

# Clustering Systems of Phylogenetic Networks

Marc Hellmuth

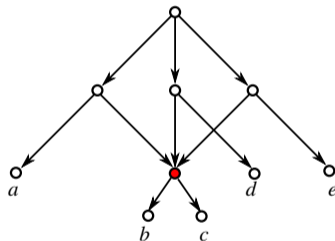
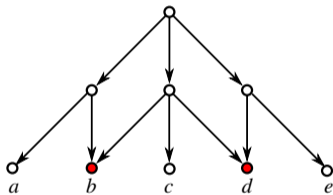
Department of Mathematics  
Faculty of Science  
Stockholm University

*TBI Winterseminar, 2024*

# Basics

All networks  $N = (V, E)$  considered here are

1. DAGs with a single root
2. phylogenetic (= no indegree 1 and outdegree 1 vertices)

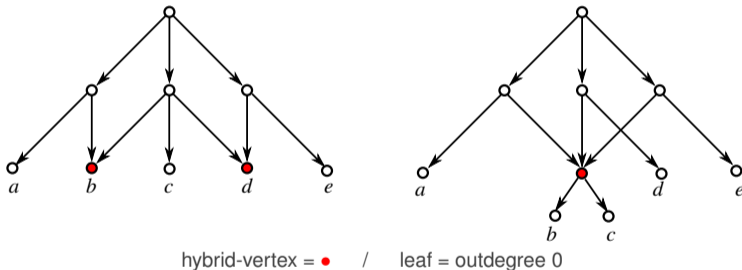


hybrid-vertex = ● / leaf = outdegree 0

# Basics

All networks  $N = (V, E)$  considered here are

1. DAGs with a single root
2. phylogenetic (= no indegree 1 and outdegree 1 vertices)



## (In)comparable vertices

$u \preceq_N v$  if  $v$  is an ancestor of  $u$ , i.e., there is a directed path from  $v$  to  $u$ .

$u$  and  $v$  are  **$\preceq_N$ -comparable** if  $u \preceq_N v$  or  $v \preceq_N u$

Otherwise,  $u$  and  $v$  are  **$\preceq_N$ -incomparable**

## Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

## Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

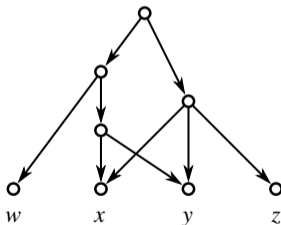
## Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .



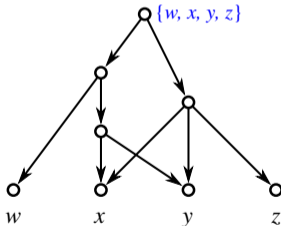
## Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .



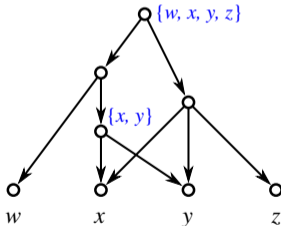
## Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .





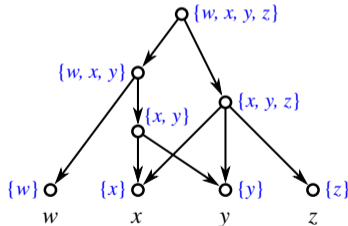
## Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .



$$\mathcal{C}_N = \left\{ \{w, x, y, z\}, \{w, x, y\}, \{x, y, z\}, \{x, y\}, \{w\}, \{x\}, \{y\}, \{z\} \right\}$$

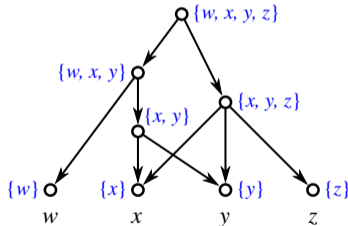
## Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .



$$\mathcal{C}_N = \left\{ \{w, x, y, z\}, \{w, x, y\}, \{x, y, z\}, \{x, y\}, \{w\}, \{x\}, \{y\}, \{z\} \right\}$$

How much information about  $N$  is contained in  $\mathcal{C}_N$  ?

# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

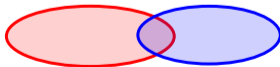
$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

**Folklore (Trees):**

- $C, C' \in \mathcal{C}$  **overlap**, if  $C \cap C' \notin \{C, C', \emptyset\}$ .

A clustering system is a **hierarchy** if it does not contain pairwise overlapping sets.



OVERLAP

# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

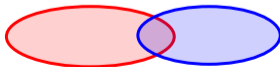
**Folklore (Trees):**

- $C, C' \in \mathcal{C}$  **overlap**, if  $C \cap C' \notin \{C, C', \emptyset\}$ .

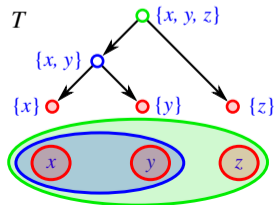
A clustering system is a **hierarchy** if it does not contain pairwise overlapping sets.

