# On the Elusiveness of Clusters

Steven M. Kelk   Celine Scornavacca   Leo van Iersel

Lukas Friedlos 22.04.21
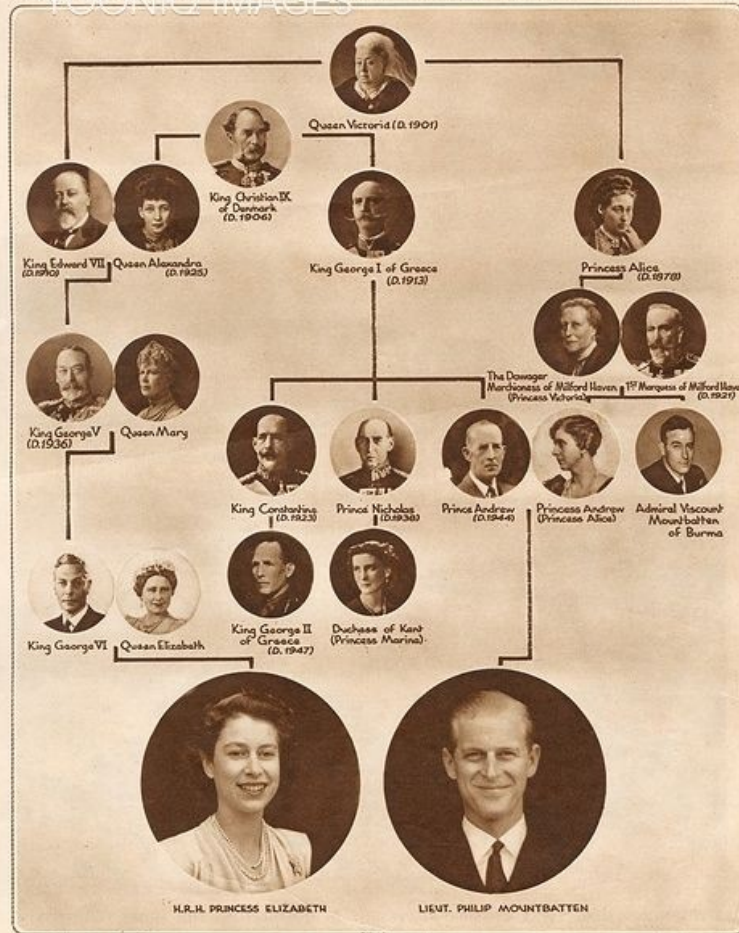
GammaCoV

KC243390.1_BatCoV_/8-724/Pip_pyg/ROU/2009
KC869678.4_BatCoV_Neoromicia/PML-PHE1/RSA/2011
KJ713299.1_MERS-CoV_KSA-CAMEL-376
JX869059.2_MERS-CoV_Human
NC_022643.1_Betacoronavirus_Erinaceus/VMC/DEU/2012
EF065505.1_BatCoV_HKU4-1
EF065509.1_BatCoV_HKU5-1
AY585228.1_HumanCoV_OC43
DQ011855.1_PorcineHemagglutinatingEncephalomyelitis
EF446615.1_Equine_CoV_strainNC99
KU558922.1_Betacoronavirus_1_BuffaloCoV_B1-24F
FJ425188.1_Sambar_deer_CoV_US/OH-WD388-TC/1994
FJ938064.1_Bovine_CoV_E-AH187-TC
FJ938065.1_Bovine_respiratory_CoV_AH187
FJ938068.1_Rat_coronavirus
AC_000192.1_Murine_hepatitis
KJ020608.1_BatCoV/B55740/S.kuh/CB/Tha/6/2012
KJ020604.1_BatCoV_B55700-3/Cyn_sph/CB/Tha/5/2012
HQ728482.1_Eidolon_BatCoV_Kenya/KY24/2006
AB918719.1_BatCoV_FB2012-8F
AB918717.1_BatCoV_IFB2012-13F
AB918718.1_BatCoV_IFB2012-17F
HM211099.1_BatCoV_HKU9-5-2
GU065421.1_Kenya_BatCoV_BtKY77
HQ728483.1_Rousettus_BatCoV_Kenya/KY06/2006
EF065513.1_BatCoV_HKU9-1
KX520654.1_BatCoV_strainRK059
AB539081.1_BatCoV_Philippines/Diliman1525G2/2008
NC_014470.1_BatCoV_BM48-31/BGR/2008
WH-Human_1IChinaI2019-Dec
DQ071615.1_Bat_SARSr-CoV_Rp3
AY278741.1_SARS-CoV_Urbani
AY304486.1_SARS-CoV_SZ3
JX993988.1_BatCoV_Bat_Cp/Yunnan2011
DQ022305.2_Bat_SARSr-CoV_HKU3-1
DQ412042.1_Bat_SARSr-CoV_Rf1
HQ166910.1_Zaria_BatCoV
FJ710043.1_BatCoV_Hipposideros/GhanaBoo/348/2008
EU769557.1_BatCoV_Trinidad/1FY2BA/2007
HQ728484.1_Miniopterus_BatCoV_Kenya/KY27/2006
BatCoV_Miniopterus_2_1B
BatCoV_Miniopterus_1A
DQ666339.1_BatCoV_HKU7_strainWCF88
EU420139.1_BatCoV_HKU8_strainAFCD77
HQ728481.1_Chaerephon_BatCoV_Kenya/KY41/2006
AY567487.2_HumanCoV_NL63
AF304460.1_HumanCoV_229E
EF203065.1_BatCoV_HKU2_strainHK/46/2006
HQ728480.1_Cardioderma_BatCoV_Kenya/KY43/2006
JQ989273.1_Hipposideros_BatCoV_HKU10
NC_022103.1_BatCoV_CDPHE15/USA/2006
AF353511.1_PEDV_strainCV777
GU190216.1_BatCoV_NM98-62/GER/2008
DQ648858.1_BatCoV_BtCoV/512/2005
AY994055.1_FelineInfectiousPeritonitis
DQ811787.1_PRCV_ISU-1

DeltaCoV

0.3

**Preliminary maximum likelihood phylogenetic analysis of novel Wuhan, China human CoV** in red, GenBank (accession MN908947)
Novel CoV seq data from: http://virological.org/t/initial-genome-release-of-novel-coronavirus/319. The Shanghai Public Health Clinical Center & School of Public Health, in collaboration with the Central Hospital of Wuhan, Huazhong University of Science and Technology, the Wuhan Center for Disease Control and Prevention, the National Institute for Communicable Disease Control and Prevention, Chinese Center for Disease Control, and the University of Sydney, Sydney, Australia.

PhyML tree based on partial RdRp gene sequence (410bp), aligned with representative human and animal CoV sequences from Genbank compiled by Alice Latinne; tree by Kevin Olival.
Analysis by EcoHealth Alliance - 11 Jan 2020 (12:30pm EST)

EcoHealth Alliance

# Cross-Hybridisation



- Hybrid speciation (through sexual reproduction)
  - Homo Sapiens ↔ Neanderthal
  - Polar Bear ↔ Brown Bear
  - Most common in plants

