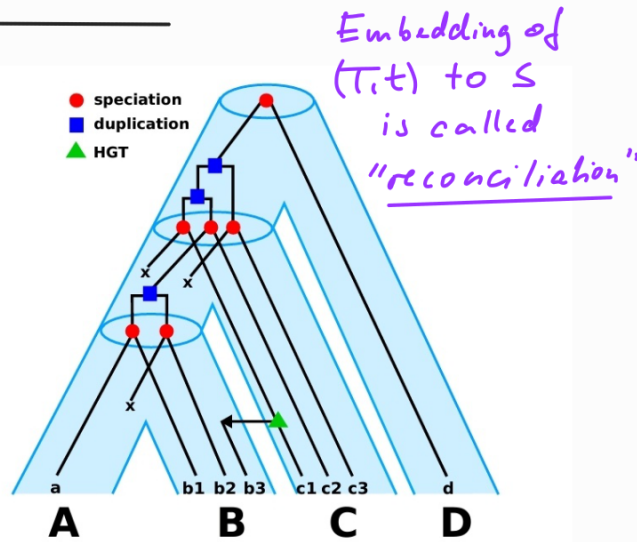


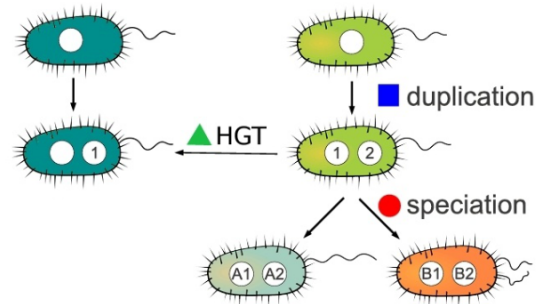
# Homology - Relations

# Detailed Evolutionary Scenarios:

- ▶ species are characterized by its genome:  
a "bag of genes"
- ▶ "Genes" evolve along a *rooted tree* with unique coloring  
 $t: V^0 \rightarrow M = \{\bullet, \blacksquare, \blacktriangle\} \cong (T, t)$
- ▶ "x" = gene loss



- **Gene duplication** : an offspring has two copies of a single gene of its ancestor
- **Speciation** : two offspring species inherit the entire genome of their common ancestor
- ▲ **HGT** : transfer of genes between organisms in a manner other than traditional reproduction and across different species



- Two genes are homologous if they share common ancestor.
- **Homology-Relations** are binary relations between genes.  
[are important: gene function, find new genes, mechanism that act on genes]

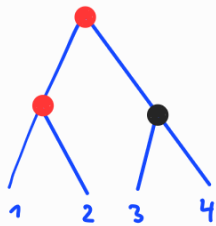
Here we investigate in more detail the structure of 2 Homology relations: ORTHOLOGY  
XENOLOGY

# ORTHOLOGY

- Two genes  $x, y$  are orthologous if they were separated by speciation event.

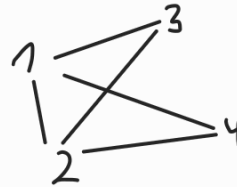
$\iff$  in given gene tree  $(T, t)$  with  $t: V \setminus L \rightarrow \{\bullet, \circ\}$   
we have  $t(\text{lca}_T(x, y)) = \bullet = \text{speciation}$ .

$R_\bullet$  = binary relation that comprises all pairs of orthologous genes.



As graph: (undirected)

$$R_\bullet = \{ (1,2), (2,1), (1,3), (3,1), (1,4), (4,1), (2,3), (3,2), (2,4), (4,2) \}.$$



## 2 classical ways to infer orthology:

### Tree-based inference

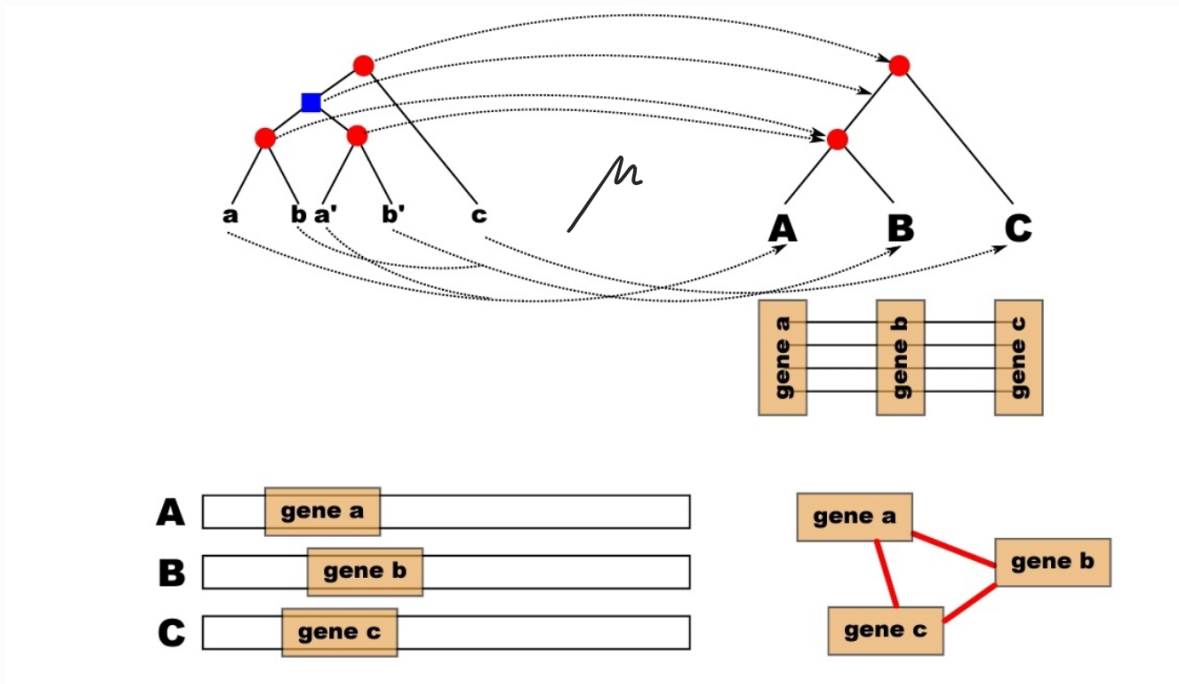
- ▶ construct gene and species trees and find reconciliation map  $\mu$  between them
- ▶ based on the placing of vertices in gene tree to species tree on infers speciation events

### Graph-based inference

Typically run in two phases:

- ▶ a **graph construction phase**, in which pairs of orthologous genes are inferred and connected by edges
- ▶ a **clustering/clean-up phase**, in which (groups of) orthologous genes are constructed/extracted based on the structure of the graph

# TREE-BASED (IDEA)



1

## Compute Species Tree:

- ▶ Find 1:1-orthologs  
= collection of genes such that from each species one gene and each gene is ortholog to all other genes in this collection
  - ▶ Select families of genes that rarely exhibit duplications (e.g. rRNAs, ribosomal proteins)
- ▶ Alignments of protein or DNA sequences and standart techniques yield gene tree with speciation-events only  
This history is believed to be congruent to that of the respective species.

