Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
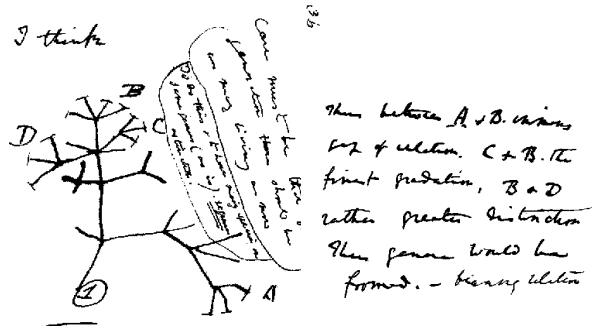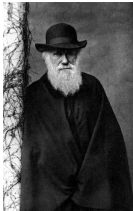00000000

ParaPhylo
000000000

# Bioinformatics
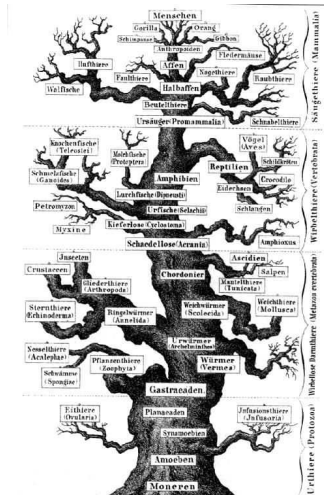## (Phylogenetic Tree Reconstruction)

Marc Hellmuth

# Phylogenetic Reconstruction



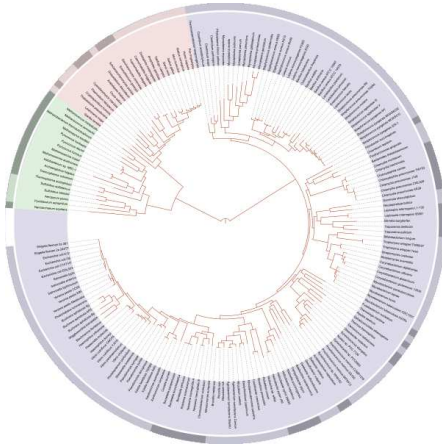"I think" by Charles Darwin (1837) - One of the first evolutionary trees.

# Tree of Live - A Better Picture



Ernst Haeckel, 1879

# Tree of Live - A Better Picture*

Relationship between species with sequenced genomes as of 2006.



center = last universal ancestor of all life on earth.
three domains of life:
eukaryota (animals, plants and fungi);
bacteria;
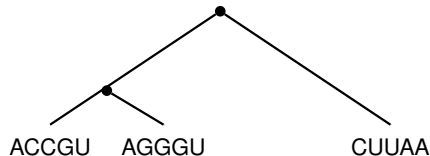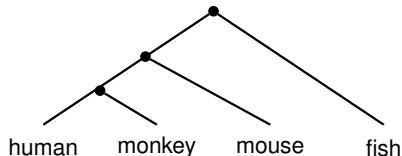archaea.

---

*Ciccarelli, FD (2006). "Toward automatic reconstruction of a highly resolved tree of life.". Science; Letunic, I (2007). "Interactive Tree Of Life (iTOL): an online tool for phylogenetic tree display and annotation.". Bioinformatics

**Aim:** Assemble a tree representing a hypothesis about the evolutionary history of a set of genes, species or other taxa.

Trees are "good" approximation (does not work if one has hybridization)

A phylogenetic tree on set of taxa $X$ is tupel $(T, \lambda)$ s.t. $T = (V, E)$ is unordered tree with unique labels $\lambda(v) \in X$ for all leaves $v \in L \subseteq V$.

# Rooted vs. Unrooted



Unrooted tree (right) "displays" all three rooted trees on three leaves.

Depending on the application, phylogenetic trees may:

- be rooted or unrooted
- have weighted or unweighted edges
- have bounded degree
  (maximum nr of children of each internal node)

The problem in practise:

- Inference of the gene or species tree *T* is a classical problem of molecular phylogenetics.
  In practice it can only be solved approximately.

- Only the subset of leaves of the species or gene tree corresponding to extant (currently living) species or genes in extant (currently living) species is observable.

- All internal nodes (and the event labeling *t*) in the gene tree must be inferred from data.
  events: duplication, speciation (Later!)

## Lemma

*There are* $(2n-3)!! = 1 \cdot 3 \cdot \cdots \cdot (2n-3)$ *rooted trees with n leaves, and* $(2n-5)!!$ *unrooted trees with n leaves*

| | *n* | 3 | 4 | 5 | 6 | 10 | 20 |
|---|---|---|---|---|---|---|---|
| Exmpl: | unrooted | 1 | 3 | 15 | 105 | 2'027'025 | $2.22 \cdot 10^{20}$ |
| | rooted | 3 | 15 | 105 | 945 | 34'459'425 | $8.20 \cdot 10^{21}$ |

Intro
0000000●

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

**Aim:** Assemble a tree representing a hypothesis about the evolutionary history of a set of genes, species or other taxa.

**Methods:**

- Distance Based e.g.:

    - Ultrametric Tree Reconstruction
    - Additive Tree Reconstruction

- Character Based e.g.:

    - Parsimony Methods
    - Maximum Likelihood

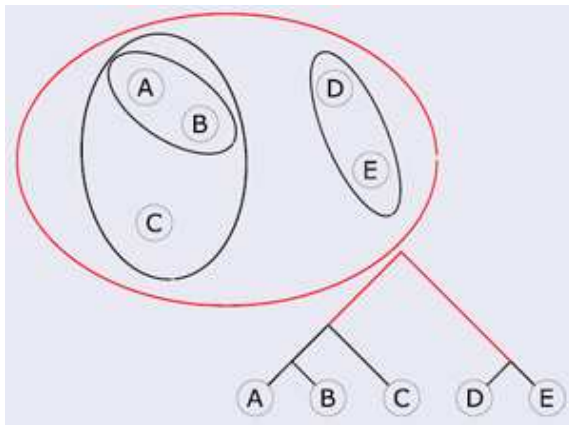- Consensus Methods e.g.:

    - BUILD

# UPGMA

**U**nweighted **P**air **G**roup **M**ethod with **A**rithmetic Mean

- Assume "constant moleculare clock":
  one assumes that mutations always appear with the same probability
  independent from time, location, kind of mutation (mutation = bygone past
  time)

- The two sequences with with the shortest evolutionary distance between
  them are assumed to have been the last that diverged, and represented by
  the most recent internal node.

- Cluster the data and at each step merge clusters.

- Distances between clusters:

$$D(C_i, C_j) = \frac{1}{|C_i||C_j|} \sum_{x \in C_i, y \in C_j} D_{x,y}$$

- Moreover, compute "ultrametric trees".
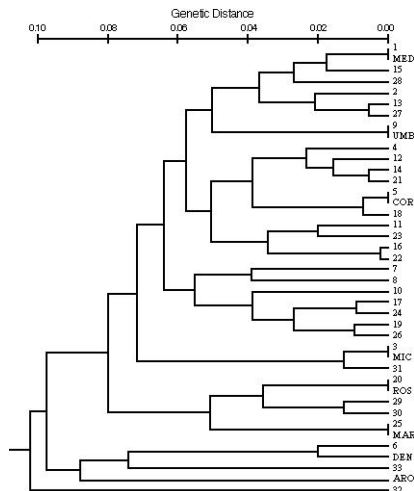
# UPGMA - Idea



It works correctly, if the underlying "distance-matrix" is an ultrametric

A metric $D$ on $M = \{1, \ldots, n\}$ is an ultrametric if for all $x, y, z \in M$ holds

$$D_{xy} \leq \max\{D_{xz}, D_{zy}\}.$$
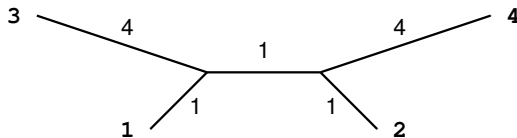
# Example: Ultrametric Tree [†]



_____

[†]taken from: Evolution of polyploid agamic complexes with examples from Antennaria (Asteraceae), RJ Bayer, Opera Bot, 1996

## Neighbor Joining and Additive Tree

For a given $n \times n$ distance matrix $D$ an additive tree $T$ for $D$ is an unrooted tree with

1. $T$ is binary, having $n$ leaves (bijectively labeled by $1, \ldots, n$)

2. each edge $(x, y)$ of $T$ is (positive) weighted with branch length $b_{xy}$

3. For any pair of leaves $i, j$ it holds: $D_{ij} =$ sum of edge weights $b_{xy}$ along path from $i$ to $j$ in $T$.

$$D = \begin{pmatrix} 0 & 3 & 5 & 6 \\ & 0 & 6 & 5 \\ & & 0 & 9 \\ & & & 0 \end{pmatrix}$$

Intro
0000000

Distance Based
0000●

Consensus Methods
0000000000000

Phylo with Event Relations
00000

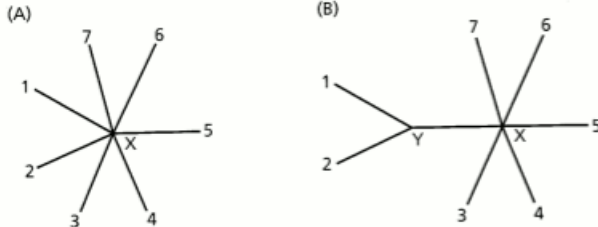Phylo with Event Relations II
00000000

ParaPhylo
000000000

# Neighbor Joining (NJ)

NJ does not assume constant molecular clock.

