

The path of trial and tribulation towards  
homology based gene/species tree inference

N. Wieseke

Swarm Intelligence and Complex Systems Group

Leipzig University

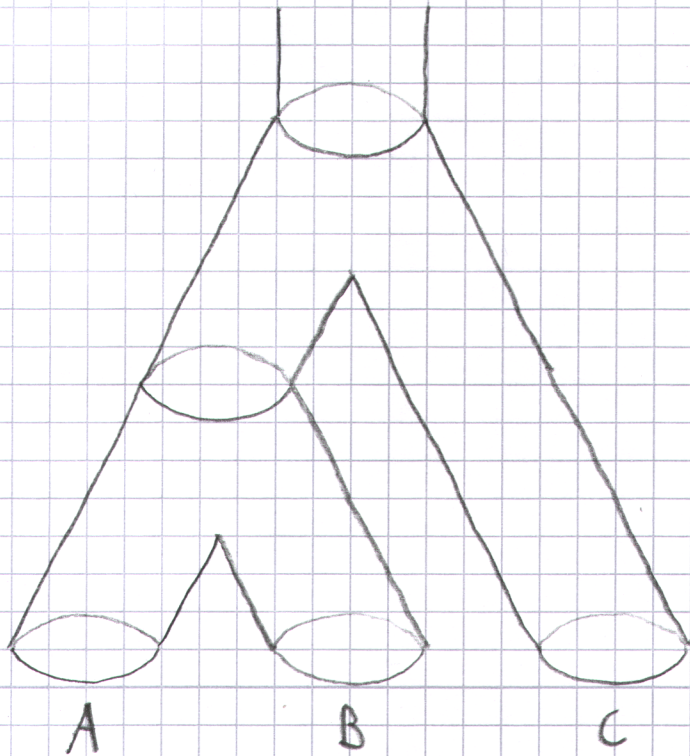
joint work with

Martin Middendorf, SICs, Leipzig

Peter F. Stadler, Bioinf, Leipzig

Marc Hellmuth, Bioinf, Gießenwald

# Problems in Phylogenetics



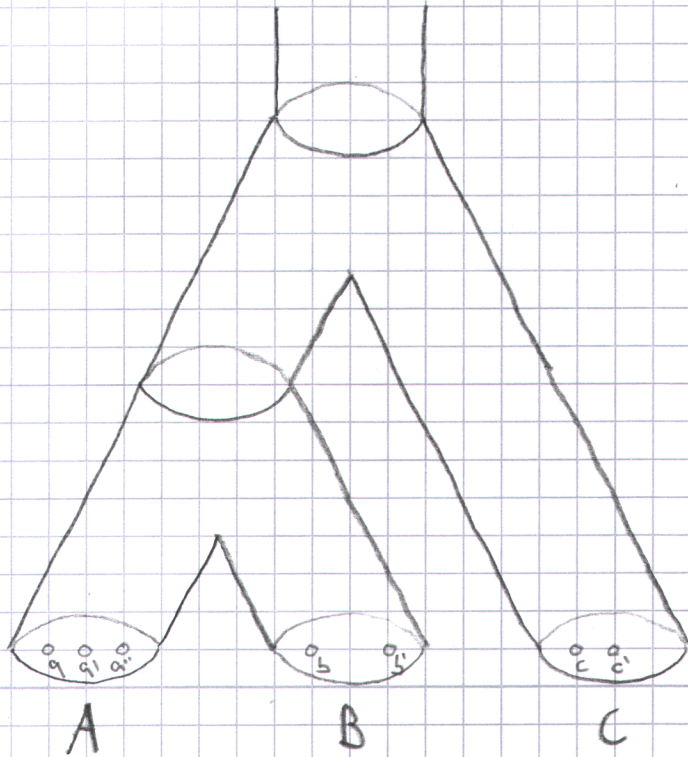
A

B

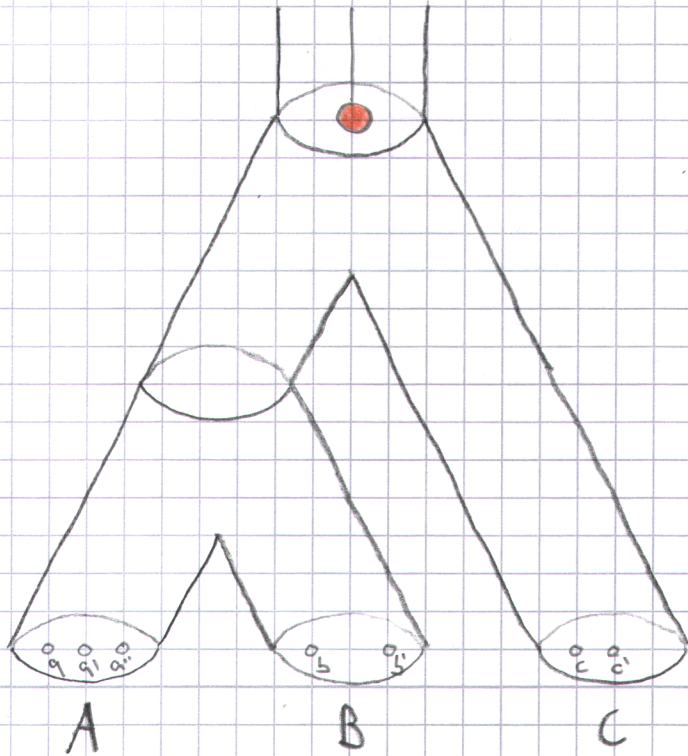
C



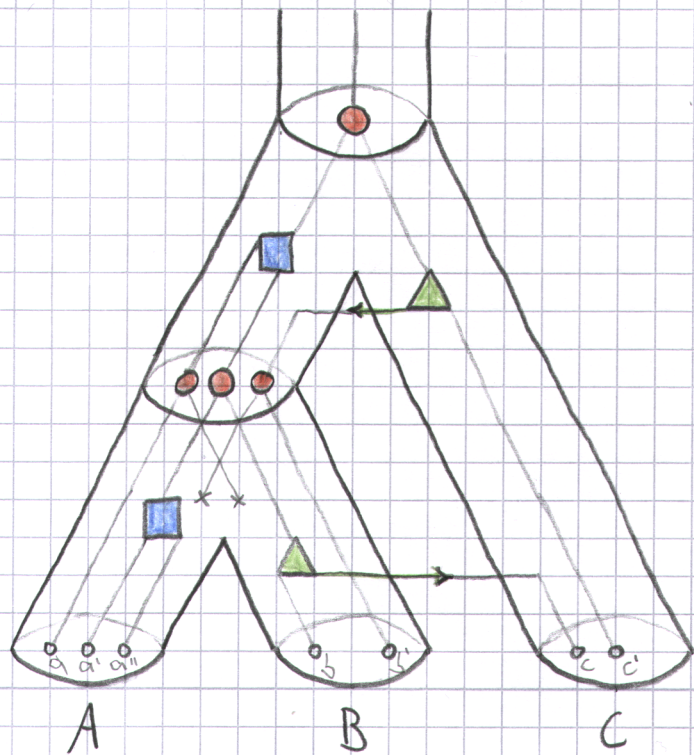
# Problems in Phylogenetics



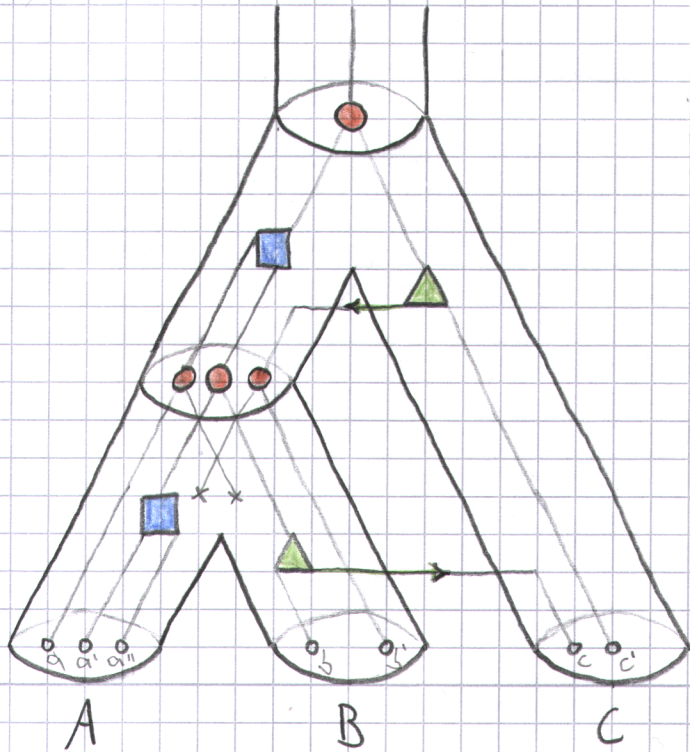
# Problems in Phylogenetics



# Problems in Phylogenetics



# Problems in Phylogenetics



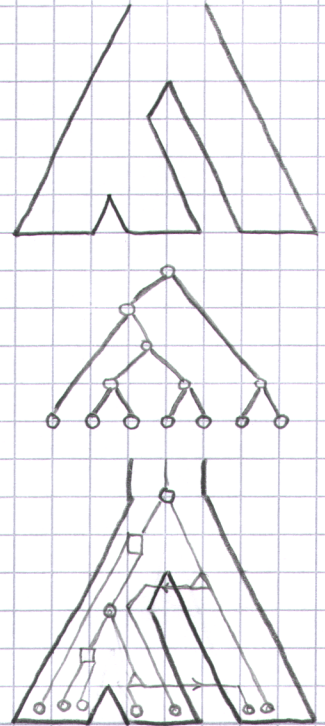
Phylogenetics asks for

- species tree
- gene trees
- evolutionary events
- reconciliation

using only the sequences  
from extant genes

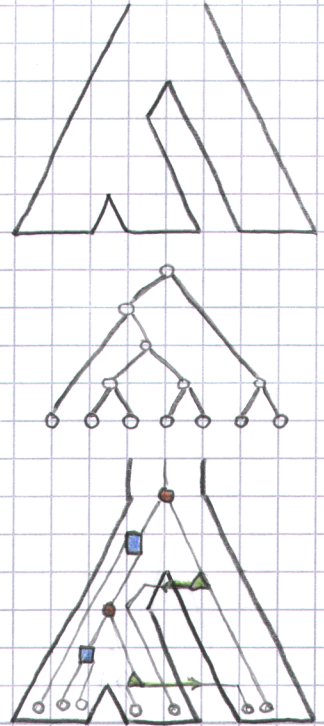
# Classical approach

- 1 infer species tree
  - using 1-to-1 orthologs
  - multiple sequence alignments
  - NJ, ML, or Bayesian methods
- 2 infer gene tree
- 3 reconcile gene tree with species tree
  - using event costs or probabilities
- 4 obtain event labels



# Classical approach

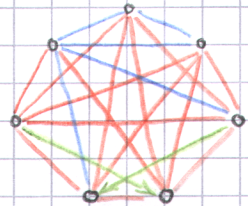
- 1 infer species tree
  - using 1-to-1 orthologs
  - multiple sequence alignments
  - NJ, ML, or Bayesian methods
- 2 infer gene tree
- 3 reconcile gene tree with species tree
  - using event costs or probabilities
- 4 obtain event labels



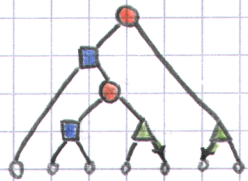


## New approach

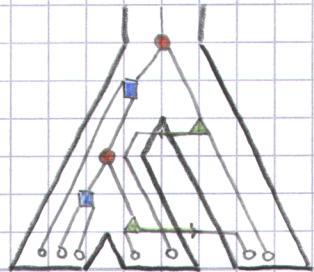
- 1 infer event labels
  - o using pairwise sequence comparisons



