

# TOWARDS THE OPTIMIZATION OF LARGE-SCALE PHYLOGENETIC TREES

**CÁTIA VAZ and ALEXANDRE FRANCISCO**

**INESC-ID;**

**Instituto Superior de Lisboa, Instituto Politécnico de Lisboa**

**COMPSTAT 2022**

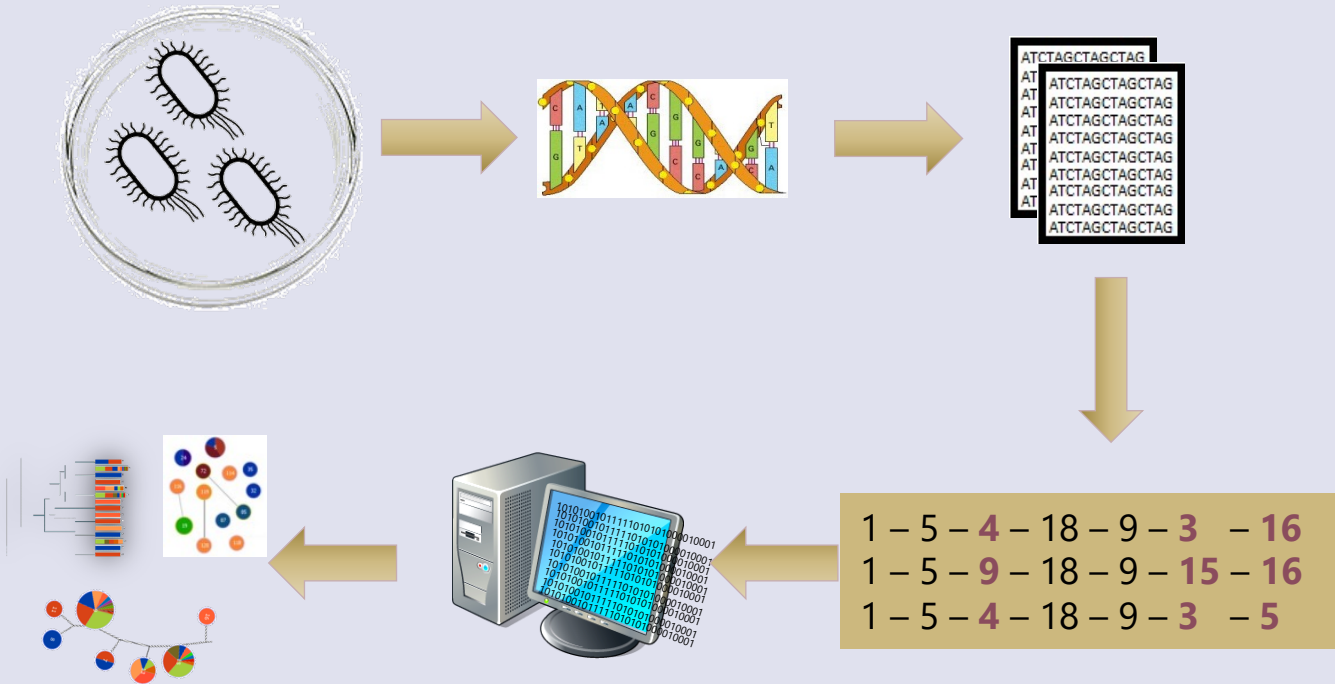
# PHYLOGENETIC INFERENCE

**PHYLOGENETICS** is the study of the evolutionary history and relationships among individuals or groups of organisms (e.g.. species, or populations) by grouping them based on some ***(dis)similarity criterion*** that underlies some ***evolution model***.

- The result is a phylogenetic tree or network

Fundamental in epidemiological and genetic studies for the surveillance of infectious diseases and about its evolution.

# PHYLOGENETIC INFERENCE



The evolutionary relationships between different species or taxa are usually inferred through known phylogenetic analysis techniques

- **Sequencing** the isolate data
- **Assembly** the sequences, comparing the draft genome to a database of gene alleles
- **Profiling** - Given the assembly results we can create an allelic profile characterizing the strain
  - Usually abstracted to categorical indexes.

# Phylogenetic tree

- After profiling, the reconstruction of a phylogenetic tree can be done by several different inference methods or algorithms.
- Most of them rely on applying clustering techniques
  - the main difference among them resides on how is defined **cluster proximity (similarity)** and on which **optimization criterion** is used.
  - Both *cluster proximity (similarity)* and *optimization criteria* rely often on a **model of evolution**
- Phylogenetic trees can be mainly built using either **distance-based methods** or **character state methods**.
- **Distance-based methods** infer the relationship between individuals as the number of genetic differences between pairs of sequences

# Cluster Proximity/Similarity/Dissimilarity

- Similarities are described according to several models like connectivity models, centroid models, distribution models...
- Distance based methods are based on *connectivity models*.
- **Connectivity models** define the similarity between elements as their distance
  - elements are more similar with nearby elements than with elements farther away)

Algorithms that implement this model differ from one another by the way distances are computed.