Basis of NJ is concept of minimum evolution, that is, the "true" tree will be that for which the total branch length is shortest.

**Idea:** Start with "star" tree and separate stepwisely vertices that are together "quite" close and also "quite" far away from the rest until a fully resolved tree has been built. (Note, these two vertices are not necessarily the nearest ones).



It works correctly, if the underlying "distance-matrix" is additive
A metric $D$ on $M = \{1, \ldots, n\}$ is additive if for all $x, y, a, b \in M$ holds

$$D_{xy} + D_{ab} \leq \max\{D_{xa} + D_{yb}, D_{xb} + D_{ya}\}.$$

# Consensus Methods[‡]

Assume a set T of phylogenetic trees has already been constructed.

Aim: Summarize the information in T in the "best way".

"best way" := find largest subtree, find supertree, ...

---

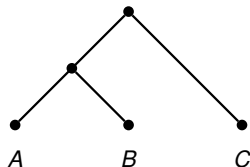[‡]parts of this section are based on talk by Jesper Jansson (2010 MSP Annual Convention)

# Supertree

Aim: Merge a given set of (possibly conflicting) phylogenetic trees into one tree.
Keep as much branching information as possible!
Motivation:

- Combine many trees constructed from different data sets.
  $\rightarrow$ more reliable answers.

- Computationally expensive methods can yield highly accurate trees for small, overlapping subsets of the objects.

- Most individual studies investigate relatively few species. Supertrees allow us to deduce new evolutionary relationships.

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# Rooted Triples

Rooted triplet= rooted binary phylogenetic tree with exactly three leaves.



For three leaves $A, B, C$ in $T$ we write $((A, B), C)$ if the path from $A$ to $B$ does not intersect the path from $C$ to the root $\rho$.
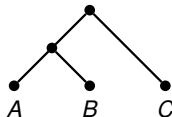
That is the unique rooted triplet with

$$lca(A, B) \prec lca(A, C) = lca(B, C)$$

Any rooted phylogenetic tree can be represented by a set of rooted triplets.

Intro
0000000
Distance Based
00000
Consensus Methods
0000●00000000
Phylo with Event Relations
00000
Phylo with Event Relations II
00000000
ParaPhylo
000000000

## Combining Rooted Triples

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



Consensus Tree "displays" all rooted triples:

Intro
0000000
Distance Based
00000
Consensus Methods
0000●00000000
Phylo with Event Relations
00000
Phylo with Event Relations II
00000000
ParaPhylo
000000000

## Combining Rooted Triples

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



Consensus Tree does not always exist!!

# Consistence



For three leaves $A, B, C$ in $T$ we write $((A, B), C)$ if the path from $A$ to $B$ does not intersect the path from $C$ to the root $\rho$.

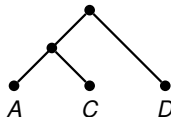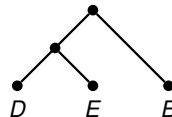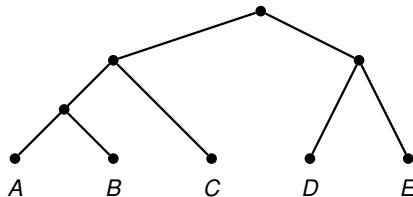That is the unique rooted triplet with

$$lca(A, B) \prec lca(A, C) = lca(B, C)$$

$T$ and an arbitrary triple $((A, B), C)$ are consistent iff

$$lca(A, B) \prec lca(A, C) = lca(B, C)$$

$T$ displays $((A, B), C)$.

## BUILD

Theorem (Aho, Sagiv, Szymanski, Ullman - 1981; Semple & Steel - 2003)

*Let $\mathscr{R}$ by a collection of rooted triples with leaf set $\mathscr{L}$. Then there is an $O(|\mathscr{R}||\mathscr{L}|)$ time algorithm – called* BUILD *– that either*

- *constructs a phylogenetic tree $T_{|\mathscr{R}}$ that displays each member of $\mathscr{R}$*

  *or*

- *recognizes $\mathscr{R}$ as inconsistent.*

Intro
0000000

Distance Based
00000

Consensus Methods
0000000●00000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

## BUILD

**Idea of this recursive, top-down approach:** Partition $\mathscr{L}$ into blocks according to $\mathscr{R}$. Output a tree consisting of a root whose children are roots of the trees obtained by recursing on each block.

# BUILD

Let $\mathscr{R}$ be a set of triples defined on a leaf set $\mathscr{L}$.

For any $L \subseteq \mathscr{L}$ define $\mathscr{R}_{|L} = \{((x,y)z) \in \mathscr{R} \mid x,y,z \in L\}$.

To find blocks use auxiliary graph $G(\mathscr{R}_{|L}, L) = (L, E)$ with $(x,y) \in E$ iff there is a triple $((x,y)z) \in \mathscr{R}_{|L}$

# BUILD

Let $\mathscr{R}$ be a set of triples defined on a leaf set $\mathscr{L}$.

For any $L \subseteq \mathscr{L}$ define $\mathscr{R}_{|L} = \{((x,y)z) \in \mathscr{R} \mid x,y,z \in L\}$.

To find blocks use auxiliary graph $G(\mathscr{R}_{|L}, L) = (L, E)$ with $(x,y) \in E$ iff there is a triple $((x,y)z) \in \mathscr{R}_{|L}$

Exmpl: $L = \{A, B, C\}$, $\mathscr{R} = ((A,B)C)$, $G(\mathscr{R}_{|L}, L)$

# BUILD

Let $\mathscr{R}$ be a set of triples defined on a leaf set $\mathscr{L}$.

For any $L \subseteq \mathscr{L}$ define $\mathscr{R}_{|L} = \{((x,y)z) \in \mathscr{R} \mid x,y,z \in L\}$.

To find blocks use auxiliary graph $G(\mathscr{R}_{|L}, L) = (L, E)$ with $(x,y) \in E$ iff there is a triple $((x,y)z) \in \mathscr{R}_{|L}$

Exmpl: $L = \{A, B, C\}$, $\mathscr{R} = ((A,B)C)$, $G(\mathscr{R}_{|L}, L)$



**Crucial observation:** If $((xy)z)$ is consistent with a tree $T$ then the leaves labeled by $x$ and $y$ cannot descend from two different children of the root of $T$, i.e., $x$ and $y$ must belong to the same block.

# BUILD

Let $\mathscr{R}$ be a set of triples defined on a leaf set $\mathscr{L}$.

For any $L \subseteq \mathscr{L}$ define $\mathscr{R}_{|L} = \{((x,y)z) \in \mathscr{R} \mid x,y,z \in L\}$.

To find blocks use auxiliary graph $G(\mathscr{R}_{|L}, L) = (L, E)$ with $(x,y) \in E$ iff there is a triple $((x,y)z) \in \mathscr{R}_{|L}$
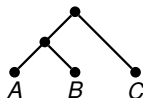
Exmpl: $L = \{A, B, C\}$, $\mathscr{R} = ((A,B)C)$, $G(\mathscr{R}_{|L}, L)$
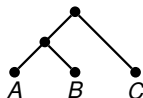


**Crucial observation:** If $((xy)z)$ is consistent with a tree $T$ then the leaves labeled by $x$ and $y$ cannot descend from two different children of the root of $T$, i.e., $x$ and $y$ must belong to the same block.
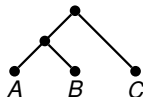
Therefore, the algorithm defines the partition of $L \subseteq \mathscr{L}$ by:
Blocks of leaves iff connected components in $G(\mathscr{R}_{|L}, L)$

## BUILD

### Lemma (Aho, Sagiv, Szymanski, Ullman (1981), Bryant & Steel (1995))

*A given triple set $\mathscr{R}$ on a leaf set $\mathscr{L}$ is consistent if and only if for all $L \subseteq \mathscr{L}$ with $|L| > 1$ the graph $G(\mathscr{R}_{|L}, L)$ is disconnected.*