- **There is a 1-to-1 correspondence between rooted phylogenetic trees and hierarchies.**

[keyword: Hasse-diagram]



OVERLAP



# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

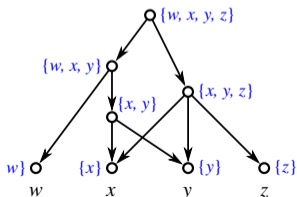
For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

## Observations about networks:

- $\mathcal{C}_N$  is usually *not* a hierarchy for general networks



Here:  $\mathcal{C} = \{\{w, x, y, z\}, \{w, x, y\}, \{x, y, z\}, \{y, z\}, \{w\}, \{x\}, \{y\}, \{z\}\}$

# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

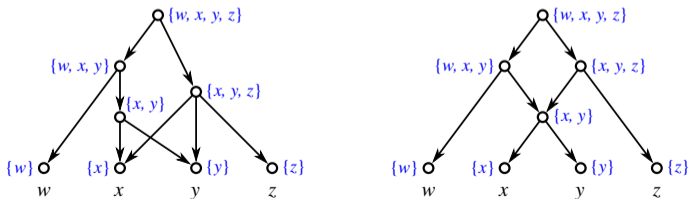
For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

## Observations about networks:

- $\mathcal{C}_N$  is usually *not* a hierarchy for general networks



Here:  $\mathcal{C} = \{\{w, x, y, z\}, \{w, x, y\}, \{x, y, z\}, \{y, z\}, \{w\}, \{x\}, \{y\}, \{z\}\}$

# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

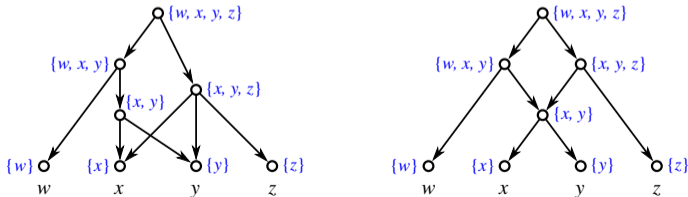
For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

## Observations about networks:

- $\mathcal{C}_N$  is usually *not* a hierarchy for general networks
- No 1-to-1 correspondence between networks and clustering systems.



Here:  $\mathcal{C} = \{\{w, x, y, z\}, \{w, x, y\}, \{x, y, z\}, \{y, z\}, \{w\}, \{x\}, \{y\}, \{z\}\}$

# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

**Central Questions:** Given  $\mathcal{C}$ .

- When is there a network  $N$  of a **given type** such that  $\mathcal{C} = \mathcal{C}_N$ ?  
type = tree, level-k, binary, galled tree, tree-child, tree-based, regular, normal, ...



# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

**Central Questions:** Given  $\mathcal{C}$ .

- When is there a network  $N$  of a **given type** such that  $\mathcal{C} = \mathcal{C}_N$ ?  
type = tree, level-k, binary, galled tree, tree-child, tree-based, regular, normal, ...
- Which type of networks are uniquely determined by  $\mathcal{C}_N$ ?

# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

**Central Questions:** Given  $\mathcal{C}$ .

- When is there a network  $N$  of a **given type** such that  $\mathcal{C} = \mathcal{C}_N$ ?  
type = tree, level-k, binary, galled tree, tree-child, tree-based, regular, normal, ...
- Which type of networks are uniquely determined by  $\mathcal{C}_N$ ?

# Clustering systems

A **clustering system** on  $X$  is a set  $\mathcal{C} \subseteq 2^X \setminus \{\emptyset\}$  such that  $X \in \mathcal{C}$  and  $\{x\} \in \mathcal{C}$  for all  $x \in X$ .

For all  $N = (V, E)$ , there is a *unique* clustering system

$$\mathcal{C}_N := \{C(v) \mid v \in V\}$$

where  $C(v)$  denotes the set of leaves “below”  $v$ .

**Central Questions:** Given  $\mathcal{C}$ .

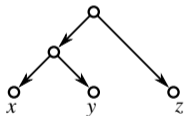
- When is there a network  $N$  of a **given type** such that  $\mathcal{C} = \mathcal{C}_N$ ?  
type = tree, level-k, binary, galled tree, tree-child, tree-based, regular, normal, ...
- Which type of networks are uniquely determined by  $\mathcal{C}_N$ ?

Let us try to answer some of the questions for level-1 networks.

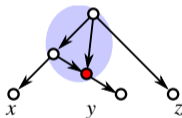
## Level- $k$ networks

A **block  $B$**  in a network is a maximal biconnected subgraph.

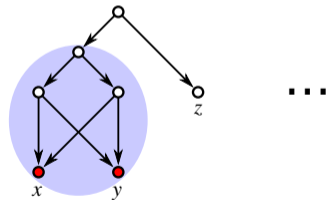
A network  $N$  is **level- $k$**  if each block  $B$  in  $N$  contains  $\leq k$  hybrid-vertices  
(distinct from root  $p_B$  of  $B$ ).



level-0 = tree



level-1

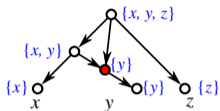
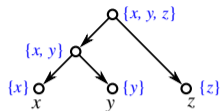


level-2

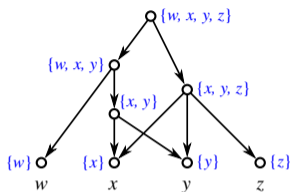
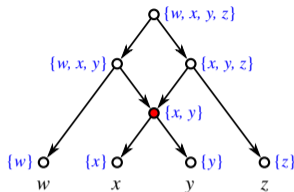
...

