# Topics of today:

# Breakpoint graph of canonical genomes

Genomes $\mathbb{A}_{\triangleright}^{f}$ and $\mathbb{B}_{\triangleright}^{f}$ are canonical, with $\mathcal{F}_{\star} = \mathcal{G}_{\star} = \{1, 2, \ldots, n\}$
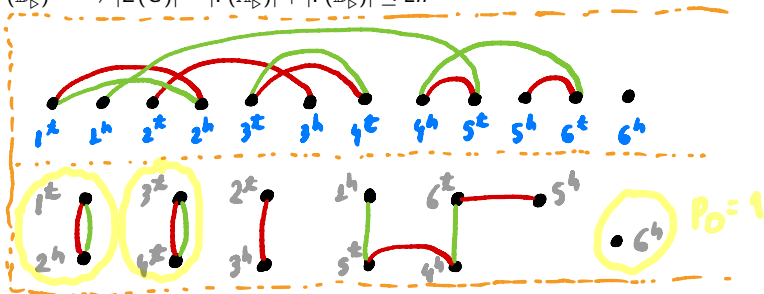
**Breakpoint graph** of $\mathbb{A}_{\triangleright}^{f}$ and $\mathbb{B}_{\triangleright}^{f}$ $\rightarrow$ $G = BG(\mathbb{A}_{\triangleright}^{f}, \mathbb{B}_{\triangleright}^{f})$:

1. Set of vertices $V(G) = \bigcup_{\mathbb{X} \in \mathcal{G}_{\star}} \{\mathbb{X}^h, \mathbb{X}^t\}$ $\Rightarrow |V(G)| = 2n$

2. Set of edges $E(G) = \Gamma(\mathbb{A}_{\triangleright}^{f}) \cup \Gamma(\mathbb{B}_{\triangleright}^{f})$ $\Rightarrow |E(G)| = |\Gamma(\mathbb{A}_{\triangleright}^{f})| + |\Gamma(\mathbb{B}_{\triangleright}^{f})| \leq 2n$



$A = \begin{bmatrix} \bar{1} & \bar{2} & \bar{3} & 4 & 5 & 6 \end{bmatrix}$

$B = \begin{bmatrix} 2 & 1 & 5 \end{bmatrix} \begin{bmatrix} \bar{6} & \bar{4} & 3 \end{bmatrix}$

Each vertex has degree 1 or 2: collection of $p$ paths and $c$ (even) cycles $\qquad$ ( $p = \kappa(\mathbb{A}_{\triangleright}^{f}) + \kappa(\mathbb{B}_{\triangleright}^{f})$ )

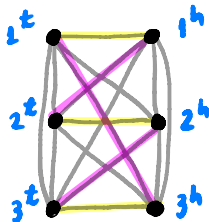length of a component: number of edges [alternating]

$d_{\mathrm{BP}}(\mathbb{A}_{\triangleright}^{f}, \mathbb{B}_{\triangleright}^{f}) = n - c_2 - \frac{p_0}{2} \begin{cases} c_2 = \text{number of 2-cycles in } BG(\mathbb{A}_{\triangleright}^{f}, \mathbb{B}_{\triangleright}^{f}) \text{ (common adjacencies)} \\ p_0 = \text{number of 0-paths in } BG(\mathbb{A}_{\triangleright}^{f}, \mathbb{B}_{\triangleright}^{f}) \text{ (common telomeres)} \end{cases}$

# Complete graph of a set of genes

**Complete graph** $\mathfrak{G}$ of $\mathcal{A} = \{1, 2, \ldots, n\}$:

Set of vertices $V(\mathfrak{G}) = \bigcup_{x \in \mathcal{A}} \{x^h, x^t\} \quad \Rightarrow |V(\mathfrak{G})| = 2n$

$Ex: n = 3$



$C_{\bigcirc} = (1)(2)(3)$

$C_{\bullet} = (1\ 2\ 3)$

$\rightarrow$ Set of non-incident edges covering all vertices

A perfect matching $M$ in $\mathfrak{G}$ corresponds to $|M| = n$ adjacencies and, consequently, defines a circular singular genome $\mathbb{C}$, with $\Gamma(\mathbb{C}) = M$.

# SCJ median of $k$ canonical genomes

Given $k$ canonical genomes $\mathbb{C}_1^f$, $\mathbb{C}_2^f$, ... $\mathbb{C}_k^f$, find another canonical genome $\mathbb{M}^f$ that minimizes the sum:

$$\mathsf{s}_{\mathrm{SCJ}}(\mathbb{M}^f) \quad = \quad \mathsf{d}_{\mathrm{SCJ}}(\mathbb{M}^f, \mathbb{C}_1^f) + \mathsf{d}_{\mathrm{SCJ}}(\mathbb{M}^f, \mathbb{C}_2^f) + ... + \mathsf{d}_{\mathrm{SCJ}}(\mathbb{M}^f, \mathbb{C}_k^f)$$

$$= \quad |\Gamma(\mathbb{C}_1^f)| + |\Gamma(\mathbb{C}_2^f)| + ... + |\Gamma(\mathbb{C}_k^f)| \quad + \quad \omega(\mathbb{M}^f)$$

For computing the median, we need to minimize:

$$\omega(\mathbb{M}^f) \ = \sum_{xy \in \Gamma(\mathbb{M}^f)} \omega(xy)$$

where $\omega(xy) = k - 2 \cdot \phi(xy, \mathbb{C}_{1..k}^f) \in \{-k, -k+2, ..., +k-2, +k\}$.