# BUILD

1: **INPUT:** Set of triples in $\mathscr{R}$, leaf set $\mathscr{L}$.
2: **OUTPUT:** A rooted, phylog. tree distinctly leaf-labeled by $\mathscr{L}$ consistent with all rooted triplets in $\mathscr{R}$, if one exists; otherwise *null*.
3: compute $G(\mathscr{R}, \mathscr{L})$
4: compute connected components $C_1, \ldots, C_s$ of $G(\mathscr{R}, \mathscr{L})$
5: **if** $s = 1$ and $|\mathscr{L}| = 1$ **then**
6:     return tree $\simeq K_1$
7: **else if** $s = 1$ and $|\mathscr{L}| > 1$ **then**
8:     return *null*
9: **else**
10:     **for** i = 1, \ldots s **do**
11:        $T_i = \text{BUILD}(\mathscr{R}_{|V(C_i)}, V(C_i))$
12:     **end for**
13:     **if** $T_i \neq null$ for all $i = 1, \ldots s$ **then**
14:        attach all of these trees to a common parent node and let $T$ be the resulting tree; else $T = null$.
15:     **end if**
16: **end if**

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# BUILD - Example

$\mathcal{R} = \{((AB)C), ((AC)D), ((DE)B)\}$

$G(\mathcal{R}, \mathcal{L})$ :



$\text{BUILD}(\mathcal{R}, \mathcal{L} = \{A, B, C, D, E\})$



$C_1 := \text{BUILD}(\mathcal{R}_{|\mathcal{L}}, \mathcal{L} = \{A, B, C\})$
$C_2 := \text{BUILD}(\mathcal{R}_{|\mathcal{L}}, \mathcal{L} = \{D, E\})$

# BUILD - Example

$\mathscr{R} = \{((AB)C), ((AC)D), ((DE)B)\}$

$C_1 := \text{BUILD}(\mathscr{R}_{|\mathscr{L}}, \mathscr{L} = \{A, B, C\})$
$\mathscr{R}_1 := \{((AB)C)\}$

$C_2 := \text{BUILD}(\mathscr{R}_{|\mathscr{L}}, \mathscr{L} = \{D, E\})$
$\mathscr{R}_2 := \emptyset$

$G(\{A, B, C\}):$



**B** —— **A**

**C**

$\text{BUILD}(\mathscr{R}, \mathscr{L} = \{A, B, C, D, E\})$



$C_1$

$C_2$

$C_{11}$

$C_{12}$    $C_{21}$    $C_{22}$

$G(\{D, E\}):$

**E**    **D**

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000●0

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# BUILD - Example

$\mathscr{R} = \{((AB)C), ((AC)D), ((DE)B)\}$

$C_1 := \mathtt{BUILD}(\mathscr{R}_{|\mathscr{L}}, \mathscr{L} = \{A, B, C\})$

$C_2 := \mathtt{BUILD}(\mathscr{R}_{|\mathscr{L}}, \mathscr{L} = \{D, E\})$

$C_{11} := \mathtt{BUILD}(\mathscr{R}_{|\mathscr{L}}, \mathscr{L} = \{A, B\})$
$C_{12} := \mathtt{BUILD}(\emptyset, \{C\})$
$C_{21} := \mathtt{BUILD}(\emptyset, \{D\})$
$C_{22} := \mathtt{BUILD}(\emptyset, \{E\})$

$G(\{A, B, C\}):$



$\mathtt{BUILD}(\mathscr{R}, \mathscr{L} = \{A, B, C, D, E\})$



$G(\{D, E\}):$

Intro
0000000

Distance Based
00000

Consensus Methods
00000000000●0

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# BUILD - Example



$\texttt{BUILD}(\mathscr{R}, \mathscr{L} = \{A, B, C, D, E\})$

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000●

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# BUILD - Example



Consensus Tree does not always exist!!

$G(\mathscr{R}, \mathscr{L})$ :

Phylogenetics with Evolutionary Event Relations

Intro
○○○○○○○

Distance Based
○○○○○

Consensus Methods
○○○○○○○○○○○○○○

**Phylo with Event Relations**
○●○○○○

Phylo with Event Relations II
○○○○○○○○

ParaPhylo
○○○○○○○○○

# The "true" evolutionary History

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
0●0000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# The "true" evolutionary History

- species are characterized by its genome: a "bag of genes"

- "Genes" evolve along a *rooted* tree with unique *event labeling* $t : V^0 \rightarrow M = \{\bullet, \blacksquare, \blacktriangle\}$

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000000

Phylo with Event Relations
0●0000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# The "true" evolutionary History



- species are characterized by its genome: a "bag of genes"

- "Genes" evolve along a *rooted* tree with unique *event labeling* $t : V^0 \rightarrow M = \{\bullet, \blacksquare, \blacktriangle\}$

- ■ 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

## The Problem in Practice



**a          b1 b2 b3     c1 c2 c3      d**
**A           B           C            D**

- Only the subset of leaves of the gene tree corresponding to genes in extant (currently living) species is observable.

- All internal nodes and the event labelling $t$ in the gene tree must be inferred from data.

- We cannot observe and reconstruct all events (losses).

- The events and the topology of the gene tree can be used (under several constraints) to infer the species tree (Reconciliation)

## The Problem in Practice



- Only the subset of leaves of the gene tree corresponding to genes in extant (currently living) species is observable.

- All internal nodes and the event labelling *t* in the gene tree must be inferred from data.

- We cannot observe and reconstruct all events (losses).

- The events and the topology of the gene tree can be used (under several constraints) to infer the species tree (Reconciliation)

## The Problem in Practice



- Only the subset of leaves of the gene tree corresponding to genes in extant (currently living) species is observable.

- All internal nodes and the event labelling *t* in the gene tree must be inferred from data.

- We cannot observe and reconstruct all events (losses).

- The events and the topology of the gene tree can be used (under several constraints) to infer the species tree (Reconciliation)

# The Problem in Practice



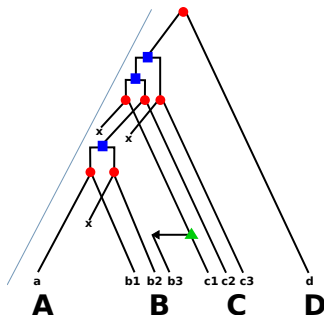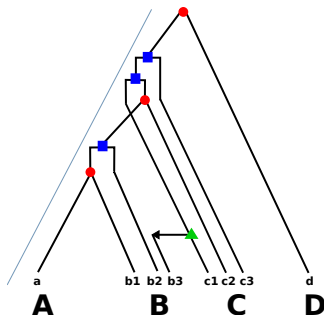- Only the subset of leaves of the gene tree corresponding to genes in extant (currently living) species is observable.

- All internal nodes and the event labelling $t$ in the gene tree must be inferred from data.

- We cannot observe and reconstruct all events (losses).

- The events and the topology of the gene tree can be used (under several constraints) to infer the species tree (Reconciliation)

Intro
○○○○○○○

Distance Based
○○○○○

Consensus Methods
○○○○○○○○○○○○○○○○

Phylo with Event Relations
○○○●○

Phylo with Event Relations II
○○○○○○○○

ParaPhylo
○○○○○○○○○

# State-of-the-Art Tree Reconstruction



- Find 1:1-orthologs.
  - Paralogs = dangerous nuisance that has to be detected and removed.
  - Select families of genes that rarely exhibit duplications
    (e.g. rRNAs, ribosomal proteins)

# State-of-the-Art Tree Reconstruction



- Find 1:1-orthologs.

    - Paralogs = dangerous nuisance that has to be detected and removed.
    - Select families of genes that rarely exhibit duplications
      (e.g. rRNAs, ribosomal proteins)

- Alignments of protein or DNA sequences and standart techniques yield
  evolutionary history that is believed to be congruent to that of the respective
  species.

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000000

Phylo with Event Relations
000●0

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# State-of-the-Art Tree Reconstruction



Pitfalls:

- Information of evolutionary events as paralogs or xenologs is ignored, although they might contain valuable information about the evolutionary history of the species.

- The set of usable gene sets is strongly restricted ($\leq 10\%$).

# State-of-the-Art Tree Reconstruction



Pitfalls:

- Information of evolutionary events as paralogs or xenologs is ignored, although they might contain valuable information about the evolutionary history of the species.
- The set of usable gene sets is strongly restricted ($\leq 10\%$).

Thus, to get a better picture of the species evolution we try to include also the information of paralogs and xenologs.

# Tree-Representable Sets of Binary Relations



An ordered pair $(x, y)$ of two genes comprises

- orthologs if $\mathrm{lca}(x, y) = \bullet = speciation$
- paralogs if $\mathrm{lca}(x, y) = \blacksquare = duplication$
- xenologs if $\mathrm{lca}(x, y) = \blacktriangle = HGT$ and $\blacktriangle$ "points from" $x$ to $y$ in $T$

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000●

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# Tree-Representable Sets of Binary Relations



The gene-tree determines three distinct relations

- $R_{\bullet}$, the orthologs (lca$(x, y) = \bullet$)
- $R_{\blacksquare}$, the paralogs (lca$(x, y) = \blacksquare$)
- $R_{\blacktriangle}$, the xenologs ( lca$(x, y) = \blacktriangle$, $\blacktriangle$ "points from" $x$ to $y$ in $T$)

# Tree-Representable Sets of Binary Relations



Orthologs, Paralogs (and to some extent HGT) can be estimated without inferring a gene- or species trees.

