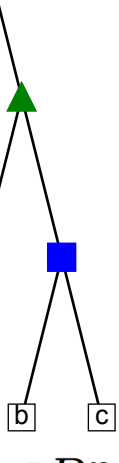
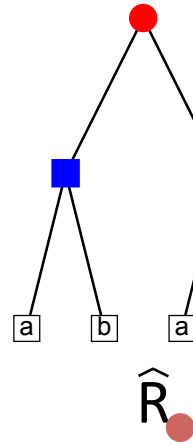


Graph-based Adjustment of Orthology-Relations

Paul Klemm

University of Greifswald, Germany

32th TBI Winterseminar, Bled 2017



$$\prod_{i=1}^{n-1} \binom{numST_i + (a_i - 1)}{a_i}$$

Introduction

Phylogenetic, Orthology relation

Cographs

Graph-based orthology inference

Artificial Data Analysis

Enumeration of all gene trees

Results

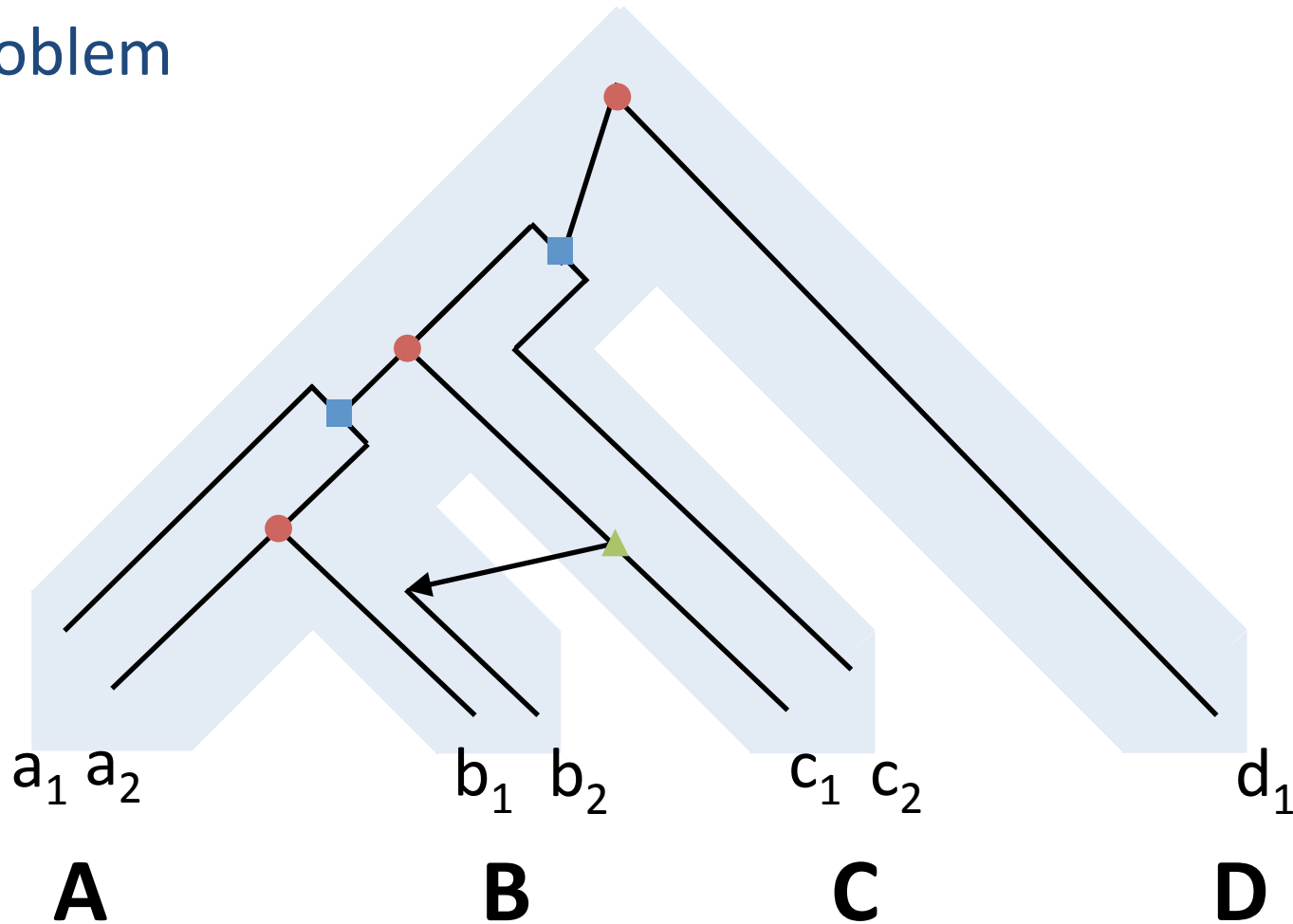
Summary and Outlook

Introduction

Artificial
Data
Analysis

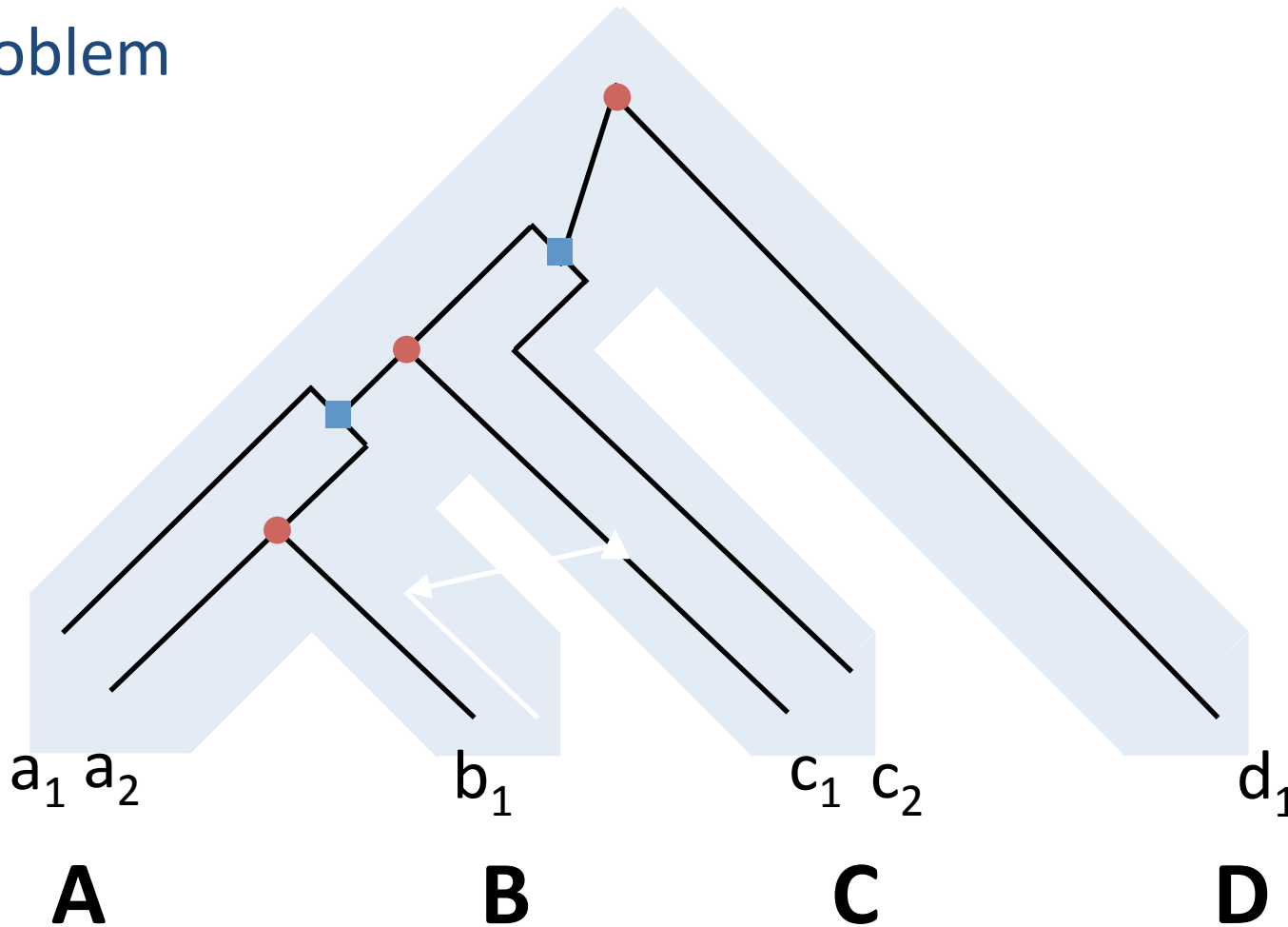
Summary
and
Outlook

The problem



Introduction

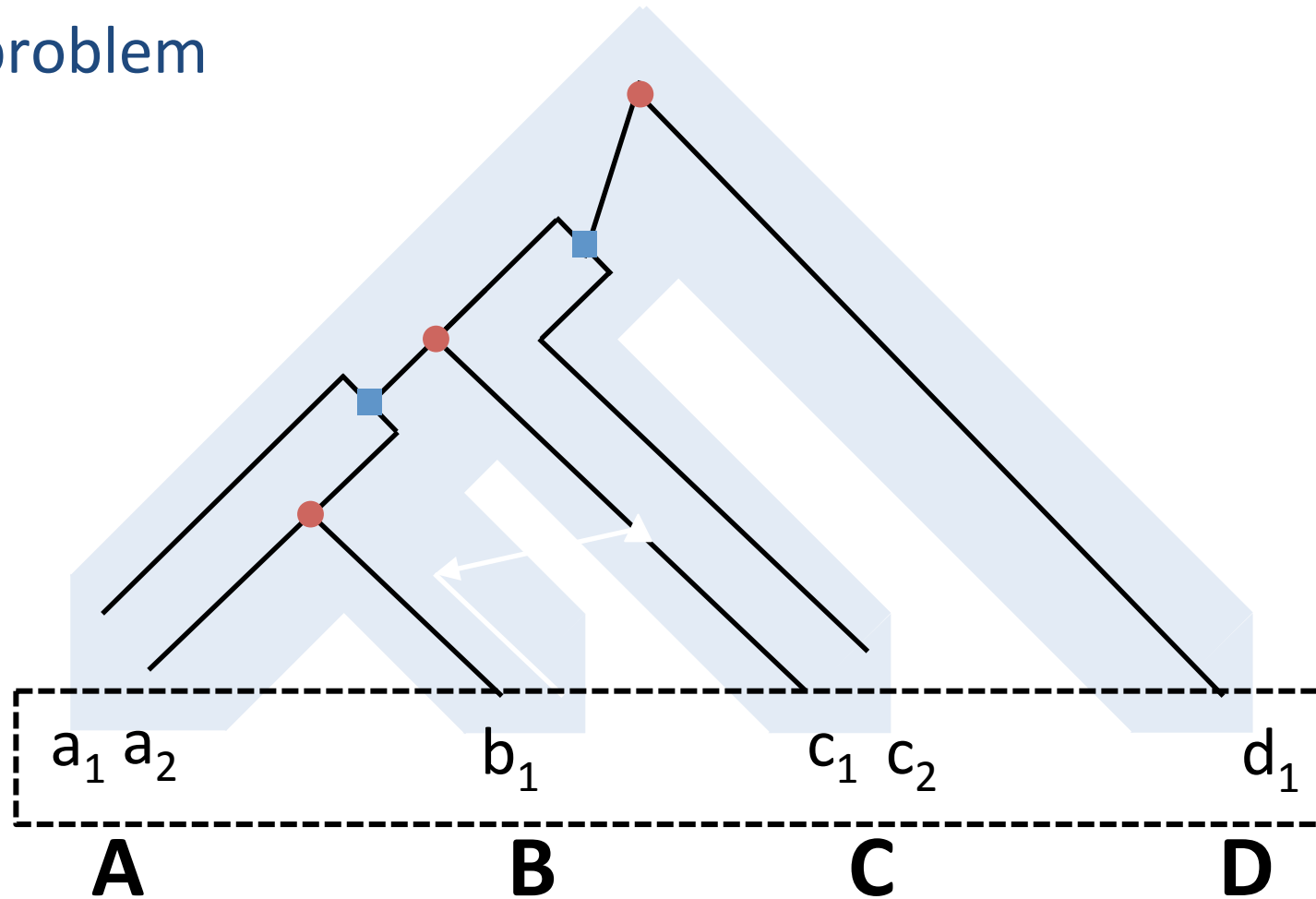
The problem



→ Simplification: no HGT-event

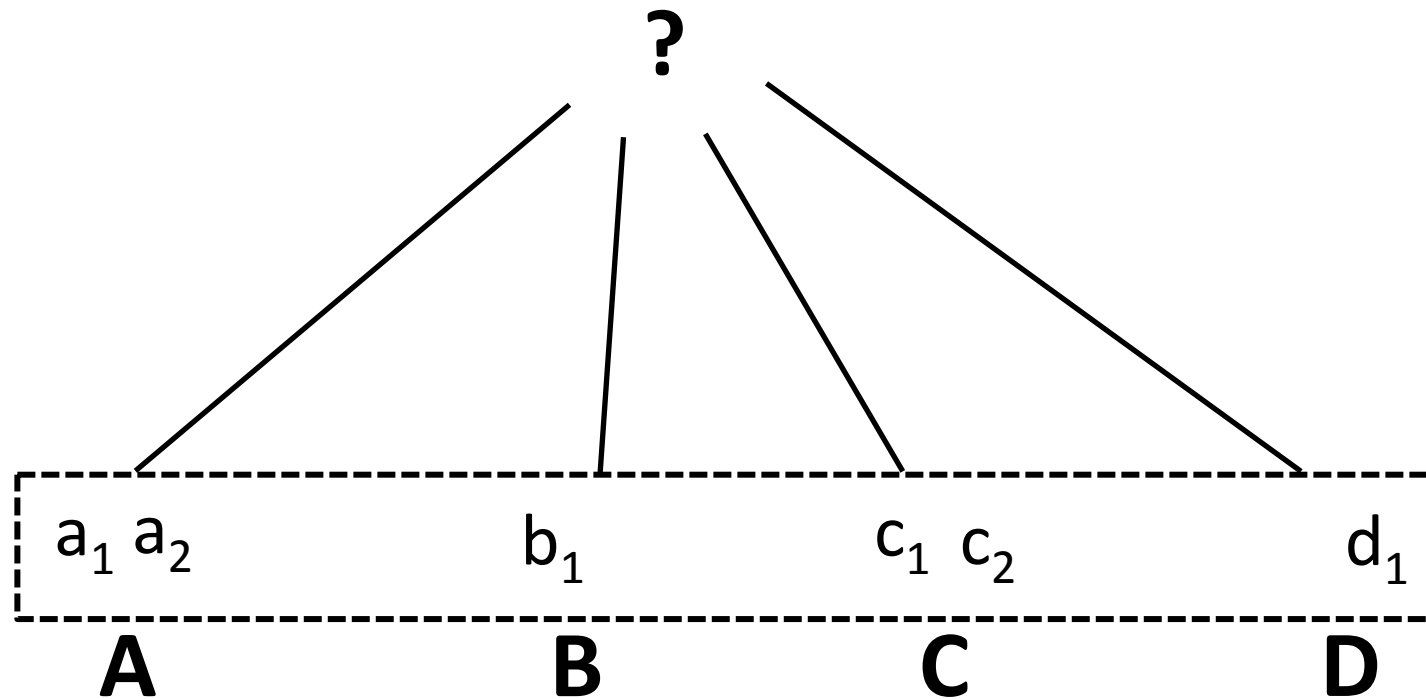
Introduction

The problem

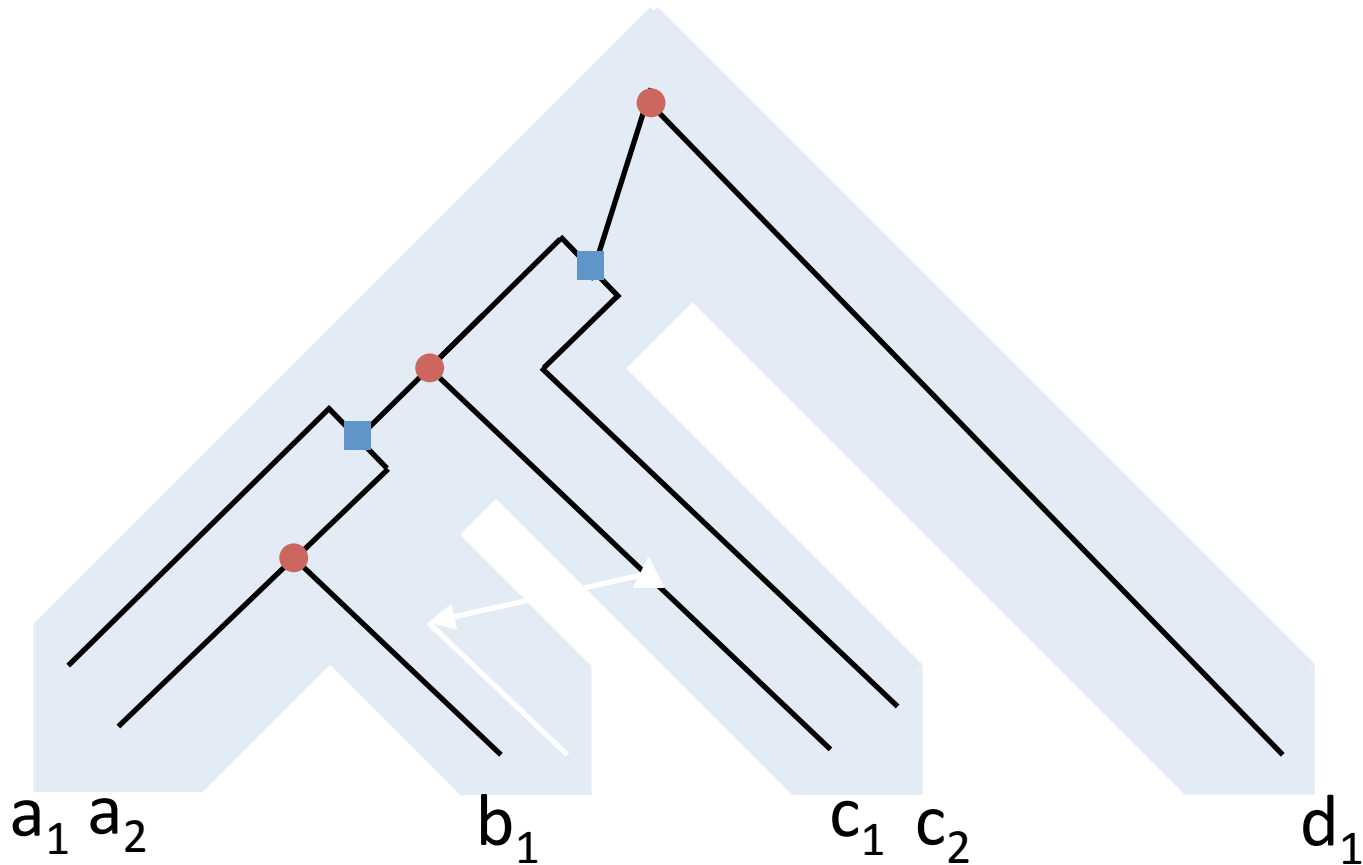


→ Simplification: no HGT-event

The problem



→ Simplification: no HGT-event



We call two genes x and y

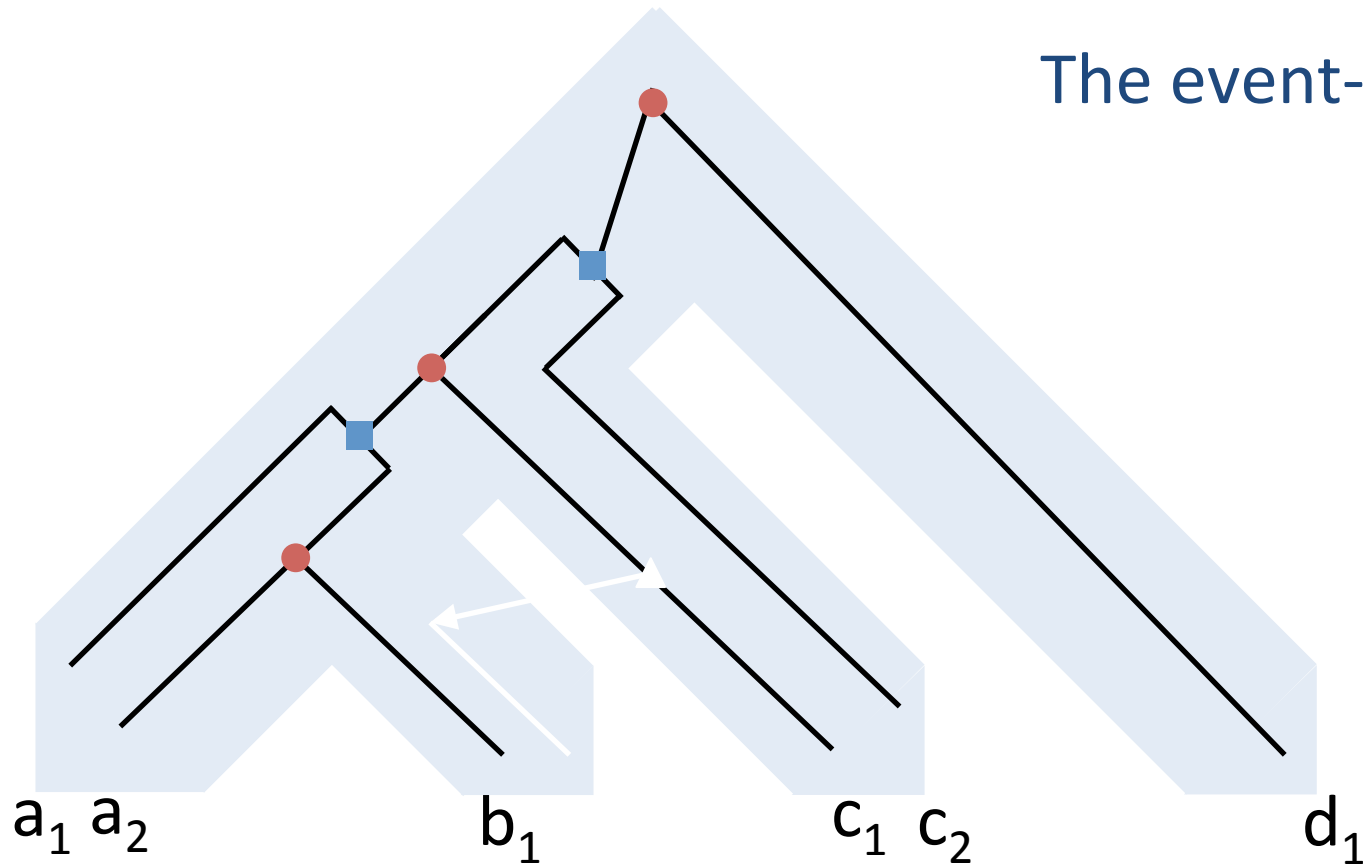
- **ortholog** if $\text{lca}(x,y) = \bullet$ (speciation)
- **paralog** if $\text{lca}(x,y) = \blacksquare$ (duplication)
- **xenolog** if $\text{lca}(x,y) = \blacktriangle$ (HGT)

