

# Phylogenetic tree

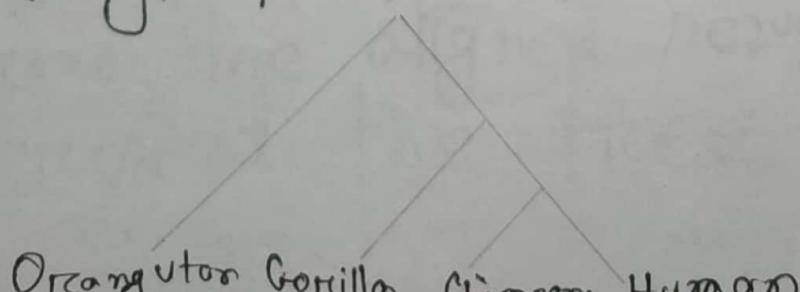
## Phylogenetic tree / Phylogeny:

A phylogenetic tree also known as a phylogeny that depicts the line of evolutionary descent of different species, organisms or genes from a common ancestor.

- \* The origin and evolution of species

has always been a mystery & a fascinating research area.

- \* Represents the evolutionary history of a group of organisms.



NB:

- \* every internal node  $\rightarrow$  hypothetical ancestor
- \* every leaf node  $\rightarrow$  species / taxa

When to say two individuals are different species?

Reproduction between two different species is not possible.

Phylogeny reconstruction:

Phylogenetic tree can be constructed from

1. Morphological data
2. Biomolecular data/sequences

\* tree can be constructed from  
various types of data

→ distance based }

→ character "

→ DNA/protein sequences

→ gene order

## Morphological data:

→ Form and structure of organisms  
and their specific structural  
features.

→ Number of legs, color of the eye  
etc.

## Biomolecular data/sequences

DNA, RNA, amino acids etc.

## Steps to construct phylogenetic tree (from sequence)

- \* acquiring a set of homologous DNA  
or protein sequences
- \* align those sequences (MSA)
- \* construction of phylogenetic tree  
from the aligned sequences
- \* present the tree

## ■ Applications of phylogenetic tree:

- gene function identification
- infer the amino acid sequence
- extinct proteins
- track the evolution of disease
- investigation of criminal case

■ Gene tree: A phylogenetic tree that depicts how a single gene has evolved in a group of related species.

■ Species tree  
Pattern of branching of species lineages via speciation.  
It shows the relationship between

different species

■ Discordance:

Gene tree don't necessarily  
Show the same branching pattern  
as their containing species tree.

■ Reason for gene tree discordance

→ Deep coalescence or incomplete  
Lineage Sorting (ALS)

→ Gene duplication/ extinction

→ Horizontal gene transfer  
(HGT)

→ ~~Estimation error~~ may also  
introduce discordance

# Speciation (outgroups)

## ■ Coalescent point

Point in the past in which

all gene copies arise from  
mutting a single gene copy.

## ■ Deep coalescence / incomplete lineage

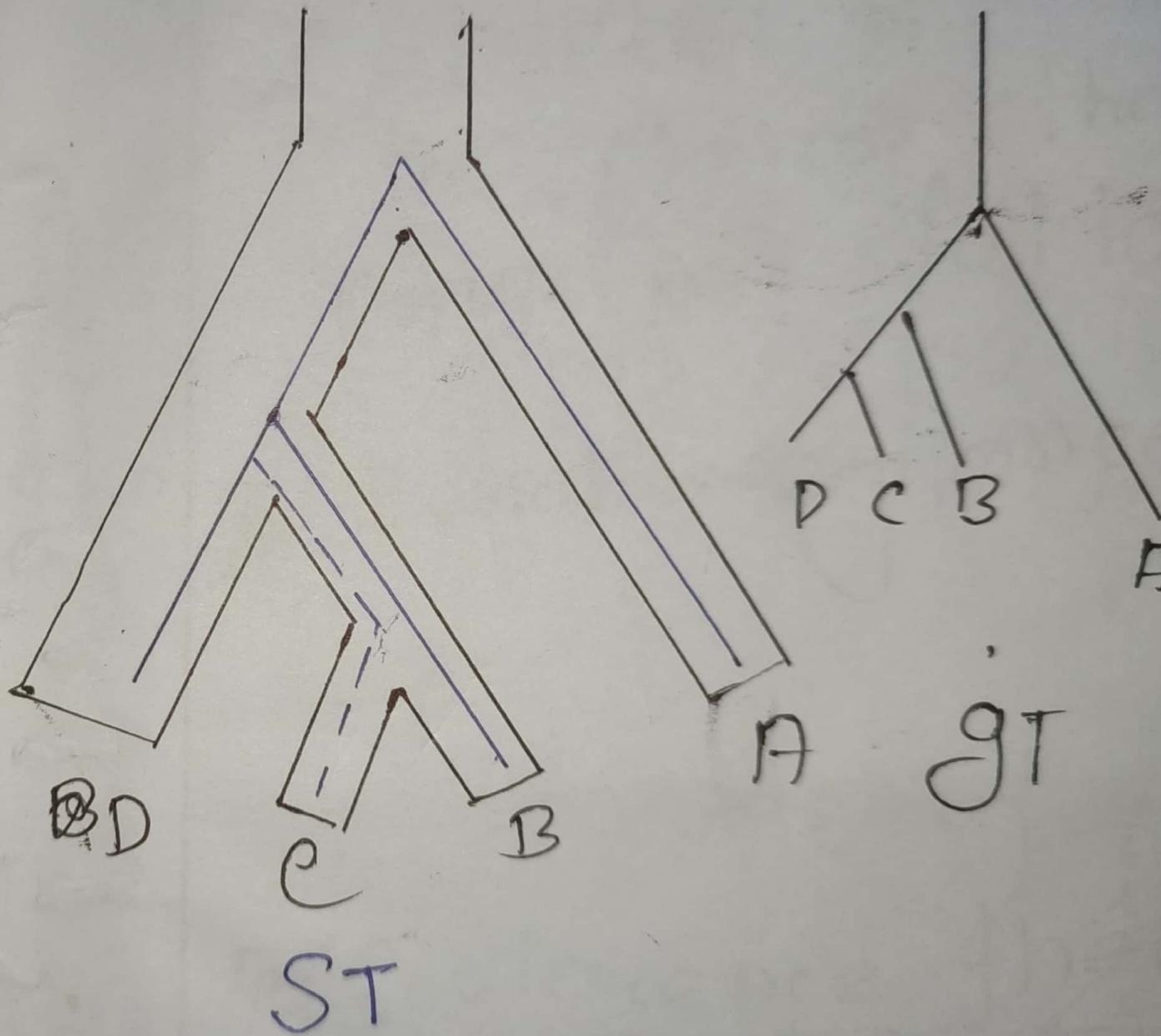
Gene copies fail to coalesce  
in the speciation point.

Genes copies at a single locus extends  
deeper than the speciation events.

## ■ things that influence the coalescent points

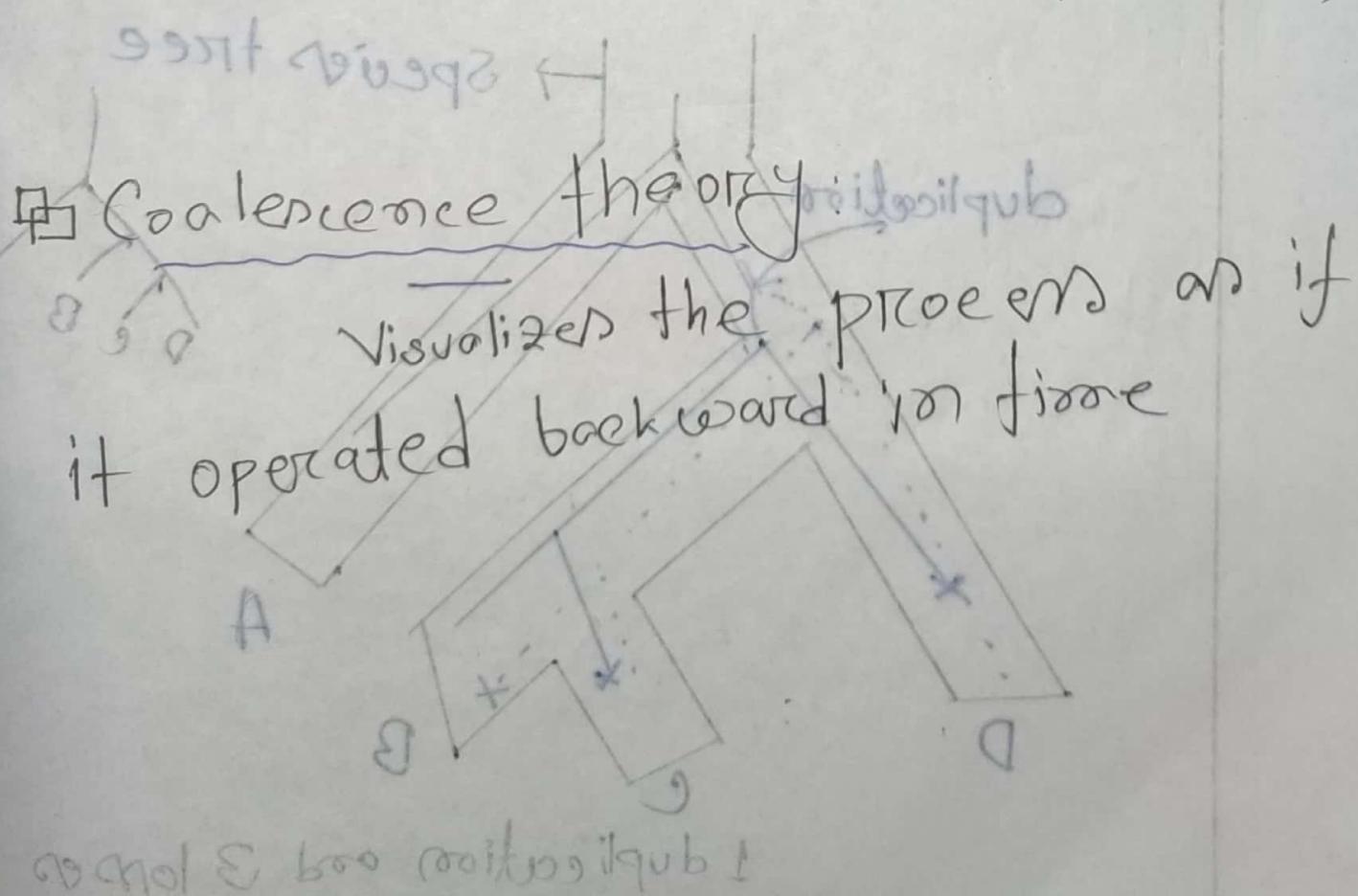
### Population size:

\* Smaller population size and longer branches increases the



## Possibility of coalescence

- \* Larger population size and shorter branches increase the chances that gene copies will fail to coalesce.
- 2) Evolutionary mechanisms (drift, type of selection)

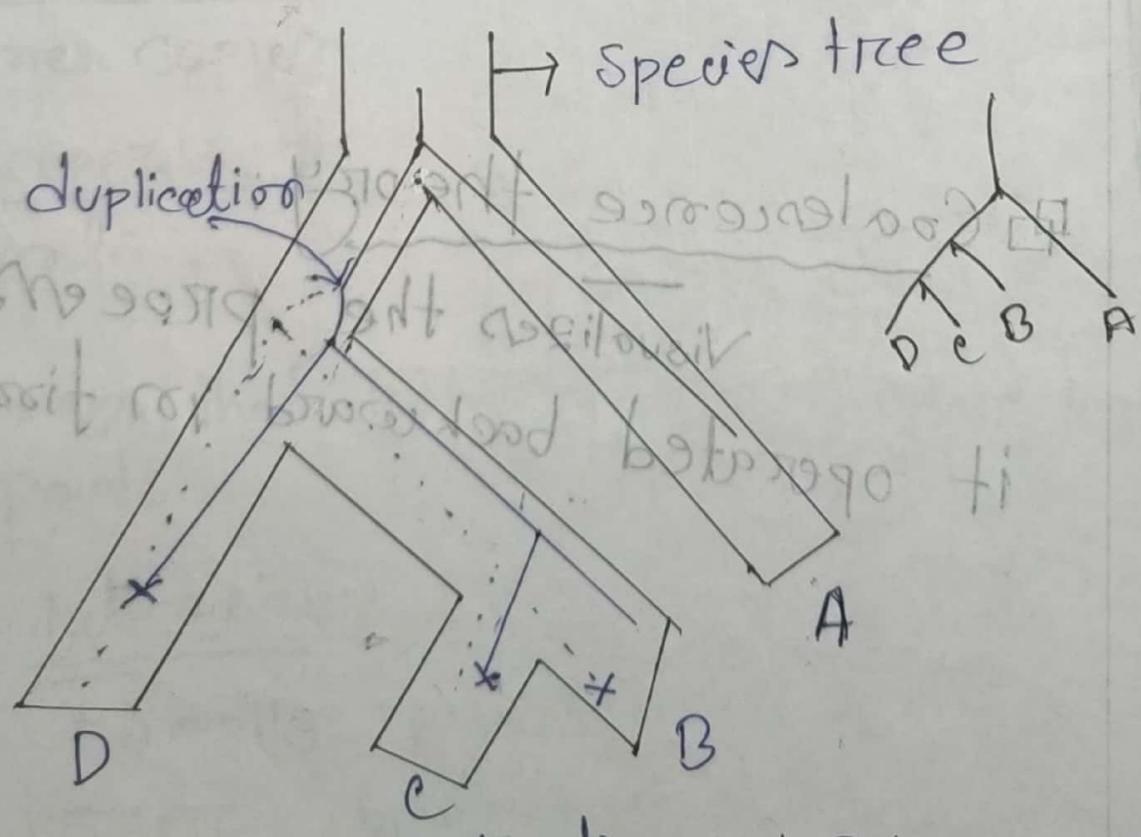


Gene duplication/loss → Duplicatio



A gene might get duplicated and both copies descend and evolve independently.

Discordance can occur if some sampled copies come from one locus and others come from another locus.



Confine: limit / Space

## Horizontal gene transfer

Renegade genes somehow break the confines of the

Species lineages and moved horizontally across the phylogeny

NB network instead of tree

Networks reflect co-type distribution  
of sites horizontally  
across populations



2023/term : unit 00

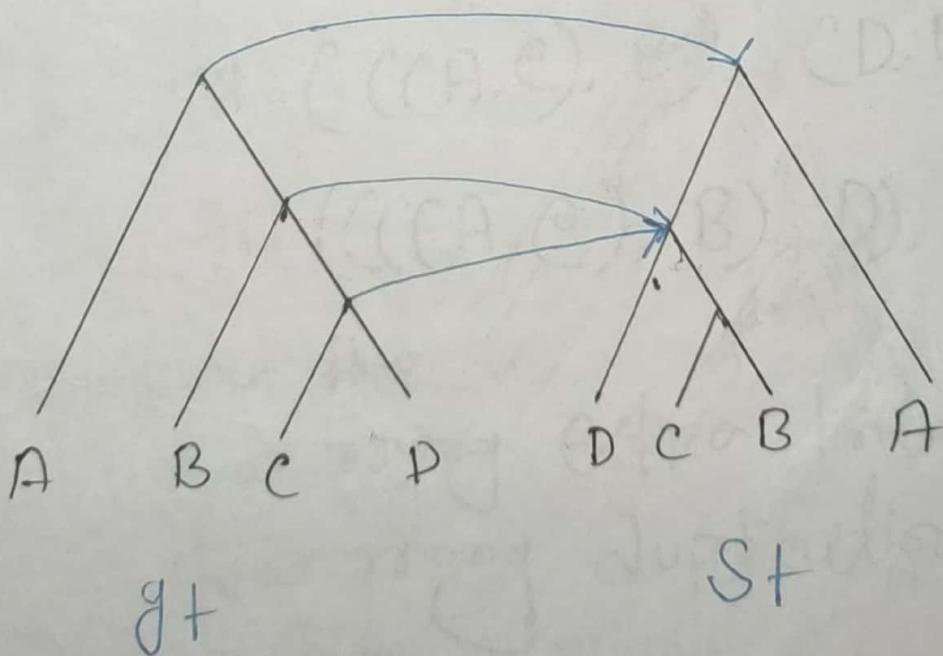
### ■ Reconciliation:

An approach to connect the history of two or more coevolving biological entities.

A phylogenetic tree representing the evolution of an entity can be drawn within another phylogenetic tree representing encompassing entity to reveal their interdependence and the evolutionary events.



- Most recent common ancestor (MRCA) /
- \* Last common ancestor (LCA)
- MRCA/LCA mapping Cgt → St
- \* An algorithm to find optional reconciliation:
  - find number of duplicate (minimum)
  - find the location of duplicate



Theorem / duplication node:  $\forall v \in V$

An internal node  $v$  of  $G$  is a  
duplication node if and only if

$$M(v) = M(\omega) \text{ for some child } \omega$$

of  $v$ .

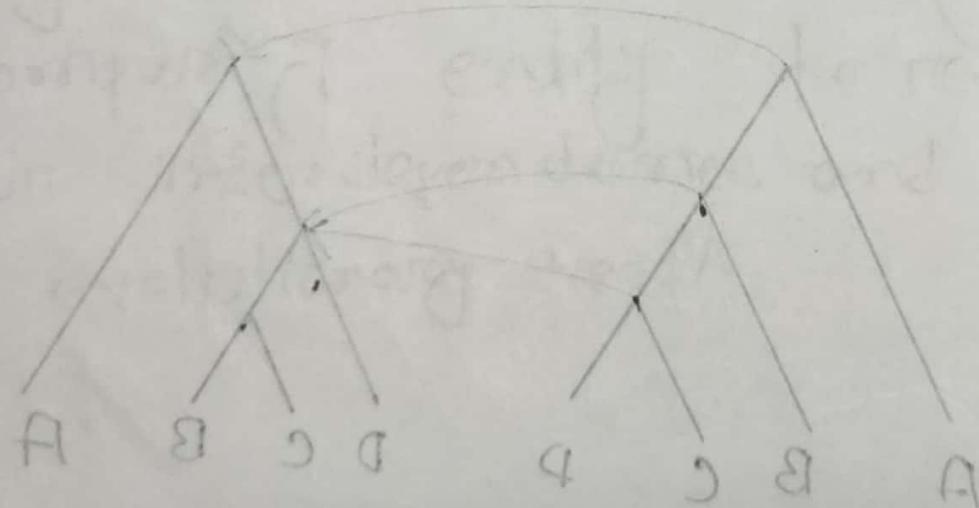
mapping of  $v$

is possible to obtain both

together free

the same entity and

the other



+2

+B

Q Consider the following gt and st  
Specification.

P-1

$$gt = (((CB,D), C), A)$$
$$st = (((A,B), C), D)$$

Newick notation

P-2

$$st = ((A,B), ((C,D), E))$$
$$gt = (((CA,C), (D,E)), B))$$

P-3

$$st = (((CA,C), B), (CD, E))$$

$$gt = (((((A,C), B), D), E),)$$

Find out, the

How many extra lineage?

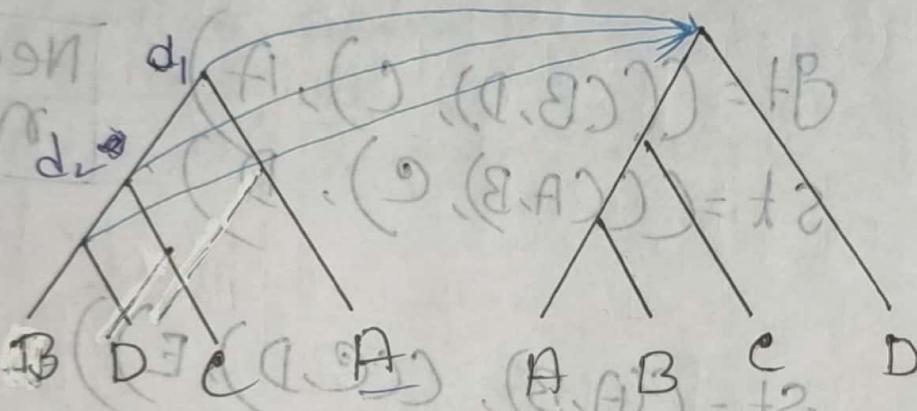
How many duplications and losses?

Soln

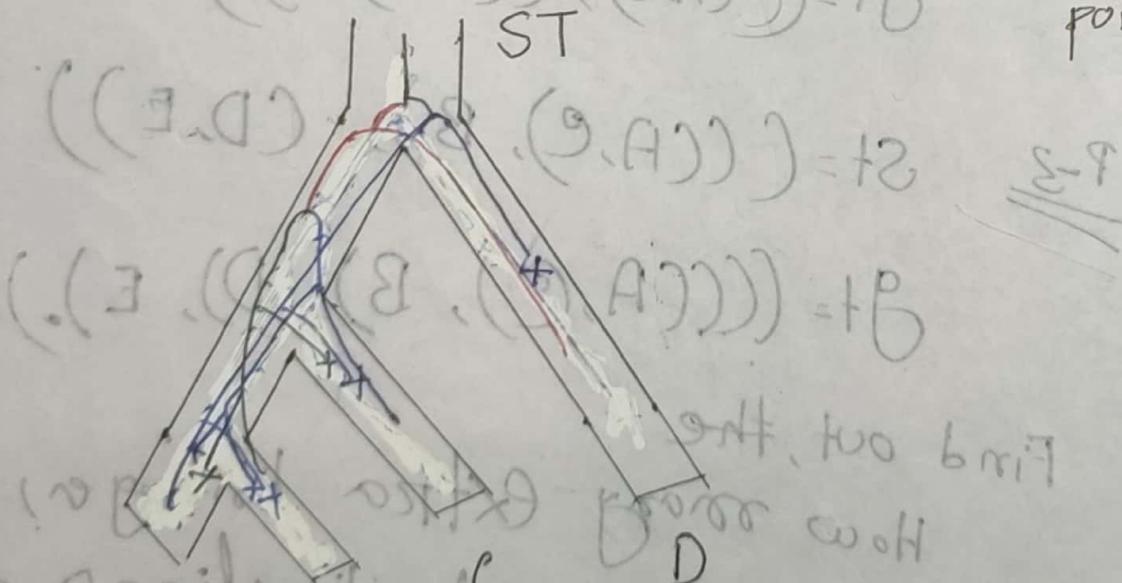
Soln  $\rightarrow$  B is a subset of  $A \cap C$   $\Rightarrow$  B is a subset of  $C$

~~P-1~~ / g+

St. George's



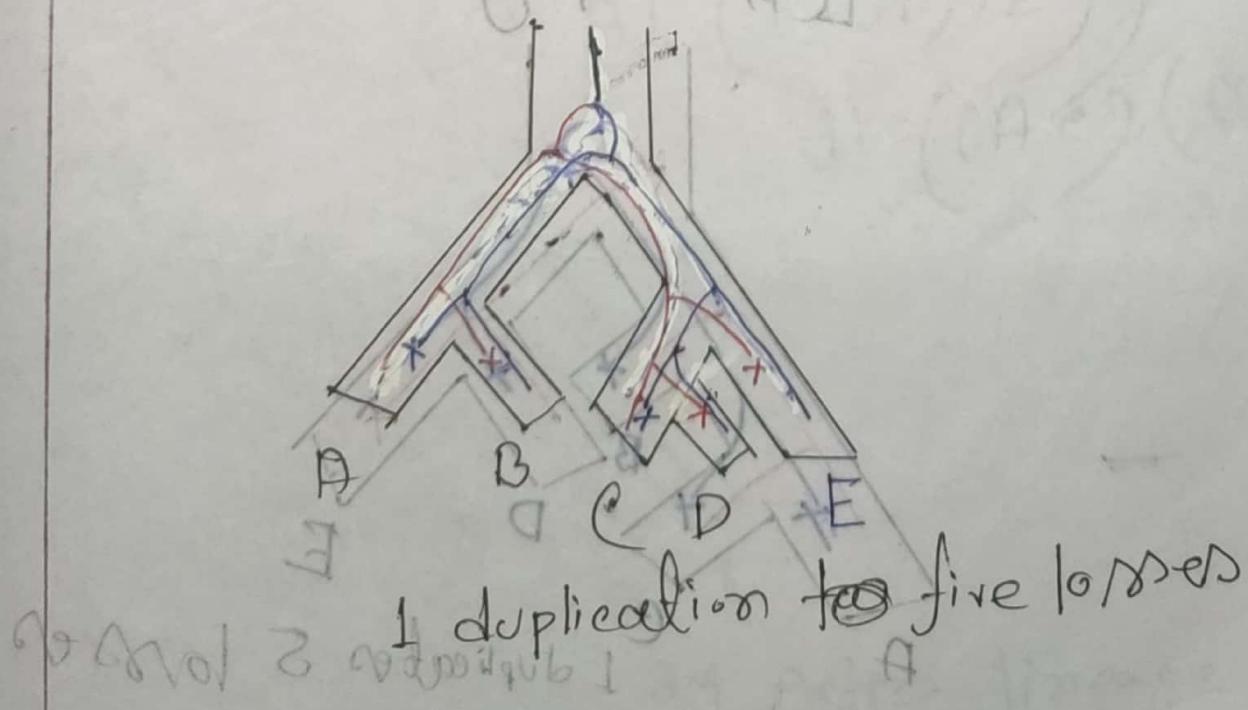
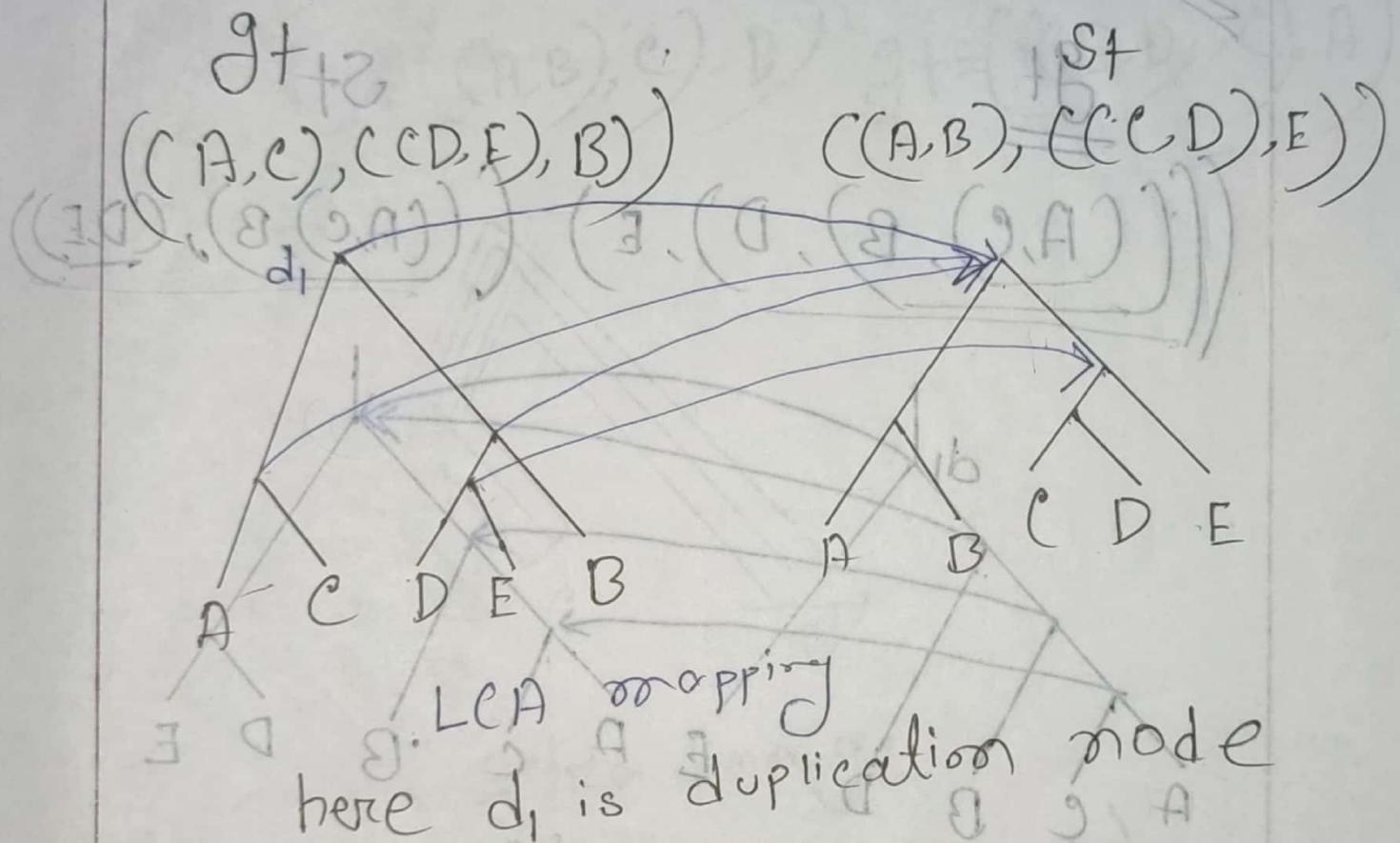
\* Here  $d_1$ ,  $d_2$  are two duplication points.



2 duplicate ~~and~~ ~~for~~ passes.