- Horizontal gene transfer

Horizontal gene transfer in bacteria
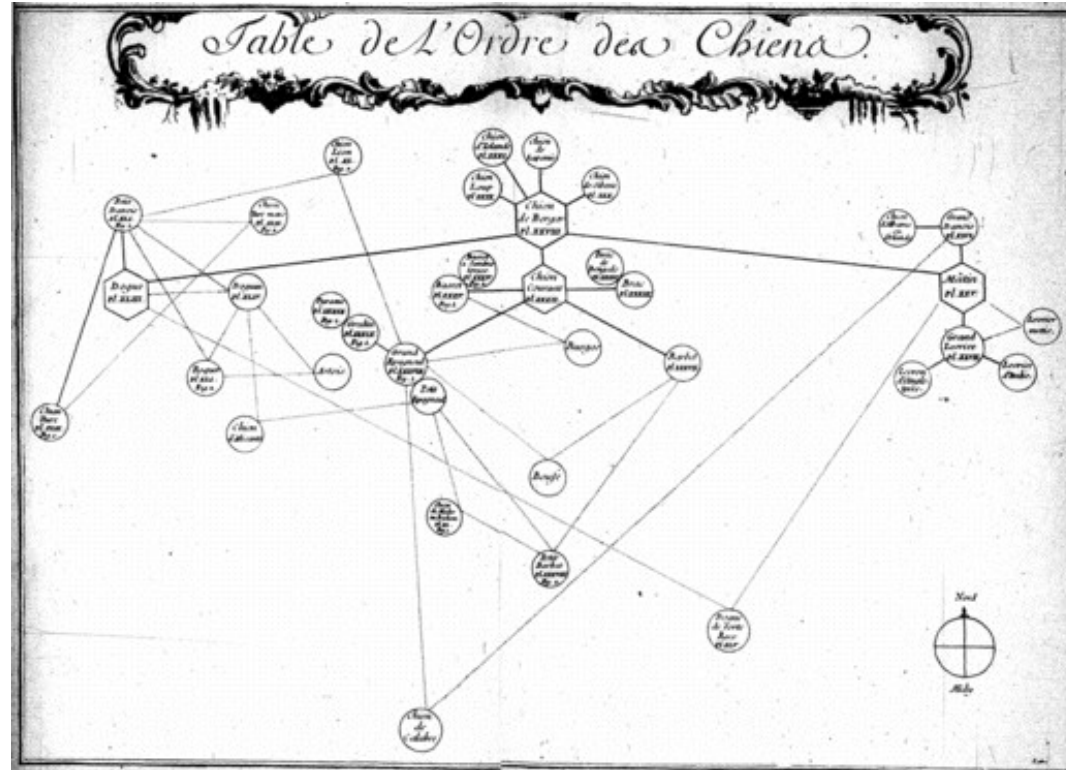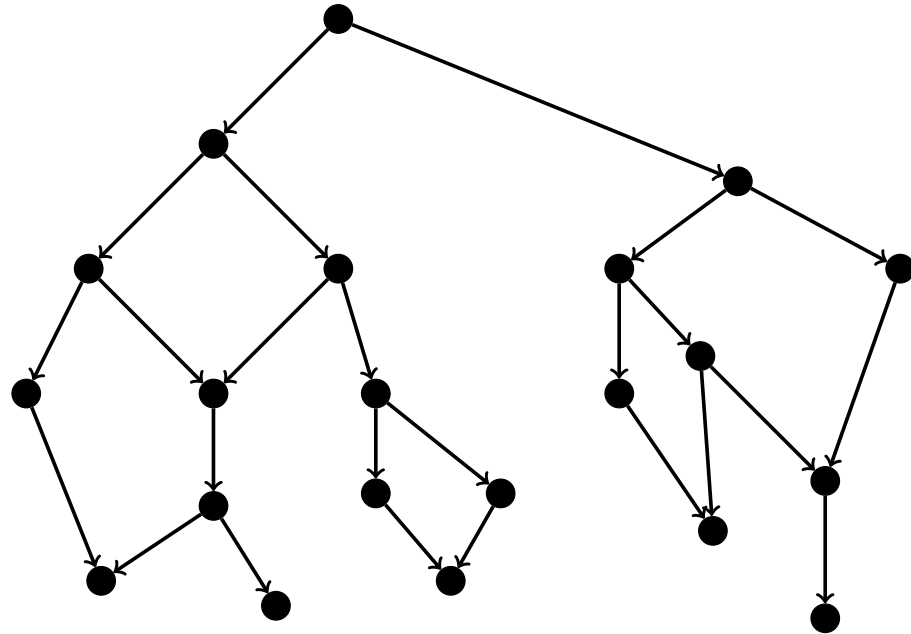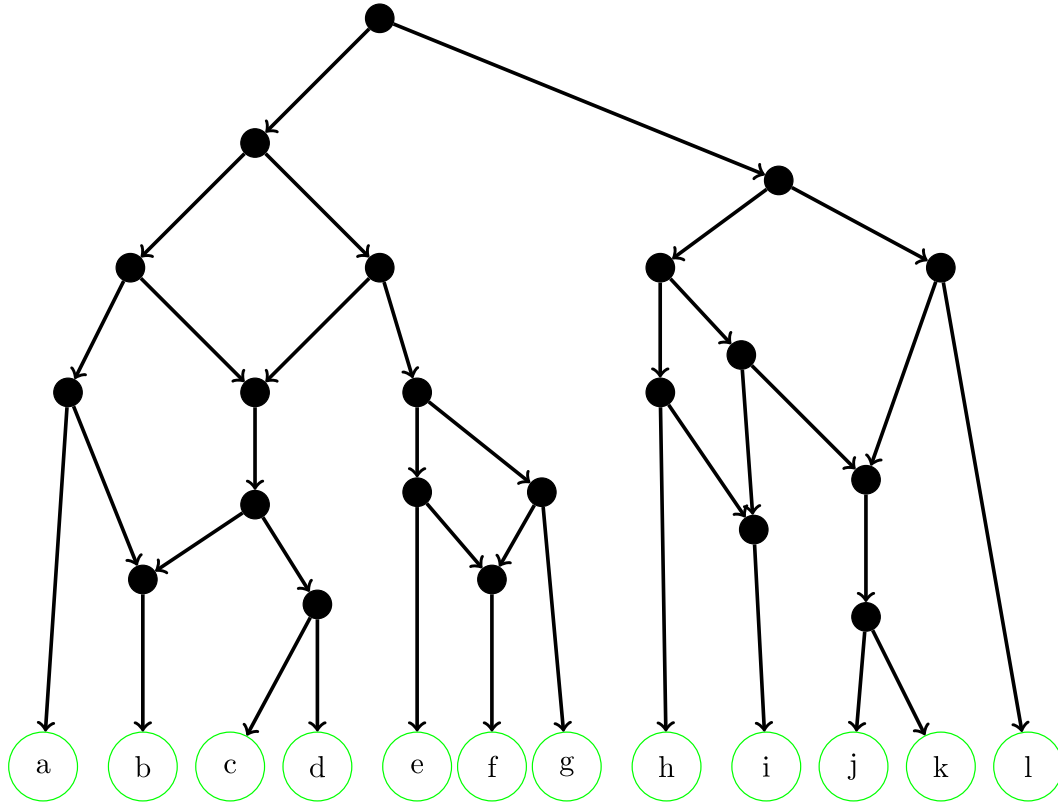


