# Reconstruction of ancestral gene orders
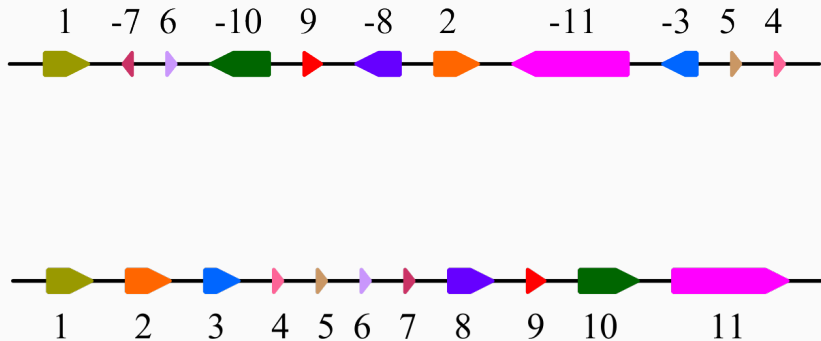
Pedro Feijao
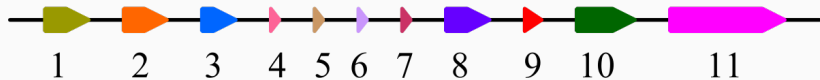
Banff, February 17, 2017

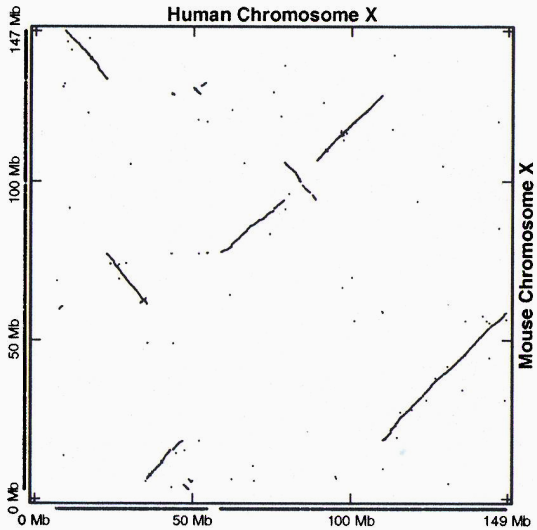SFU SIMON FRASER UNIVERSITY
ENGAGING THE WORLD

Universität Bielefeld Technische Fakultät
Genominformatik

# Introduction

**Mouse X-Chromosome**

1    -7  6    -10    9    -8    2    -11    -3  5  4

**Human X-Chromosome**

1    2    3    4  5  6  7    8    9    10    11

(Pevzner and Tesler, 2003)

Human Chromosome X / Mouse Chromosome X
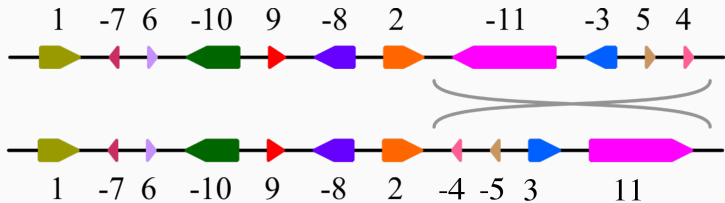
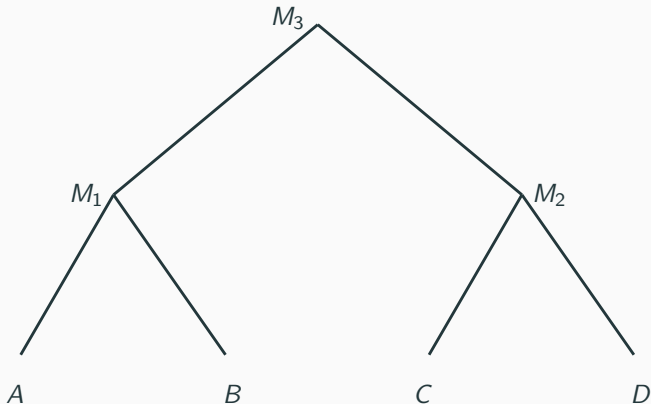- **Distance**: Minimum # of rearrangements from $A$ to $B$?

- **Distance**: Minimum # of rearrangements from $A$ to $B$?