Assume we have *estimated* binary relations $R_1, \ldots, R_k$ s.t.

$$(xy) \in R_i \text{ iff } \mathrm{lca}(xy) = i \text{ in ordered tree } T$$

Thus, it is important to understand, when these estimates $R_1, \ldots, R_k$ can be "represented" in a single tree — thus, the edge-colored graph-representation.

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
●0000000

ParaPhylo
000000000

# Sketch: Estimating $R_\bullet$ directly from the Data

- Simplify: No losses, No HGT // $T$ gene tree, $S$ species tree

- Let $d_S(A, B)$ be divergence time of species $A, B$.

- $y \in B$ is orthologous to $x \in A$, if

  1. $A \neq B$,

     orthologs are never found in the same
     species

  2. $d_T(x, y) = d_S(A, B)$,
     divergence time of $x$ and $y$ must be equal to
     $d_S(A, B)$.



● speciation
■ duplication

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
●0000000

ParaPhylo
000000000

# Sketch: Estimating $R_\bullet$ directly from the Data

- Simplify: No losses, No HGT // $T$ gene tree, $S$ species tree

- Let $d_S(A, B)$ be divergence time of species $A, B$.

- $y \in B$ is orthologous to $x \in A$, if

  

  ● speciation
  ■ duplication

  1. $A \neq B$,

     orthologs are never found in the same species

  2. $d_T(x, y) = d_S(A, B)$,
     divergence time of $x$ and $y$ must be equal to $d_S(A, B)$.

- If no losses, then for each $x \in A$ there is an orthologous gene $y \in B$.

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
●0000000

ParaPhylo
000000000

# Sketch: Estimating $R_\bullet$ directly from the Data

- Simplify: No losses, No HGT // $T$ gene tree, $S$ species tree

- Let $d_S(A, B)$ be divergence time of species $A, B$.

- $y \in B$ is orthologous to $x \in A$, if

  - ● speciation
  - ■ duplication

  1. $A \neq B$,

     orthologs are never found in the same species

  2. $d_T(x, y) = d_S(A, B)$,
     divergence time of $x$ and $y$ must be equal to $d_S(A, B)$.



- If no losses, then for each $x \in A$ there is an orthologous gene $y \in B$.

- If no HGT, then $d_T(x, y) < d_S(A, B)$ is not possible.

# Sketch: Estimating $R_\bullet$ directly from the Data

- Simplify: No losses, No HGT // $T$ gene tree, $S$ species tree

- Let $d_S(A, B)$ be divergence time of species $A, B$.

- $y \in B$ is orthologous to $x \in A$, if

  ● speciation
  ■ duplication

  1. $A \neq B$,

     orthologs are never found in the same
     species

  2. $d_T(x, y) = d_S(A, B)$,
     divergence time of $x$ and $y$ must be equal to
     $d_S(A, B)$.



- If no losses, then for each $x \in A$ there is an orthologous gene $y \in B$.

- If no HGT, then $d_T(x, y) < d_S(A, B)$ is not possible.

- If $d_T(x, y) > d_S(A, B)$, then $x, y$ must be paralogs

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
●0000000

ParaPhylo
000000000

# Sketch: Estimating $R_\bullet$ directly from the Data

- Simplify: No losses, No HGT // $T$ gene tree, $S$ species tree

- Let $d_S(A, B)$ be divergence time of species $A, B$.

- $y \in B$ is orthologous to $x \in A$, if

  ● speciation
  ■ duplication

  1. $A \neq B$,

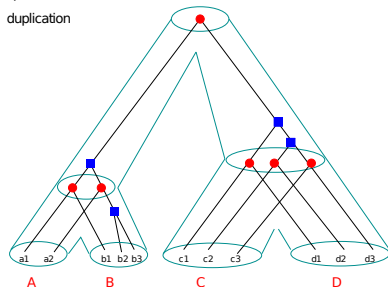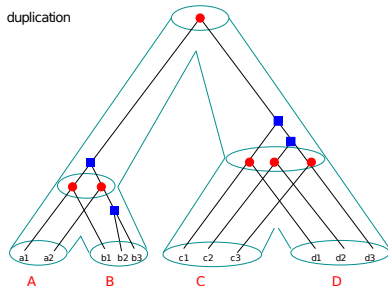     orthologs are never found in the same species

  2. $d_T(x, y) = d_S(A, B)$,
     divergence time of $x$ and $y$ must be equal to $d_S(A, B)$.



Set of orth. genes in $B$ for $x \in A$

$$R_\bullet(x, B) = \{y \in B \mid d_T(x, y) = \min_{z \in B} d_T(z, x)\}$$

For all $x \in A$, $y \in B$

$$y \in R_\bullet(x, B) \iff x \in R_\bullet(y, A), \text{then } (x, y) \in R_\bullet$$

## Sketch: Estimating $R_\bullet$ directly from the Data

- We don't know the true divergence time $\Rightarrow$ genetic distance / similarity scores

- We know the assignment of genes to species and we can measure similarity $s(x, y)$ of two genes using sequence alignments and `blast` bit scores

- $y \in B$ is a (putative) ortholog of $x \in A$, in symbols $(x, y) \in \widehat{R_\bullet}$, if

  1. $A \neq B$,
     orthologs are never found in the same species

  2. $s(x, y) \approx \max\limits_{z \in B} s(x, z) \approx \max\limits_{z \in A} s(z, y)$,
     if $x$ and $y$ are orthologs, then they do not have (much) closer relatives in the two species.



● speciation
■ duplication

# Orthologs and Paralogs

$\Rightarrow$ we get an estimate $\widehat{R}_{\bullet}$ of the true relation $R_{\bullet}$

# Orthologs and Paralogs

$\Rightarrow$ we get an estimate $\widehat{R}_\bullet$ of the true relation $R_\bullet$

An estimate $\widehat{R}_\bullet$ is *valid* iff there is a tree-representation $T$ with

- $\mathrm{lca}(x, y) = \bullet = $ *speciation* for all $(x, y) \in \widehat{R}_\bullet$ and

- $\mathrm{lca}(x, y) = \blacksquare = $ *duplication* for all $(x, y) \in \widehat{R}_\blacksquare \sim \overline{\overline{\widehat{R}_\bullet}}$

# Orthologs and Paralogs

$\Rightarrow$ we get an estimate $\widehat{R}_{\bullet}$ of the true relation $R_{\bullet}$

An estimate $\widehat{R}_{\bullet}$ is *valid* iff there is a tree-representation $T$ with

• $\text{lca}(x,y) = \bullet = speciation$ for all $(x,y) \in \widehat{R}_{\bullet}$ and

• $\text{lca}(x,y) = \blacksquare = duplication$ for all $(x,y) \in \widehat{R}_{\blacksquare} \sim \overline{\widehat{R}_{\bullet}}$



$G_{\widehat{R}_{\bullet}}$ with edge set $\widehat{R}_{\bullet} = \{(v0,v2),(v0,v4),(v2,v3),(v3,v4)\}$

# Orthologs and Paralogs

$\Rightarrow$ we get an estimate $\widehat{R}_\bullet$ of the true relation $R_\bullet$

An estimate $\widehat{R}_\bullet$ is *valid* iff there is a tree-representation $T$ with

- $\mathrm{lca}(x,y) = \bullet = $ *speciation* for all $(x,y) \in \widehat{R}_\bullet$ and

- $\mathrm{lca}(x,y) = \blacksquare = $ *duplication* for all $(x,y) \in \widehat{R}_\blacksquare \sim \overline{\overline{\widehat{R}_\bullet}}$

**Question:** When are estimates $\widehat{R}_\bullet$ valid?

**Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

# Orthologs and Paralogs

$\Rightarrow$ we get an estimate $\widehat{R}_\bullet$ of the true relation $R_\bullet$

An estimate $\widehat{R}_\bullet$ is *valid* iff there is a tree-representation $T$ with

- $\text{lca}(x,y) = \bullet = speciation$ for all $(x,y) \in \widehat{R}_\bullet$ and

- $\text{lca}(x,y) = \blacksquare = duplication$ for all $(x,y) \in \widehat{R}_\blacksquare \sim \overline{\widehat{R}_\bullet}$

**Question:** When are estimates $\widehat{R}_\bullet$ valid?

## Theorem (2013)

*The estimate $\widehat{R}_\bullet$ (and $\widehat{R}_\blacksquare$) is valid* $\quad \Leftrightarrow \quad G_{\widehat{R}_\bullet}$ *is $P_4$-free = Cograph*

---

**Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
000●0000

ParaPhylo
000000000

Look at all possible gene trees that encode $R_\bullet, R_\blacksquare$ on on some set $X$, $|X| = 4$

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
000●0000

ParaPhylo
000000000

Look at all possible gene trees that encode $R_\bullet, R_\blacksquare$ on on some set $X$, $|X| = 4$



All symmetric relations $R_\bullet, R_\blacksquare$ have a tree-representation, except:



$A - B, B - C, C - D \in R_\bullet$

$A - C, A - D, B - D \in R_\blacksquare \sim \overline{R_\bullet}$

# Cograph (=Complement reducible graph)

Corneil et al., 1981:

Cographs are defined recursively (blackboard)

$G$ is Cograph IFF $G$ is "induced $P_4$-free"

Forbidden:



Allowed:





**Complement reducible graphs**, Corneil DG, Lerchs H, Steward Burlingham L, *Discr. Appl. Math.*, 1981

Intro
0000000
Distance Based
00000
Consensus Methods
000000000000
Phylo with Event Relations
00000
Phylo with Event Relations II
00000●000
ParaPhylo
000000000

# Cograph (=Complement reducible graph)

Corneil et al., 1981:

Cographs are defined recursively (blackboard)

*G* is Cograph IFF *G* is "induced $P_4$-free"

Every Cograph is associated with a unique Cotree.