- 2 obtain event labeled gene trees



- 3 infer species tree and corresponding reconciliation



# LCA Relations

# Homology Relations

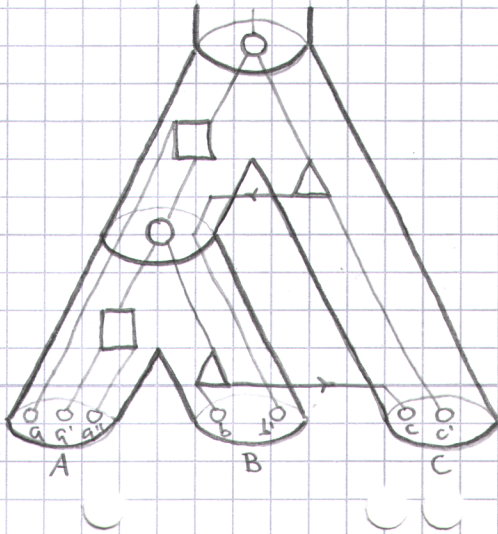
● lca-orthologs  $\Theta_L$

■ lca-paralogs  $\Pi_L$

△ lca-xenologs  $\chi_L$

○ lca-donor-xenologs  $\chi_L^D$

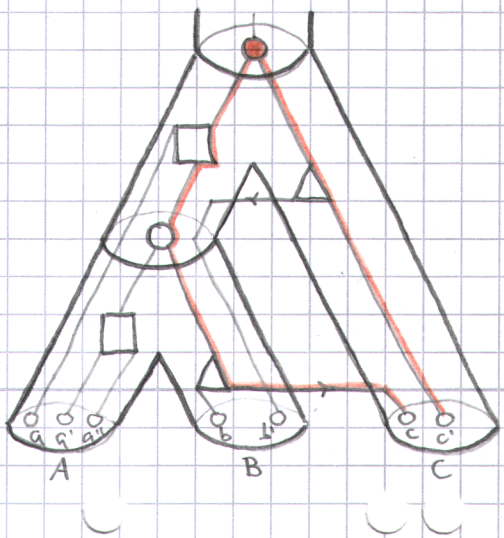
○ lca-acceptor-xenologs  $\chi_L^A$



# LCA Relations

# Homology Relations

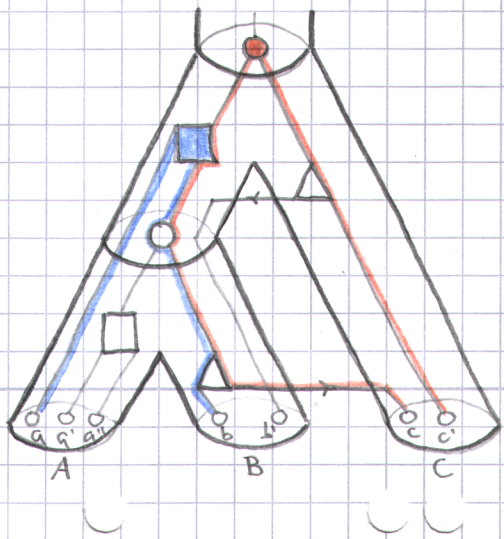
- lca-orthologs  $\Theta_L$
- lca-paralogs  $\Pi_L$
- ▲ lca-xenologs  $\chi_L$ 
  - lca-donor-xenologs  $\chi_L^D$
  - lca-acceptor-xenologs  $\chi_L^A$



# LCA Relations

# Homology Relations

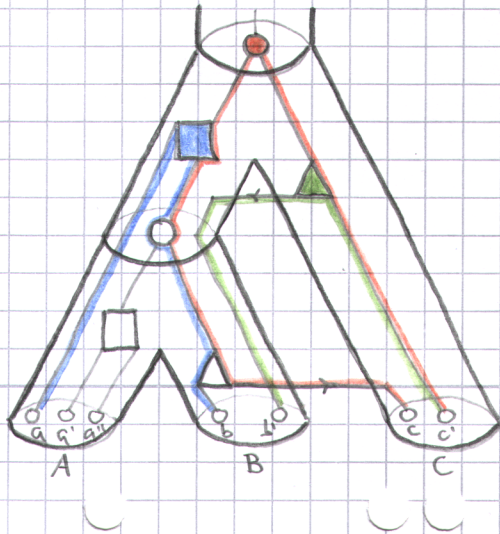
- lca-orthologs  $\Theta_L$
- lca-paralogs  $\Pi_L$
- ▲ lca-xenologs  $\chi_L$ 
  - lca-donor-xenologs  $\chi_L^D$
  - lca-acceptor-xenologs  $\chi_L^A$



# LCA Relations

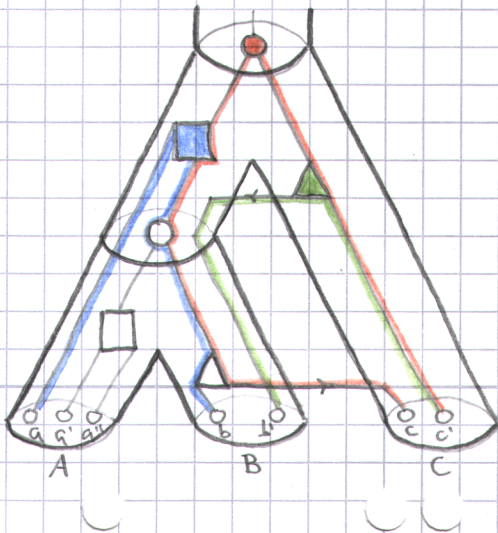
# Homology Relations

- lca-orthologs  $\Theta_L$
- lca-paralogs  $\Pi_L$
- ▲ lca-xenologs  $\chi_L$ 
  - lca-donor-xenologs  $\chi_L^D$
  - lca-acceptor-xenologs  $\chi_L^A$



# LCA Relations

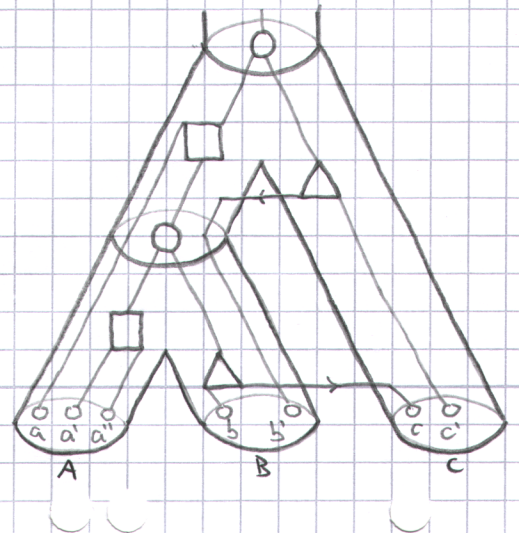
- lca-orthologs  $\Theta_L$
- lca-paralogs  $\Pi_L$
- ▲ lca-xenologs  $\chi_L$ 
  - lca-donor-xenologs  $\chi_L^D$
  - lca-acceptor-xenologs  $\chi_L^A$



# Homology Relations

# Fitch Relations

- Fitch-orthologs  $\Theta_F$   
 $(a,b) \in \Theta_F \iff n(\text{lca}_F(a,b)) = \text{lca}_S(A,B)$
- ▲ Fitch-xenologs  $\chi_F$   
 $(a,b) \in \chi_F \iff$  transfer on path  $a-b$ 
  - Fitch-donor-xenologs  $\chi_F^D$
  - Fitch-acceptor-xenologs  $\chi_F^A$

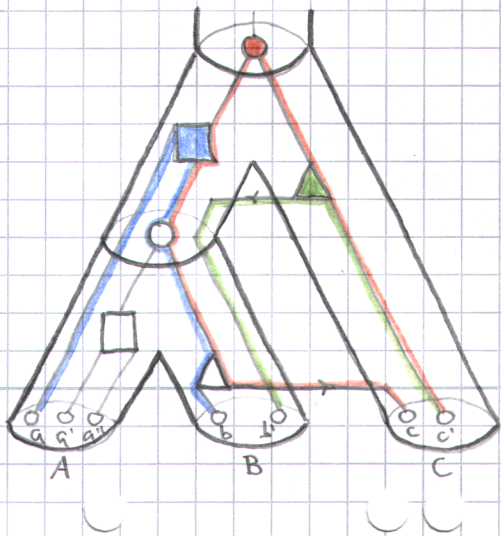


# LCA Relations

- lca-orthologs  $\Theta_L$
- lca-paralogs  $\Pi_L$
- ▲ lca-xenologs  $\chi_L$ 
  - lca-donor-xenologs  $\chi_L^D$
  - lca-acceptor-xenologs  $\chi_L^A$

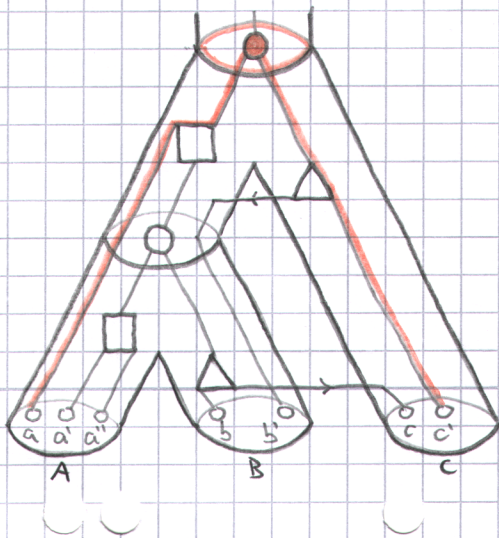
# Homology Relations

$$\Theta_L \supseteq \Theta_F$$



# Fitch Relations

- Fitch-orthologs  $\Theta_F$   
 $(a,b) \in \Theta_F \iff m(lca_F(a,b)) = lca_S(A,B)$
- ▲ Fitch-xenologs  $\chi_F$   
 $(a,b) \in \chi_F \iff$  transfer on path  $a-b$ 
  - Fitch-donor-xenologs  $\chi_F^D$
  - Fitch-acceptor-xenologs  $\chi_F^A$



# LCA Relations

- lca-orthologs  $\Theta_L$
- lca-paralogs  $\Pi_L$
- ▲ lca-xenologs  $\chi_L$ 
  - lca-donor-xenologs  $\chi_L^D$
  - lca-acceptor-xenologs  $\chi_L^A$

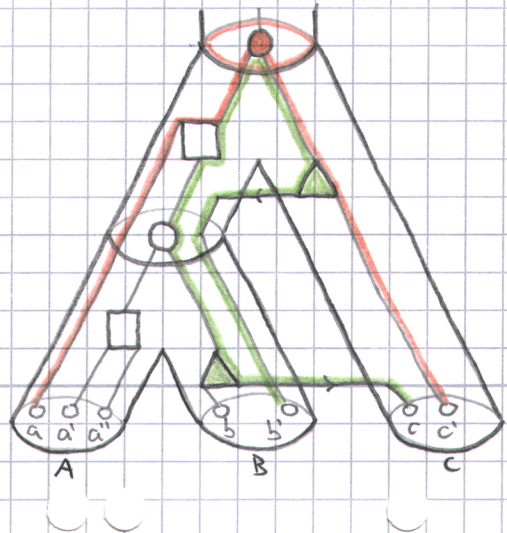
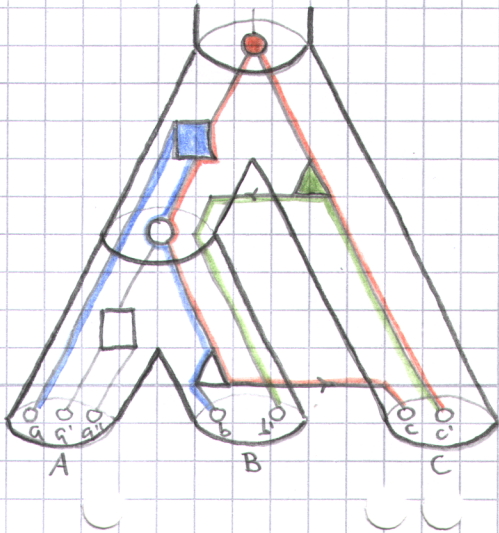
# Homology Relations

$$\Theta_L \supseteq \Theta_F$$

$$\chi_L \subseteq \chi_F$$

# Fitch Relations

- Fitch-orthologs  $\Theta_F$   
 $(a,b) \in \Theta_F \iff \mu(lca_F(a,b)) = lca_S(A,B)$
- ▲ Fitch-xenologs  $\chi_F$   
 $(a,b) \in \chi_F \iff$  transfer on path  $a-b$ 
  - Fitch-donor-xenologs  $\chi_F^D$
  - Fitch-acceptor-xenologs  $\chi_F^A$