e.g. a_2 and b_1

e.g. c_1 and c_2

e.g. b_2 and c_1

The event-Relations



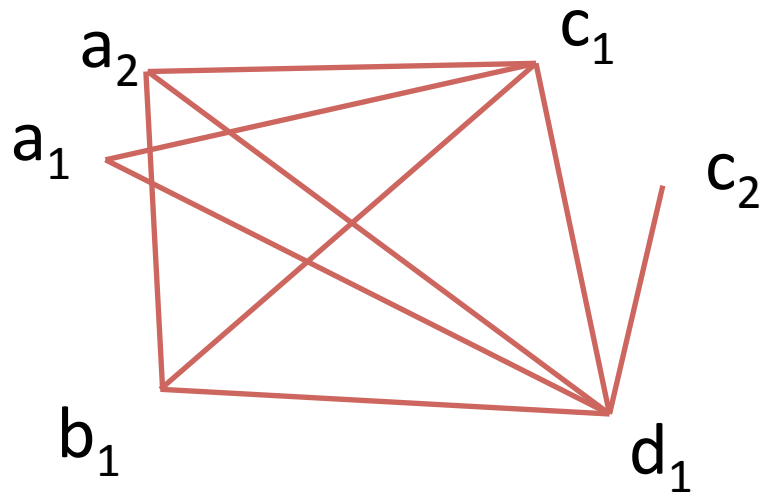
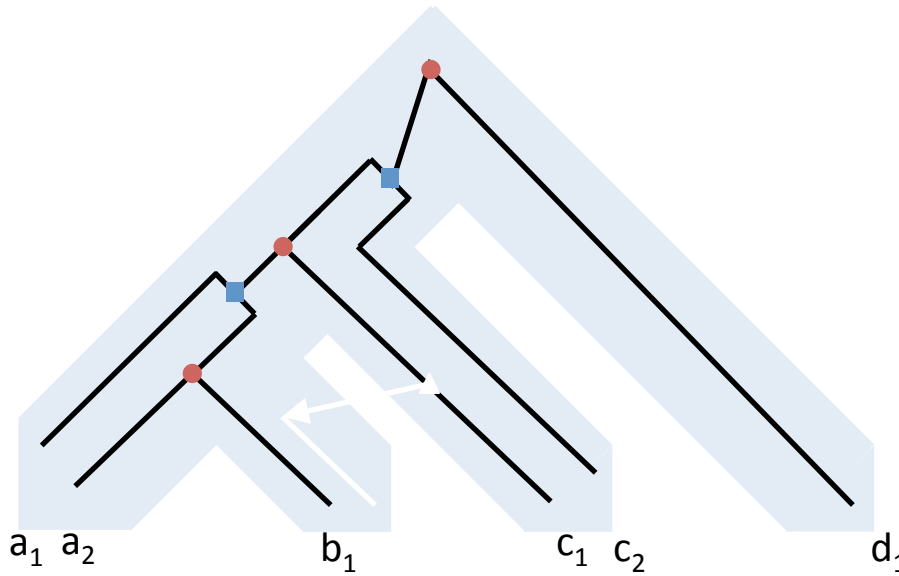
The (distinct, binary, symmetric) **event-Relations**:

- R_{\bullet} = the set of all (x,y) with $\text{lca}(x,y) = \bullet$ (speciation)
- R_{\blacksquare} = ... with $\text{lca}(x,y) = \blacksquare$ (duplication)
- R_{\blacktriangle} = ... with $\text{lca}(x,y) = \blacktriangle$ (HGT)

The event-Relations

Simplification: no HGT-events

$$\rightarrow R_{\bullet} = \overline{R_{\blacksquare}}$$

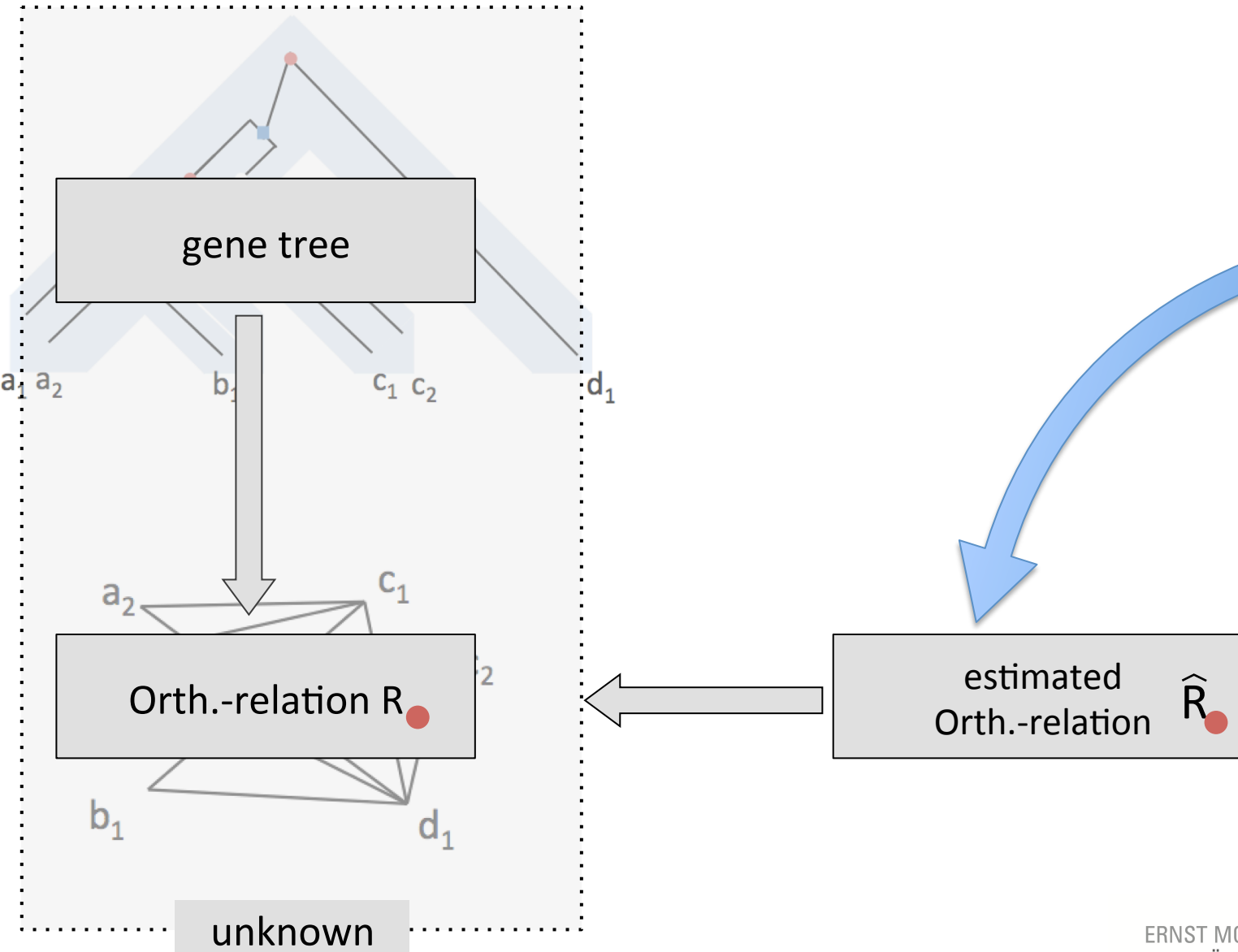


\rightarrow graph representation of

R_{\bullet} and R_{\blacksquare}

Introduction

Estimation

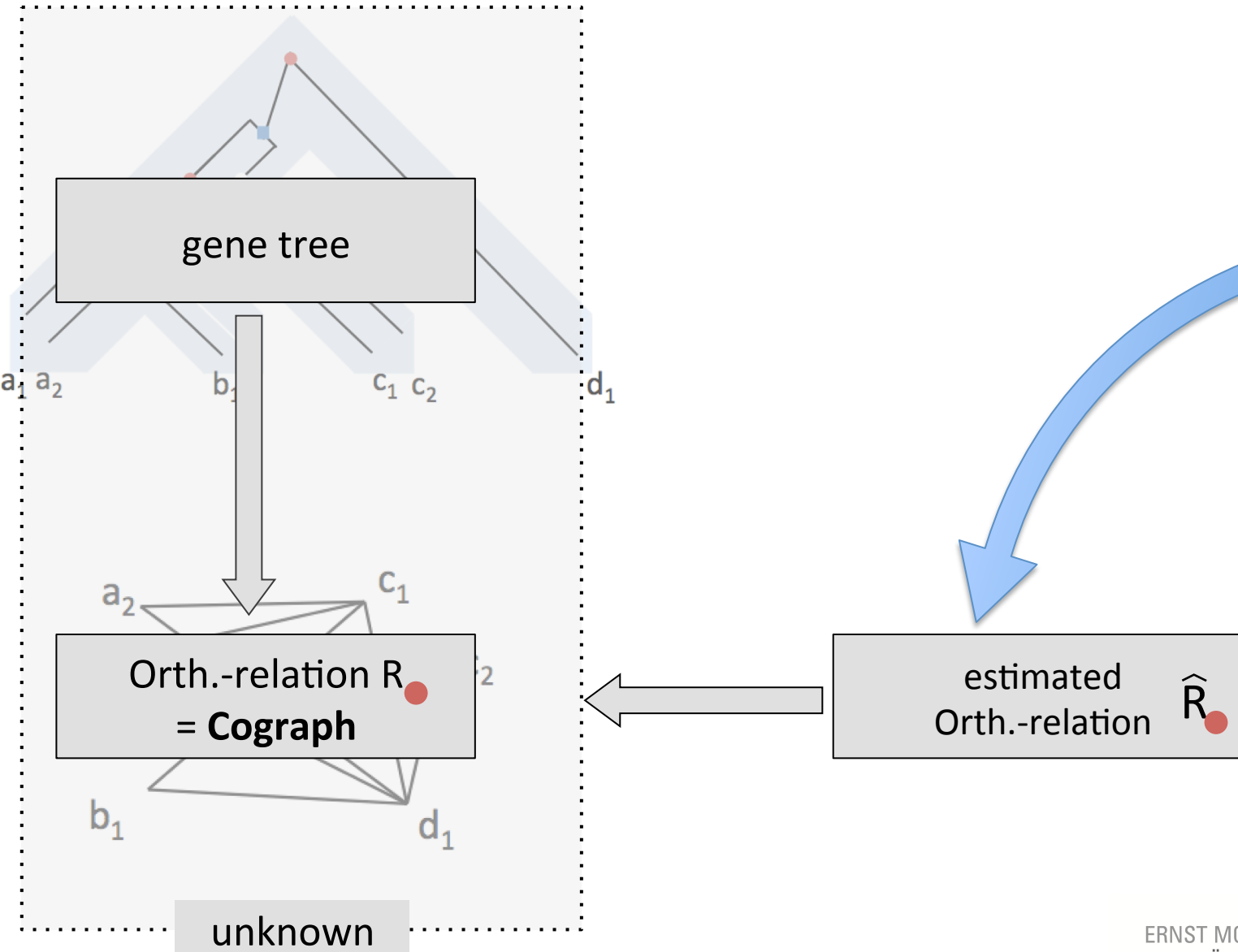


Introduction


Artificial
Data
Analysis

Summary
and
Outlook


Estimation

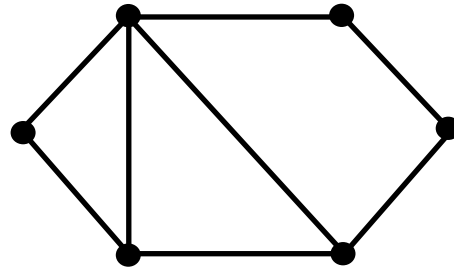


Cographs


A graph is called **cograph** if and only if there exists no induced subgraph on 4 nodes that is a P_4 ()

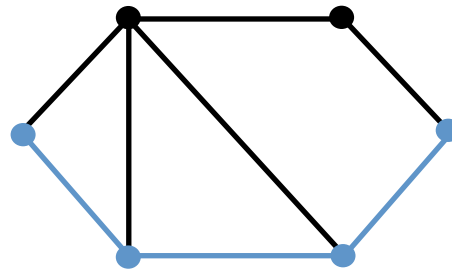
Cographs

A graph is called **cograph** if and only if there exists no induced subgraph on 4 nodes that is a P_4 ()




Cographs

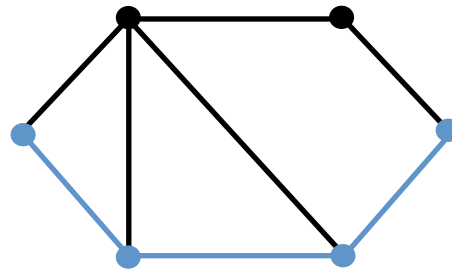
A graph is called **cograph** if and only if there exists no induced subgraph on 4 nodes that is a P_4 ()



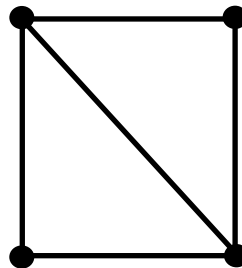
not a Cograph

Cographs

A graph is called **cograph** if and only if there exists no induced subgraph on 4 nodes that is a P_4 ()



not a Cograph



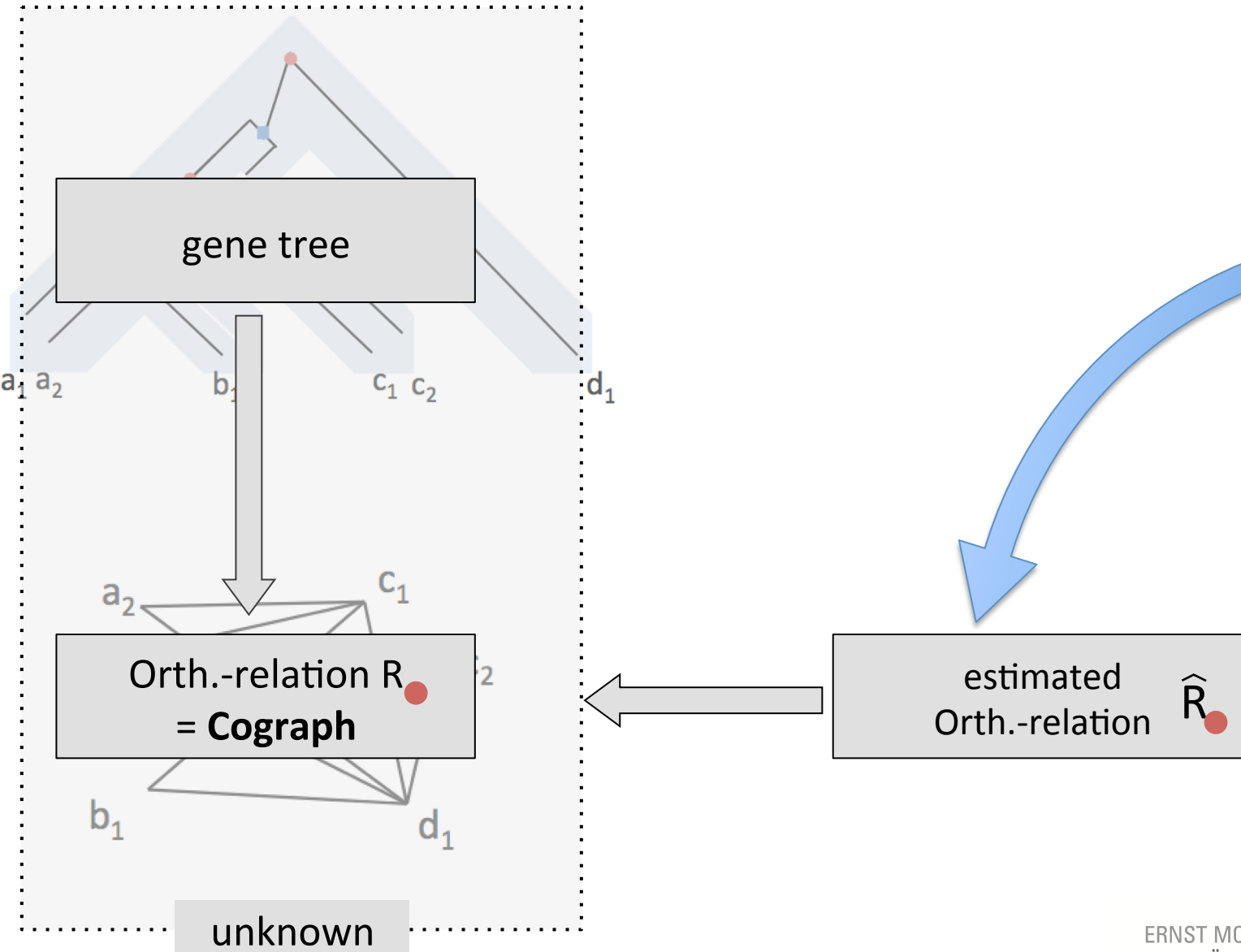
Cograph

Introduction

Artificial
Data
Analysis

Summary
and
Outlook

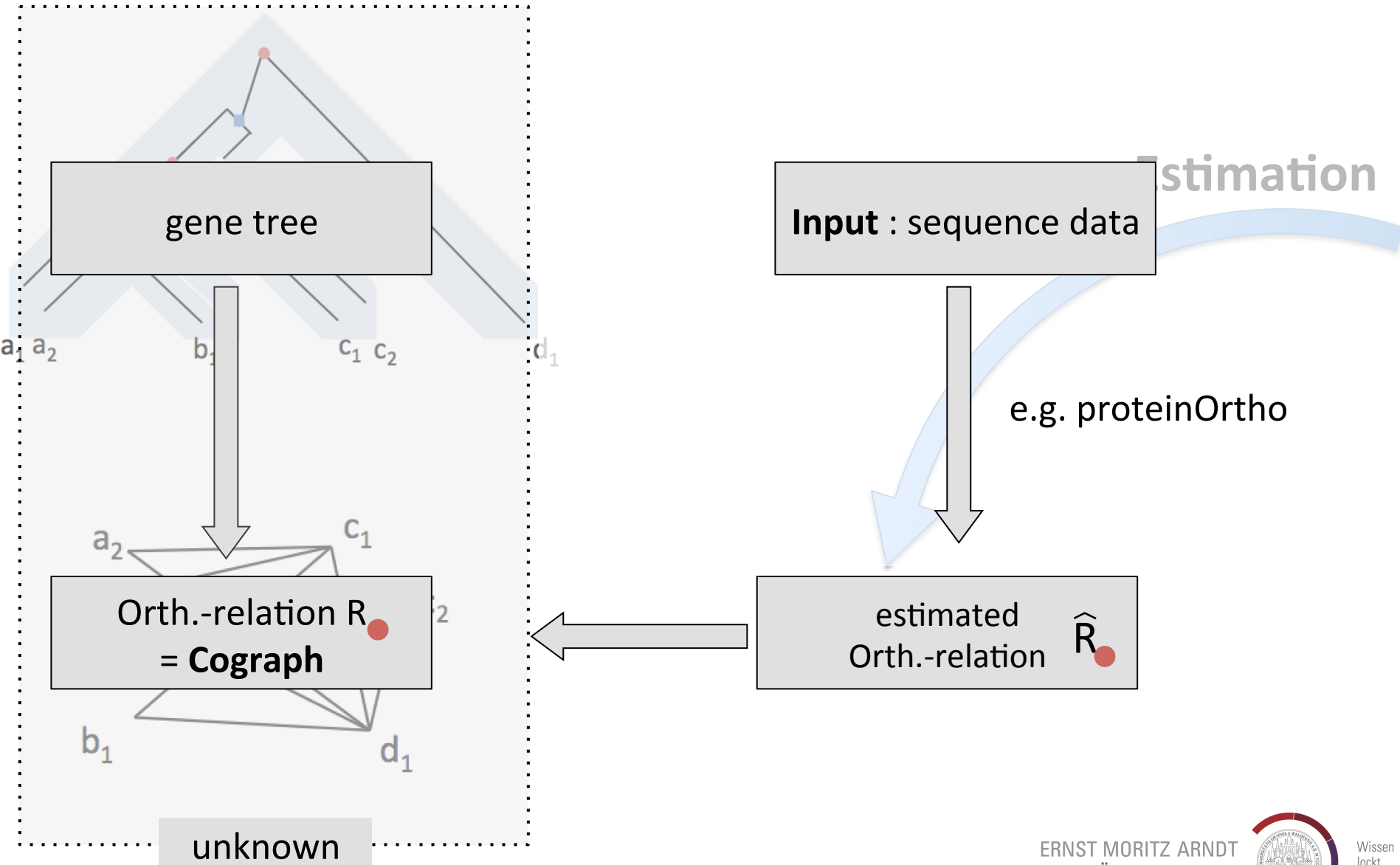
Estimation



Introduction

Artificial
Data
Analysis

Summary
and
Outlook



Graph-based orthology inference (e.g. proteinOrtho)

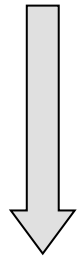
Input : sequence data

Local alignment search
(e.g. BLAST)

Sequence similarity $s(x,y)$

Graph-based orthology inference (e.g. proteinOrtho)

Input : sequence data



Local alignment search
(e.g. BLAST)

Sequence similarity $s(x,y)$

Then two genes x in X and y in Y are **estimated orthologs** if:

- i. they are from different species and
- ii. $s(x,y)$ is the ~best score compared to $s(x',y)$ and $s(x,y')$ where x' in X , y' in Y

Graph-based orthology inference (e.g. proteinOrtho)

Then two genes x and y are **estimated orthologs** if:

- i. they are from different species and
- ii. $s(x,y)$ is the \sim best score compared to $s(x',y)$ and $s(x,y')$ where x' in X , y' in Y

Simplification: leaf distance $d(x,y)$ of the „true“ gene tree:

$d(x,y)$ = length of the shortest path between x and y

Graph-based orthology inference (e.g. proteinOrtho)