P-2

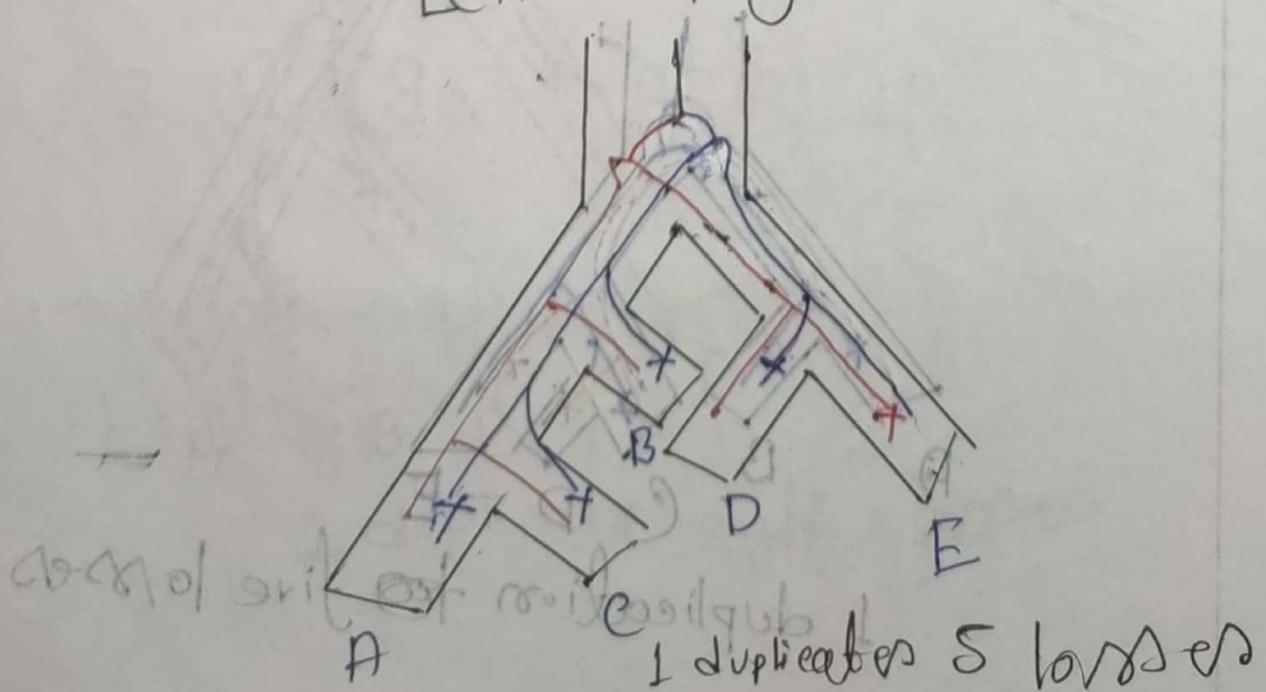
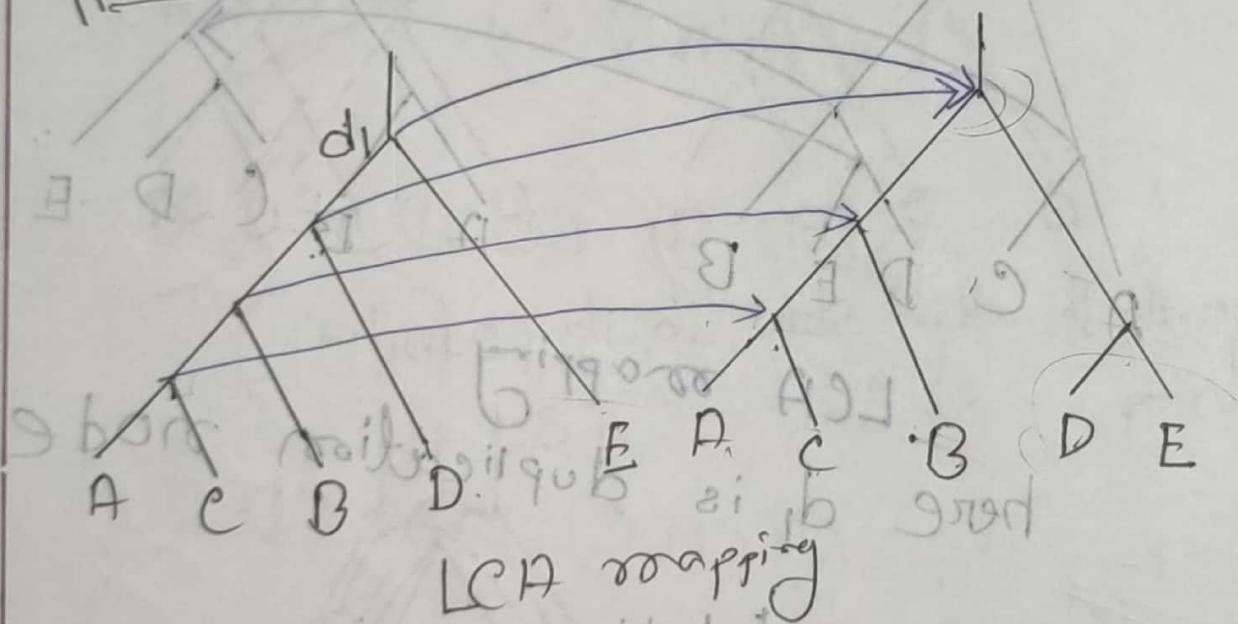
S-9



P-3

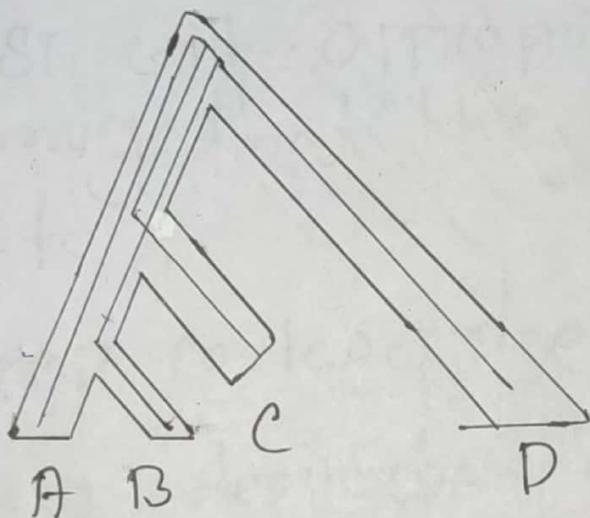
S-9

$$\begin{array}{c} \text{gt} \\ \text{gt} \\ ((\underline{\underline{((A,B),C)}},D),E) \end{array} \quad \begin{array}{c} \text{gt} \\ \text{gt} \\ ((B,(\underline{\underline{C,D}})),(\underline{\underline{E,A}})) \end{array} \quad \begin{array}{c} \text{gt} \\ \text{gt} \\ (((A,C),B),(\underline{\underline{D,E}})) \end{array}$$



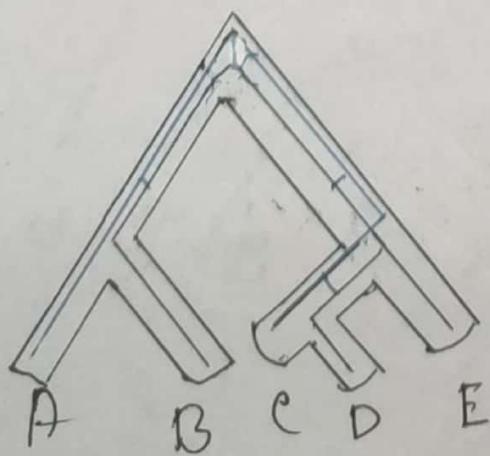
P1

$$st = ((A \cdot B), C) \cdot D \quad gt = ((C \cdot D), C) \cdot A$$



2 extra lineage

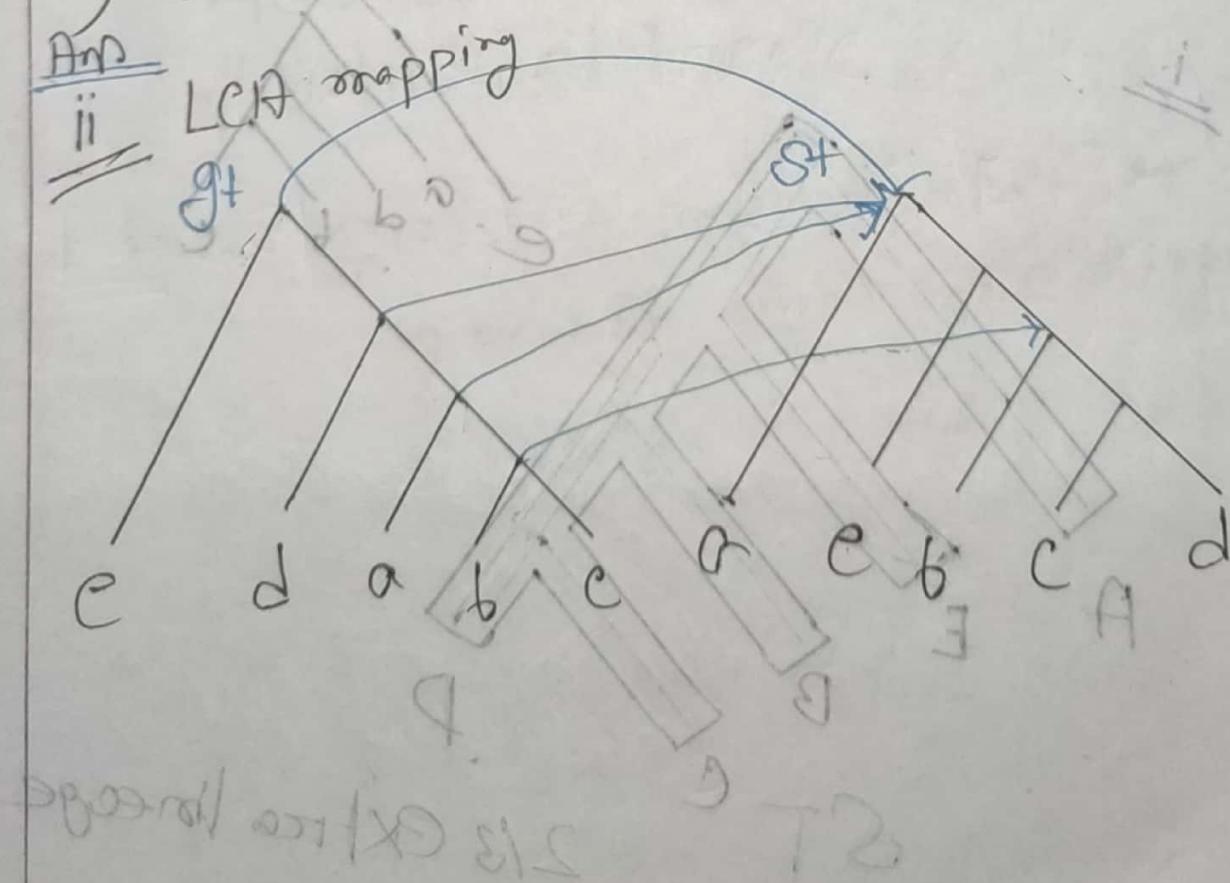
P2     $st = (CA, B), ((CC, D), E)$   
 $gt = ((A, e), ((D, E), B))$

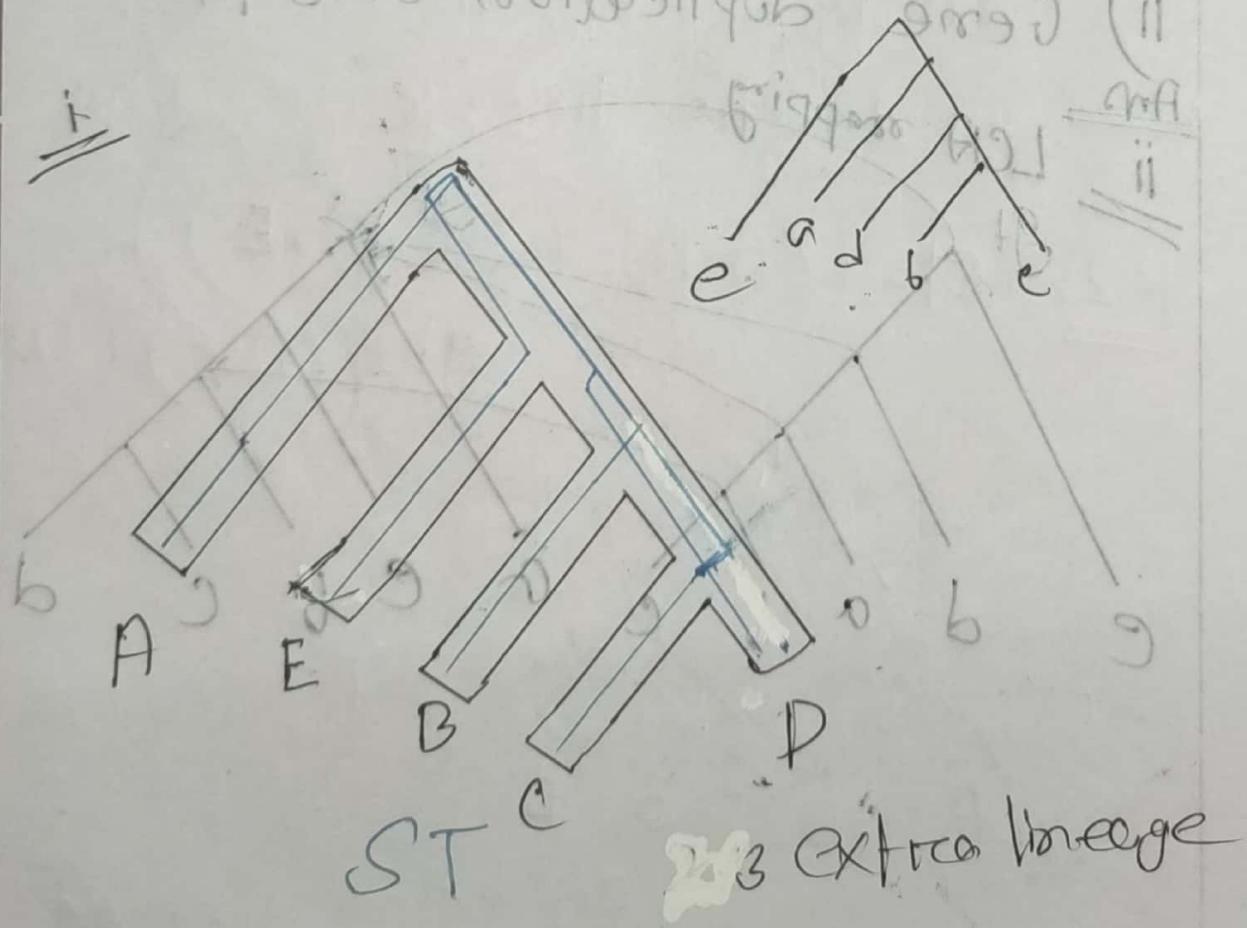
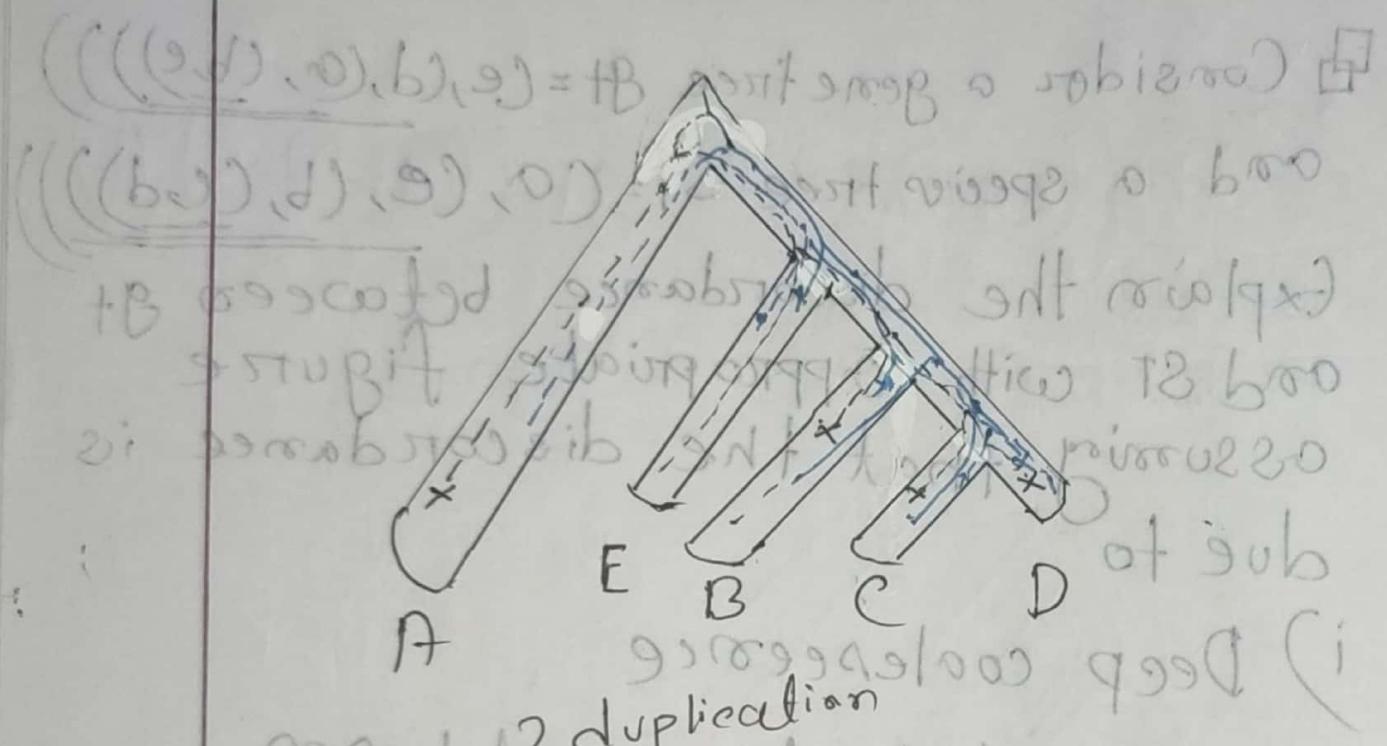


3/4 extra lineage

Consider a gene tree  $gt = (e, (d, (\underline{a}, (\underline{b}, e))))$   
 and a species tree  $ST = (a, (e, (b, (c, d))))$   
 Explain the discordance between  $gt$  and  $ST$  with appropriate figure  
 assuming that the discordance is due to

- Deep coalescence
- Gene duplication and losses





## Tree:

Mathematically a tree is a graph  $G = (V, E)$ , (where  $V$  is the vertex set, and  $E$  is the edge set) that is connected and acyclic.

## # types of tree

\* Rooted tree

\* Unrooted tree

# Tree can be represented by Newick representation/notation

→ Newick representation

→ Cladogram

## Rooted tree:

In a rooted tree  $T$ , we can orient the edges in the direction of the root  $r = r(T)$  so that all vertices other than  $r$  have outdegree one.

\*  $r$ oot is denoted by  $r(T)$

\* set of vertices of  $T$  "  $V(T)$

" of edges of  $T$  "  $E(T)$

" of leaves of  $T$  "  $L(T)$

For all nodes  $v \neq r$  there is an unique vertex  $\omega$  such that

$v \rightarrow \omega$  is an arc,

$\omega$  is called the parent of  $v$

$v$  is " the child of  $\omega$

siblings: two or more vertex sharing vertices

and the same parent are called siblings

leaf/tip: A vertex without any children

is called a leaf / tip

internal nodes: All the remaining nodes with-

ut leaf.

polytomy: A vertex with more than

two children

binary - bifurcating or fully resolved:

A rooted tree that has no

polytomy

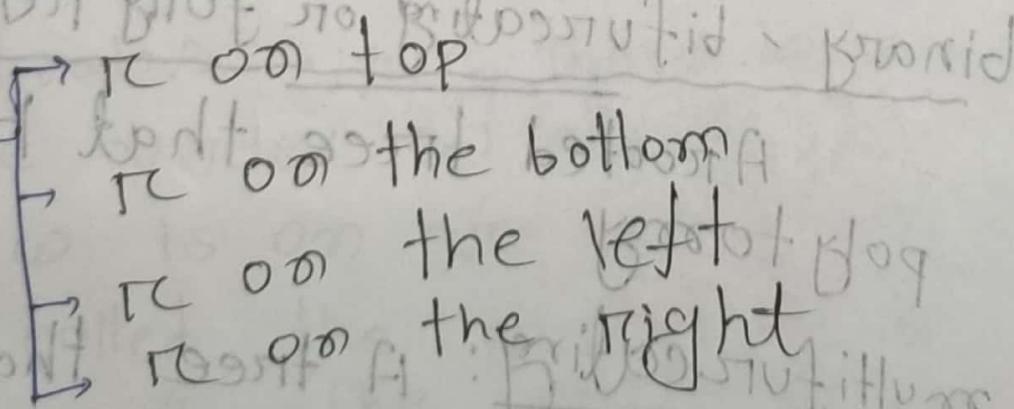
multifurcating: A tree that is not

binary.

Leaves: In a phylogenetic tree, the leaves represent the taxa of interest.

Ancestors: Internal nodes represent the ancestors of the taxa at the leaves:

Note: There are multiple ways to draw the topology of a rooted phylogenetic tree



Two way to represent rooted tree

~~A root is a node which has no parent~~

~~The root of tree is a node which has no parent~~

~~Rooted tree is a tree where every node has a parent~~

(A, B, C)

Access to children is from A to B

B " " to "

$$(A \cdot B) = (B \cdot A)$$

Note:

~~if two cars meet the road will~~

~~surface to ground~~

~~for a root position design~~

~~float foot~~

## Newick notation

The Newick notation for a rooted binary tree with subtrees A and B is given by  $(A', B')$ , where,

$A' \rightarrow$  newick notation of subtree A

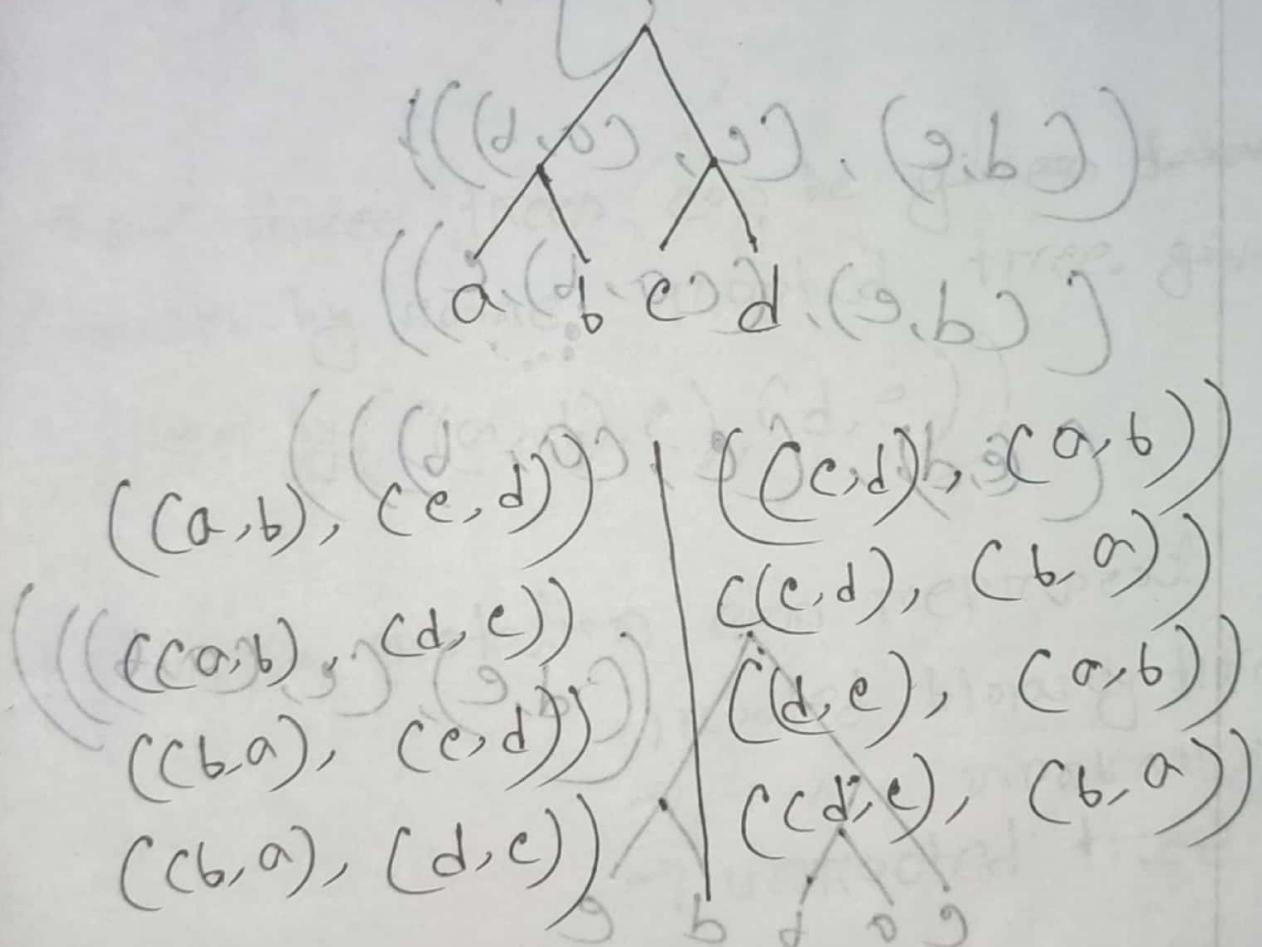
$B' \rightarrow$  " " of " B

Note:  $(A', B') = (B', A')$

→ don't care about the left-to-right ordering of subtree;

→ Newick notation for a leaf is leaf itself

Find Newick notation of following rooted binary tree



N.B.: There is more than one newick string for a (binary) rooted tree.

Rooted trees for the following newick string:

~~((c,d,e), (c,a,(c,b))))~~

~~((c,d,e), ((c,a,b), c)))~~

~~((c,d,e), (c,c, (c,b))))~~

~~((c,d,e), ((c,b), (d,a)))~~

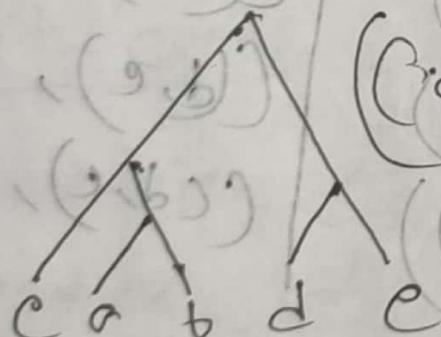
~~((c,d,e), ((b,g), (a,f)))~~

~~((c,d,e), ((g,h), (f,i)))~~

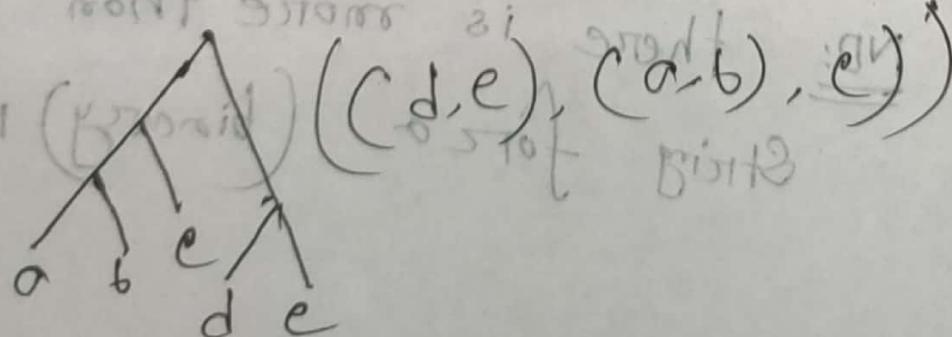
~~((c,d,e), (c,c, (c,b))))~~

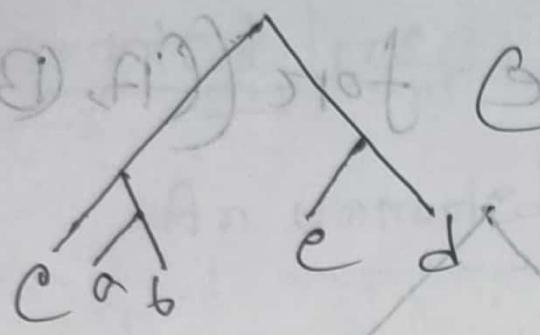
~~((c,d,e), ((c,b), (a,f)))~~

~~((c,d,e), ((g,h), (f,i)))~~



Rooted tree for the following newick string:

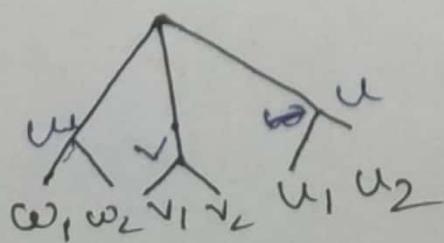




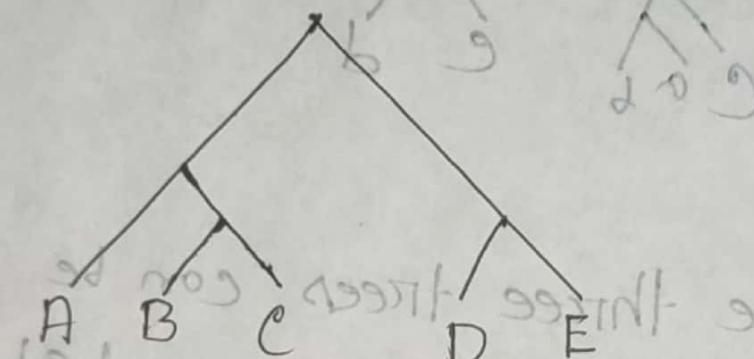
these three trees can be given by  
 shown by some rooted tree given  
 given by  $((a, b), c), (d, e))$

- \* Newick notation can represent
  - rooted binary tree
  - ... nonbinary
  - unrooted tree

■ rooted tree for the Newick string  
 $((v_1, w_1), (v_1, w_2), (w_1, w_2))$



~~Rooted tree for  $((A, B, C), (D, E))$~~



$((g, b), (g, (g, g)))$  is morib

to transport and visitors when \*

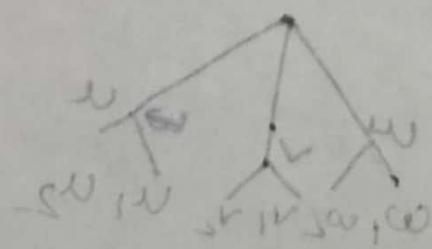
use broad based

\* Broad base

base below

is not good for broad base

$((w, w), (w, v), (w, v))$ .



## Unrooted tree:

An unrooted phylogenetic tree is

a type of phylogenetic tree,

that only describes the relatedness of a group of organisms without

making assumptions about ancestry.

It can be generated from a rooted tree by simply omitting the root.

If the root of T has two children, then these two children are made adjacent to each other.

- \* unrooted  $\rightarrow$  rooted (one to many)
- \* rooted  $\rightarrow$  unrooted (one to one)

NB

\* each rooted tree has an unique unrooted edge

so each sibling relationship is A version

\* each unrooted tree has multiple rooted version.

function: analysis of tree root

Outgroup to group A to

★ Taxa that are clearly

not closely related to the rest of the input taxa as the remaining taxa are to each other.

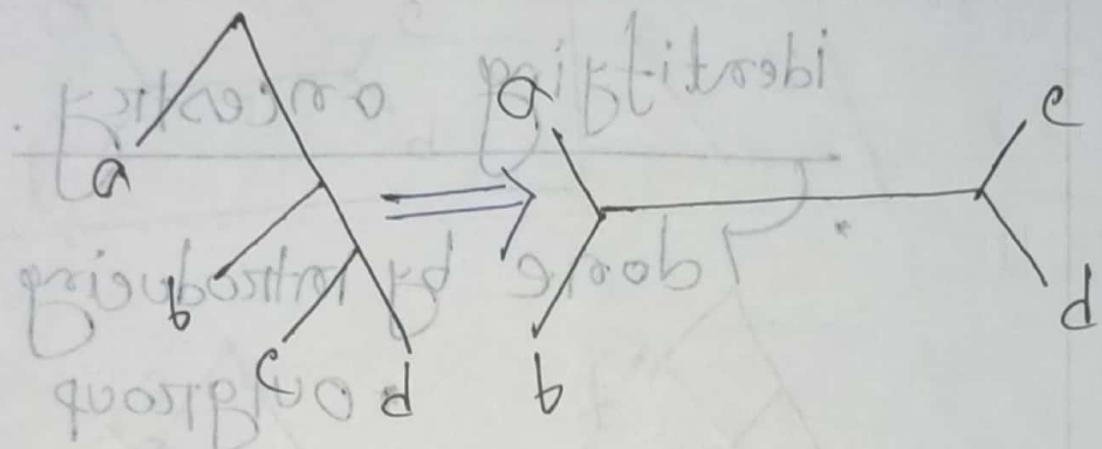
out and T to root of tree

and out event root

root of tree

▪ Unrooted tree of rooted tree  
(Rooted  $\rightarrow$  unrooted)

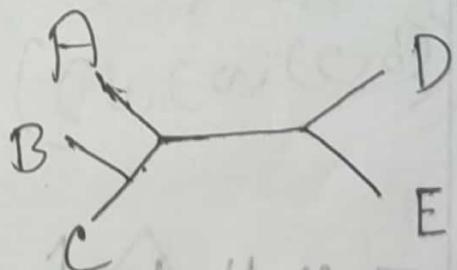
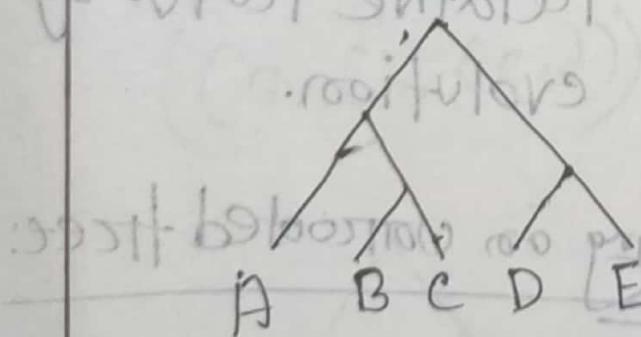
to remove some writing



$((a, (b, (c, d))))$

$((a, b), (c, d))$

to remove some writing



$((a, (b, c)), (d, e))$

remove some writing

remove some writing

remove some writing

remove some writing

田 Untrooted tree  $\Rightarrow$  Rooted tree

\* requires some means of identifying ancestry.

\* done by introducing an outgroup

(\* additional assumption about the relative rates of evolution.)