	A	B	C
A	-	17	11
B	17	-	28
C	11	28	-

Hamming Distance

A	1	5	4	18	9	3	16
B	1	5	9	18	9	15	16
C	1	5	4	18	9	3	5

Profile data

# Common PHYLOGENETIC Analysis workflow

**After Profiling**, the phylogenetic analysis workflow using *distance based methods* can be summarized into four consecutive steps:

- distance calculation,
  - distance correction,
  - inference algorithm,
  - local optimization steps.
- 
- **Distance Calculation:** e.g. Hamming (normalized or unnormalized), Euclidian, etc;
  - **Distance Correction:** based on some model of evolution, e.g.: Jukes-Cantor, Kimura, Felsenstein, etc.
  - **Inference algorithms** and **local optimization** steps uses different different principles/criteria e.g.: Minimum evolution, minimum deviation, etc.

# Inference algorithms/methods

Different methods can built different trees.

The strategy that each method uses to find the best tree topology can be classified as:

- exhaustive search methods
- completely bifurcating tree search methods
- stepwise clustering method

The phylogenetic tree produced by an inference method is not unique (Backeljau et al, 1996; Teixeira et al, 2015), depending on the inference method or even on the dataset input order.

Inferred trees might also not necessarily represent the best tree for the underlying evolution model

# goeBURST algorithm

goeBURST follows a graph theoretic approach to better estimate a true phylogeny, reflecting the principle of minimum evolution.

The problem is similar to the **minimum spanning tree problem**, trying to identify a subset of relationships ("links") that connects the profiles or OTUs ("points") without any cycles and with the minimum possible total link length.

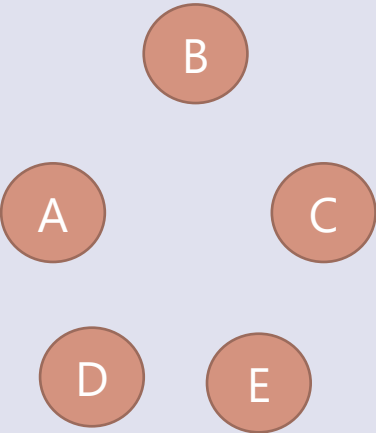
It is a stepwise algorithm, which interactively process:

- The **cluster-pair selection**: selects links connecting two OTUs with the minimum distance.
  - **Tie break rules**: select the links between OTUs with higher number of:
    1. SLVs
    2. DLVs
    3. TLVs
    4. OTU's frequency
    5. OTU's ID
- **Cluster-pair joining**: Join those two OTUs.



# goeBURST algorithm

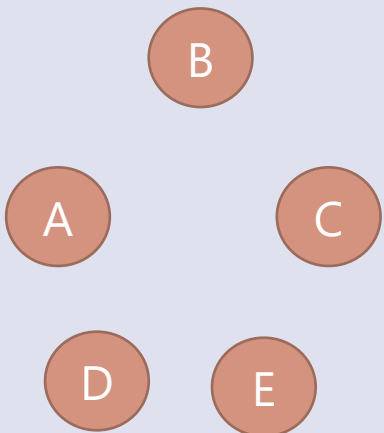
Profile	A	B	C	D	E
A	-	2	7	7	6
B	-	-	7	7	6
C	-	-	-	5	5
D	-	-	-	-	3
E	-	-	-	-	-



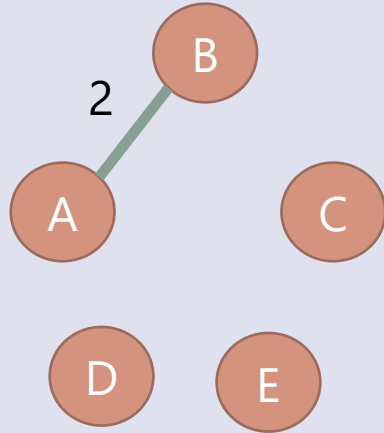
LEVEL = 1 (SLV)

# goeBURST algorithm

Profile	A	B	C	D	E
A	-	2	7	7	6
B	-	-	7	7	6
C	-	-	-	5	5
D	-	-	-	-	3
E	-	-	-	-	-



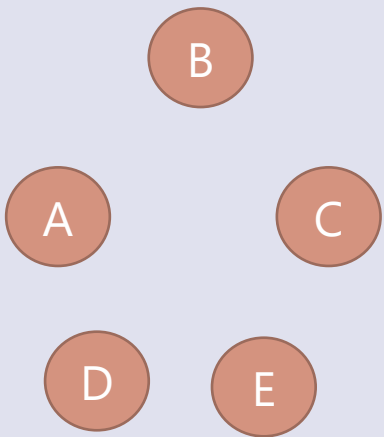
LEVEL = 1 (SLV)



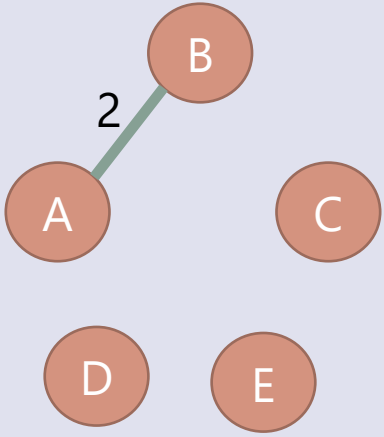
LEVEL = 2 (DLV)