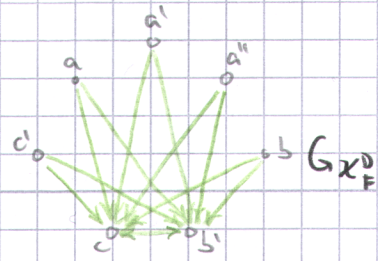
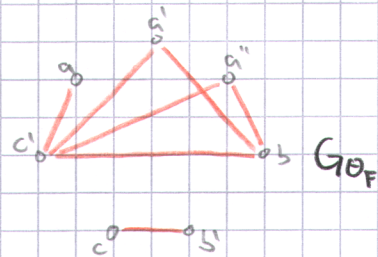
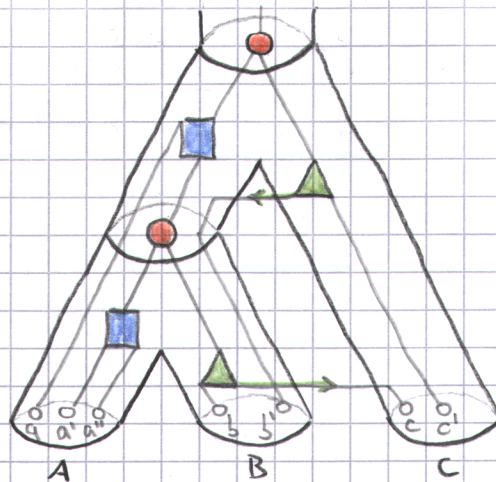
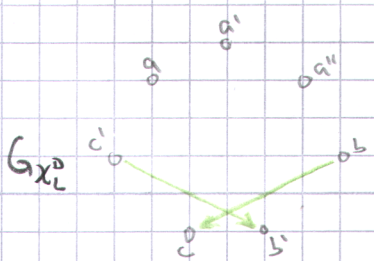
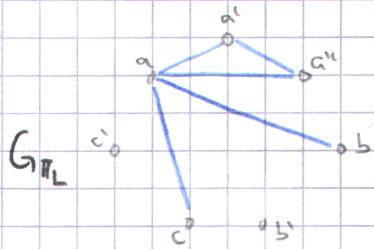
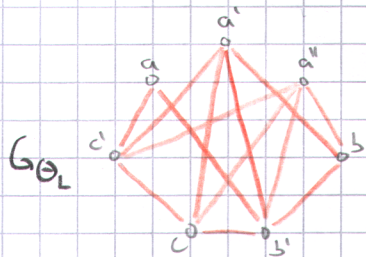




# LCA Relations

# Graph Representation

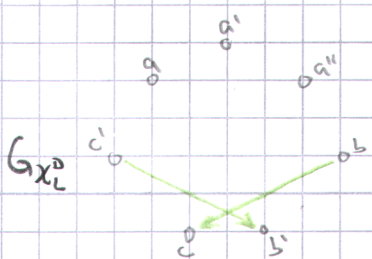
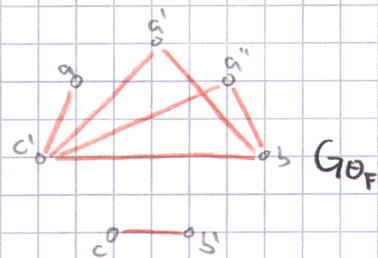
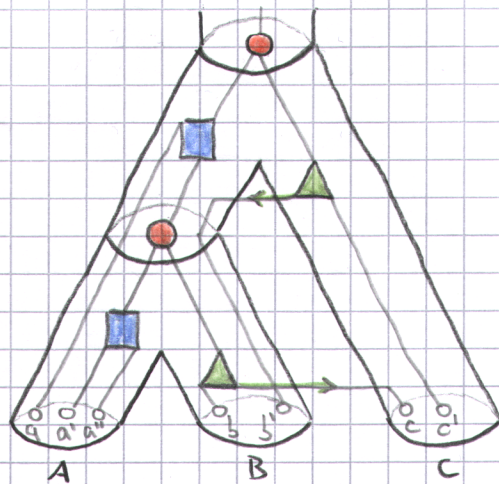
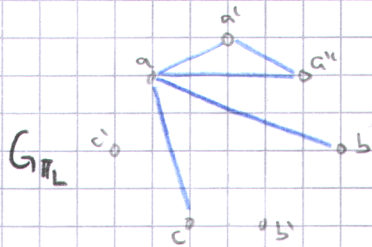
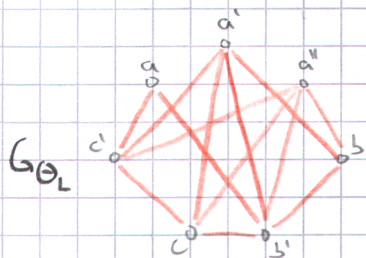
# Fitch Relations



# LCA Relations

# Graph Representation

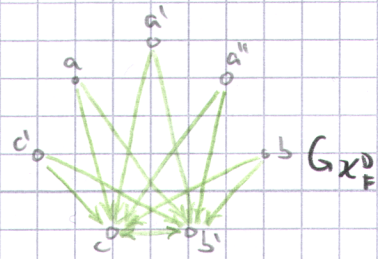
# Fitch Relations



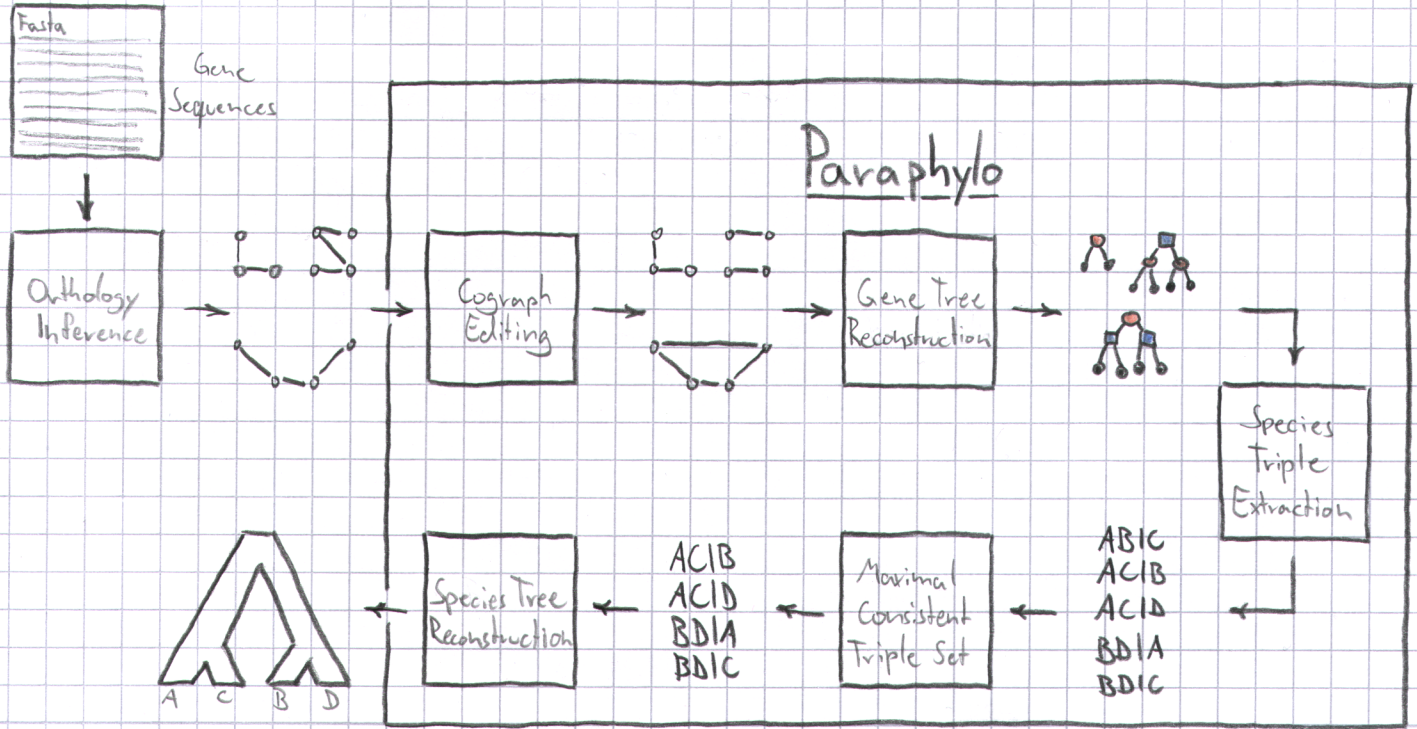
In case of no HGT events

$$G_{\Theta_L} \stackrel{!}{=} G_{\Theta_F}$$

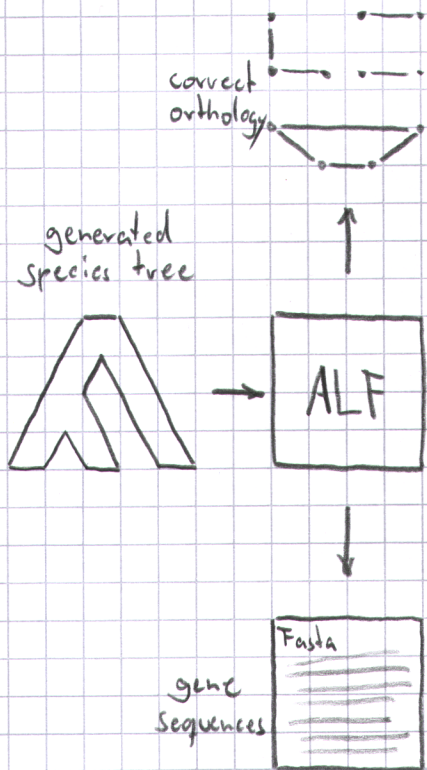
$$G_{\Pi_L} \stackrel{!}{=} G_{\Theta_F}$$