田 Methods for rooting an unrooted tree:

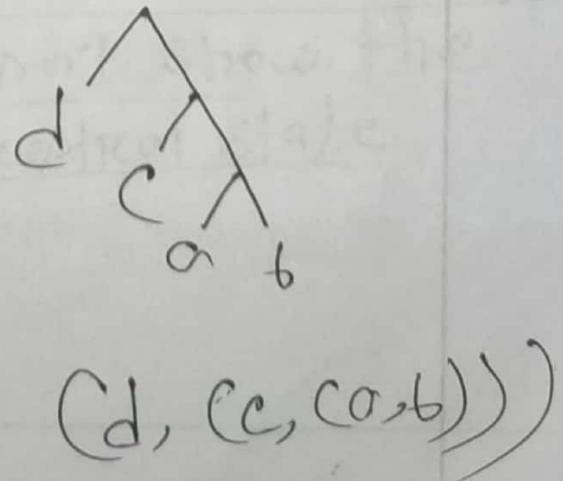
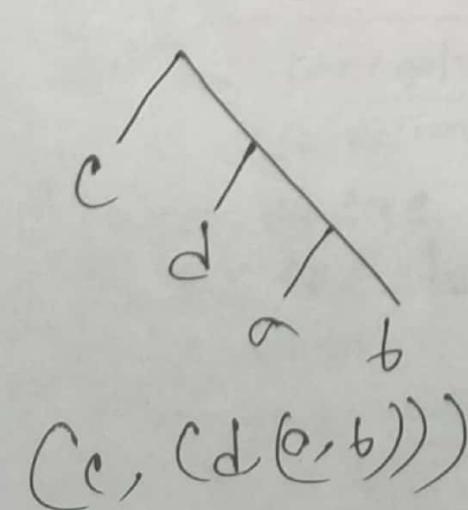
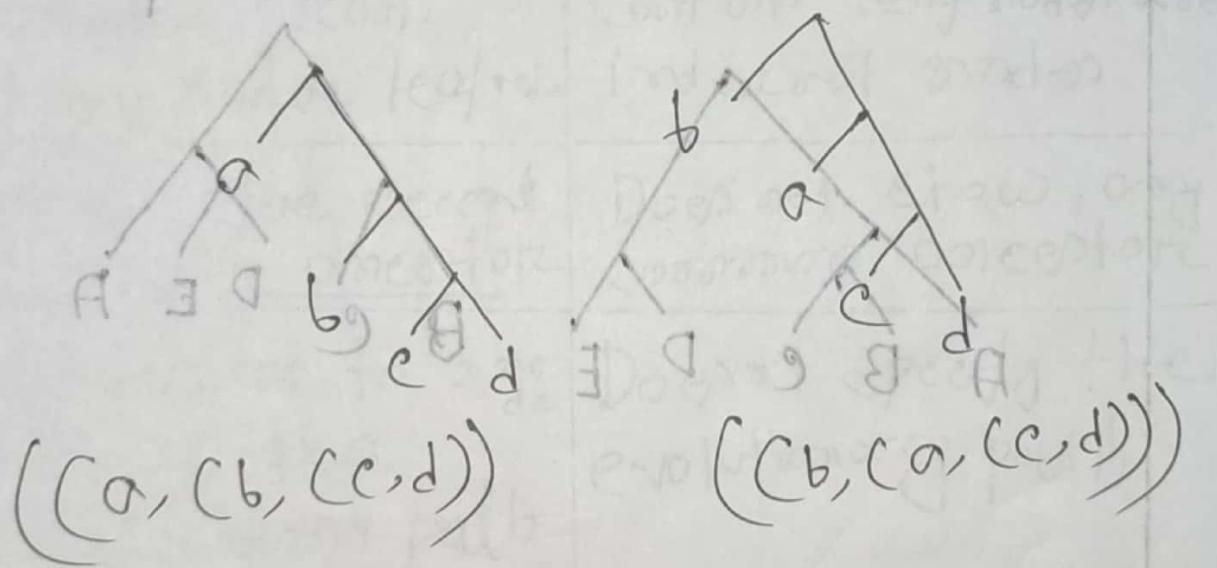
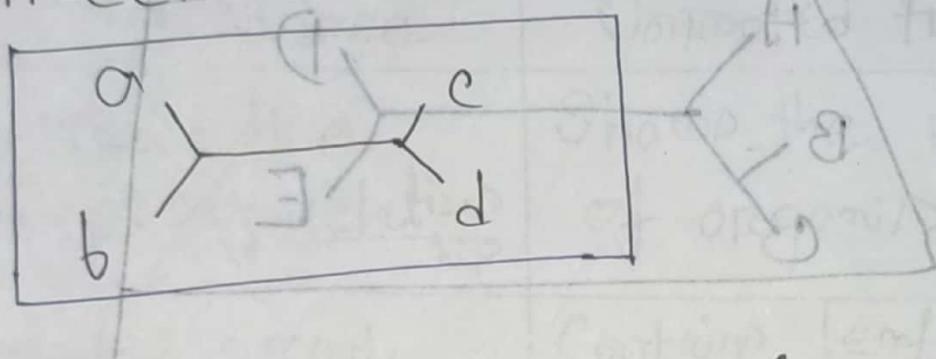
→ outgroup rooting

→ midpoint rooting

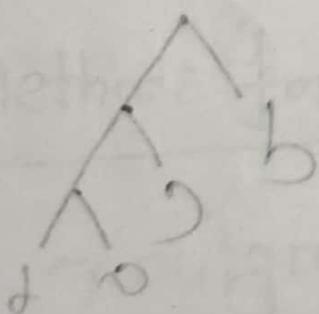
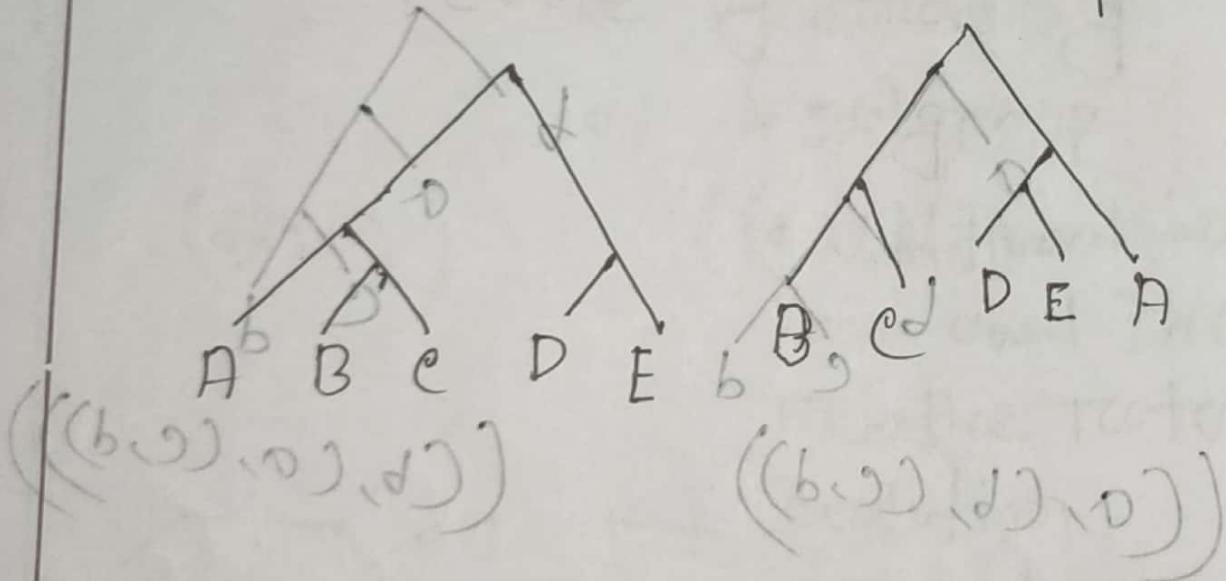
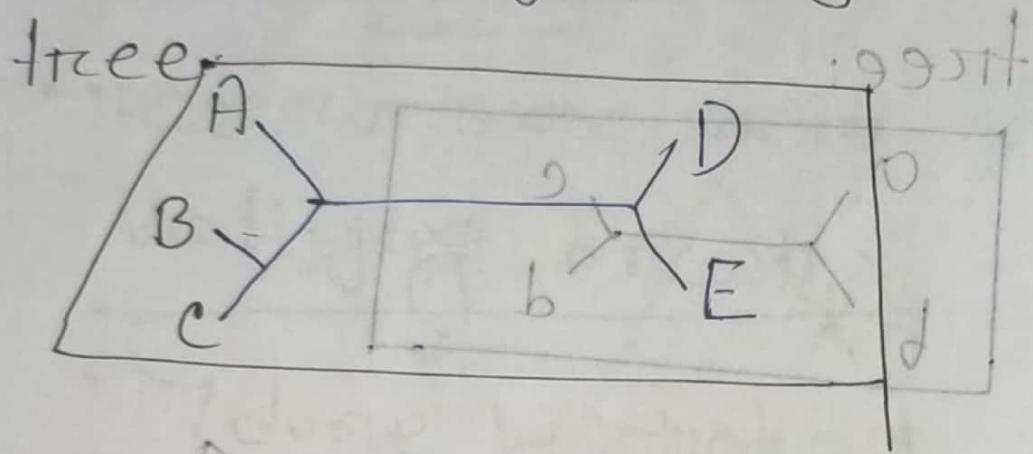
→ molecular clock rooting

→ Bayesian molecular clock rooting

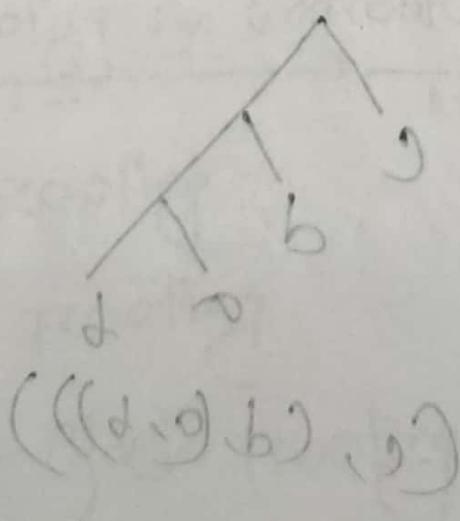
Rooted version of following unrooted tree:



Q Rooted tree of following unrooted tree



$((d, o), g), b)$



$((((d, o), b), g))$

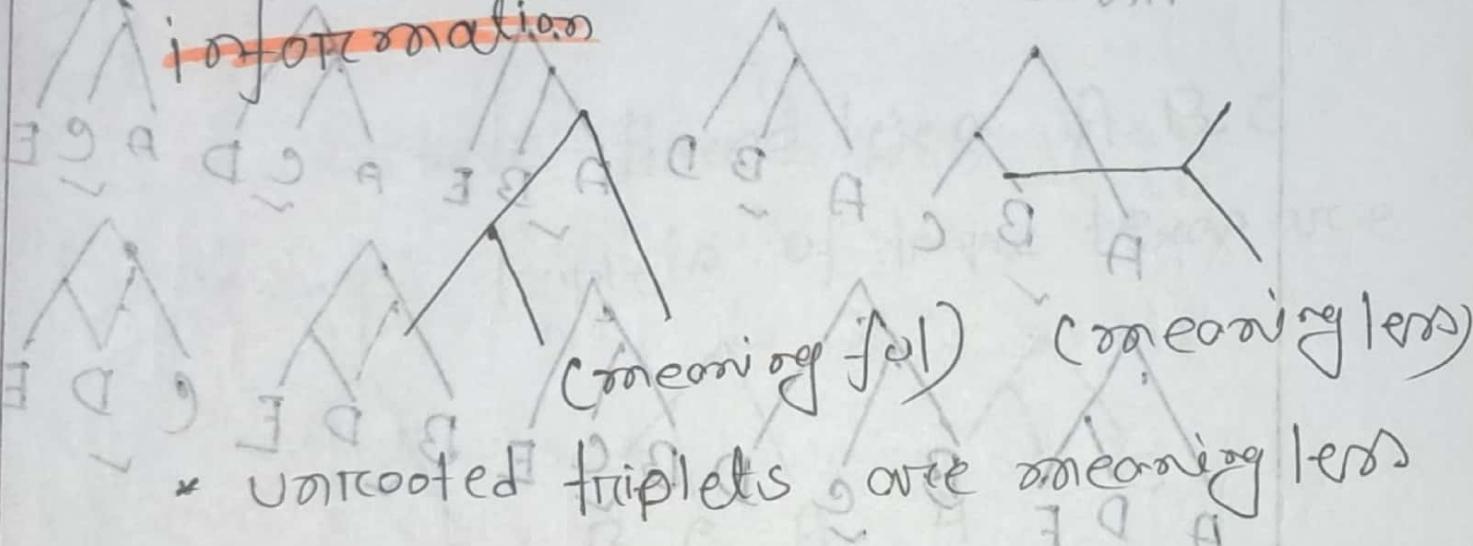
## ★ Rooted tree Vs Unrooted tree

<u>Rooted tree</u>	<u>Unrooted tree</u>
It shows the <u>ancestry relationship</u>	Showsthe relatedness of organisms
Contains a root, internal nodes, leaf node	Contains leaf nodes and internal nodes
Shows the <u>recent common ancestor</u>	Does not show any common ancestor
Each root to a <sup>age</sup> node Shows the <u>evolutionary path</u>	Doesn't specify the evolutionary path
Shows the <u>ancestral state</u> of organisms at the gene at the bottom to the terminal branches.	Doesn't show the <u>ancestral state</u>

■ Triplets: Rooted tree with three leaves  
(species)

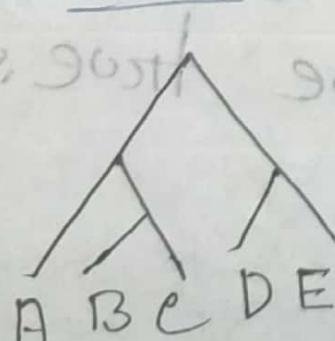
\* The most basic piece of phylogenetic

information



■ Decomposing input gene tree into HD

induced triplets:



$$\frac{n!}{3!(n-3)!}$$

NB:

→ Number of induced gene tree (triplet)

$$= \frac{n!}{3!(n-3)!}$$

# + exch

A B C

A B D

A B E

A C D  
A C E

A D E

B C D  
B C E

B D E

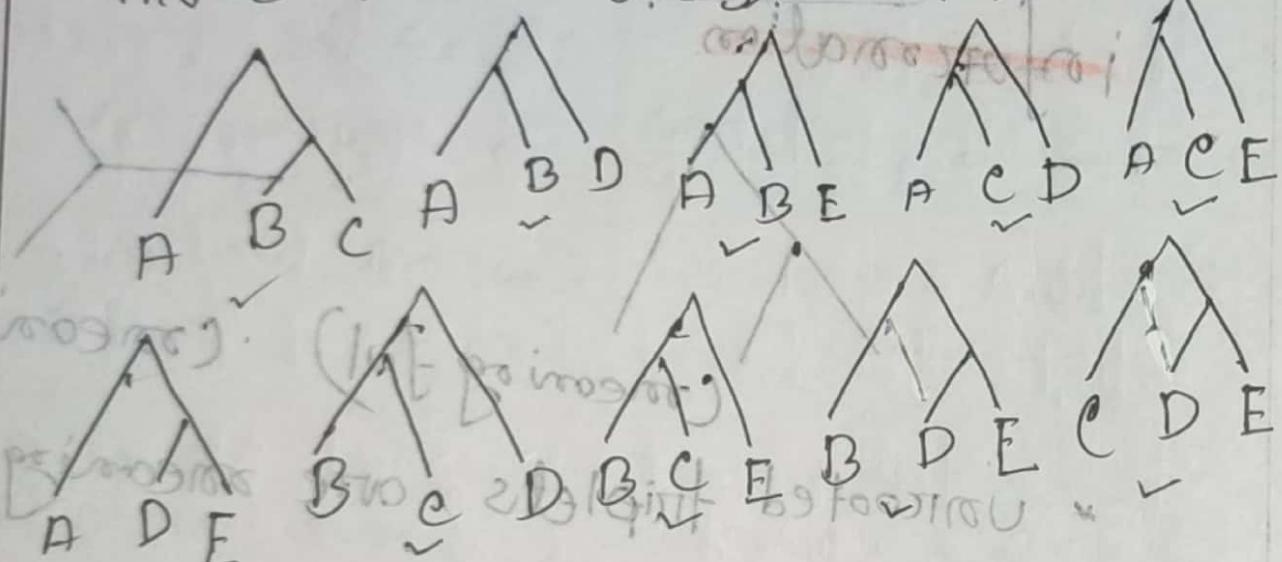
C D E

There are

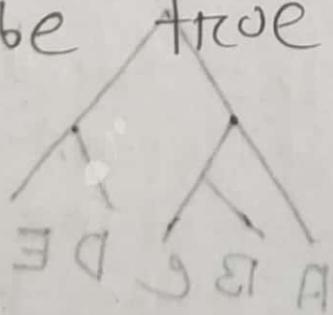
$$\frac{5!}{3!(5-3)!}$$

$$\frac{5!}{3!2!}$$

$$\frac{4 \times 5}{2} = 10$$



N.B: Triplet which come maximum time that will be true species (TST) tree

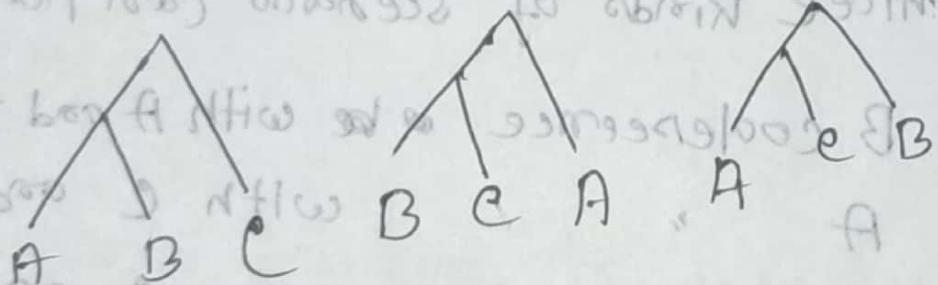


\* Triplets and quartets are statistically consistent estimators.

■ Triplets are statistically consistent estimators of true species tree.

Proof:

Let we have three taxa A, B, C. Indeed triplets of these taxa are



Probability of B's coalescence with A is  $P_{node}$

Probability of B's coalescence with A is  $1 - P_{node}$

$$E(P_{node}) = P_{node} + (1 - P_{node}) \cdot \frac{1}{2} = \frac{1}{2}$$

$$E(1 - P_{node}) = 1 - P_{node} - (1 - P_{node}) \cdot \frac{1}{2} = \frac{1}{2}$$

$$E(1 - P_{node}) = (1 - P_{node}) \cdot \frac{1}{2} + P_{node} \cdot \frac{1}{2} = \frac{1}{2}$$

Blloositeitote gwo ab~~form~~ baa ab~~gint~~ +  
retorites frakiaro)

For second case of B.

If B coalescence with A at node  $\times$  then  
species gene tree is (ABC)

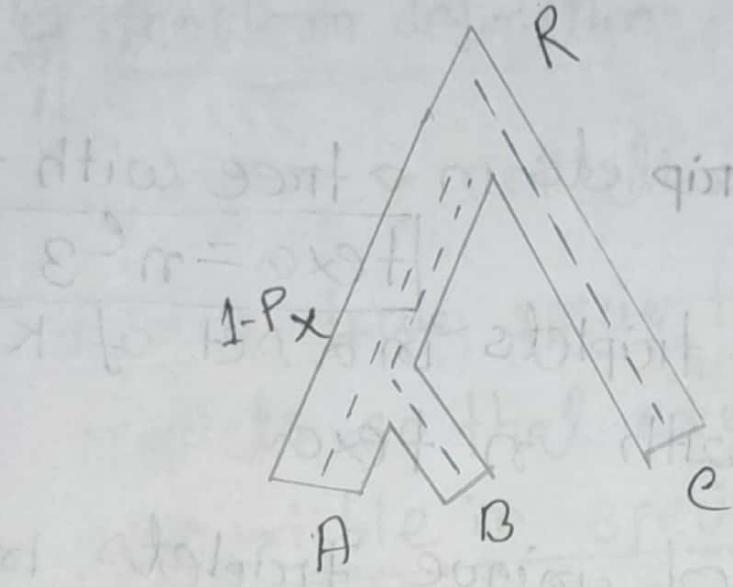
If B doesn't coalescence with A  
at node  $\times$  then if coalescence  
three kinds of scenario can happen.

If B coalescence with A and then C (ABC)  
A , A G with C and then B (ACB)  
with C and then A (BCA)  
gene trees are (CBA) ; (ACB)  
and (CBA) respectively and their  
coalescence probability is  $(1-P)/3$

$$\therefore \text{probability of } (ABC) = P + (1-P)/3$$

$$\therefore \text{of } (ACB) = (1-P)/3$$

$$\therefore \text{of } (BCA) = (1-P)/3$$



4

$$\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8}$$

$$\frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8}$$

- \* Number of triplets in a tree with  $n$  nodes =  $n^3$
- \* Number of triplets in a set of  $K$  trees with  $n$  nodes
- \* Number of unique triplets in a set of  $K$  trees with  $n$  nodes
- \* Number of unrooted binary tree on  $n$  nodes =  $(2n-5)!!$
- \* Number of rooted binary tree on  $n$  nodes =  $(2n-3)!!$

## Problem definition:

Given a collection of gene tree triplets  
find a true species tree by joining these  
triplets so that every gene trees  
compatible in species tree.

## Algorithm:

- If the number of taxa in  $\mathcal{X}$  is less than 3, just return the tree in  $\mathcal{X}$ .
- Else,
  - Find a sibling pair  $a, b$  (which is always possible)
  - If no such pair exists, return "No tree".
  - Else, continue
  - Remove all rooted triples that induced include  $a$  or  $b$  from the set  $\mathcal{X}$

→ Recursively compute a tree on the

Reduced set  $X'$  of rooted triple

→ Invert  $a$  into the tree by making  
it subtling to  $b$ .

→ Construct a tree on the leaf

\* set  $\{a, b, c, d, e\}$  which is inconsistent

with each of the following triples?

You have to show the intermediate  
steps of your algorithm.

able able a b

Sol 9

Step 1:

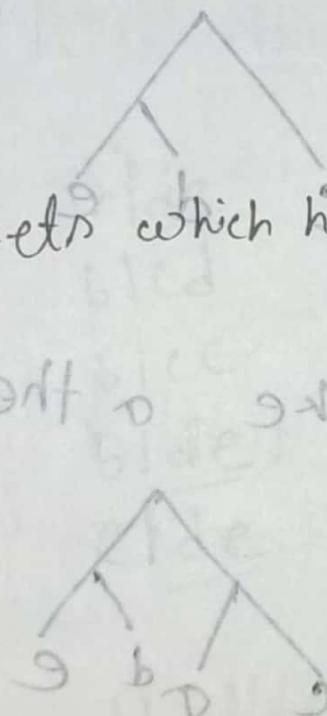
pair of siblings which are always grouped together are, bc, de.

choose bc

Step 2

Remove triplets which has b and remaining triplets are

olde  
aeld  
aele  
elde



Step-3

pair of siblings which are always grouped together are ac, de and choose ae

Step-4

Remove triplets with a and

108

remaining triplets are cde

about two sets are going to pair  
Step 5

Step 6

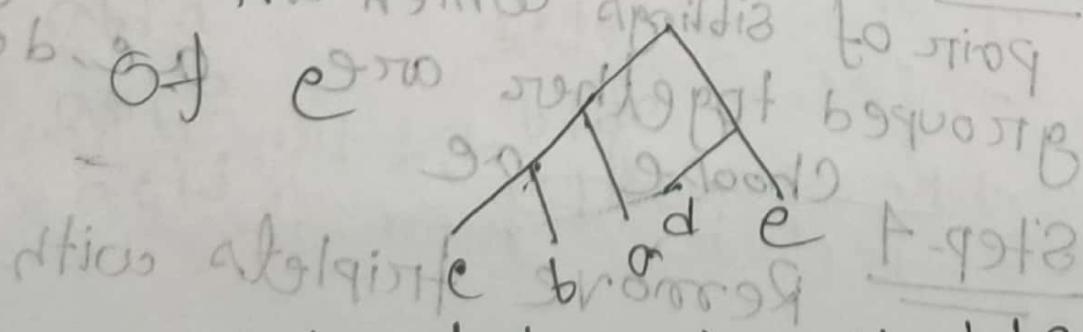
Step 6

Now make a the sister of e.

```
graph TD; c --- a; c --- d; a --- e;
```

Step 7

Now make b the sister of c



∴ final tree is (c, (b, a), (d, e))

Given a tree on leaf tree labeled  
which is consistent with each of the  
following 6 triplets. You have to  
show the intermediate steps of your  
algorithm.

albe	alde	ei work pot
albd	bled	
albe	bice	ppk
aled	blde	
alce	clde	

Sol:

Step 1: Sibling that is always grouped  
together is de.

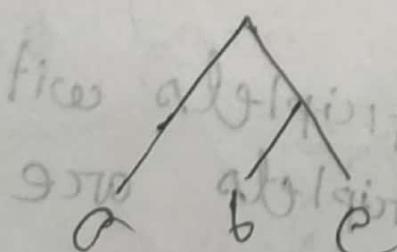
Step 2: Remove triplets with d and  
remaining 6 triplets are

Step 3: Sibling with that is always grouped together is ce

#### Step-4

Remove triplets with ce and remaining triplets are bd

#### Step-5



Step-6

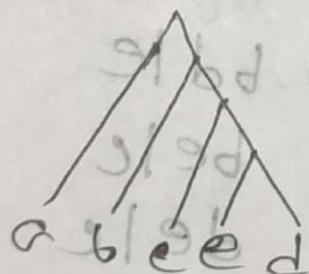
make c as the sister of e  
 about pricoll of ent to dos dices  
 rabi ent words of end up  
 mthioblo subg b e aqet2 strobzne

```

graph TD
    c[c] --> a[a]
    c --> b[b]
    a --> d[d]
    a --> e[e]
    b --> f[f]
    b --> g[g]
  
```

Step-7

make make bd as the sister of e



: 103

group word start with qibdi2 19512  
 .bd, go one by one  
 .go 69906 C

Construct a tree on the leaf set  
 $\{a, b, c, d, e\}$  which is consistent  
with each of the following triplets.

You have to show the intermediate steps of your algorithm.

actb  
falso  
bela  
dela  
acid

ac|e  
bdle  
belc  
deles

F-9913

Sol<sup>n</sup>:

Step 1: Siblings that are always grouped are ac, bd.  
choose ac.

Step-2

Remove triplets with a ~~b~~ and ~~c~~ remaining  
and ~~d~~ left.

bdic

bde

bele

dele.

Step-3

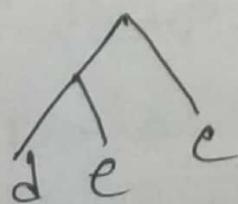
Sibling that is always grouped  
together is bd.

Re

Step-4 : Remove triplets with b and  
remaining are.

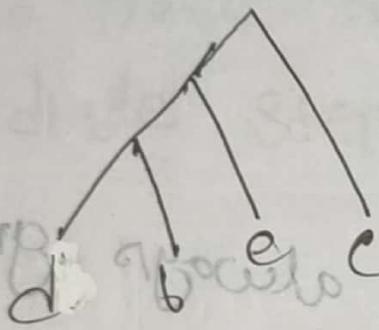
deic

Step-5



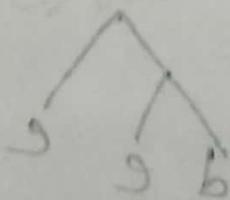
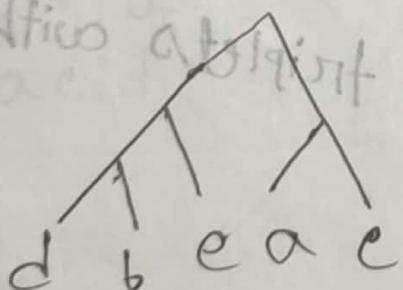
Step-6

make b as the sibling of / sister of d



Step-7

make a as the sister of c.



↗ group of organisms that share common

### Clades

$L(T)$  is called the leafset of  $T$ .  
 $T_v$  denotes the subtree of  $T$  below node  $v$ . (The clades of  $T$  are the subset of  $L(T)$  that are equal to the  $L(T_v)$  for some vertex  $v$ .)

### Cluster:

Set of taxa in a subtree

### types of clades:

trivial clade:

non-trivial clade

trivial clade: Let  $T$  be a rooted tree

on leafset  $S$ .

All the singleton clades and set  $S$  are trivial clades.

are trivial clades.

~~non-trivial clades~~

20b010

All the clades that are not

~~trivial~~

## Properties of clade to be compatible



→ containment:

→ disjoint

20f010

→ One clade's leaf set fully contain  
the leaf set of another clade

→ Clade that doesn't contain only  
singleton clade

20f010

20f010

20f010

Q trivial clades and non-trivial clades of tree

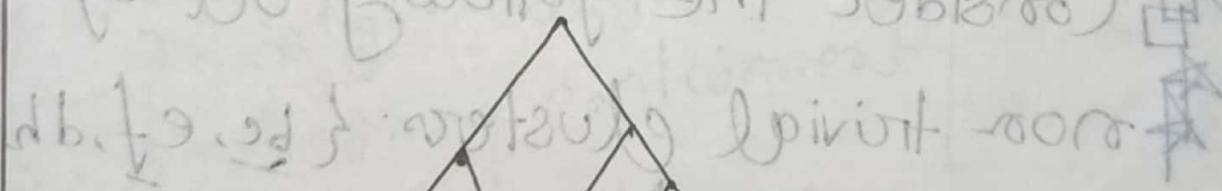
ans state  $C((a,b)(c,d,e)))$ .

trivial clades are  $\{a\}, \{b\}, \{c\}, \{d\}, \{e\}$ ,

$\{a, b, c, d, e\}$

non-trivial clades are  $\{a, b\}, \{c, d, e\}, \{d, e\}$

$\therefore \text{clades } (T) = \{\{a\}, \{b\}, \{c\}, \{d\}, \{e\},$   
 $\{a, b, c, d, e\}, \{a, b\}, \{c, d, e\}, \{c, d, e\}\}$



## E Compatibility:

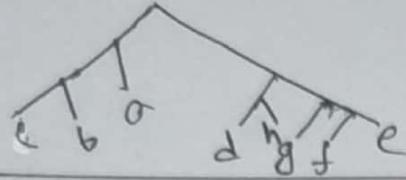
Two clades/clusters are compatible if they can "co-exist" in a binary rooted tree.

A clade can't overlap another

Clade to become compatible

Consider the following set of non-trivial clusters: {bc, ef, dh, abc, efg, defgh}. Is this a compatible set of clusters?

If yes show the corresponding tree. If not explain why.



Sol<sup>n</sup>

Now find out the relation between every possible pair of non-trivial cluster

to be, ef → disjoint

be, dh → "

be, abe → containment

be, efg → disjoint

ef, dh → disjoint

ef, abe → "containment

ef, efg → containment

ef, defgh → "

be, defgh → "

dh, abe → disjoint

dh, efg → disjoint

dh, defgh → containment

\* dh, defgh → containment

abe, efg → disjoint

abe, defgh →

efg, defgh →

containment

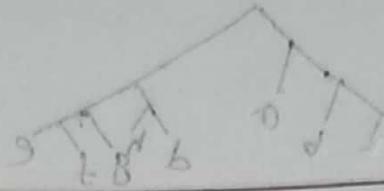
So there is no overlapping relation between any pair of non-trivial

cluster.

So given set is compatible.

~~Appendix~~

2020



Q10  
Consider the following set of non-trivial clusters  $\{be, fg, abc, defg, edfg, bedefg\}$ . Is this a compatible set of clusters? If yes, show the corresponding tree. If not why.  
Sol.

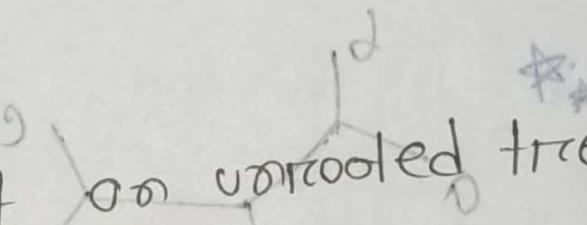
Now find out the "relation between every possible pair of non-trivial clusters".

$be, fg \rightarrow$ disjoint	$fg, abc \rightarrow$ disjoint
$be, abc -$ containment	$fg, defg \rightarrow$ contains
$be, defg -$ disjoint	$fg, edfg \rightarrow ..$
$be, bedefg \rightarrow ..$	$fg, bedefg \rightarrow ..$
$bc, bedefg \rightarrow$ contains	

$abc, defg \rightarrow$ disjoint	$defg, abc \rightarrow$ contains
$abc, edfg \rightarrow$ "	$defg, bedefg \rightarrow$ "
$abc, bedefg \rightarrow$ overlap	
	$edfg, fedefg \rightarrow$ containment

As there is a overlapping relationship  
so given set is not compatible.

### Bipartition:

Bipartition of  on unrooted tree  
formed by taking each edge (in turn)  
and writing down the two sets of  
leaves that would be formed by deleting  
the edge.

$\Pi(e)$  used to denote bipartition by

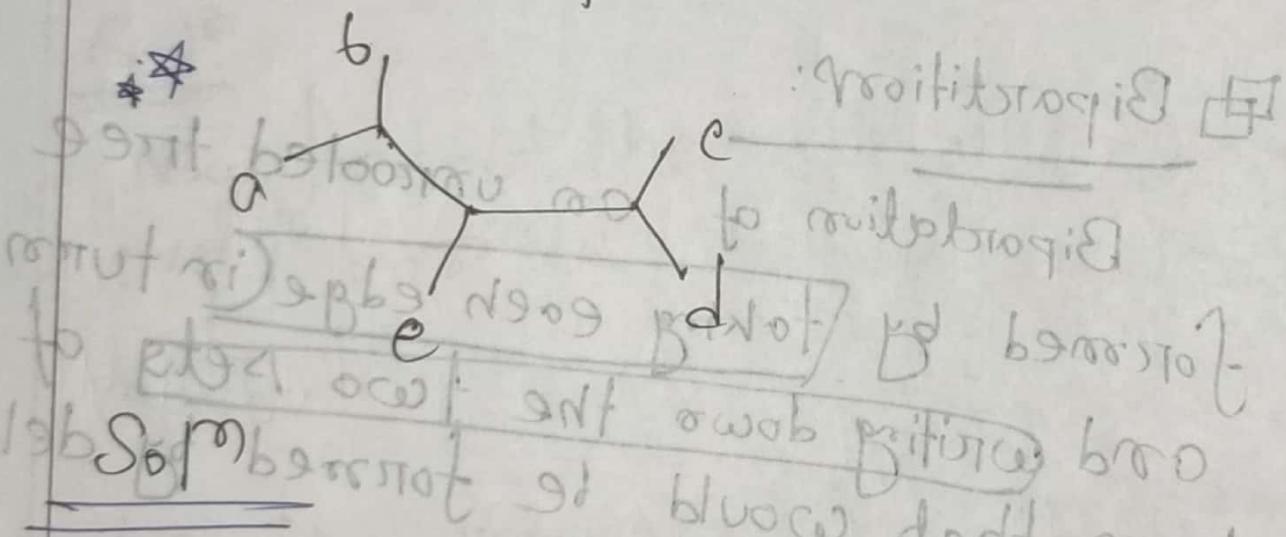
edge  $e$  with  $\Pi(e) = A \sqcup B$

where  $A$  is one half of the bipartition  
and  $B$  is other

two types of bipartitions:

1. trivial bipartition
2. non-trivial bipartition

Find out the trivial and non-trivial bipartitions of tree



non-trivial bipartition:

$$\{a, b\} \mid \{c, d, e\}$$

$$\{a, b\} \mid \{c, d, e\}$$

$$\{a, b, e\} \mid \{c, d\}$$

Trivial bipartitions:

$$\{\emptyset\} \mid \{b, c, d, e\}$$

$$\{b\} \mid \{\emptyset, c, d, e\}$$

$$\{c\} \mid \{a, b, d, e\}$$

$$\{d\} \mid \{a, b, c, e\}$$

$$\{e\} \mid \{a, b, c, d\}$$

Bipartition

not both or both a set

coincides with elements and pairs

either pairs or sets of two  
elements must form a pair

single

## Anomalous gene tree (AGT)

~~most probable gene trees with a topology~~  
~~Gene trees with a topology~~  
~~different from species tree that~~  
~~are more probable to observe~~  
~~than congruent gene tree.~~

NB: No anomaly zone in 3 species triplets.