- **Scenario**: Which rearrangements? (also called *Sorting*)

- **Distance**: Minimum # of rearrangements from $A$ to $B$?

- **Scenario**: Which rearrangements? (also called *Sorting*)

- **Phylogeny**: How did the genomes evolve?

- **Distance**: Minimum # of rearrangements from $A$ to $B$?

- **Scenario**: Which rearrangements? (also called *Sorting*)

- **Phylogeny**: How did the genomes evolve?

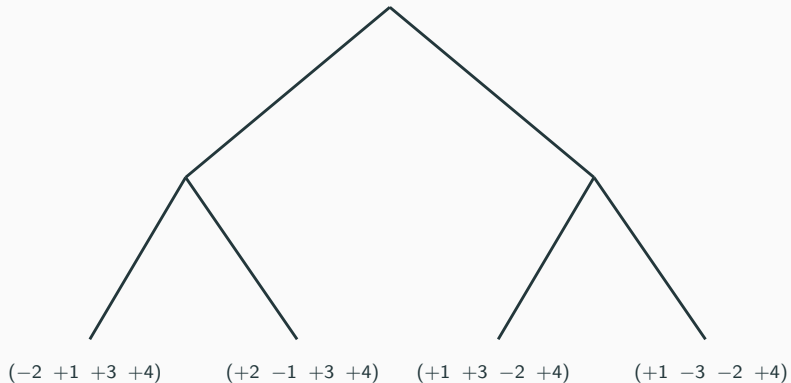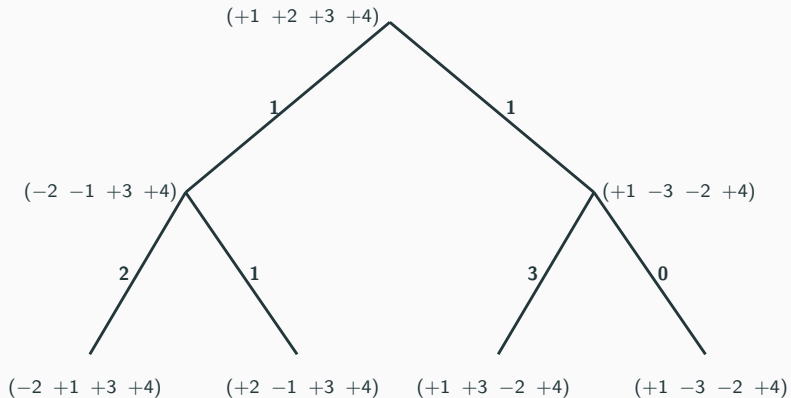- **Ancestral Reconstruction**: How do the ancestors look like?
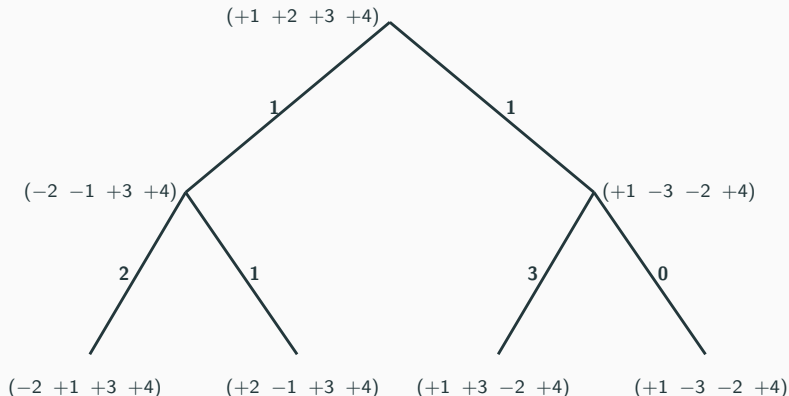
# Ancestral Reconstruction