## Compute Gene Tree *without events*:

- ▶ Alignments of protein or DNA sequences and standart techniques

2

## Compute Gene Tree *without events*:

- ▶ Alignments of protein or DNA sequences and standart techniques

3

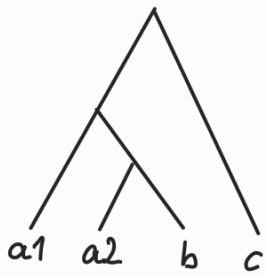
## Compute Events of Gene Tree:

- ▶ Find reconciliation map  $\mu$  w.r.t. certain optimization criteria (e.g. minimize number of losses and duplications)



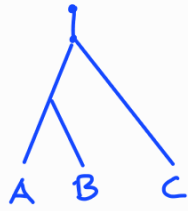
2

T (gene tree)



1

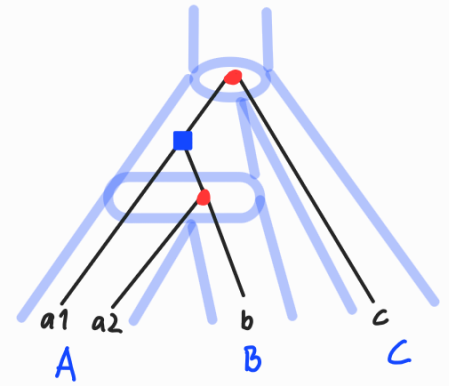
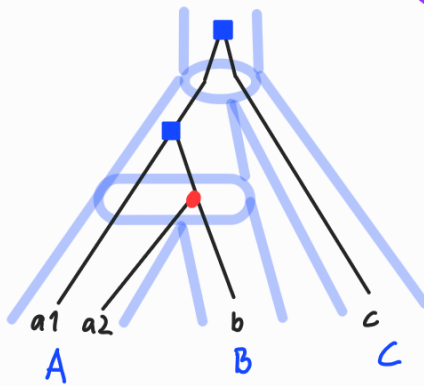
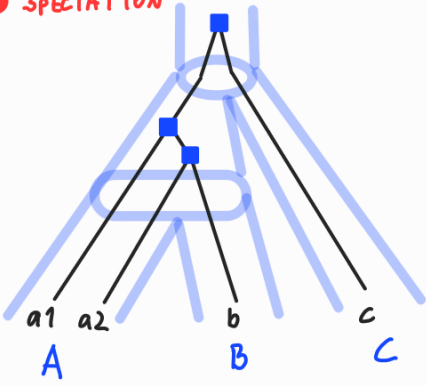
S (species tree)



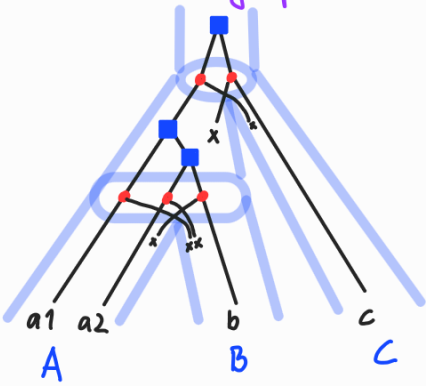
possible reconciliations

■ DUPLICATION  
● SPECIATION

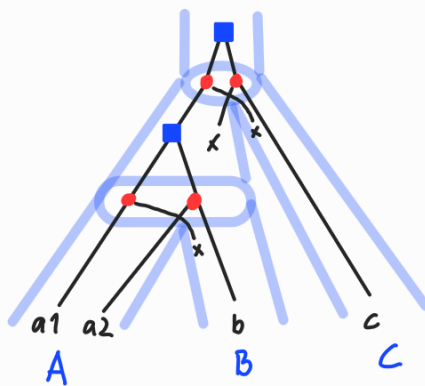
3



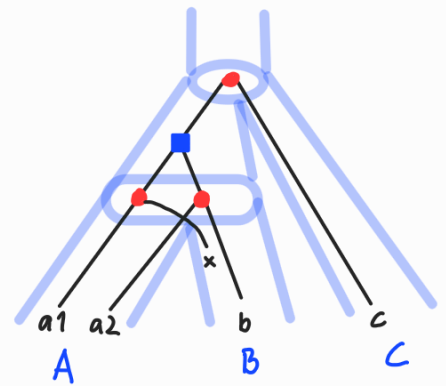
Resulting putative true evol. histories:



3 Dupl. / 5 losses



2 Dupl / 3 losses



1 Dupl / 1 loss

In practice: Find most parsimonious reconciliation, i.e., one with a min. nr of losses & duplications

## Observation (TREE-BASED)

### ► Compute Species Tree

- some orthologs must already be known!
- since only 1:1 orthologs are used, ~ 90% of the genetic sequence material remains unused

### ► Compute Gene Tree + Reconciliation

- Methods that allow to reconstruct the history of arbitrary genes rely on “restrictive” evolutionary models (e.g. event probabilities, maximum parsimony)

### This reveals a circular problem:

Reconstruction of species trees requires identifying **events** of the family evolution

Reconstruction of **event-labeled** gene trees requires a known species trees

**Accuracy** strongly depends on the predicted gene tree and the used methods (together with underlying evolutionary model) to reconcile gene and species tree.