Then two genes x and y are **estimated orthologs** if:

- i. they are from different species and
- ii. $s(x,y)$ is the \sim best score compared to $s(x',y)$ and $s(x,y')$ where x' in X , y' in Y

Simplification: leaf distance $d(x,y)$ of the „true“ gene tree:

$d(x,y)$ = length of the shortest path between x and y

$$d \sim 1/s$$

Graph-based orthology inference (e.g. proteinOrtho)

Then two genes x and y are **estimated orthologs** if:

- i. they are from different species and
- ii. $s(x,y)$ is the \sim best score compared to $s(x',y)$ and $s(x,y')$ where x' in X , y' in Y

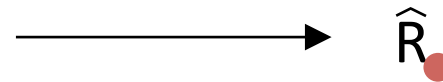
Simplification: leaf distance $d(x,y)$ of the „true“ gene tree:

$d(x,y)$ = length of the shortest path between x and y

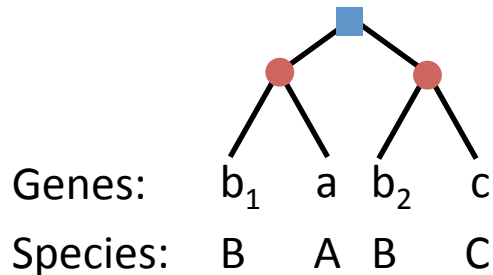
$$d \sim 1/s$$

→ **estimated orthologs** if:

- i. they are from different species and
- ii. the distance $d(x,y)$ is „small“



Graph-based orthology inference



R :

\hat{R} :

Simplification: leaf distance $d(x,y)$ of the „true“ gene tree:

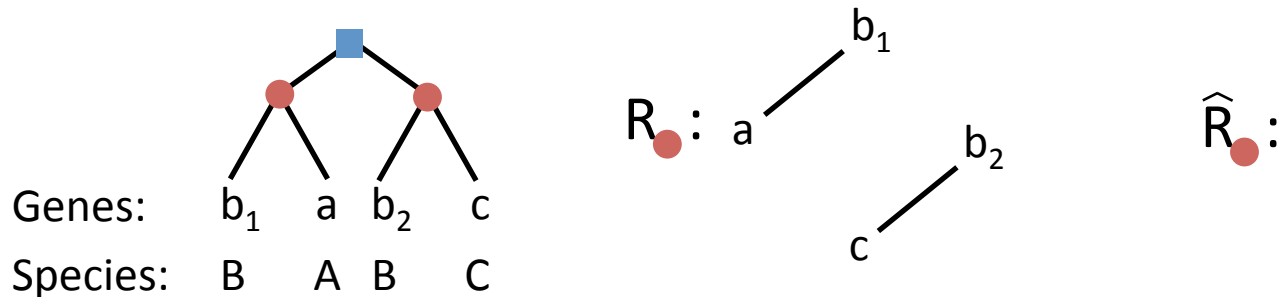
$d(x,y)$ = length of the shortest path between x and y

→ **estimated orthologs** if:

- i. they are from different species and
- ii. the distance $d(x,y)$ is „small“

→ \hat{R}

Graph-based orthology inference

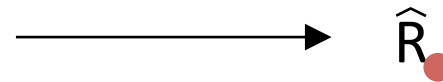


Simplification: leaf distance $d(x,y)$ of the „true“ gene tree:

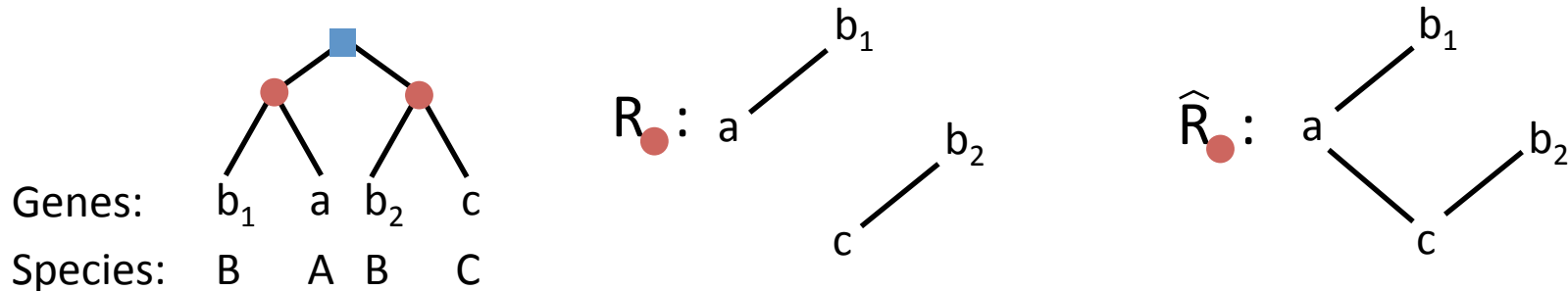
$d(x,y)$ = length of the shortest path between x and y

→ **estimated orthologs** if:

- i. they are from different species and
- ii. the distance $d(x,y)$ is „small“



Graph-based orthology inference

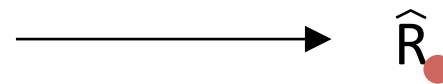


Simplification: leaf distance $d(x,y)$ of the „true“ gene tree:

$d(x,y)$ = length of the shortest path between x and y

→ **estimated orthologs** if:

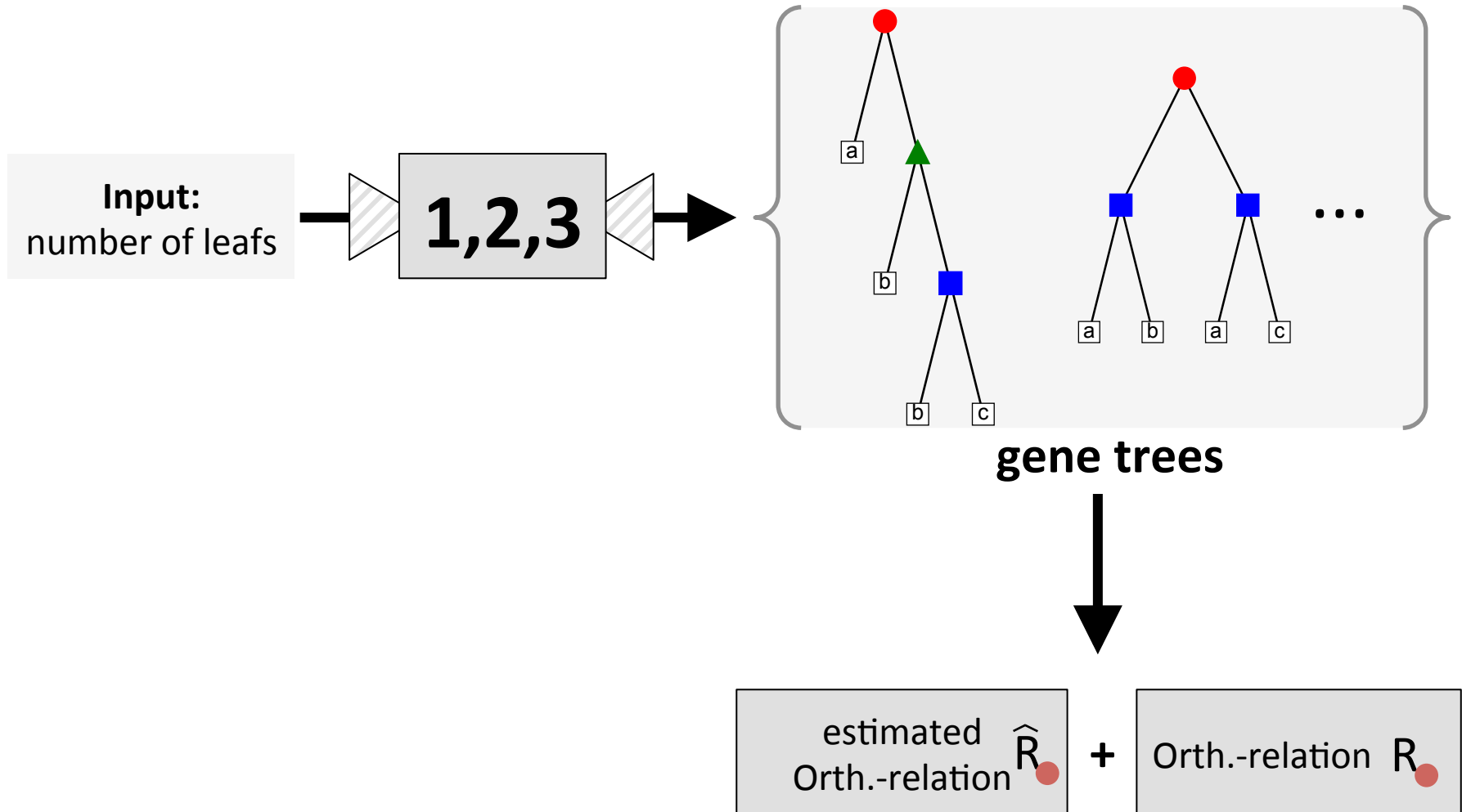
- i. they are from different species and
- ii. the distance $d(x,y)$ is „small“



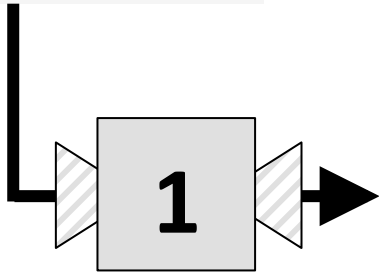
Question:

How much information can we infer from non-cograph relations ?

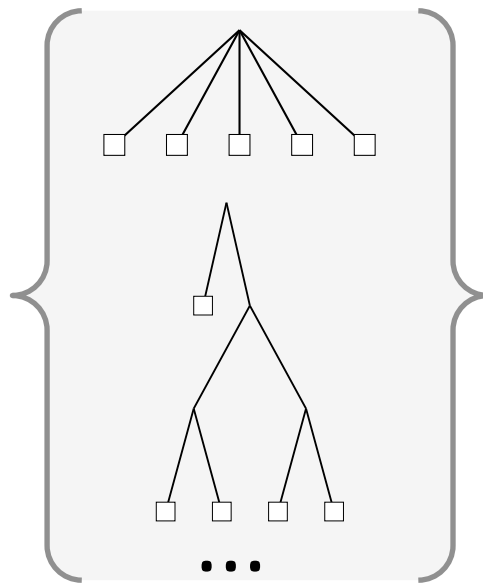
How much information can we infer from non-cograph relations ?



Input:
number of leafs



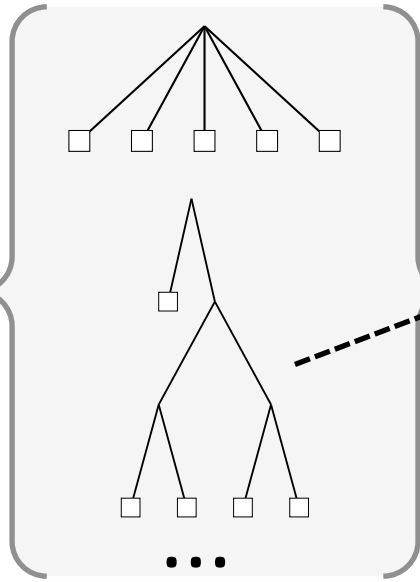
skeleton tree



Input:
number of leaves

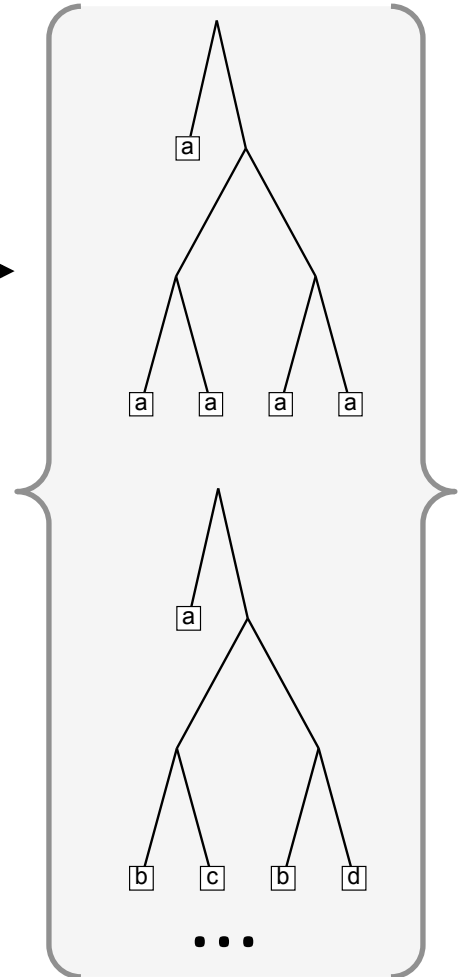
1

skeleton tree



2

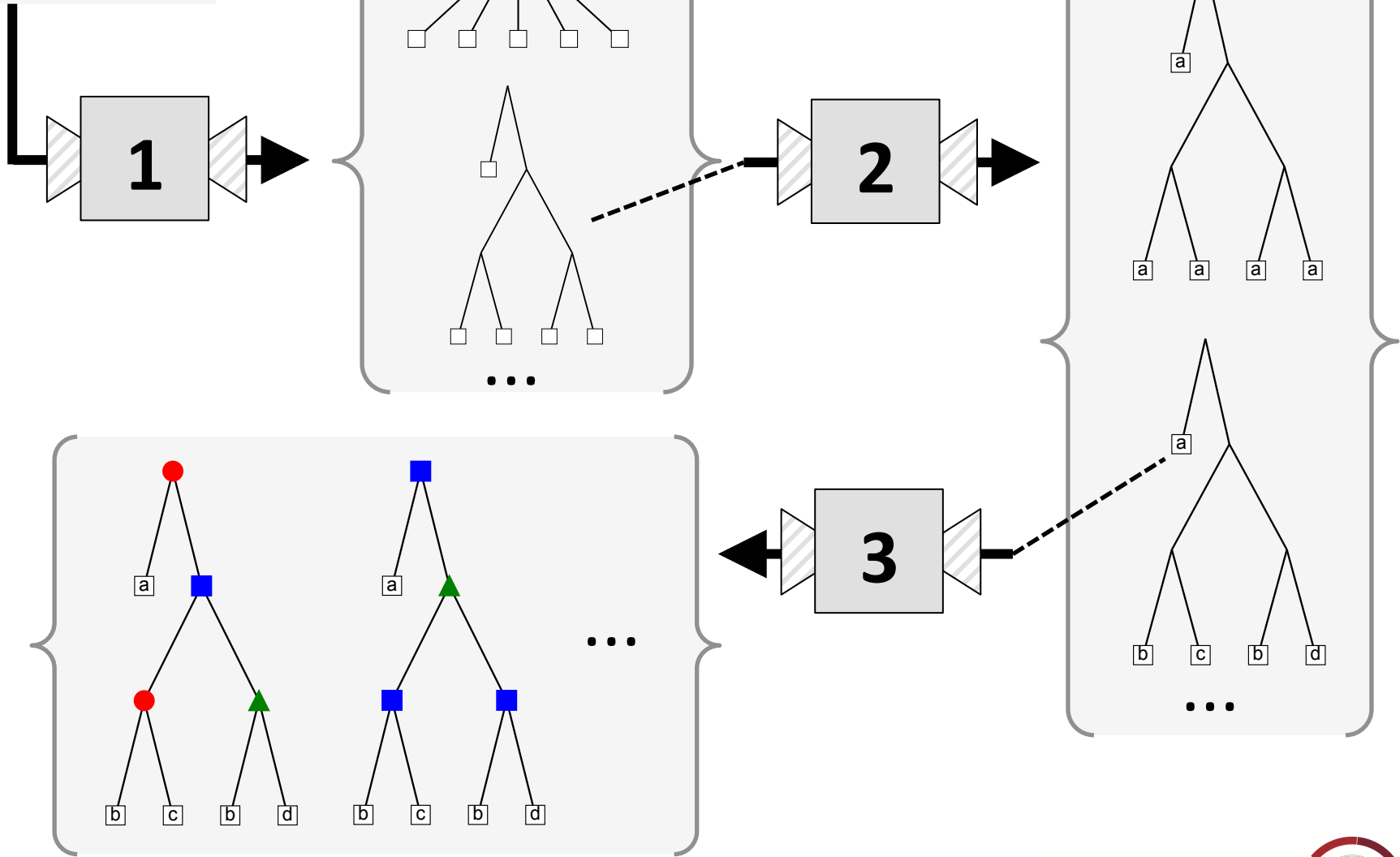
labeled
skeleton tree



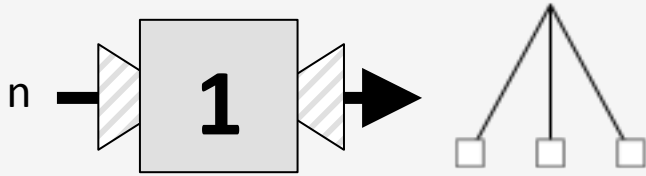
Input:
number of leafs

skeleton tree

labeled
skeleton tree



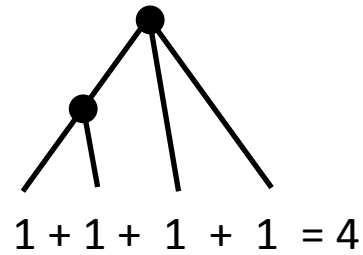
gene trees

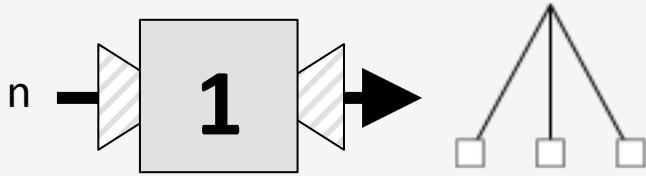


1. Generate skeleton trees

We want:

1. Rooted trees with n leaves

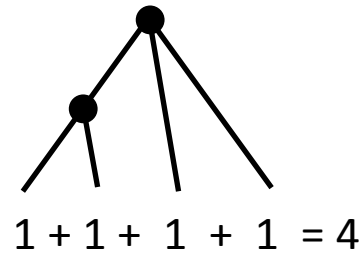




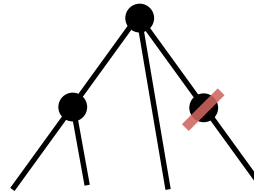
1. Generate skeleton trees

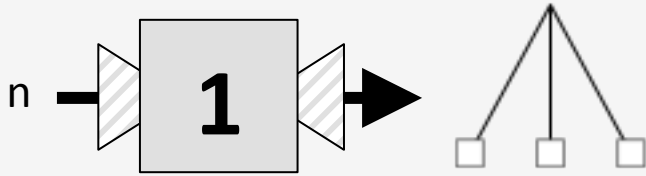
We want:

1. Rooted trees with n leafs



2. All inner nodes should have at least 2 childs

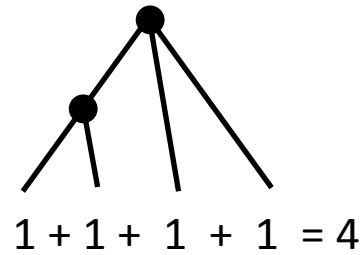




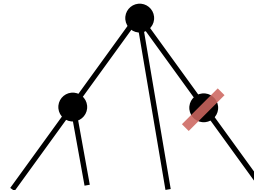
1. Generate skeleton trees

We want:

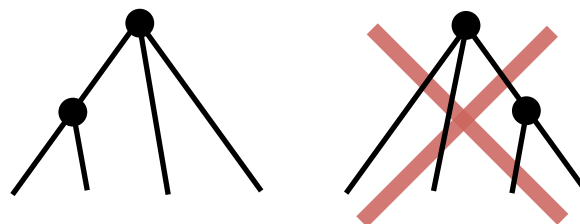
1. Rooted trees with n leaves

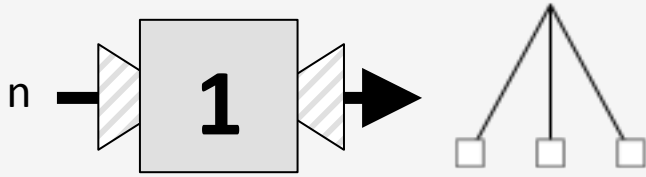


2. All inner nodes should have at least 2 childs

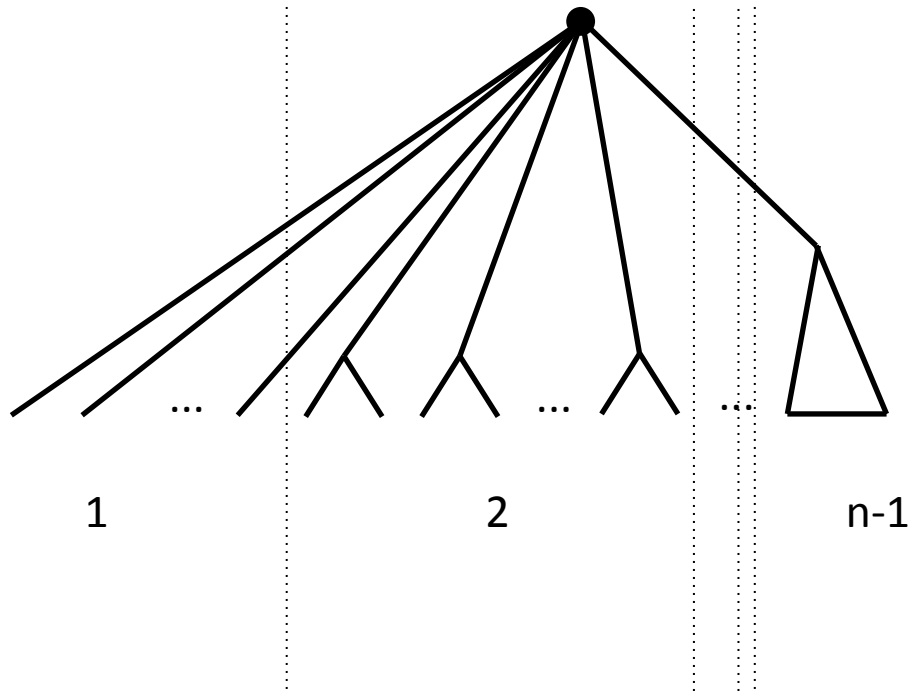


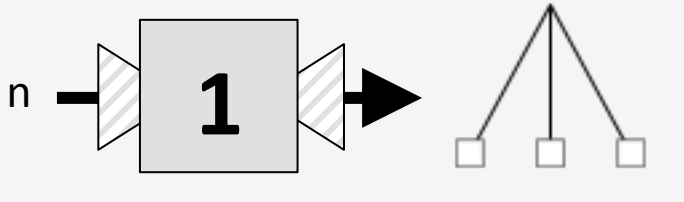
3. Only one representative of the isomorphism classes