# Pipeline for datasets without HGT

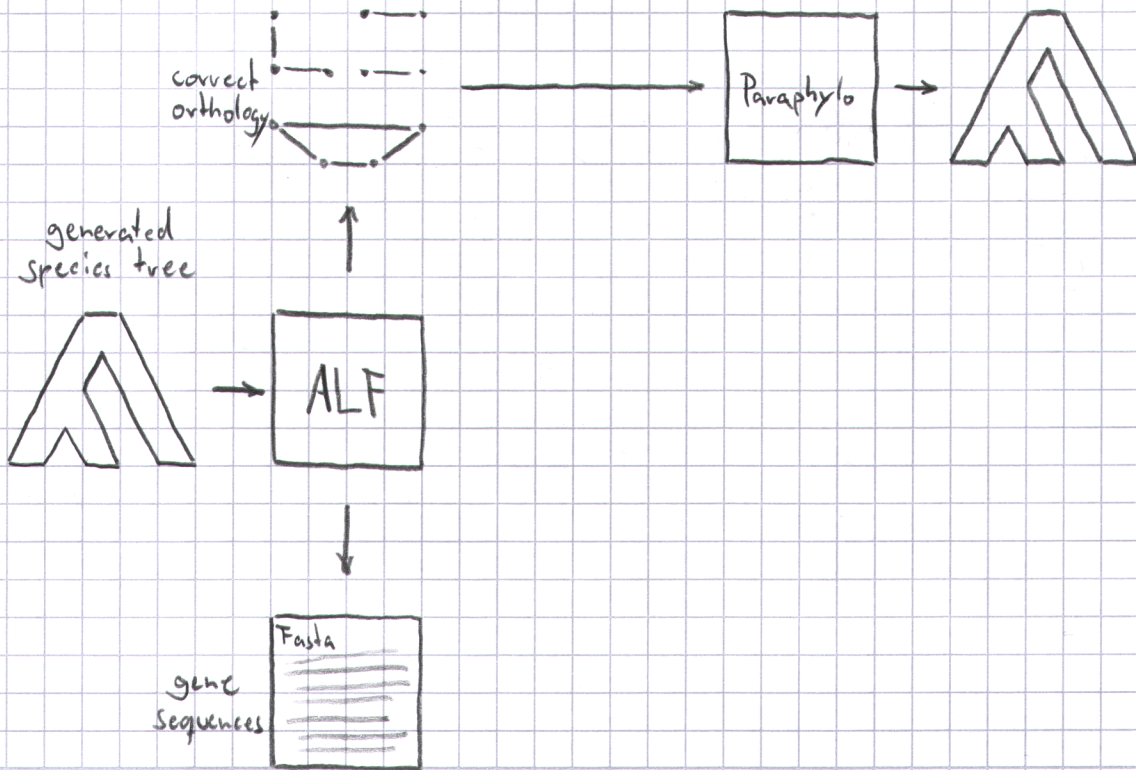


# Simulations without HGT



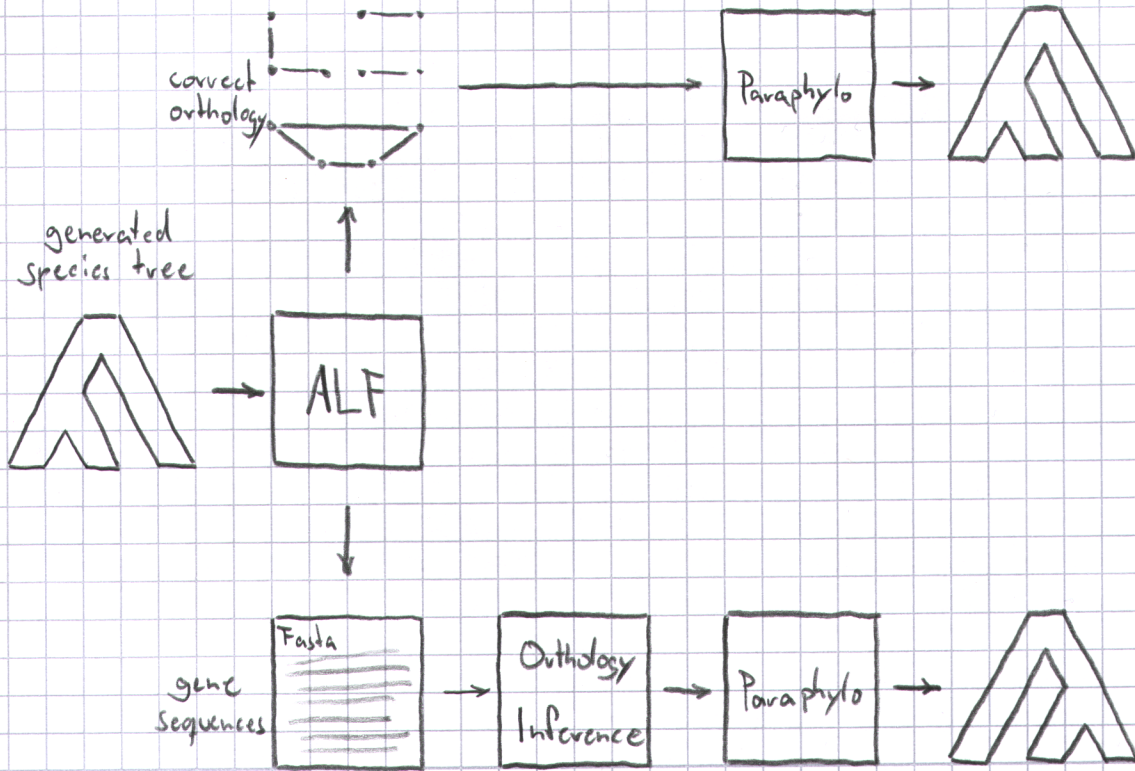
# Simulations without HGT

reconstructed species tree

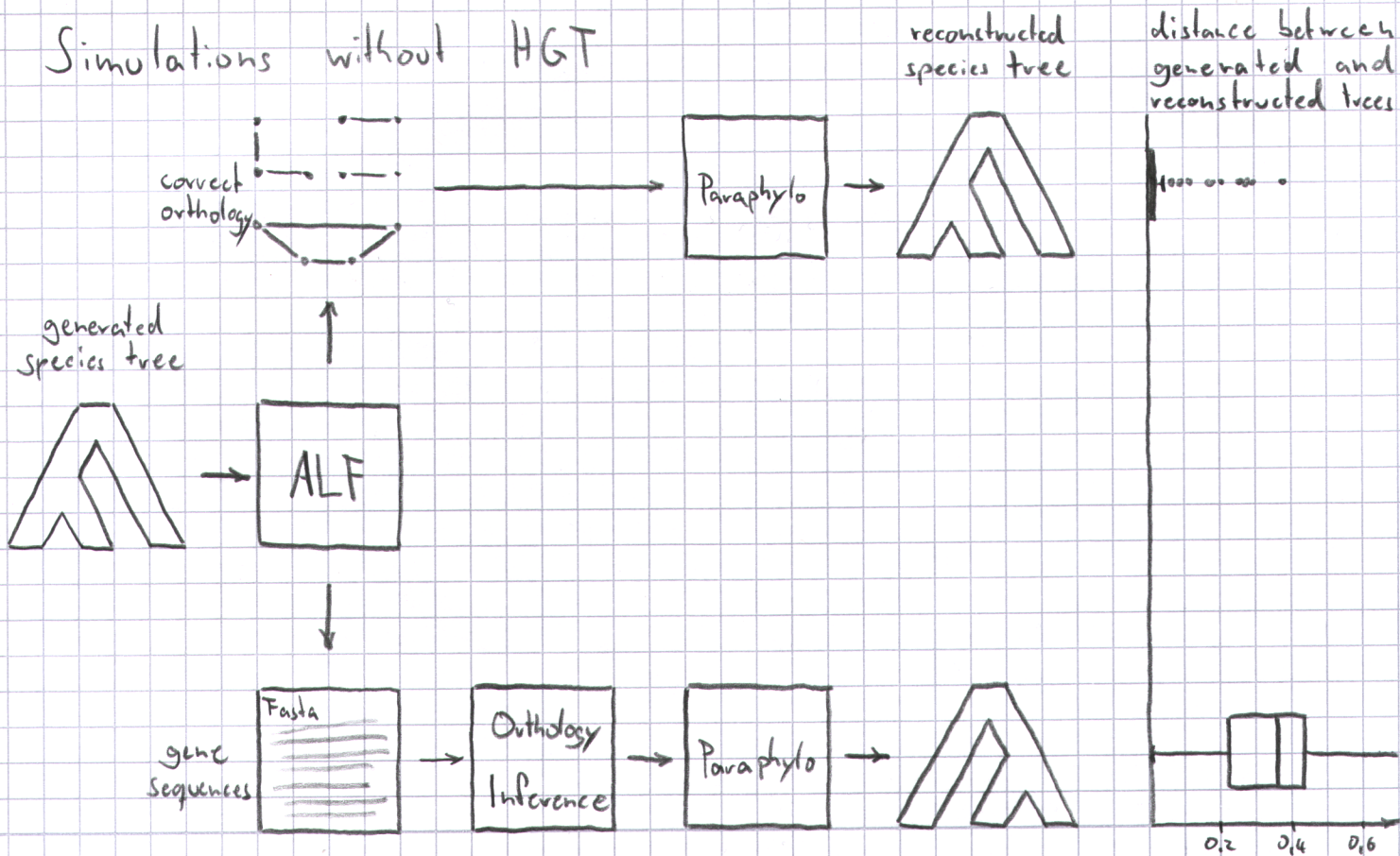


# Simulations without HGT

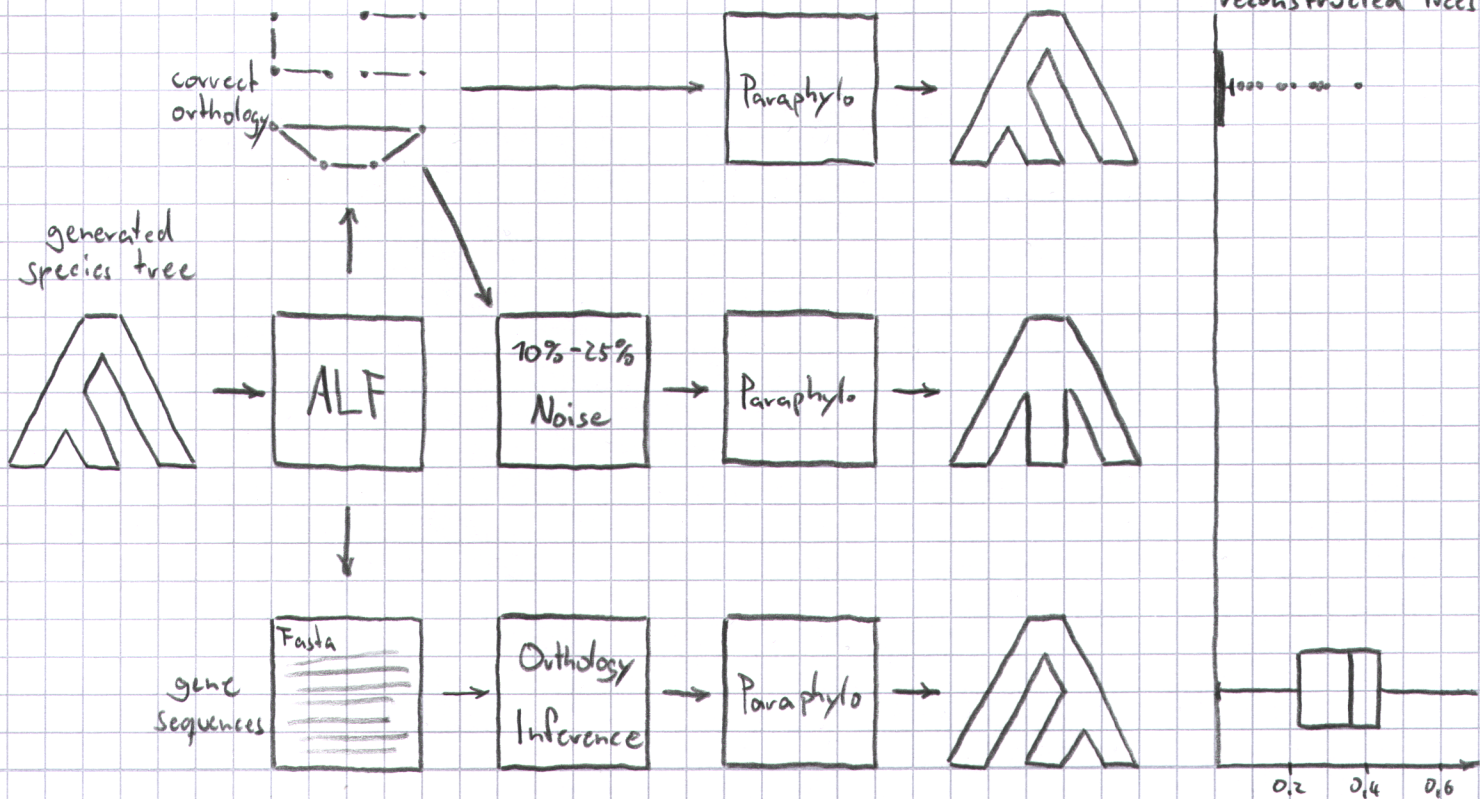
reconstructed species tree



# Simulations without HGT

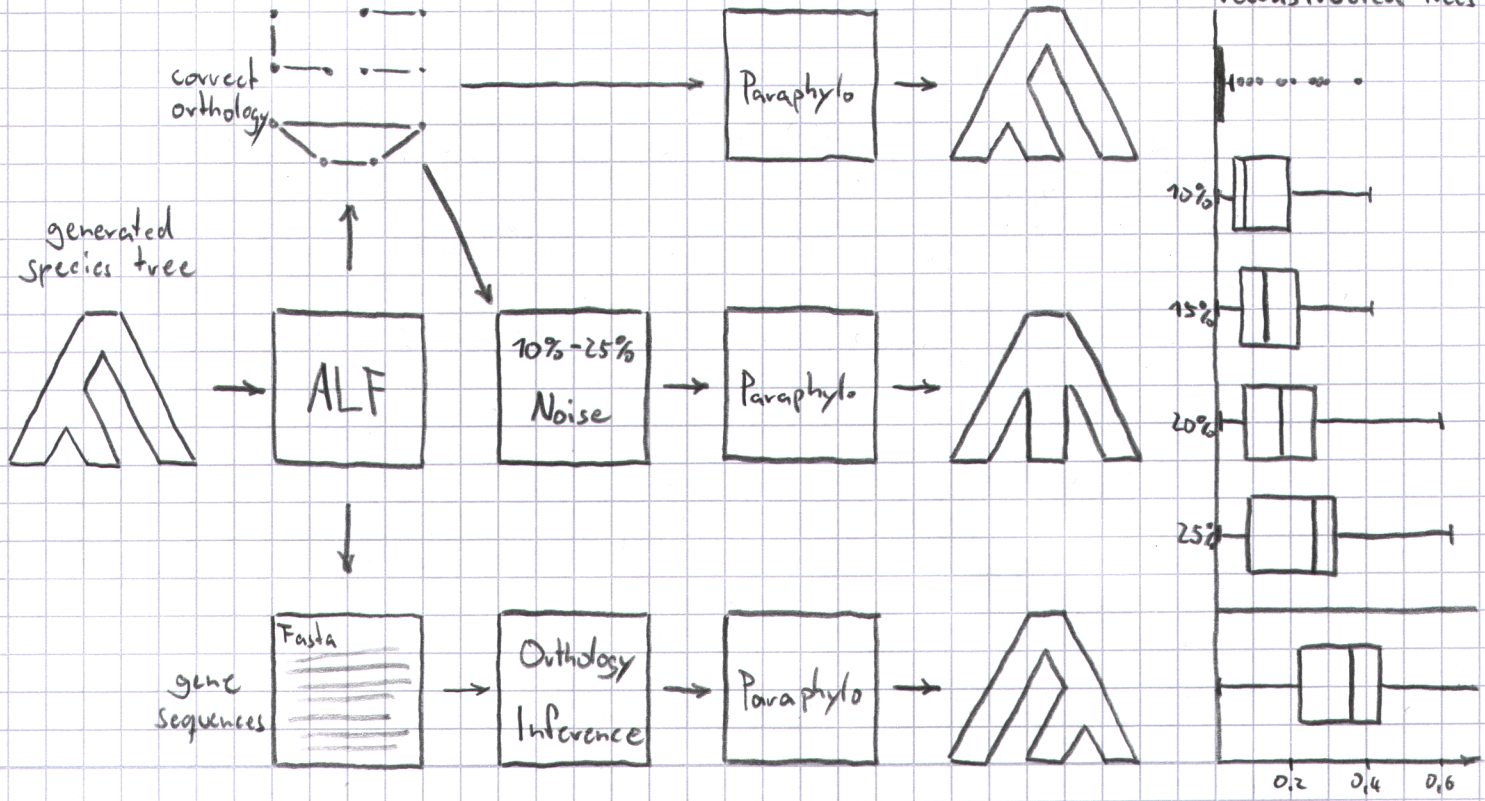


# Simulations without HGT





# Simulations without HGT



Orthology Inference - when does it go wrong

□ no molecular clock-like data

Orthology Inference - when does it go wrong

- no molecular clock-like data
- horizontal gene transfer
- gene losses

# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance

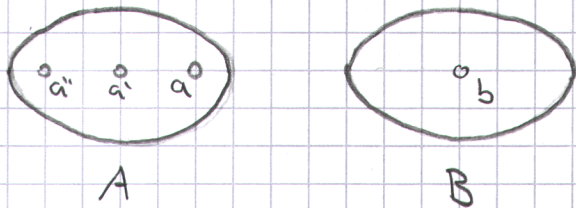
$d(a, b)$  between all genes  $a, b$

- determine reciprocal best matches

$(a, b)$  is a RBM iff  $\forall a' \in A, b' \in B$

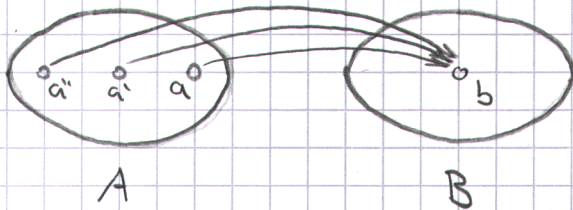
$d(a, b) \leq d(a', b)$  and  $d(a, b) \leq d(a, b')$

- consider RBMs as orthologs



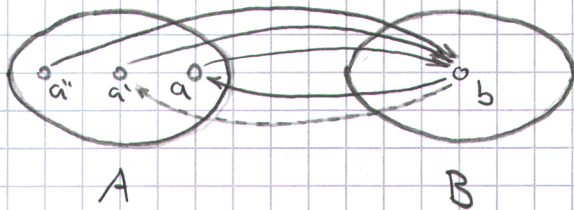
# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance  
 $d(a,b)$  between all genes  $a,b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs



# Orthology Inference - based on pairwise sequence comparisons

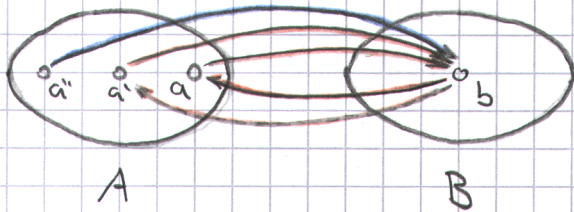
- compute evolutionary distance  $d(a,b)$  between all genes  $a,b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs



$$d(a,b) \leq d(a',b) \leq d(a'',b)$$

# Orthology Inference - based on pairwise sequence comparisons

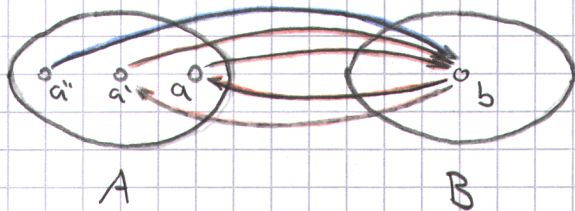
- compute evolutionary distance  
 $d(a,b)$  between all genes  $a,b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs



$$d(a,b) \leq d(a',b) \leq d(a'',b)$$

# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance  $d(a,b)$  between all genes  $a,b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs



$$d(a,b) \leq d(a',b) \leq d(a'',b)$$

time

0

$a''$

$a$

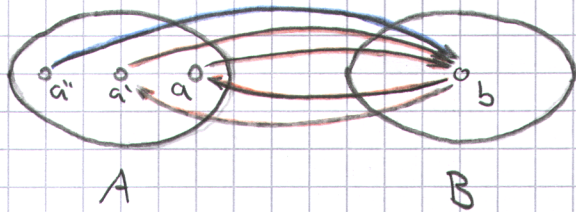
$a'$

$b$

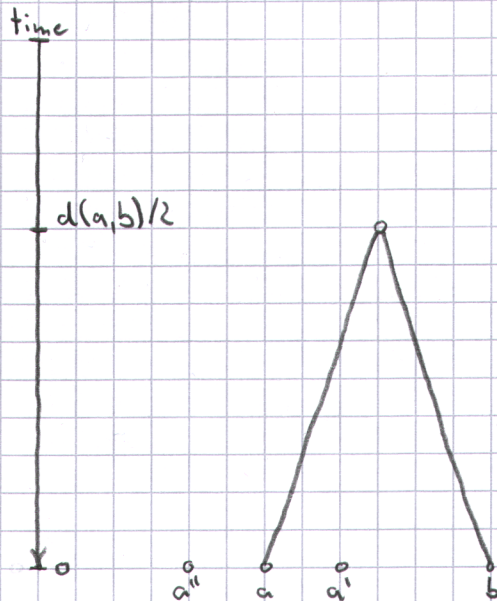


# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance  $d(a,b)$  between all genes  $a,b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs

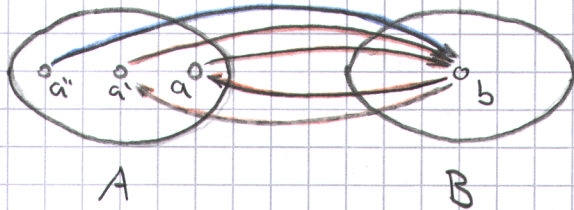


$$d(a,b) \approx d(a',b) \ll d(a'',b)$$

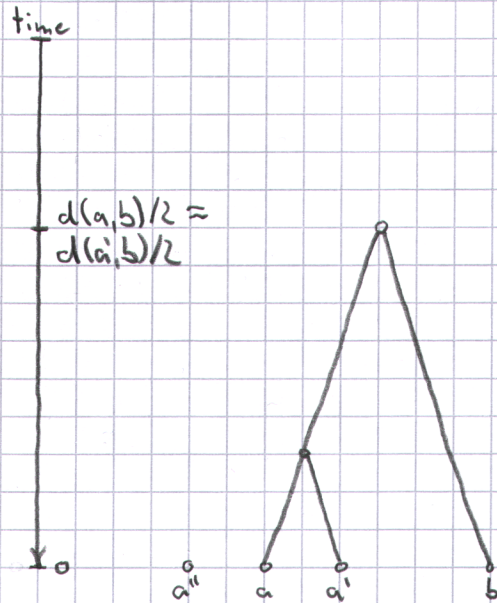


# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance  $d(a,b)$  between all genes  $a, b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs

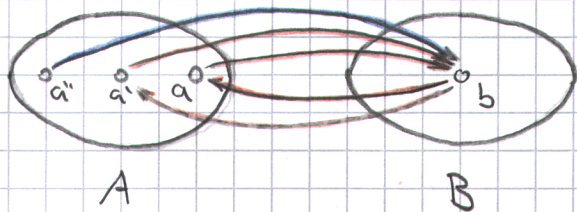


$$d(a,b) \approx d(a',b) \ll d(a'',b)$$

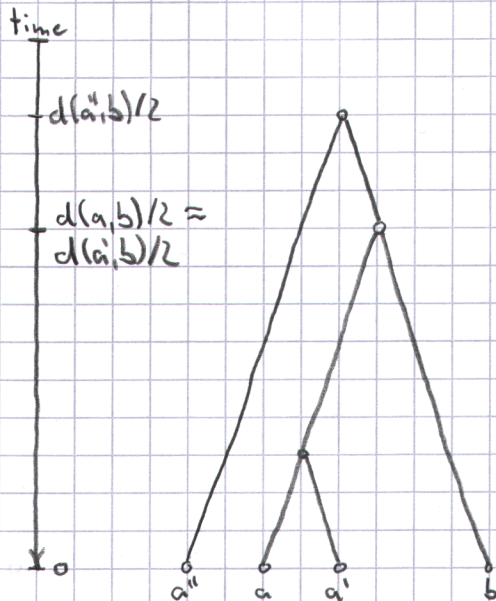


# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance  $d(a, b)$  between all genes  $a, b$
- determine reciprocal best matches  
 $(a, b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a, b) \leq d(a', b)$  and  $d(a, b) \leq d(a, b')$
- consider RBMs as orthologs

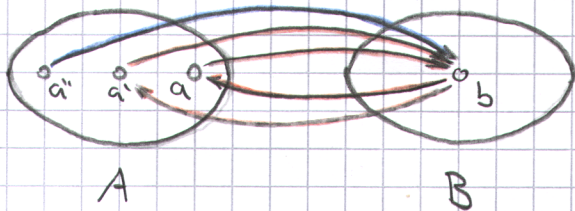


$$d(a, b) \approx d(a', b) \ll d(a'', b)$$

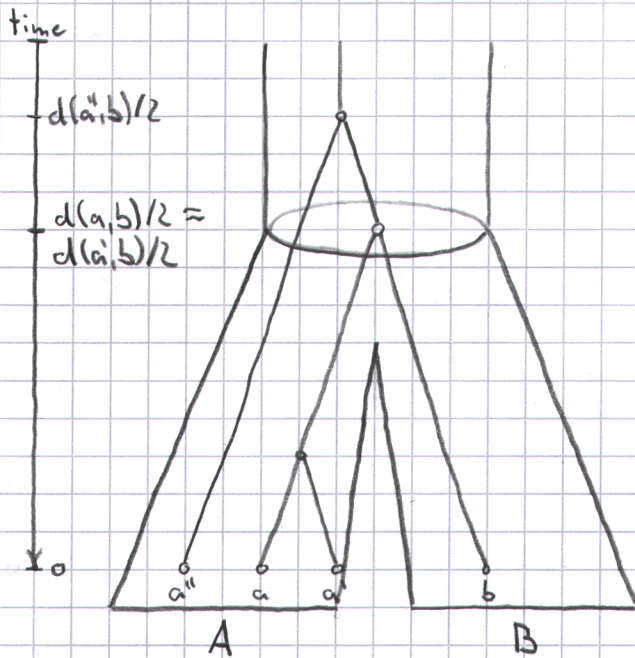


# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance  $d(a,b)$  between all genes  $a,b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs

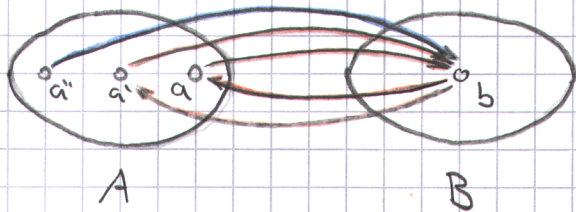


$$d(a,b) < d(a',b) < d(a'',b)$$

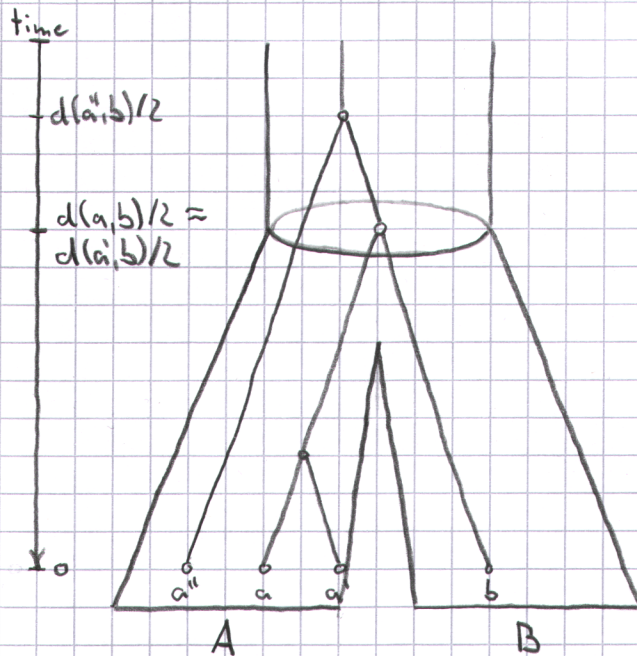


# Orthology Inference - based on pairwise sequence comparisons

- compute evolutionary distance  $d(a,b)$  between all genes  $a,b$
- determine reciprocal best matches  
 $(a,b)$  is a RBM iff  $\forall a' \in A, b' \in B$   
 $d(a,b) \leq d(a',b)$  and  $d(a,b) \leq d(a,b')$
- consider RBMs as orthologs



$$d(a,b) \approx d(a',b) \ll d(a'',b)$$

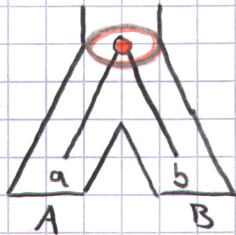


If  $(a,b)$  is a RBM then assume  
 $d(a,b) \approx d(A,B)$

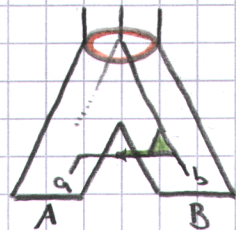
# Orthology Inference - when does it go wrong

Let  $(a,b)$  be a RBM

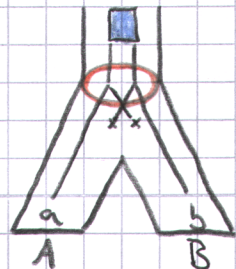
$$\square d(a,b) \hat{=} d(A,B)$$



$$\square d(a,b) \hat{<} d(A,B)$$



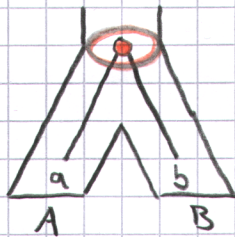
$$\square d(a,b) \hat{>} d(A,B)$$



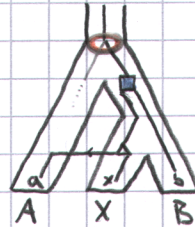
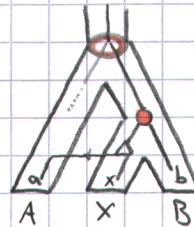
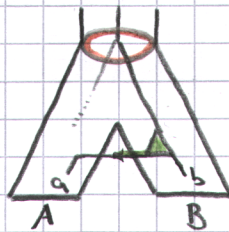
# Orthology Inference - when does it go wrong

Let  $(a,b)$  be a RBM

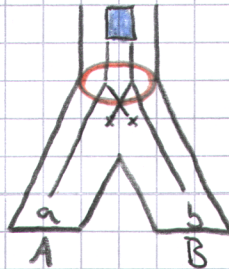
□  $d(a,b) \cong d(A,B)$



□  $d(a,b) \hat{=} d(A,B)$

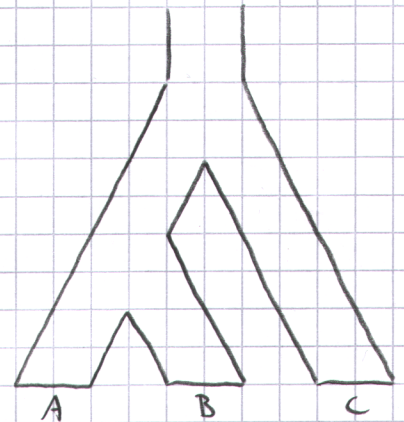
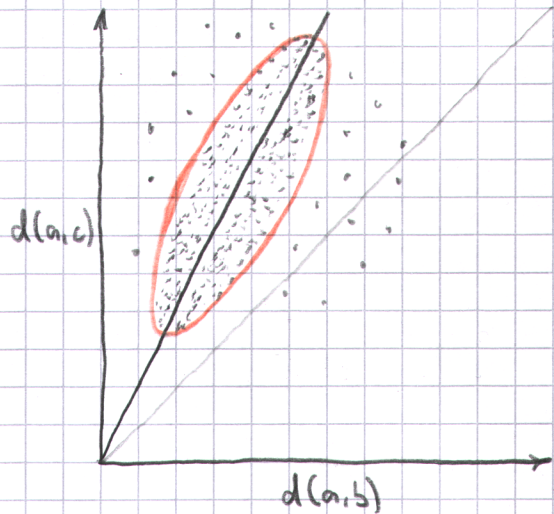


□  $d(a,b) \hat{>} d(A,B)$



# Improved Orthology Inference

- use RBM strategy to obtain an initial orthology estimation

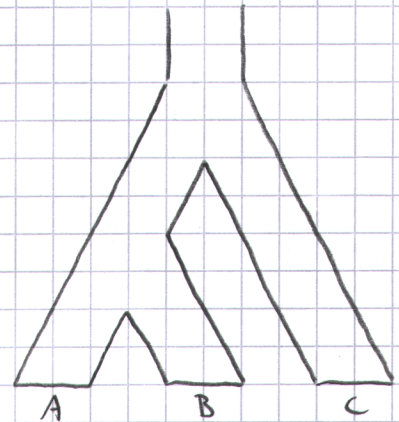
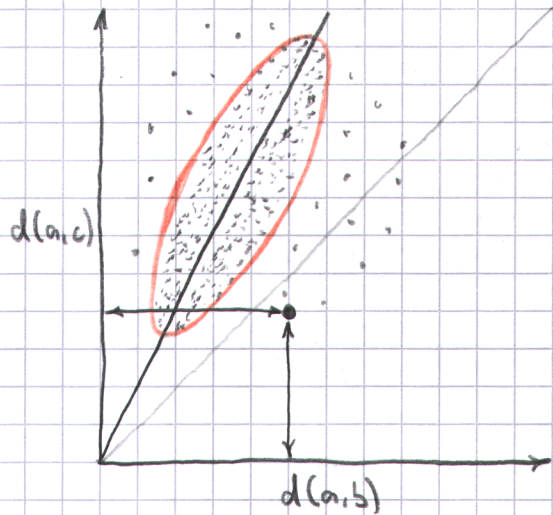


- remove  $(a,b), (a,c)$  for all outliers  $(d(a,b), d(a,c))$



# Improved Orthology Inference

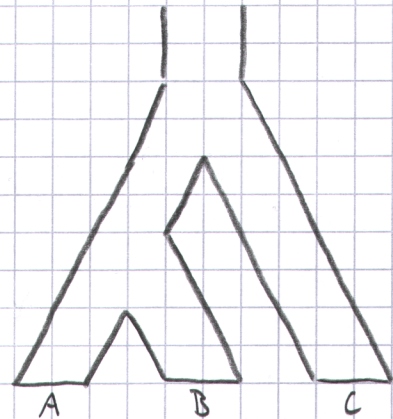
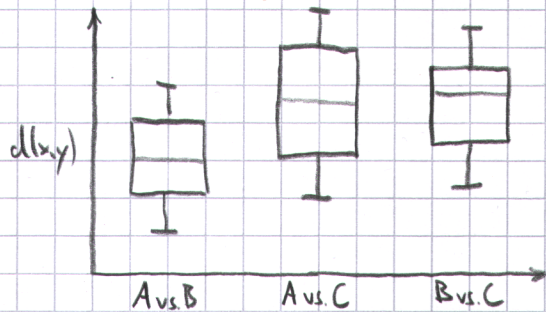
- use RBM strategy to obtain an initial orthology estimation