*⇒ alternative!*

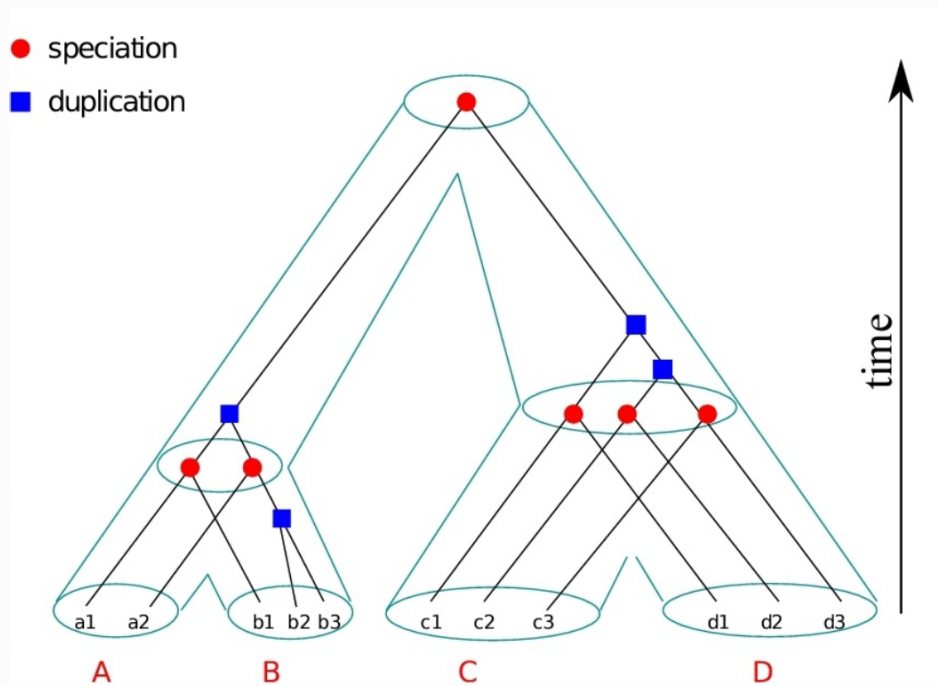
## GRAPH-BASED (IDEA)

### 2 Phases:

- 1 ► a **graph construction phase**, in which pairs of orthologous genes are inferred and connected by edges
- 2 ► a **clustering/clean-up phase**, in which (groups of) orthologous genes are constructed/extracted based on the structure of the graph

1

Assume we know true history [no HGT]



Observation:

Orthologs tend to be the homologs that diverged least. Why?

If no HGT, orthologs branched by definition at the latest possible time point—the speciation between the two genomes in question.

IDEA:

- ▶  $T$  gene tree,  $S$  species tree
- ▶  $t_S(X, Y)$  = divergence time of species  $X, Y$ .
- ▶  $y \in Y$  is orthologous to  $x \in X$ , if

- 1  $X \neq Y$ ,  
orthologs are never found in the same species
- 2  $t_T(x, y) \simeq t_S(X, Y)$ , divergence time of  $x$  and  $y$  in  $T \simeq t_S(X, Y)$ .

- ▶  $T$  gene tree,  $S$  species tree
- ▶  ~~$t_S(X, Y)$  = divergence time of species  $X, Y$ .~~
- ▶  $y \in Y$  is a candidate orth. of  $x \in X$ , if

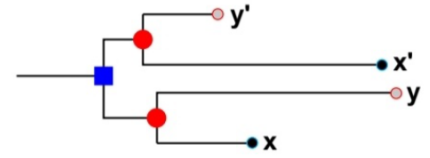
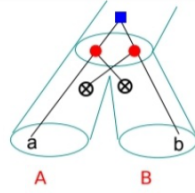
- 1  $X \neq Y$ ,  
orthologs are never found in the same species
- 2  $sim(x, y) \gtrsim sim(x, y') \forall y' \in Y$  and  
 $sim(x, y) \gtrsim sim(x', y) \forall x' \in X$ .

Don't know divergence time!  
But can measure sequence similarities

if  $x$  and  $y$  are orthologs, then they do not have (much) closer relatives in the two species.

**Not too weird mutation rates:**  $t_T(x, y) \leq t_T(x, y') \iff \text{sim}(x, y) \geq \text{sim}(x, y')$   
 “closer related, higher similarity”

This cannot work perfectly:



But we can “approximate” Orthology from such sequence similarities

AKA: (Reciprocal) Best Hits/  
Best matches.



How can we trust such estimates  $\hat{R}_\bullet$  of the true  $R_\bullet$ ?

**The least task we can do:**

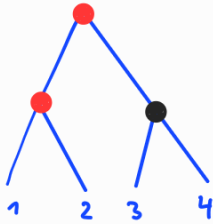
Ask for an event-labeled gene tree that supports our observation.

An estimated orthology relation  $\hat{R}_\bullet$  is **feasible** if there is a tree  $T = (V, E)$  with coloring  $t: V^0 \rightarrow \{\bullet, \bullet\}$  such that

$$t(\text{lca}_T(x, y)) = \bullet \iff (x, y) \in \hat{R}_\bullet \text{ for all distinct } x, y \in X.$$

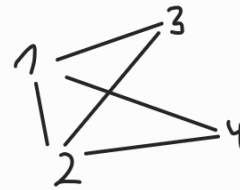
Can we mathematically characterize **feasible** estimates  $\hat{R}_\bullet$ ?

# The structure of Orthology



As graph: (undirected)

$$R_{\bullet} = \{ (1,2), (2,1), (1,3), (3,1), (1,4), (4,1), (2,3), (3,2), (2,4), (4,2) \}$$



$R_{\bullet}$  is symmetric ( $\text{lca}(xy) = \text{lca}(yx)$ )  
 but not transitive (eg. 3 & 1 are orthol.  
 1 & 4 are orthol.  
 But 3 & 4 are not orthol.)

We can represent  $R_{\bullet}$  as an undirected graph  $G = (V, E)$

$$\{x, y\} \in E \Leftrightarrow (xy), (yx) \in R_{\bullet}$$

To understand  $R_{\bullet}$  we can investigate its graph structure.



Thus we consider the equivalent problem:

Given a graph  $G$

Is there a tree  $(T, t)$  with leaf set  $L(T) = V(G)$   
 & labeling  $t: V^0 \rightarrow \{0, 1\}$  with  $V^0 := V(T) \setminus L(T)$

st

$$\{x, y\} \in E(G) \Leftrightarrow t(\text{lca}(x, y)) = 1$$

In the following we characterize feasible graphs

To this end, we need some extra notation:

DEF [complement]:  $\bar{G}$  of  $G: V(\bar{G}) = V(G)$   
 $E(\bar{G}) = \{\{x, y\} \mid \{x, y\} \notin E(G), x \neq y\}$

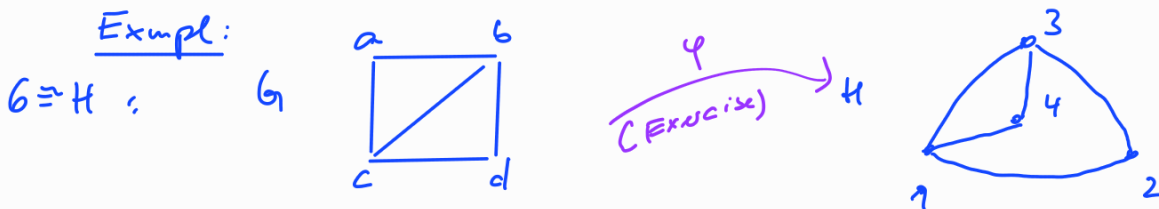
DEF (isomorphism)