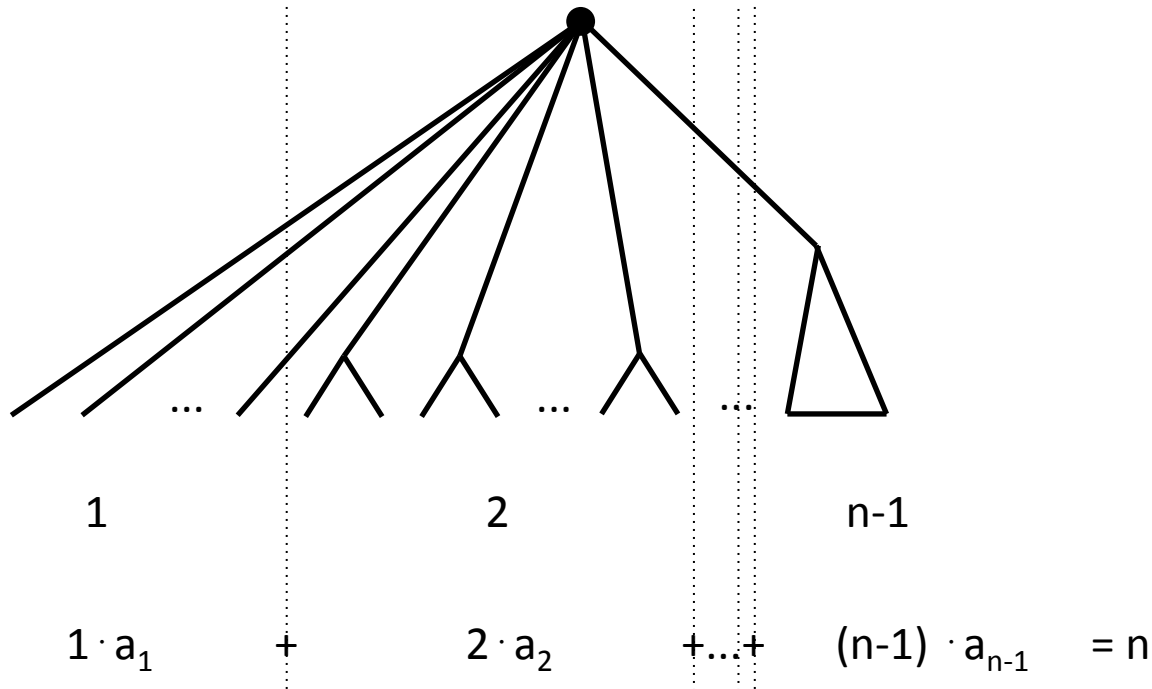


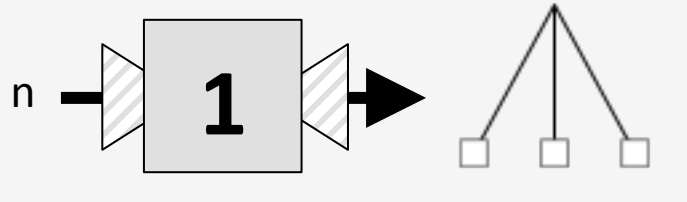
1. Generate skeleton trees



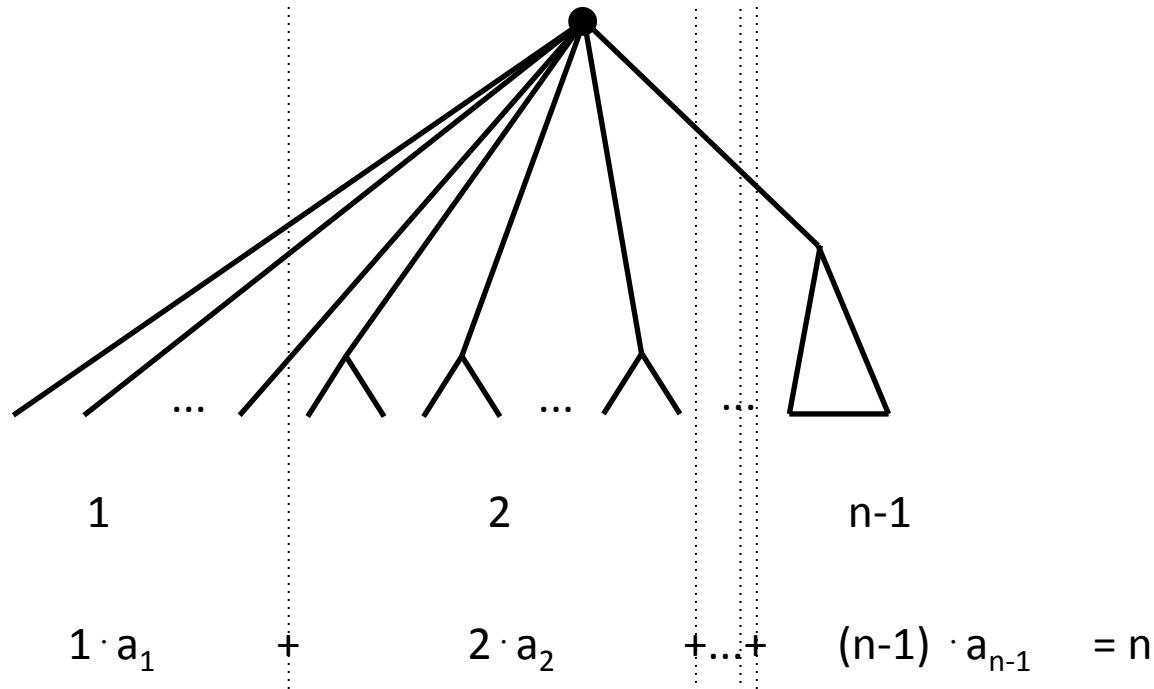


1. Generate skeleton trees





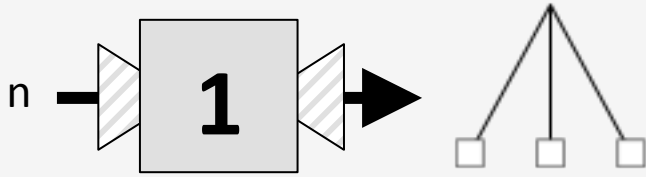
1. Generate skeleton trees



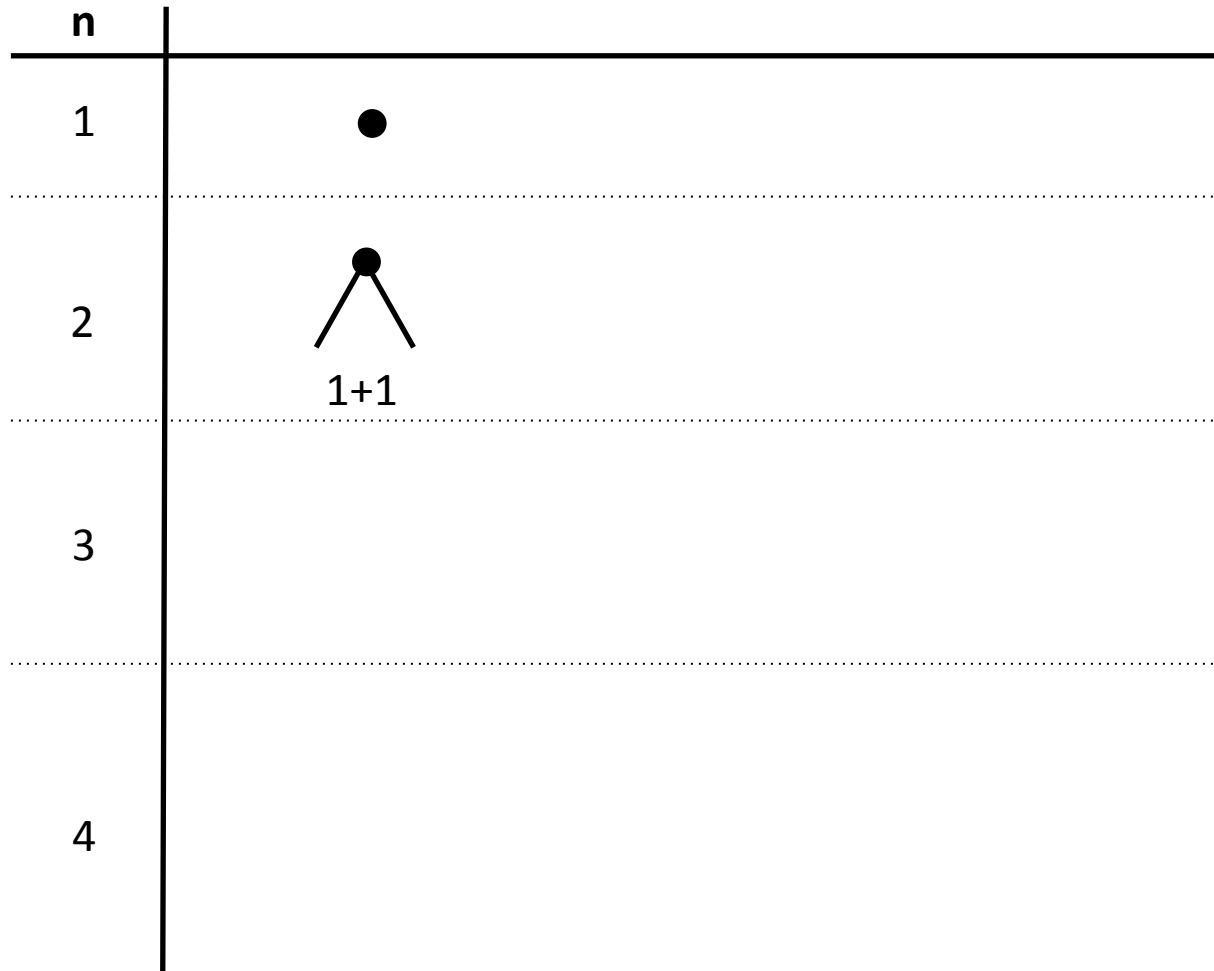
Partition problem:

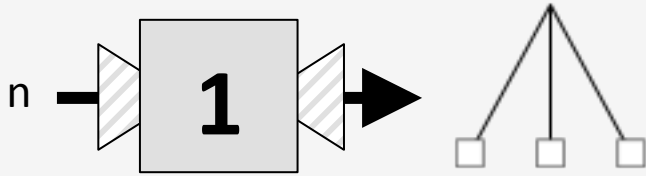
given an integer n

Question : Obtain all possible ways to write n as sum of $1, \dots, n-1$

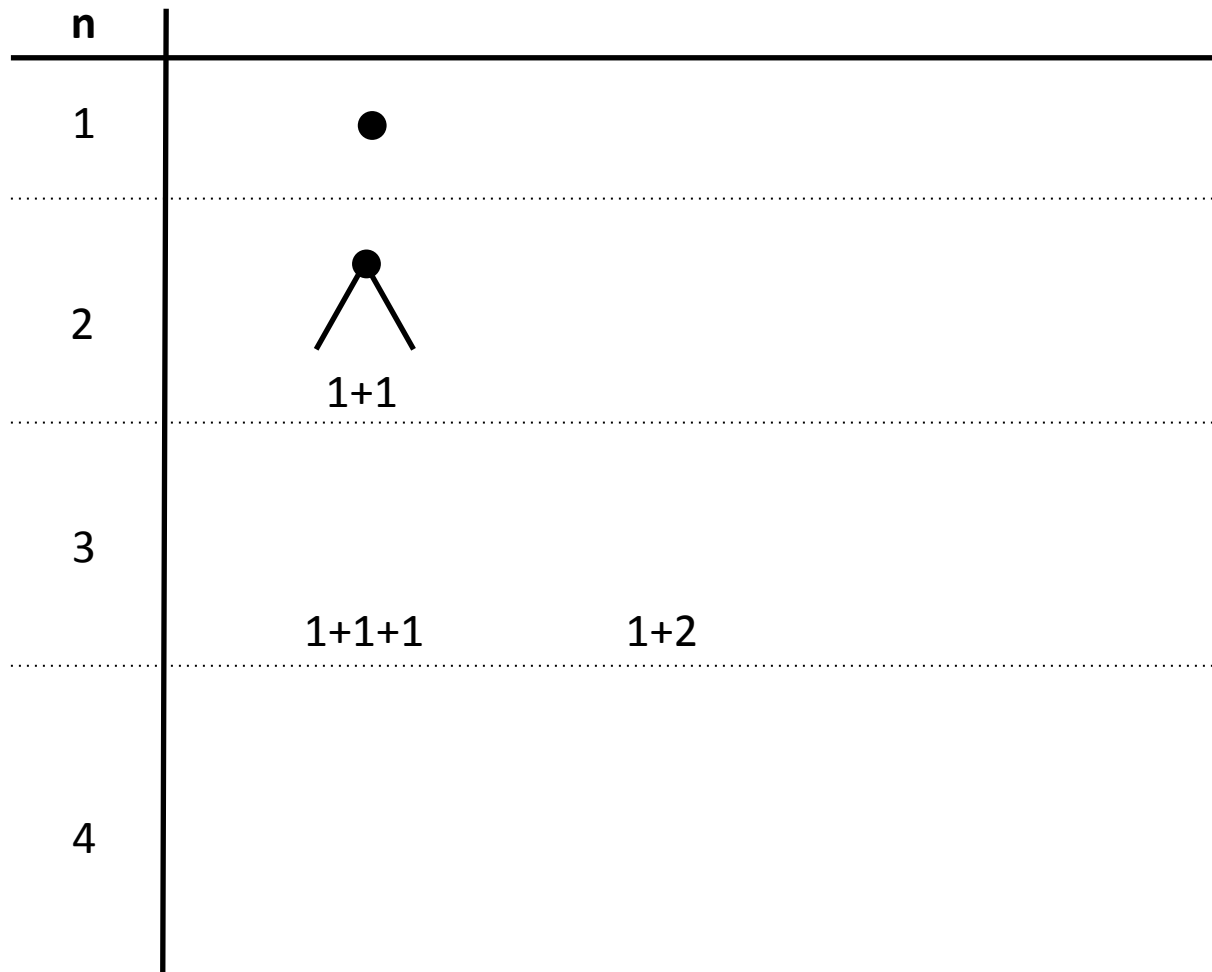


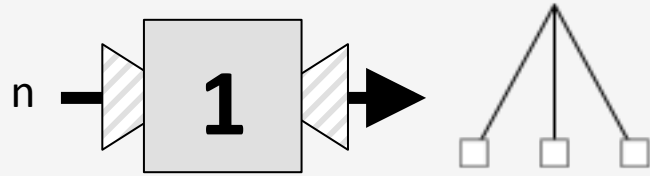
1. Generate skeleton trees



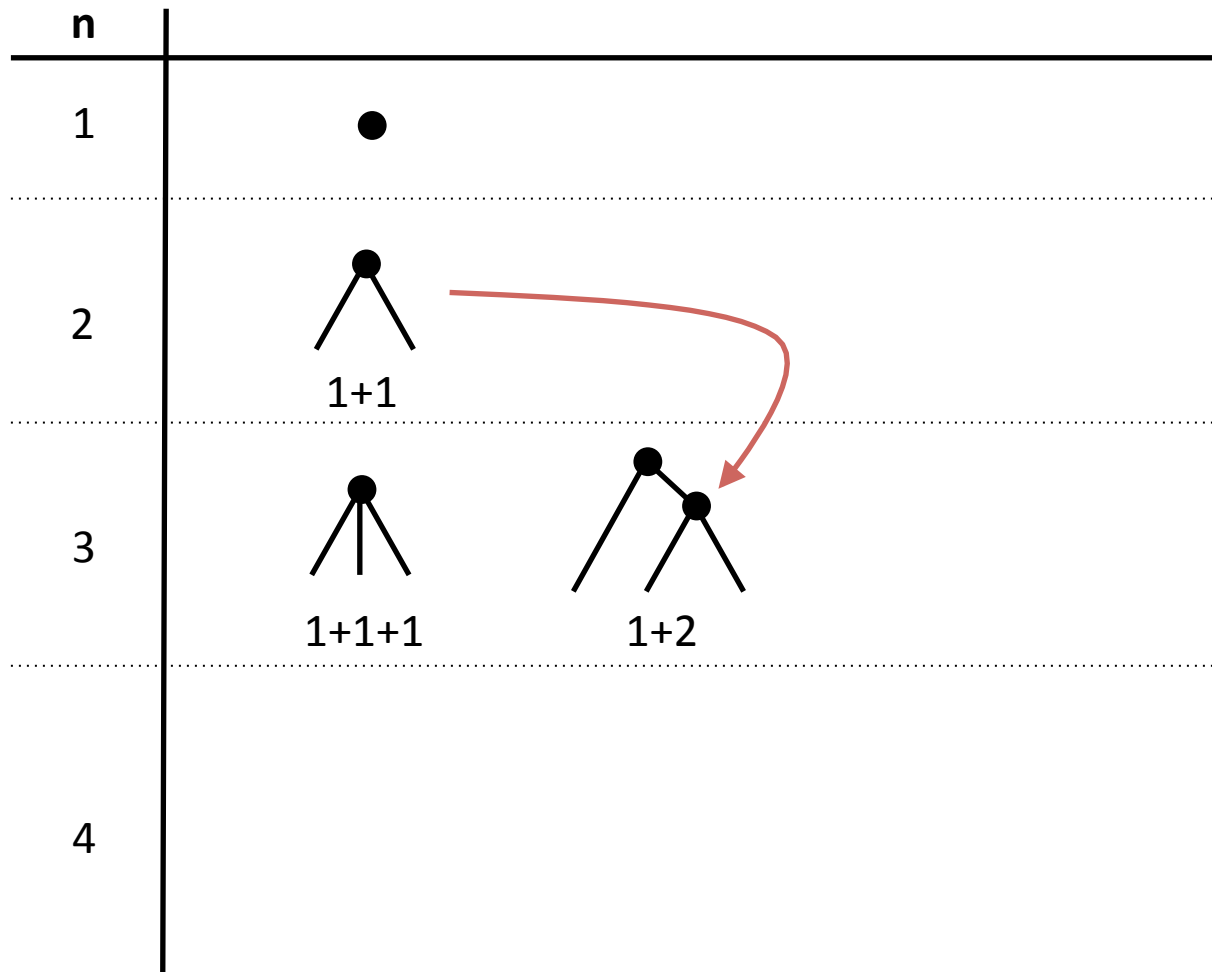


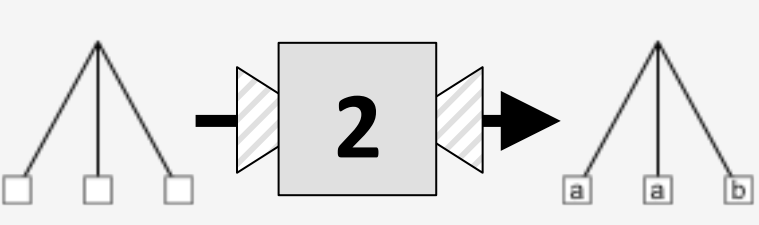
1. Generate skeleton trees





1. Generate skeleton trees

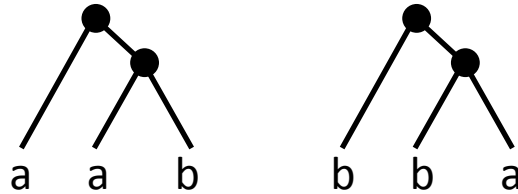


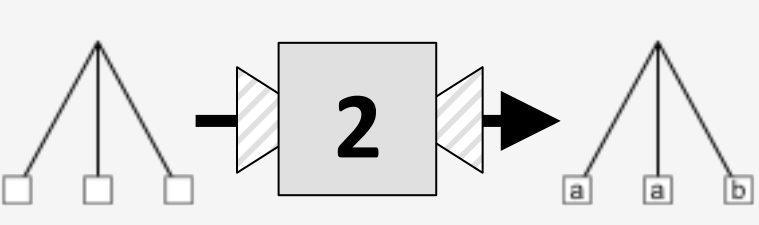


2. Generate labels for skeleton trees

We want:

1. Only one representative over all permutations

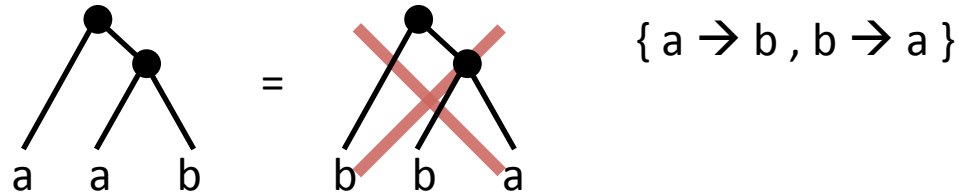


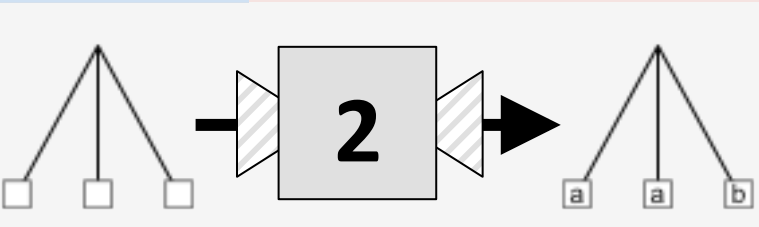


2. Generate labels for skeleton trees

We want:

1. Only one representative over all permutations

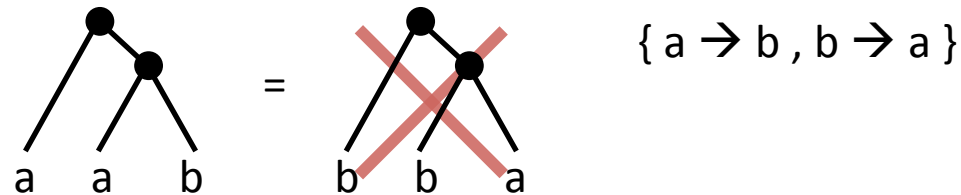




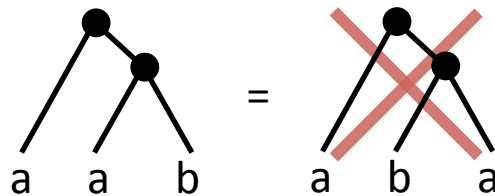
2. Generate labels for skeleton trees

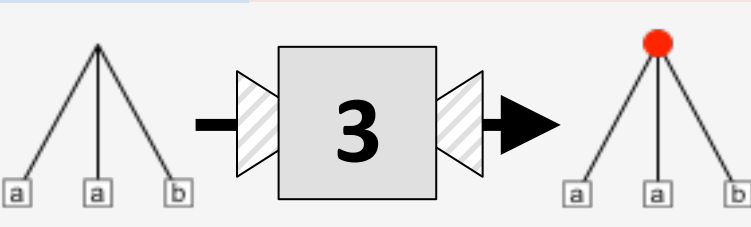
We want:

1. Only one representative over all permutations



2. Only one representative of the isomorphism classes

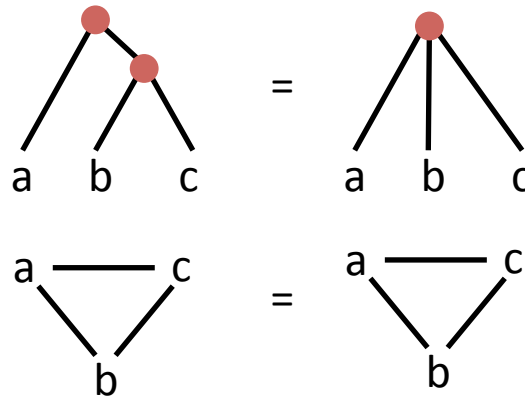


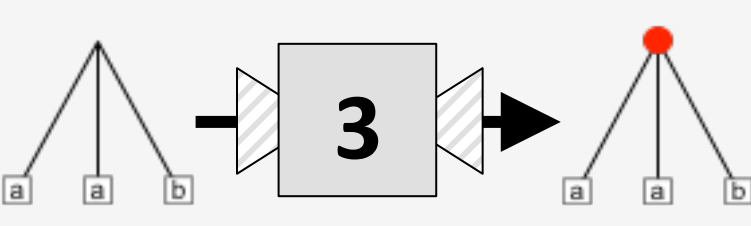


3. Coloring the labeled skeleton trees

We want:

No edge with the same event (color) on both ends

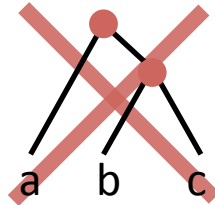




3. Coloring the labeled skeleton trees

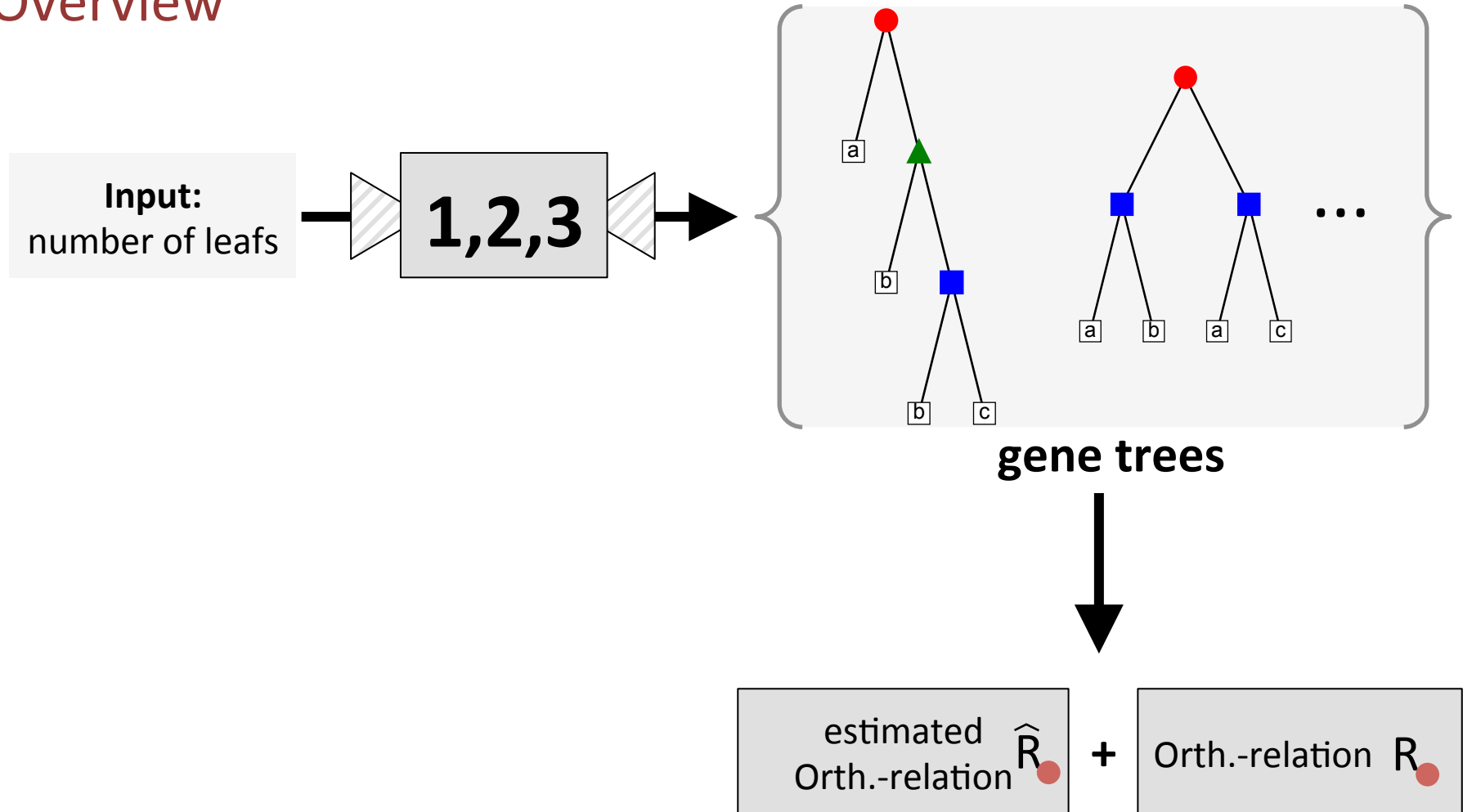
We want:

No edge with the same event (color) on both ends



= 3-coloring of the inner nodes

Overview



Results for n=4 leafs

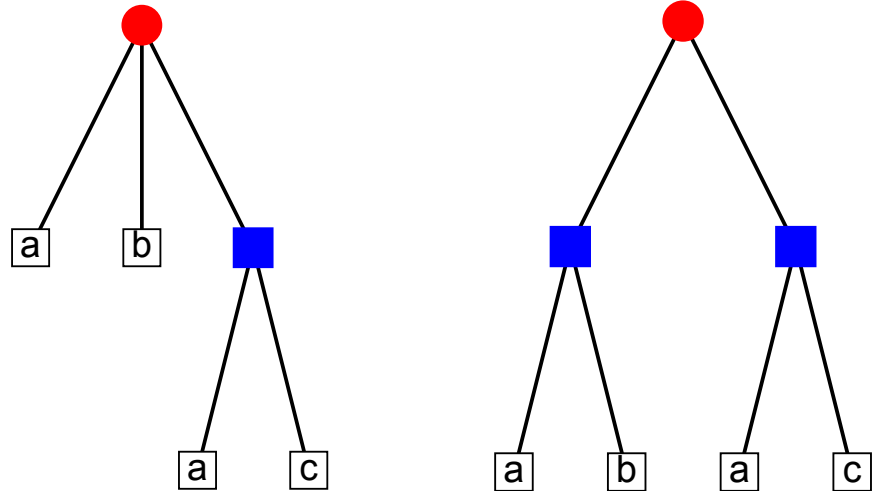
Total number of gene trees : **324**

Non-cograph cases : **27**

Results for n=4 leafs

Total number of gene trees : **324**

Non-cograph cases : **27**

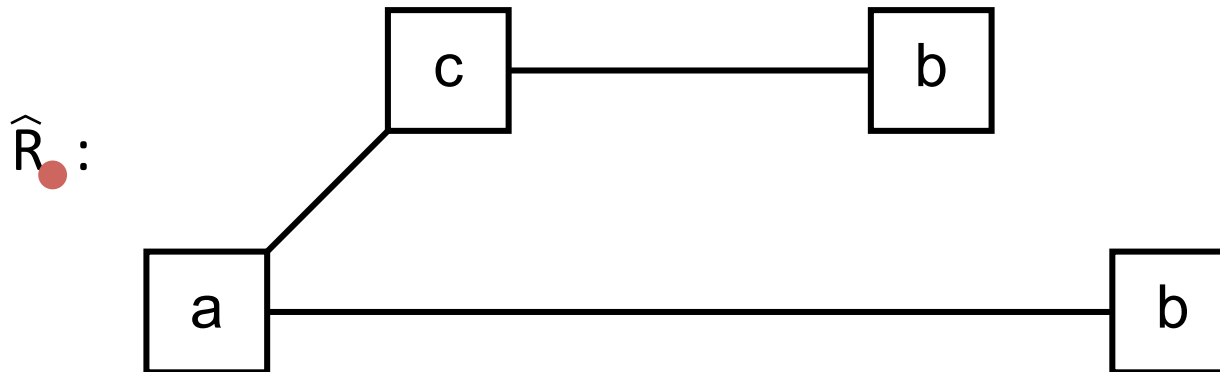
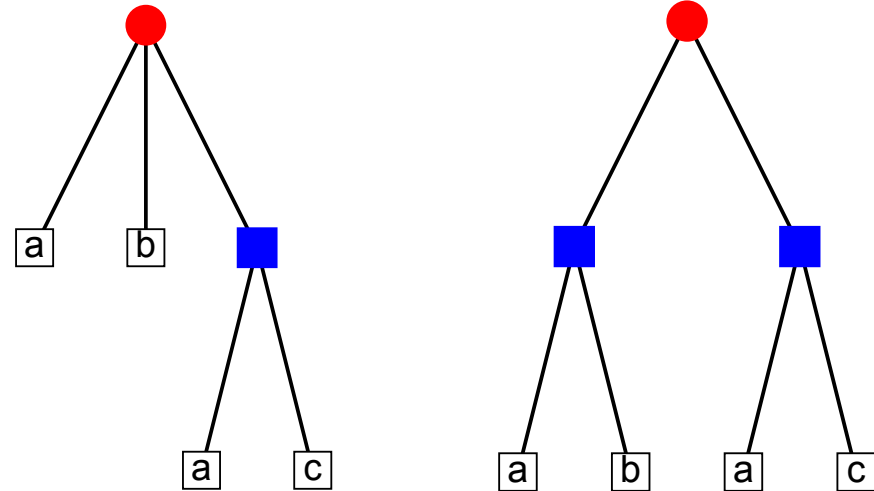


Results for n=4 leafs

Total number of gene trees : **324**

Non-cograph cases : **27**

Number of problemclasses : **1**



problemclass #1
27 generating gene trees

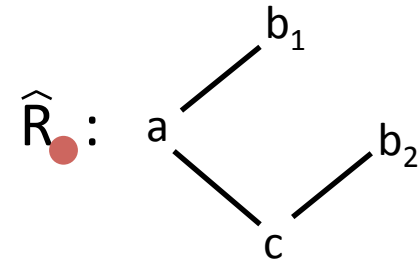
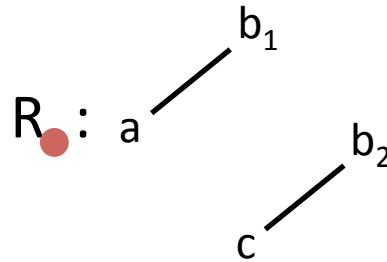
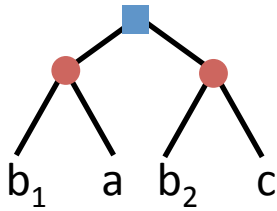
Results for n=5 leafs

Total number of gene trees : **3543**

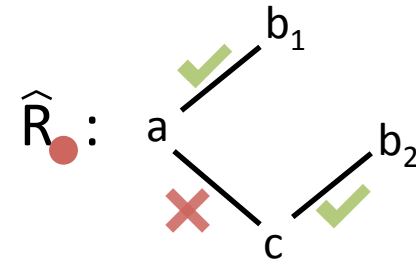
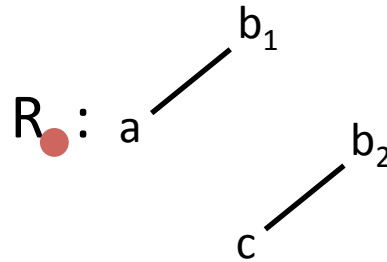
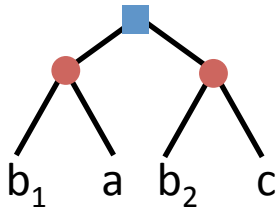
Non-cograph cases : **822**

Number of problemclasses : **9**

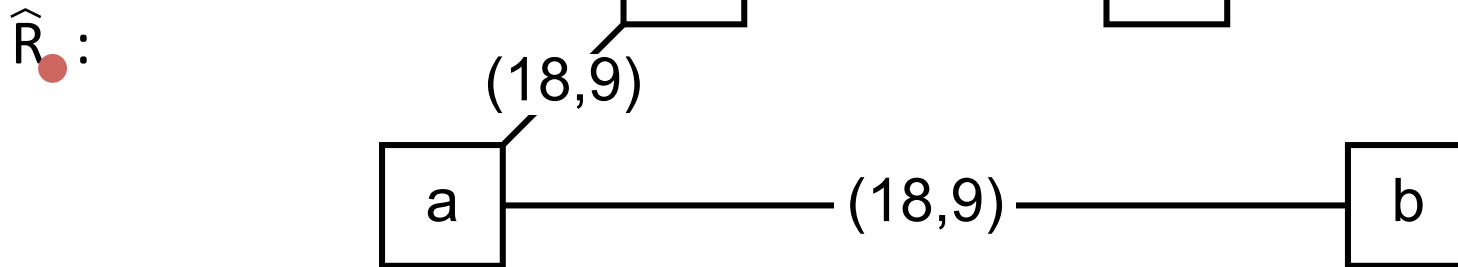
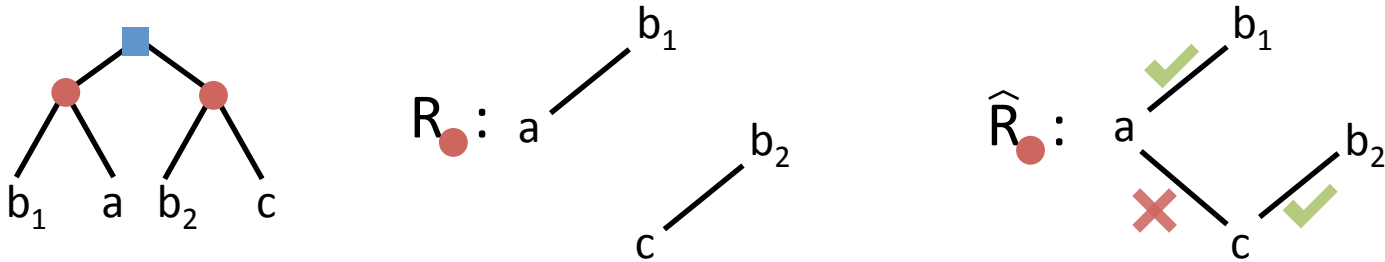
Question 1 : Is there a preference of right/wrong edges?



Question 1 : Is there a preference of right/wrong edges?



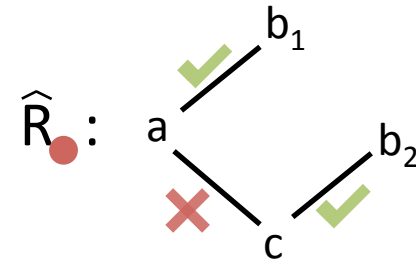
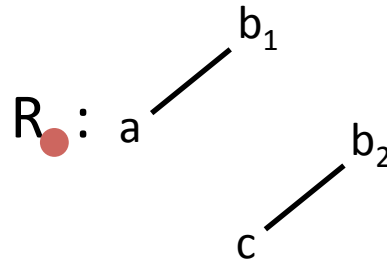
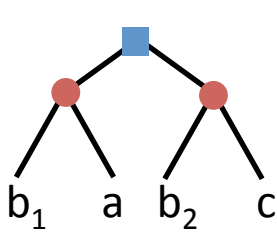
Question 1 : Is there a preference of right/wrong edges?



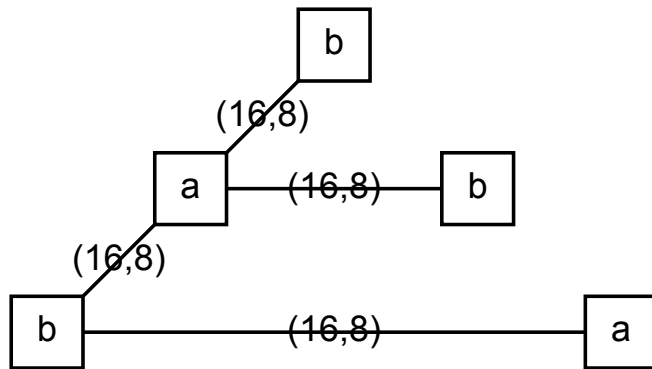
problemclass #1

27 generating gene trees

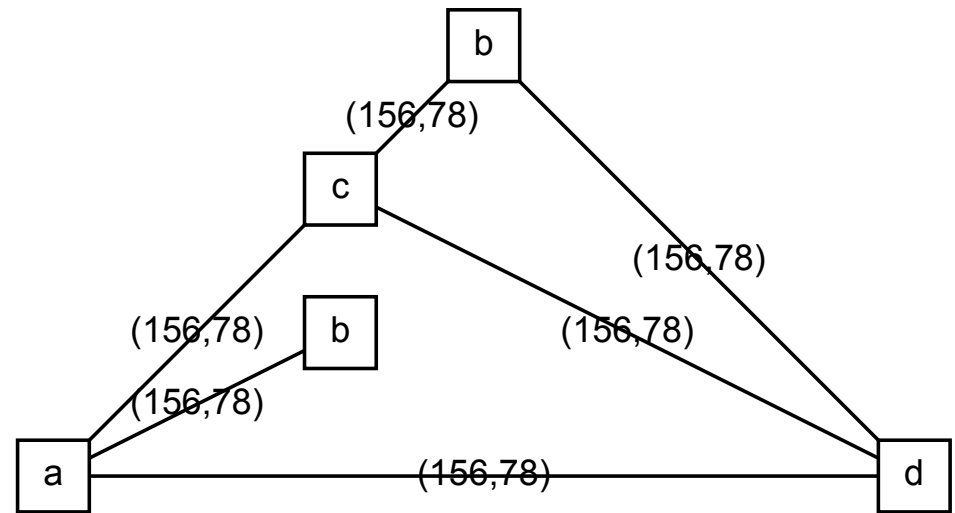
Question 1 : Is there a preference of right/wrong edges?



\hat{R} :

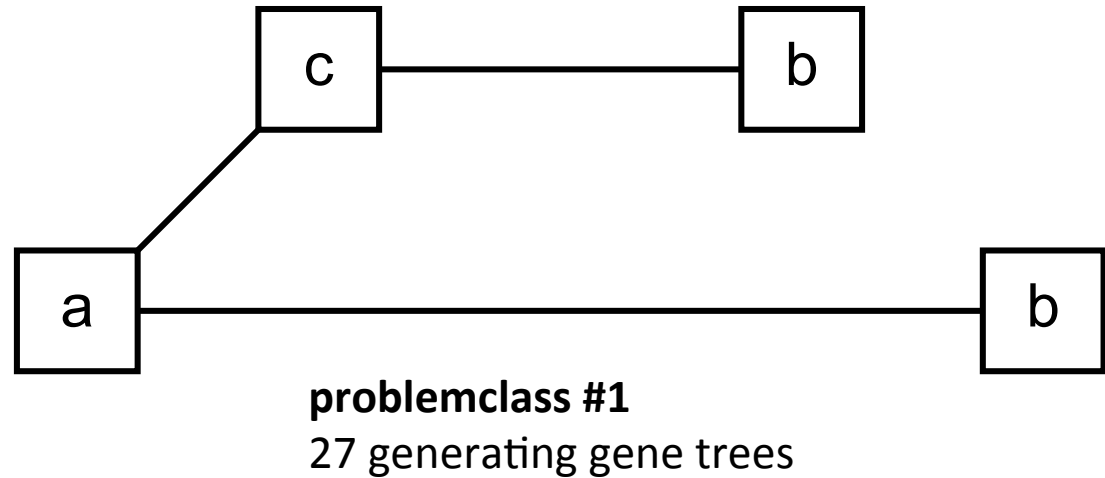


problemclass #3
24 generating gene trees

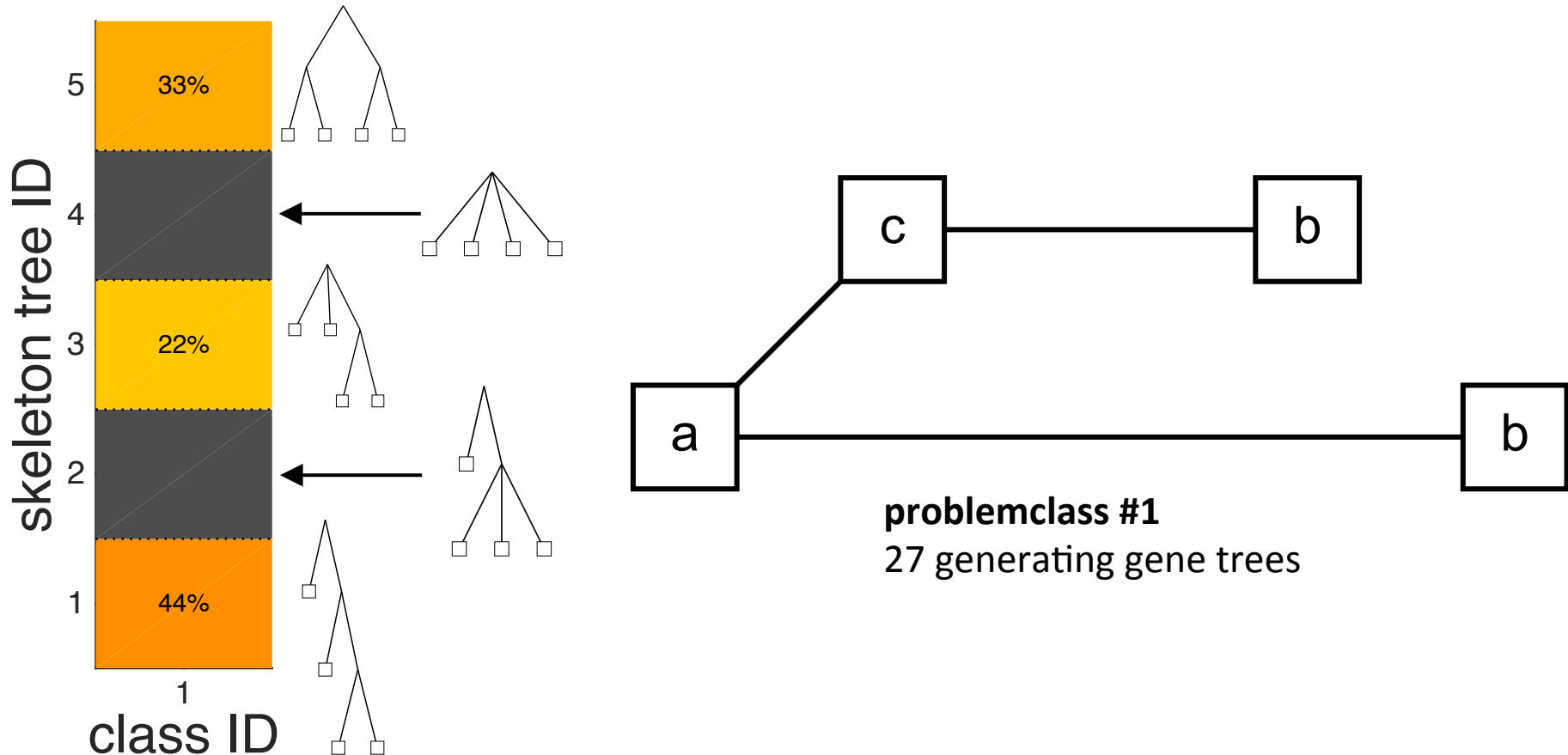


problemclass #7
234 generating gene trees

Question 2 : Can one infer information about the skeleton?