**Input**: Tree and genomes $A, B, C, D$

**Ouput**: Ancestral genomes $(M_1, M_2, M_3)$
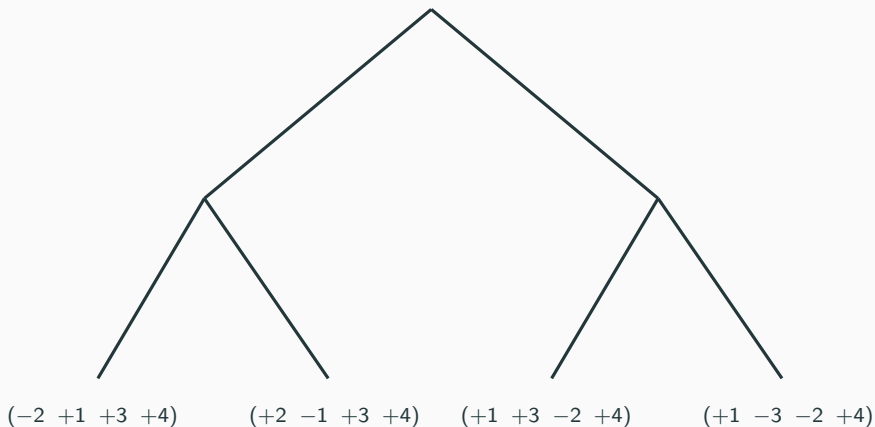
- Distance-based Methods

- Homology-based Methods

$$(-2 \ +1 \ +3 \ +4) \qquad (+2 \ -1 \ +3 \ +4) \qquad (+1 \ +3 \ -2 \ +4) \qquad (+1 \ -3 \ -2 \ +4)$$

(+1 +2 +3 +4)

1     1

(−2 −1 +3 +4)     (+1 −3 −2 +4)

2   1     3   0

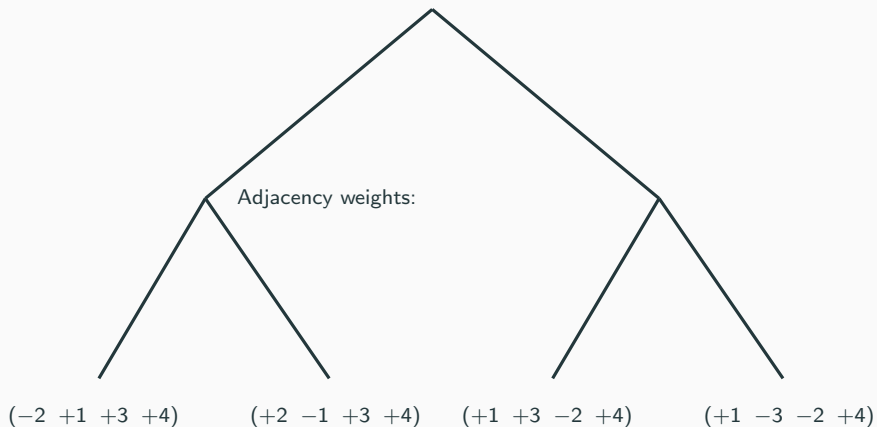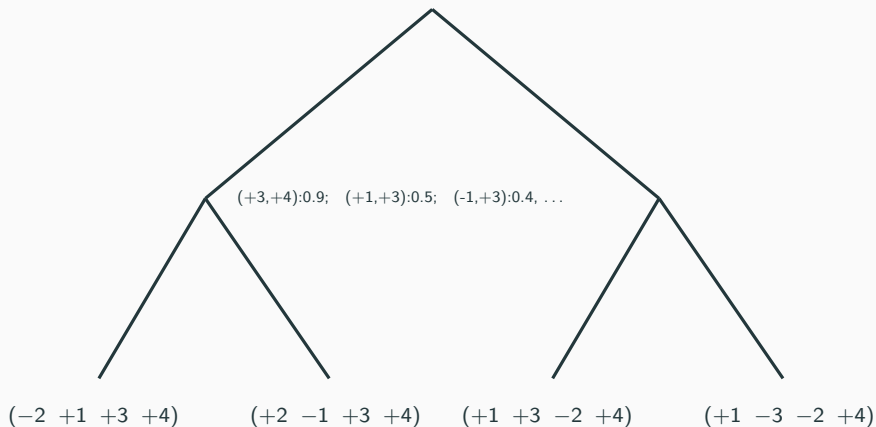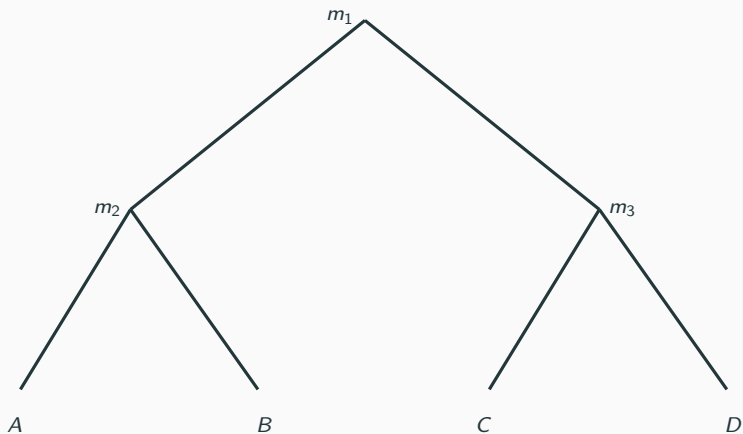(−2 +1 +3 +4)   (+2 −1 +3 +4)   (+1 +3 −2 +4)   (+1 −3 −2 +4)