Heliconius heurippa

Parental 1    H. heurippa    Parental 2

Jaguma          Prizzly

# Phylogenetic Networks

# Phylogenetic Networks



Definition
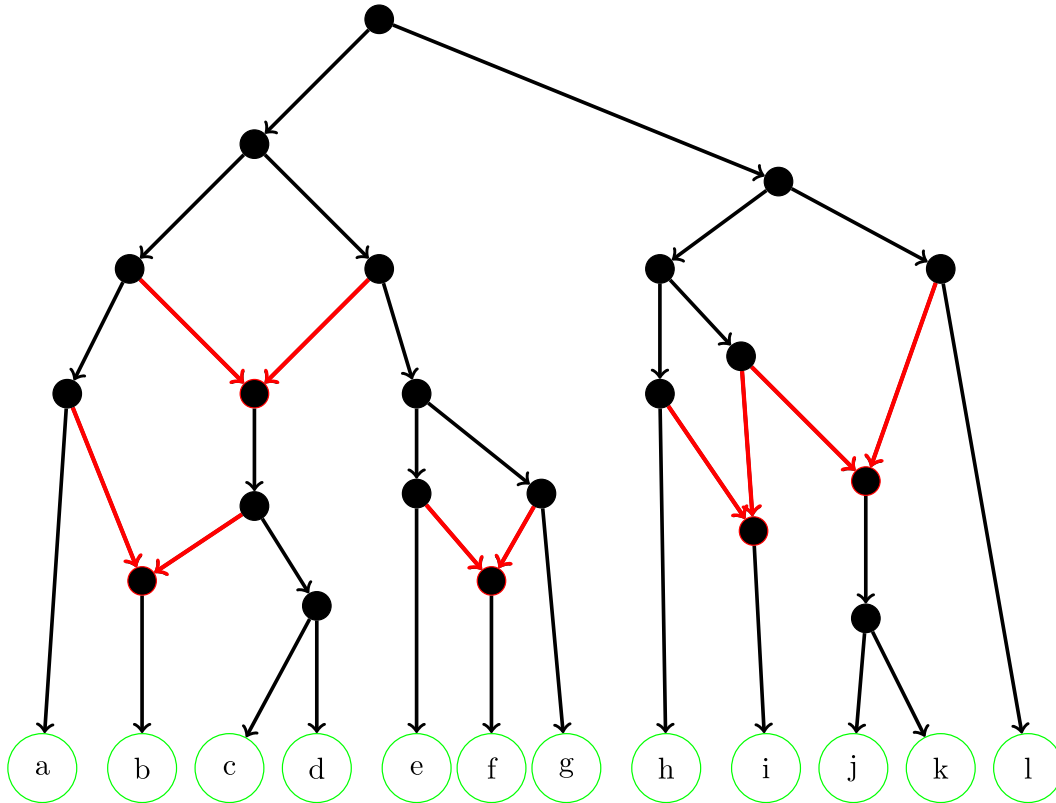
- rooted directed acyclic graph

  Henceforth network $N$

# Phylogenetic Networks
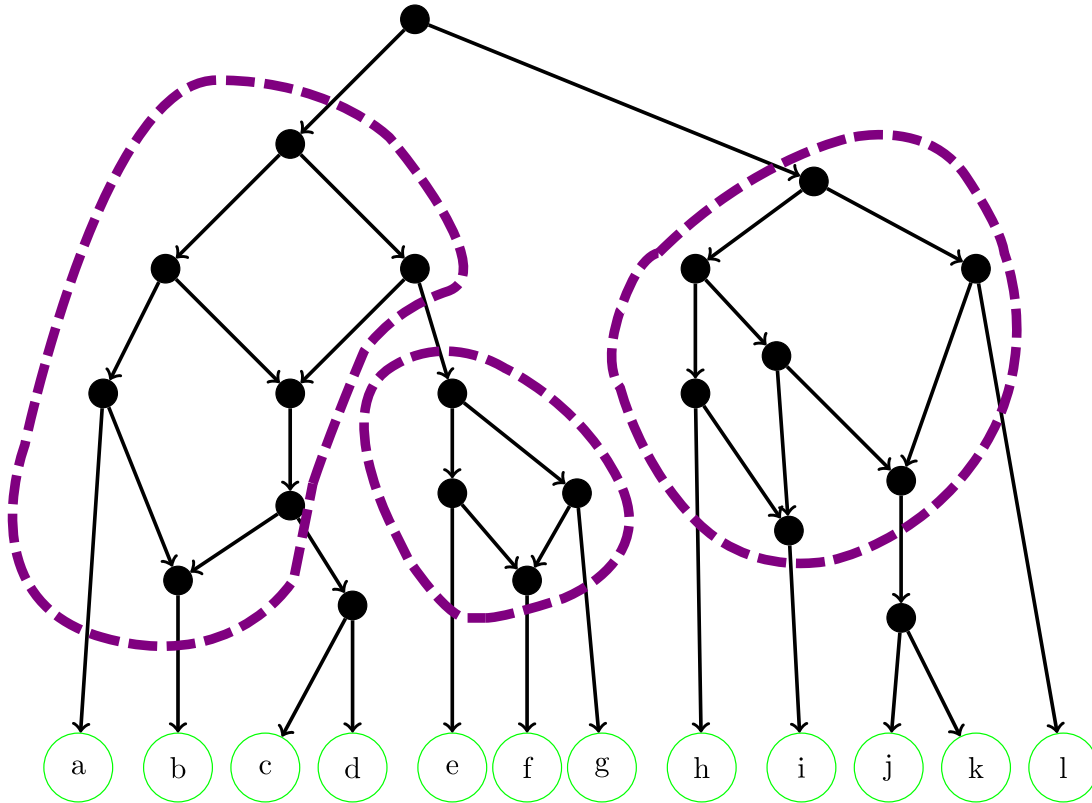


Definition

- rooted directed acyclic graph
  Henceforth network $N$

- leaves representing a set of taxa
  $\mathcal{X} = \{a, b, c, d, e, f, g, h, i, j, k, l\}$

# Phylogenetic Networks



Definition

- rooted directed acyclic graph
  Henceforth network $N$

- leaves representing a set of taxa
  $$\mathcal{X} = \{a, b, c, d, e, f, g, h, i, j, k, l\}$$

Terminology

- reticulation
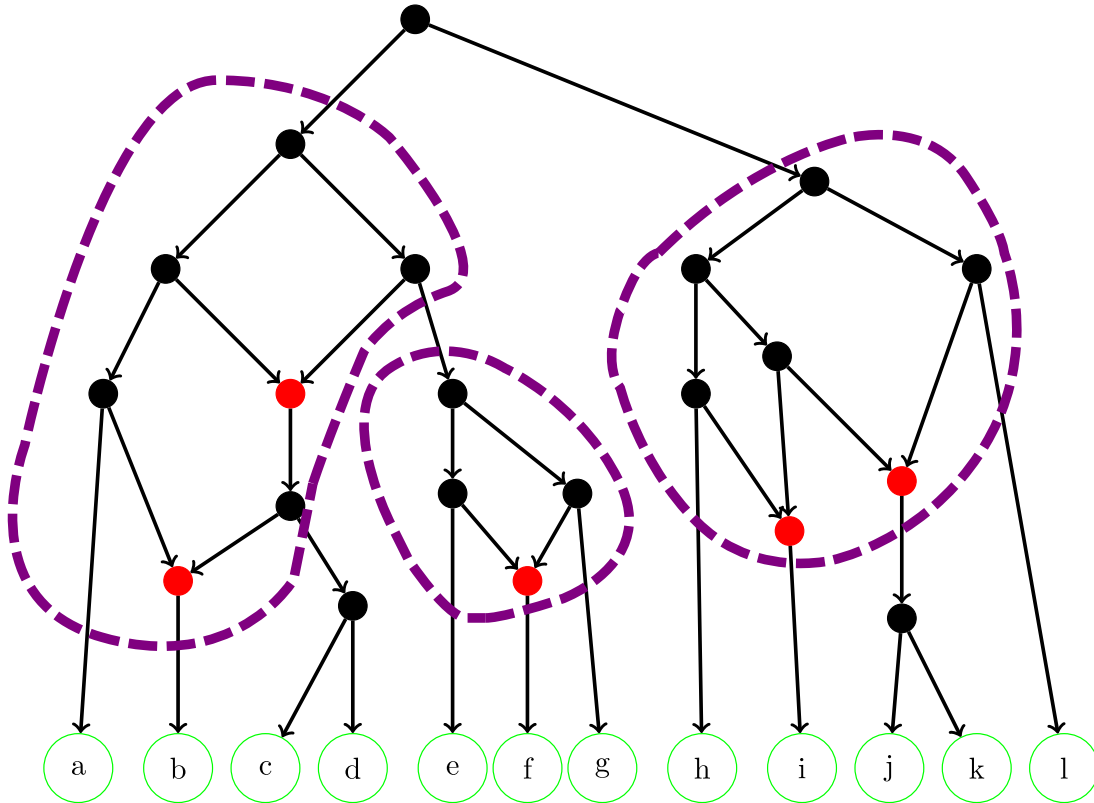  $$\text{in-deg}(v) \geq 2$$