# Cograph (=Complement reducible graph)

Corneil et al., 1981:

Cographs are defined recursively (blackboard)

$G$ is Cograph IFF $G$ is "induced $P_4$-free"

Every Cograph is associated with a unique Cotree.

# Cograph (=Complement reducible graph)

Corneil et al., 1981:

Cographs are defined recursively (blackboard)

*G* is Cograph IFF *G* is "induced $P_4$-free"

Every Cograph is associated with a unique Cotree.

# Cograph (=Complement reducible graph)

Corneil et al., 1981:

Cographs are defined recursively (blackboard)

*G* is Cograph IFF *G* is "induced $P_4$-free"

Every Cograph is associated with a unique Cotree.

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000●000

ParaPhylo
000000000

# Cograph (=Complement reducible graph)

Corneil et al., 1981:

Cographs are defined recursively (blackboard)

*G* is Cograph IFF *G* is "induced $P_4$-free"

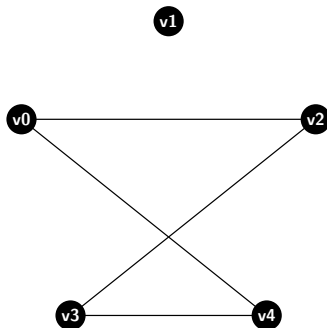Every Cograph is associated with a unique Cotree.

# Cograph (=Complement reducible graph)

Corneil et al., 1981:

Cographs are defined recursively (blackboard)

$G$ is Cograph IFF $G$ is "induced $P_4$-free"

Every Cograph is associated with a unique Cotree.



$$(x, y) \in E(G_{\widehat{R}_\bullet}) \text{ if and only if } lca(x, y) = 1 = \bullet$$

# Orthologs and Paralogs

An estimate $\widehat{R}_\bullet$ is valid iff there is a tree-representation $T$ (with event-label $t$) with

- $t(\text{lca}(x, y)) = \bullet = \textit{speciation}$ for all $(x, y) \in \widehat{R}_\bullet$ and

- $t(\text{lca}(x, y)) = \blacksquare = \textit{duplication}$ for all $(x, y) \in \widehat{R}_\blacksquare \sim \overline{\widehat{R}_\bullet}$

## Theorem (2013)

*The estimate $\widehat{R}_\bullet$ (and $\widehat{R}_\blacksquare$) is valid* $\quad \Leftrightarrow \quad$ $G_{\widehat{R}_\bullet}$ *is a Cograph ($P_4$-free)*

‡**Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

# Orthologs and Paralogs

An estimate $\widehat{R}_\bullet$ is valid iff there is a tree-representation $T$ (with event-label $t$) with

- $t(\mathrm{lca}(x,y)) = \bullet = $ *speciation* for all $(x,y) \in \widehat{R}_\bullet$ and

- $t(\mathrm{lca}(x,y)) = \blacksquare = $ *duplication* for all $(x,y) \in \widehat{R}_\blacksquare \sim \overline{\widehat{R}_\bullet}$

## Theorem (2013)

*The estimate $\widehat{R}_\bullet$ (and $\widehat{R}_\blacksquare$) is valid* $\quad\Leftrightarrow\quad$ $G_{\widehat{R}_\bullet}$ **is a Cograph** *($P_4$-free)*

The cotree (= least resolved gene tree) can then be computed in linear time.

[‡] **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

**a**          **b1 b2**    **c1 c2 c3**    **d**
**A**          **B**        **C**           **D**

Given valid relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ (there is no HGT)

Given valid relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ (there is no HGT) $\rightarrow$ event-labeled gene tree

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

Given valid relations $\widehat{R}_{\bullet}$ and $\widehat{R}_{\blacksquare}$ (there is no HGT) $\rightarrow$ event-labeled gene tree

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
0000000●0

ParaPhylo
000000000

Given valid relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ (there is no HGT) $\rightarrow$ event-labeled gene tree

**Question:** When does there exist a species tree for a given gene tree

**Answer:** BLACKBOARD + next slides

Intro
○○○○○○○

Distance Based
○○○○○

Consensus Methods
○○○○○○○○○○○○○○○

Phylo with Event Relations
○○○○○

**Phylo with Event Relations II**
○○○○○○●○

ParaPhylo
○○○○○○○○○

Given valid relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ (there is no HGT) $\rightarrow$ event-labeled gene tree

**Question:** When does there exist a species tree for a given gene tree

**Answer:** BLACKBOARD + next slides

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

Given valid relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ (there is no HGT) $\rightarrow$ event-labeled gene tree

**Question:** When does there exist a species tree for a given gene tree and a reconciliation map $\mu$ between them?

**Answer:** BLACKBOARD + next slides

## Triples

For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
0000000●0

ParaPhylo
000000000

## Triples

For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.



$$\mathscr{R}(T) = \{ab_1|x \text{ with } x = b_2, c_1, c_2, c_3, d;$$
$$ab_2|x \text{ with } x = c_1, c_2, c_3, d;$$
$$b_1b_2|x \text{ with } x = c_1, c_2, c_3, d;$$
$$\dots\}$$

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000000

Phylo with Event Relations
00000

**Phylo with Event Relations II**
00000000

ParaPhylo
000000000

## Triples

For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.

$$
\begin{aligned}
\mathscr{R}(T) = \{ & ab_1|x \text{ with } x = b_2, c_1, c_2, c_3, d; \\
& ab_2|x \text{ with } x = c_1, c_2, c_3, d; \\
& b_1b_2|x \text{ with } x = c_1, c_2, c_3, d; \\
& \dots \}
\end{aligned}
$$

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

## Triples

For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.

We write ab|c$^\bullet$ if ab|c $\in \mathcal{R}(T)$

$$\text{lca}(a, b, c) = \bullet = \text{"speciation"}$$

$$\mathcal{R}(T) = \{\text{ab}_1|\text{x with } x = b_2, c_1, c_2, c_3, d;$$
$$\text{ab}_2|\text{x with } x = c_1, c_2, c_3, d;$$
$$\text{b}_1\text{b}_2|\text{x with } x = c_1, c_2, c_3, d;$$
$$\dots\}$$

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
0000000●0

ParaPhylo
000000000

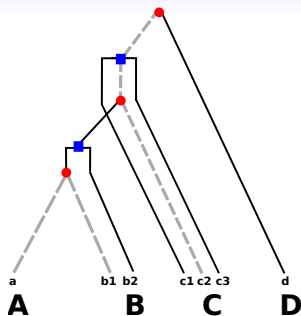## Triples

For three leaves $a, b, c$ in $T$ we write $ab|c$ if the path from $a$ to $b$ does not intersect the path from $c$ to the root.

We write $ab|c^{\bullet}$ if $ab|c \in \mathscr{R}(T)$

$$\mathrm{lca}(a, b, c) = \bullet = \text{``speciation''}$$

Examples: $ab_1|c_2^{\bullet}, ab_1|d^{\bullet}, b_2c_3|d^{\bullet} \ ac_2|d^{\bullet}, \ldots$

Intro
○○○○○○○

Distance Based
○○○○○

Consensus Methods
○○○○○○○○○○○○○

Phylo with Event Relations
○○○○○

**Phylo with Event Relations II**
○○○○○○●○

ParaPhylo
○○○○○○○○○

## Triples

For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.

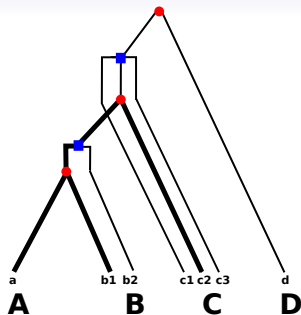We write ab|c$^\bullet$ if ab|c $\in \mathscr{R}(T)$

$$\text{lca}(a, b, c) = \bullet = \text{``speciation''}$$

We know the assignment of genes to the species in which they occur. This gives us the triple set:

$$\mathbb{S} = \{(AB|C : \exists \ \text{ab|c}^\bullet \ \text{with} \ a \in A, b \in B, c \in C\}$$

Examples: $ab_1|c_2^\bullet$, $ab_1|d^\bullet$, $b_2c_3|d^\bullet$ $ac_2|d^\bullet$, ...