Find ancestral genomes that **minimize events** on the tree
→ **Small Parsimony Problem**

$(-2\ +1\ +3\ +4)$ $\quad$ $(+2\ -1\ +3\ +4)$ $\quad$ $(+1\ +3\ -2\ +4)$ $\quad$ $(+1\ -3\ -2\ +4)$

Adjacency weights:

$(-2 \; +1 \; +3 \; +4)$    $(+2 \; -1 \; +3 \; +4)$    $(+1 \; +3 \; -2 \; +4)$    $(+1 \; -3 \; -2 \; +4)$

$(+3,+4):0.9;$   $(+1,+3):0.5;$   $(-1,+3):0.4,$ . . .

$(-2\ +1\ +3\ +4)$     $(+2\ -1\ +3\ +4)$     $(+1\ +3\ -2\ +4)$     $(+1\ -3\ -2\ +4)$
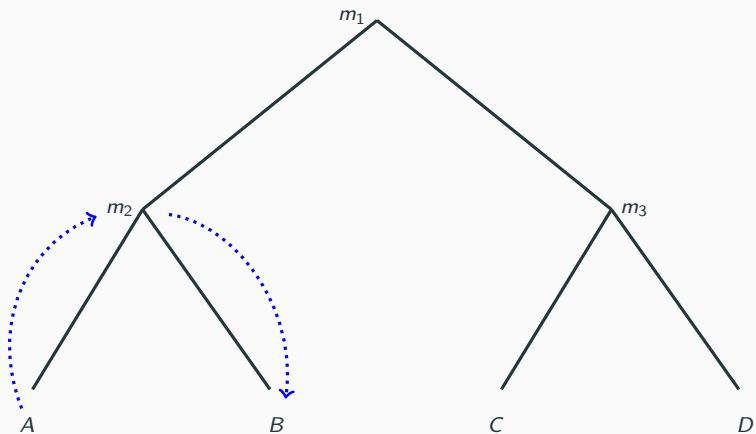
- Distance-based methods:

  - Assume a rearrangement model

  - Minimize branch lengths

- Homology-based methods:

  - Find conserved structures

  - Maximize some weight/probability function

**Ancestral Reconstruction** where internal nodes

are

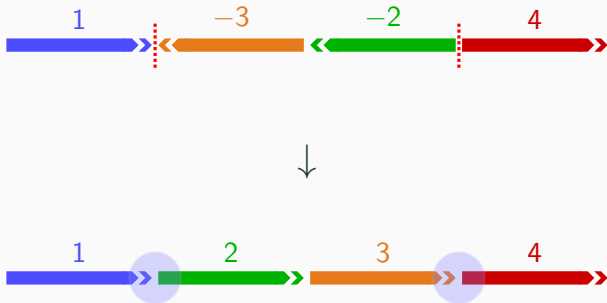**Intermediate Genomes** of its children.

# Definitions

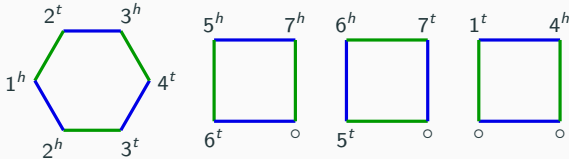$$A = \{\circ 1^t, 1^h 2^t, 2^h 3^h, 3^t 4^t, 4^h \circ\}$$