## Bootstrapping:

It is a statistical method for calculating the sampling distribution of an estimator by sampling with replacement from original sample.

This is any test or metric that uses random sampling with replacement and falls under the broad broader class of resampling

NB : Branch length in a phylogenetic tree indicate

- \* confidence
- \* mutation rate
- \* time

longer isn't most  
stap

and longer branch = longer

tree = at best more

and) or worse. o wait does

biology this branch

biology a lot of stuff

Estimating branch support by using  
bootstrapping method

→ input sequence alignment is used to generate a large number boot strap replicate dataset

→ boot strap replicate is data matrix with the same dimensions of original matrix

→ column of boot strap replicate is obtained by sampling with replacement

→ from the original data

\* suppose a original input has 1 species

→ sequence length = 200 nt

\* each time a column is chosen randomly with replacement.

\* After 200 times replacement

## Calculating bootstrap support

A new input set is generated.

- \* 100 or more input data is generated
- \* A phylogeny is estimated on each bootstrap dataset and original dataset using same method

\* characterize the support for edge in the tree  $T$ .

- \* An edge  $e$  in  $T$  has bipartition  $A|A'$
- \* From bootstrap tree for each edge  $e$  of  $T$  determine the % of trees have some bipartition like original tree.

\* bootstrap support/posterior probabilities:

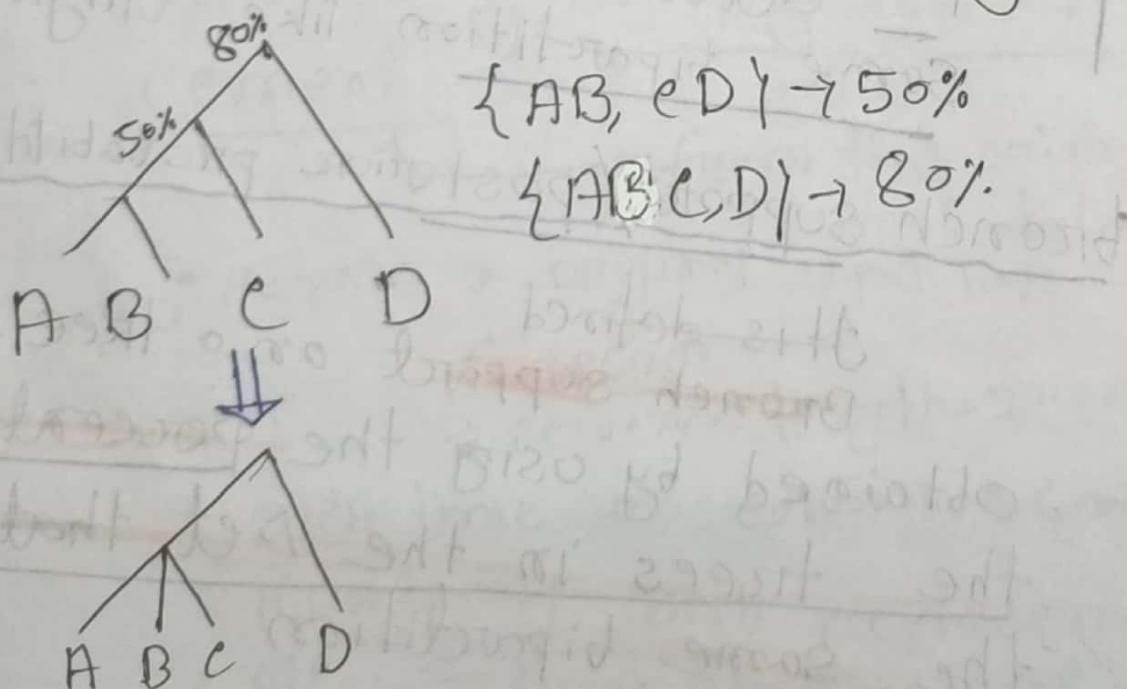
This defined  
~~Bootstrap support~~ on a tree is the percentage of the trees in the set that include the same bipartition

NB

High bootstrap support can be misleading when the estimation method is not statistically consistent.

- \* support value  $< 50\%$  : unreliable
- \* support value  $> 95\%$  : reliable
- \* support value ( $50\% - 95\%$ ): opinions differ as to the reliability of edges.

# It may generate non-binary tree



## Tree construction methods:

→ distance based method

\* UPGMA method

\* Neighbor Joining

→ parsimony based method

\* Fitch's algorithm

\* Sankoff's algorithm

→ maximum likelihood

between no

← CN is

## Distance based estimation

2 step processes

1. Computing distance matrix

2. Computing a tree from a distance matrix

two algorithms

1. UPGMA → produce a rooted tree

2. NJ → " an unrooted "

■ UPGMA (Unweighted pair group method with arithmetic mean)

→ finds a pair  $\{x, y\}$  that have the smallest distance and make them siblings

→ these two are then replaced by the cluster  $\{\{x, y\}, z\}$

→ the distance from  $\{x, y\}$  to every other taxon  $z$  is defined to be the average of  $d(x, z)$  and  $d(y, z)$

$$d(\{\{x, y\}, z\}) = \frac{d(x, z) + d(y, z)}{2}$$

→ the process is repeated until all the taxa are merged into a single cluster

## Requirement

★ ★

→ Ultrametric matrix

→ Additive

\* UPGMA fail when the distance does not obey strict molecular

clock (not ultrametric)

$d(\text{root} \rightarrow \text{leaf})$  are same

↳ Ultrametric matrix

An  $n \times n$  matrix  $M$  corresponding to distances between the leaves in a rooted edge-weighted tree  $T$  where the sum of the

edge weights in the path from the root to any leaf of  $T$  does not depend on the selected leaf.

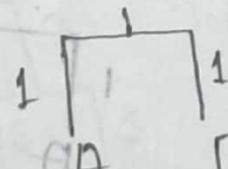
OR soon

By Constructing tree from following distance matrix

	A	B	C	D	AB
A	0	2	16	16	
B		0	16	16	
C			0	10	
D				0	

Step-1: As  $d(A, B)$  is smallest value  
 $d(A, B) = 2$

cluster  $\{A, B\}$



Step-2

Replace  $A, B$  with cluster  $AB$ .

$$\therefore d(A, B, C) = \frac{d(A, C) + d(B, C)}{2} = \frac{16 + 16}{2} = 16$$

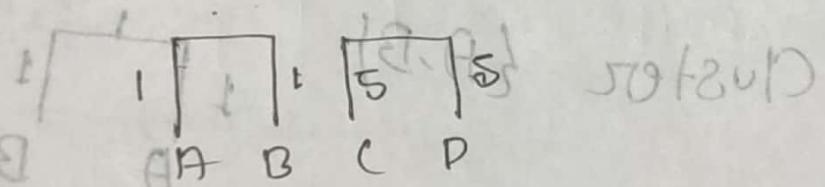
$$d(AB, D) = \frac{d(A, D) + d(B, D)}{2} = \frac{16 + 16}{2} = 16$$

Step 2 : Build a matrix of cost function

AB	C	D	Total	
AB	16	16	B	A
C	16	10	S	A
D	16	0	0	B
	61	0		S

Step 3 : Choose C,D and or d(C,D)

is smallest and make cluster  
 $S = \{C, D\}$



Step 4

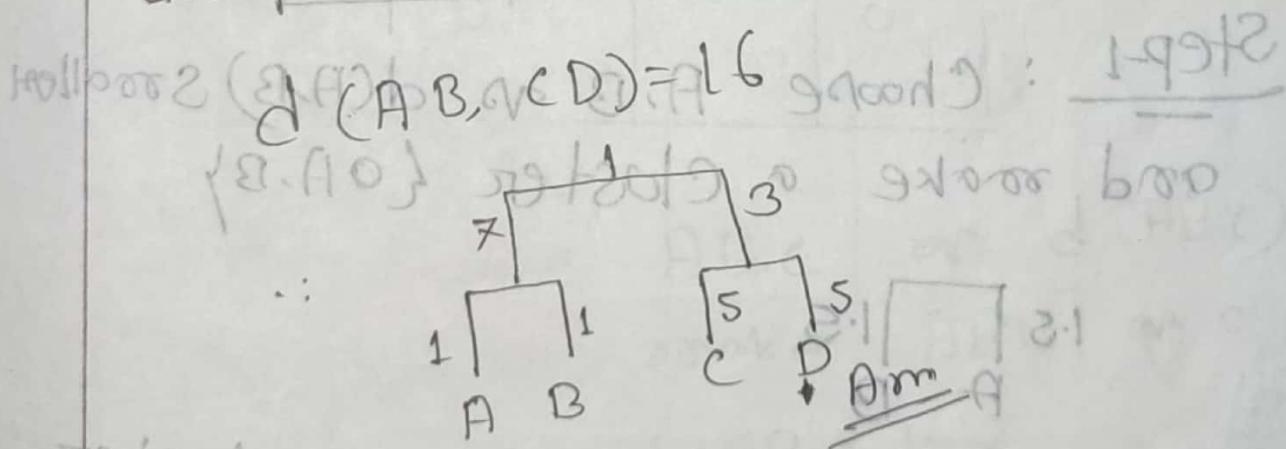
Replace {C,D} with {C,D}.

$$d = d_1 + d_2 : d(AB, CD) = \frac{d(A, C) + d(B, D)}{2}$$

$$= \frac{16 + 16}{2} = 16$$

AB	CD	16
CD	E D C B A	Y S O B A

Step-5 choose AB, CD and make a cluster ABCD as there is only one dis possible cluster.



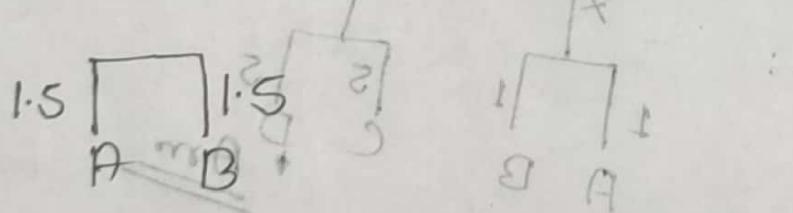
~~so k=10~~  $b = \frac{(2 \cdot B) + (2 \cdot A)}{5} = (2 \cdot BA) b$  Step-2

$$F = \frac{F_1 + F_2}{2}$$

■ Constructing tree from following  
distance matrix,

	A	B	C	D	E
A	0	3	7	10	10
B		0	7	10	10
C			0	10	10
D				0	8
E					0

Step 1 : Choose A, B as d(A,B) smaller  
and make a cluster {A,B}



Step 2 : Replace d(A,B) with cluster

$$d(A,B,C) = \frac{d(A,C) + d(B,C)}{2}$$

$$= \frac{7+7}{2} = 7$$

$$d(A, B, D) = \frac{d(A, D) + d(B, D)}{2} = \frac{10 + 10}{2}$$

$$d(A, B, E) = \frac{d(A, E) + d(B, E)}{2} = \frac{10 + 10}{2} = 10$$

	AB	CD	D	E
AB	0	7	10	10
C	0	10	10	10
D	5	0	8	
E	8		0	

Step 3 choose ABC or  $d(ABC)$  is smallest and make ABC as a cluster

$$d(AB, C) = 7$$



follows (E.D.B, 1.5) 2.5 3.5 2.5  
follows a parent node and we have

Step-4

Replace  $AB, C$  with  $ABe$

$$d(ABe, D) = \frac{d(AB, D) + d(Be, D)}{2}$$

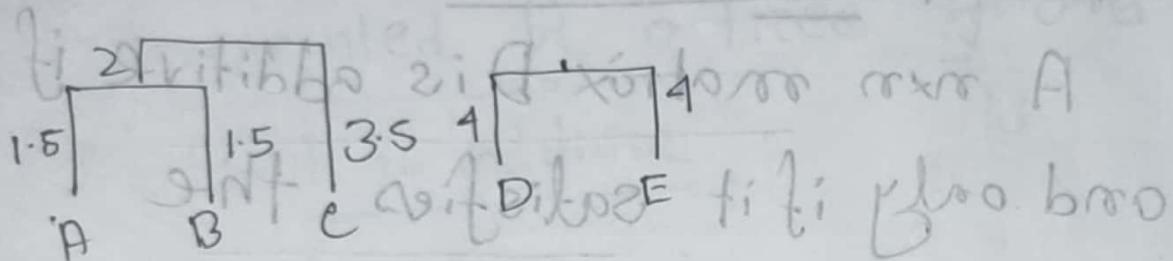
$OI =$

$$OI = \frac{10 + 10}{2} = 10$$

$$d(CABe, E) = \frac{d(CAB, E) + d(Be, E)}{2}$$

$ABe$	$D$	$E$
90	10	10
ABC	ABe	ABe
D	Be	Be

Step-5 choose  $D, E$  as  $d(D, E)$  smallest  
and make them as a cluster.

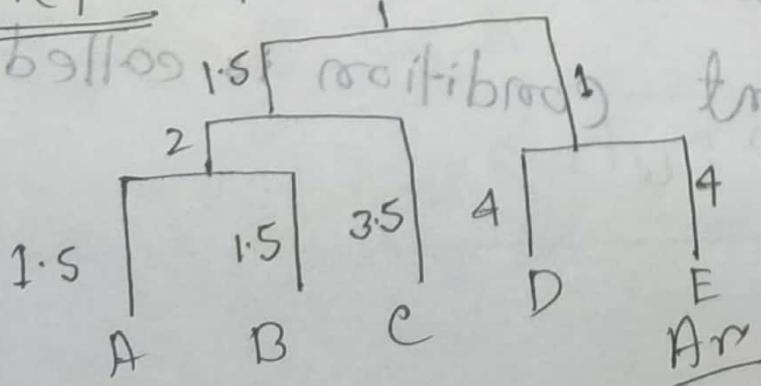


Step-6 Replace  $D, E$  with  $DE$

$$d(ABe, DE) = \frac{d(ABe, D) + d(ABe, E)}{2} = \frac{10 + 10}{2} = 10$$

	ABe	DE
ABe		10
DE		

Step-7 make  $ABe, DE$  as a cluster.



## Four point theorem

A  $n \times n$  matrix  $D$  is additive if and only if it satisfies the "four point condition" for all indices  $i, j, k, l$  which is that the median and largest of the following three values are the same.

$$D_{i,j} + D_{k,l}$$

$$D_{ik} + D_{jk}$$

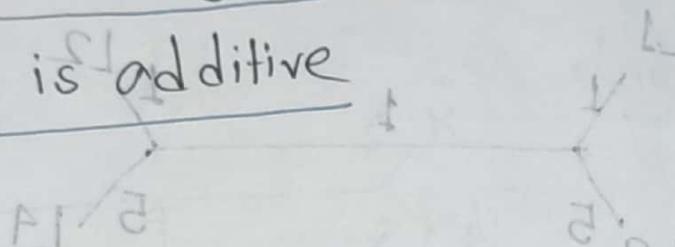
$$D_{i,k} + D_{j,k}.$$

	DE	ABG	ABC
10			
			DE

Additive: A matrix that satisfies the

four point condition is called additive

Theorem: A distance matrix can be represented by a tree if and only if it is additive.



Theorem:

Ultrametric: A dissimilarity map  $S$  on  $X$  is an ultrametric if for every three

distance  $i, j, k \in X$

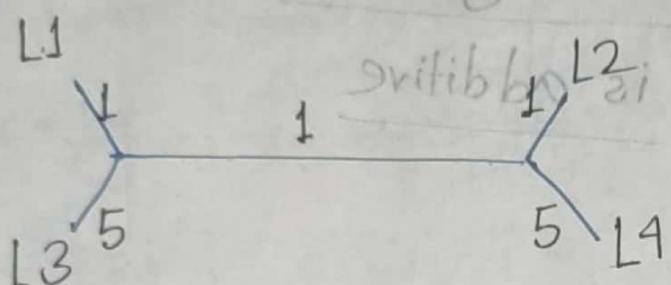
$$S(i,j) \leq \max\{S(i,k), S(j,k)\}$$

triangle inequality

An  $n \times n$  matrix  $d$  is said to satisfy the triangle inequality if for all  $i, j, k$ ,

$$d_{ik} \leq d_{ij} + d_{kj}$$

Example of additive matrix



$$\lambda = D_{(L_1, L_3)} + D_{(L_2, L_4)} = (1+5) + (1+5) = 6+6=12$$

$$\mu = D_{(L_1+L_2)} + D_{(L_3, L_4)} = (1+1+1) + (5+1+5) = 3+11=14$$

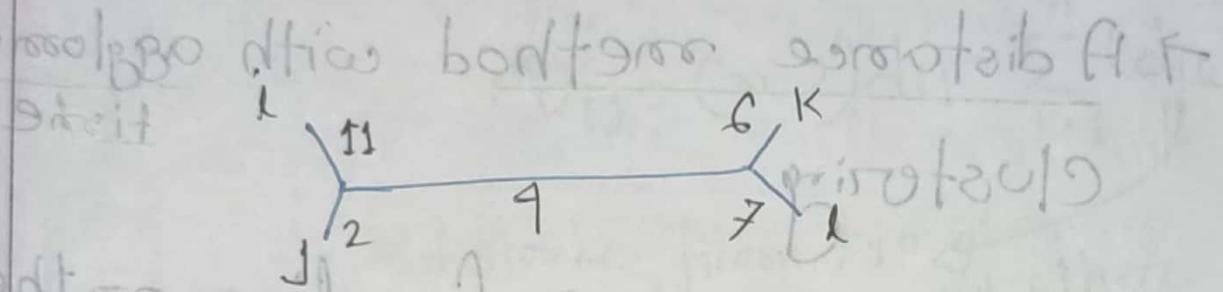
$$\nu = D_{(L_1, L_4)} + D_{(L_3, L_2)} = (1+1+5) + (5+1+1) = 7+7=14$$

As  $\lambda, \mu, \nu > \lambda$  so, it is an additive matrix.

$$ab + ab \geq ab$$

Determine if the following is additive or  
not

(CH) ~~Binded~~ ~~addition~~



$$u = D_{ij} + D_{kl} = (11+2) + (6+7) = 13 + 13 = 26$$

$$j = D_{ik} + D_{jl} = (11+9+6) + (2+7) = 21 + 13 = 34$$

$$z = D_{il} + D_{jk} = (11+9+7) + (2+6+4) = 22 + 12 = 34$$

$y = 2 \gamma u$

So it is an additive matrix

P. 158

most condition ~~isnt~~ ~~between~~ : ~~two~~

as . 1 ball/3 dol

## Neighbour Joining (NJ)

→ A distance method with agglomerative clustering

→ Compute the  $Q$  function on the given distance matrix

→ Clustering is similar to UPGMA, but computing the distance from a clustered node to the others is different from UPGMA

Input:  $n \times n$  dissimilarity matrix  $d$  with  $n > 1$

Output: Unrooted tree with  $n$  leaves labelled  $1 \dots n$

Initialization: Compute the  $n \times n$  matrix  $Q$ , defined by

$$Q_{ij} = (n-2)d_{ij} - \sum_{k=1}^n (d_{ik} + d_{jk})$$

total distance from  $i$  to  $j$  +   
 all other

" " " " (  $i$  to  $j$  )

$$D \rightarrow D^*(Q) \rightarrow D'$$

recursively apply neighbor joining on  $D'$

(Distance  $\rightarrow$  NJ matrix)

$$D_{ij}^* = (n-2) D_{ij} - \text{Total distance}_d(i) - \text{Total distance}_d(j)$$

$$\Delta_{ij} = \frac{\text{Total distance}_d(i) - \text{Total distance}_d(j)}{(n-2)}$$

$$\text{Link length } (i) = \frac{1}{2} (D_{ij} + \Delta_{ij})$$

$$\text{Link length } (j) = \frac{1}{2} (D_{ij} - \Delta_{ij})$$

Let cluster  $m = \{i, j\}$

$$D_{k,m} = \frac{1}{2} (D_{ki} + D_{kj} - D_{ij}) \quad (\text{NJ matrix} \rightarrow \text{Distance})$$

## Nearest Joining theorem

Given an additive matrix  $D$ ,  
the smallest element  $D_{ij}$  of its  
nearest-joining matrix  $D^*$  corresponds  
to a pair of neighboring leaves  $i$  and  
 $j$  in  $\text{Tree}(D)$ .

$$D \leftarrow D^* + D$$

$$(c_i^A + c_j^A) \frac{1}{2} = c_{ij}^A$$

$$\frac{(c_i^A + c_j^A) \frac{1}{2} - c_{ij}^A}{(5-1)} = c_i^A$$

$$(c_i^A + c_j^A) \frac{1}{2} - c_{ij}^A = (c_i^A + c_j^A) \frac{1}{2}$$

$$(c_i^A - c_{ij}^A) \frac{1}{2} = c_j^A$$

$$(c_i^A - c_{ij}^A) \frac{1}{2} = c_j^A$$

$$D_{kj} = \frac{1}{2} (D_{ii} + D_{jj}) \frac{1}{2} = \min D$$

Constructing tree from the following  
Distance matrix (Neighbour Joining)

	j	i	k	l	m	n	t	i
j	0	13	21	22			83	0
i	13	0	12	13			83	0
k	21	12	0	13			83	0
l	22	13	13	0			83	0

Step-1: Distance matrix to Neighbour-Joining matrix  
Calculate total distance

$$TD_i = 56 \quad i \{l, i\} = m$$

$$TD_j = 38$$

$$TD_k = 46$$

$$TD_l = 48$$

t	k	m	
51	01	0	m
31	0	01	k
01	11	11	1

Step 2

$$D \rightarrow D^*$$

$$n=4$$

$$\therefore (n-2) = 4-2 = 2$$

	$i$	$j$	$k$	$\lambda$	$\tau$
$i$	0	-68	-60	-60	
$j$	-68	0	-60	-60	
$k$	-60	-60	0	-68	
$\lambda$	-60	-60	-68	0	

$$D_{ij} = 2 \times D_{ij} - TD_i - TD_j$$

$$D_{ij} = 2 \times 13 - (56 + 38)$$

$$D_{ij} = 26 - 94$$

$$D_{ij} = -68$$

$$D_{ik} = 2 \times D_{ik} - TD_i - TD_k$$

$$D_{ik} = 2 \times 21 - 56 - 16$$

$$D_{ik} = 42 - 102$$

-68 is smallest value

$\therefore$  construct a cluster  $m$

$$m = \{i, j\}$$

$m = \{i, j\}$

Step 3  $D^* \rightarrow D'$

	$m$	$k$	$\lambda$	$\tau$
$m$	0	10	11	
$k$	10	0	13	
$\lambda$	11	13	0	

$$D_{mk} = \frac{1}{2} (21 + 12 - 13)$$

$$D_{mk} = \frac{1}{2} (21 - 1) = \frac{1}{2} 20 = 10$$

$$D_{m\lambda} = \frac{1}{2} (22 + 18 - 13) = \frac{1}{2} 27 = 13.5$$

Step-4

Calculate TD<sub>m</sub> and D'

$$TD_m = 21$$

$$TD_K = 23$$

$$TD_L = 29$$

Step-5

$D' \rightarrow D''$

m	K	L
m	0 -34 -34	

K -34	0 -34	
-------	-------	--

L -34 -34	0	
-----------	---	--

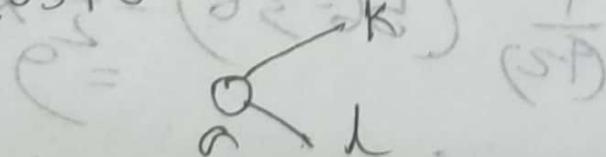
$$m=3 \\ m-2=(3-2)=1$$

$$D''_{mk} = 1 \times D'_{mk} - (TD_m + TD_K) \\ = 1 \times 10 - (21 + 23) \\ = 10 - 44 \\ = -34$$

$$D''_{ml} = 1 \times D'_{ml} - (TD_m + TD_L) \\ = 1 \times 11 - (21 + 29) \\ = 11 - 50 \\ = -39$$

choose (K, L)

8 Cluster (K, L) = {K, L}



$$= -34$$

$$D''_{kl} = (1 \times 13) - (23 + 29) \\ = 13 - 52 \\ = -39$$

$$(e+b) \frac{1}{3} = (l+d+c) \frac{1}{3} = (i,j,k) \frac{1}{3} \\ 55 \frac{1}{3} = (e+60) \frac{1}{3} = 62 \frac{1}{3} \\ 11 = 62 \frac{1}{3} - 55 \frac{1}{3} = 62 - 55 \\ = 7$$

Step 6

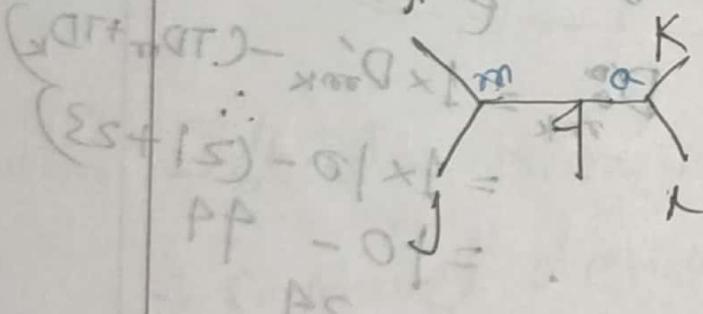
P-SPR

$$D' - D''$$

$$\begin{matrix} m & o & n \\ m & o & 4 \\ o & 4 & o \end{matrix}$$

$$\begin{aligned} IS &= CT \\ ES_{max} &= \frac{1}{n-2} (D_{m,n} + D_{max} - D_{KL}) \\ PS &= \frac{1}{2} (10 + 11 - 13) \\ &= \frac{1}{2} (10 - 2) \\ &= \frac{1}{2} \times 8 = 4 \end{aligned}$$

$$I = (f + g) = 5 - 2$$



$$\begin{aligned} I &= K - m \\ PS - PE &= O - m \\ PE - O &= PS - m \end{aligned}$$

Calculating link length

$$(PS + IS) - (11 \times 4) = \frac{1}{(n-2)} (TD_i - TD_j)$$

$$24 - 11 = \frac{1}{4-2} (56 - 38) = \frac{1}{2} \times 18 = 9$$

$$(PS + ES) - (13 \times 1) = \frac{1}{2} (D_{ij} + D_{ji}) = \frac{1}{2} (B + G)$$

$$13 - 13 = \frac{1}{2} (D_{ij} - 9) = \frac{1}{2} (13 - 9) = 2$$

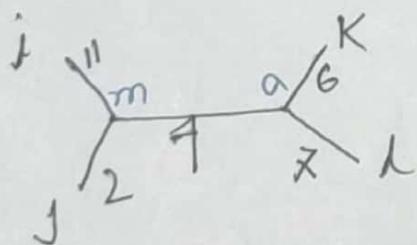
$$\therefore G = \frac{1}{2} (D_{ij} - 9) = \frac{1}{2} (13 - 9) = 2$$

$$\begin{aligned}\text{Length}(1K) &= \frac{1}{2} (D_{K1} + \Delta_{K1}) \\ &= \frac{1}{2} (13 - 1) \\ &= \frac{1}{2} 12 = 6\end{aligned}$$

$$\begin{aligned}\Delta_{K1} &= \frac{1}{(n-2)} (TD_K - TD_V) \\ &= \frac{1}{3-2} (23 - 21) \\ &= -1\end{aligned}$$

$$\begin{aligned}\text{Length}(1) &= \frac{1}{2} (D_{K1} + \Delta_{K1}) \\ &= \frac{1}{2} (13 - (-1)) \\ &= \frac{1}{2} 19 = 7\end{aligned}$$

$\therefore$  final tree



## Gene tree estimation

- \* given sequence
- \* apply multiple sequence alignment
- \* run Maximum likelihood method

NB: Discrepancy in gene tree estimation may cause gene tree discordance.

## Species tree estimation:

→ combined analysis

→ summary methods

## ~~Combined analysis~~

\* Combining

Supergene:

Combining all genes to  
create supergene

→ create supergene ( $\delta^*$ )

→ apply Bayesian based estimation  
(maximum likelihood,  
parsimony)

Dis.ad.

→ ignores gene tree discordance

Adv:

→ provide high accuracy.

## # Summary method:

- create different gene tree

multiple origins from different genes  
\* apply different technique to merge  
gene trees

- generate a species tree

merging technique

- Brin → maximum voting

→ maximum parsimony

## Disadvantages

- discordance

- anomaly zone

multiple trees from same species

↓ result with minimum

Gene tree parsimony / maximum parsimony

\* Find the tree that requires minimum number of changes to explain the data.

input: A set of rooted binary gene tree with each species having a single copy of genome.

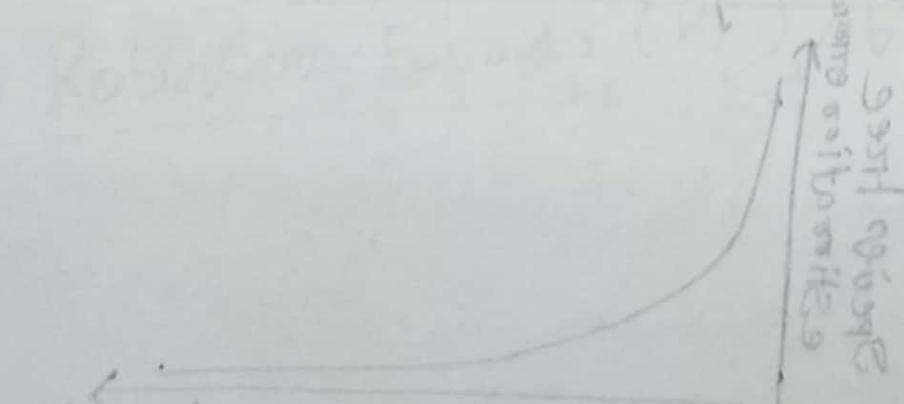
output: A species tree that optimizes a particular cost function.

- \* minimize cost w.r.t duplication
- \* maximize extra lineage

Potassium approaches usually have two different components

1. Search through many different tree topologies to find a good one

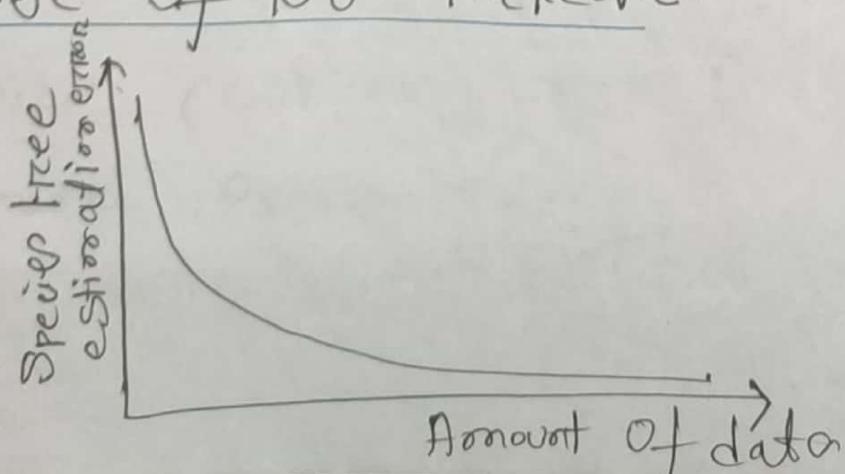
2. Find the minimum number of charges needed to explain the data for a given tree topology.