2 graphs  $G, H$  are isomorphic (in symbols  $G \cong H$ )

if  $\exists$  bijective map  $\varphi: V(G) \rightarrow V(H)$  st

$$\{x, y\} \in E(G) \Leftrightarrow \{\varphi(x), \varphi(y)\} \in E(H)$$

Exmpl:



DEF [induced subgraph]

Given  $G = (V, E)$  &  $W \subseteq V$ . The subgraph  $G[W]$   
induced by  $W$  has vertex set  $W$  & edges  $\{x, y\}$   
 for all  $x, y \in W$  &  $\{x, y\} \in E$

## DEF [path]

A path  $P \subseteq G$  is subgraph of the form:



[ formal:  $V(P) = \{v_1, \dots, v_n\}$ ,  $n \geq 1$  [i.e.  $v_i \neq v_j \ \forall i \neq j$ ]  
 $E(P) = \{ \{v_i, v_{i+1}\} \mid 1 \leq i < n \}$  ]

length of path  $P = |V(P)|$

Paths of length  $n$  are denoted by  $P_n$ .

## DEF [disjoint union / join]

Given graphs  $G, H$  with  $V(G) \cap V(H) = \emptyset$ .

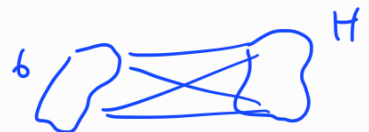
Their disjoint union  $G + H$ :



$$V(G + H) := V(G) \cup V(H)$$

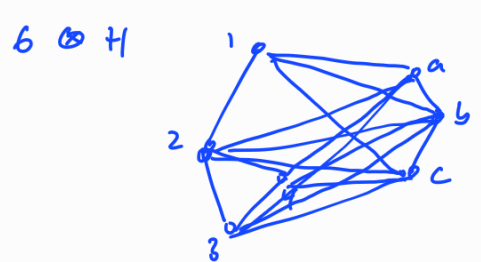
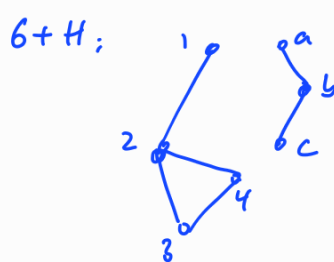
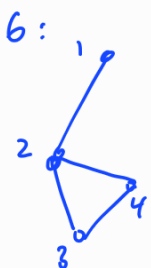
$$E(G + H) := E(G) \cup E(H)$$

Their join  $G \otimes H$ :



$$V(G \otimes H) := V(G) \cup V(H)$$

$$E(G \otimes H) := E(G + H) \cup \{ \{x, y\} : x \in V(G), y \in V(H) \}$$





DEF: Cograph recursively defined:

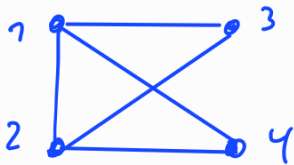
1)  $K_1$  is cograph

2) if  $b_1$  &  $b_2$  cograph

$\Rightarrow b_1 + b_2$  &  $b_1 \otimes b_2$  cograph.

[ $K_1 \cong \bullet$  single vertex graph]

Is this graph a cograph?



[Exercise start with if  $b$  cograph

$\Rightarrow b = b_1 \otimes b_2 / b_1 \cup b_2$

what is then  $b_1$  &  $b_2$ ?

recurse on  $b_1$  &  $b_2$ ]

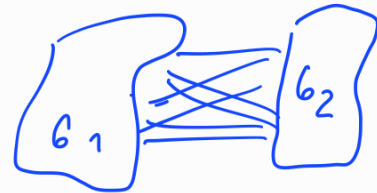
Obs. if  $b$  cograph  $\Rightarrow b \cong K_n$  or either  $b$  or  $\bar{b}$  disconnected,

$\Rightarrow b$  cograph implies  $\bar{b}$  cograph.

Why?

$b = b_1 + b_2$  disconn.  $\checkmark$

$b = b_1 \otimes b_2$



$\Rightarrow b$  connected

&

$\bar{b}$ :



$\bar{b} = \bar{b}_1 + \bar{b}_2$

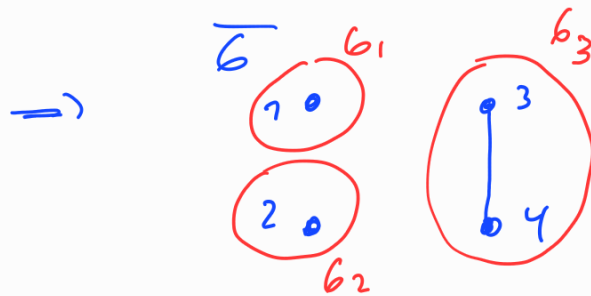
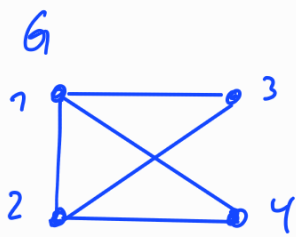
$\Rightarrow$  IF  $G$  cograph :  $G = G_1 + G_2$

or  $G = G_1 \otimes G_2 \Leftrightarrow \overline{G} = \overline{G_1 \otimes G_2} = \overline{G_1} + \overline{G_2}$

with  $G_1/G_2$  cographs.

& either  $G$  or  $\overline{G}$  disconnected.

$\Rightarrow$  check components of disconnected graph  
in  $\{G, \overline{G}\}$   
whether they are cographs!