$$A = \{\circ 1^t, 1^h 2^t, 2^h 3^t, 3^h 4^t, 4^h \circ, \circ 5^t, 5^h 6^t, 6^h 7^t, 7^h \circ\}$$
$$B = \{1^h 2^h, 2^t 3^h, 3^t 4^t, 4^h 1^t, \circ 6^t, 6^h 5^t, 5^h 7^h, 7^t \circ\}$$
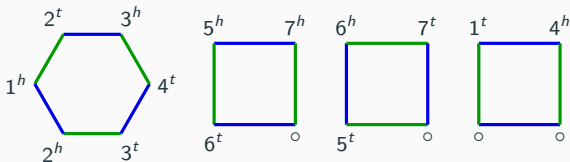
$$A = \{\circ 1^t, 1^h 2^t, 2^h 3^t, 3^h 4^t, 4^h \circ, \circ 5^t, 5^h 6^t, 6^h 7^t, 7^h \circ\}$$
$$B = \{1^h 2^h, 2^t 3^h, 3^t 4^t, 4^h 1^t, \circ 6^t, 6^h 5^t, 5^h 7^h, 7^t \circ\}$$



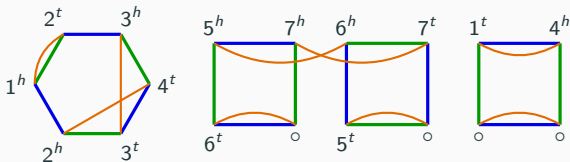$A$-edges are drawn in green, and $B$-edges in blue.

$$d_{\mathrm{DCJ}}(A, B) = N - C$$

where $N$ is the number of genes and $C$ is the number of cycles in $BP(A, B)$.
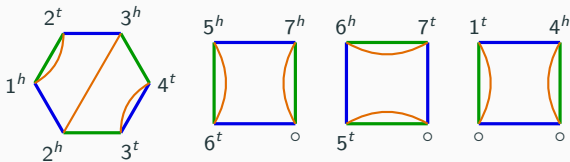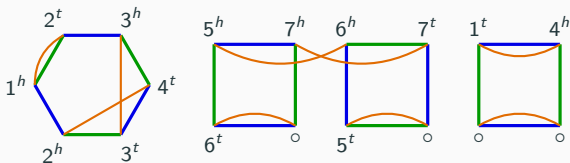
(Bergeron et al, 2006)

Genomes are **matchings** in the BP graph:

Genomes are **matchings** in the BP graph:

Genomes are **matchings** in the BP graph:

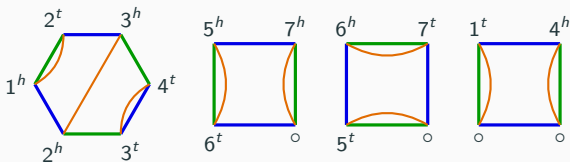Edges are non-crossing chords in the cycles of $BP(A, B)$

$\Rightarrow$

Intermediate Genome of $A$ and $B$

Edges are non-crossing chords in the cycles of $BP(A, B)$
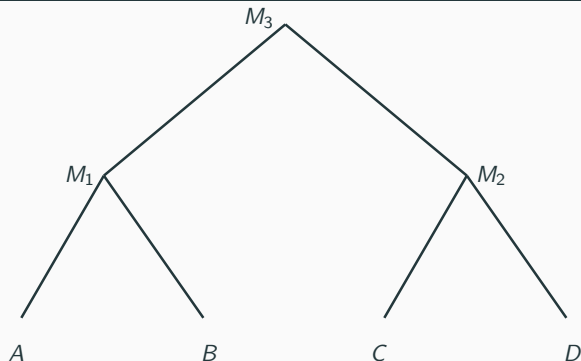
$\Rightarrow$

Intermediate Genome of $A$ and $B$

- Very easy to detect (linear time)

- Very easy to detect (linear time)
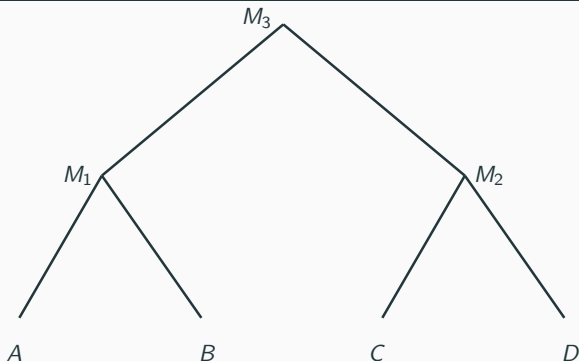
- Reduces the search space. In the example:

- Very easy to detect (linear time)

- Reduces the search space. In the example:

  - 34,459,425 possible genomes

- Very easy to detect (linear time)

- Reduces the search space. In the example:

  - 34,459,425 possible genomes

  - Only 40 intermediate genomes between A and B.