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

## Triples
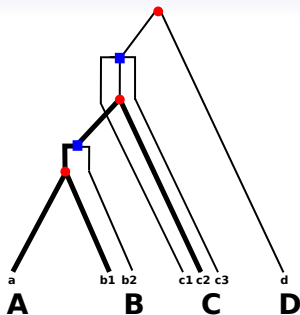
For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.

We write ab|c$^\bullet$ if ab|c $\in \mathscr{R}(T)$

$$\text{lca}(a, b, c) = \bullet = \text{``speciation''}$$

We know the assignment of genes to the species in which they occur. This gives us the triple set:

$$\mathbb{S} = \{(AB|C : \exists \text{ ab|c}^\bullet \text{ with } a \in A, b \in B, c \in C\}$$

Examples: $ab_1|c_2{}^\bullet$, $ab_1|d^\bullet$, $b_2c_3|d^\bullet$ $ac_2|d^\bullet, \ldots$

$$\mathbb{S} = \{(AB|C, AB|D, BC|D, AC|D\}$$

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

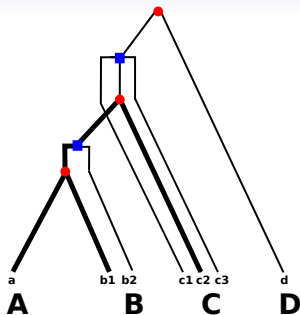**Phylo with Event Relations II**
00000000

ParaPhylo
000000000

## Triples

For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.
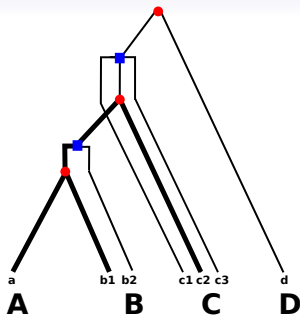
We write ab|c$^\bullet$ if ab|c $\in \mathscr{R}(T)$

$$\text{lca}(a, b, c) = \bullet = \text{``speciation''}$$

We know the assignment of genes to the species in which they occur. This gives us the triple set:

$$\mathbb{S} = \{(AB|C : \exists \text{ ab|c}^\bullet \text{ with } a \in A, b \in B, c \in C\}$$

$$\mathbb{S} = \{(AB|C, AB|D, BC|D, AC|D\}$$

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

**Phylo with Event Relations II**
00000000

ParaPhylo
000000000

## Triples

For three leaves $a, b, c$ in $T$ we write ab|c if the path from $a$ to $b$ does not intersect the path from $c$ to the root.

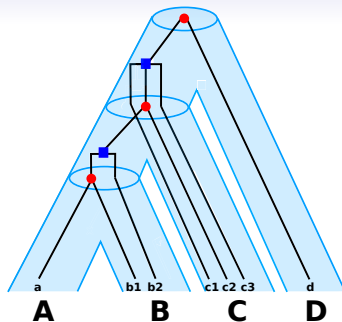We write ab|c$^\bullet$ if ab|c $\in \mathscr{R}(T)$

$$\text{lca}(a, b, c) = \bullet = \text{``speciation''}$$

We know the assignment of genes to the species in which they occur. This gives us the triple set:

$$\mathbb{S} = \{(AB|C : \exists \; ab|c^\bullet \text{ with } a \in A, b \in B, c \in C\}$$

## Theorem (2012)

*There is a species tree $S$ for the gene tree $T \iff$ the triple set $\mathbb{S}$ is consistent*
*(can be tested efficiently).*

*A reconciliation map $\mu$ from $T$ to $S$ can be constructed in polynomial time.*

---

**From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

## Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent. The reconciliation map $\mu : T \rightarrow S$ is then "for free".

‡ **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

‡ **The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

‡ **From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

# Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid
   iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent.
   The reconciliation map $\mu : T \to S$ is then "for free".

Generalizations to non-disjoint non-symmetric relations have recently been published (characterization via uniformly non-prime 2-structures and di-cographs)

‡**Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

‡**The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016
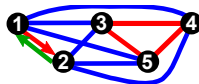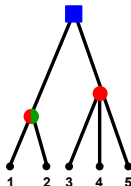
‡**From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

## Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid
   iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent.
   The reconciliation map $\mu : T \rightarrow S$ is then "for free".

Estimated relations usually don't have a tree-representation
(noise in the data, inference methods, ...)

$\longrightarrow$ Find "closest" valid event-relations (NP-hard).

‡ **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

‡ **The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

‡ **From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

# Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent. The reconciliation map $\mu : T \to S$ is then "for free".

Estimated relations usually don't have a tree-representation
(noise in the data, inference methods, ... )

$\longrightarrow$ Find "closest" valid event-relations (NP-hard).

Design of heuristics is work in progress.

‡**Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

‡**The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

‡**From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

## Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid
   iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent.
   The reconciliation map $\mu : T \to S$ is then "for free".

‡ **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

‡ **The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

‡ **From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

## Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid
   iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent.
   The reconciliation map $\mu : T \to S$ is then "for free".

The gene tree provides a lot of structural information of the species tree.

---

[‡] **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

[‡] **The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

[‡] **From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

# Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid
   iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent.
   The reconciliation map $\mu : T \to S$ is then "for free".

The gene tree provides a lot of structural information of the species tree.

The species triple set $\mathbb{S}$ is usually not consistent
(noise in the data, HGT, …)

$\longrightarrow$ Find max-consistent triple set of $\mathbb{S}$ (NP-hard).

Design of heuristics is work in progress.

---

‡ **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

‡ **The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

‡ **From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

# Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid
   iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent.
   The reconciliation map $\mu : T \rightarrow S$ is then "for free".

In the presence of HGT, we have characterized tree-representable event-relations,

---

[‡] **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013

[‡] **The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

[‡] **From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012

# Intermediate Summary and Open Problems

Characterization in the absence of HGT:

1. The two complementary estimated relations $\widehat{R}_\bullet$ and $\widehat{R}_\blacksquare$ are valid
   iff $G_{\widehat{R}_\bullet}$ is a cograph

2. There is a species tree $S$ for a gene tree $T$ iff the triple-set $\mathbb{S}$ is consistent.
   The reconciliation map $\mu : T \rightarrow S$ is then "for free".

In the presence of HGT, we have characterized tree-representable event-relations,
but an axiomatic framework for the reconciliation between gene trees (with HGT) and species
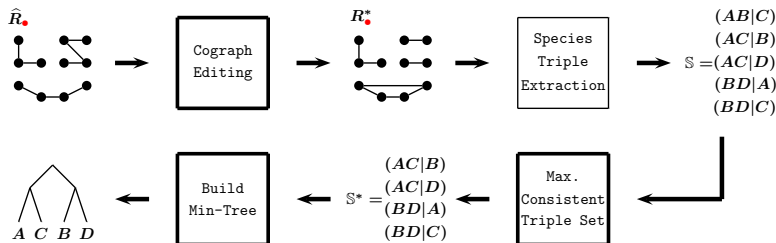tree/networks is missing

Work in progress.

‡ **Orthology Relations, Symbolic Ultrametrics, and Cographs**, *Hellmuth M*, H.-Rosales M, Huber K, Moulton V, Stadler PF, Wieseke N, *J. Math. Biol.*, 2013
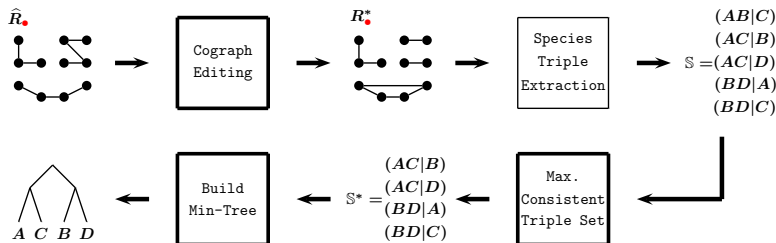
‡ **The Mathematics of Xenology: Di-cographs, Symbolic Ultrametrics, 2-structures and Tree-representable Systems of Binary Relations**, *Hellmuth M*, Stadler PF, Wieseke N, (accepted) *J. Math. Bio.*, 2016

‡ **From Event-Labeled Gene Trees to Species Trees.**, H.-Rosales M, *Hellmuth M*, Huber K, Moulton V, Wieseke N, Stadler PF, *BMC Bioinformatics*, 2012
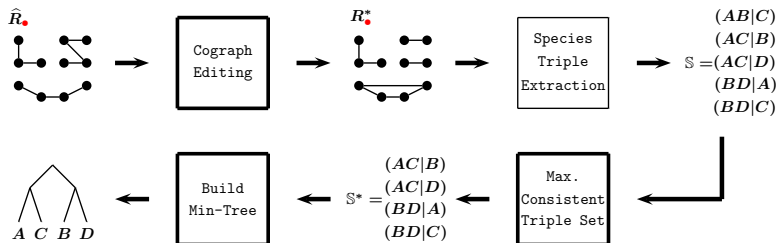
## Workflow **ParaPhylo**