# goeBURST algorithm

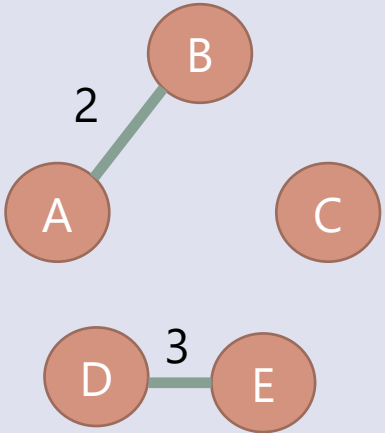
Profile	A	B	C	D	E
A	-	2	7	7	6
B	-	-	7	7	6
C	-	-	-	5	5
D	-	-	-	-	3
E	-	-	-	-	-



LEVEL = 1 (SLV)



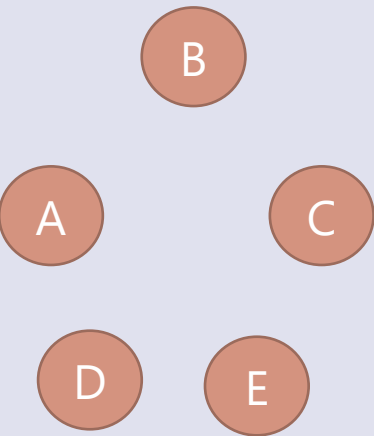
LEVEL = 2 (DLV)



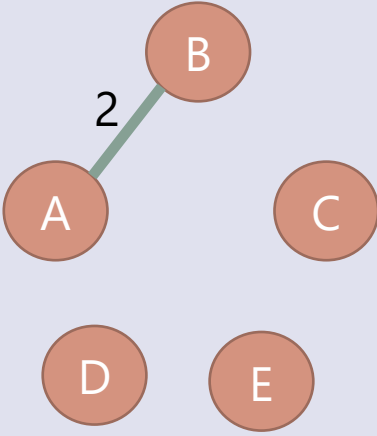
LEVEL = 3 (TLV)

# goeBURST algorithm

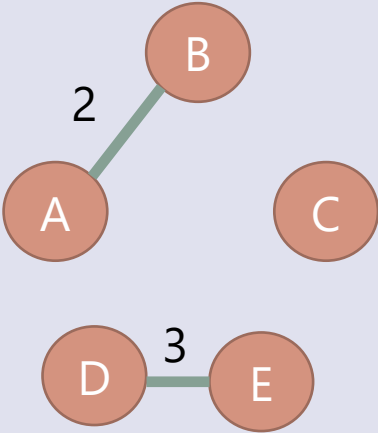
Profile	A	B	C	D	E
A	-	2	7	7	6
B	-	-	7	7	6
C	-	-	-	5	5
D	-	-	-	-	3
E	-	-	-	-	-



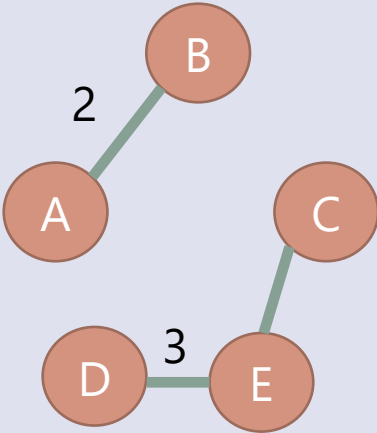
LEVEL = 1 (SLV)



LEVEL = 2 (DLV)



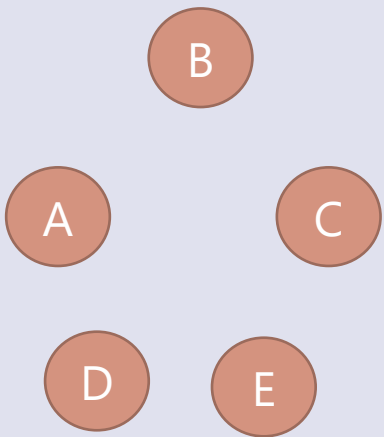
LEVEL = 3 (TLV)



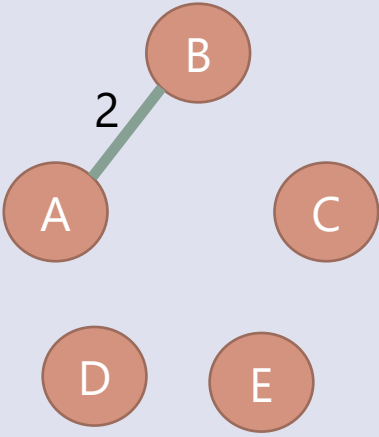
LEVEL = 5,  
assuming #E>#D

# goeBURST algorithm

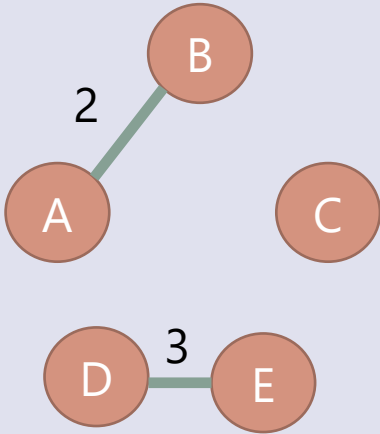
Profile	A	B	C	D	E
A	-	2	7	7	6
B	-	-	7	7	6
C	-	-	-	5	5
D	-	-	-	-	3
E	-	-	-	-	-