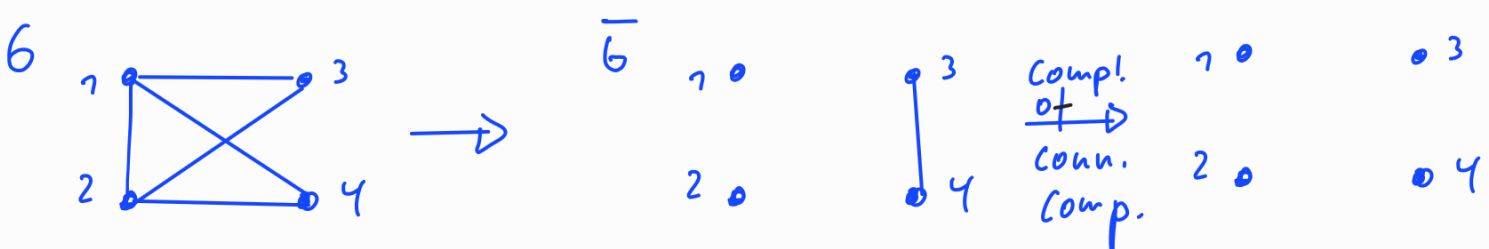
$$\overline{G} = G_1 + G_2 + G_3$$

since  $G_1 \cong G_2 \cong K_1$  they are cographs

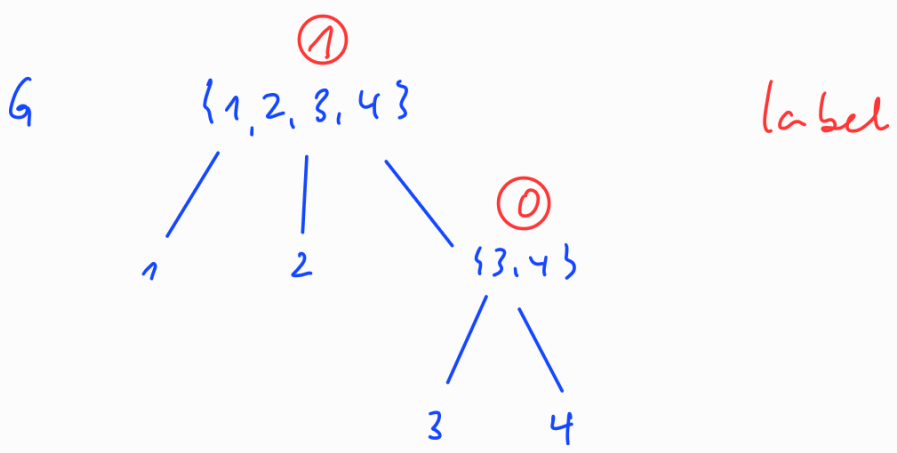
For  $G_3 = \overset{3}{\cdot} \otimes \overset{4}{\cdot}$  cograph by def.

$$\Rightarrow G = (\overset{1}{\cdot} \otimes \overset{2}{\cdot}) \otimes (\overset{3}{\cdot} + \overset{4}{\cdot})$$

Cograph = complement reducible graph,  
 i.e. stepwise complementation  
 of conn. components yield  
 complete edge-less graph.



In this way one can compute tree for cographs:



$$t(\text{LCA}(x,y)) = 1 \iff \{x,y\} \in E(G).$$

ALGO (Input cograph  $G$ )

init: Add root  $g$  to empty tree  $T$

IF ( $G$  connected)

$t(g) = 1, B = 0$ , call  $COTREE(\overline{G}, g, B)$

ELSE  $t(g) = 0, B = 1$ , call  $COTREE(G, g, B)$

COTREE( $G, v, B$ )

IF ( $|V(G)| = 1$ )

    return rooted tree " $\bullet x$ " //  $V(G) = \{x\}$

ELSE //  $G$  disconnected since cograph

    let  $C_1 \dots C_k$  conn. of  $G$

    add vertices  $v_1 \dots v_k$  to  $T$

    add edges  $\{v, v_i\}$  to  $T$ ,  $1 \leq i \leq k$ .

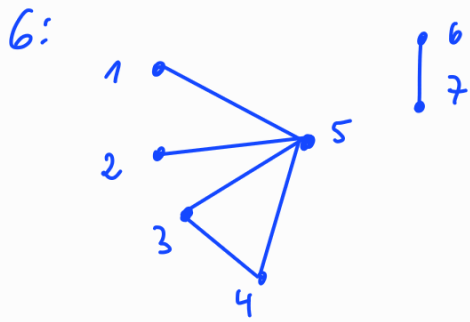
    Label  $t(v_i) = B$

    FOR ( $i = 1 \dots k$ )

        call  $COTREE(\overline{G[C_i]}, v_i, \overline{B})$  //  $\overline{B} = \begin{cases} 1, B=0 \\ 0, B=1 \end{cases}$

Exmpl:

label



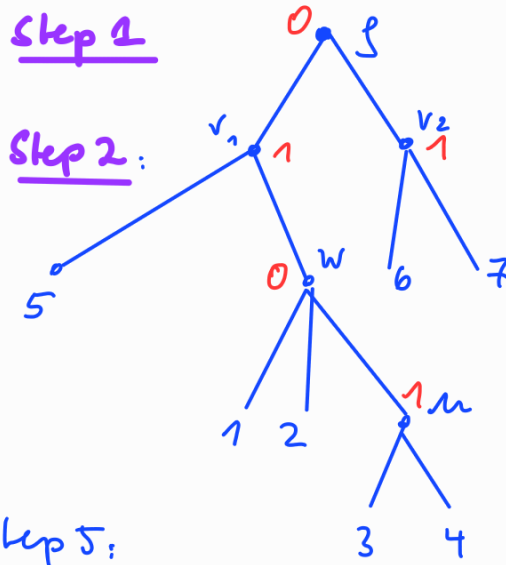
Step 1

Step 2:

Step 3

Step 4:

Step 5:



Step 1 init

call  $(G, s, B=1)$

Step 2 6 disconnected with 2 conn. comp.

$C_1 = \{1 \dots 5\}$

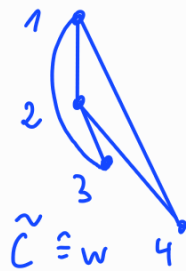
$C_2 = \{6, 7\}$

→ call  $(\overline{G[C_1]}, v_1, 0)$   
call  $(\overline{G[C_2]}, v_2, 0)$

$\overline{G[C_i]}$  must be disconnected!