# Phylogenetic Networks



Definition

- rooted directed acyclic graph
  Henceforth network $N$

- leaves representing a set of taxa
  $$\mathcal{X} = \{a, b, c, d, e, f, g, h, i, j, k, l\}$$

Terminology

- reticulation
  $$\text{in-deg}(v) \geq 2$$

- biconnected components

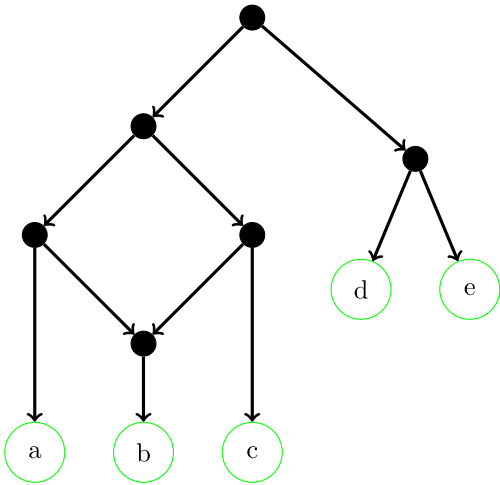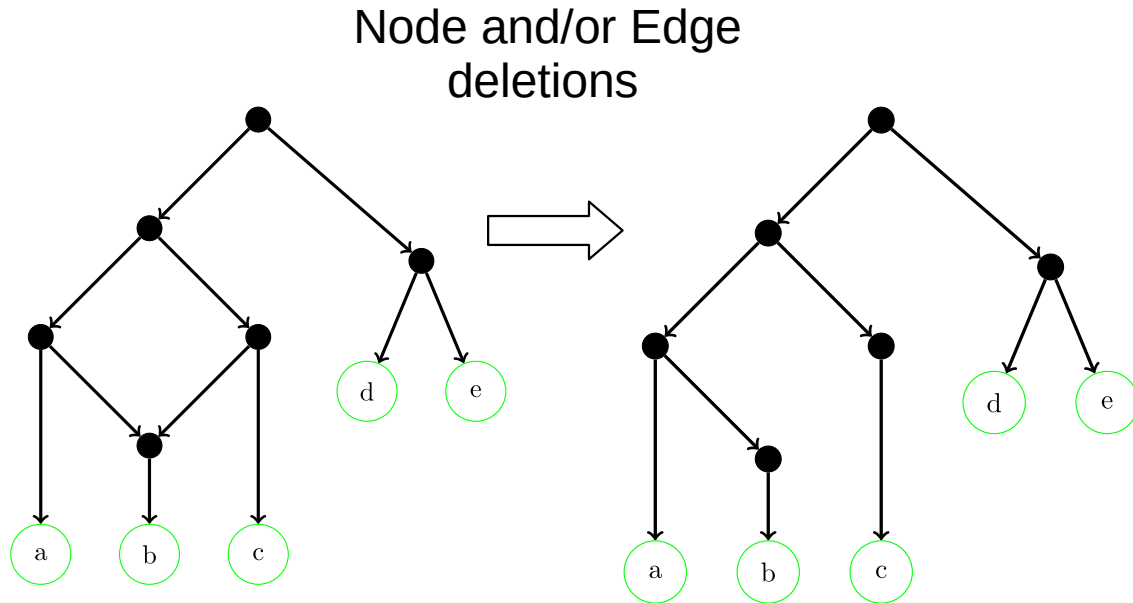# Phylogenetic Networks



Definition

- rooted directed acyclic graph
  Henceforth network $N$

- leaves representing a set of taxa
  $\mathcal{X} = \{a, b, c, d, e, f, g, h, i, j, k, l\}$

Terminology

- reticulation
  $$\text{in-deg}(v) \geq 2$$

- biconnected components

- $k$ -level

# Network 'Displaying' a Tree

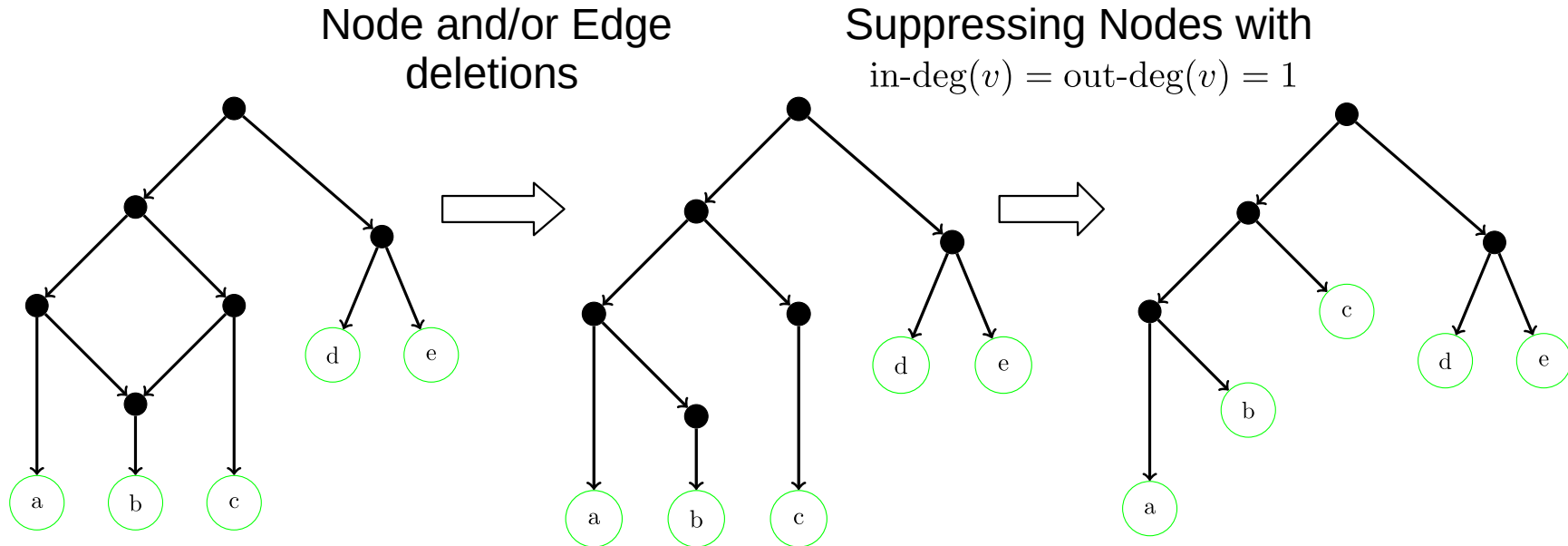A Network $N$ *displays* a tree $T$, if it can be obtained from $N$ via

# Network 'Displaying' a Tree

A Network $N$ *displays* a tree $T$, if it can be obtained from $N$ via

Node and/or Edge
deletions

# Network 'Displaying' a Tree

A Network $N$ *displays* a tree $T$, if it can be obtained from $N$ via

Node and/or Edge deletions

Suppressing Nodes with
$$\text{in-deg}(v) = \text{out-deg}(v) = 1$$

# Clusters

Given a set of taxa $\mathcal{X} = \{a, b, c, d, e\}$

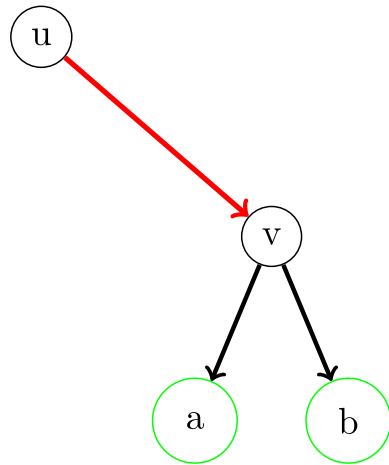A cluster $C \subset \mathcal{X}$ is a proper subset of all taxa, e.g. $\{a, b, c\}$