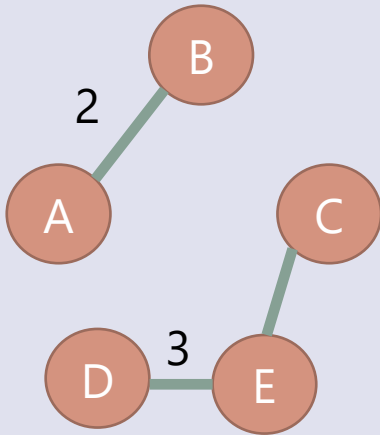
LEVEL = 1 (SLV)



LEVEL = 2 (DLV)



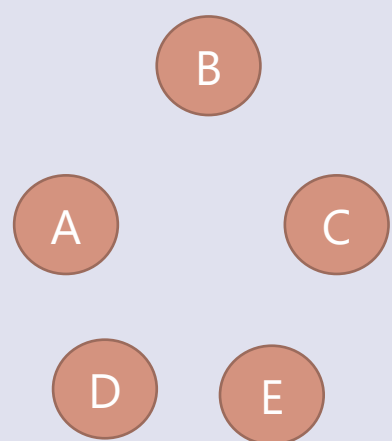
LEVEL = 3 (TLV)



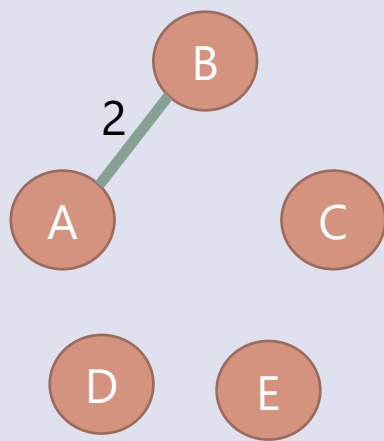
LEVEL = 5,  
assuming #E > #D

# goeBURST algorithm

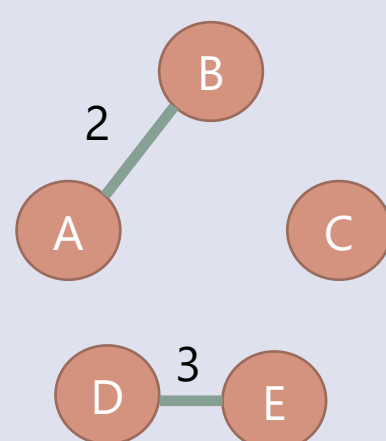
Profile	A	B	C	D	E
A	-	2	7	7	6
B	-	-	7	7	6
C	-	-	-	5	5
D	-	-	-	-	3
E	-	-	-	-	-



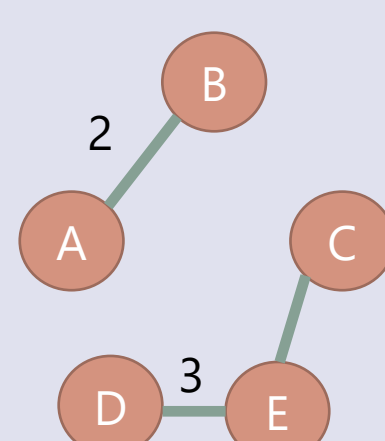
LEVEL = 1 (SLV)



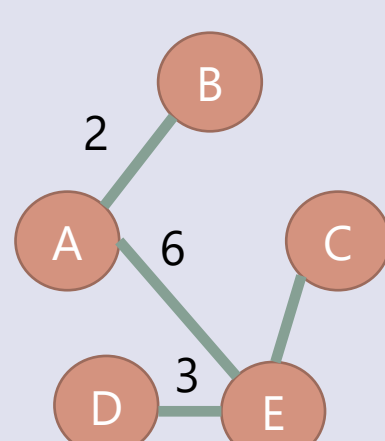
LEVEL = 2 (DLV)



LEVEL = 3 (TLV)



LEVEL = 5,  
assuming #E>#D



LEVEL = 6, assuming  
ID tiebreak

# Local Branch Recrafting

The phylogenetic tree does not necessarily represent true phylogenetic relationships between sequences

- allelic distances do not always correlate with divergence time.

To correct, we may apply a **Local Branch Recrafting (LBR)**

**Input:** A phylogenetic tree **T** over a set of elements **S**.

**Initialization:** Initialize the set **E** with the edges of the tree **T**.

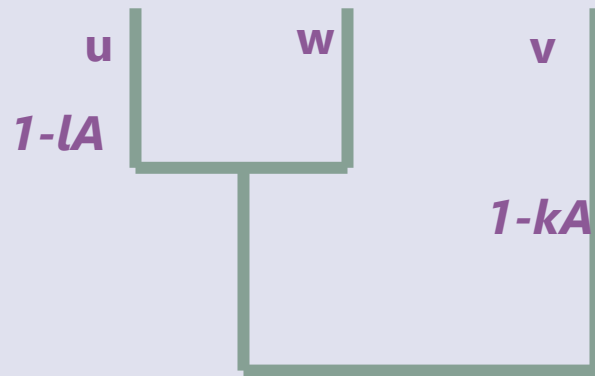
**Loop:** While  $|E| > 0$  do:

1. **Selection:** Select an edge  $(u \rightarrow v)$  of the set **E** and remove it from the tree **T**, dividing it into two sub-trees **T<sub>u</sub>** (containing **u**) and **T<sub>v</sub>** (containing **v**).
  2. **Joining:** Find two vertices **w** and **z** that best connect the two sub-trees by an edge  $(w \rightarrow z)$ .
  3. **Reduction:** Remove the edge  $(u \rightarrow v)$  from **E** and add the edge  $(w \rightarrow z)$  to **E**.
1. **Finalization:** Return the tree **T**.