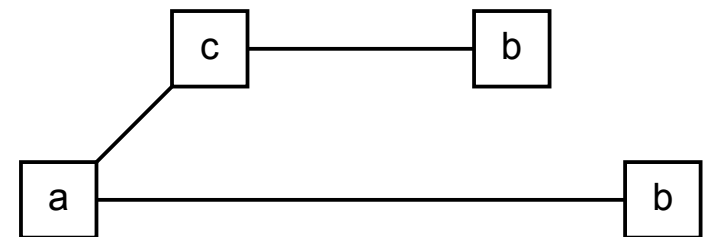
Question 2 : Can one infer information about the skeleton?



Question 2 : Can one infer information about the skeleton?

Consider \hat{R} with leaf distance d as additional information.

$d(x,y) =$
length of the shortest path
between x and y

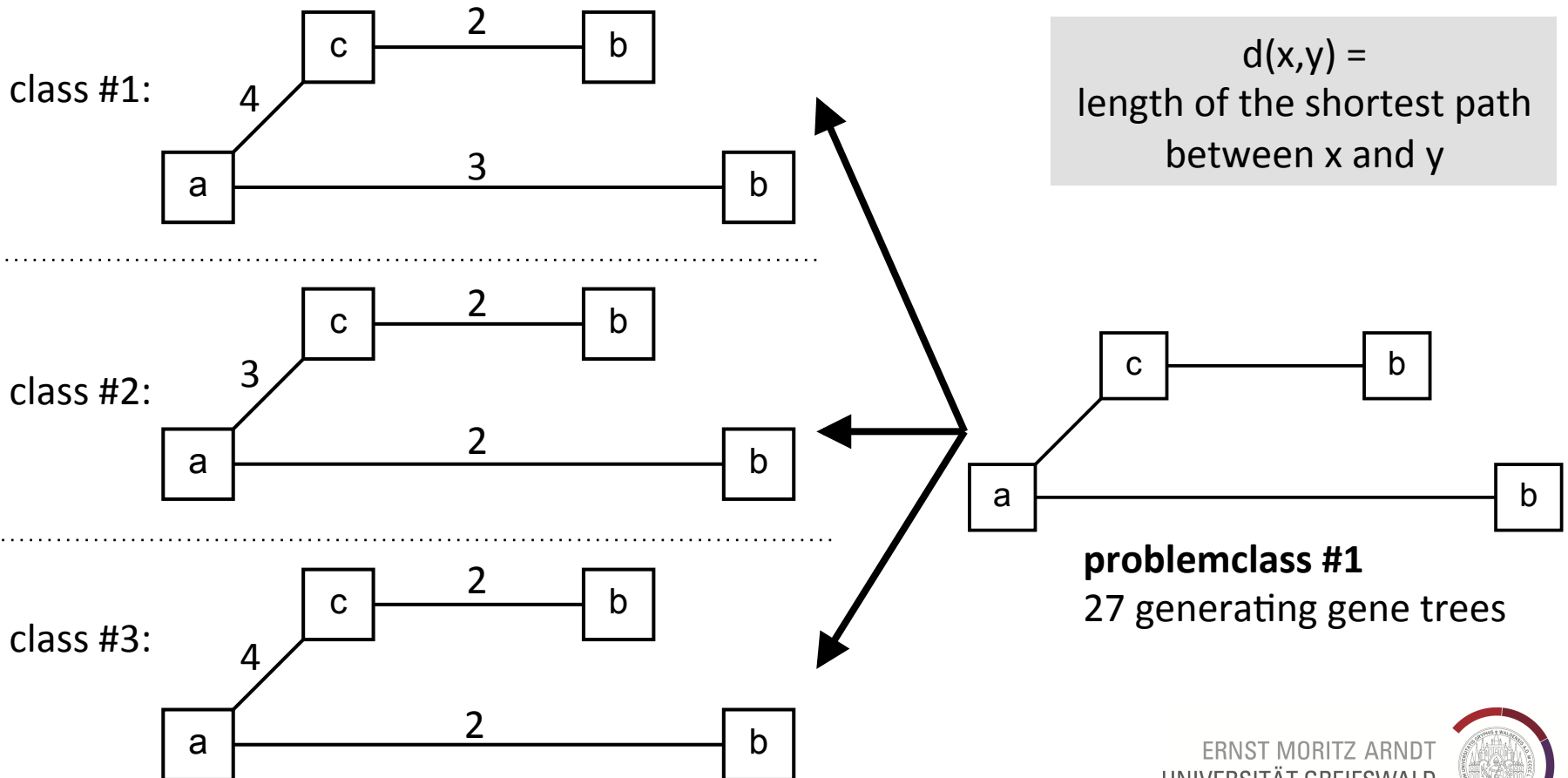


problemclass #1

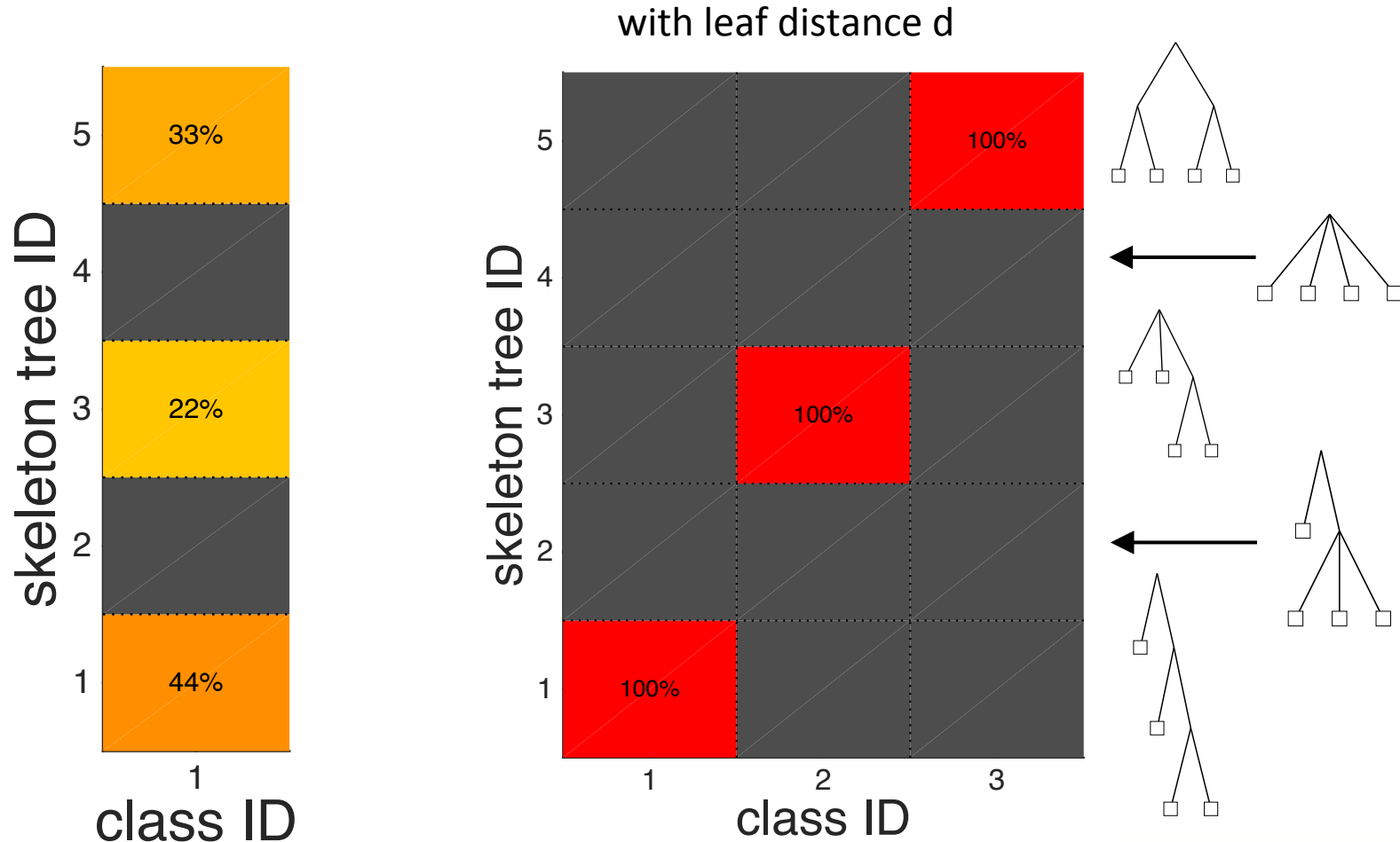
27 generating gene trees

Question 2 : Can one infer information about the skeleton?

Consider \hat{R}_\bullet with leaf distance d as additional information.

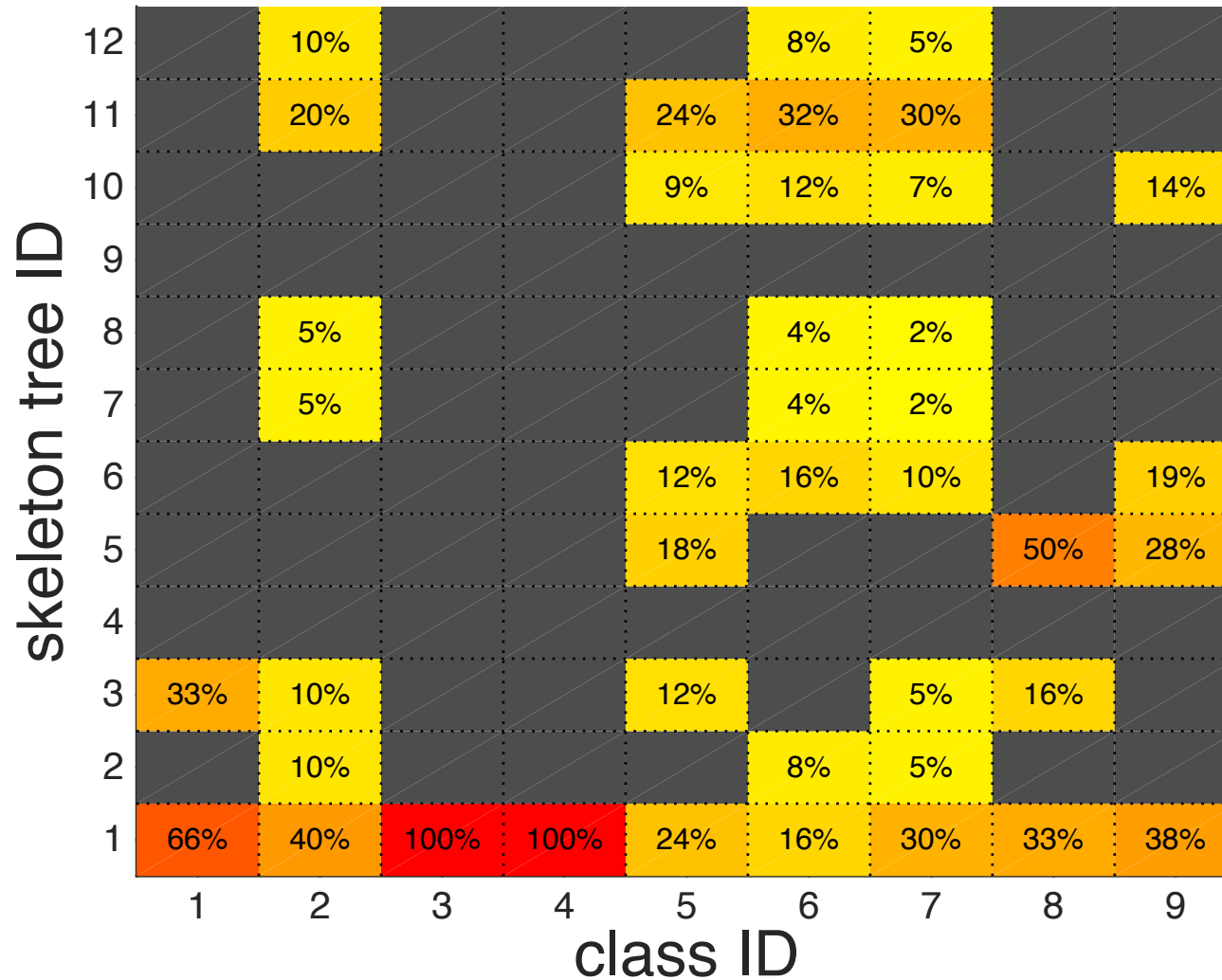


Question 2 : Can one infer information about the skeleton?

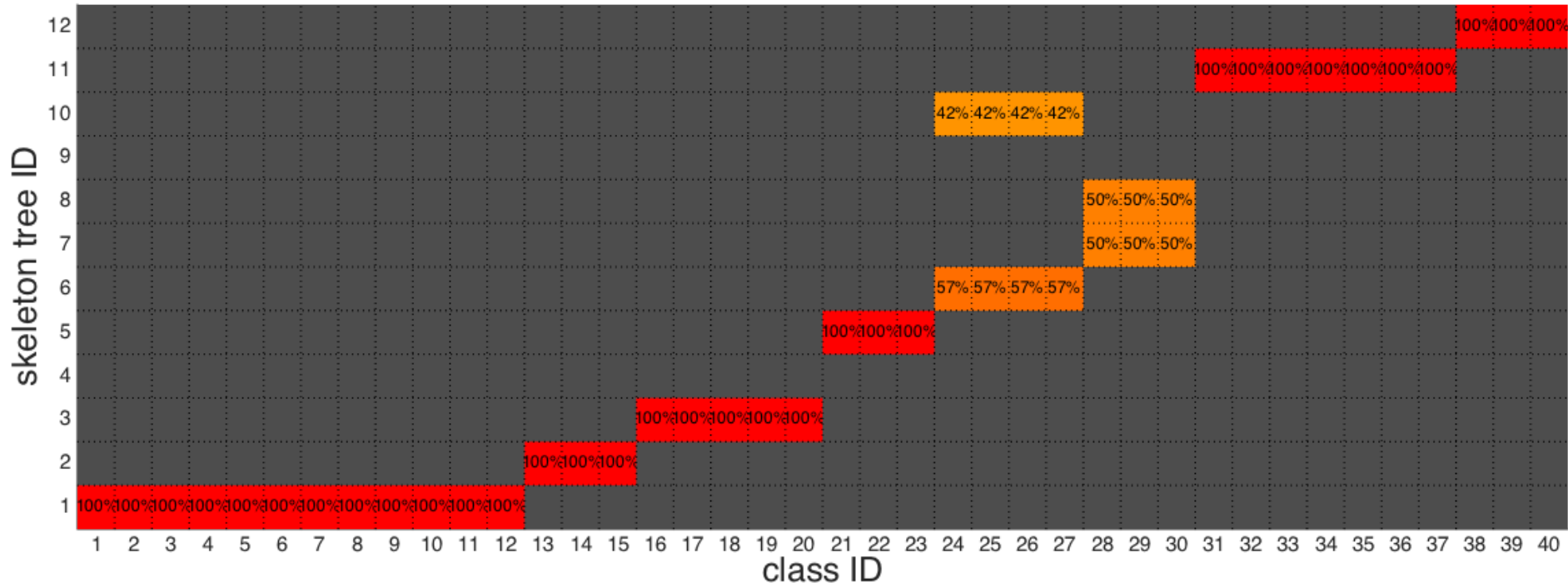


Question 2 : Can one infer information about the skeleton?

n=5



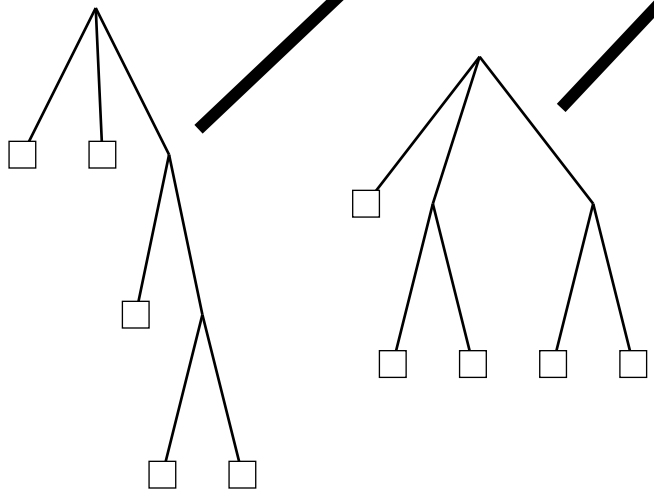
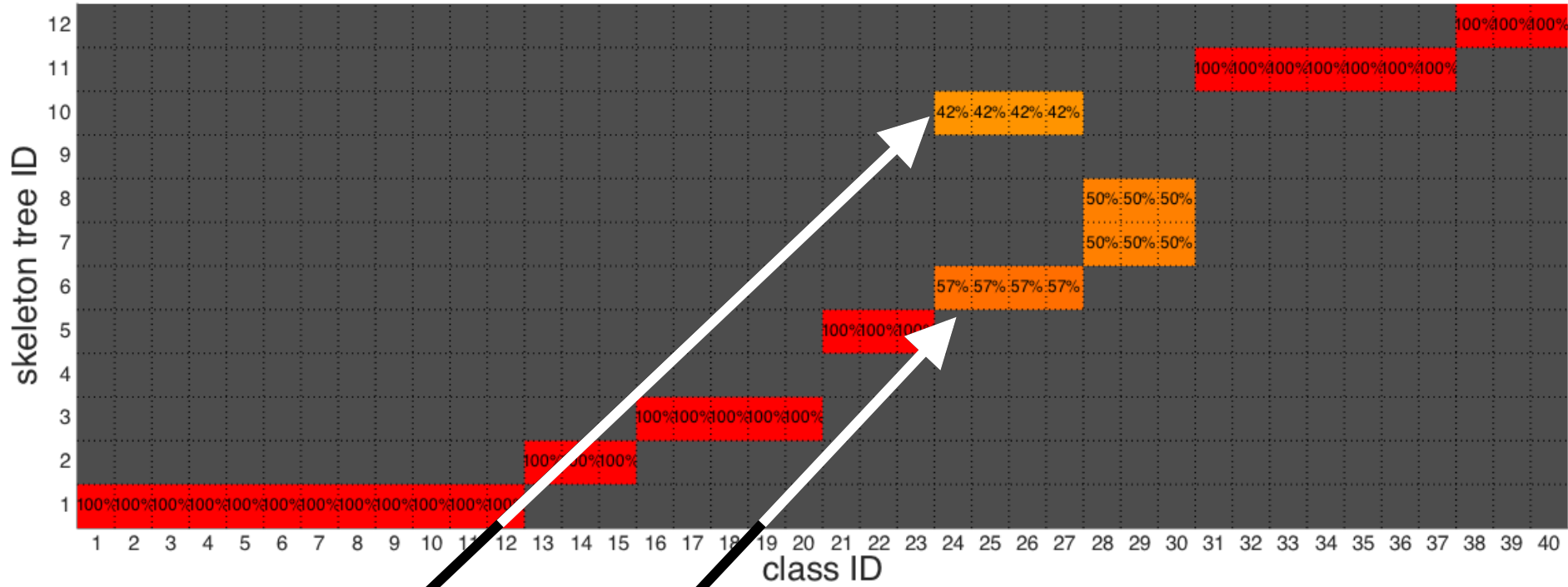
Question 2 : Can one infer information about the skeleton?