Taxa are in a cluster $\implies$ Taxa have a common ancestor

A set of clusters $\mathcal{C}$ on a set of taxa, e.g. $\left\{ \{a, b, c\}, \{a, b\}, \{d, e\} \right\}$
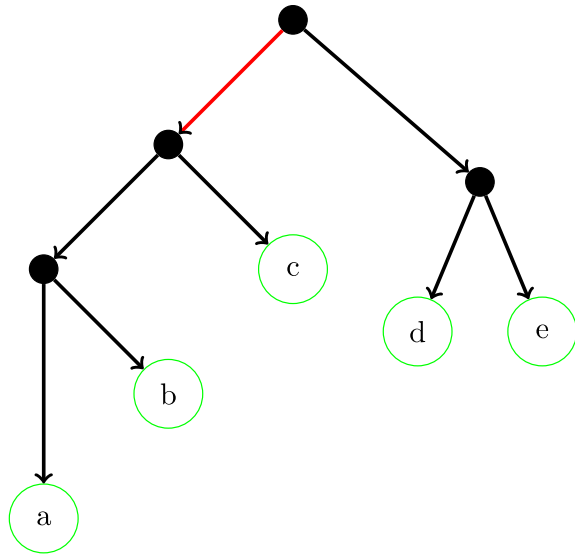
**Networks represent Clusters**

# Representing Clusters

We say that an **edge** $(u, v)$ *represents* a cluster $C$
if $C$ is the set of leaf descendants of $v$.



$$C = \{a, b\}$$
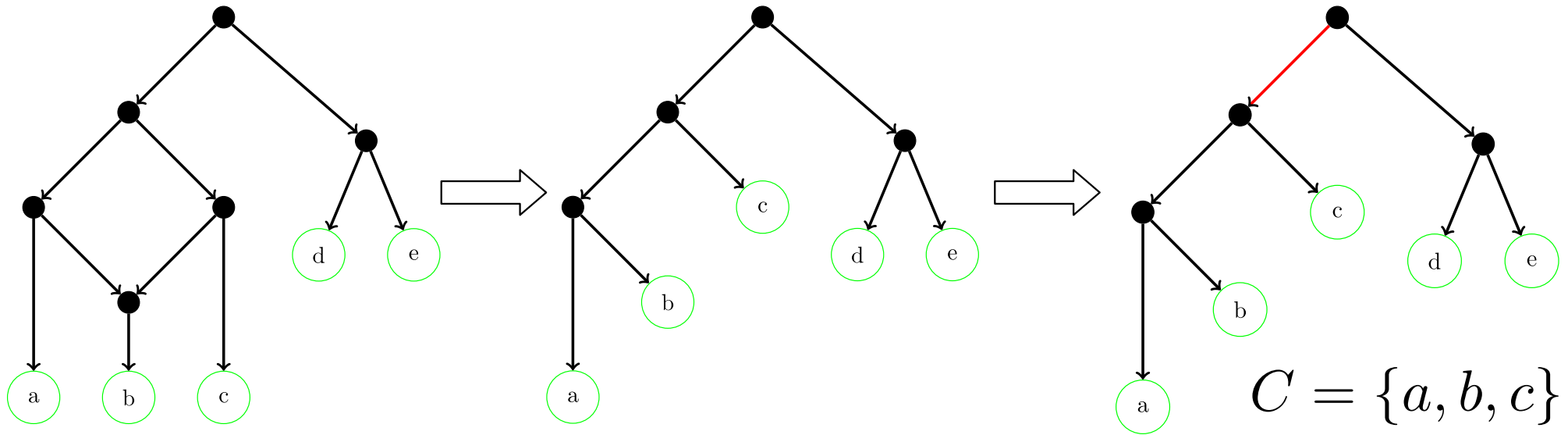
# Representing Clusters

We say that a tree $T$ *represents* a cluster $C$
if $T$ has an edge that represents $C$.



$$C = \{a, b, c\}$$

# Representing Clusters

We say that a network $N$ *represents* a cluster $C$ if $N$ displays a tree that represents $C$.
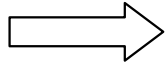
$$C = \{a, b, c\}$$

# A Theoretical Polynomial-Time Algorithm for Constructing Level-k Networks

# Goal

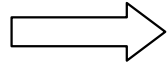Given: • set of clusters $\mathcal{C}$
       • fixed $k \geq 0$

# Goal

Given:
- set of clusters $\mathcal{C}$
- fixed $k \geq 0$

$\Longrightarrow$

Construct level-$k$ network $N$ representing $\mathcal{C}$, if such a network exists.

# Goal

Given:
- set of clusters $\mathcal{C}$
- fixed $k \geq 0$

$\Longrightarrow$

Construct level-$k$ network $N$ representing $\mathcal{C}$, if such a network exists.
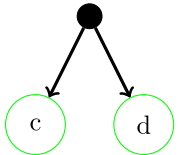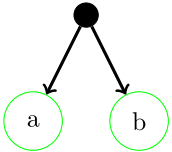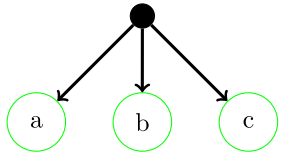
**Possible in Polynomial-Time!**

# From Clusters to Networks (Trivial)

E.g. $\mathcal{X} = \{a, b, c, d\}$, $\mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}$

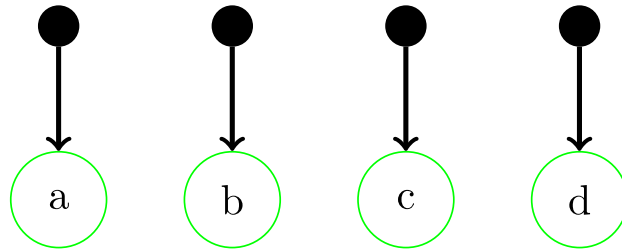# From Clusters to Networks (Trivial)

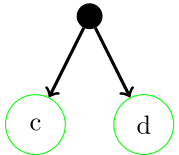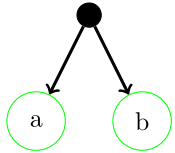E.g. $\mathcal{X} = \{a, b, c, d\}$, $\mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}$

# From Clusters to Networks (Trivial)

E.g. $\mathcal{X} = \{a, b, c, d\}$, $\mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}$

# From Clusters to Networks (Trivial)

E.g. $\mathcal{X} = \{a, b, c, d\}$, $\mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}$

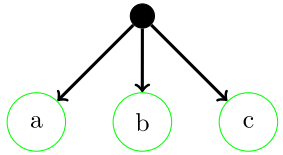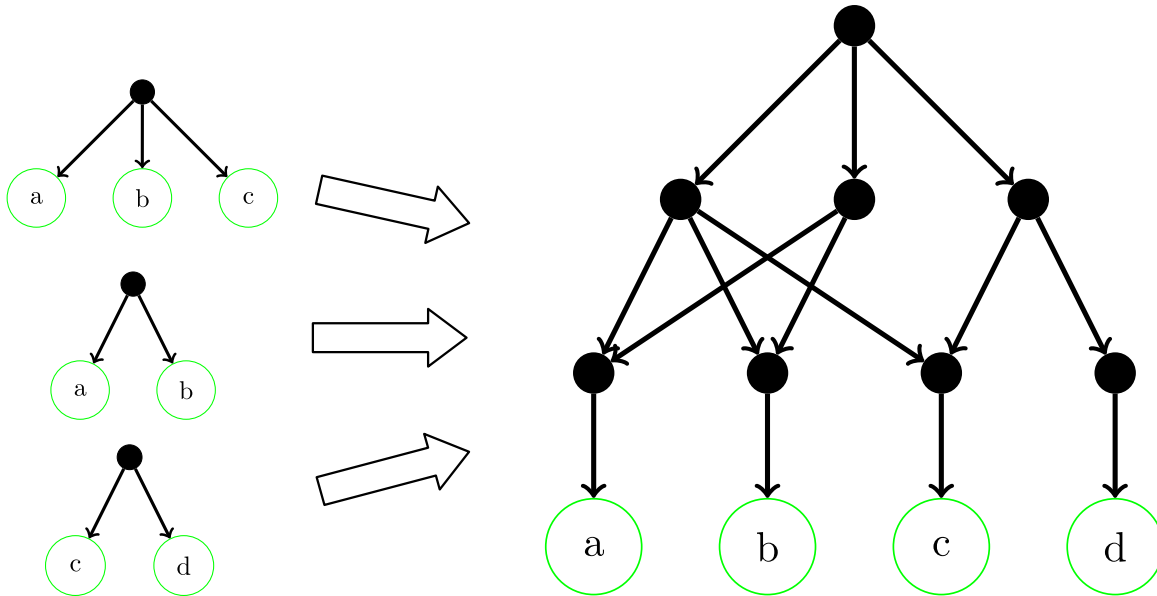E.g. $\mathcal{X} = \{a, b, c, d\}$, $\mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}$