Step 3

$\overline{G[C_1]}$



$\overline{G[C_2]} = \begin{matrix} \circ & 6 \\ & \circ & 7 \end{matrix}$

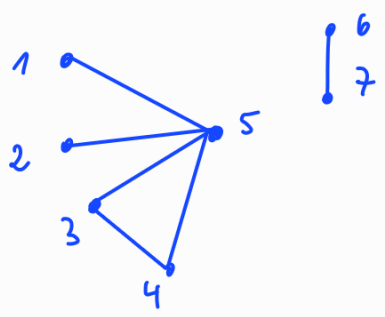
call  $(\overline{G[C]}, w, 1)$ , call  $(\overline{G[\{5\}]}, 5, 1)$   
 $(\overline{G[\{6\}]}, 6, 1)$   
 $(\overline{G[\{7\}]}, 7, 1)$  } Stop in next iteration

Step 4

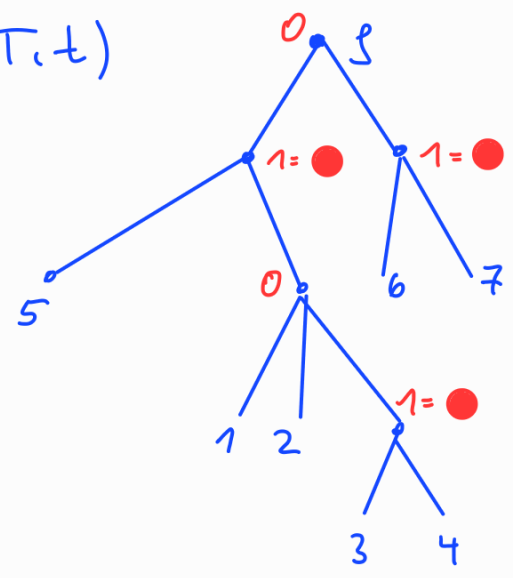


→ step 5 ... and so on.

6



$(T, t)$



leaves of  $T$  are vertices of  $G$ .

By construction:

$$\{x, y\} \in E(G) \Leftrightarrow t(\text{lca}(x, y)) = 1$$

"G is explained by  $(T, t)$ "

Thm:  $G$  cograph  $\Leftrightarrow \exists (T, t)$  that explains  $G$

[without proof  
but idea  
clear]

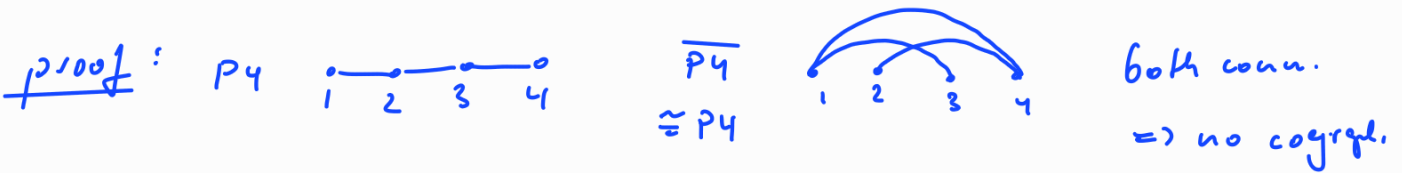
Recap:  $\{x, y\} \in G_\Theta \Leftrightarrow t(\text{lca}(x, y)) = \bullet$

Thm [2013]  $G_\Theta$  valid  $\Leftrightarrow G_\Theta$  is cograph.

⊗ If  $G$  cograph  $\Rightarrow$  every induced subgraph  $H \subseteq G$  is cograph

[exercise, hint look at underlying tree structure]

Thm:  $G$  cograph  $\Leftrightarrow \nexists$  induced  $P_4$ .



$\otimes$   
 $\Rightarrow$  IF  $G$  cograph then it cannot contain induced  $P_4$ .

Now assume  $G$  does not contain induced  $P_4$ .

By induction on  $|V(G)|$  with  $|V(G)| \leq 3 \Rightarrow G$  cograph.  
 Assume statement holds for all  $G$  with  $|V(G)| < n$ , i.e.:  
 every graph with less than  $n$  vertices & without induced  $P_4$  is cograph.

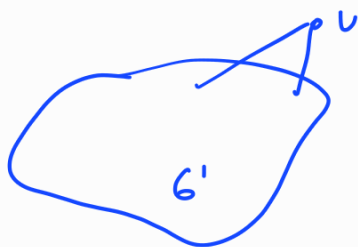
Consider now  $G$  with  $|V(G)| = n \geq 4$  & no induced  $P_4$ .

IF  $G$  disconnected  $\Rightarrow G$  can be written as  $G = G_1 \cup G_2$

$\Rightarrow |V(G_1)|, |V(G_2)| < n$   
 $\stackrel{\text{Ind hyp}}{\Rightarrow} G_1, G_2$  cogr.  $\stackrel{\text{Dct.}}{\Rightarrow} G_1 \cup G_2 = G$  cogr.

IF  $G$  connected, then consider  $G' = G - v$

(remove  $v$  & all incid. edges from  $G$ )

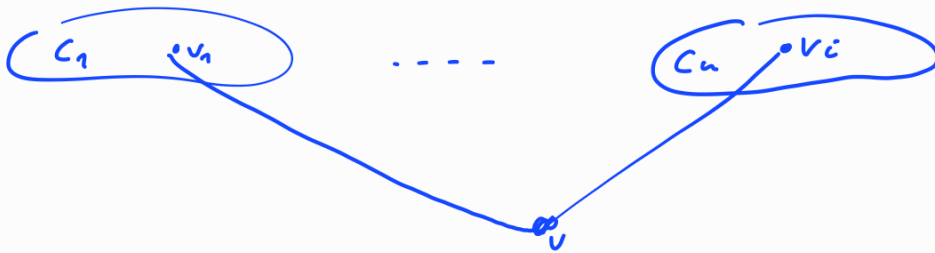


$G'$  is cograph (by Ind hyp).



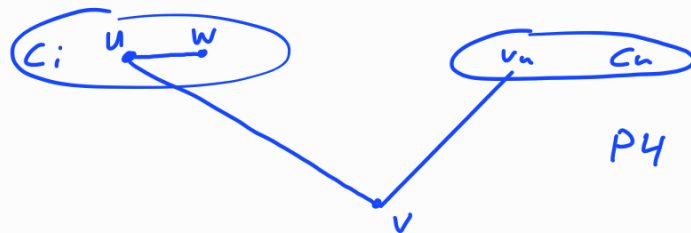
a) if  $G'$  disconnected it has conn. comp.  $C_1 \dots C_n$

Since  $G$  connected  $\Rightarrow v$  must be incident to some  $v_i \in C_i \ \forall i \in \{1 \dots n\}$ .



Assume, for contradiction,  $v$  not adj to all  $w \in V(G')$ .

$\Rightarrow \exists$  conn. comp  $C_i$  of  $G'$  & vertices  $w, u$  st.



$P_4 \nmid$  since  $G$  contains no  $P_4$

$\Rightarrow G = G' \otimes K_1$  conn. per Def.

b) IF  $G'$  connected.

since  $G'$  conn.  $\Rightarrow \overline{G'}$  disconn. & conn.