# Methods

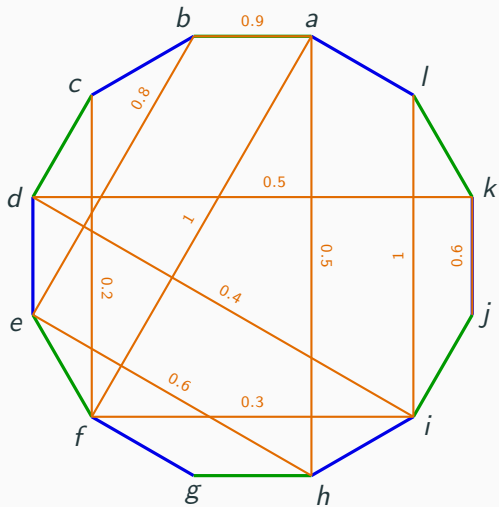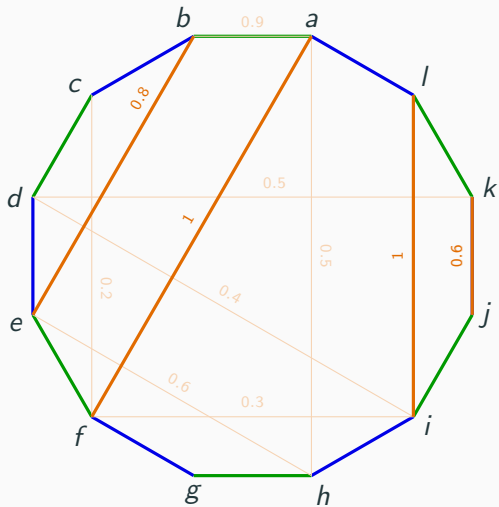- Small parsimony with the restriction that internal nodes are IG's of the children.

- Small parsimony with the restriction that internal nodes are IG's of the children.
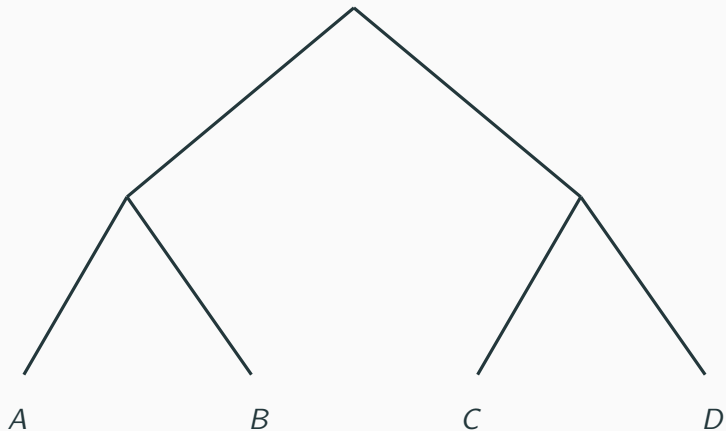
- Still NP-hard

- Given adjacency weights, can we find an IG with maximum weight?

- **Maximum Weight Independent Set**: Polynomial Time
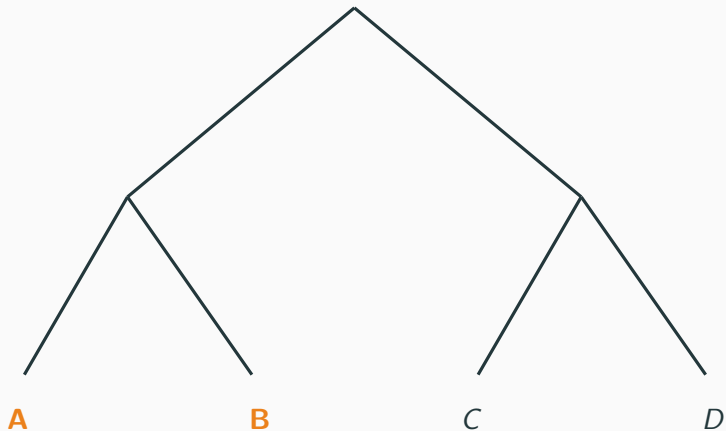
SFU