# Local Branch Recrafting – Best connection?

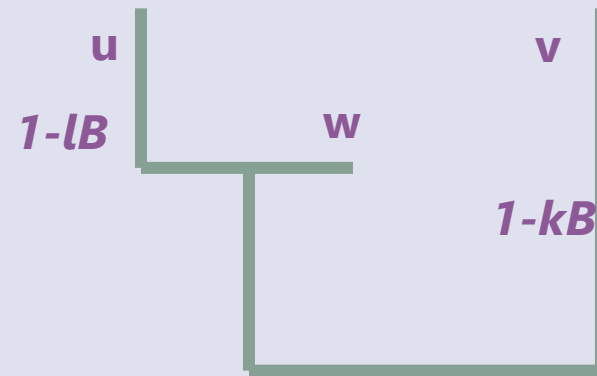
1. **Joining:** For each edge  $(u \rightarrow v)$  Find two vertices  $w$  and  $z$  that **best connect** the two sub-trees  $T_u$  and  $T_v$  by an edge  $(w \rightarrow z)$ .

**Contemporary model (MA):** nodes  $u$  and  $w$  are contemporary sisters that diverged from a hypothetical common ancestor.



$l_A \times l_A$  -> invariant alleles proportion/probability between  $u$  and  $w$ ;  
 $k_A \times k_A$  -> invariant alleles proportion between  $(u,w)$  and  $v$

**Ancestor-Descendent model (MB):** node  $w$  is the direct ancestor of  $u$ .



$l_B \times l_B$  -> invariant alleles proportion between  $u$  and  $w$ ;  
 $k_B \times k_B$  -> invariant alleles proportion between  $(u,w)$  and  $v$



# Local Branch Recrafting – Best connection?

- $1 - I_{Ax|A}$  -> **variant** alleles proportion/probability between u and w;
- $1 - I_{AxKA}$  -> **variant** alleles proportion/probability between (u,w) and v;

## Best connection?

- **Step 1:** Calculate the likelihoods of a contemporary model versus an ancestor-descendent model

Given L, the set of distinct alleles

$$\underset{0 \leq I_A \leq 1, 0 \leq k_A \leq 1}{\operatorname{argmax}} \log P(MA \mid I_A, k_A) = \underset{0 \leq I_A \leq 1, 0 \leq k_A \leq 1}{\operatorname{argmax}} \log P(u \rightarrow w \mid I_A) P(u \rightarrow v \mid I_A, k_A) P(w \rightarrow v \mid I_A, k_A)$$

# Local Branch Recrafting – Best connection?

- $(1-IA)xIA$  -> **variant** alleles probability between u and w;
- $(1-IA)xkA$  -> **variant** alleles propability between (u,w) and v;

- L, the set of distinct alleles
- $\log(a*b) = \log(a) + \log(b)$
- Distances are normalized

$$\underset{0 \leq IA \leq 1, 0 \leq kA \leq 1}{\operatorname{argmax}} \log P(MA \mid IA, kA) = \underset{0 \leq IA \leq 1, 0 \leq kA \leq 1}{\operatorname{argmax}} \log P(u \rightarrow w \mid IA) P(u \rightarrow v \mid IA, kA) P(w \rightarrow v \mid IA, kA)$$

$$= \underset{0 \leq IA \leq 1, 0 \leq kA \leq 1}{\operatorname{argmax}}. \log P(u \rightarrow w \mid IA) + \log P(u \rightarrow v \mid IA, kA) + \log P(w \rightarrow v \mid IA, kA)$$

The likelihood of a branch is:

(branch length) <sup>(number of different alleles)</sup> \* (1-branch length) <sup>(number of identical alleles)</sup>

Thus, for instance, branch  $u \rightarrow w$

$$|L|d(u,v) \log(1-IA^2) + |L|(1-d(u,v)) \log(1-IA^2)$$

# Local Branch Recrafting – Best connection?

$$\begin{aligned} & \operatorname{argmax}_{0 \leq |A| \leq 1, 0 \leq k_A \leq 1} \log P(MA \mid |A, k_A) = \dots \\ &= \operatorname{argmax}_{0 \leq |A| \leq 1, 0 \leq k_A \leq 1} \log P(u \rightarrow w \mid |A) + \log P(u \rightarrow v \mid |A, k_A) + \log P(w \rightarrow v \mid |A, k_A) \\ &= \operatorname{argmax}_{0 \leq |A| \leq 1, 0 \leq k_A \leq 1} |L|d(u \rightarrow w) \log(1 - |A|^2) + |L|(1 - d(u \rightarrow w)) \log(|A|^2) + \\ & \quad |L|d(u \rightarrow v) \log(1 - |A|k_A) + |L|(1 - d(u \rightarrow v)) \log(|A|k_A) + \\ & \quad |L|d(w \rightarrow v) \log(1 - |A|k_A) + |L|(1 - d(w \rightarrow v)) \log(|A|k_A) \end{aligned}$$