$\Rightarrow \overline{G'}$   $P_4$ -free

now apply analogous arguments as in case a) for  $\overline{G'}$

## Summary

Thm:

R. feasible  $\Leftrightarrow$  its graph representation  
is a cograph

$\Leftrightarrow \nexists$  induced  $P_4$  

$\Rightarrow$  estimated orthology that violates the  
"cograph property" must be corrected

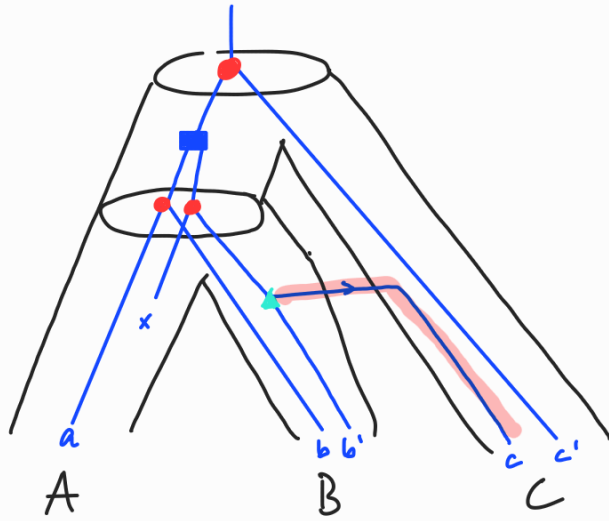
[cograph-editing is NP hard].

# Xenology

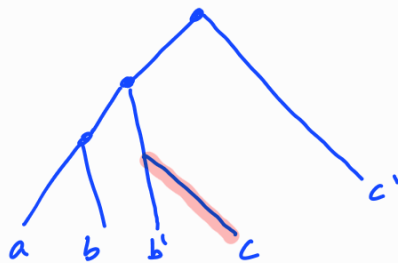
Examples for HBT

& IDEA of inference methods  
→ slides.

# HBT



- losses cannot be observed!
- gene tree with specified HBT-edges

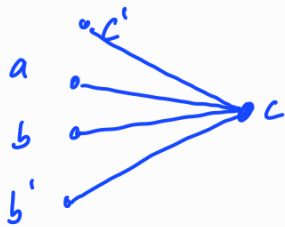


$(T, \lambda)$

$\lambda: E \rightarrow \{0, 1\}$ ,  $\lambda(e) = 1 \Leftrightarrow e$  is HBT edge.

$F(T, \lambda) = (V, E)$  with

$V := L(T)$ ,  $E := \{ \{x, y\} \mid \exists \text{ edge } e \text{ with } \lambda(e) = 1 \text{ along unique } xy\text{-path in } T \}$



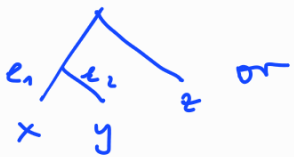
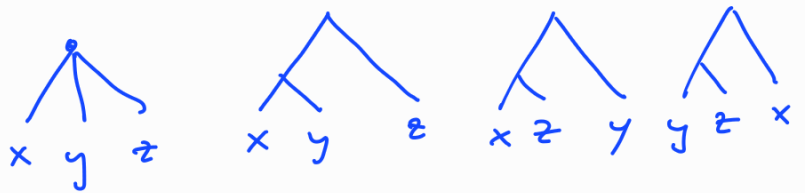
Given  $G$ , when exist  $(T, \lambda)$  st  $G \cong F(T, \lambda)$ ?  
 In this case,  $(T, \lambda)$  explains  $G$ .

Lemma:  $F(T, d)$  does not contain  $K_1 + K_2 \cong " \bullet \text{---} \bullet "$  as an induced subgraph.

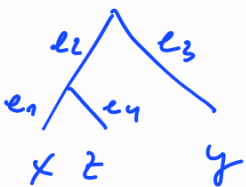
Proof: Assume, for contradiction, it contains  $K_1 + K_2$  as induced subgraph.



Possible trees:



$\Rightarrow$  at least one of  $e_1$  or  $e_2$  has label 1  
 $\Rightarrow \exists$  edge  $\{x, z\}$  or  $\{y, z\} \notin E$



$\Rightarrow d(e_4) = 0$  [since  $z$  not adj to  $y$  &  $x$ ]  
 $d(e_1) = 0$  [since  $z$  not adj to  $x$ ]  
 $d(e_2) = d(e_3) = 0$  [since  $z$  not adj to  $y$ ]  
 $\Rightarrow$  no HBT edges  $\Rightarrow \frac{1}{2}$  to  $\{x, y\} \in E$

[androg case  $\frac{1}{2}$  to  $\{y, z, x\}$ ]

$\Rightarrow F(T, d)$  cannot contain  $K_1 + K_2$   $\square$

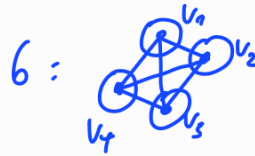
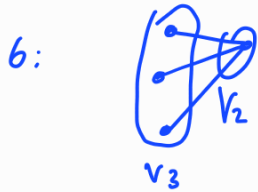
DEF:

$G = (V, E)$  is complete multipartite, if  
exist partition  $V_1 \dots V_k, k \geq 1$  of  $V$   
st

$G[V_i]$  contains no edges,  $1 \leq i \leq k$

$\forall x \in V_i, y \in V_j, i \neq j: \{x, y\} \in E.$

Exmpl:



Lemma


$G$  complete multipartite

$\Leftrightarrow$  does not contain  $K_{k_1+k_2}$  as an induced subgraph

proof:

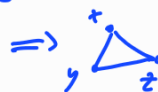
" $\Rightarrow$ " let  $V_1 \dots V_k, k \geq 1$  be partition of  $V$   
Take 3 vertices  $x, y, z \in V.$

cases:  $x, y, z \in V_i \Rightarrow G[V_i]$  edges  
 $\Rightarrow G[\{x, y, z\}] \cong \cdot \cdot \cdot$

$x, y \in V_i, z \in V_j, i \neq j \Rightarrow$  

[analog if  $xz \in V_i, y \in V_j$   
 $yz \in V_i, z \in V_j$ ]

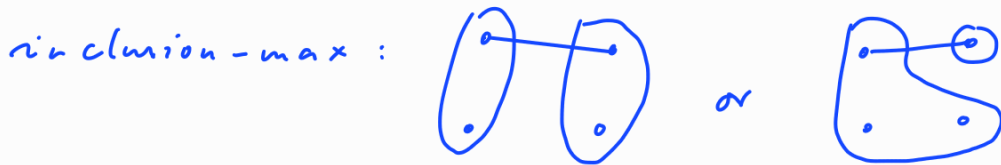
$x \in V_i, y \in V_j, z \in V_k, i, j, k$  pairw. distinct