Universität Bielefeld

- DeClone (Chauve et al., 2015)

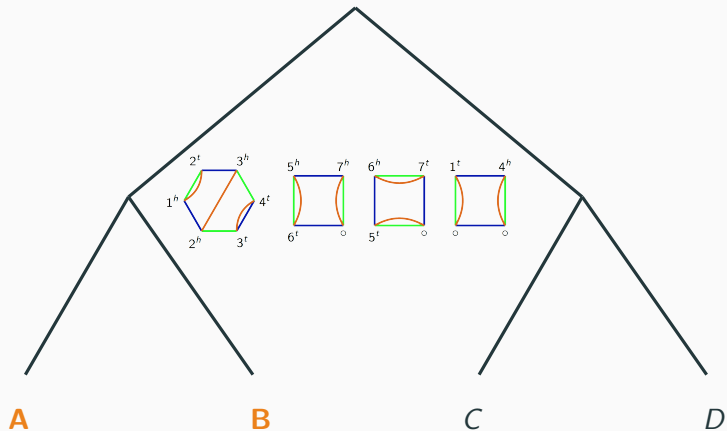- New proposed algorithm based on InferCARs (Ma et al., 2006).
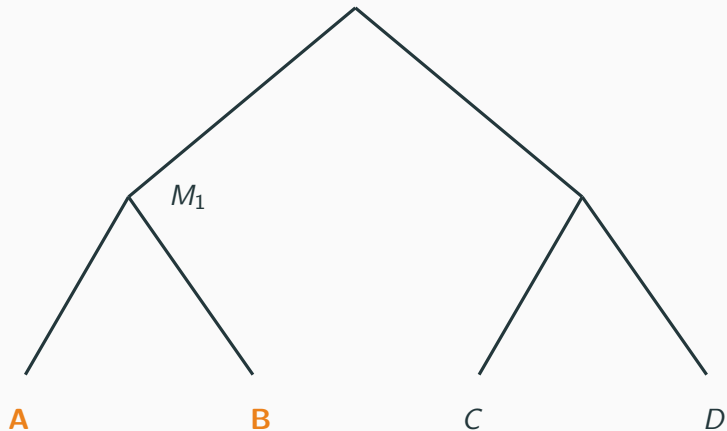
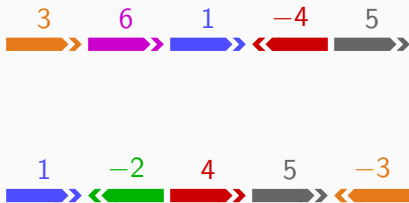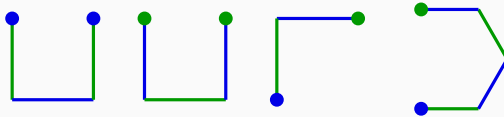(+3,+4):0.9;  (+1,+3):0.5;  (-1,+3):0.4, ...
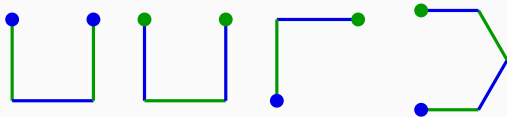
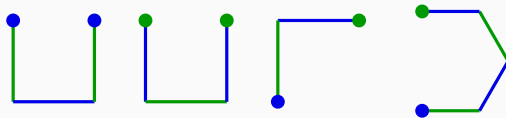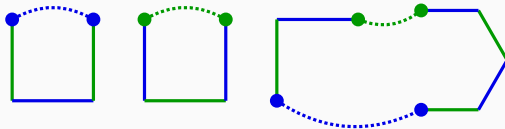A          B          C          D

DCJ InDel Model (Braga et al., 2010; Compeau, 2012)

**New components**: AA-, BB-, AB- , A-, and B-paths.

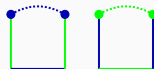**New components**: AA-, BB-, AB- , A-, and B-paths.