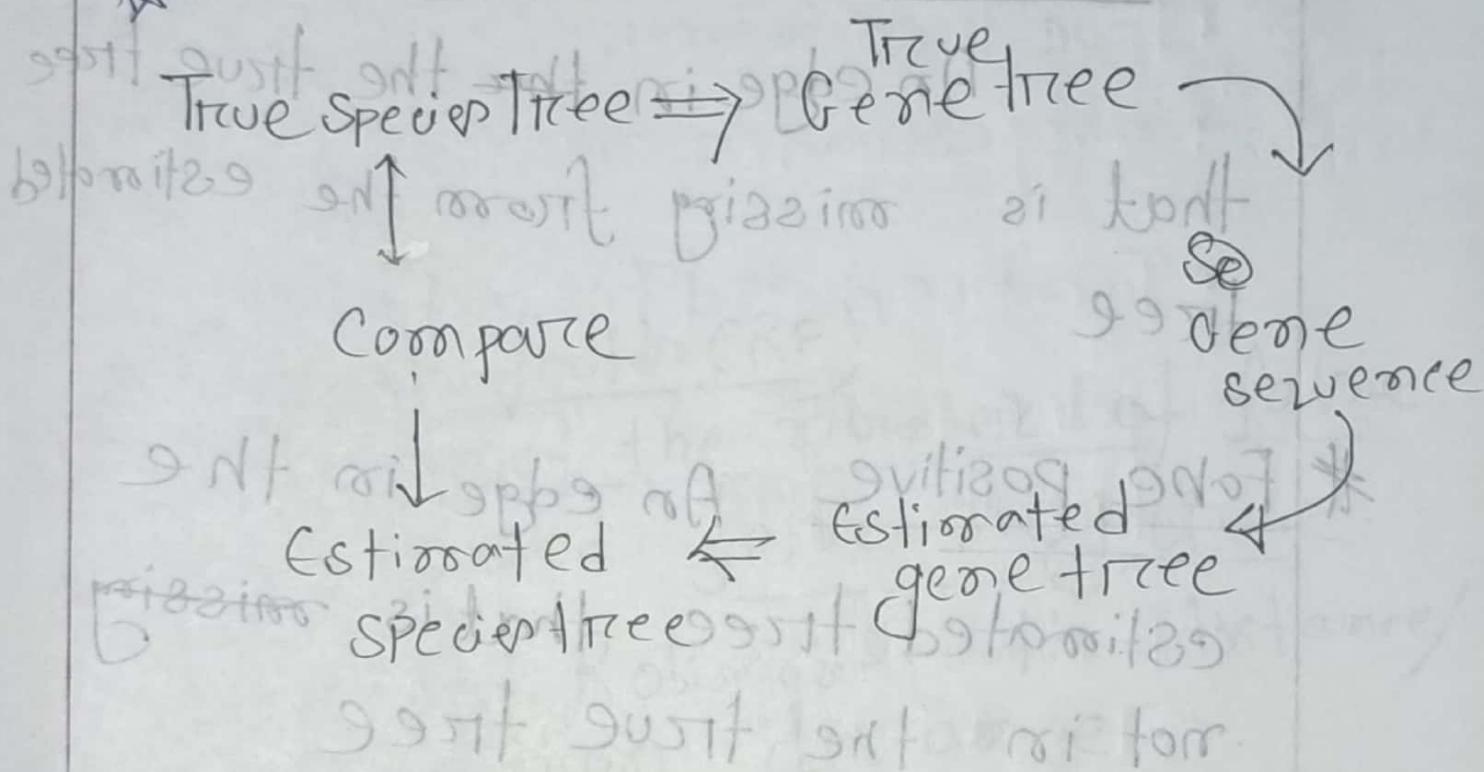
## ■ Statistical consistency: (Species tree estimation)

A species tree estimation method is called statistically consistent if the probability of returning true species tree converges to one as the amount of data increases.

- \* we usually assume both the number of sites/ loci and the number of loci increase.



## Evaluation procedures for SNPs



## Metrics

- False negative (FN) → false negative rate
- False positive (FP) → false positive rate
- Robinson-Foulds (RF) / bipartition distance

# False negative

An edge in the true tree  
that is missing from the estimated  
tree

# False positive An edge in the  
estimated tree that is missing  
not in the true tree

# False negative rate :

Number

Proportion of the missing  
edges + (99)

# FN edges

# internal edges in the  
true tree

# false positive rate:

# Robinson-Foulds (RF)

sum of the number of false  
positive and false negative.

Also called Robinson-Foulds distance/  
Bipartition distance.

$$RF = \frac{\# FP + \# FN}{2^n - 6}$$

number of leaves

of even new branch for all trees

for each node

RF

Q

stop svitlog slot #

- \* true tree is binary (assumed)
- \* many estimated tree are not binary
- \*  $T_0$  is binary and  $T$  is not binary  
false positive error rate will be less than false negative error rate
- \*  $T_0$  is not binary
- \*  $T_0$  and  $T$  both is not binary  
false negative error rate could be smaller than the false positive error rate.
- \* So Robinson-Foulds is not equal to the average of all negative and false positive
- \* if  $T$  is not binary we have to report both FP and FN instead of RF.

Nov 2020 / April 2020

Suppose you are trying to construct a species tree on 10 different species. You have sampled 250 genes from each of these 10 species. Your supervisor has asked you to use a method called GT-est for constructing trees from sequence alignments, and SP-est (which is a summary method) for estimating species trees from gene trees.

- i) How many times you need to run GT-est and SP-est to estimate a species tree by summarizing gene trees.
- ii) How many times do you need to run GT-est and SP-est to estimate a species tree by "combined analysis".

Ans

Ans:

i) Species tree by summarizing gene trees:

GT-est: 250 times

SP-est: 100 times

ii) Species tree by "Combined Analysis" of (combined

mitochondrial

and nuclear

GT-est: 1 time

SP-est: 0 time

of bears up to 2000 years old

stratigraphic levels of the 92 bear fossils

from 1000 to 2000 years old

Nov 2021 / April 2021

Suppose you are trying to construct a species tree on 15 different species.

You have sampled 300 gene sequences from each of these 15 species. Your supervisor has asked you to use a method called

GT-est for constructing trees from sequence alignment and SP-est (which is a summary method) for estimating species trees from gene trees.

- i) How many times you need to run GT-est and SP-est to estimate species tree by summarizing gene trees?
- ii) How many times you need to run GT-est and SP-est to estimate species tree if you use "Naive Binning" with 10 bins?

iii) How many times do you need to run GTree and SP-est to estimate species tree if you use (Statistical Binning) assuming that SB returns 12 bins

iv) How many times do you need to run GTree and SP-est to estimate a species tree if you use weighted Statistical Binning (WSB) assuming that WSB returns 12 bins

v) How many times do you need to run GTree and SP-est to estimate a species tree using combined analysis?

vi) How many times do you need to run GTree and SP-est to estimate a species tree using combined analysis?

Ans:

i) Species tree by summarizing genetree:

SP-GT-est: 300 times

SP-est: 1 time

v) Species tree by using combined analysis:

GT-est: 1 time

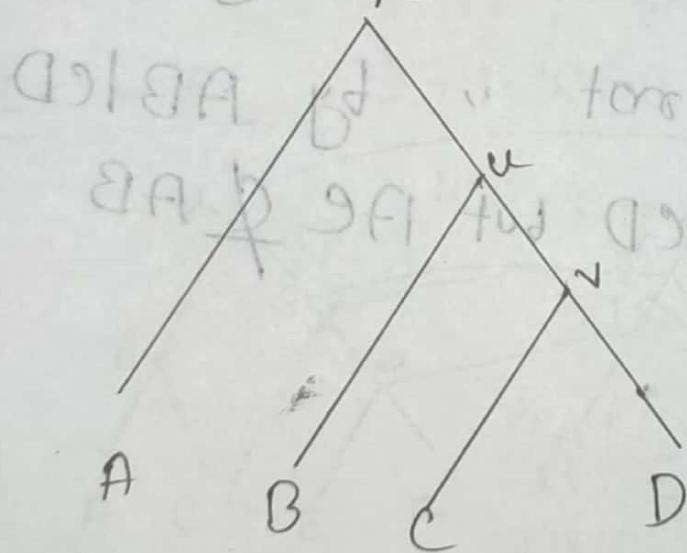
SP-est: 0 time

## Subtree-bipartition:

For an internal node  $u$  in a binary rooted tree  $T$ ,

$$SBP(u) = \text{cluster}(T_L) \sqcup \text{cluster}(T_R)$$

Example: Subtree bipartition of following tree



for node  $u$        $B \sqcup C \sqcup D$

,,    "       $v$        $C \sqcup D$

,,    "       $u$        $A \sqcup B \sqcup D$

## Domination

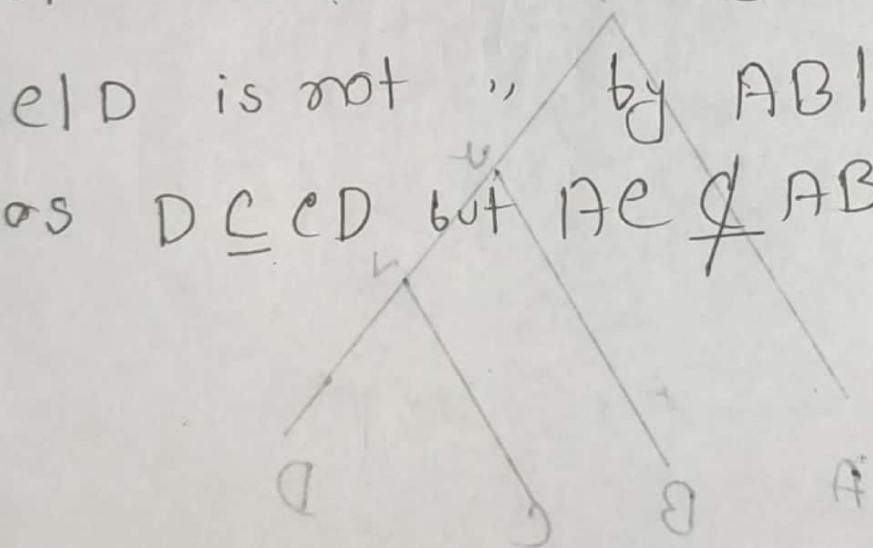
$X \sqsubset Y$  is dominated by  $P \sqsubset Q$

if  $X \subseteq P$  and  $Y \subseteq Q$



### Example

- \*  $A \sqsubset D$  is dominated by  $AB \sqsubset D$
- \*  $A \sqsubset D$  is not " by  $AB \sqsubset D$   
as  $D \subseteq CD$  but  $A \notin AB$

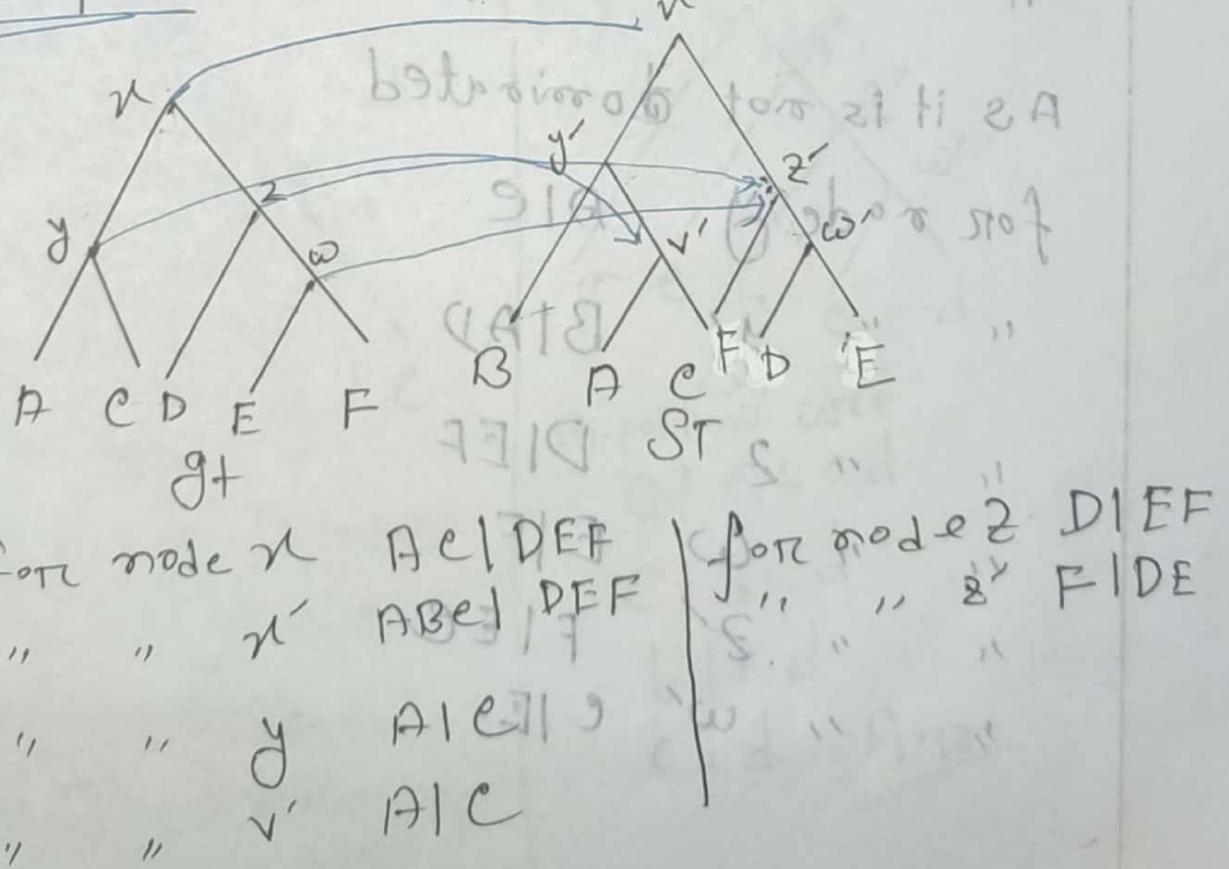


$D \sqsubset B$        $\sqsubset$  shows set  
 $D \sqsubset C$       " "  
 $B \sqsubset A$       "

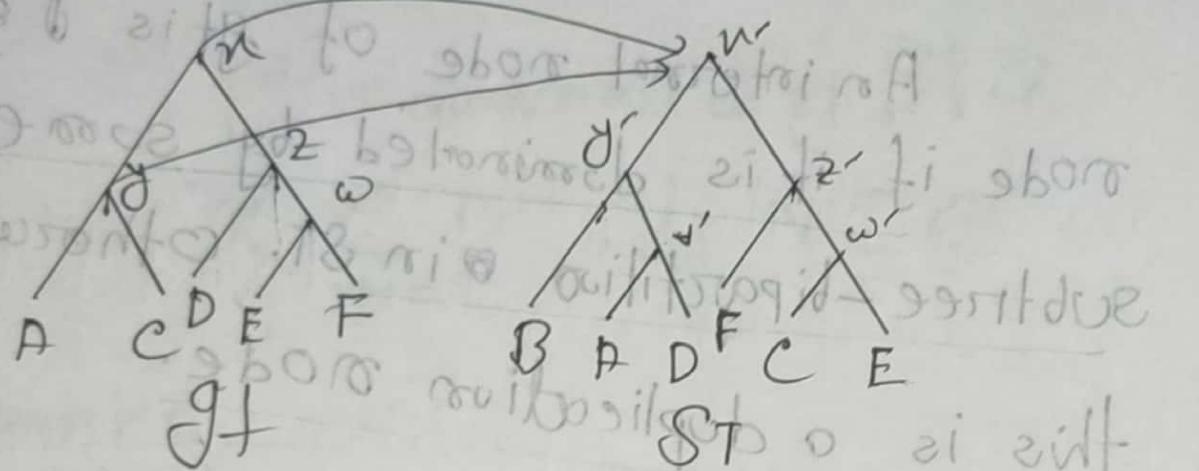
## Duplication (alternative defn)

An internal node of  $gt$  is a speciation node if it is dominated by some subtree-bipartition in ST. Otherwise, this is a duplication node.

### Example



#



for node  $n$  Ael DEF

" "  $n'$  BADI EFC

As it is not dominated

for node  $y$  AIE

" "  $y'$  BIAD

" "  $z$  DIEF

E/F/A

F/EC/A

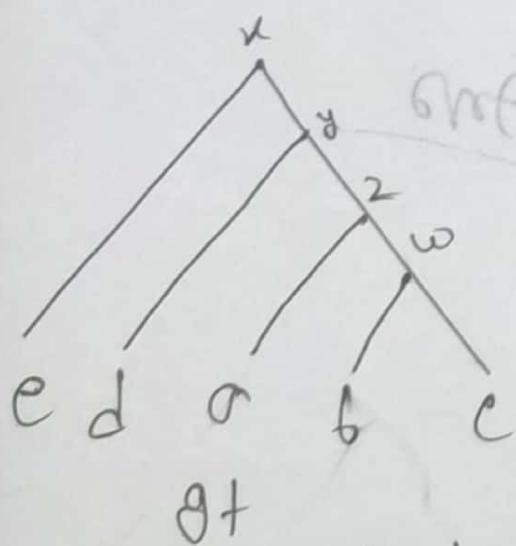
C/IE/A

April  
6/1

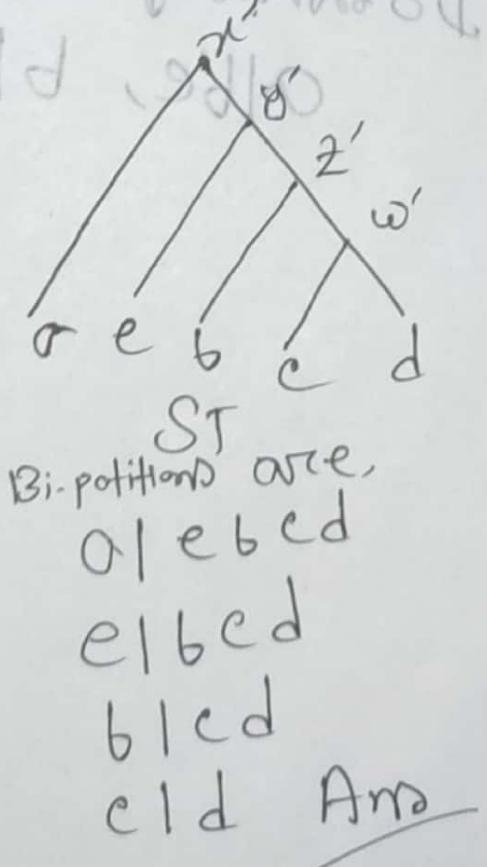
Consider the following  $G_t$  and  $ST$ . Find the sets of subtree-bipartitions in  $G_t$  and  $ST$ . Find the subtree-bipartitions in  $G_t$  that are dominated by some subtree bipartitions in  $ST$ .

$$G_t = (e, (d, (a, (b, c))))$$

$$ST = (a, (e, (b, (c, d))))$$



Bipartitions are:  
el|d|ab|e  
dl|ab|e  
al|b|e  
bl|e  
Ans



Bipartitions are:  
a|eb|cd  
el|b|cd  
b|l|cd  
cl|d  
Ans

for  $a \leq b$  and  $a \leq c$

$+B$  is  $a \leq c$ ,  $b \leq c$  to  $+B$

$\therefore a \leq c$  is dominant

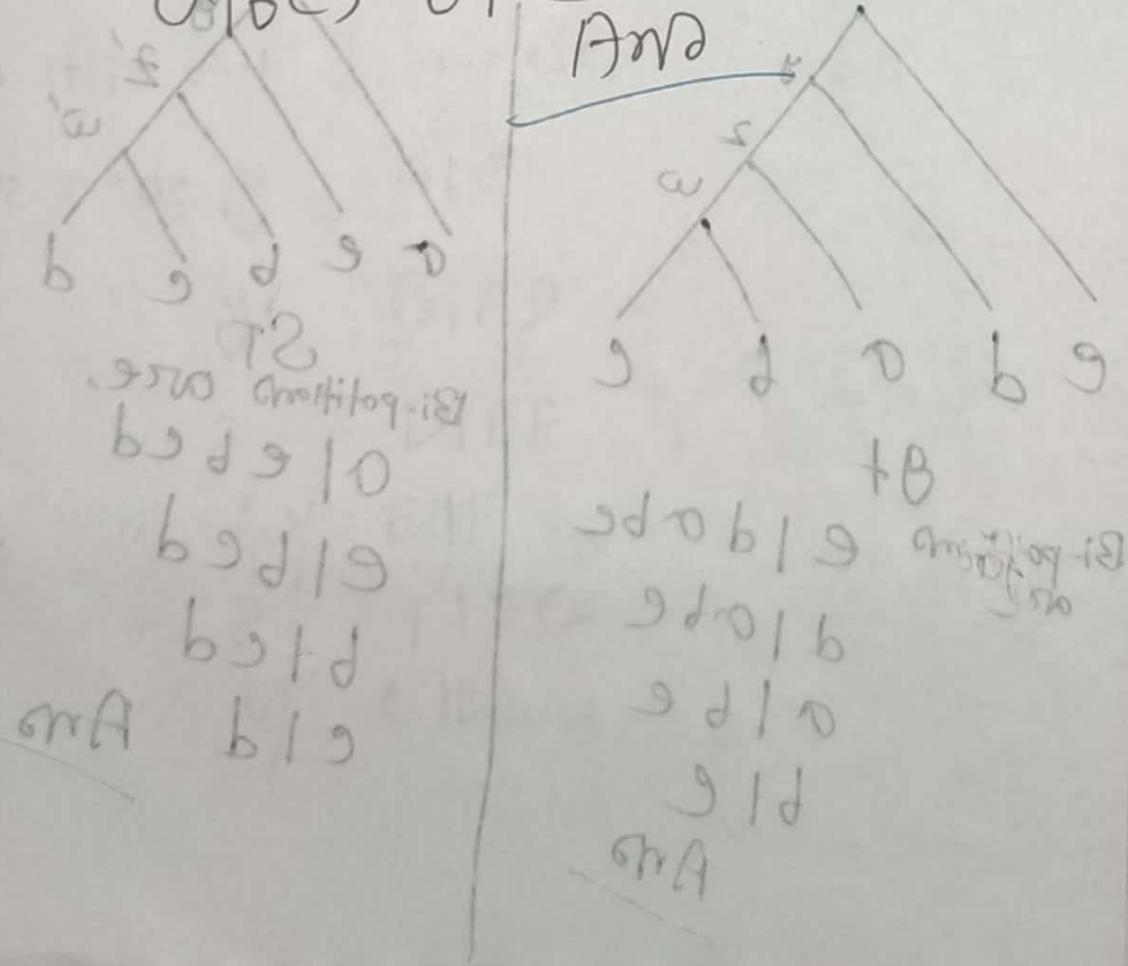
for  $b \leq c$  and  $b \leq d$

$b \leq c \leq d$

$\therefore b \leq c$  is dominant

$\therefore$  Dominant of sub-tree bipartition are

$a \leq b$ ,  $b \leq c$



## Multiple Sequence Alignment (MSA)

### MSA!

\* It reflects evolutionary processes operating on sequences so that the homology can be inferred.

Homology: Two letters in two sequences are homologous if they descend from a letter in a common ancestor.

NB

- \* MSA → estimating tree
- \* errors in MSA estimation tend to produce errors in estimated tree
- \* MSA introduces "—" in the sequences sometimes treated as an additional state

5 states for nucleotide alignment  
(A, G, C, T, N)

\* 21 " " protein "

- Column (sites) within the alignment

among homologous sequences define the homology

benefit of using homologous

### Evolutionary Processes operating on MSA:

- Substitution (mutation)
- Insertion and deletion (Indels)
- Rearrangement
  - inversions
  - transpositions

→ Duplication of regions

Note: Most alignment methods

stretch out sequences so that

they line up well and also only

addresses substitutions and indels

## ★ Inversion and decomposition

### Inversion:

$$S_1 = \text{ACC} \xrightarrow{\text{obri}} \text{AGTCA} \xleftarrow{\text{spodho}} \text{CCTTA}$$

$$S_2 = \text{ACC} \xleftarrow{\text{KTF}} \text{ACTGA} \xrightarrow{\text{spodho}} \text{CCTTA}$$

Sequence segments are same but they change their direction between two sequences.

### Decomposition:

$$S_1 = \text{ACC} \xrightarrow{\text{obri}} \text{AGTCA} \xleftarrow{\text{spodho}} \text{CCTTA}$$

$$S_2 = \text{AC} \xleftarrow{\text{spodho}} \text{CC} \xrightarrow{\text{obri}} \text{AGTCA}$$

Two segments of sequence swap their location between two sequences.

## Example of recombination

Turning Cabbage into Turnip

→ Share a common ancestry  
but they look and taste different

→ 99% similarity between genes

→ Differ in gene order

~~FITGAGCTGAAGA = 12~~

~~ATGAAATGGGA = 12~~

just swap → allows out

new and mutated variants

~~Example~~

Evolving one sequence into another

Sequence through a combination of insertion, deletion and substitution:

~~Deletion~~ → deletion      ~~Substitution~~ → substitution  
...ACGGTGCA~~G~~T~~A~~CCA...

...ACCA~~G~~T~~C~~AC~~T~~A...  
  insertion

Allignment/Evolving

$s_1 = \text{ACGGTGCA}G\text{TACCA}...$

$s_2 = \text{AC}---\text{CAGTCACCTA}...$

Space is reserved to add a long sequence to both ends.

~~Ex ample~~

Evolving sequence from Birla

$$S_1 = ACAT$$

$S_2 = AGAT$

~~so it fulfills~~ ~~so it is broken~~ → deleting C  
First pairwise offligament by inserting G

AC - AT  
A-GAT

Second pointwise alignment

A - e AT

$\text{ATGCAATTAGCTA} = \text{P}$

HG - AT

- ATGGAATCTAG --- GA = 2

There is more than one representation as a pairwise alignment, each of which accurately reflects true homology.

→ Conserved regions don't necessarily mean that mutation didn't happen.

## Importance of MSA



→ Shared ancestry (Homology):

Comparative Genomics

→ Prerequisite for sequence based phylogeny

→ Gene finding

→ Genome assembly

→ Protein attribute prediction

→ Conserved regions, motifs can be identified

Highly conserved DNA sequences  
are thought to have functional

Lethal mutation: A type of mutation

in which the effects can result in death or  
reduce significantly the expected longevity  
of an organism carrying mutation.

mutation occurs  
but it  
causes  
death

~~more fibroblasts~~ Fibroblasts migrate towards growth  
support fibroblast infiltration test

## Types of mutation: ~~in~~ to ~~most~~ most

→ Good

→ Bad

→ Silent

Good: Mutation that enhances the organism's function

Bad: Mutation that cause harmful traits (Huntington disease)

Silent: Mutation that cause no difference in the function of the organism

NB: Lethal dominant genes are rare because the organism dies.

types of sequence alignment: implemented by

DP

→ Local alignment → Smith - waterman

→ Global " → Needleman - Wuneh

Local alignment:

match a query with a substring  
(portion) of your subject

→ finds regions of similarity in parts of the region

→ searching for local similarities  
within large sequence

Global alignment

End-to-end alignment of input sequences

→ Find the best alignment across the whole sequences

→ Comparing two gene with similar functions

Two sequences  $s$  and  $t$

## Needleman-Wunck algorithm (Global alignment)

For two sequences  $s$  and  $t$

$$v_{i,j} = \max \left\{ \begin{array}{l} v_{i-1,j} + S(s_i, -) \\ v_{i,j-1} + S(-, t_j) \\ v_{i-1,j-1} + S(s_i, t_j) \end{array} \right.$$

$$\left. \begin{array}{l} (-, -) \\ (-, s_i) \\ (s_i, -) \end{array} \right\} + 1$$

Alphabet:

DNA: {A, T, C, G}

Scoring scheme:

$$S(s_i, t_j) = 1 \rightarrow \text{match}$$

$$S(s_i, -) = -2 \rightarrow \text{indel (insertion/deletion)}$$

$$S(-, t_j) = -1 \rightarrow \text{gap}$$

$$S(s_i, t_j) = -1 \rightarrow \text{mismatch}$$

S and t are input sequences

Needleman-Wunsh CS) pseudo code:

$m \leftarrow \text{length}(s)$

$n \leftarrow \text{length}(t)$

$V[0,0] \leftarrow 0$

for i in 1:m

$$V[i, 0] = V[i-1, 0] + \delta(s_i, -)$$

for j in 1:n

$$V[0, j] = V[0, j-1] + \delta(s_0, t_j)$$

for i in 1:m

for j in 1:n

$$V[i, j] = \max \begin{cases} V[i-1, j-1] + \delta(s_i, t_j) \\ V[i-1, j] + \delta(s_i, -) \\ V[i, j-1] + \delta(s_0, t_j) \end{cases}$$

elsewise

$$V[i, j] = \min \begin{cases} V[i-1, j-1] + \delta(s_i, t_j) \\ V[i-1, j] + \delta(s_i, -) \\ V[i, j-1] + \delta(s_0, t_j) \end{cases}$$

Q Align and find the aligned sequences of following two sequences S and t with the help of Needham - Wunsch algorithm.

$$S = \text{A A A C}$$

$$t = \text{A G C}$$

$$S(n) = 1$$

$$S(n-1) = -2$$

$$S(n-2) = -1$$

	A	A	A	G	C
	0	-2	-4	-6	
A	-2	1	-1	-3	
A	-4	-1	0	-2	
A	-6	-3	-2	-1	
C	-8	-5	-4	0	

to : Alignment score = -1 first two gap A  $\square$   
 deletion + loss 2 reasons of cost of insertion  
Higher costs Aligned sequence

Aligned Sequence with respect to given

First one

A G C  
 S-C  
 A A A C  
 I-B-1-2-3-1  
 2-3=-1

GAGA = 2  
 GGA = 1

Second one

A - G C

A A A C

I - 2 - 1 1 2 - 3 = -1

Third one

S- A G C

A A A C

I - 2 1 - 1 1 2 - 3 = -1

S- A G C

A A A C

I - 2 1 - 1 1 2 - 3 = -1

S- A G C

A A A C

I - 2 1 - 1 1 2 - 3 = -1

S- A G C

A A A C

I - 2 1 - 1 1 2 - 3 = -1

## Local alignment:

A min global alignment between substrings of the original sequences

Goal: Find the best local alignment between two strings

Input: String  $s$  and  $t$ , as well as a scoring matrix  $S$

Output: Alignment of substrings of  $s$  and  $t$  whose alignment score is maximum among all possible alignments of all possible substrings of  $s$  and  $t$ .

\* DP to solve the local alignment problem

\* Smith Waterman algorithm

$$\begin{aligned} & \text{(L.i.v)} S + H.i.v \\ & \text{(L.i.v)} S + H.v.i \\ & \text{(L.i.v)} S + H.v.v \end{aligned}$$

Naive approach:

Find the path with maximum score between every pair of matched vertices  $(i, j)$  and  $(i', j')$ .  
Select the pair with maximum score.

Smith-Waterman

+ find a path with the maximum score from  $(0, 0)$  to every other vertex.  
+ find odd edges of weight 0 in the edit graph.

Recurrence: For two sequences  $s$  and  $t$

$$v_{i,j} = \max \left\{ \begin{array}{l} v_{i-1,j} + \delta(s_i, -) \\ v_{i,j-1} + \delta(-, t_j) \\ v_{i-1,j-1} + \delta(s_i, t_j) \end{array} \right\}$$

~~Ques~~ Find the best local alignment of the following two sequences using a gap penalty -2

b) ~~Ques~~ with mismatch penalty -1 and match score = +1.

$$S_1 = CCG$$

$$S_2 = ATCCG$$

$(7 \times 5)$  matrix

$$S(s_i, t_j) = S_{mn}$$

$$S(n, j) = +1$$

$$S(n, j) = -1$$

$$S(-, m) = S_{mn} = -2$$

<del>S1</del>	-	C	C	G
<del>S2</del>	0 0 0 0			
A	0 0 0 0			
T	0 0 0 0			
C	0 1 1 0			
C	0 1 2 0			
G	0 0 0 3			
T	0 1 0 0			
T	0 0 0 0			
		C C G T		
		C C G		

~~Example~~

Find the best local alignment from the following two sequences by using match score +1 mismatch penalty -1 and gap penalty -2

(S1) S = G T E G G T		Here, $S(C, C) = 1$	
$\text{I}^+ =$		$S(C, N) = 1$	
$\text{I}^- =$		$S(C, -, ) = -2$	
$(7 \times 7)$ matrix		$S(G, C) = -2$	

S1 S2	-	T	C	G	G	T
-	0 0 0 0 0 0					
A	0 0 0 0 0 0					
C	0 0 1 0 0 0					
G	0 0 0 2 1 0					
T	0 1 0 0 1 2					

Alignment-1  
C G  
C G

Alignment-2  
G T  
G T

## Example

Given that,

$$\text{DNA, } \Sigma = \{A, C, G, T\}$$

Input:

$$S = TGTTACGG$$

~~$$T = TGCTTCACTA$$~~

Scoring scheme:

~~$$S(C, C) = -3$$~~

~~$$S(C, A) = 3$$~~

~~$$S(C, U) = S(C, A) = 2$$~~

S	-	T	G	T	T	A	C	G	G
-	0	0	0	0	0	0	0	0	0
C	0	0	3	1	0	0	0	3	3
G	0	0	3	1	0	0	0	3	6
T	0	3	1	6	4	2	0	9	9
T	0	3	1	4	9	7	5	13	2
G	0	4	6	1	7	6	4	8	6
A	0	0	4	3	5	10	8	6	5
C	0	0	2	1	3	8	13	11	9
T	0	3	1	5	9	8	11	10	8
A	0	1	0	3	2	7	9	8	7

best alignment score: 13

example

G T T - A C  
G T T G A C  
3 6 9 7 10 13

length

= 3

: tuqat

ATGATTGT = 2

ATGA-GTT DP approach

time complexity of sequence alignment

- \* time complexity / time grow with the number of sequences

\* for two sequences  $O(n^2)$

\* " three "  $O(n^3)$

\* " k "  $O(n^k)$

### 8 Alternative approaches

→ branch and bound

→ Heuristics

→ Dynamic programming table centred along the diagonal.