...

## Level-1 Networks



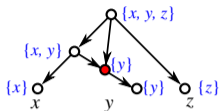
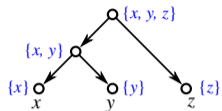
**Level-1**



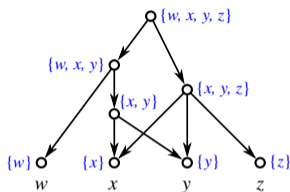
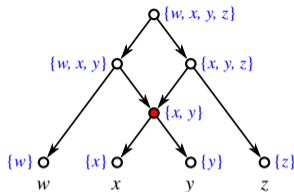
**not Level-1**

In trees  $T$ :  $u \preceq_T v \iff C(u) \subseteq C(v)$

## Level-1 Networks



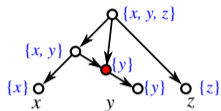
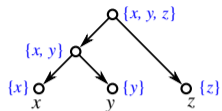
**Level-1**



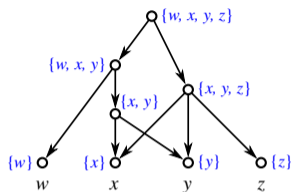
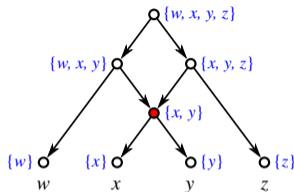
**not Level-1**

In trees  $T$ :  $u \preceq_T v \iff C(u) \subseteq C(v)$

## Level-1 Networks



**Level-1**

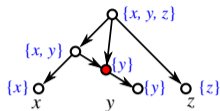
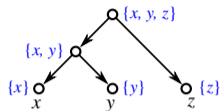


**not Level-1**

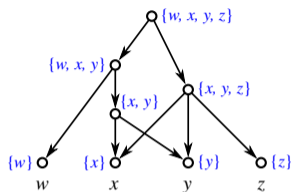
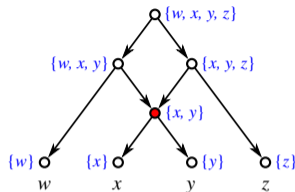
In trees  $T$ :  $u \preceq_T v \iff C(u) \subseteq C(v)$

In networks  $N$ :  $u \preceq_N v \implies C(u) \subseteq C(v)$  (converse not true in general)

## Level-1 Networks



**Level-1**



**not Level-1**

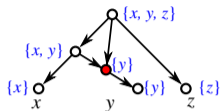
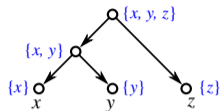
In trees  $T$ :  $u \preceq_T v \iff C(u) \subseteq C(v)$

In networks  $N$ :  $u \preceq_N v \implies C(u) \subseteq C(v)$  (converse not true in general)

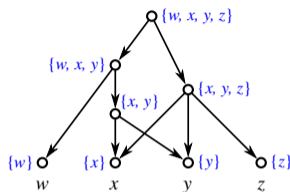
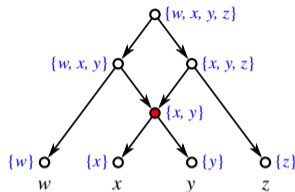
Hence, if  $C(u)$  and  $C(v)$  overlap or are disjoint, then  $u$  and  $v$  are  $\preceq_N$ -incomparable in  $N$



## Level-1 Networks



Level-1



not Level-1

In trees  $T$ :  $u \preceq_T v \iff C(u) \subseteq C(v)$

In networks  $N$ :  $u \preceq_N v \implies C(u) \subseteq C(v)$  (converse not true in general)

Hence, if  $C(u)$  and  $C(v)$  overlap or are disjoint, then  $u$  and  $v$  are  $\preceq_N$ -incomparable in  $N$

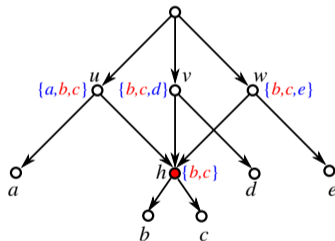
**Lemma** (H., Schaller, Stadler, 2023)

In level-1 networks  $N$ :  $u$  and  $v$  are  $\preceq_N$ -incomparable in  $N \iff C(u) \cap C(v) \in \{\emptyset, C(h_B)\}$  where  $h_B \neq u, v$  is the unique hybrid in block  $B$  that contains  $u$  and  $v$ .

## Property (L)

If  $C(u)$  overlaps with  $C(v)$ , then  $C(u) \cap C(v) = C(h_B)$  with  $B$  being the block containing  $u$  and  $v$ .

[this property cannot be observed when just looking at  $\mathcal{C}_N$ !]



