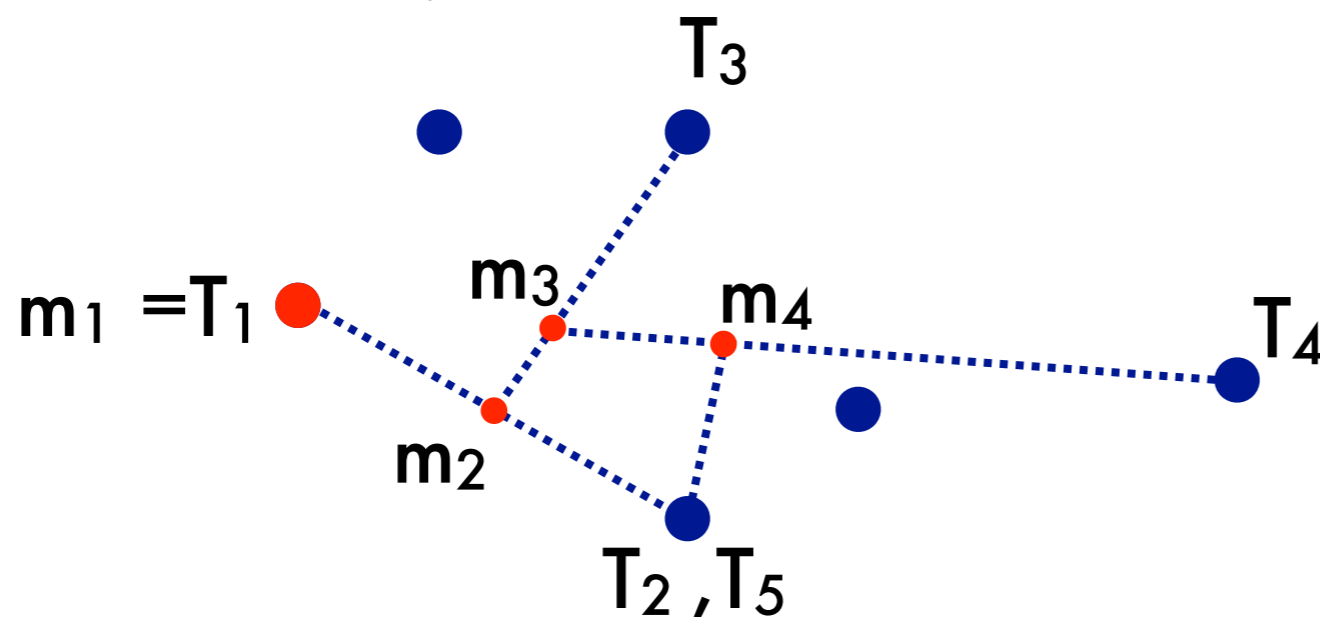
- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )



# Computing Mean

**Theorem** (Sturm, 2003): the following algorithm converges to the mean tree:

- $m_1 = T_1$
- $i^{\text{th}}$  iteration :
  - randomly choose tree  $T_i$  from tree set with replacement
  - $m_i = \frac{1}{i}$  (geodesic from  $m_{i-1}$  to  $T_i$ )