Similarly, in Model B

$$\begin{aligned} & \operatorname{argmax}_{0 \leq |B| \leq 1, 0 \leq k_B \leq 1} \log P(MB \mid |B, k_B) = \operatorname{argmax}_{0 \leq |B| \leq 1, 0 \leq k_B \leq 1} \log P(w \rightarrow u \mid |B)P(u \rightarrow v \mid |B, k_B)P(w \rightarrow v \mid |B, k_B) \\ &= \operatorname{argmax}_{0 \leq |B| \leq 1, 0 \leq k_B \leq 1} |L|d(w \rightarrow u) \log(1 - |B|) + |L|(1 - d(w \rightarrow u)) \log(|B|) + \\ & \quad |L|d(u \rightarrow v) \log(1 - |B|k_B) + |L|(1 - d(u \rightarrow v)) \log(|B|k_B) + \\ & \quad |L|d(w \rightarrow v) \log(1 - k_B) + |L|(1 - d(w \rightarrow v)) \log(k_B) \end{aligned}$$

# Local Branch Recrafting – Best connection?

The solution to parameters is:

$$I_A = \sqrt{1 - d(u \rightarrow w)}$$

$$k_A = \frac{1 - (1/2)(d(u \rightarrow v) + d(w \rightarrow v))}{I_A}$$

$$I_B = 1 + \frac{x d(w \rightarrow u)}{d(u \rightarrow v) - 2x}$$

$$k_B = 1 + \frac{x d(w \rightarrow v)}{d(u \rightarrow v) - 2x}$$

**with**

$$x = 1 - \frac{(1 - d(w \rightarrow u))(1 - d(w \rightarrow v)) + (1 - d(u \rightarrow v))}{2}$$

## Best connection?

- **Step 1:** Calculate the likelihoods of a contemporary model versus an ancestor-descendent model
- **Step 2:** the joining criterion consists of finding two nodes that have the minimum harmonic distance if the contemporary model has a higher or equal likelihood to the ancestor-descendent model

# Best connection

**Input:** Initial edge  $(\mathbf{u} \rightarrow \mathbf{v})$  , Tree  $\mathbf{t}(\mathbf{u}) \in \mathbf{F}$ , harmonic tiebreaker  $\mathbf{ht}$

Output: New edge  $(\mathbf{u}' \rightarrow \mathbf{v}')$

1: Initialize  $\mathbf{u}' = \mathbf{u}$

2: for each node  $\mathbf{w} \in \mathbf{TargetNodes}(\mathbf{t}(\mathbf{u}'))$  do

3:  $\mathbf{P}(\mathbf{MA}), \mathbf{P}(\mathbf{MB}) = \mathbf{ModelSelection}(d(\mathbf{u}' \rightarrow \mathbf{w}), d(\mathbf{w} \rightarrow \mathbf{u}'), d(\mathbf{u}' \rightarrow \mathbf{v}), d(\mathbf{w} \rightarrow \mathbf{v}))$

4: if  $\mathbf{P}(\mathbf{MA}) \geq \mathbf{P}(\mathbf{MB})$  and  $\mathbf{ht}(\mathbf{u}') > \mathbf{ht}(\mathbf{w})$  then

5:  $\mathbf{u}' = \mathbf{w}$

6: else if  $\mathbf{P}(\mathbf{MA}) < \mathbf{P}(\mathbf{MB})$  and  $\mathbf{d}(\mathbf{u}' \rightarrow \mathbf{v}) \geq \mathbf{d}(\mathbf{w} \rightarrow \mathbf{v})$  then