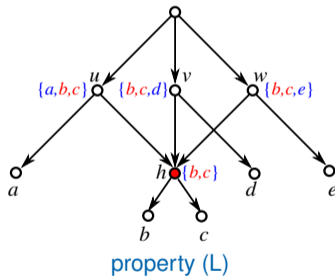
## Property (L)

If  $C(u)$  overlaps with  $C(v)$ , then  $C(u) \cap C(v) = C(h_B)$  with  $B$  being the block containing  $u$  and  $v$ .

[this property cannot be observed when just looking at  $\mathcal{C}_N$ ]

A clustering system  $\mathcal{C}$  satisfies **property (L)** if

$C \cap C_1 = C \cap C_2$  for all  $C, C_1, C_2 \in \mathcal{C}$  where  $C$  overlaps both  $C_1$  and  $C_2$ .



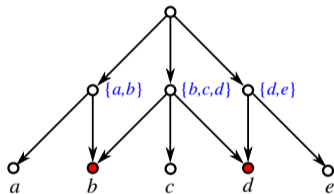
## Property (L)

If  $C(u)$  overlaps with  $C(v)$ , then  $C(u) \cap C(v) = C(h_B)$  with  $B$  being the block containing  $u$  and  $v$ .

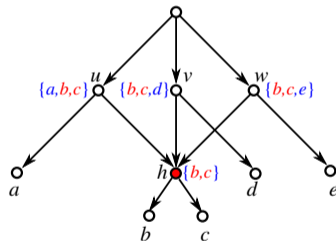
[this property cannot be observed when just looking at  $\mathcal{C}_N$ ]

A clustering system  $\mathcal{C}$  satisfies **property (L)** if

$C \cap C_1 = C \cap C_2$  for all  $C, C_1, C_2 \in \mathcal{C}$  where  $C$  overlaps both  $C_1$  and  $C_2$ .



no property (L)



property (L)

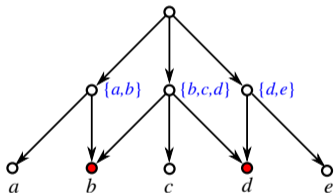
## Property (L)

If  $C(u)$  overlaps with  $C(v)$ , then  $C(u) \cap C(v) = C(h_B)$  with  $B$  being the block containing  $u$  and  $v$ .

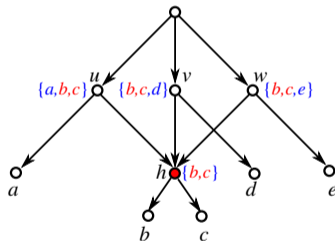
[this property cannot be observed when just looking at  $\mathcal{C}_N$ ]

A clustering system  $\mathcal{C}$  satisfies **property (L)** if

$C \cap C_1 = C \cap C_2$  for all  $C, C_1, C_2 \in \mathcal{C}$  where  $C$  overlaps both  $C_1$  and  $C_2$ .



no property (L)



property (L)

### Lemma (H., Schaller, Stadler, 2023)

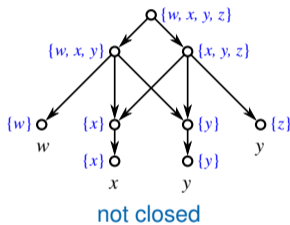
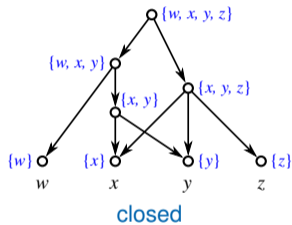
For all level-1 networks  $N$ , the set  $\mathcal{C}_N$  satisfies property (L).

## Closed clustering systems

A clustering system  $\mathcal{C}$  is **closed** if  $A \cap B \in \mathcal{C}$  for all  $A, B \in \mathcal{C}$  with  $A \cap B \neq \emptyset$ .

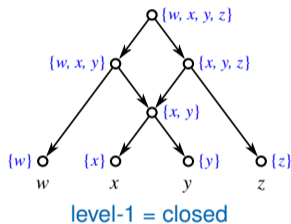
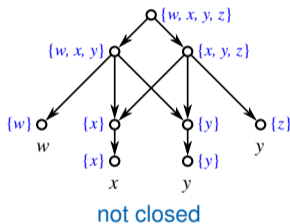
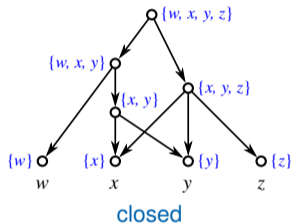
## Closed clustering systems

A clustering system  $\mathcal{C}$  is **closed** if  $A \cap B \in \mathcal{C}$  for all  $A, B \in \mathcal{C}$  with  $A \cap B \neq \emptyset$ .



# Closed clustering systems

A clustering system  $\mathcal{C}$  is **closed** if  $A \cap B \in \mathcal{C}$  for all  $A, B \in \mathcal{C}$  with  $A \cap B \neq \emptyset$ .



**Lemma** (H., Schaller, Stadler, 2023)

For level-1 networks  $N$ , the set  $\mathcal{C}_N$  is always closed.



## Characterization for level-1 Networks

A clustering system  $\mathcal{C}$  ...

... is **closed** if  $A \cap B \in \mathcal{C}$  for all  $A, B \in \mathcal{C}$  with  $A \cap B \neq \emptyset$ .