## # Heuristic approaches to MSA:

### → Heuristic methods

↳ static alignment

↳ progressive alignment methods

\* CLUSTALW

\* T-coffe

\* MUSCLE

↳ heuristic variant of DP approaches

↳ genetic algorithm

↳ Gibbs sampler

↳ Branch & Bound

### → Basic techniques

\* progressive

\* iterative

\* profile based

\* divide & conquer

## Recurrence for 3 sequences



Align  $(s_i^1, s_j^2, s_k^3)$

in sequences  
 $\rightarrow (2^n - 1)$  recurrence

$$* 2^3 - 1 = 8 - 1 = 7$$

equation

$$V_{i-1, j-1, k-1} + S(a_i, a_j, a_k)$$

$$V_{i-1, j-1, k-1} + S(a_i, -, a_k)$$

$$V_{i-1, j-1, k-1} + S(a_i, a_j, -)$$

$$V_{i-1, j-1, k-1} + S(a_i, -, a_k)$$

$$V_{i-1, j-1, k-1} + S(a_i, a_j, -)$$

$$V_{i-1, j-1, k-1} + S(a_i, a_j, -)$$

$$V_{i-1, j-1, k-1} + S(a_i, a_j, a_k)$$

For K sequences dynamic programming

table will have size  $\underline{\underline{n^k}}$

$$\frac{A+3}{2} \quad \frac{S+8}{2} \quad \frac{T+9}{2} \quad \frac{V+3}{2}$$

## Sum of pairs (SP) score:

\* Consider all pairs of letters in each column

\* add the scores

$$SP \text{ score} = \left( \begin{array}{c} A \\ V \\ V \end{array} \right) = 2 \times \text{Score}(A, V) + \text{Score}(V, V) + \text{Score}(A, -) + 2 \times \text{Score}(V, -)$$

NB K Dvenues gives  $\frac{K(K-1)}{2}$  addends

$$\text{Score}(-, -) = 0$$

$$+ (A, A) 25 + (-A) 25 = 25$$

$$(A, A) 25 + (V, V) 25$$

$$+ (S, S) 25 + (T, T) 25 = 50$$

$$+ (F, F) 25 + (P, P) 25 = 50$$

$$+ (L, L) 25 = 25$$

$$(W, W) 25 = (W, V) 25 + (W, S) 25 + (W, T) 25 + (W, F) 25 + (W, L) 25 = 125$$

$$S = 25 \times 5 = 125$$

### # Example

\* For the following alignment determine the SP score. consider  
 match score = 2 and mismatch/gap penalty = -2

$$S_1 = A C G \text{ -- } G A G A$$

$$S_2 = \text{-- } C G T T G A C A$$

$$S_3 = A C H T - G A - A$$

$$S_4 = C C G T T C A C \text{ -- }$$

Ans

For column 1 (-, -)

$$\text{SP score} = 2S(A, -) + 2S(Ac) +$$

$$S(C) + S(A, A)$$

$$= (2 \times -2) + 2(-2) + (-2) + 2$$

$$= (-4) + (-4) - 2 + 2$$

$$= -8$$

For column 2

$$\text{SP score}(C, C, C) = 6S(C, C)$$

$$= 6 \times 2 = 12$$

~~9x3~~

$\frac{3 \times 2}{2}$

$\frac{2 \times 2}{2}$

$$SP\text{-Score}(G, G, -, G) = 3S(G, G) + 3(G, -)$$

$$(S(G, G) + S(G, -)) = 3(2) + 3(-2) \\ = 6 - 6 = 0$$

$$SP\text{-Score}(-, T, T, T) = 3S(T, T) + 3(-, T) \\ = 3(2) + 3(-2) \\ = 6 - 6 = 0$$

$$SP\text{-Score}(-, T, -, T) = 2(S(G, T) + S(-, -)) + \\ + 4(T, -) \\ = 2(1) + 0 + (-8) \\ = -6$$

$$SP\text{-Score}(G, G, G, e) = 3S(G, e) + 3S(G, e) \\ = 3(2) + 3(-2) \\ = 6 - 6 \\ = 0$$

$$SP\text{-Score}(A, A; A; A) = 6S(A, A) \\ = 6 \times 2 \\ = 12$$

$$SP\text{-Score}(G, e, -, e) = 2S(G, e) + 8(e, e) \\ + 2S(G, -) + S(G, -) \\ = (2)(1) + 2(1) + (-1) + (-2) = -8$$

~~3+2~~

$$\text{SP-Score}(A, A, A, -) = 3(A, A) + 3(A, -)$$

$$(5 \times 8 + 5 \times 8) = 6 + 3(5 \times 2)$$

$$0.5(8 - 2) = 6 - 6 = 0$$

$$(T, -)S + (T, T)R S = (T, T, T, -)$$

$$\text{Total SP-Score} = -8 + 12 + 0 + 0 - 6 + 0 + 12 - 8 + 0$$

$$0 = 2 - 2 = 24 - 22$$

$$(-T)P +$$

$$(8 - 1 + 0 + 0) = -1$$

$$(9.0)S + (9.0)R S = (9.0, 0, -1)$$

$$(5 \times 8 + 5 \times 8) =$$

$$8 - 2 =$$

$$0 =$$

$$(A, A)S = (A, A, A, A)$$

$$5 \times 2 =$$

$$51 =$$

$$(9.0)S + (9.0)R S = (9.0, 0, 0)$$

$$(-1)S + (-1)R +$$

$$8 = (8) + (8) + 15 + (-1) =$$

## ■ Evaluating alignment method:



Given an estimated alignment and a true alignment, compute various statistics (based on homology)

- \* SPFN: Sum of the false negative homology pairs
- \* SPFP: Sum of the false positive homology pairs
- \* TC: Total column score. Number of columns that are identical in the two alignments
- \* Comprehension: Ratio of the estimated alignment length, to true alignment length.

$$\rightarrow \text{SP. score} = 1 - \text{SPFN}$$

measure of recall

$$\rightarrow \text{Modeler score} = 1 - \text{SPFP}$$

measure of precision

Error rate: SPFN, SPFP

Accuracy rate: The score, SP score,  
model score

British sign language: 97.92%

British sign language

Chinese sign language: 87%

American sign language: 65%

British sign language: 81%

Chinese sign language: 80.1%

American sign language: 78%

Percentage of success: 97.92 - 1 = 96.92%

Percentage of failure: 97.92 - 1 = 2.08%

tree → alignment  
(iteratively)

### • Progressive alignment:

- Start with a binary tree (merge tree/guide tree) from the input set of sequences.
- multiple sequence alignment is built from the bottom-up.
- Starting with internal nodes that have only leaves as children.
- Interactively align the alignments to obtain an alignment for larger number of sequences.

### Example

- Ex:
    - \* Clustal
    - \* ClustalW
    - \* Clustal-Omega
    - \* PICOOK → PIAGAN → Canopy
  - Phylogeny
  - \* PICOOK-aware
  - \* MAFFT
  - \* PIAGAN
- improve the accuracy of PICOOK, PIAGAN

Problem: Create a bias free

that produce bias alignment.

\* Doesn't realign the sequences

\* Not guaranteed to coverage

converge to the optimal solution.

→ the final alignment will be less accurate

if we have poor initial alignment.

Gap penalty

gap-opening

gap-extension

match/mismatch

match/mismatch

gap

gap-opening → gap-extension

gap-opening → gap-extension

## Iterative alignment:

\* Generate an initial MSA using a method like progressive alignment.

\* Iteratively improves the MSA.

(Example)

SATE, SATe<sup>-2</sup>, PASTA

process terminate

after fixed number of

iterations

of iterations

Complexity

DTAAST

For N sequences, running time

$$O(c^{2n}) \cdot n^2$$

DTAAST = 2

DTASA = 2

$$S = (N \times D)^2 \cdot 2 \cdot 2 = (N \times D)^2 \cdot 4$$

$$S = (N \times D)^2 \cdot 2$$

April 2020

Find the optimal local alignment of the following two sequences using -2 as the gap penalty, -1 as the mismatch penalty, and 2 as the score for a match. You have to implement the Smith-Waterman algorithm. Show the DP table and mark the path, which corresponds to the alignment in the DP table.

TCAATG  
ACATAG  
S<sub>1</sub>: T C A A T G      S<sub>2</sub>: A C A T G  
              | | | | |  
              C i r c l e      Given

$$S_1 = \text{TCATG}$$

$$S_2 = \text{ACATG}$$

$$S_{(i,j)} = 2 \quad S_{(i,j)} = -1$$

$$S_{(i-1,j)} = -2$$

Longest

<del>S<sub>11</sub></del>	A	C	A	T	G
T	0 0	0	0	0 0	0 0
C	0 0	2	0	0 1	0 1
A	0 2	0 1	1	2 1	2 0
A	0 2	1	2	3 1	1
T	0 1	1	0 1	4 2	2
G	0 0	0	0	2 6	6

highest alignment score = 6

ACTT  
 ACTT  
 ACTT  
 ACTT

ATG  
 ATG  
 ATG  
 ATG

By April 2021

Find the best (in terms of the alignment score) shortest and longest local alignments of the following two sequences using DP  
Penalty of -2, a mismatch penalty of -1 and a match score of 2.

You have to use the Smith-Waterman Algorithm, and show the corresponding dynamic programming (DP) table.  
Please Show the alignments and  
marks the paths that correspond  
to the alignments in the DP table

CCAGTC  
CTAT

Given that,

$$S_1 = \text{CCATGCG}$$

$$S_2 = \text{CTATGG}$$

$$\begin{aligned} S(n) &= -1 \\ S(n, n) &= 2 \\ S(n-) &= -2 \end{aligned}$$

Improbable bad match 2

Improbable bad match

TAATG

TGAGG

$\Sigma = S - S + S$

$S_1 \downarrow$	C	T	A	T
/	0	0	0	0
C	0	2	0	0
C	0	2	1	0
A	0	0	1	3
G	0	0	0	1
T	0	0	2	0
G	0	0	0	1
C	0	2	0	0

Bent: Alignment score = 3

Shortest best alignment:

$$CTA \text{ A } G T S = 2$$

$$\begin{matrix} C & C & A & G & T & S \\ 2 & -1 & 2 & = 3 \end{matrix}$$

Longest best alignment

$$CTA \text{ } \underline{G} \text{ } T$$

$$CCAGT$$

$$2 \text{ } -1 \text{ } 2 \text{ } -2 \text{ } 2 = 3$$

$$I = (G, N) 2$$

$$S = (N, W) 2$$

$$S = (-W) 2$$

T	A	T	G		58 ↓R
O	O	O	O	O	
O	O	O	S	O	S
O	O	D	S	O	S
I →	E	I	O	O	A
I	I	O	O	S	
S	I	O	O	S	
E	O	S	O	O	T
I	I	O	O	O	S
O	O	O	S	O	S

Score: 58

## Genome sequencing:

Determining the order of nucleotides or bases in a DNA molecule/gene.

## Genome assembly

Process of putting a large number of short DNA sequences (read) back together to recreate the original chromosome from which the DNA originated.

OTL, it refers to aligning and merging fragments from a longer DNA sequence in order to reconstruct whole genome

- Human genome: 3 billion base pairs
- small bacterial genome: few million base pairs

Q Why challenging?

Ans: ~~Sequencing machines can't read whole genome at a time~~  
~~It can determine small fragments called reads.~~

# Read: ~~Single observation of the (partial) sequence of a DNA molecule~~

# Shotgun: ~~Random fragmentation of the whole genome, like if it was fired from a shotgun~~

# Contig: ~~Contiguous stretch of sequence often derived from multiple reads~~

# Scaffold: ~~Linearly ordered non-contiguous oriented group of contigs~~

## Steps in genome sequencing:

- \* Start with many copies of genome
- \* Fragment them and read one or both ends
- \* Find overlapping reads
- \* Merge them into contigs
- \* Scaffold contigs using paired read ends.

Coverage: It refers to the number of reads covering a particular position in the genome.

Average coverage: Average number of reads

covering a position in the genome

↳ Some people use average coverage per position

(approx) Reads per base pair

## Types of coverage:

- small coverage (shallow sequencing)
- High coverage (Deep sequencing)

## Basic challenges in genome assembly:

- \* two main challenges are:
  - getting sufficient coverage of the genome (read length, number of reads, size of genome etc.)
  - Assembling the reads into a complete genome.

## Issues:

- ↳ Coverage (more is good)
- ↳ Errors in reads (error by sequencing machine)
- ↳ genomes are double stranded.  
Reads vary from very short (35bp) to quite long (800bp)

↳ Repeats ( $\text{repeat length} > \text{read length}$ )

↳ Running time and memory

↳ Double strand input of loups

Mistakes made by reference matching

Servicing errors: flufflib

Mistakes made by reference matching

Repeat: Some portion of genome can be

repeated more than one time in the whole genome

Problems when

\* repeats are longer than read lengths

Advantages: repeats are present in many copies

\* Fast analysis  
- short reads can map to several places

Disadvantages:

\* needs large reference genome  
\* repeats are problematic

## Read length

Assembly should be

- \* Reads much longer than repeats → Assembly should be easy
- \* Reads equal to the repeats → Assembly computationally difficult (NP-hard)
- \* Read shorter than repeats → Assembly undetermined

\* Human genome is  $\sim 5\%$  repetitive

## Coverage

More coverage is good for better assembly

## Summary of challenges

## Types of genome assembly:

→ De novo assembly → Do everything from scratch

→ Reference based / Comparative assembly

If we have a reference genome

Comparative assembly → easier than De novo

→ take the reads

→ map them onto the reference genome

→ collect all overlapping reads

→ produce a MSA

→ produce consensus sequence

### Advantages:

→ highly accurate even when

have errors

\* Fast

\* Short reads can map to several places  
(if they have errors)

### Disadvantages:

\* Needs close reference genome

\* Repeats are problematic.

■ De novo assembly:

- Much easier to do with long reads
- Need very good coverage
- Necessary when you don't have a closely related correctly assembled reference genome

■ De Novo Assembly paradigms:

- \* Overlap graph
- \* K-mer graph / De Bruijn graph

→ Overlap layout consensus methods

- ① greedy (TGICR Assembler, phrap)
- ② graph-based (Celera assembler, Arachne)

## Shortest common superstring (SCS)

input: give collection of strings  $S$

output: find SCS

\* without requirement of shortest:

\* just concatenate input strings

## # String composition problem: (String $\rightarrow$ K-mers)

Given a string that its K-mer composition is a collection of all K-mer substrings of Text. (includes repeated K-mers)

Example:

Composition<sub>3</sub>(TATGCGGTGC) = {TAT, ATG, TGG, GGG, GGC, GCT, GTG, TCC}

There will be  $(\text{len}(\text{Text}) - K + 1)$  K-mers.

NB: genome assembly is inverse problem

$\xrightarrow{\text{# K-mer from a text}}$

Forward problem :: Substrings

## # String reconstruction problem

(K-mers → string) tip: fugri

Reconstructing a string from its K-mer composition.

Input: an integer  $k$  and a collection of patterns of  $k$ -mers

Output: A string  $text$  with  $k$ -mers

\* Connect a pair of  $k$ -mers if they

overlap in  $k-1$  positions

Example:

$$S = \{ATA, ATG, CTT, TAA, TGT\}$$

Simple concatenation

• Superstring length =  $3 \times 5 - 15$

TAA

AATG

CTT

TAA

TGT

GTT

TAA TGT TT

∴ Superstring length = 7

## Example 2 :

$S = \{AAT, ATG, ATG, ATG, CAT, CCA, GAT, GCG, GGA, GGG, GTT, TAA, TGC, TGG, TGT\}$

best prefix = AAT ATG ATG ATG CAT CCA GAT GCG GGA GGG GTT TAA TGC TGG TGT

TAA

AAT

ATG to ATG ATG ATG CAT CCA GAT GCG GGA GGG GTT TAA TGC TGG TGT

We can extend it using TGC, TGG, TGT

TAA

AAT

ATG

TGT

GTT

If we extend by using TGT there is no other extension other than GTT ← TGT

Extend by using TGC

TAA

AAT

ATG

TGC

GCC

CCA

{CAT TAA} :)

ATG

TGG {J.D.}

CCA

GAT

ATG  
TGT

## Overlap graph:

This is a directed graph in which

- vertex is a read, a directed
- edge is on overlap between suffix of source and prefix of sink.

Weight of edges are numbers of overlap.

In summary,

vertex → read  
edge → overlap between suffix of source to prefix of sink

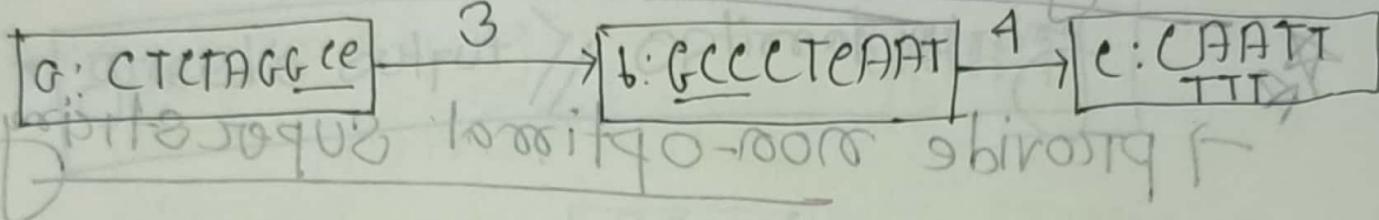
weight → number of overlap.

## Example:

$$V = \{a: CTCAGGCG, b: GCGCTCAAT,$$

$$c: CAATTATTT\}$$

$$\text{edge} = \{(b, b); (b, c)\}$$



# How can we solve it? (finding SES from overlap graph)

- \* modify overlap graph where,

$$\text{cost} = -(\text{length of overlap})$$

- \* SES corresponds to a path that visits every node once, minimizing total cost along path.

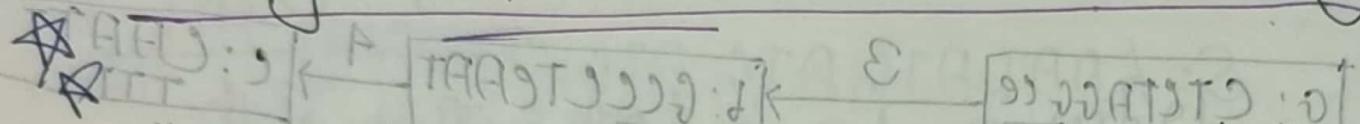
- \* So it is nothing but a Traveling Salesman Problem (TSP)

$\therefore$  It is NP hard

If we disregard edge weights and just look for a path that visits all the nodes exactly once

visit  $\Rightarrow$  Hamilton path problem (NP complete)

## Greedy Shortest Common Superstring



→ Provide non-optimal superstring

→ greedy algorithm chooses longest overlap between two reads.

→ Ties for best overlap arbitrarily

→ Concatenate two strings

→ Update graph

## Failure scenario of Greedy SES

\* not guaranteed to choose overlap

yielding SES

\* It has a good approximation ratio

(ratio ~ 2.5)

o Superstring by greedy SES is

(it won't be more than ~2.5 times

longer than the true SES

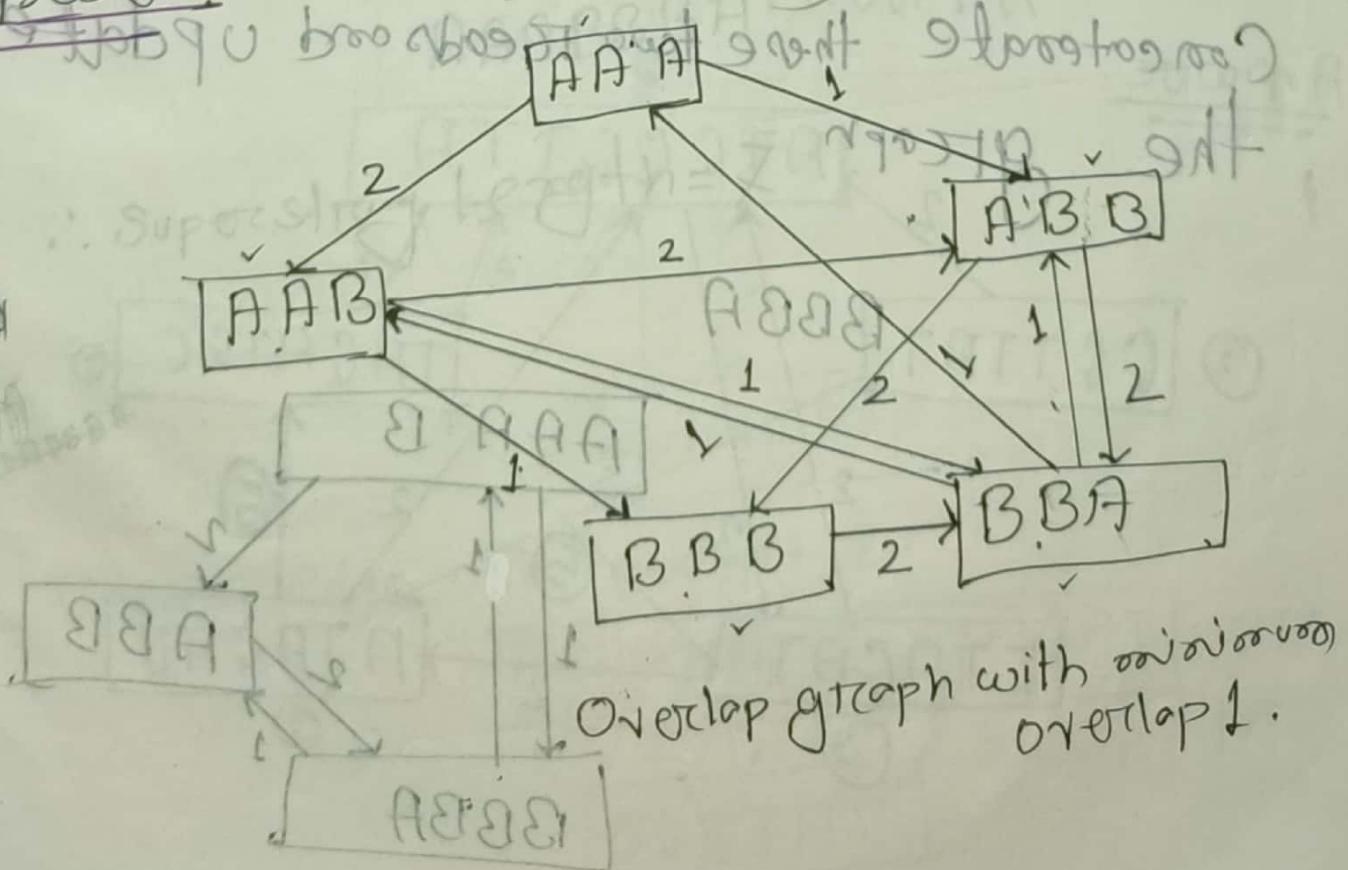
With many results more than one fragment

\* greed's output // optional strategy

Example: Construct overlap graph with minimum overlap 1 and find a shortest common superstring from that graph.

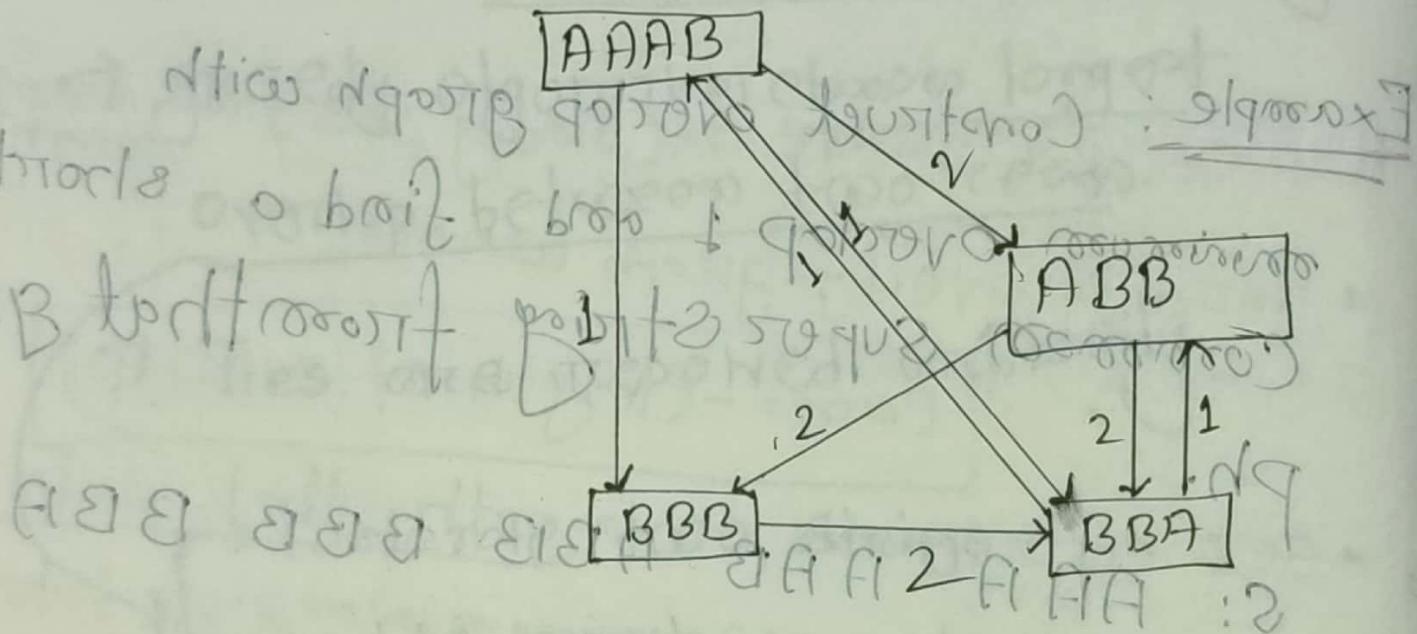
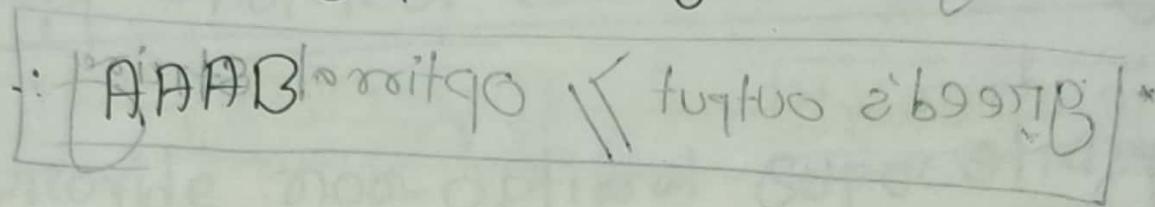
Ph:  
S: AAA AAB ABA BBA BBA

Answer:



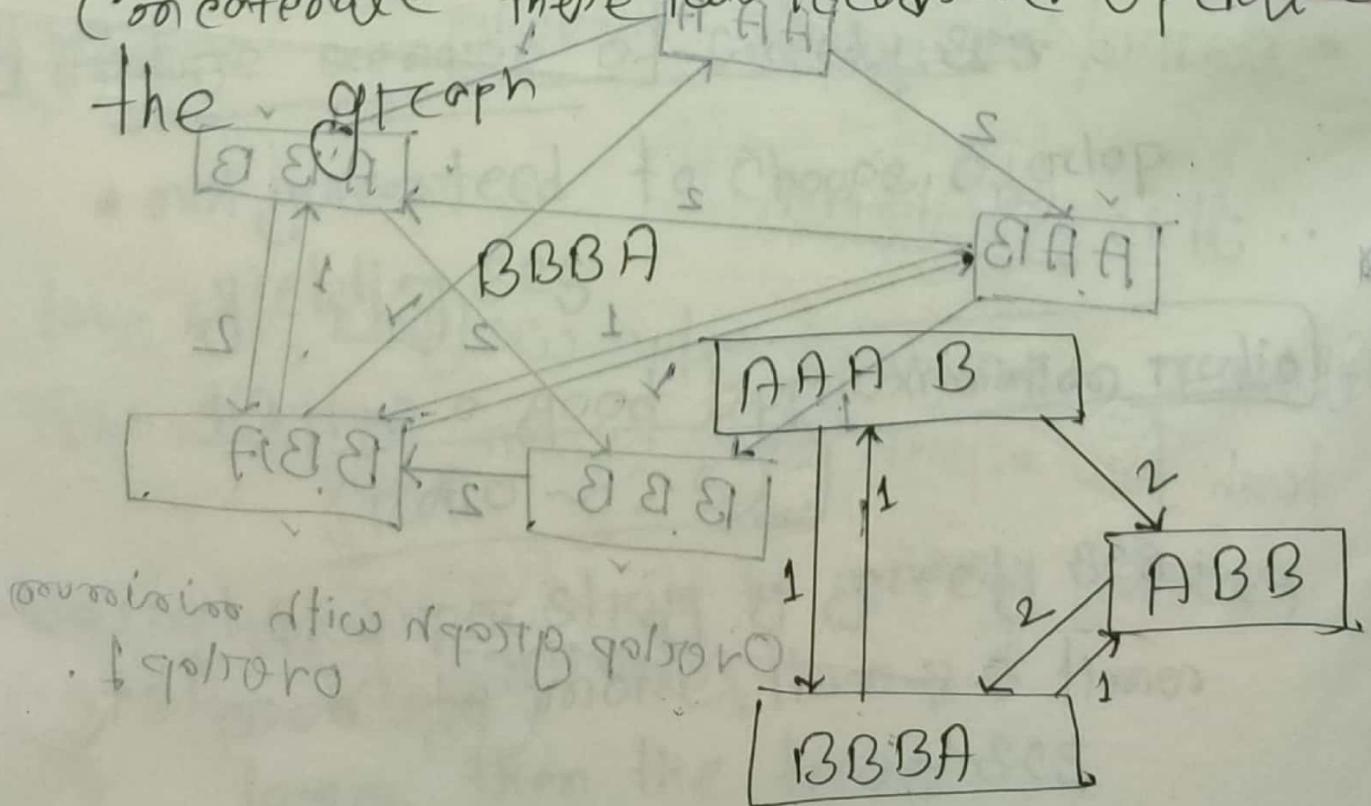
Choose edge (AAA  $\rightarrow$  AAB)

Concatenate these two strings and update the graph

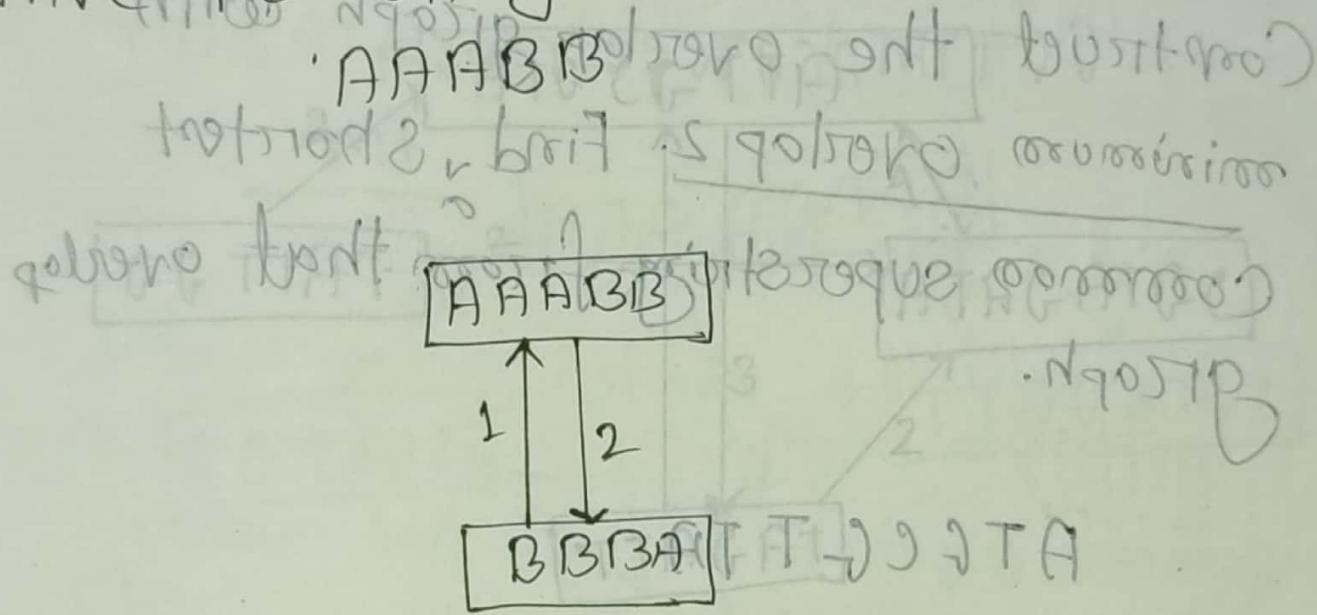


Choose edge (BBB  $\rightarrow$  BBA)

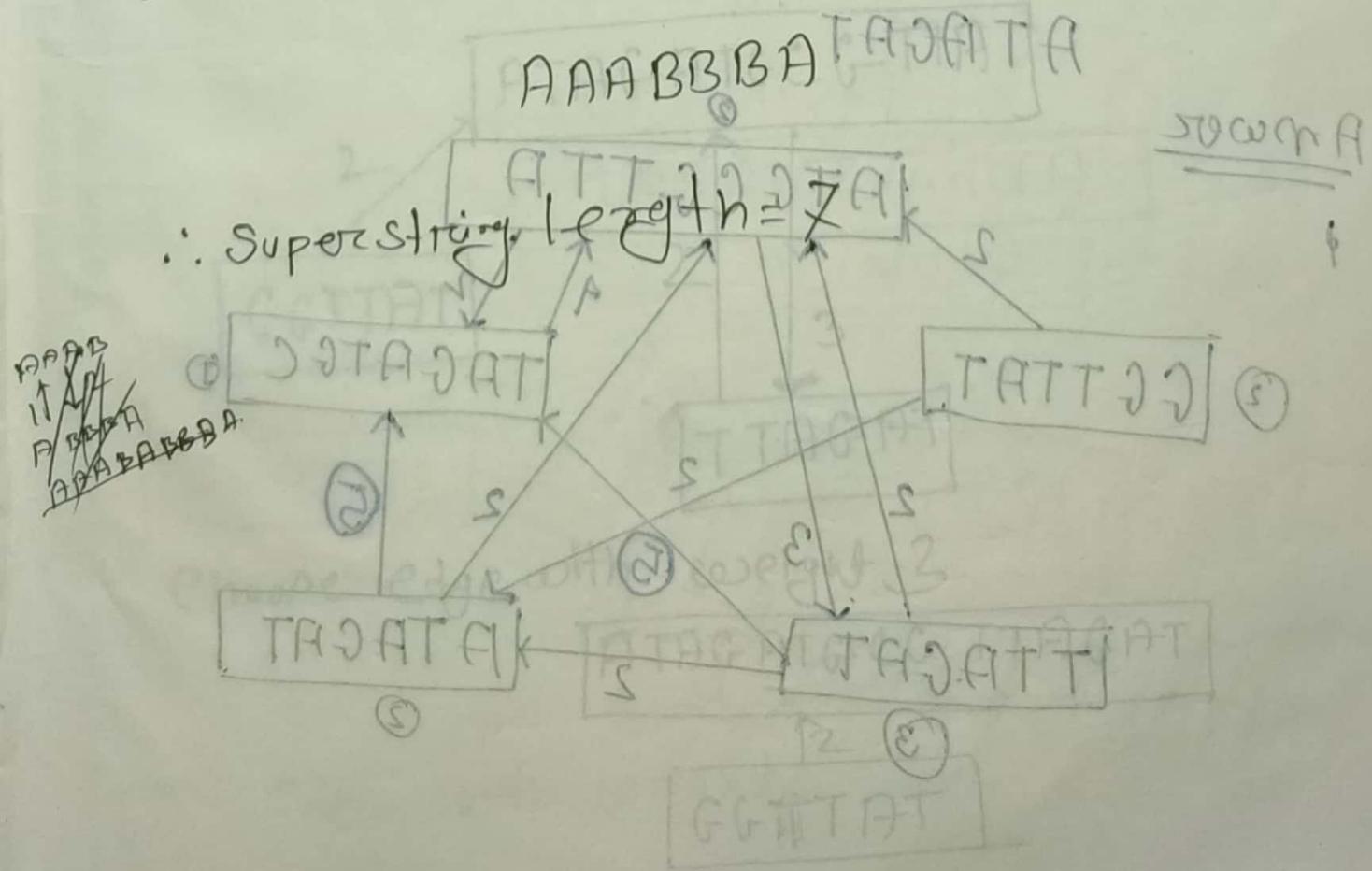
Concatenate these two edges and update the graph



Choose edge ( $AAAAB \xrightarrow{?} ABB$ ) and concatenate these two strings.