# Workflow `ParaPhylo`



To demonstrate the potential of the approach without confounding it with computational approximations, we formulated all NP-hard problems (CE, MCT, LRT) as Integer Linear Program (ILP):

$$\min F(x) \text{ s.t. } Ax \leq b$$

---

[‡] **Phylogenomics with Paralogs**, *Hellmuth M*, Wieseke N, Lechner M, Lenhof HP, Middendorf M, Stadler PF, *PNAS*, 2015
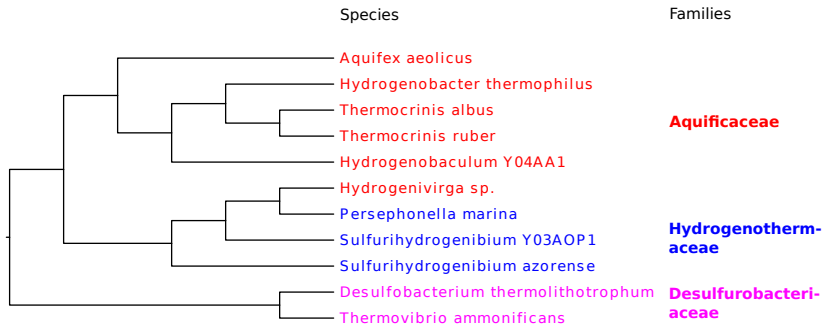
# Workflow `ParaPhylo`



The entire worflow as ILP is implemented in the Software `ParaPhylo`
using IBM ILOG CPLEX™ Optimizer 12.6.

It is freely available from
stubber.math-inf.uni-greifswald.de/~hellmuth/paraphylo

‡**Phylogenomics with Paralogs**, *Hellmuth M*, Wieseke N, Lechner M, Lenhof HP, Middendorf M, Stadler PF, *PNAS*, 2015

# Results - Real Life Data



- Class of bacteria that live in harsh environmental settings, e.g., hot springs, sulfur pools, . . .

- 11 Aquificales species with 2887 gene families
  (1372 - 3809 genes per species)

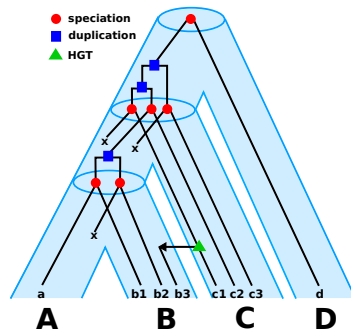- `ProteinOrtho` → **ParaPhylo** $\xrightarrow{34sec}$ Species Tree

‡ **ProteinOrtho: Detection of (Co)orthologs in large-scale analysis.**, Lechner M, Findeiß S, Steiner L, Marz M, Stadler PF, Prohaska SJ, *BMC Bioinformatics*, 2011

Intro
○○○○○○○

Distance Based
○○○○○

Consensus Methods
○○○○○○○○○○○○○

Phylo with Event Relations
○○○○○

Phylo with Event Relations II
○○○○○○○○

ParaPhylo
○○●○○○○○○○

# Results - Simulated Data

Artificial data generated with `ALF`:

Simulation of "true" evol. history

- generate binary species tree
- simulate dupl./loss/HGT history of
  gene sequences (within species tree)

**Output:** Species tree with embedded
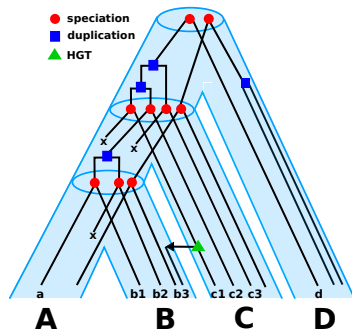gene trees and gene-sequences



‡ **ALF-a simulation framework for genome evolution.**, Dalquen et al., *Mol. Biol. Evol.*, 2012

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000●00000

# Results - Simulated Data

### Artificial data generated with `ALF`:

Simulation of "true" evol. history

- generate binary species tree
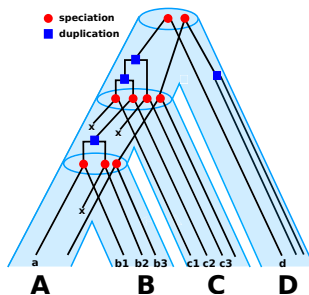- simulate dupl./loss/HGT history of
  gene sequences (within species tree)

**Output:** Species tree with embedded
gene trees and gene-sequences



‡**ALF-a simulation framework for genome evolution.**, Dalquen et al., *Mol. Biol. Evol.*, 2012
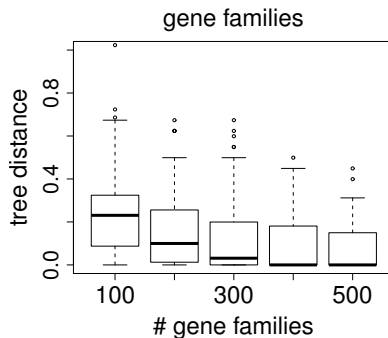
# Results - Simulation without HGT



ALF (no HGT)

$\longrightarrow$ The cograph $G_{R_\bullet}$ is directly accessible

$\longrightarrow$ Compute cotree of $G_{R_\bullet}$

$\longrightarrow$ Extract the species triples set $\mathbb{S}$ (consistent)

$\longrightarrow$ Compute least resolved species tree and compare it with initial species tree

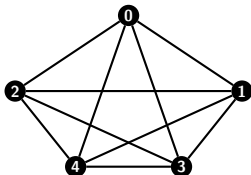# Results - Simulation without HGT

Accuracy of reconstructed species trees (20 species)
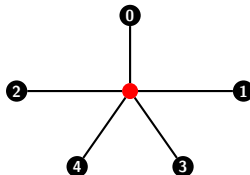as function of number of independent gene families:



Simulation with `ALF` with duplication/loss rate 0.005
($\sim 8\%$ duplications) and no HGT.

# Results - Simulation without HGT

Since no HGT, we have $(x, y) \in R_\bullet$ iff $(x, y) \notin R_\blacksquare$
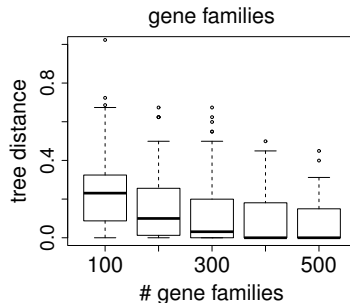


$$G_{R_\bullet}$$

$$T$$

If $\nexists$ paralogs $\rightarrow G_{R_\bullet}$ is a clique $\rightarrow$ gene tree is a star $\rightarrow$ no species triples can be inferred.

To obtain fully resolved species trees, a sufficient number of gene duplications must have occurred, since the phylogenetic information utilized by our approach is entirely contained in the duplication events.

# Results - Simulation without HGT

Accuracy of reconstructed species trees (20 species)
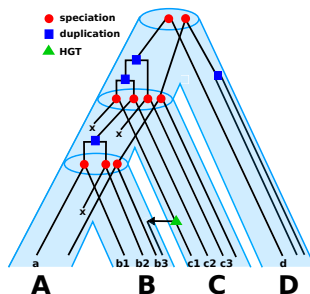as function of number of independent gene families:



More genefamilies (*incl. paralogs*) → more accurate species trees.

Fewer gene families → less duplicated genes → species trees less resolved.

Deviations from perfect reconstructions are exclusively explained by a lack of perfect
resolution.