... satisfies **property (L)** if  $C \cap C_1 = C \cap C_2$  for all  $C, C_1, C_2 \in \mathcal{C}$  where  $C$  overlaps both  $C_1$  and  $C_2$ .

## Characterization for level-1 Networks

A clustering system  $\mathcal{C}$  ...

... is **closed** if  $A \cap B \in \mathcal{C}$  for all  $A, B \in \mathcal{C}$  with  $A \cap B \neq \emptyset$ .

... satisfies **property (L)** if  $C \cap C_1 = C \cap C_2$  for all  $C, C_1, C_2 \in \mathcal{C}$  where  $C$  overlaps both  $C_1$  and  $C_2$ .

### Theorem (H., Schaller, Stadler, 2023)

*For  $\mathcal{C}$ , there is a level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is closed and satisfies (L).*

## Characterization for level-1 Networks

A clustering system  $\mathcal{C}$  ...

... is **closed** if  $A \cap B \in \mathcal{C}$  for all  $A, B \in \mathcal{C}$  with  $A \cap B \neq \emptyset$ .

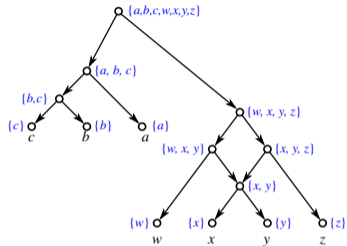
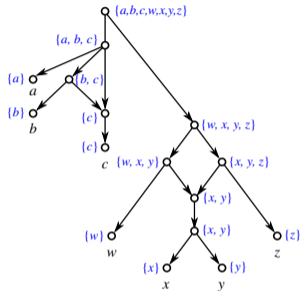
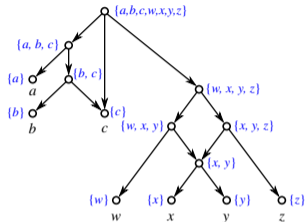
... satisfies **property (L)** if  $C \cap C_1 = C \cap C_2$  for all  $C, C_1, C_2 \in \mathcal{C}$  where  $C$  overlaps both  $C_1$  and  $C_2$ .

### Theorem (H., Schaller, Stadler, 2023)

*For  $\mathcal{C}$ , there is a level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is closed and satisfies (L).*

*Recognition of such  $\mathcal{C}$  and reconstruction of such a level-1 network can be done in polynomial time.*

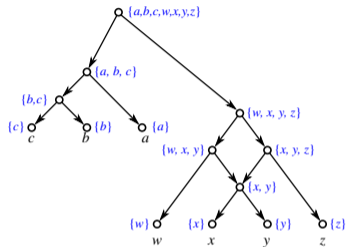
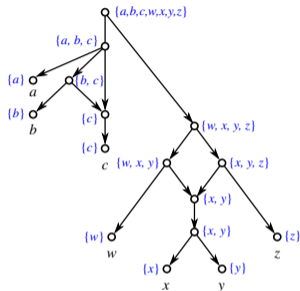
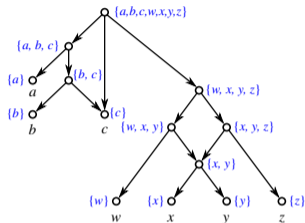
# Uniqueness Results



Hasse diagram

In general, several different level-1 networks may represent the same  $\mathcal{C}$

# Uniqueness Results



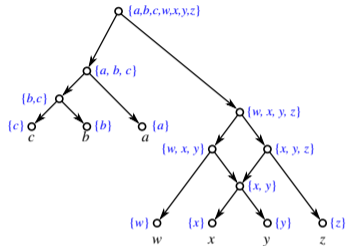
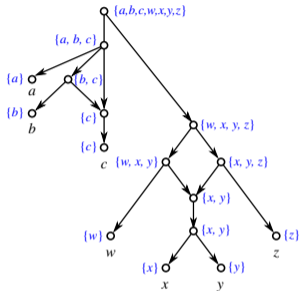
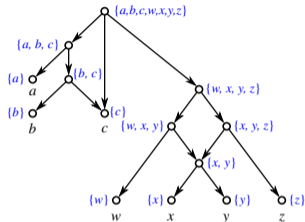
Hasse diagram

In general, several different level-1 networks may represent the same  $\mathcal{C}$

## Proposition (H., Schaller, Stadler, 2023)

Let  $\mathcal{C}$  be a closed clustering system that satisfies (L). Then, the Hasse-diagram of  $\mathcal{C}$  is the unique least-resolved level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C}$ .