Choose edge ( $AAAAB \xrightarrow{?} BBBB$ ) and find final string.



April 2021

Example Consider the following five reads.

Construct the overlap graph with minimum overlap 2. Find shortest common superstring from that overlap graph.

ATGCGTTA

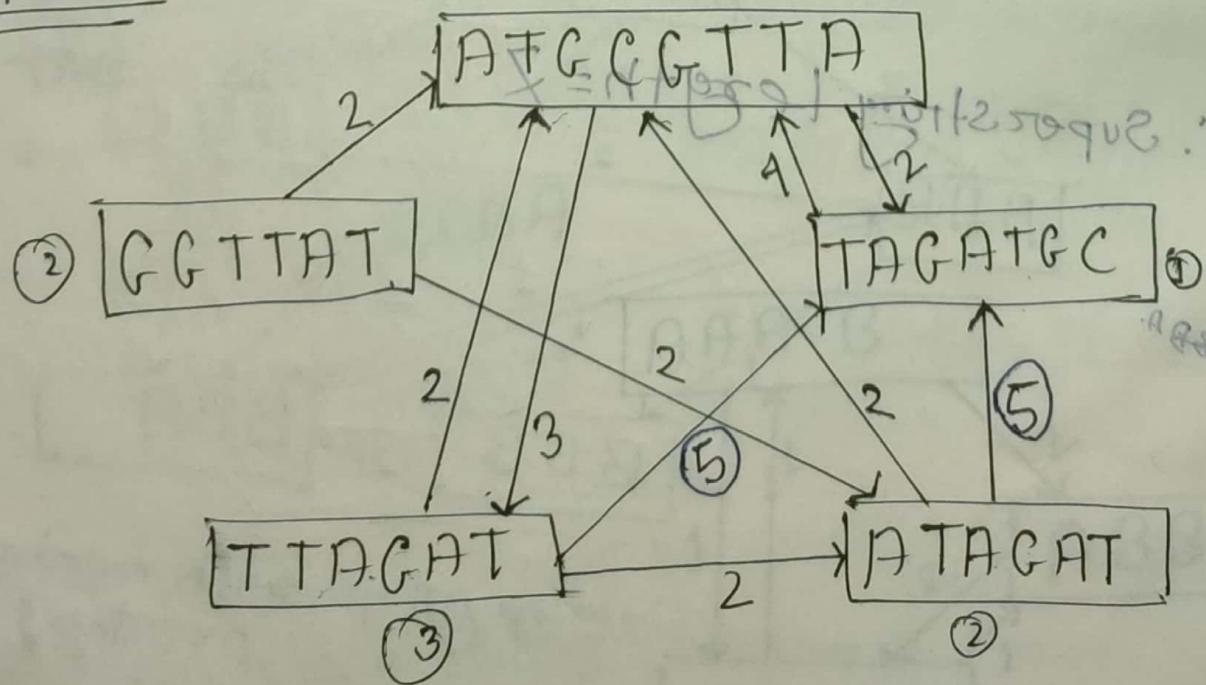
GGTTAT

TAGATGC

TTAGAT

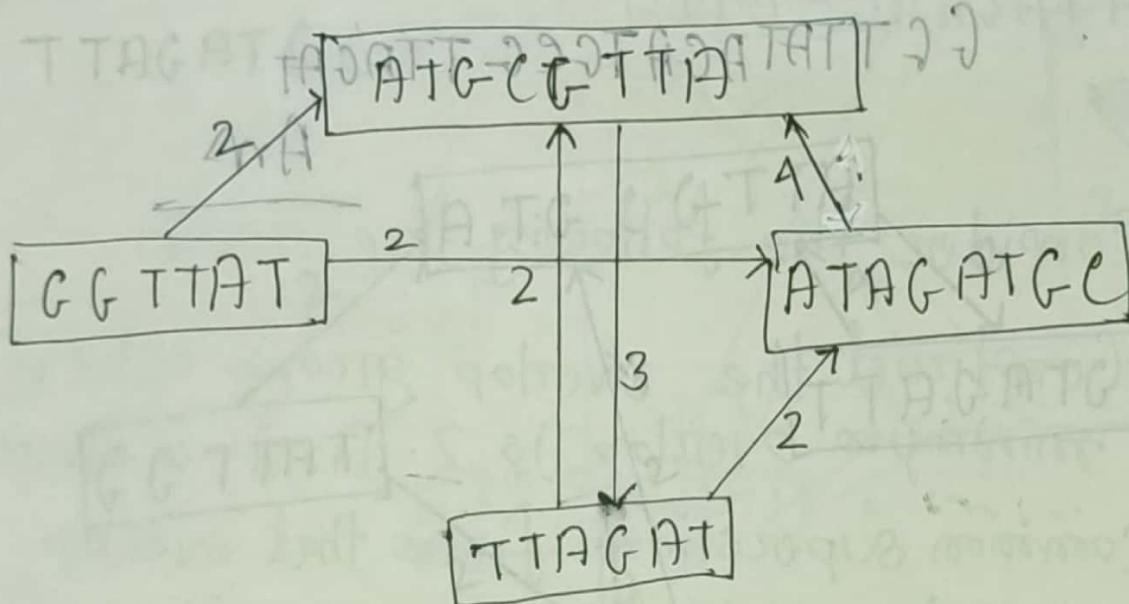
ATAGAT

Answer

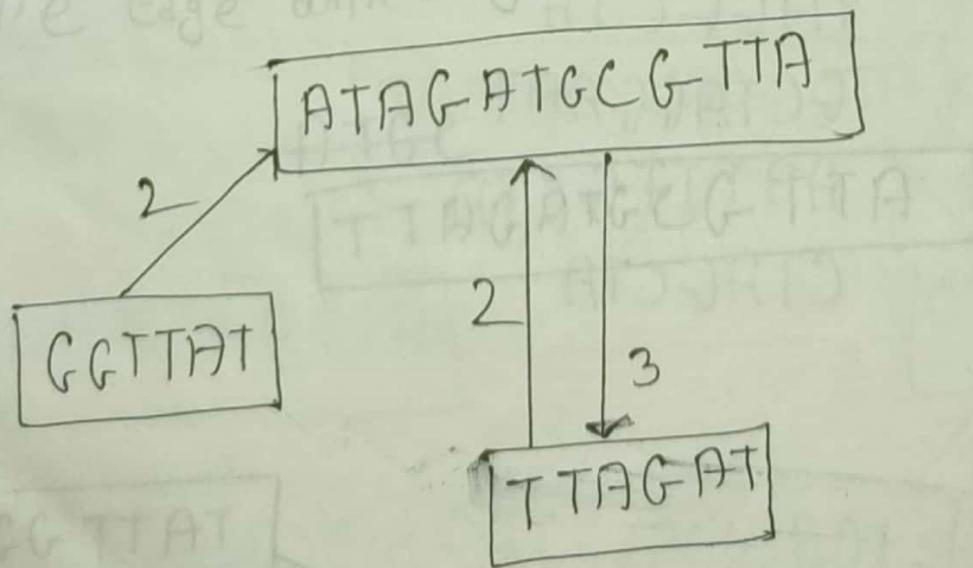


Approach -1

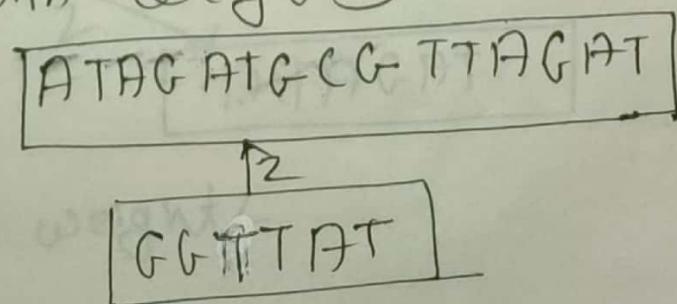
Choose edge ATAGAT  $\xrightarrow{5}$  TAGATGC



Choose edge with weight 1



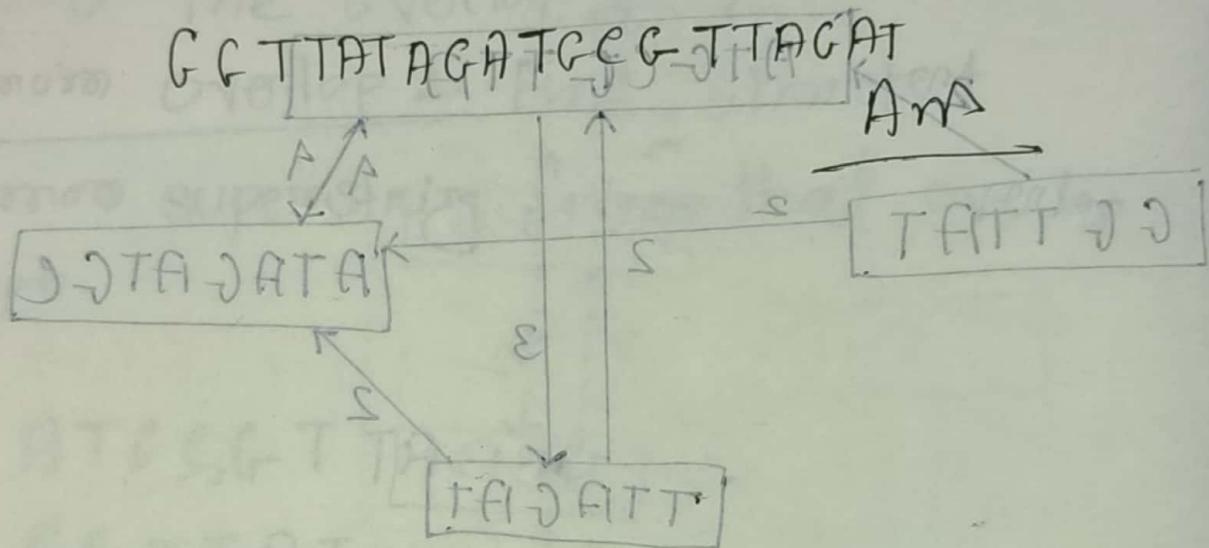
Choose edge with weight 3



Appal 1

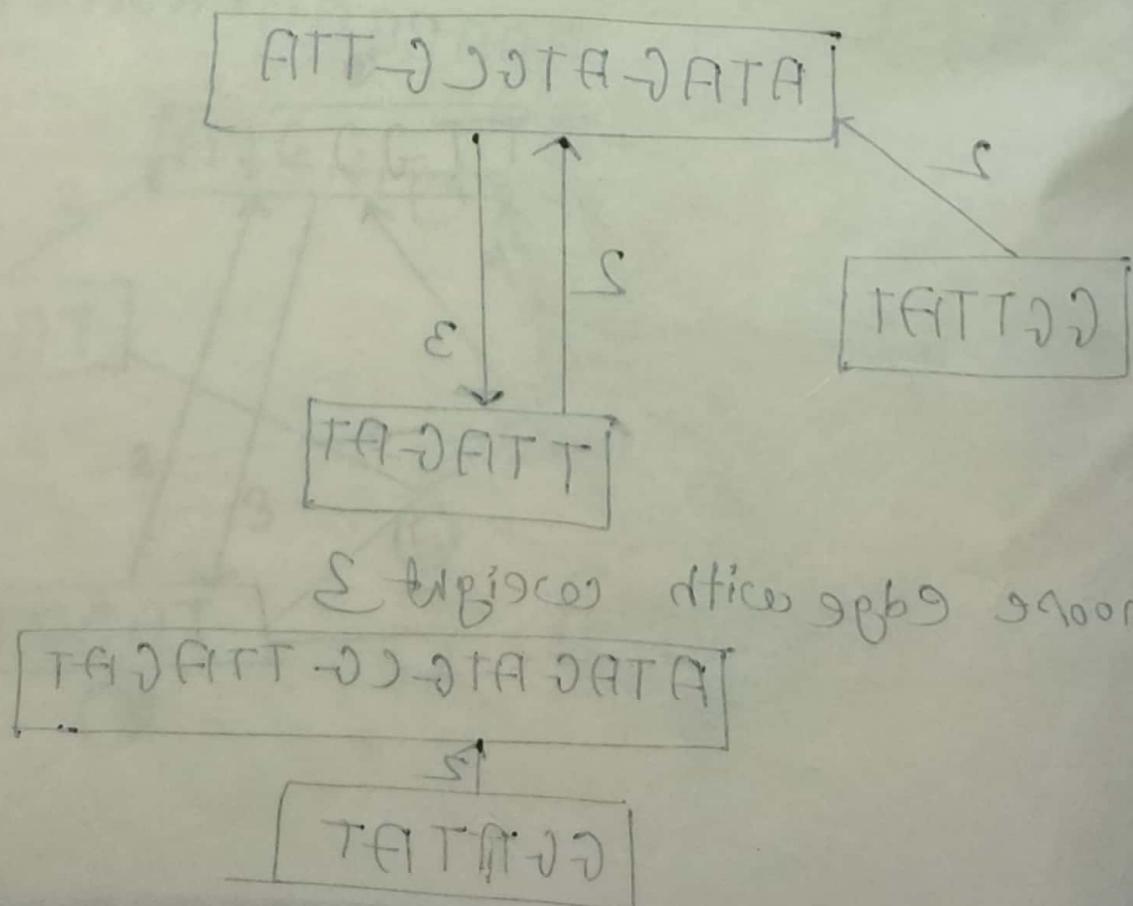
Exo

Choose edge with weight 2



Ans

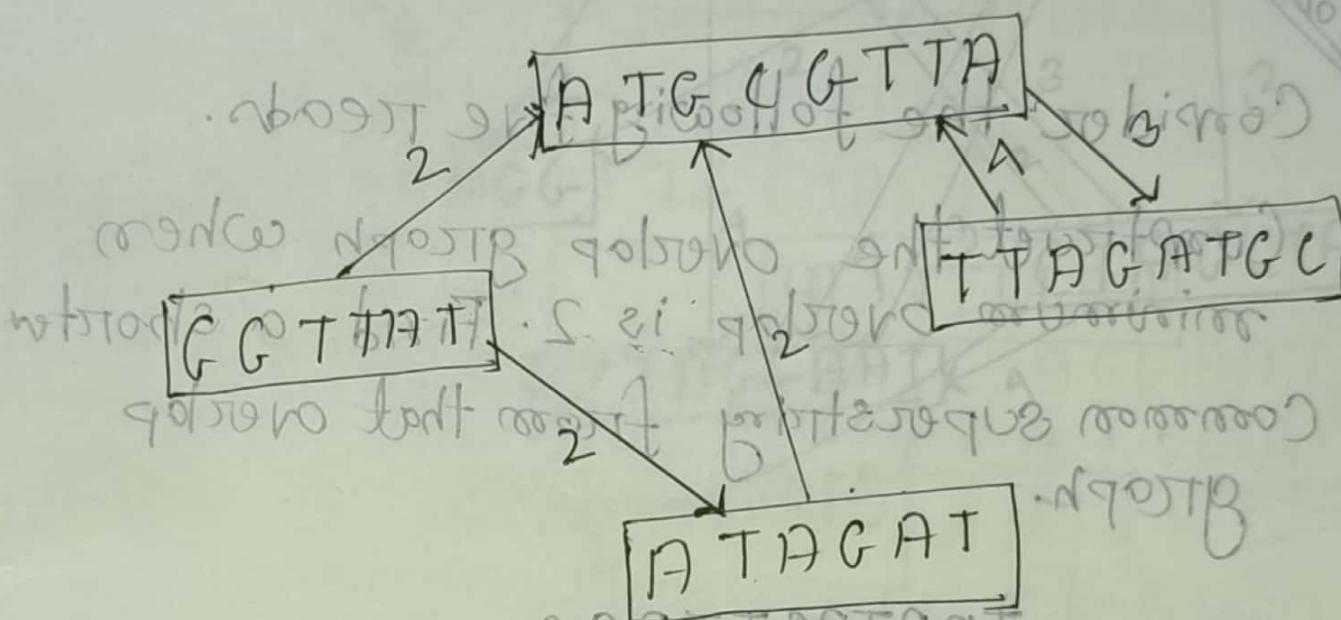
②



Approach 2  
Choose edge ( $TTAGAT \xrightarrow{?} TAGATGTC$ )

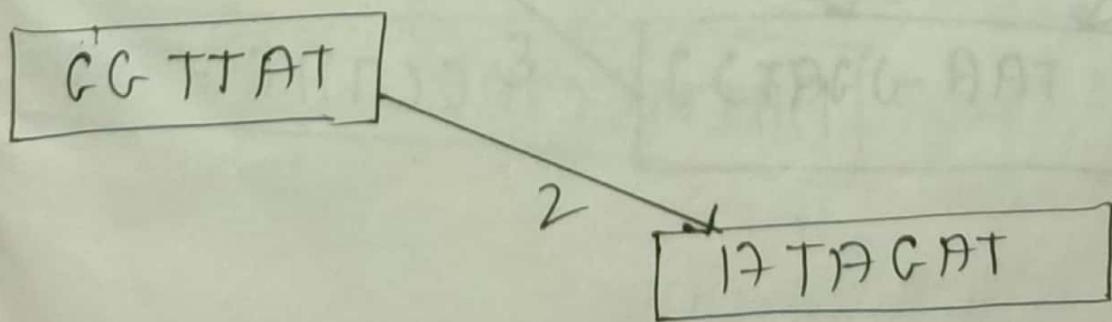
TTAGATGC

ATTCGGCTAATT



Choose edge with weight 1 and concatenate.

ATCGC  
ATCGATGCGATA



Choose edge with weight 2

CGTT ATAGAT TATT  
TTAGATGCCCTTA

APRIL 2020

Consider the following five reads.

Construct the overlap graph when minimum overlap is 2. Find the shortest common superstring from that overlap graph.

TAATATA

TAATACTTAGG

TAGCTA

GCTAGGAAT

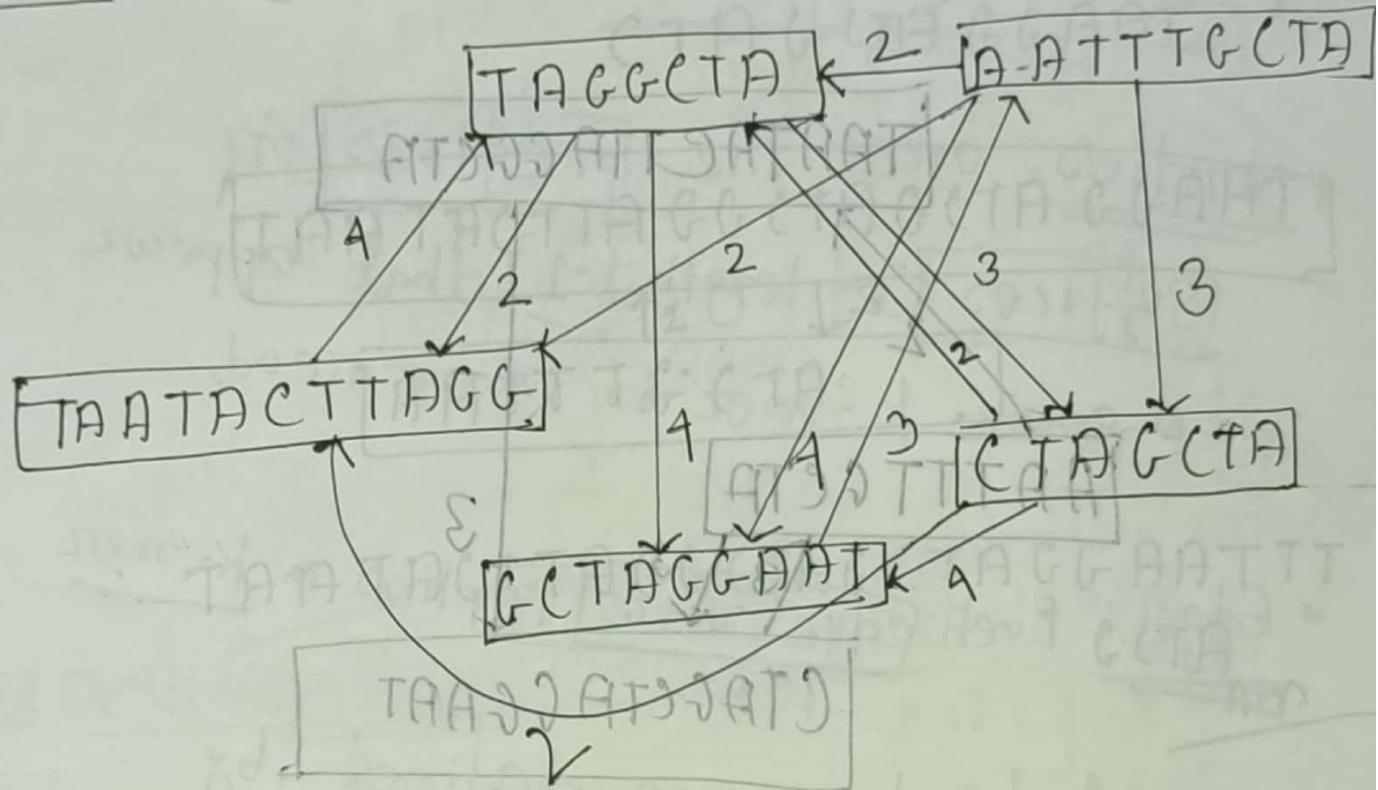
AATTGGTA

CTAGCTA

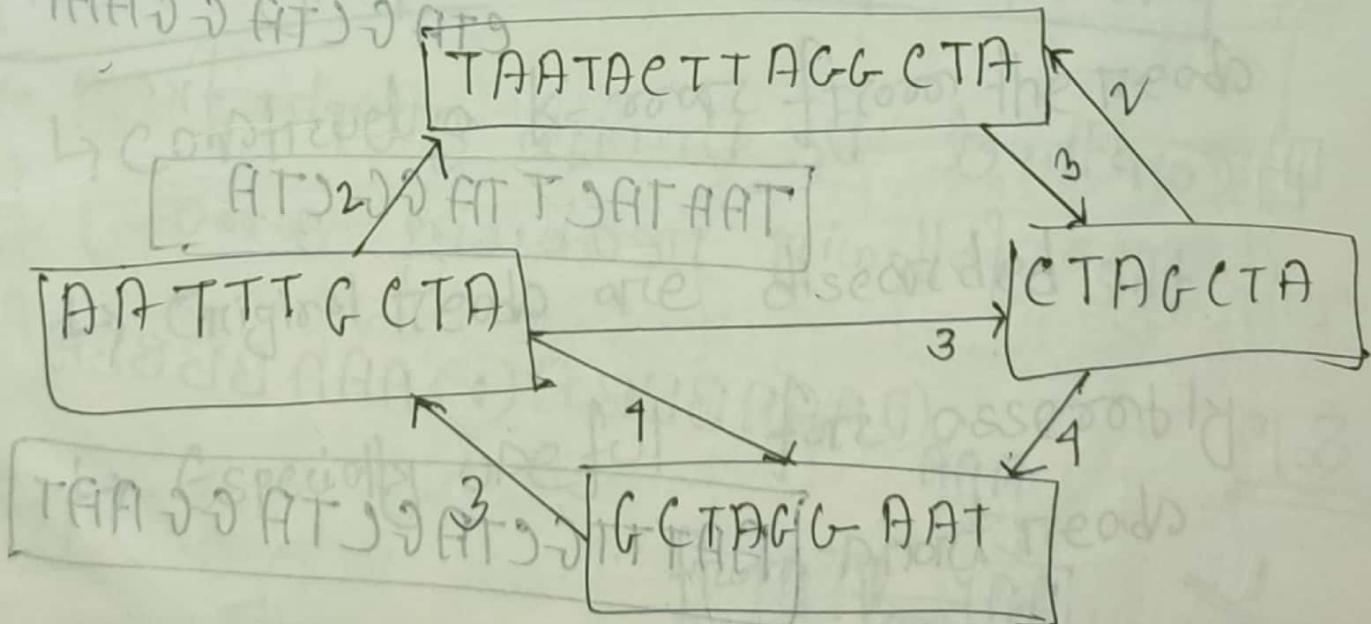
TATT

TAATATA

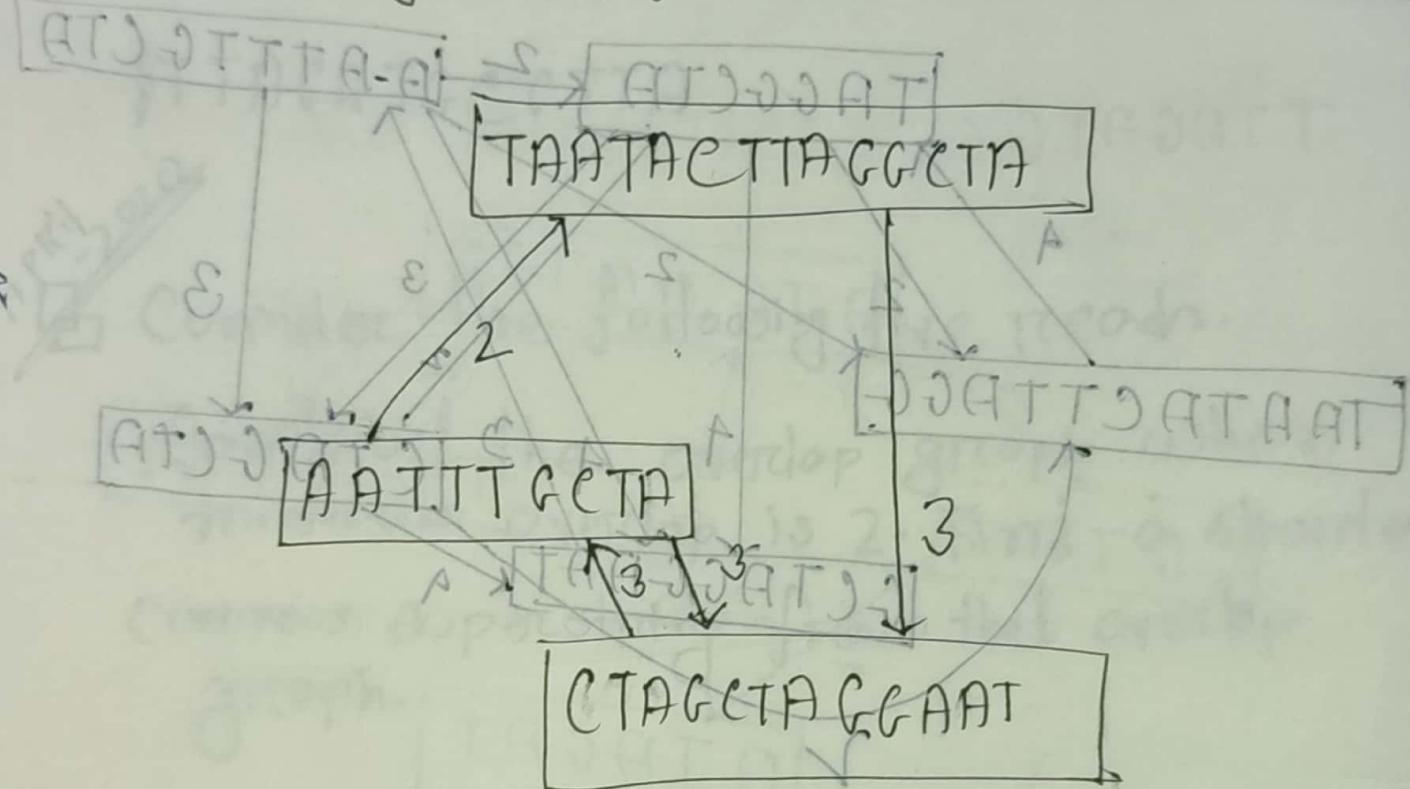
Answer



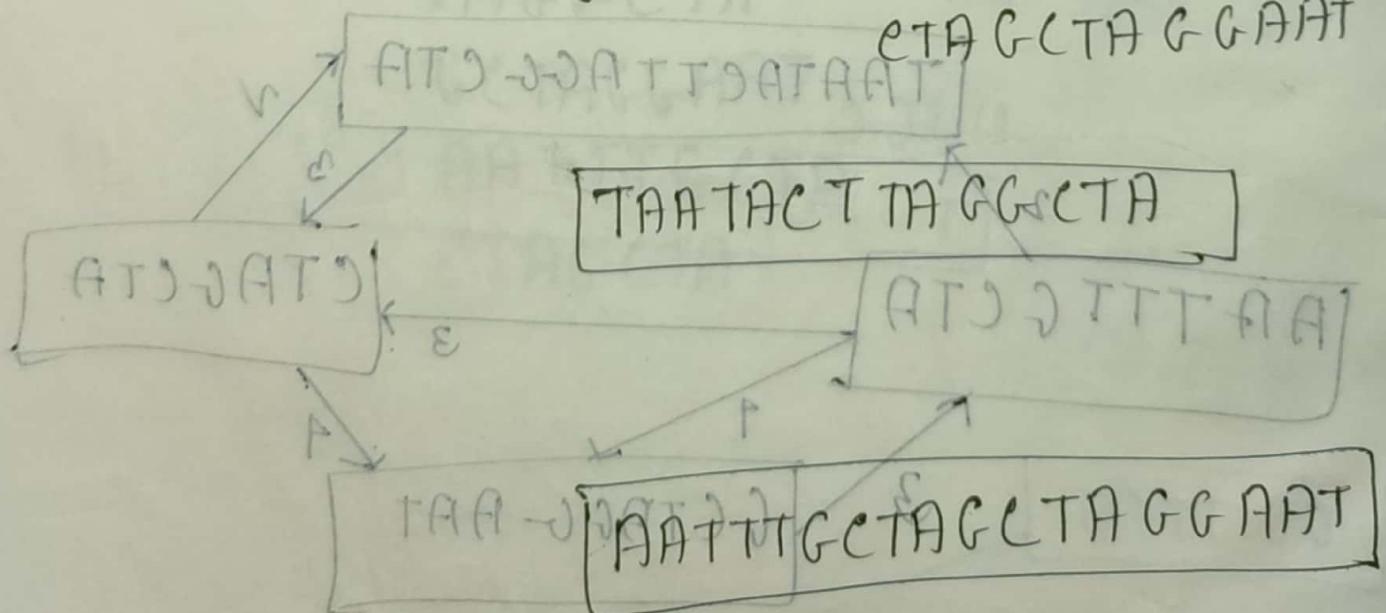
Chosen  $\rightarrow$  TAATACCTTACG  $\xrightarrow{4}$  TAGGCTA



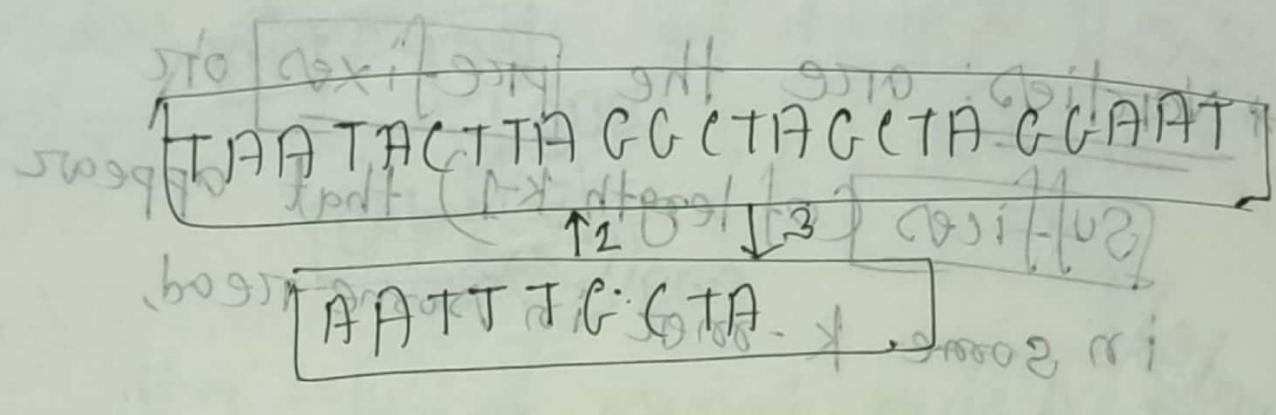
Choose edge  $\text{CTAGCTA} \xrightarrow{A} \text{CTAGCTAGGAAT}$



Choose edge with  $\text{AAATTGCTA} \xrightarrow{3} \text{CTAGCTAGGAAT}$



Choose edge with TAA TA CTTA GG CTAT<sup>3</sup>  
CTA G CTAT GG AAT



TAA TACTTA CG CTAG CTA CG AAT  
AATT TG GTAA

### K-mer graph

↳ Constructing K-mer from the reads

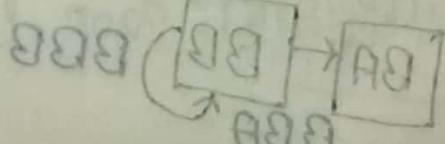
↳ Original reads are discarded

↳ Especially useful for assembly

from short reads

Eulerian path

Path



Endpoint of  
multiple paths

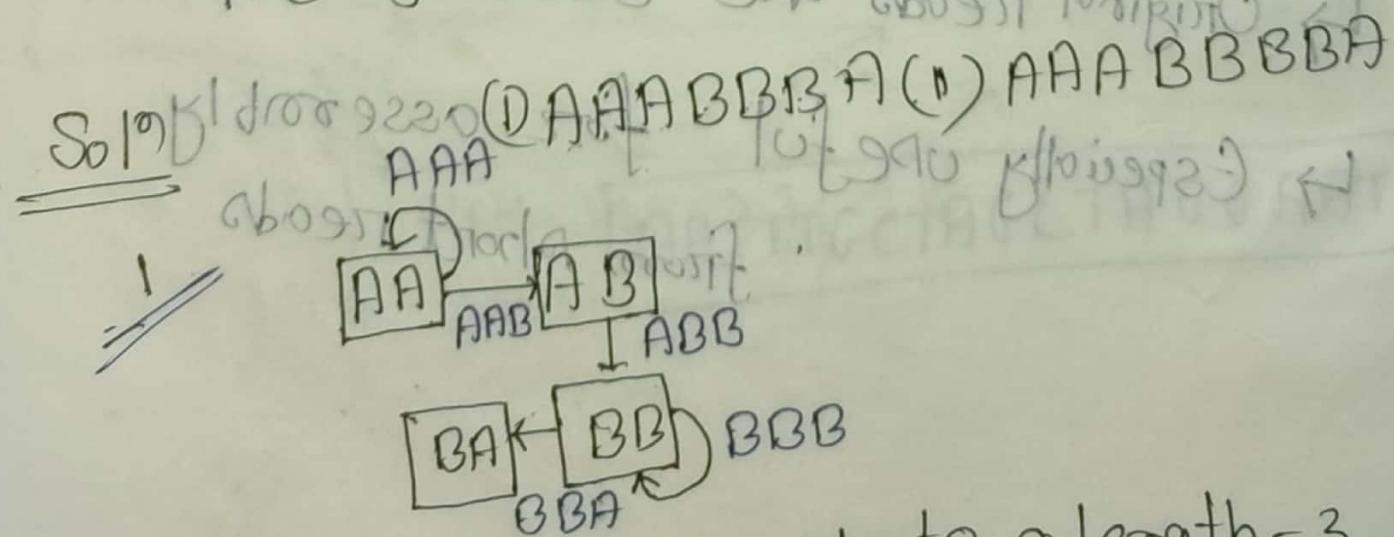
## De Bruijn graph

- \* Vertices: are the prefixes or suffixes (of length,  $k-1$ ) that appear in some  $k$ -mer in some read.