Find an **optimal completion**

**New components**: AA-, BB-, AB- , A-, and B-paths.

Find an **optimal completion**
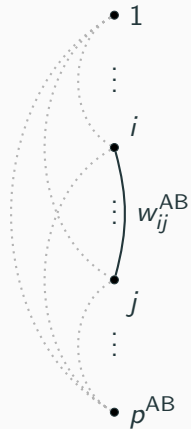
- *AA*- and *BB*- components are closed
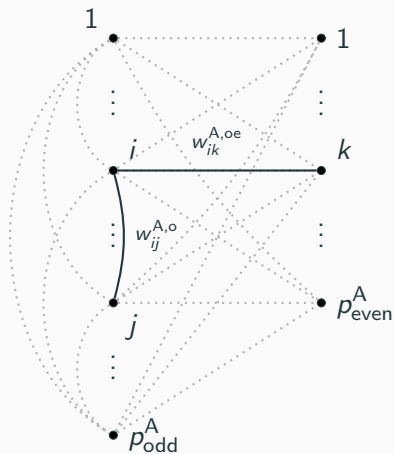
- *AB*- are paired

- *A*- and *B*- paths are paired, with opposing parity.

- Sometimes *A*-, *B*- and *AB*- paths are joined in a **triplet**.

- How to find a completion with maximum weight?

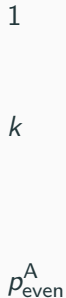- Calculate all possible pairings and solve a *Maximum Weight Matching*
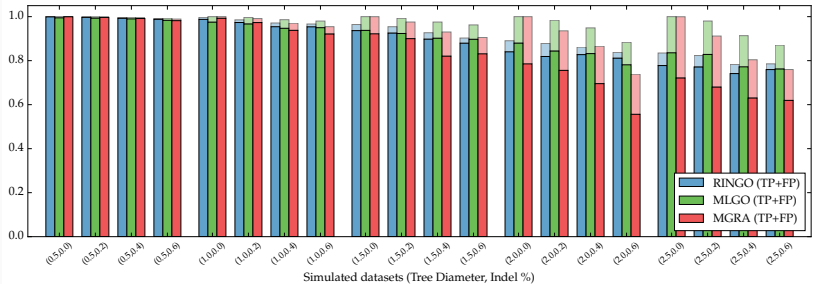
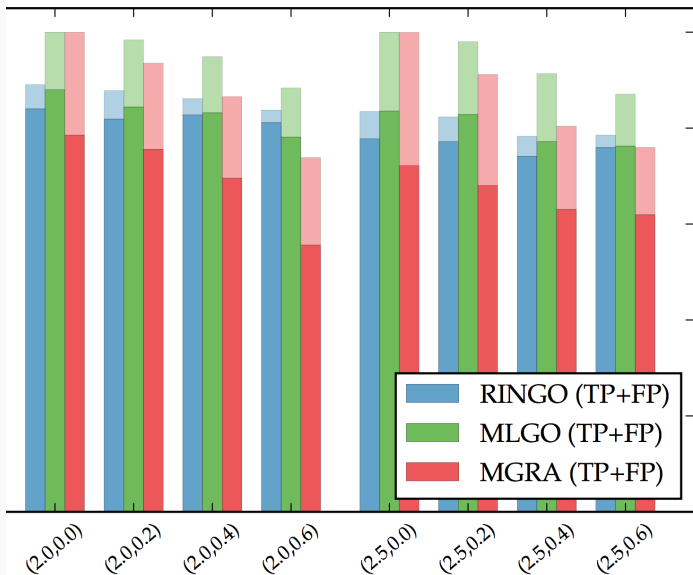$AB$-paths     odd $A$-paths     even $A$-paths

- Sometimes *A*-, *AB*-, *B*- triplets are possible.

- Triple matching is usually NP-hard, but it is still open in this case.

# Results

- RINGO - ancestral **R**econstruction with **IN**termediate **G**en**O**mes (Feijao and Araujo, 2016)

- MGRA2 (Avdeyev et al., 2016)

- MLGO (Hu et al., 2014)

| Dataset | $I = 1$, unitary indels | | | | |
|---|---|---|---|---|---|
| **Diameter (D)** | $0.5n$ | $1n$ | $1.5n$ | $2n$ | $2.5n$ |
| RINGO | 3s | 3s | 5s | 7s | 7s |
| MLGO | 1m6s | 1m10s | 1m7s | 1m9s | 1m16s |
| MGRA | 7s | 1m46s | 12m12s | 56m55s | 2h2m41s |

- Duplicated genes

- Statistical models

- Elói Araújo (UFMS, Brazil)



- Jens Stoye (Bielefeld University, Germany)

**Thanks!**