n=5 with leaf-distant d

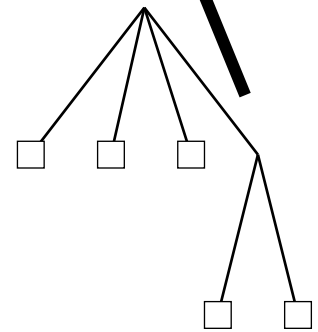
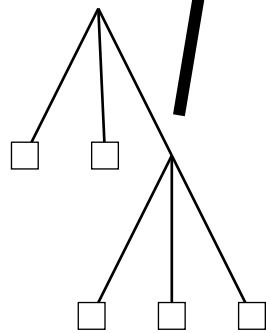
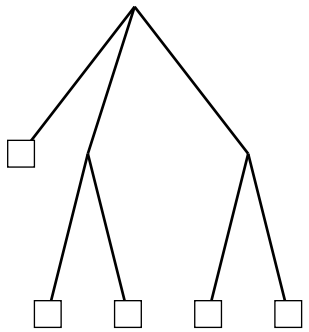
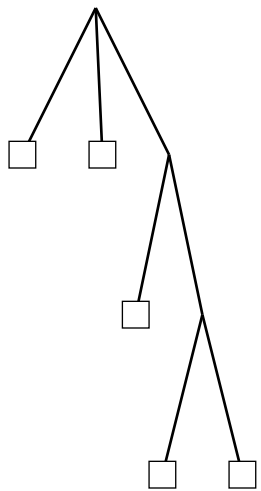
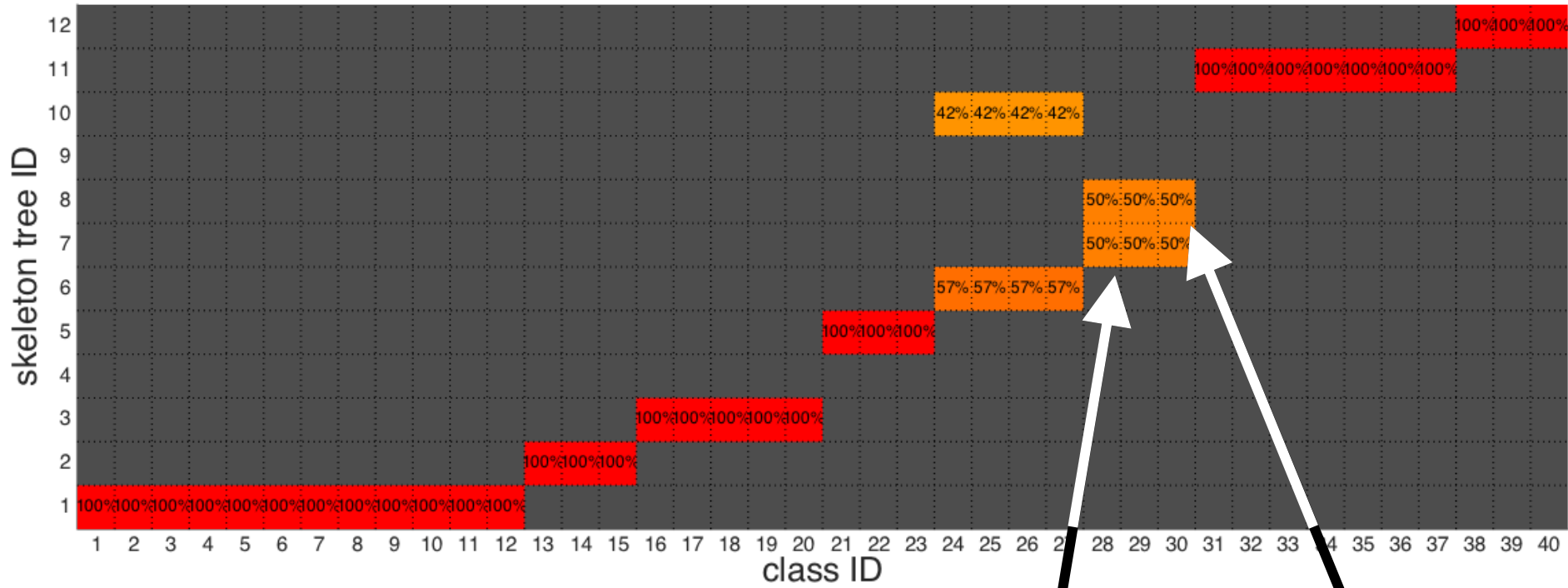


Question 2 : Can one infer information about the skeleton?



n=5 with leaf-distant d

Question 2 : Can one infer information about the skeleton?



n=5 with leaf-distant d



Summary and Outlook

Summary:

- Generated all possible gene trees for up to 7 leafs
- No information about right or wrong placed edges in the problemclasses
- The leaf distance d gives us the unrooted skeleton tree

For the future:

- Find a more efficient way to generate the leaf-labeling
- Define and investigate more questions

Summary and Outlook

Summary:

- Generated all possible gene trees for up to 7 leafs
- No information about right or wrong placed edges in the problemclasses
- The leaf distance d gives us the unrooted skeleton tree

For the future:


- Find a more efficient way to generate the leaf-labeling
- Define and investigate more questions

THANK YOU!

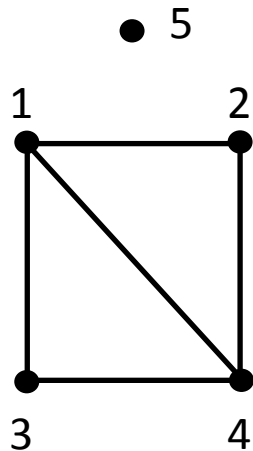
Summary and Outlook

THANK YOU!


P4-free Graphs / Cographs

A graph is called **P4-free / cograph** if and only if there exists no induced subgraph on 4 nodes that is a P4 ()

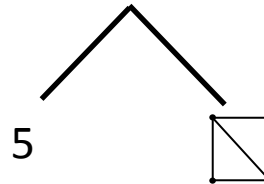
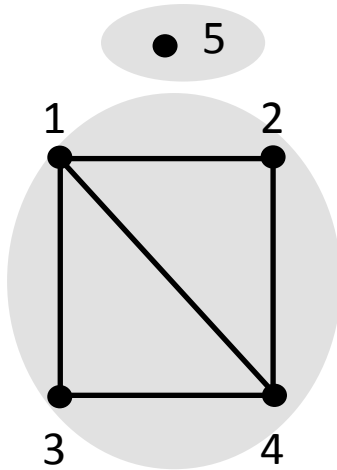
Every cograph is uniquely associated with a **cotree**
= construction instruction to create the cograph




P4-free Graphs / Cographs

A graph is called **P4-free / cograph** if and only if there exists no induced subgraph on 4 nodes that is a P4 ()

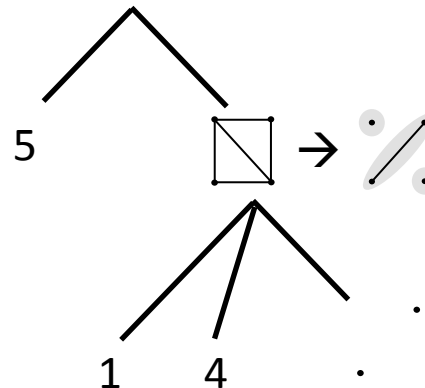
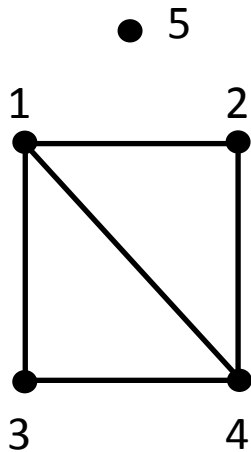
Every cograph is uniquely associated with a **cotree**
= construction instruction to create the cograph




P4-free Graphs / Cographs

A graph is called **P4-free / cograph** if and only if there exists no induced subgraph on 4 nodes that is a P4 ()

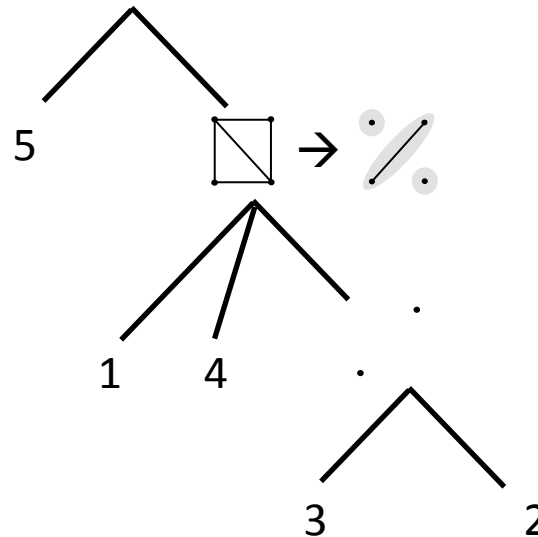
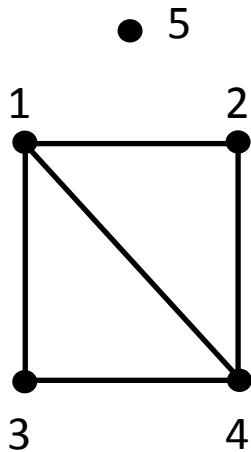
Every cograph is uniquely associated with a **cotree**
= construction instruction to create the cograph




P4-free Graphs / Cographs

A graph is called **P4-free / cograph** if and only if there exists no induced subgraph on 4 nodes that is a P4 ()

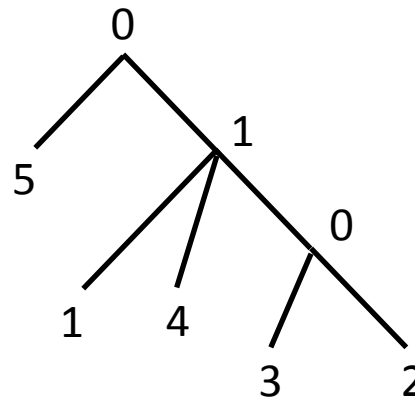
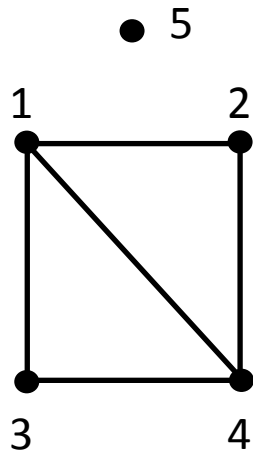
Every cograph is uniquely associated with a **cotree**
= construction instruction to create the cograph




P4-free Graphs / Cographs

A graph is called **P4-free / cograph** if and only if there exists no induced subgraph on 4 nodes that is a P4 ()

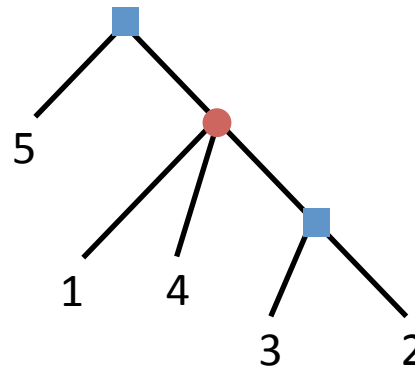
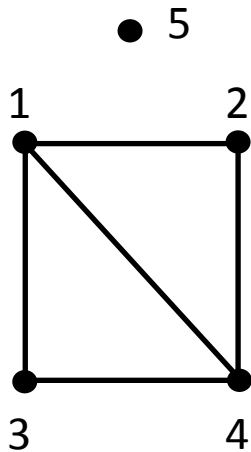
Every cograph is uniquely associated with a **cotree**
= construction instruction to create the cograph



P4-free Graphs / Cographs

A graph is called **P4-free / cograph** if and only if there exists no induced subgraph on 4 nodes that is a P4 ()

Every cograph is uniquely associated with a **cotree**
= construction instruction to create the cograph

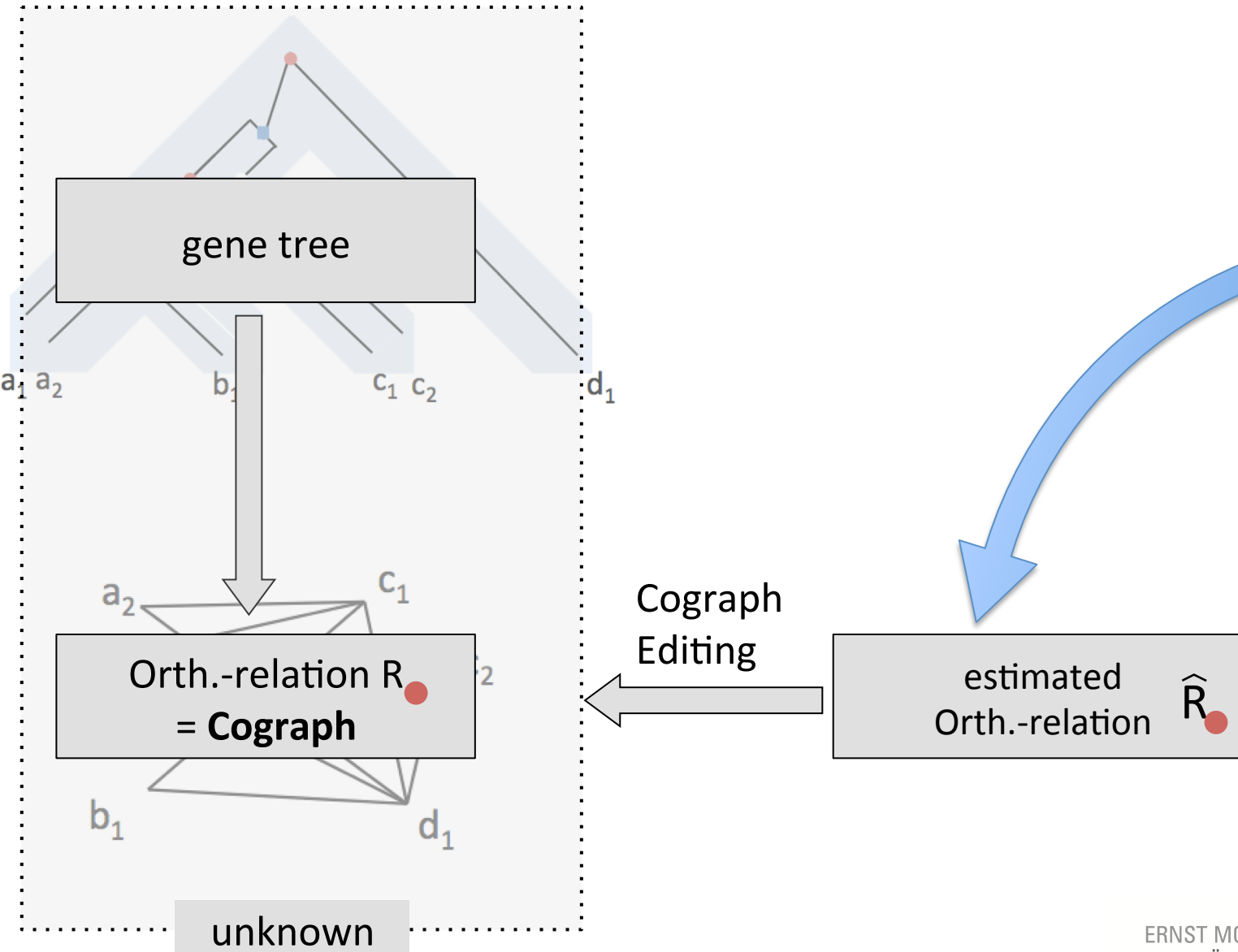


Introduction

Artificial
Data
Analysis

Summary
and
Outlook

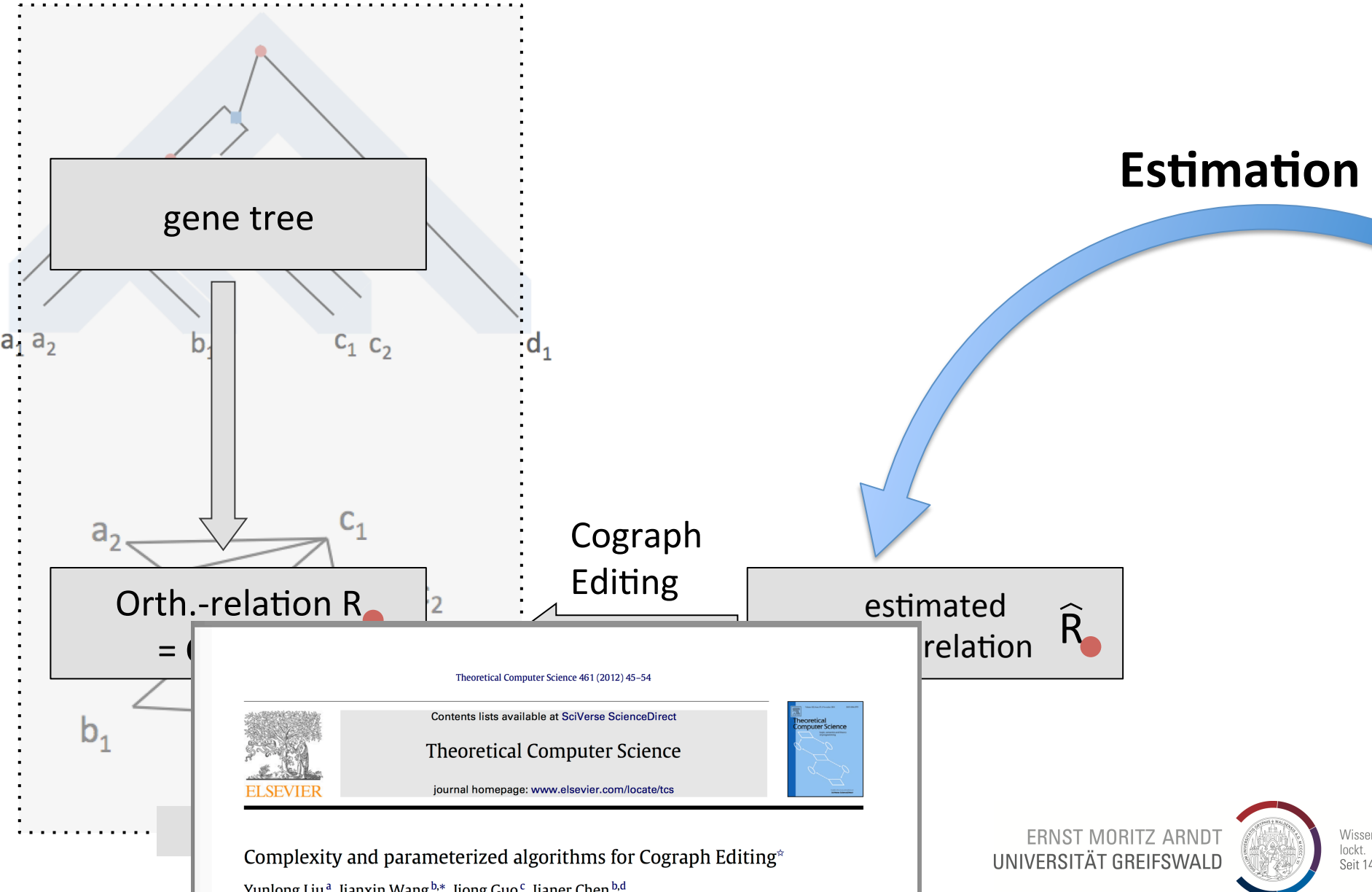
Estimation



Introduction

Artificial
Data
Analysis

Summary
and
Outlook

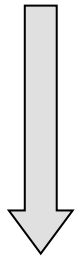


Complexity and parameterized algorithms for Cograph Editing*

Yunlong Liu^a, Jianxin Wang^{b,*}, Jiong Guo^c, Jianer Chen^{b,d}

Graph-based orthology inference (e.g. proteinOrtho)

Input : sequence data



Local alignment search
(e.g. BLAST)

Sequence similarity $s(x,y)$

Then two genes x (from species A) and y (from species B) are **estimated orthologs** if:

- i. $A \neq B$
- ii. The sequence similarity $s(x,y)$ is greater than the one between x to all other genes of B and y to all other genes of A

Graph-based orthology inference (e.g. proteinOrtho)

x (from specie A) and y (from specie B) are **estimated orthologs** if:

- i. $A \neq B$
- ii. The sequence similarity $s(\mathbf{x}, \mathbf{y})$ is greater than the one between x to all other genes of B and y to all other genes of A

Simplification: leaf distance $d(x,y)$ of the (unknown) gene tree:

$d(x,y)$ = length of the shortest path between x and y

Graph-based orthology inference (e.g. proteinOrtho)

x (from specie A) and y (from specie B) are **estimated orthologs** if:

- i. $A \neq B$
- ii. The sequence similarity $s(\mathbf{x}, \mathbf{y})$ is greater than the one between x to all other genes of B and y to all other genes of A

Simplification: leaf distance $d(\mathbf{x}, \mathbf{y})$ of the (unknown) gene tree:

$d(\mathbf{x}, \mathbf{y})$ = length of the shortest path between x and y

$$d \sim 1/s$$

Graph-based orthology inference (e.g. proteinOrtho)

x (from specie A) and y (from specie B) are **estimated orthologs** if:

- i. $A \neq B$
- ii. The sequence similarity $s(x,y)$ is greater than the one between x to all other genes of B and y to all other genes of A

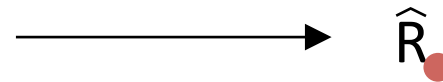
Simplification: leaf distance $d(x,y)$ of the (unknown) gene tree:

$d(x,y)$ = length of the shortest path between x and y

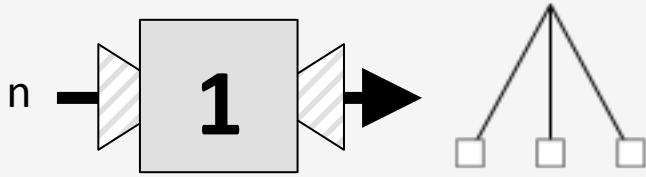
$$d \sim 1/s$$

→ **estimated orthologs** if:

- i. $A \neq B$
- ii. $d(x,y)$ is smaller than the distance between x to all other genes of B and y to all other genes of A