$$k = 3$$

# From Clusters to Networks (Trivial)

E.g. $\mathcal{X} = \{a, b, c, d\}$, $\mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}$
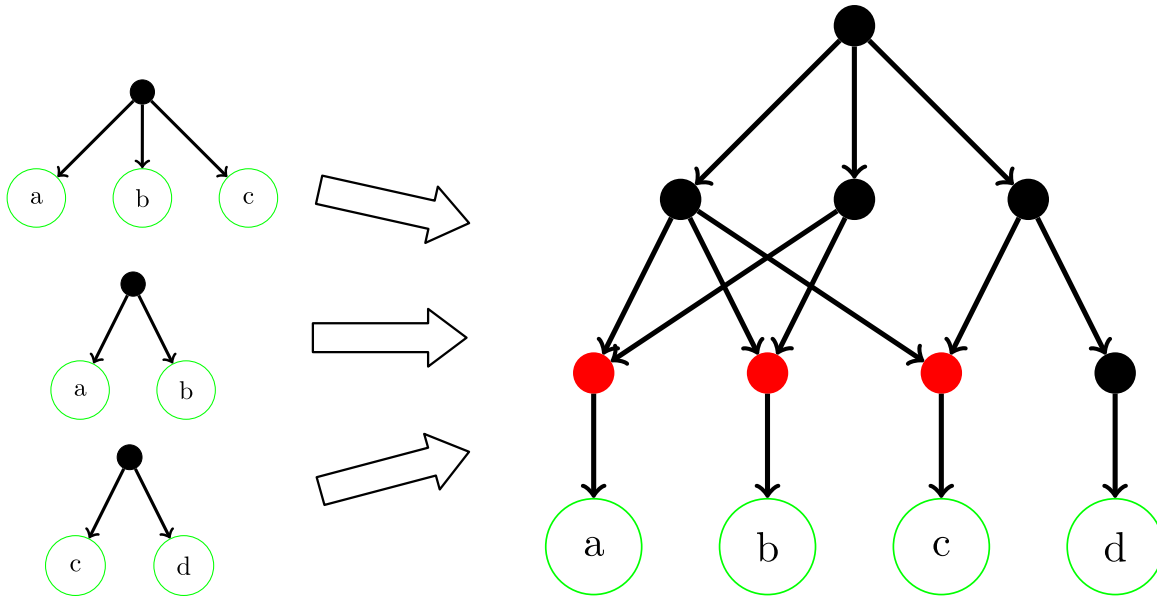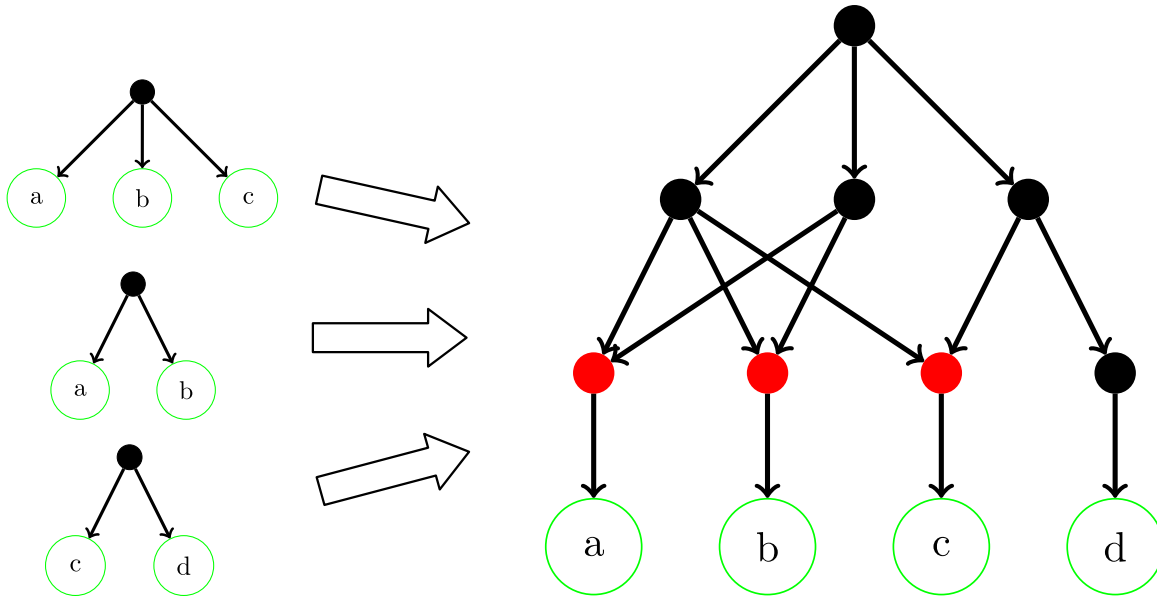


$k = 3$

**But:**

$k = 1$

# From Clusters to Networks

- Constructing a network representing clusters is trivial

# From Clusters to Networks

- Constructing a network representing clusters is trivial

- Constructing a network representing clusters with minimal level is NP-hard

# Generators

Any level-$k$ network can be reduced to a level-$k$ *generator*



Edges and nodes with $\text{out-deg} = 0$ are called *sides*

# Algorithm Outline

- Guess a generator for the network

- Build the network from the generator with guesses

**Disclaimer:** This algorithm is purely theoretical and
doesn't lend itself for practical implementation

# Algorithm

$$\mathcal{X} = \{a, b, c, d\},\ \mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\},\ k = 1$$

1. Guess a generator

# Algorithm

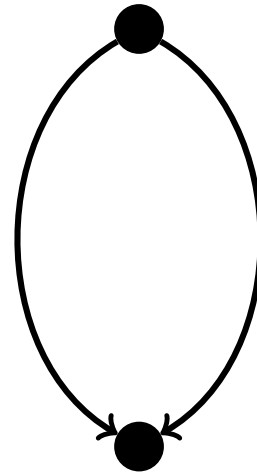$$\mathcal{X} = \{a, b, c, d\}, \, \mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}, \, k = 1$$

1. Guess a generator

# Algorithm

$$\mathcal{X} = \{a, b, c, d\}, \, \mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\}, \, k = 1$$

1. Guess a generator
2. For each side, guess if it has
   $0, 1, 2$, or more leaves on its side

# Algorithm

$$\mathcal{X} = \{a, b, c, d\},\ \mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\},\ k = 1$$

1. Guess a generator
2. For each side, guess if it has
   $0, 1, 2$, or more leaves on its side

# Algorithm

$$\mathcal{X} = \{a, b, c, d\},\ \mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\},\ k = 1$$