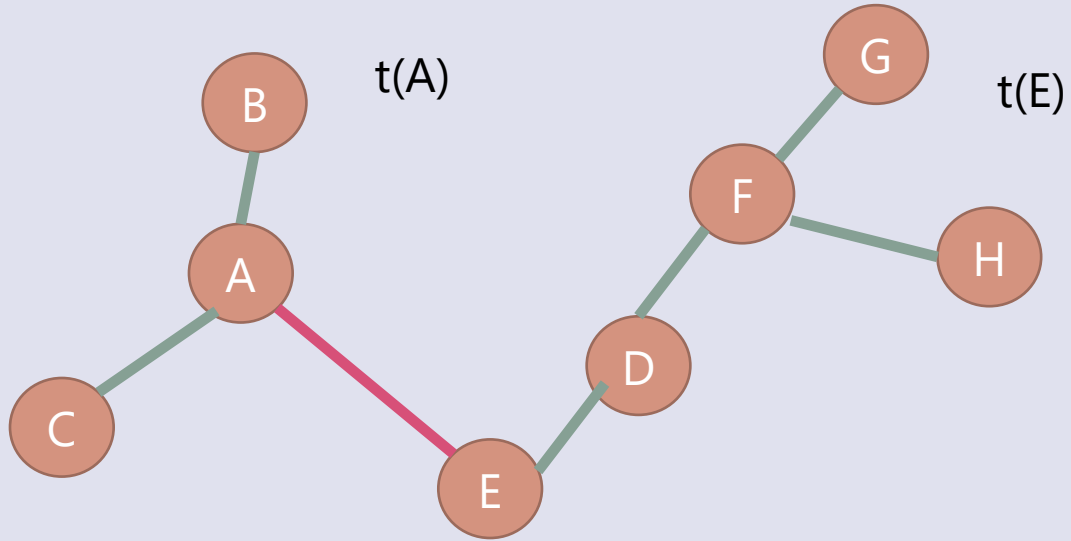
7:  $\mathbf{u}' = \mathbf{w}$

8: Repeat **(1–7)** on  $\mathbf{v}$  to obtain  $\mathbf{v}'$

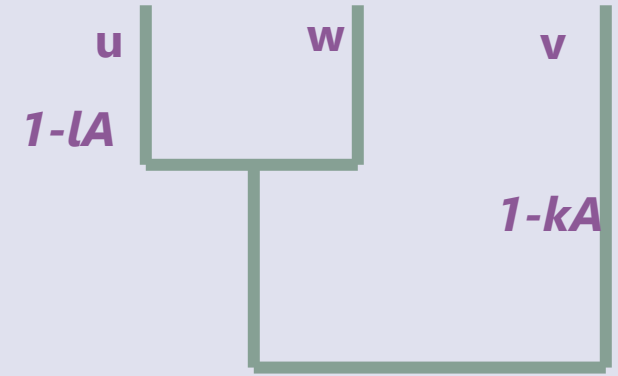
9: Return  $(\mathbf{u}' \rightarrow \mathbf{v}')$

**harmonic mean** gives less weight to vertices that are close to the vertex of interest than to those far away,

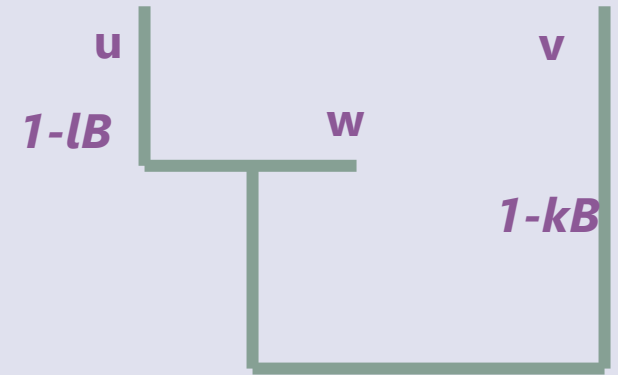
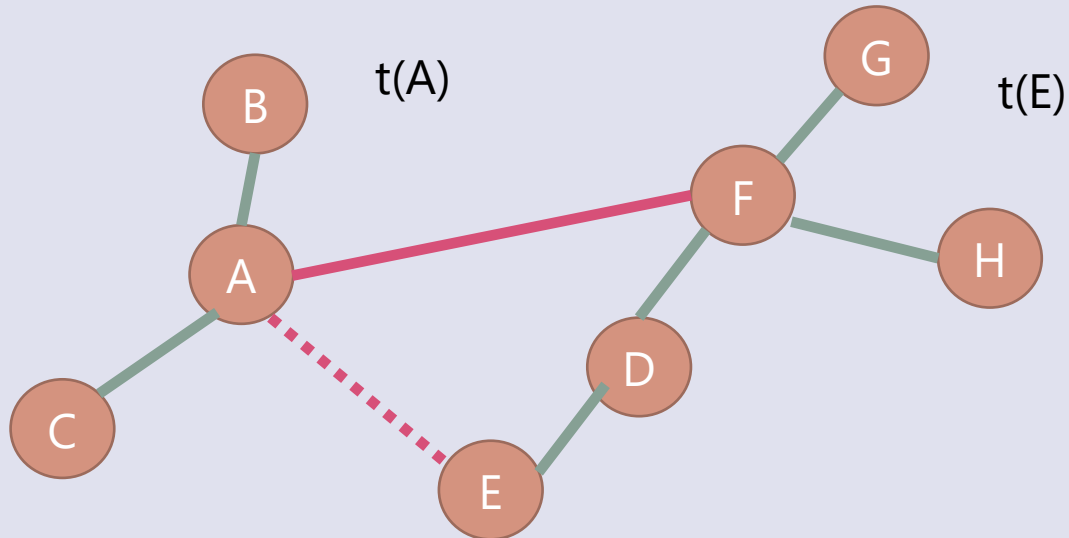
# Example



( $A \rightarrow E$ ) is compared with branch ( $A \rightarrow F$ ), where node F has the lowest harmonic average distance to other nodes and Model A has higher probability