\hat{R}

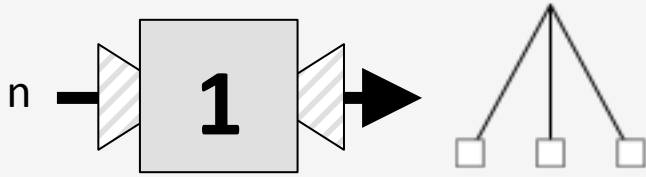


1. Generate skeleton trees

$n = 3$

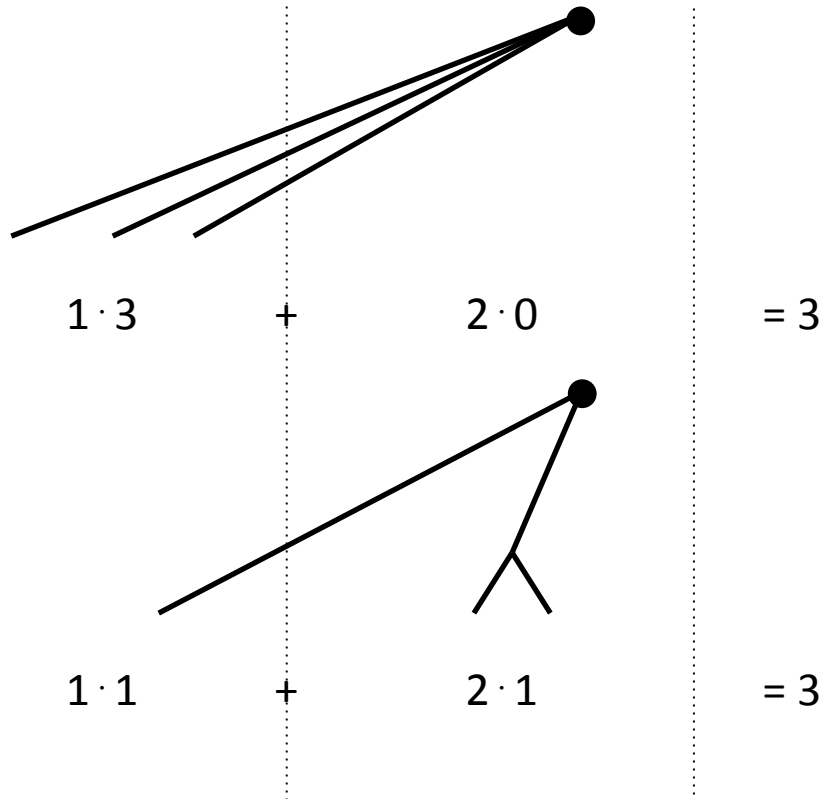
$$1 \cdot 3 + 2 \cdot 0 = 3$$

$$1 \cdot 1 + 2 \cdot 1 = 3$$

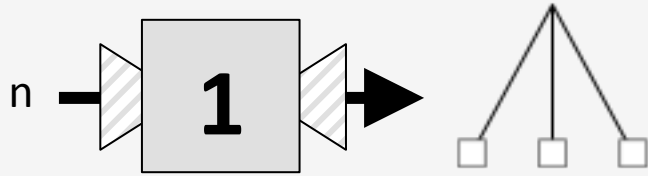


1. Generate skeleton trees

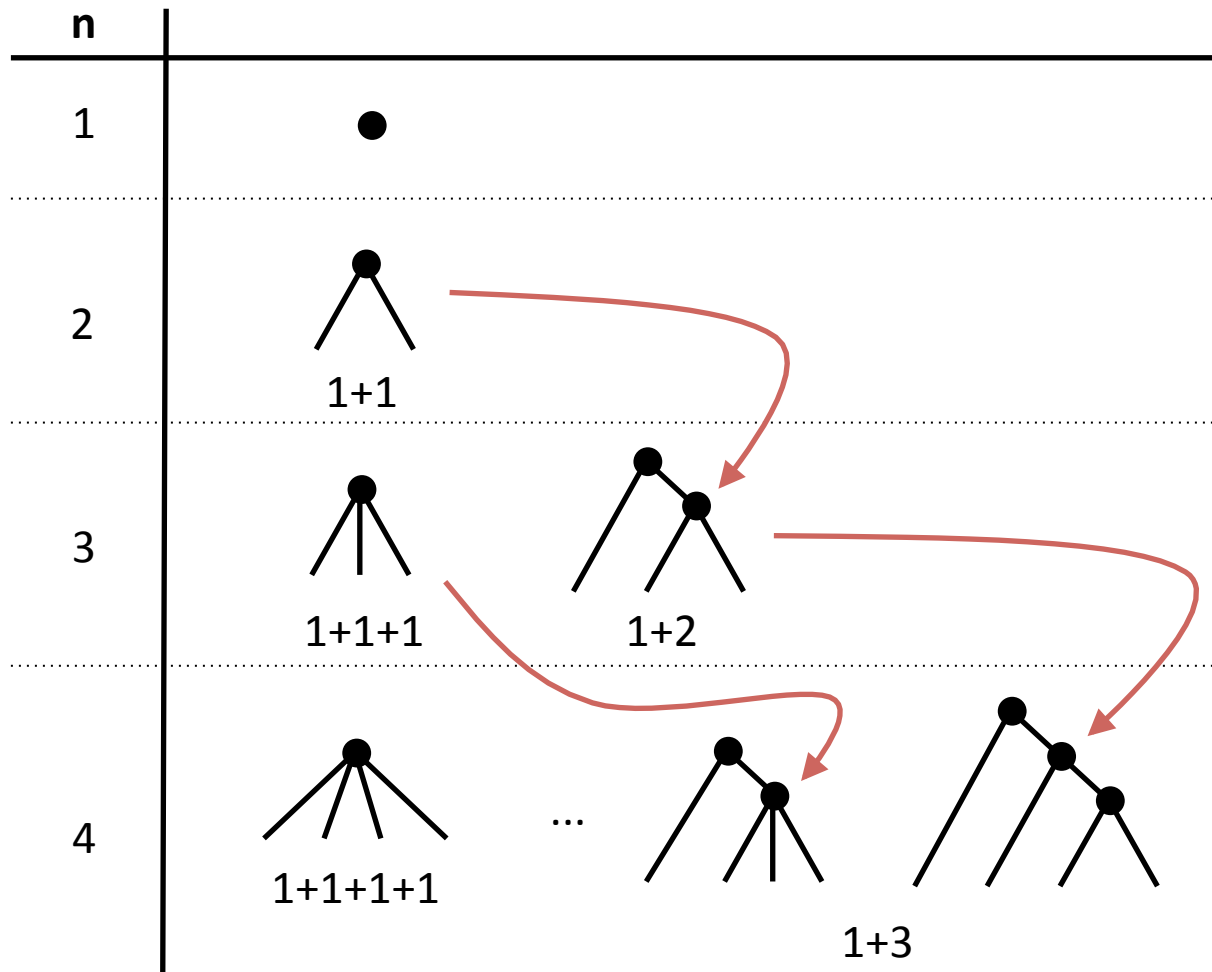
$n = 3$

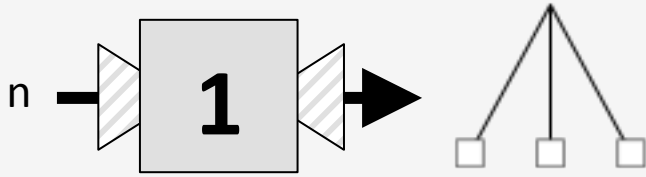


Artificial Data Analysis

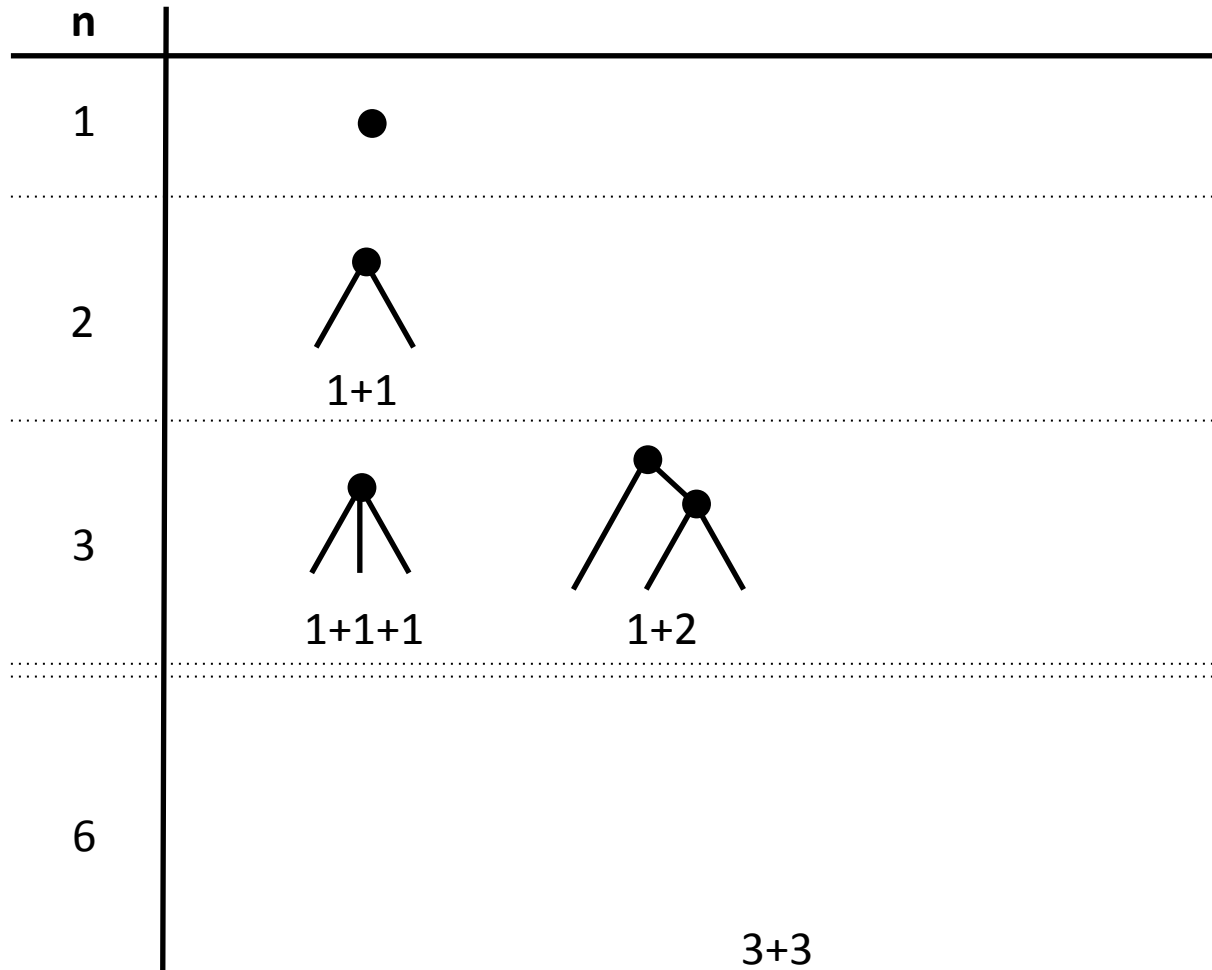


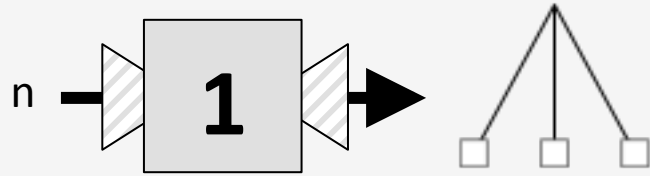
1. Generate skeleton trees



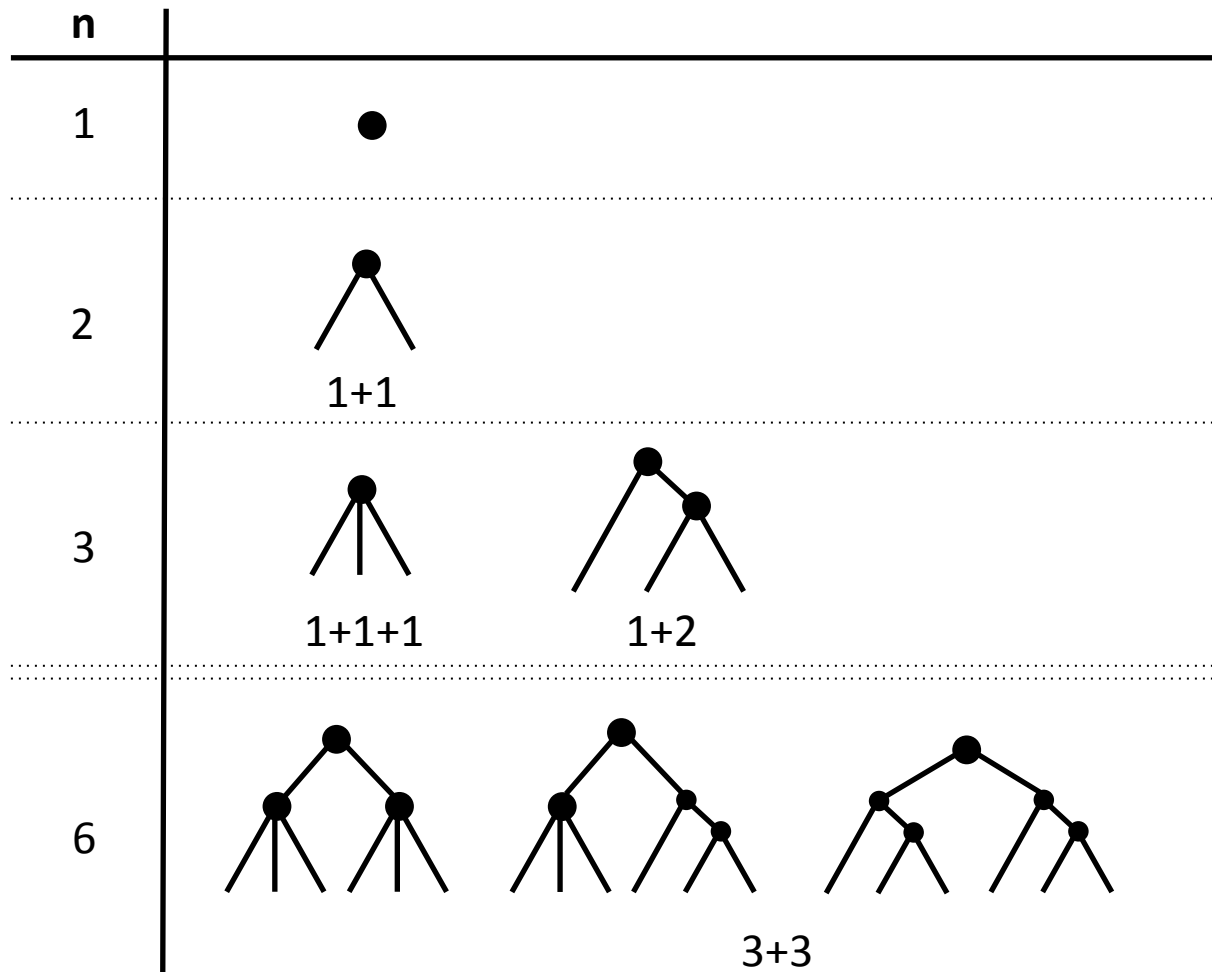


1. Generate skeleton trees

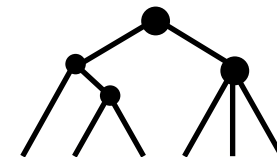


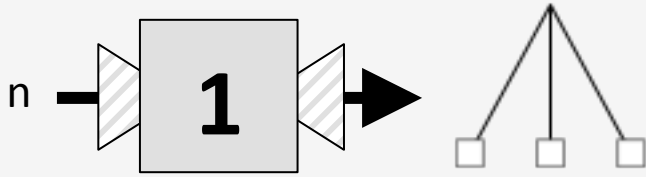


1. Generate skeleton trees



we do not want:

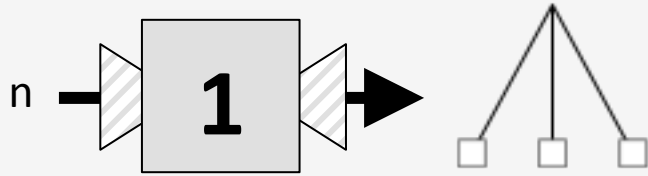




1. Generate skeleton trees

→ Choosing a_i **unordered** elements of the skeleton trees with i leafs **with** replacement

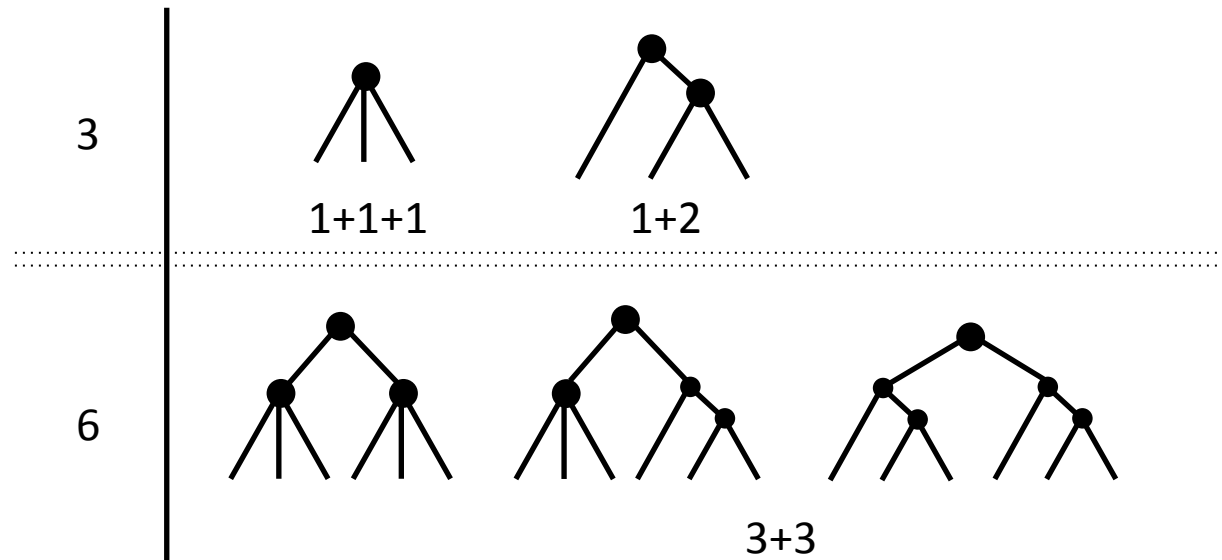
$$\binom{numST_i + (a_i - 1)}{a_i}$$



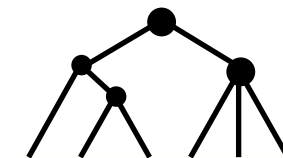
1. Generate skeleton trees

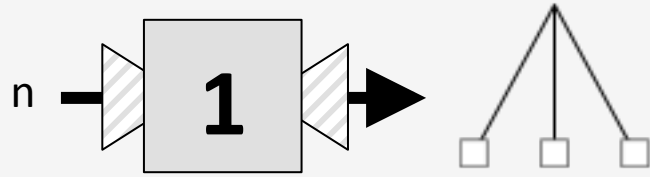
- Choosing a_i **unordered** elements of the skeleton trees with i leaves **with** replacement
- Choosing 2 **unordered** elements of the skeleton trees with 3 leaves **with** replacement

$$\binom{numST_i + (a_i - 1)}{a_i} = 3$$



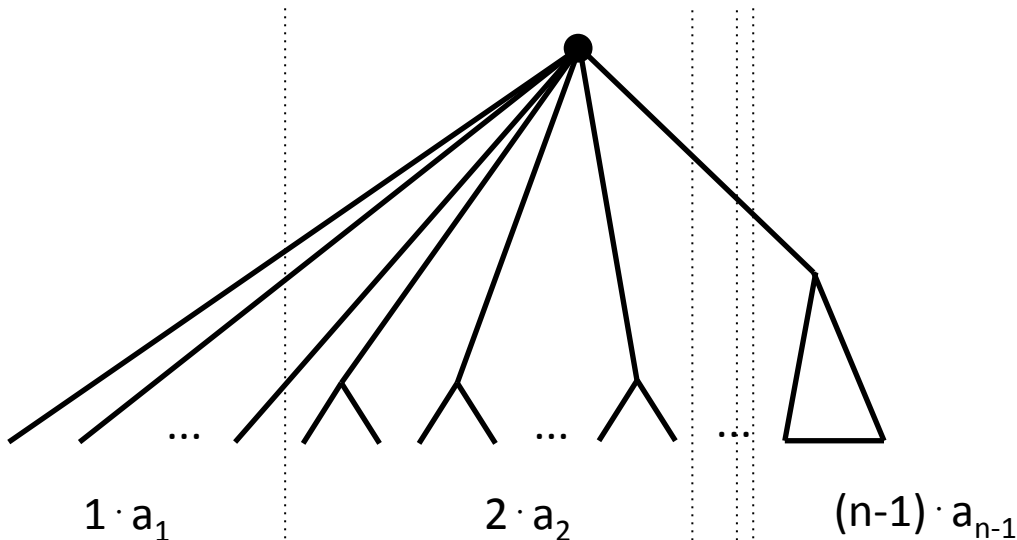
we do not want:





1. Generate skeleton trees

$$\text{num}ST_n = \sum_{(a_1, a_2, \dots, a_{n-1}) \in P^n} \prod_{i=1}^{n-1} \binom{\text{num}ST_i + (a_i - 1)}{a_i}$$



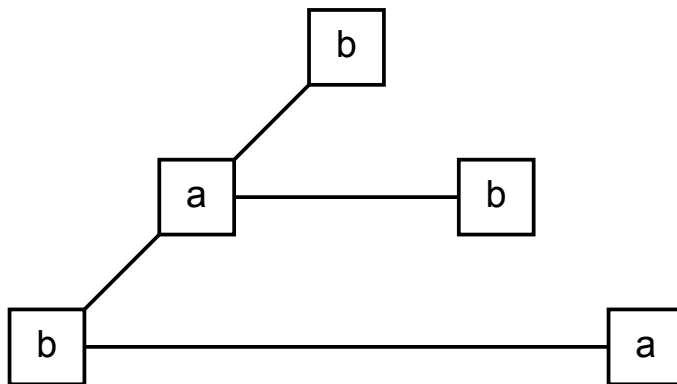
Results for n=5 leafs

Total number of gene trees : **3543**

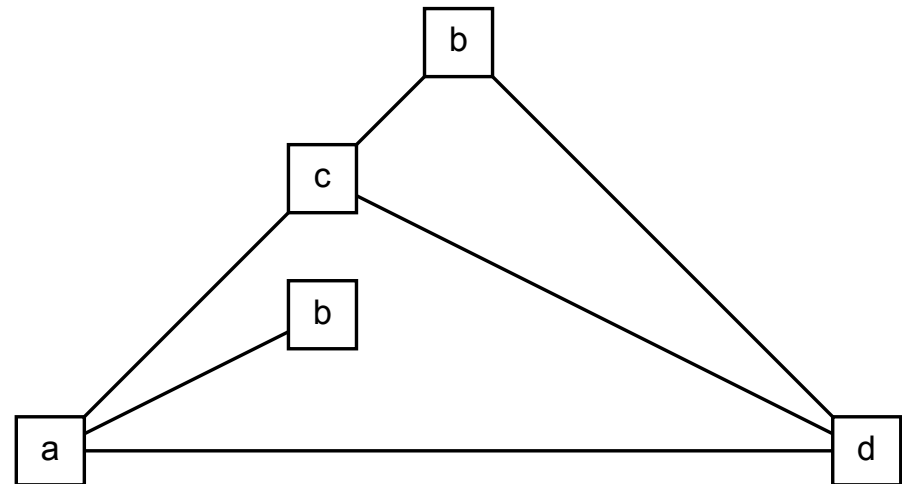
Non-cograph cases : **822**

Number of problemclasses : **9**

\hat{R} :



problemclass #3
24 generating gene trees



problemclass #7
234 generating gene trees

Question 2 : Can one infer information about the skeleton?

Theorem:

Given the leaf distance d_1 of the skeleton T_1 and d_2 of T_2 then:

$$d_1 = d_2 \Rightarrow \text{Unroot}(T_1) \sim \text{Unroot}(T_2)$$

$\text{Unroot}(T)$: unrooted version of the tree T