# Uniqueness Results



Hasse diagram

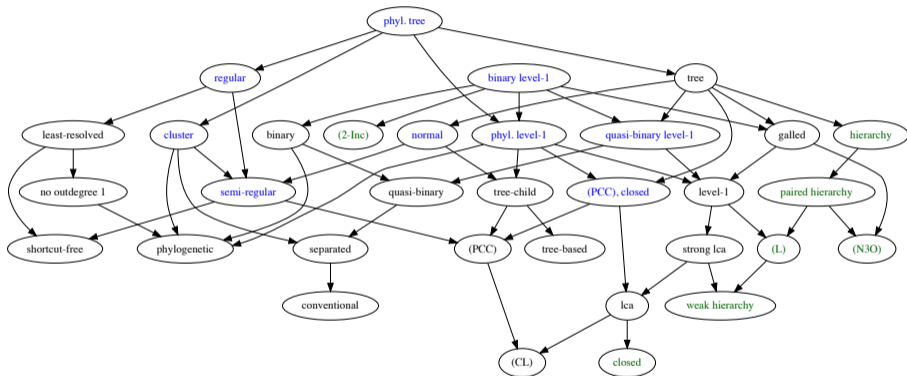
In general, several different level-1 networks may represent the same  $\mathcal{C}$

## Proposition (H., Schaller, Stadler, 2023)

Let  $\mathcal{C}$  be a closed clustering system that satisfies (L). Then, the Hasse-diagram of  $\mathcal{C}$  is the unique least-resolved level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C}$ .

Every level-1 network  $N$  is a refinement (=adding shortcuts + expand vertices) of the Hasse-diagram of  $\mathcal{C}_N$ .

## Other types of networks ...



The latter results and plenty of other characterizations for dozens of other network types can be found in *Hellmuth, Stadler, Schaller, Clustering Systems of Phylogenetic Networks, Theory in Biosciences (142), 301-358, 2023*

The initial motivation to investigate the clustering systems of networks:

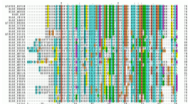
- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology



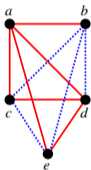
The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology

*genomic sequence data +  
similarity information*



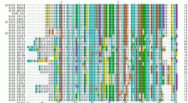
*pairwise relationships  
between genes*



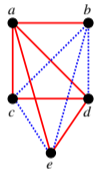
The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology

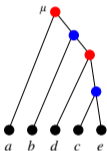
*genomic sequence data +  
similarity information*



*pairwise relationships  
between genes*



*phylogenetic tree  
representing relationships*

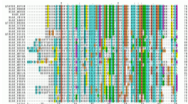


color of  $\text{lca}(x,y)$   
= relationship of  $x$  and  $y$

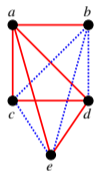
The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology

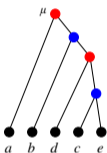
*genomic sequence data +  
similarity information*



*pairwise relationships  
between genes*



*phylogenetic tree  
representing relationships*



color of  $\text{lca}(x,y)$   
= relationship of  $x$  and  $y$

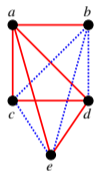
Given a relationship  $R$  that can be represented by trees  $T$

The initial motivation to investigate the clustering systems of networks:

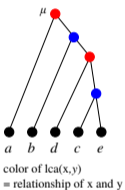
- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology



*pairwise relationships  
between genes*



*phylogenetic tree  
representing relationships*

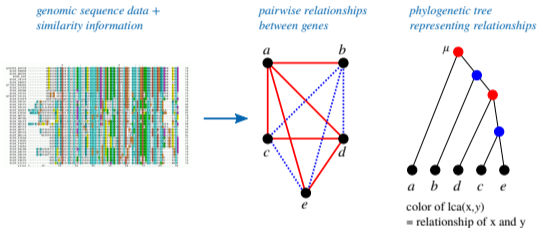


Given a relationship  $R$  that can be represented by trees  $T$

- Clustering system of  $T$  is determined by subgraphs (modules) of  $R$

The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology

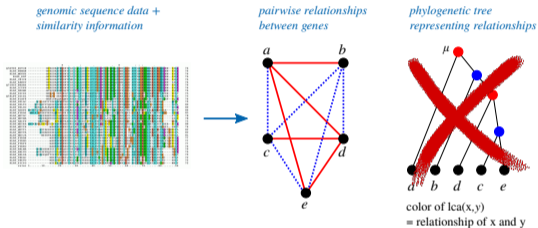


Given a relationship  $R$  that can be represented by trees  $T$

- Clustering system of  $T$  is determined by subgraphs (modules) of  $R$
- Colors of pairwise lca's determine relationship

The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology

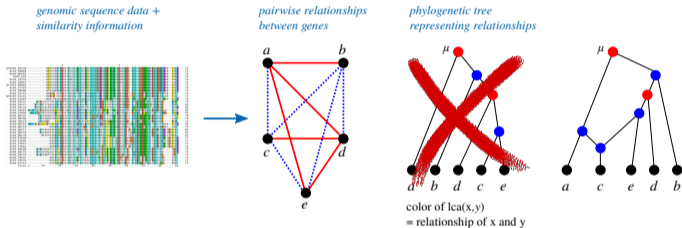


Given a relationship  $R$  that can be represented by trees  $T$

- Clustering system of  $T$  is determined by subgraphs (modules) of  $R$
- Colors of pairwise lca's determine relationship
- Noise in the data or NON-tree-like evolution  $\implies$  **cannot expect trees !**