# Results - Simulation with HGT



ALF with HGT (10 Species, 1000 Gene Families):

(1) we get simulated sequences:
   ProteinOrtho → **ParaPhylo** → Species Tree

(2) we get $R_\bullet, R_\blacktriangle, R_\blacksquare$ from the gene tree
   But **ParaPhylo** can only deal with $R_\bullet$ and $\overline{R_\bullet}$, so-far
   Thus, we use $\widehat{R}_\bullet = R_\bullet \cup \mathscr{R}$, where $\mathscr{R} \subseteq R_\blacksquare \cup R_\blacktriangle$.
   Graph $G_{\widehat{R}_\bullet} \to$ **ParaPhylo** → Species Tree

# Results - Simulation with HGT



ALF with HGT (10 Species, 1000 Gene Families):

  (1)  we get simulated sequences:
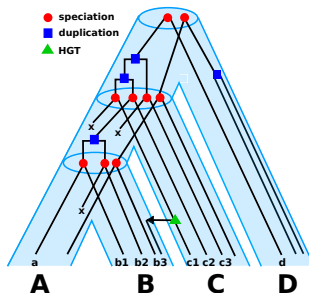      ProteinOrtho → **ParaPhylo** → Species Tree

  (2)  we get $R_\bullet, R_\blacktriangle, R_\blacksquare$ from the gene tree
      But **ParaPhylo** can only deal with $R_\bullet$ and $\overline{R_\bullet}$, so-far
      Thus, we use $\widehat{R}_\bullet = R_\bullet \cup \mathscr{R}$, where $\mathscr{R} \subseteq R_\blacksquare \cup R_\blacktriangle$.
      Graph $G_{\widehat{R}_\bullet} →$ **ParaPhylo** → Species Tree

# Results - Simulation with HGT



ALF with HGT (10 Species, 1000 Gene Families):

(1) we get simulated sequences:
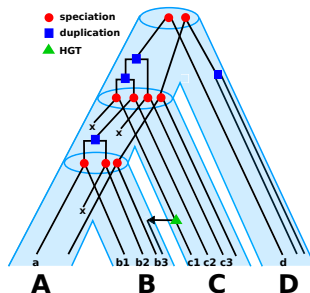ProteinOrtho → **ParaPhylo** → Species Tree

(2) we get $R_\bullet$, $R_\blacktriangle$, $R_\blacksquare$ from the gene tree
But **ParaPhylo** can only deal with $R_\bullet$ and $\overline{R_\bullet}$, so-far
Thus, we use $\widehat{R}_\bullet = R_\bullet \cup \mathscr{R}$, where $\mathscr{R} \subseteq R_\blacksquare \cup R_\blacktriangle$.
Graph $G_{\widehat{R}_\bullet} \to$ **ParaPhylo** → Species Tree

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# Results - Simulation with HGT



ALF with HGT (10 Species, 1000 Gene Families):

(1) we get simulated sequences:
    ProteinOrtho → **ParaPhylo** → Species Tree
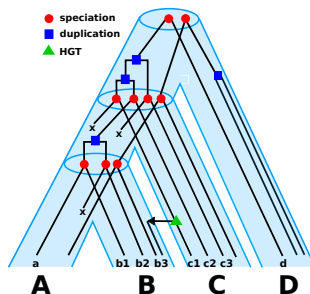
(2) we get $R_\bullet, R_\blacktriangle, R_\blacksquare$ from the gene tree
    But **ParaPhylo** can only deal with $R_\bullet$ and $\overline{R_\bullet}$, so-far
    Thus, we use $\widehat{R}_\bullet = R_\bullet \cup \mathscr{R}$, where $\mathscr{R} \subseteq R_\blacksquare \cup R_\blacktriangle$.
    Graph $G_{\widehat{R}_\bullet} \to$ **ParaPhylo** → Species Tree

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000●0

# Results - Simulation with HGT



ALF with HGT (10 Species, 1000 Gene Families):

(1) we get simulated sequences:
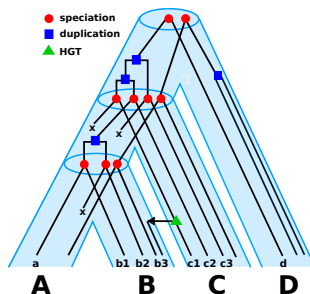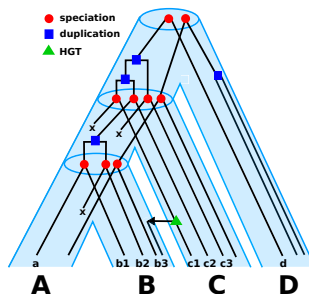ProteinOrtho → **ParaPhylo** → Species Tree

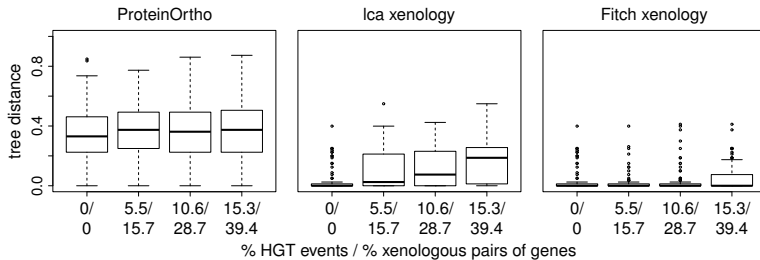(2) we get $R_\bullet, R_\blacktriangle, R_\blacksquare$ from the gene tree
But **ParaPhylo** can only deal with $R_\bullet$ and $\overline{R_\bullet}$, so-far
Thus, we use $\widehat{R}_\bullet = R_\bullet \cup \mathscr{R}$, where $\mathscr{R} \subseteq R_\blacksquare \cup R_\blacktriangle$.
Graph $G_{\widehat{R}_\bullet} \to$ **ParaPhylo** → Species Tree

Intro
0000000

Distance Based
00000

Consensus Methods
000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000

# Results - Simulation with HGT



$\mathtt{ALF}$ with HGT (10 Species, 1000 Gene Families):

(1) we get simulated sequences:
$\mathtt{ProteinOrtho} \rightarrow \mathbf{ParaPhylo} \rightarrow$ Species Tree

(2) we get $R_\bullet, R_\blacktriangle, R_\blacksquare$ from the gene tree
But $\mathbf{ParaPhylo}$ can only deal with $R_\bullet$ and $\overline{R_\bullet}$, so-far
Thus, we use $\widehat{R}_\bullet = R_\bullet \cup \mathscr{R}$, where $\mathscr{R} \subseteq R_\blacksquare \cup R_\blacktriangle$.
Graph $G_{\widehat{R}_\bullet} \rightarrow \mathbf{ParaPhylo} \rightarrow$ Species Tree

Intro
0000000
Distance Based
00000
Consensus Methods
000000000000
Phylo with Event Relations
00000
Phylo with Event Relations II
00000000
ParaPhylo
000000000●

# Results - Simulation with HGT

Accuracy of reconstructed species trees vs. intensity of HGT



left $\widehat{R}_\bullet$ = "estim." orthologs via `ProteinOrtho`

middle $\widehat{R}_\bullet$ = orthologs $R_\bullet$ + lca-xenologs $R_\blacktriangle$

*(orthology-overprediction / all paralogs are correctly identified)*

right $\widehat{R}_\bullet$ = orthologs $R_\bullet$ + all pairs of genes having at least one
HGT event on their path

*(orthology-overprediction / all paralogs that are not disturbed by HGT on their paths are correctly identified)*

# Results - Simulation with HGT

Accuracy of reconstructed species trees vs. intensity of HGT



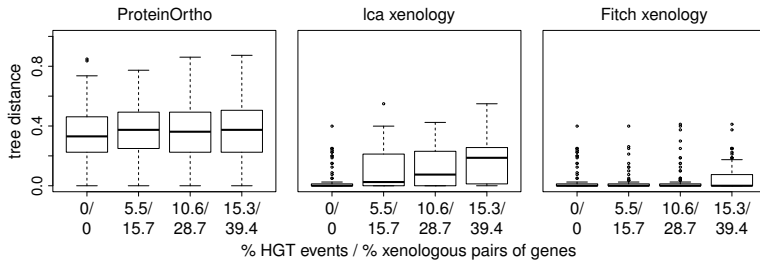left $\widehat{R}_\bullet$ = "estim." orthologs via `ProteinOrtho`

middle $\widehat{R}_\bullet$ = orthologs $R_\bullet$ + lca-xenologs $R_\blacktriangle$

*(orthology-overprediction / all paralogs are correctly identified)*

right $\widehat{R}_\bullet$ = orthologs $R_\bullet$ + all pairs of genes having at least one
HGT event on their path

*(orthology-overprediction / all paralogs that are not disturbed by HGT on their paths are correctly identified)*

Intro
0000000

Distance Based
00000

Consensus Methods
0000000000000

Phylo with Event Relations
00000

Phylo with Event Relations II
00000000

ParaPhylo
000000000●

# Results - Simulation with HGT

Accuracy of reconstructed species trees vs. intensity of HGT



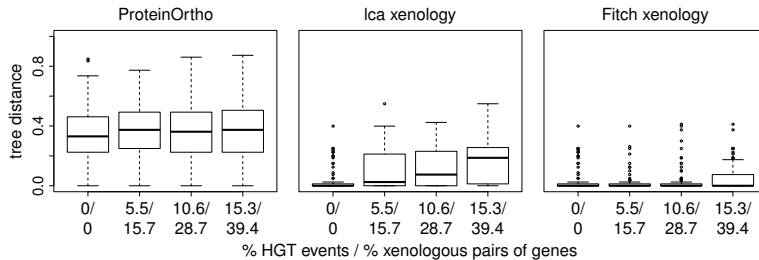left $\widehat{R}_\bullet$ = "estim." orthologs via `ProteinOrtho`

middle $\widehat{R}_\bullet$ = orthologs $R_\bullet$ + lca-xenologs $R_\blacktriangle$

*(orthology-overprediction / all paralogs are correctly identified)*

right $\widehat{R}_\bullet$ = orthologs $R_\bullet$ + all pairs of genes having at least one
HGT event on their path

*(orthology-overprediction / all paralogs that are not disturbed by HGT on their paths are correctly identified)*