1. Guess a generator
2. For each side, guess if it has
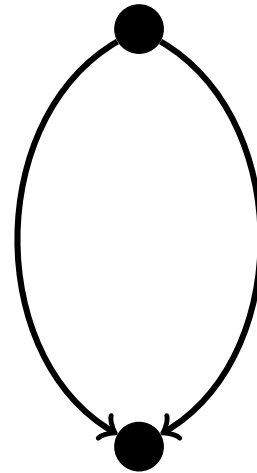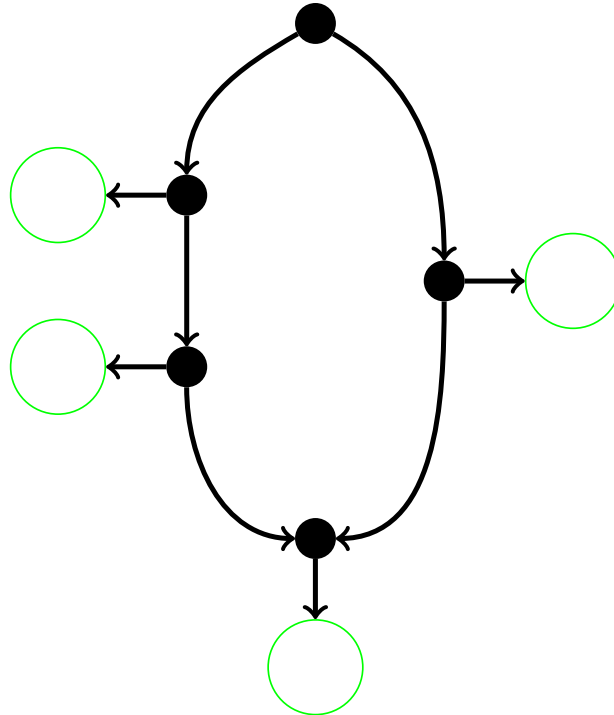   $0, 1, 2$, or more leaves on its side
3. Guess leaves of the sides with $1$ leaf

# Algorithm

$$\mathcal{X} = \{a, b, c, d\}, \; \mathcal{C} = \{\{a, b, c\}, \{a, b\}, \{c, d\}\}, \; k = 1$$

1. Guess a generator
2. For each side, guess if it has $0, 1, 2$, or more leaves on its side
3. Guess leaves of the sides with $1$ leaf

# Algorithm

$$\mathcal{X} = \{a, b, c, d\}, \; \mathcal{C} = \left\{\{a, b, c\}, \{a, b\}, \{c, d\}\right\}, \; k = 1$$

1. Guess a generator
2. For each side, guess if it has $0, 1, 2$, or more leaves on its side
3. Guess leaves of the sides with $1$ leaf
4. Guess top $s^+$ and bottom $s^-$ of the sides with $\geq 2$ leaves

# Algorithm

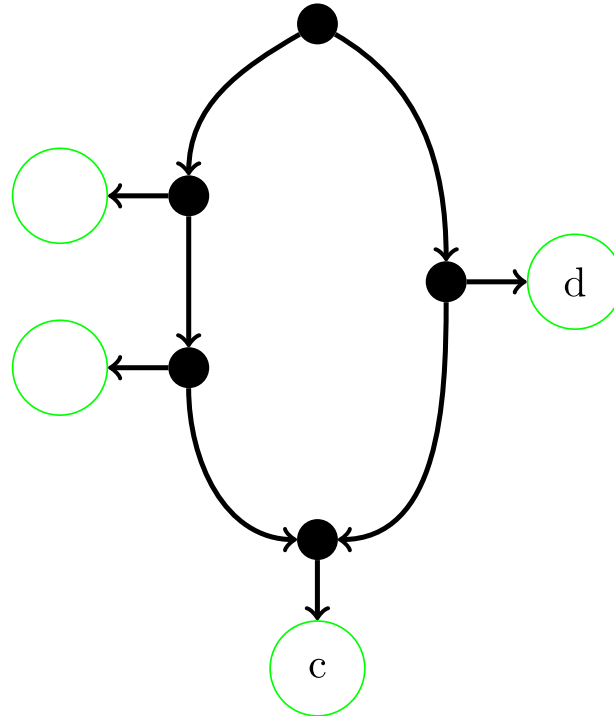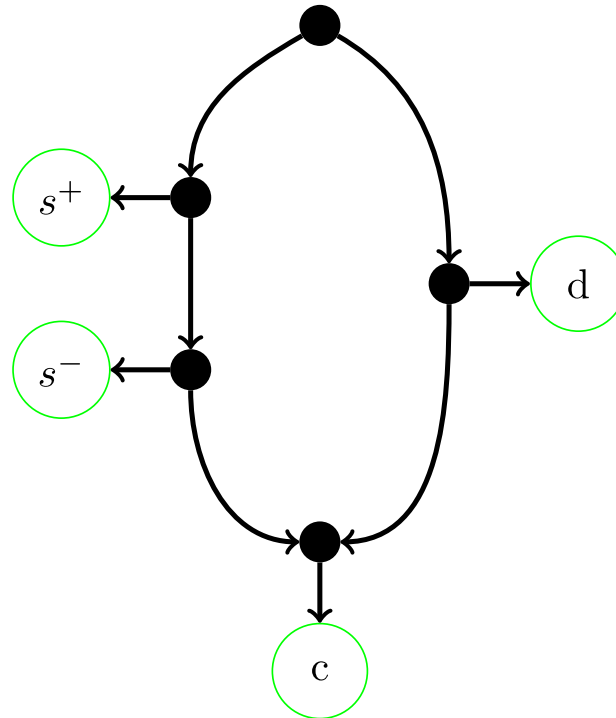$$\mathcal{X} = \{a, b, c, d\}, \ \mathcal{C} = \{\{a, b, c\}, \{a, b\}, \{c, d\}\}, \ k = 1$$

1. Guess a generator
2. For each side, guess if it has $0, 1, 2$, or more leaves on its side
3. Guess leaves of the sides with $1$ leaf
4. Guess top $s^+$ and bottom $s^-$ of the sides with $\geq 2$ leaves

# Algorithm

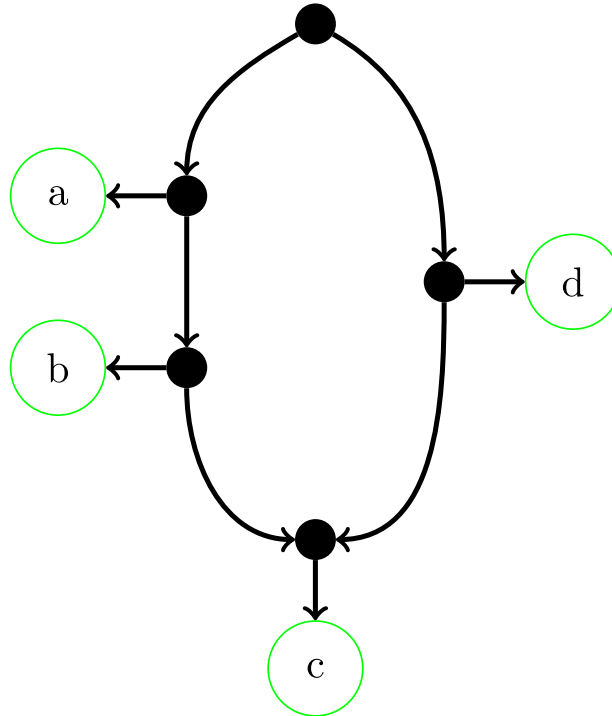$$\mathcal{X} = \{a, b, c, d\}, \ \mathcal{C} = \{\{a, b, c\}, \{a, b\}, \{c, d\}\}, \ k = 1$$