The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology

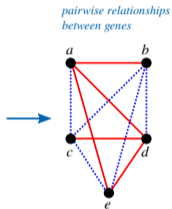


Given a relationship  $R$  that can be represented by trees  $T$

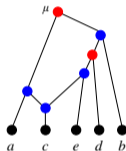
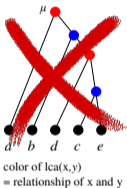
- Clustering system of  $T$  is determined by subgraphs (modules) of  $R$
- Colors of pairwise lca's determine relationship
- Noise in the data or NON-tree-like evolution  $\Rightarrow$  **cannot expect trees !**

The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology



phylogenetic tree  
representing relationships



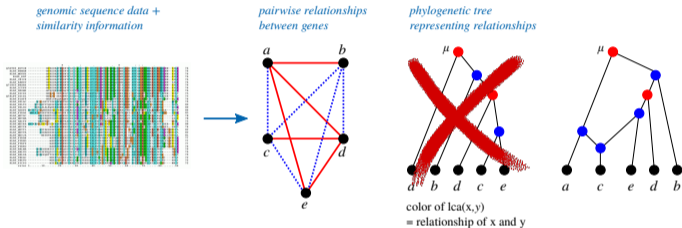
Given a relationship  $R$  that can be represented by trees  $T$

- Clustering system of  $T$  is determined by subgraphs (modules) of  $R$   
How are subgraphs in  $R$  related to clusters in  $N$  ?
- Colors of pairwise lca's determine relationship
- Noise in the data or NON-tree-like evolution  $\implies$  cannot expect trees !



The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology

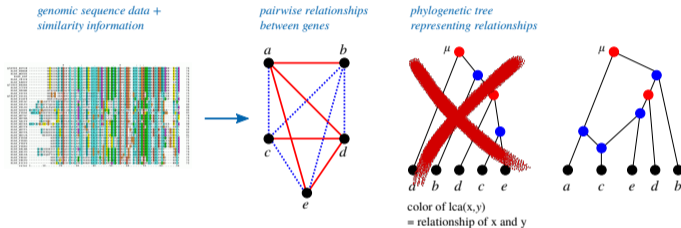


Given a relationship  $R$  that can be represented by trees  $T$

- Clustering system of  $T$  is determined by subgraphs (modules) of  $R$   
*How are subgraphs in  $R$  related to clusters in  $N$  ?*
- Colors of pairwise lca's determine relationship  
*How are networks characterized with pairwise-lca-properties?*
- Noise in the data or NON-tree-like evolution  $\implies$  **cannot expect trees !**

The initial motivation to investigate the clustering systems of networks:

- lack of results in literature
- we wanted to understand in more detail inference of horizontal gene transfer and orthology



Given a relationship  $R$  that can be represented by trees  $T$

- Clustering system of  $T$  is determined by subgraphs (modules) of  $R$   
How are subgraphs in  $R$  related to clusters in  $N$  ?
- Colors of pairwise  $\text{lca}$ 's determine relationship  
How are networks characterized with pairwise- $\text{lca}$ -properties?
- Noise in the data or NON-tree-like evolution  $\Rightarrow$  **cannot expect trees !**

A first starting point is provided by

## Theorem (Shanavas, Changat, H., Stadler, 2024)

For  $\mathcal{C}$ , there is a network  $N$  with  $\mathcal{C} = \mathcal{C}_N$  and pairwise  $\text{lca}$ -property

$\iff \mathcal{C}_N$  is **pre-binary**, i.e., there is a unique **incl.-min. cluster**  $C \in \mathcal{C}$  such that  $\{x,y\} \subseteq C$  for all  $x,y \in L(N)$

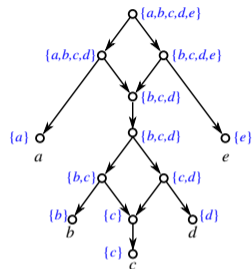
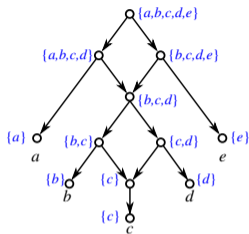
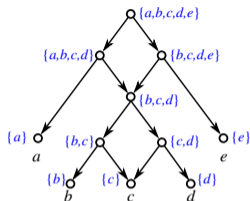
Shanavas, Changat, Hellmuth, Stadler **Unique Least Common Ancestors and Clusters in Directed Acyclic Graphs**, LNCS vol 14508, 2024

- David Schaller (Biontech, GER)
- Peter F. Stadler (Uni Leipzig, GER)
- Ameera Vaheeda Shanavas (Uni Kerala, IND)
- Manoj Changa (Uni Kerala, IND)
- Anna Lindeberg (Uni Stockholm, SWE)

- David Schaller (Biontech, GER)
- Peter F. Stadler (Uni Leipzig, GER)
- Ameera Vaheeda Shanavas (Uni Kerala, IND)
- Manoj Changa (Uni Kerala, IND)
- Anna Lindeberg (Uni Stockholm, SWE)

**Thanks!**

# Uniqueness Results



We obtain uniqueness for mild restrictions!