$\Rightarrow$  no  $K_{k_2+k_1}$  in  $V$

⇐

Let  $V_1 \dots V_k$  be partition of  $V$   
 st  $\{V_i\}$  is inclusion-maximal  
 wrt to "edge" less,  $1 \leq i \leq k$

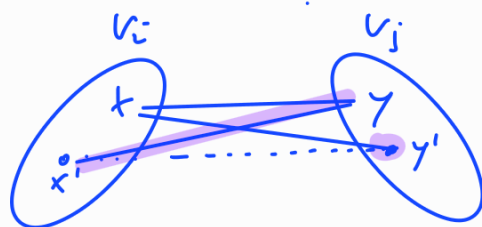
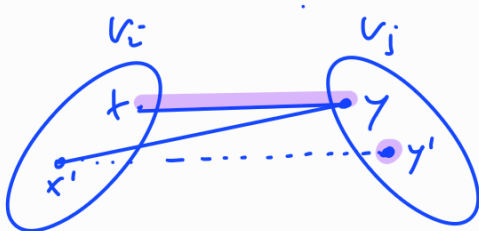
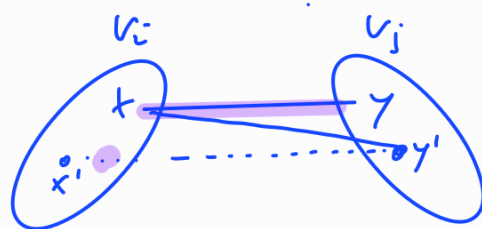
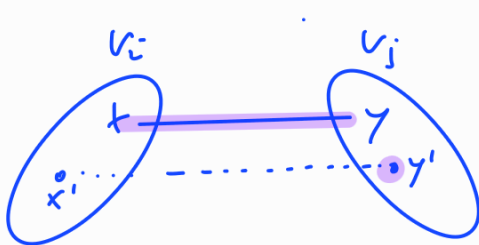


⇒ between distinct  $V_i$  &  $V_j$  there must be  
 at least one edge  $\{xy\}$ ,  $x \in V_i$ ,  $y \in V_j$

Assume  $\exists x' \in V_i$  st  $x'$  not adj to all  $y' \in V_j$



Subcases:



—  $k_2 + k_1$   
 •

⇒ contains  $k_1 + k_2$   $\hookrightarrow$

⇒  $\forall x \in V_i, y \in V_j: \{xy\} \in E \Rightarrow$  complete bipartite graph  $\square$



COR:  $F(T, d)$  is complete multipartite graph.

Theorem [2018]

$G$  can be explained by  $(T, d)$ , i.e.  $G \cong F(T, d)$   
 $\iff G$  is complete multipartite graph.

Proof:

" $\implies$ "  $G \cong F(T, d) \xrightarrow{\text{COR}} G$  complete multip.

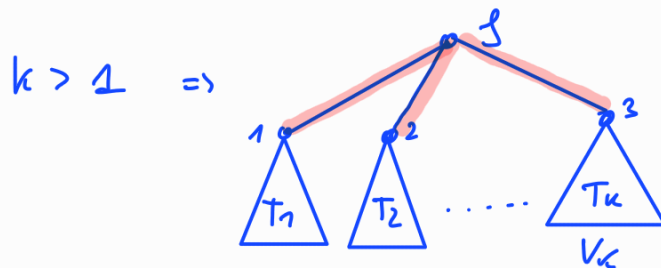
" $\impliedby$ "  $G$  complete multipartite.



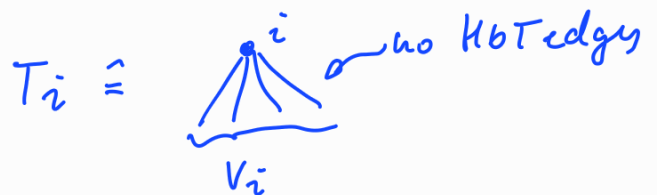
if  $k = 1 \implies G$  edge-less



with no HBT edges



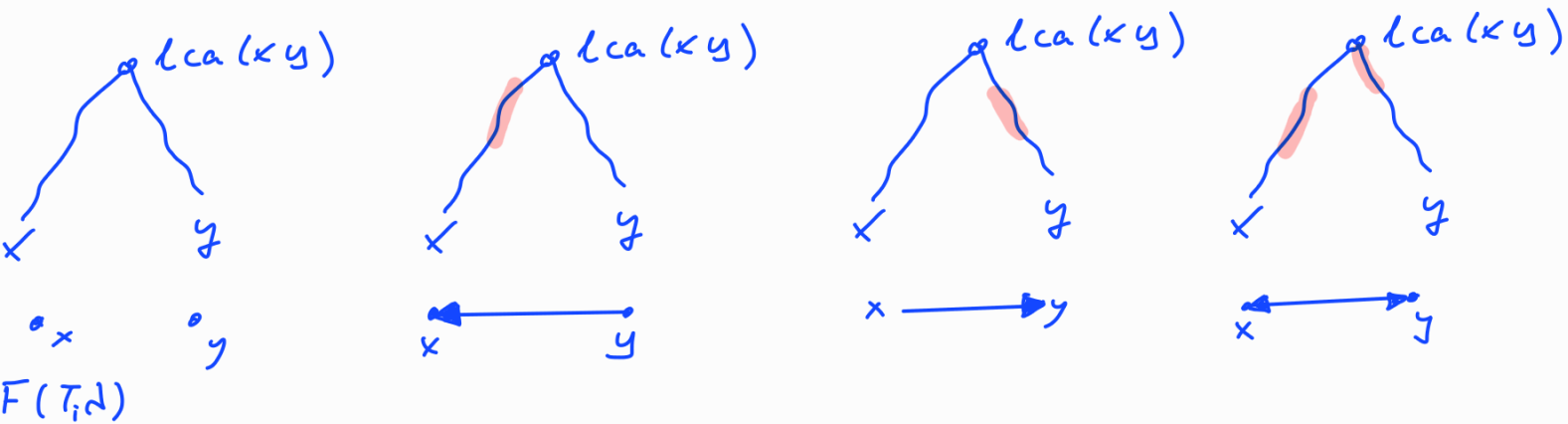
HBT edge



## Outlook:

- The latter characterization helps us to find trees & if there is no tree ( $\exists$  forbidden subgraphs) we must edit graphs so no forb. subgraphs exist (NP hard)   
 hw 17 LCS.

### ► what if direction known?



"directed version".

→ we end in further "forbidden subgraph"

► what if ortholog & paralog are known?  
when do they fit in common tree & do they constrain each other?

► what if evolution is "network" like?