- remove  $(a,b), (a,c)$  for all outliers  $(d(a,b), d(a,c))$

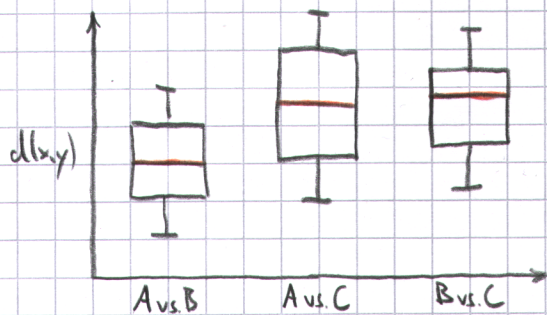
# Improved Orthology Inference

- use filtered orthology estimates

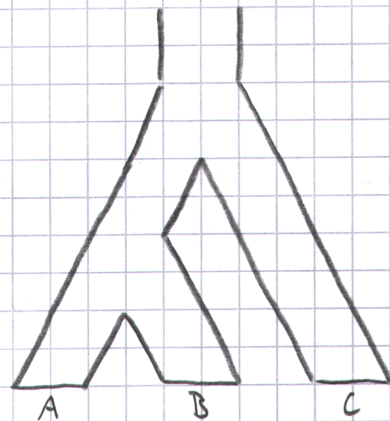


# Improved Orthology Inference

- use filtered orthology estimates

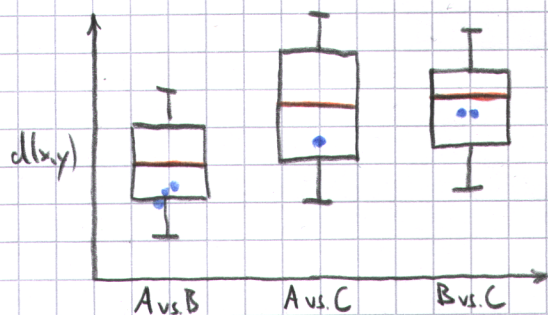


- define  $d(X,Y) := \text{median}(d(x,y))$

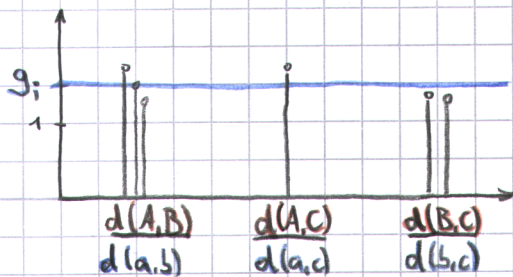


# Improved Orthology Inference

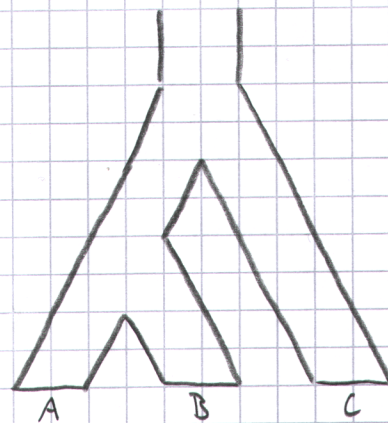
- use filtered orthology estimates



- define  $d(X,Y) := \text{median}(d(x,y))$



with  $g_i, b, c$  from gene family  $i$



## Improved Orthology Inference

$$\square d(a,b) \hat{=} d(A,B) \quad \text{translates to} \quad d(a,b) \cdot g_i \approx d(A,B)$$

$$\square d(a,b) \hat{<} d(A,B) \quad \text{--- u ---} \quad d(a,b) \cdot g_i \ll d(A,B)$$

$$\square d(a,b) \hat{>} d(A,B) \quad \text{--- i ---} \quad d(a,b) \cdot g_i \gg d(A,B)$$

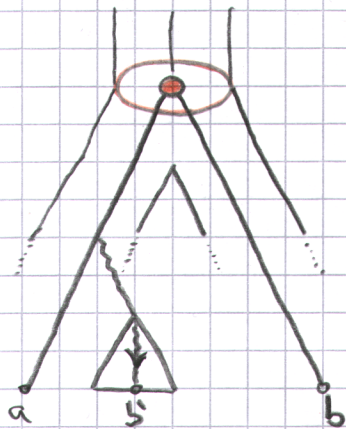
## Infer Fitch Orthologs $\Theta_F$

$$\Theta_F = \{(a,b) \mid d(a,b) \cdot g_i \approx d(A,B)\}$$

# Fitch Xenology Inference

Idea:  $(a, b') \in \chi_F^D$  if  $\exists b$  s.t.  $d(a, b') \cdot g_i \ll d(A, B) \approx d(a, b) \cdot g_i$

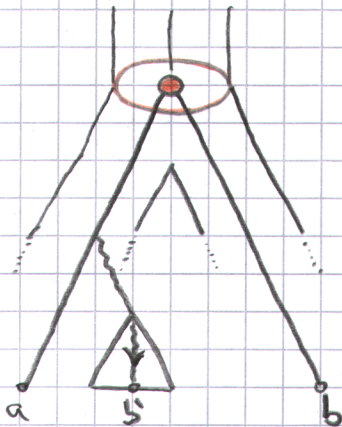
$$\chi_F^D = \left\{ (x, b') \mid \exists a, b, a \in A, b, b' \in B, d(x, a) \leq d(x, b') \text{ and } d(a, b') \cdot g_i \ll d(A, B) \approx d(a, b) \cdot g_i \right\}$$



# Fitch Xenology Inference

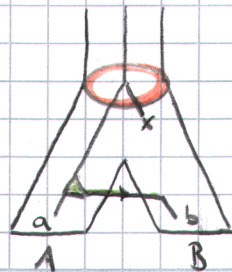
Idea:  $(a, b') \in \chi_F^D$  if  $\exists b$  s.t.  $d(a, b') \cdot g_i \ll d(A, B) \approx d(a, b) \cdot g_i$

$$\chi_F^D = \left\{ (x, b') \mid \exists a, b, a \in A, b, b' \in B, d(x, a) \leq d(x, b') \text{ and } d(a, b') \cdot g_i \ll d(A, B) \approx d(a, b) \cdot g_i \right\}$$

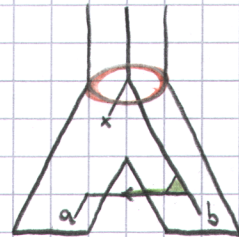


There is ambiguity!

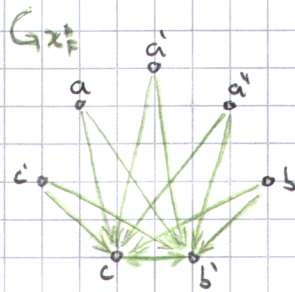
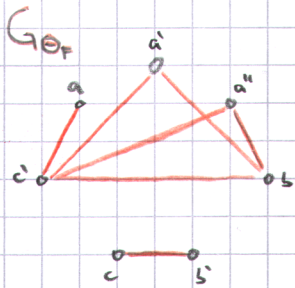
$(a, b) \in \chi_F^D$



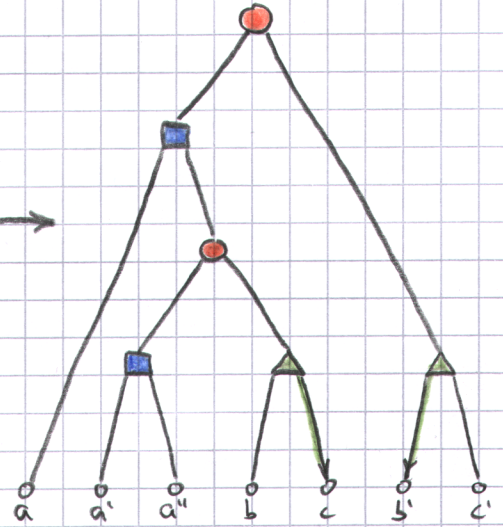
$(b, a) \in \chi_F^D$



# Reconstruction of event labeled gene trees



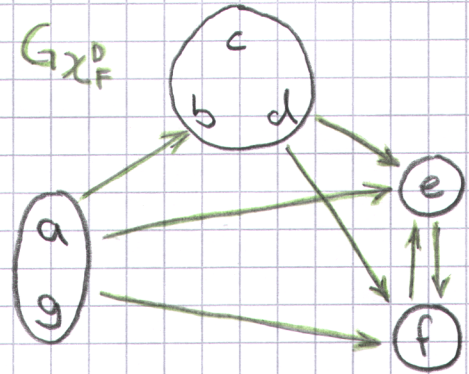
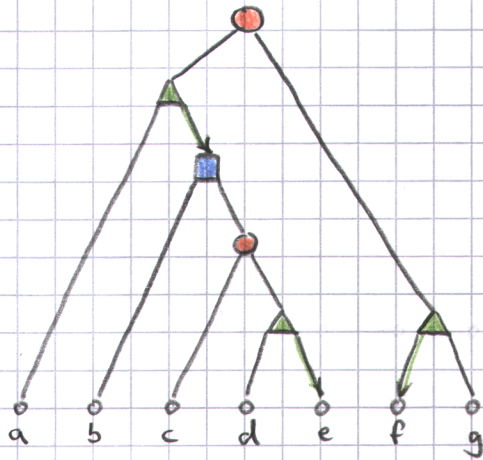
estimated Fitch Relations



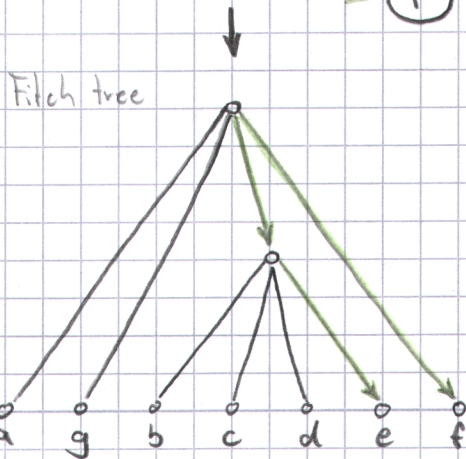
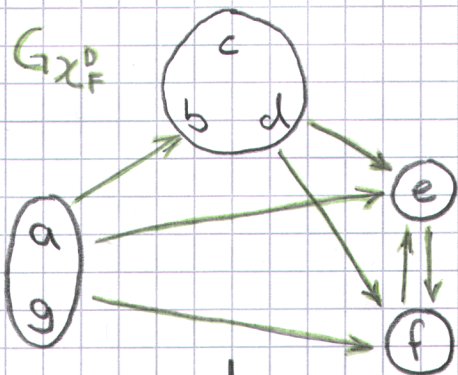
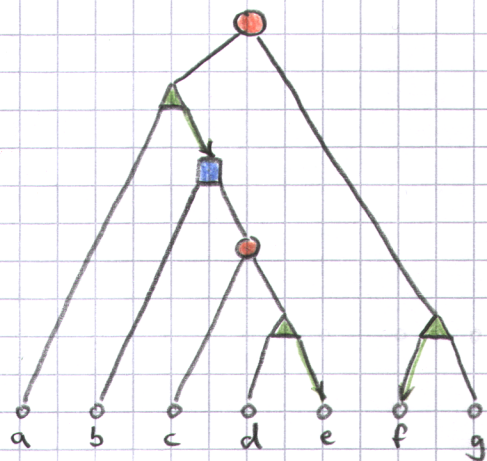
event labeled gene tree