## Theorem

Let  $\mathcal{C}$  be a closed clustering system that satisfies (L). Then, there is a unique shortcut-free level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C}$  that satisfies precisely one condition:

- $N$  contains no vertex  $v$  with outdegree 1 (=Hasse diagram)
- every leaf in  $N$  has indegree 1 and all vertices  $v$  with outdegree 1 are adjacent to leaves.
- every hybrid in  $N$  has outdegree 1 (thus all leaves have indegree 1).

## Further results

**Galled tree** = level-1 network where all non-trivial blocks correspond to “undirected cycles”

## Further results

**Galled tree** = level-1 network where all non-trivial blocks correspond to “undirected cycles”

### Theorem

*There is a galled tree  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and does not contain three distinct pairwise overlapping clusters.*

## Further results

**Galled tree** = level-1 network where all non-trivial blocks correspond to “undirected cycles”

### Theorem

*There is a galled tree  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and does not contain three distinct pairwise overlapping clusters.*

**Binary network** = network where all non-hybrids  $v$  is either a leaf or has  $\text{outdeg}(v) = 2$ , and every hybrid  $v$  satisfies  $\text{indeg}(v) = 2$  and  $\text{outdeg}(v) = 1$ .



## Further results

**Galled tree** = level-1 network where all non-trivial blocks correspond to “undirected cycles”

### Theorem

*There is a galled tree  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and does not contain three distinct pairwise overlapping clusters.*

**Binary network** = network where all non-hybrids  $v$  is either a leaf or has  $\text{outdeg}(v) = 2$ , and every hybrid  $v$  satisfies  $\text{indeg}(v) = 2$  and  $\text{outdeg}(v) = 1$ .

### Theorem

*There is a binary level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and for all clusters  $C \in \mathcal{C}$ , there are at most two inclusion-maximal clusters  $A, B \in \mathcal{C}$  with  $A, B \subsetneq C$  and at most two inclusion-minimal clusters  $A, B \in \mathcal{C}$  with  $C \subsetneq A, B$ .*

## Further results

**Galled tree** = level-1 network where all non-trivial blocks correspond to “undirected cycles”

### Theorem

*There is a galled tree  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and does not contain three distinct pairwise overlapping clusters.*

**Binary network** = network where all non-hybrids  $v$  is either a leaf or has  $\text{outdeg}(v) = 2$ , and every hybrid  $v$  satisfies  $\text{indeg}(v) = 2$  and  $\text{outdeg}(v) = 1$ .

### Theorem

*There is a binary level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and for all clusters  $C \in \mathcal{C}$ , there are at most two inclusion-maximal clusters  $A, B \in \mathcal{C}$  with  $A, B \subsetneq C$  and at most two inclusion-minimal clusters  $A, B \in \mathcal{C}$  with  $C \subsetneq A, B$ .*

A clustering system  $\mathcal{C}$  is **compatible** w.r.t.  $N$  if  $\mathcal{C} \subseteq \mathcal{C}_N$ .

## Further results

**Galled tree** = level-1 network where all non-trivial blocks correspond to “undirected cycles”

### Theorem

*There is a galled tree  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and does not contain three distinct pairwise overlapping clusters.*

**Binary network** = network where all non-hybrids  $v$  is either a leaf or has  $\text{outdeg}(v) = 2$ , and every hybrid  $v$  satisfies  $\text{indeg}(v) = 2$  and  $\text{outdeg}(v) = 1$ .

### Theorem

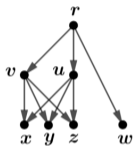
*There is a binary level-1 network  $N$  with  $\mathcal{C}_N = \mathcal{C} \iff \mathcal{C}$  is **closed** and satisfies **Property (L)** and for all clusters  $C \in \mathcal{C}$ , there are at most two inclusion-maximal clusters  $A, B \in \mathcal{C}$  with  $A, B \subsetneq C$  and at most two inclusion-minimal clusters  $A, B \in \mathcal{C}$  with  $C \subsetneq A, B$ .*

A clustering system  $\mathcal{C}$  is **compatible** w.r.t.  $N$  if  $\mathcal{C} \subseteq \mathcal{C}_N$ .

### Theorem

*$\mathcal{C}$  is compatible w.r.t. a level-1 network  $N \iff \mathcal{C}$  is satisfies **Property (L)**.*

*In this case, compute  $A \cap B$  for all overlapping  $A, B \in \mathcal{C}$  and add it to  $\mathcal{C}$  if their intersection is not present. (can be done in polynomial time)*



**Fig. 2.** Consider the DAG  $G$  with leaf set  $X = L(G)$  where  $\mathcal{C}_G = \{\{x\}, \{y\}, \{z\}, \{w\}, \{x, y, z\}, X\}$ . Here,  $\mathcal{C}_G$  satisfies **(KS)** and **(KC)** for  $k = 2$ . By definition,  $\mathcal{C}_G$  is thus pre-binary. However,  $G$  is not a pairwise lca-network since  $\text{lca}(x, y)$ ,  $\text{lca}(x, z)$ , and  $\text{lca}(y, z)$  are not defined. Moreover,  $\mathcal{C}_G$  also satisfies **(KC)** for  $k = 3$  but  $G$  is not a 3-lca-network since  $\text{lca}(x, y, z)$  is not defined.

pre-binary clustering system but not pairwise lca-property