1. Guess a generator
2. For each side, guess if it has $0, 1, 2$, or more leaves on its side
3. Guess leaves of the sides with $1$ leaf
4. Guess top $s^+$ and bottom $s^-$ of the sides with $\geq 2$ leaves
5. Add remaining leaves

# **Algorithm**

$$\mathcal{X} = \{a, b, c, d\},\ \mathcal{C} = \{\{a, b, c\}, \{a, b\}, \{c, d\}\},\ k = 1$$

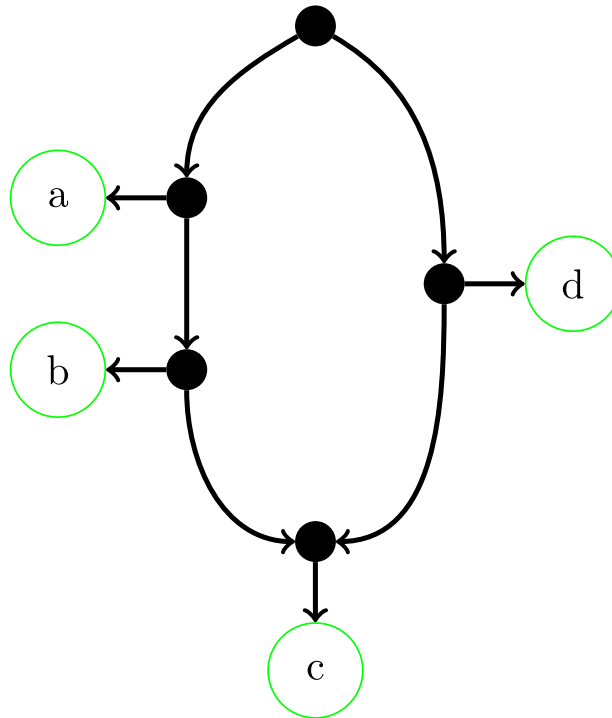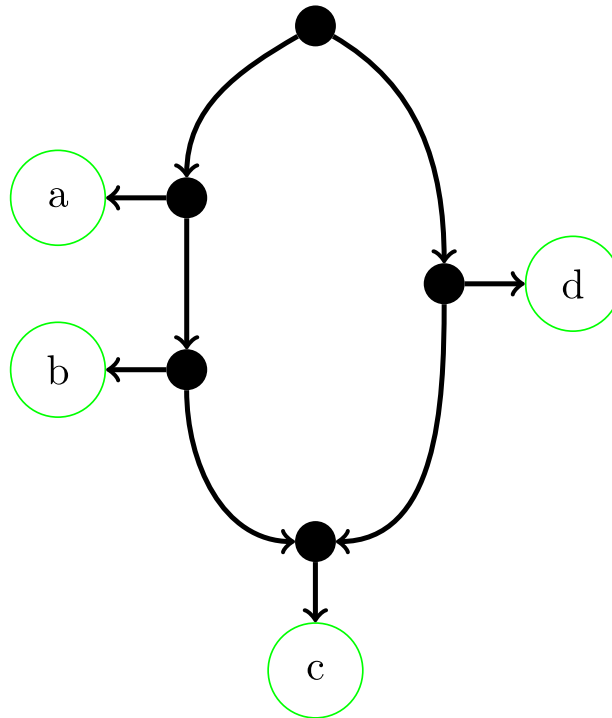1. Guess a generator
2. For each side, guess if it has $0, 1, 2$, or more leaves on its side
3. Guess leaves of the sides with $1$ leaf
4. Guess top $s^+$ and bottom $s^-$ of the sides with $\geq 2$ leaves
5. Add remaining leaves
6. Verify that the network actually represents $\mathcal{C}$

# Adding Remaining Leaves
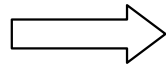
$$\mathcal{C} = \big\{ \{a, b, c\}, \{a, b\}, \{c, d\} \big\}$$

# Adding Remaining Leaves

$$\mathcal{C} = \big\{\{a, b, c\}, \{a, b\}, \{c, d\}\big\} \implies$$

$$a \to b$$
$$b \to a$$
$$c \to d$$
$$d \to c$$

$$a \to b$$
$$=$$
"Every cluster containing $a$, also contains $b$"

# Adding Remaining Leaves

$$\mathcal{C} = \{\{a, b, c\}, \{a, b\}, \{c, d\}\} \implies$$

$$a \to b$$
$$b \to a$$
$$c \to d$$
$$d \to c$$

$$a \to b$$
$$=$$

"Every cluster containing $a$, also contains $b$"

Start with the lowest side

Based on $s^+$, $s^-$ and the relationships from above, you can determine the leaf to be inserted in polynomial time

$$s^+ \to x \to s^- \quad \text{must hold}$$

# Adding Remaining Leaves

$$\mathcal{C} = \{\{a, b, c\}, \{a, b\}, \{c, d\}\} \implies$$

$$a \to b$$
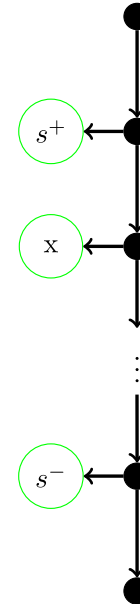$$b \to a$$
$$c \to d$$
$$d \to c$$

$$a \to b$$
$$=$$
"Every cluster containing $a$, also contains $b$"

Start with the lowest side

Based on $s^+$, $s^-$ and the relationships from above, you can determine the leaf to be inserted in polynomial time

$$s^+ \to x \to s^- \quad \text{must hold}$$

If no leaf can be added $\implies$ Start new

# Complexity Recap

- Guess from a **constant** number of generators

- For a **constant** number of sides, guess how many leaves from 4 options

- $\mathcal{O}(n^2)$ guesses for top and bottom

- Polynomial time for the remaining leaves

- Polynomial time to verify

$\Longrightarrow$ **Very polynomial-time algorithm!**

# Issues

# Issues

- # generators explodes

# Issues

- # generators explodes

$$[2^{k-1}, k!^2 50^k]$$

| $k$ | generators |
|---|---|
| 1 | 1 |
| 2 | 4 |
| 3 | 65 |
| 4 | 1993 |
| 5 | 91454 |

# Issues

- \# generators explodes

$$[2^{k-1}, k!^2 50^k]$$

| $k$ | generators |
|-----|------------|
| 1 | 1 |
| 2 | 4 |
| 3 | 65 |
| 4 | 1993 |
| 5 | 91454 |

- \# edges (and hence sides) in the generators grows linearly in $k$

# Issues

- # generators explodes
$$[2^{k-1}, k!^2 50^k]$$

| $k$ | generators |
|---|---|
| 1 | 1 |
| 2 | 4 |
| 3 | 65 |
| 4 | 1993 |
| 5 | 91454 |

- # edges (and hence sides) in the generators grows linearly in $k$

- The need to guess $0, 1, 2, > 2$ leaves per side

# Issues

- # generators explodes

$$[2^{k-1}, k!^2 50^k]$$

| $k$ | generators |
|---|---|
| 1 | 1 |
| 2 | 4 |
| 3 | 65 |
| 4 | 1993 |
| 5 | 91454 |

- # edges (and hence sides) in the generators grows linearly in $k$

- The need to guess $0, 1, 2, > 2$ leaves per side

- The need to make up to $\mathcal{O}(n^2)$ guesses per side to guess $s^+/s^-$

# Issues

- # generators explodes

  $$[2^{k-1}, k!^2 50^k]$$

  | $k$ | generators |
  |---|---|
  | 1 | 1 |
  | 2 | 4 |
  | 3 | 65 |
  | 4 | 1993 |
  | 5 | 91454 |

- # edges (and hence sides) in the generators grows linearly in $k$

- The need to guess $0, 1, 2, > 2$ leaves per side

- The need to make up to $\mathcal{O}(n^2)$ guesses per side to guess $s^+/s^-$

**In sum this makes the algorithm practically unfeasible :(**

# The End

Thank you for your attention!