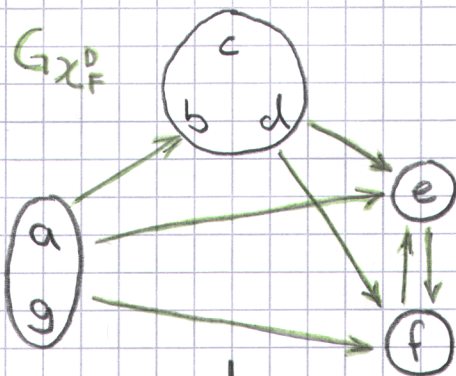
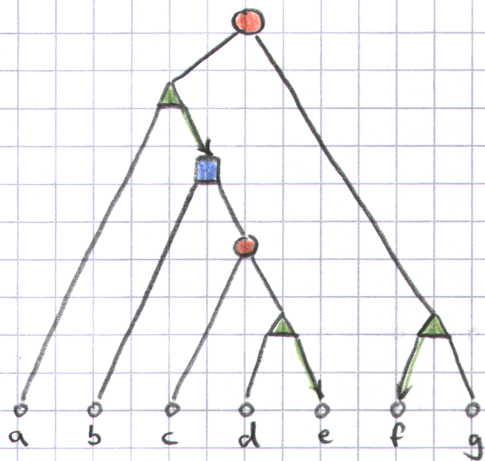
# Reconstruction of event labeled gene trees



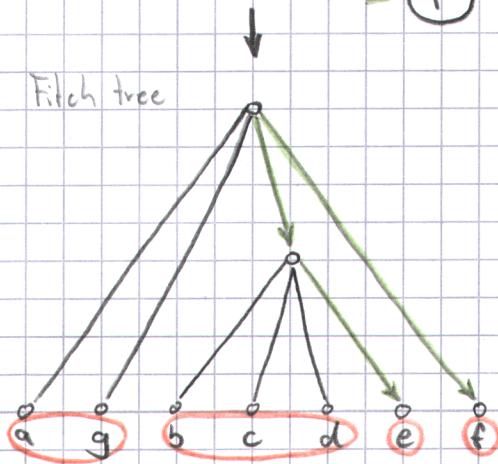
# Reconstruction of event labeled gene trees



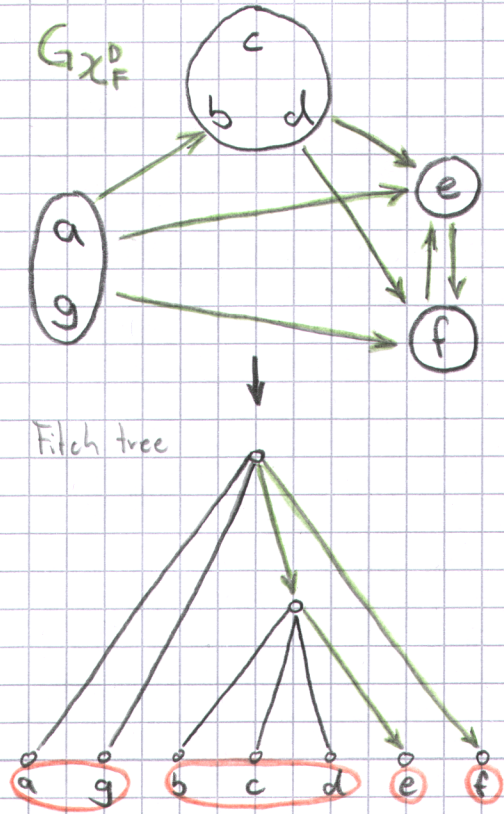
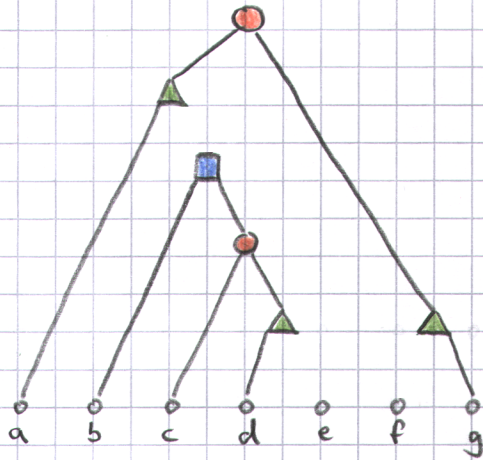
# Reconstruction of event labeled gene trees



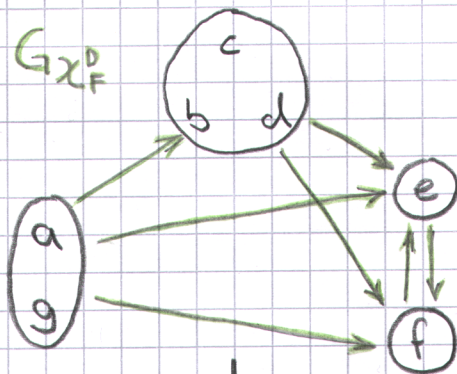
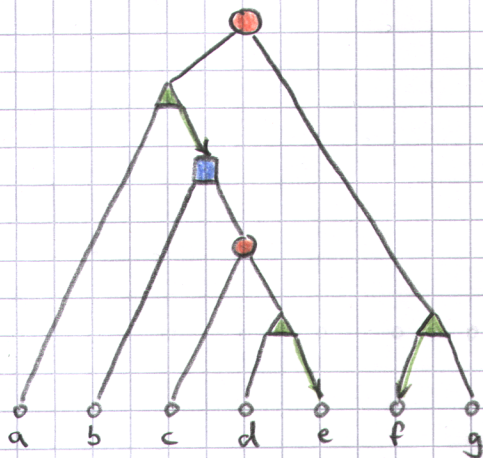
Fitch tree



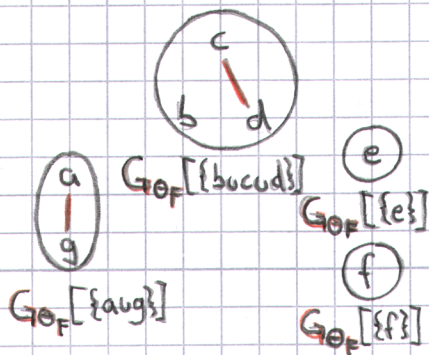
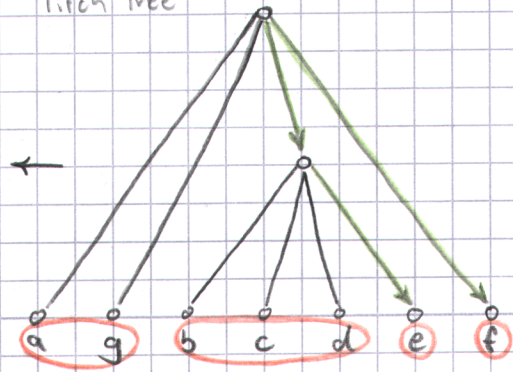
# Reconstruction of event labeled gene trees



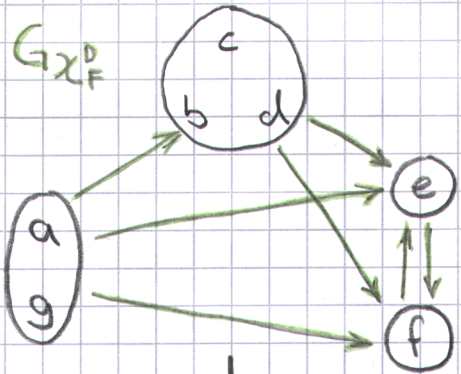
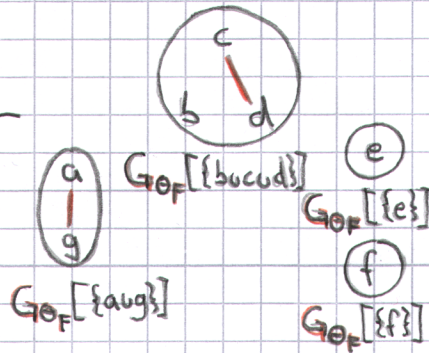
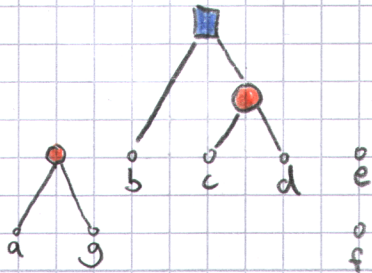
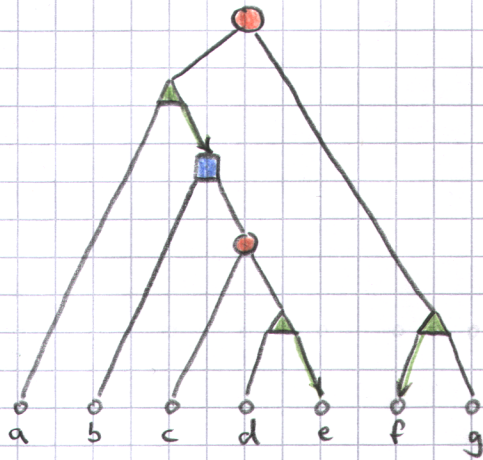
# Reconstruction of event labeled gene trees



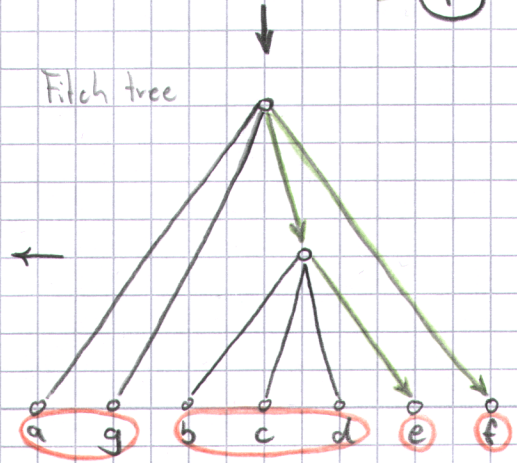
Fitch tree



# Reconstruction of event labeled gene trees

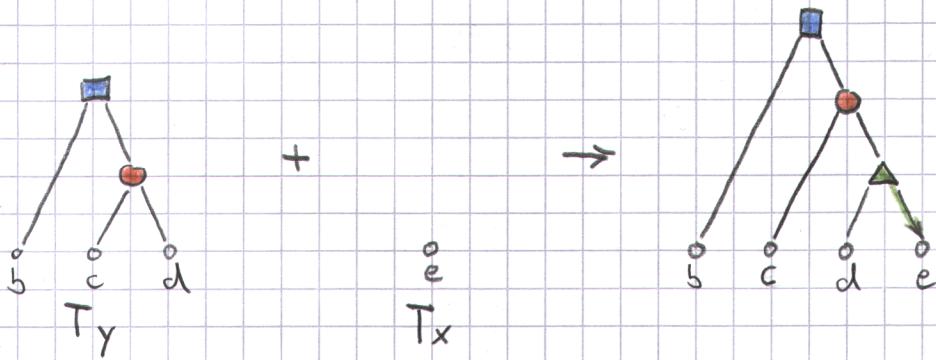


Fitch tree



# Reconstruction of event labeled gene tree

How to append tree  $T_x$  to  $T_y$



□ find subtree  $T_y(v)$  with  $d(y,x) \ll d(y',x)$  for all  $y \in L(T_y(v))$ ,  $y' \in L(T_y) \setminus L(T_y(v))$  and  $x \in L(T_x)$

□ append  $T_x$  to  $T_y$  at edge  $(u,v)$  with HGT towards  $T_x$

## Summary

- improved orthology / xenology inference
- method for constructing event labeled gene trees from Fitch relations  $\Theta_F, \chi_F^D$

## Open Problems

- method relies on molecular clock-like data
- understand the connection between  $\Theta_F, \chi_F^D$  and species trees
- infer species trees
- construct time-consistent reconciliation map



Thank You...