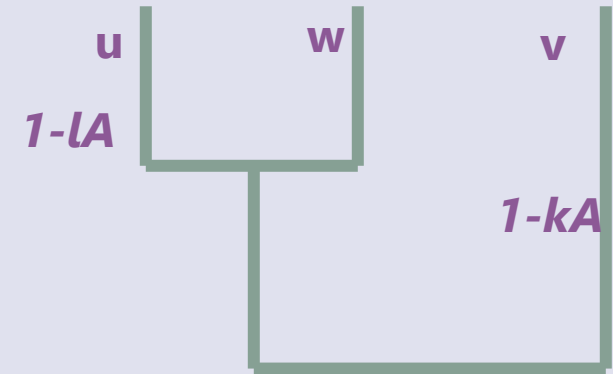
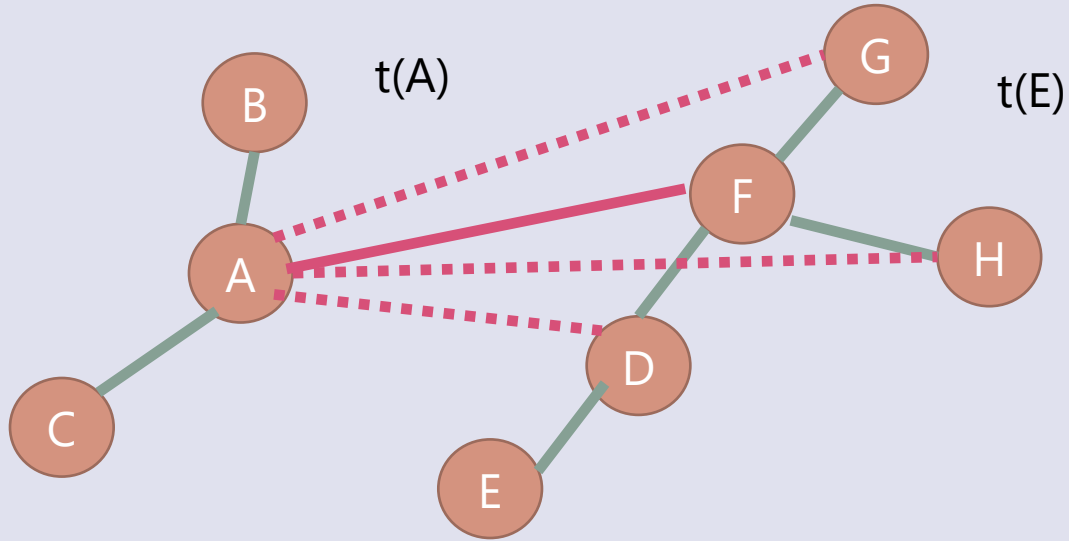
Model A



Model B

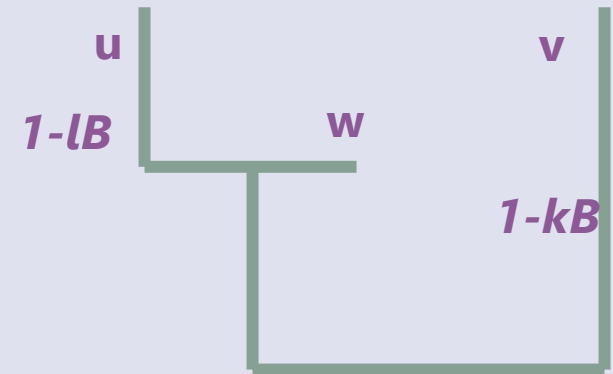
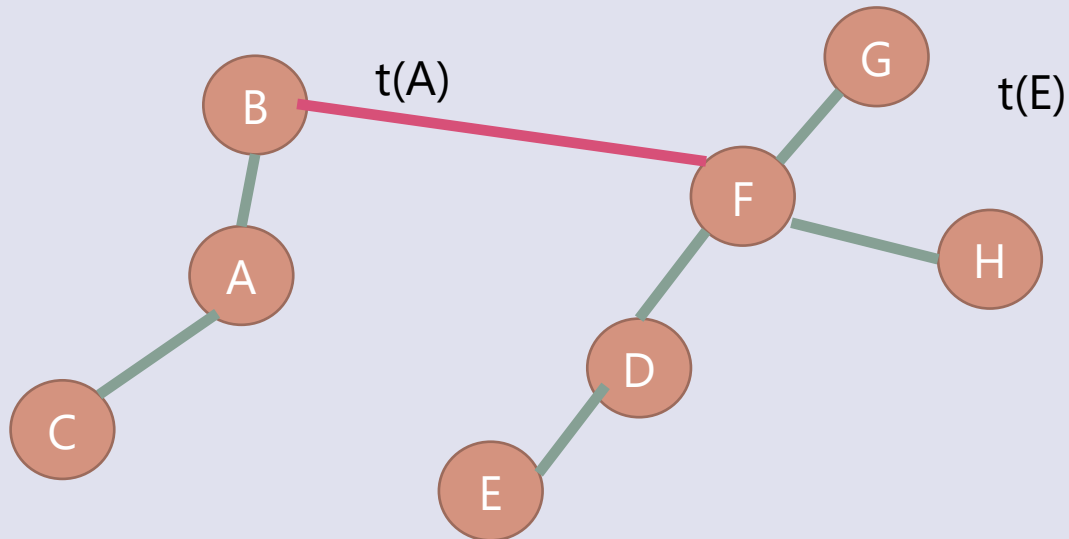
# Example

$(A \rightarrow F)$  is compared with all the nodes that are directly connected with F, but  $(A \rightarrow F)$  is still the most probable branch



Model A

The same process is performed for tree  $t(A)$ , which results in  $(B \rightarrow F)$  becoming the most probable branch.



Model B

# Final Remarks

- The proposed local optimization was implemented (in Java) in a generic way such that it can be applied to the output of goeBURST algorithm, or to the output of any other algorithm that produces a phylogenetic tree.
- In the case of goeBURST algorithm, we were able to incorporate this optimization in the comparison criteria, and it does not affect the running time of the algorithm.
- This method was tested with data from EnteroBase, and results are promising with trees being more meaningful from the biological point of view.



# Acknowledgments

1. This work was supported by:

- A research project from Polytechnic Institute of Lisbon - IPL/2021/DIVA ISEL
- A research project from Fundação para a Ciência e Tecnologia(FCT) with reference NGPHYLO PTDC/CCI-BIO/29676/2017.