- \* Edge: Each edge  $v \rightarrow w$  implies for  $k$ -mer

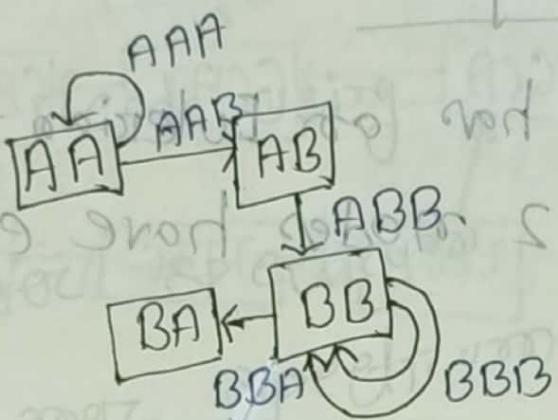
- \* Directed edges are defined by overlap of  $k-2$  nucleotides

- \* Construct De Bruijn graphs from the following reads using 3-mers



- \* Each edge corresponds to a length- $3$  input string.

Given read: AAA BBBBBA, K=3



↓  
value

and  $\text{dist} \geq 1$

and  $\text{dist} \geq 1$

Sensitive to sequencing errors

Characteristics:

### De Bruijn Graph

- \* Small value of  $K$  produce small graph

- \* Doesn't require all pairwise overlap calculations

- \* Also produces fragmented assemblies

- \* Long range connectivity information is lost for small  $K$

### Assembly / finding string from De Bruijn graph

→ by using Eulerian path.

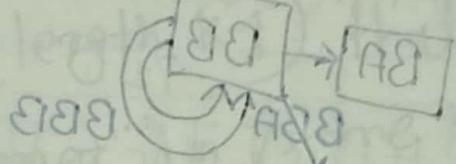
Eulerian path:

Path that goes through every edge exactly once.

## Condition of Eulerian path

~~undirected graph~~

if a graph has an Eulerian path  
the all but 2 nodes have even degree.



~~directed graph~~

A graph has an Eulerian path  
if and only if all the nodes

$$\text{indegree}(v) = \text{outdegree}(v) \quad \forall v$$

all but 2 nodes are even where

$$\text{indegree}_v = \text{outdegree}_v + 1$$

$$\text{indegree}_v = \text{outdegree}_v - 1$$

## Failure scenario of De Bruijn Graph

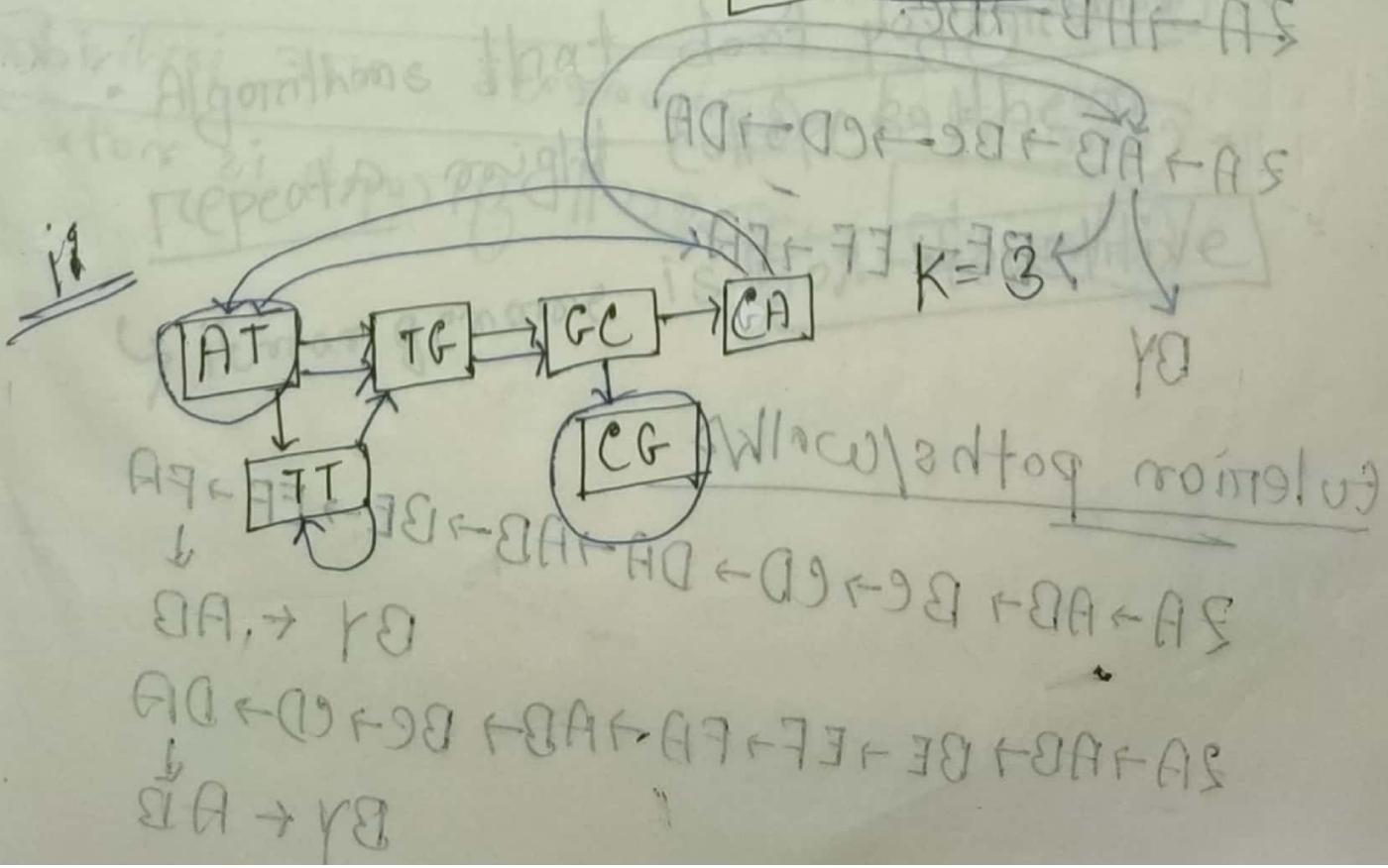
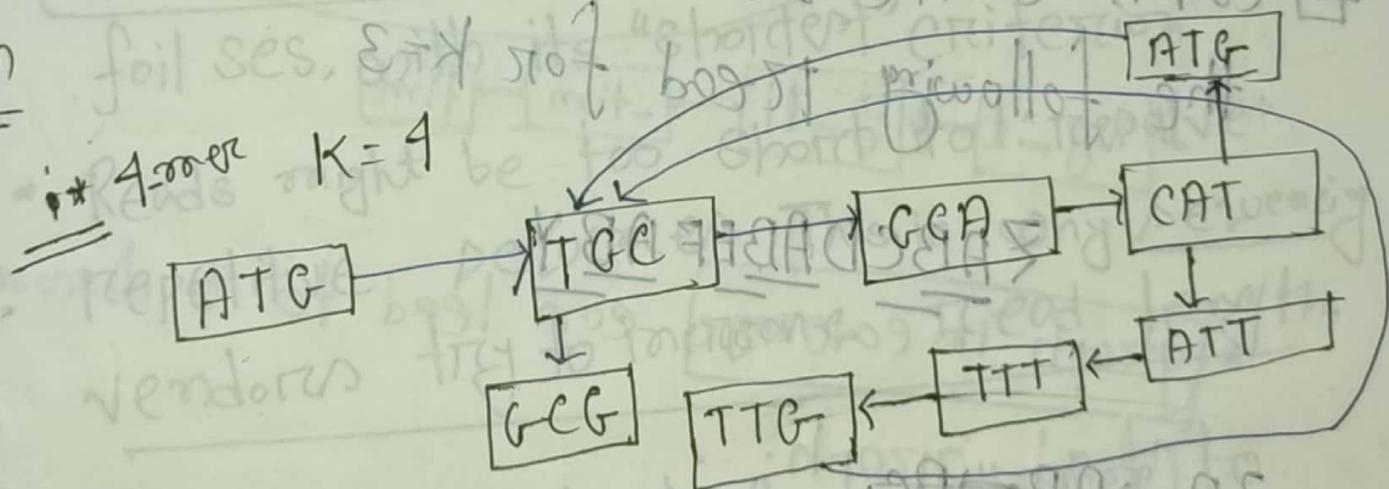
- \* Need to know the correct order of k-mers if there are repeats in the string

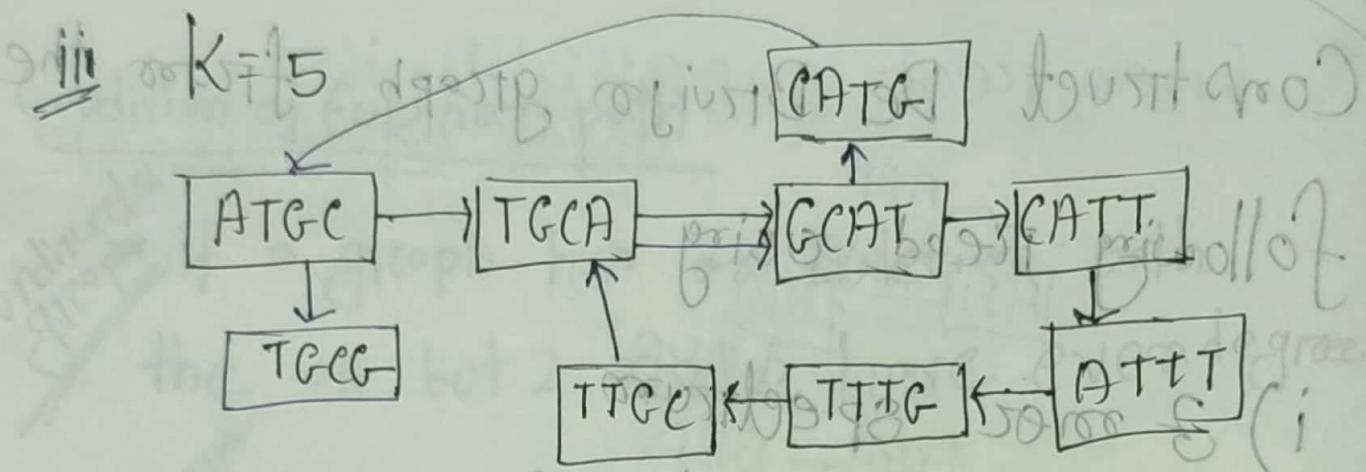
April  
2024

Construct De Bruijn graph from the following reads using i) 3-mers spectrum ii) 5-mers spectrum

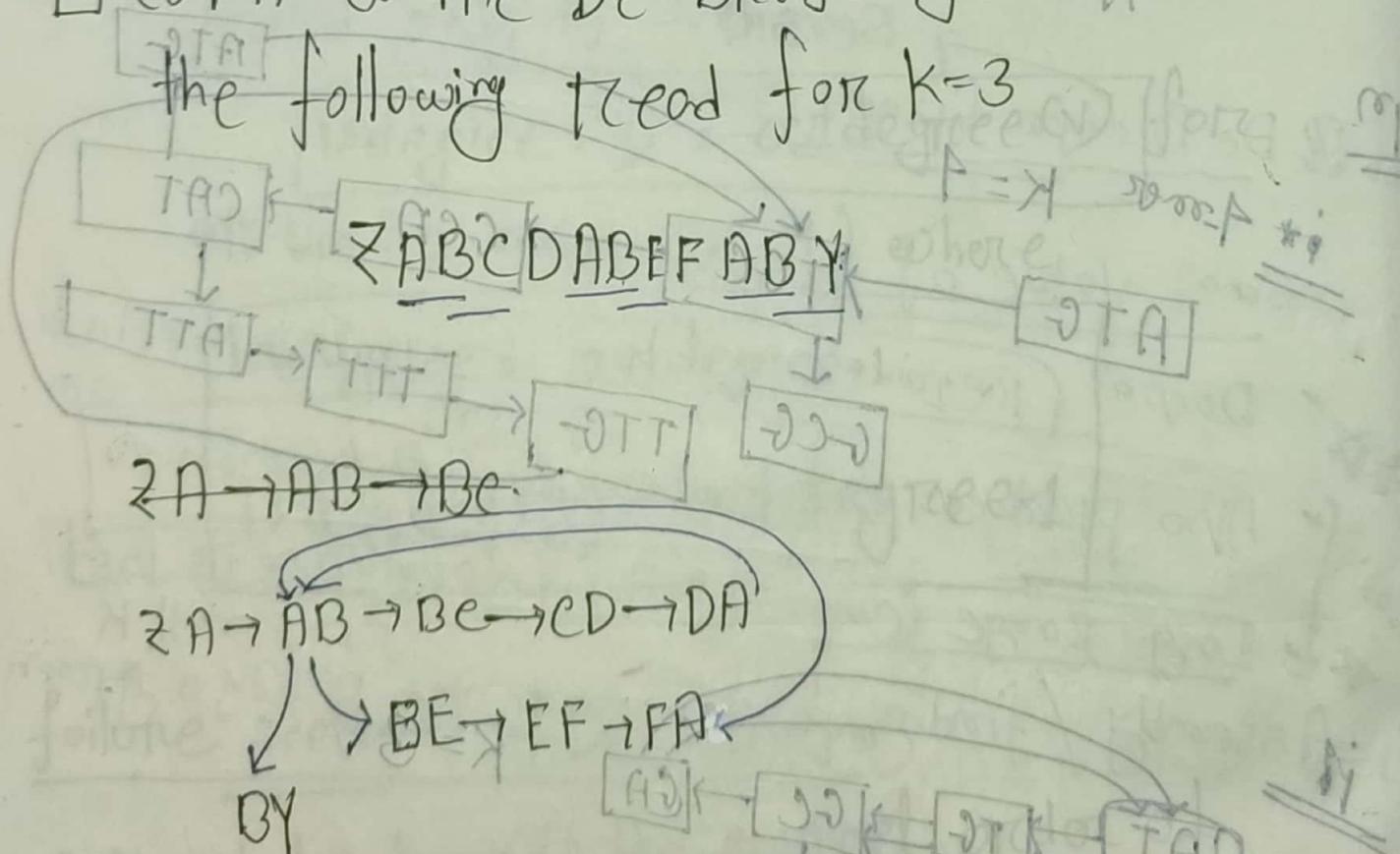
ATGCATTTGGATGCG

Soln





□ Construct the De Bruijn graph from the following read for  $K=3$



Eulerian paths/walks

$Z \rightarrow A \rightarrow B \rightarrow C \rightarrow D \rightarrow A \rightarrow B \rightarrow E \rightarrow F \rightarrow F \rightarrow A$

BY  $\leftarrow AB$

$Z \rightarrow A \rightarrow B \rightarrow E \rightarrow F \rightarrow F \rightarrow A \rightarrow B \rightarrow B \rightarrow C \rightarrow D \rightarrow D \rightarrow A$

BY  $\leftarrow A \bar{B}$

These correspond to two disjoint directed cycles joined by node ABD. Edges left to right

cycle 1: A-B-D-C-A  
cycle 2: A-B-D-C-A  
start - bro1 - bro1 - bro1 - O

### ■ Repeat

\* repeats often fail assembly. They certainly

fail with its "shortest" criterion.

\* Reads might be too short to "overlap".

smooth repetitive sequences. This is why sequencing vendors try to increase read length.

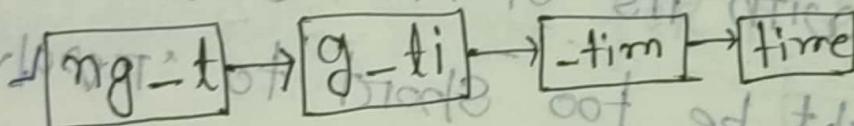
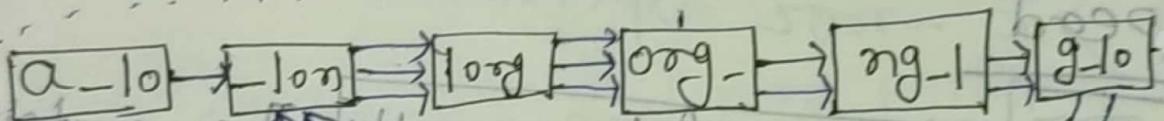
Algorithms that don't pay attention to repeats might collapse them.

For  $z_i$ ,  $P(z_i)$  might collapse.

Human genome is 50% repetitive

↳

Compute the De Bruijn graph for the following string with  $k=5$ .  
alpha-long-long-long-time



graph is connected.

Connected components are individual

graph is not.

Eulerian

eulerian paths/walks

AB → BA → AB → CD → DA → AB → BA → EF → FA → AG → BE → CD → DA

G1 +, AB

AB → DE → EF → FA → AG → BE → CD → DA

## Genome Rearrangement

Process by which the order of genes, DNA segments, or even entire chromosomes in an organism's genome are altered or rearranged (gold biologs)

→ Turnip and Cabbage share a common ancestry

→ Human and mouse share a common

- they share the same genes, but in different order

A series of rearrangements transforms one genome

most of them others

A breed of mice had similar symptoms

caused by the same type of gene or mutation

■ Mouse and Human Genome Rearrangement

from mouse genome to human genome

■ Genome Rearrangement:

Mouse and human chromosomes represented as eleven colored, dissected segments (synteny block) by STO

↳ Most human genes have a mouse counterpart

↳ But in different order

↳ Hundreds of genes often appear one after another

but may have different order in different species

# Transforming Mouse genome to Human genome (Sorting by Reversal)

## Sorting by Reversals

$$\Pi_1 = 5 \ 1 \ 4 \ 3 \ 2 \ 6 \ 8 \ 9 \ 10 \quad \text{locus}$$

$\tau(1, 3)$  is in reverse

$$\Pi_2 = 1 \ 5 \ 4 \ 3 \ 2 \ 6 \ 7 \ 8 \ 9 \ 10$$

$\tau(2, 5)$

$$\Pi_3 = 1 \ 2 \ 3 \ 5 \ 4 \ 6 \ 7 \ 8 \ 9 \ 10$$

around in between

## Common Rearrangement Practical implication:

- \* Mouse provides insight into Human Genetics  
Disorders

Wooldeaburg's Syndrome

- ↳ Gene implicated in the disease was linked to human chromosome 2q11.2 and also

- ↳ A breed of mice had similar symptoms

caused by the same type of gene  
as in human

- ↳ Scientist succeeded in identifying location → genes responsible for disorder in mice
- ↳ Finding the gene in mice gives clues to where the same gene is located in humans

⊕ GWAS (Genome Wide Association Studies)

- # Genetic Association:  
genotype → phenotype
- \* Attempts to detect how genotype affects phenotype in population

- # Key components of association studies:
  - ↳ measure genetic variation
  - ↳ measure phenotype variation
  - ↳ Quantify the association between genetic and phenotypic variations

SNP: Single Nucleotide Polymorphism  
 Located in the genome of multiple organisms or cells, etc.

### Heritability:

Extent to which a trait is predictably passed from generation to generation

\* Phenotypic variation = genetic + environmental

<ul style="list-style-type: none"> <li>→ Down's Syndrome</li> <li>→ Huntington's Disease</li> </ul>	<ul style="list-style-type: none"> <li>→ Heart disease</li> <li>→ Height</li> </ul>
<ul style="list-style-type: none"> <li>→ 3<sup>rd</sup> copy of Chromosome 21</li> </ul>	

→ Consider potential confounders

Case and Control design

take 200 people genome sequence

Case: 100 people have disease/trait

Control: 100 " " no "

\* Other traits are almost same  
(weight, age, height)

\* notified the different portion between  
two group of people genome

Confounding factors if other traits are not

similar, then one ambiguity will create

Ex:

Case: people of Bangladesh have Diabetes

Control: " " of America have no "

\* traits are genetic / environmental  
Some

\* GWAS is used for genetic variant

to bring  
IS

Manhattan Plot → Association testing  
↳ Association of each SNP with <sup>p-value</sup> phenotype  
↳ Statistical test (P-value)

→ Used in any GWAS study

→ take SNP of each chromosome

→ find out the likelihood against each SNP

→ SNP that cross a threshold are responsible <sup>Association Studies</sup>

SNP: Single Nucleotide Polymorphism

[Portion of genome] that varies between two (people) agents (A<sub>1</sub>, A<sub>2</sub>)

→ find the corresponding gene against SNP

A DNA sequence variation that occurs when a single nucleotide (A, T, C, G) in genome sequence is altered and the particular alteration is present in 1% of the population.

## BLAST (Basic Local Alignment Search Tool)

- ↳ find regions of similarity between biological sequences
- ↳ Most widely used bioinformatics programs
- ↳ Replace Smith-Waterman full alignment algorithm with faster heuristics for local alignments

## Alignment tool      Gene tree estimation

- SATE
- RAxML
- PASTA
- FastTree
- MAFFT
- MUSCLE
- Bowtie

\* Species tree from gene trees: And both H

- WQFM
- ASTRAL (2019, highly used)
- MP-EST (Highly used before 15 years ago)
- BEAST (Can outperform ASTRAL, Not significant scale)
- BUCKY
- Combined Analysis

■ DNA storage/DNA digital data storage:

※ Refers to the process of storing digital data in the base sequence of DNA

• Next generation storage

Main idea for biostorage

↳ convert file → into binary data

↳ binary data → DNA sequence

↳ DNA isotheresis (DNA sequence → biological material)

↳ store (DNA)

- ↳ read DNA seq by using sequencer
- ↳ Mapping to binary data
  - ↳ find original data for DNA synthesis

### Confounding variable

- A variable that influences both the dependent variable and independent variable, causing a spurious association.
- \* Confounding is a causal concept and such can't be described in terms of correlations or associations.
- \* They confound the true relationship between two variables.

## Phylogenomics practical challenges

Improved species tree estimation in the presence of

→ Gene tree discordance

→ Missing data

→ Poor signal per gene

→ Large dataset

\* Even the concatenated analysis provides the advantage of high level of signal.

\* Combined analysis → Statistically inconsistent

\* Summary method → generate inaccurate species tree if gene tree is "

↳ small enough genome → inaccurate gene tree

↳ discordance created by model misspecification

↳ poor signal per gene

Sol<sup>n</sup>: Removing short sequenced

Distorting

the true gene tree distribution

## # Species tree estimation methods

\* Statically inconsistent using bayesian dating to estimate

↳ CA

↳ Phylogenetic gene tree

↳ MRP

↳ GC

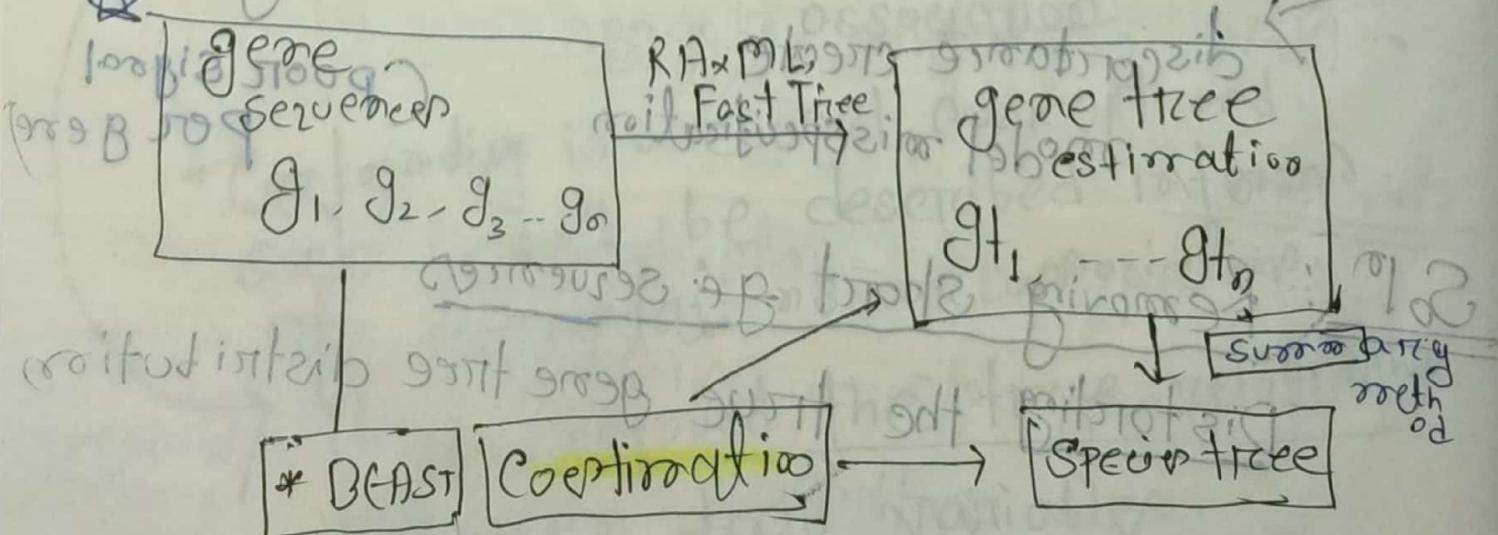
\* Statically consistent

↳ \* BEAST (Coestimation)

↳ MP-EST

↳ BUCKY-Pop

⊕ Why \*BEAST (so good)



## Gene tree accuracy

\* BEAST produces much more accurate gene trees, than the most popular likelihood methods: (RAxML, FastTree)

accuracy > accuracy

\* BEAST > RAxML, FastTree

\* Even the simplest methods like greedy consensus produce good species tree as \*BEAST given the gene trees estimated by \*BEAST

NB: / finding:

Highly accurate species tree of \*BEAST is largely due to the highly accurate gene tree:



edge of taxon

Color graph

taxon level of discordance

## Problem:

Combined analysis → Specifically inconsistent

\* BEASTOM soft → Not sealable

Summary method → poor signal

81

Combining the strengths

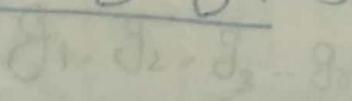
## Combined analysis

\* Evolution of Summary methods  
\* Approaches

↳ Nine binning

→ Stoptical  
T2A381

↳ Weighted statistical binning



Naive

Naive binning

Avian Genomes

Avian tree

→ Let we have  $n$  numbers of genes

→ take  $m$  numbers of gene information randomly

→ Combine them

→ generate super gene trees

→ generate species tree

\* Naive binning improves the summary methods substantially

Problems

→ Model violation (Combining gene randomly)

Sol: Statistical binning

\* Combining genes that are "similar"

→ make a graph

→ vertex → each gene tree by using incompatibility graph

→ edge → if two vertex have certain level of discordance

Color graph

- \* Use graph colors
- \* two types of bin
- \* # bin = # colors
- balanced bin size
- unbalanced binning
- \* unbalanced bins are problematic
- \* large bin size
- \* small bin

bottom would not provide increased significance

3 approaches

(1) Binning

Theorem:

Phylogenomics pipeline using statistical

bining and coalescent-based summary

method is not statistically consistent

GOALS

other outgroups

genomic features

Statistical binning + coalescent-based summary method

weighted " " " + statistically consistent  
 not statistically consistent

## Coalescent-based summary method:

- \* MP-EST
- \* ASTRAL
- \* STAR
- \* NEST
- \*\* BEAST (Coalescent)

## Sequencing technology

Technology	Read length	Throughput	Cost
Illumina	36 - 150 bp	~1 Gb/h	~\$1000
Solexa	36 - 150 bp	~1 Gb/h	~\$1000
Pyrosequencing	36 - 150 bp	~1 Gb/h	~\$1000

April 6, 2021

Suppose you are trying to construct a species tree on 15 different species.

You have sampled 300 gene sequences from each of these 15 species.

Your supervisor has asked you to

use a method called

GT-est for constructing trees from sequence alignment and SP-est (which is a summary method) for estimating species trees from gene trees.

i) How many times do you need to run

GT-est and SP-est to estimate a species tree by summarizing gene trees?

ii) How many times do you need to run

GT-est and SP-est to estimate a species tree by if you use "main binning" with 10 bins?

- iii) How many times do you need to run GT-est and SP-est to estimate a species tree if you use "Statistical Binning (SB)" assuming that SB returns 12 bins!  $\uparrow : t_{95-98}$
- iv) How many times do you need to run GT-est and SP-est to estimate a species tree if you use "Weighted Statistical Binning (WSB)" assuming that WSB returns 12 bins!
- v) How many times do you need to run GT-est and SP-est to estimate a species tree using combined analysis!  $\uparrow : t_{95-98}$

Ans: + Given 300 gene

i 300 gene  
15 species  
Species tree by summarizing gene tree

ii GT-est: 300 " (B) binning

SP-est: 1 (mid SI mutation)

iii Species tree by using "Naive Binning"

with 10 bins.

GT-est: 10

SP-est: 1

iii Species tree by using "Statistical Binning"  
assuming that SB returns 12 bins

GT-est: 12

SP-est: 1

iv Species tree by using "Weighted Statistical  
Binning" assuming that WSB returns 12  
bins:

GT-est: 300

SP-est: 1

v Species tree by using Combined  
analysis

GT-est: 1

SP-est: 0