**Solution:** take only the adjacencies with negative weight

# SCJ median of $k$ canonical circular genomes

Canonical genomes $\mathbb{C}_1^f$, $\mathbb{C}_2^f$, $\ldots \mathbb{C}_k^f$ are circular:

$$s_{\text{SCJ}}(\mathbb{M}^f) \quad = \quad |\Gamma(\mathbb{C}_1^f)| + |\Gamma(\mathbb{C}_2^f)| + \ldots + |\Gamma(\mathbb{C}_k^f)| \quad + \quad \omega(\mathbb{M}^f)$$

$$= \quad \cancel{3n} \quad + \quad \omega(\mathbb{M}^f)$$

*kn*

Again, we need to minimize $\omega(\mathbb{M}^f) = \sum_{xy \in \Gamma(\mathbb{M}_b^f)} \omega(xy)$, where $\omega(xy) = k - 2 \cdot \phi(xy, \mathbb{C}_{1..k}^f)$, but since the median is required to be circular, it may not be possible to take only the adjacencies that have a negative weight

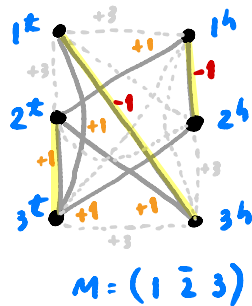**Solution:**

1. Build the complete graph $\mathfrak{G}$ of $\mathcal{G}_\star$
2. Assign weights to each edge $xy$ of $\mathfrak{G}$: $\omega(xy) = k - 2 \cdot \phi(xy, \mathbb{C}_{1..k}^f)$.

$$\omega(xy) \in \{-3, -1, +1, +3\}$$

$$(3) \quad (2) \quad (1) \quad (0)$$

Ex: $C_1 = (1\ 2\ 3)$

$C_2 = (3\ 2\ \bar{1})$

$C_3 = (3, 1, \bar{2})$



$$M = (1\ \bar{2}\ 3)$$

Perfect matching $M$ in $\mathfrak{G}$ $\Leftrightarrow$ Circular genome $\mathbb{M}^f$ ; with $\omega(M) = \omega(\mathbb{M}^f)$

A perfect matching $M_{\text{MIN}}$ with **minimum weight** gives a circular SCJ median $\mathbb{M}_{min}^f$ with **minimum weight**

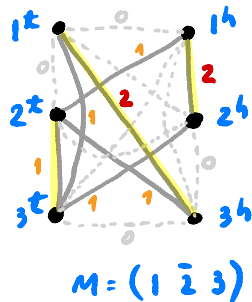# Breakpoint median of three canonical circular genomes

Given canonical circular genomes $\mathbb{C}_1^f$, $\mathbb{C}_2^f$ and $\mathbb{C}_3^f$, find a canonical circular genome $\mathbb{M}_\triangleright^f$ that minimizes the sum:

$$
\begin{aligned}
s_{BP}(\mathbb{M}_\triangleright^f) \quad &= \quad d_{BP}(\mathbb{M}_\triangleright^f, \mathbb{C}_1^f) \;+\; d_{BP}(\mathbb{M}_\triangleright^f, \mathbb{C}_2^f) \;+\; d_{BP}(\mathbb{M}_\triangleright^f, \mathbb{C}_k^f) \\
&= \quad n - \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_1^f) \;+\; n - \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_2^f) \;+\; n - \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_3^f) \\
&= \quad 3n \;-\; \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_{1..3}^f) \\
&= \quad 3n \;-\; \omega'(\mathbb{M}_\triangleright^f)
\end{aligned}
$$

Here we need to maximize $\omega'(\mathbb{M}_\triangleright^f) = \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \omega'(xy)$, where $\omega'(xy) = \phi(xy, \mathbb{C}_{1..3}^f)$.

1. Build the complete graph $\mathfrak{G}$ of $\mathcal{G}_\star$
2. Assign weights to each edge $xy$ of $\mathfrak{G}$: $\omega'(xy) = \phi(xy, \mathbb{C}_{1..k}^f)$.

$$\omega(xy) \in \{0, 1, 2, 3\} \qquad \text{Ex: } C_1 = (1\ 2\ 3)$$
$$C_2 = (3\ 2\ \overline{1})$$
$$C_3 = (3, 1, \overline{2})$$



$$M = (1\ \overline{2}\ 3)$$

Perfect matching $M$ in $\mathfrak{G}$ $\Leftrightarrow$ Circular genome $\mathbb{M}^f$ ; with $\omega'(M) = \omega'(\mathbb{M}^f)$

A perfect matching $M_{\text{MIN}}$ with **maximum weight** gives a circular BP median $\mathbb{M}_{min}^f$ with **maximum weight**

# Breakpoint median of three canonical genomes

Given canonical genomes $\mathbb{C}_1^f$, $\mathbb{C}_2^f$ and $\mathbb{C}_3^f$, find a canonical genome $\mathbb{M}_\triangleright^f$ that minimizes the sum:

$$
\begin{aligned}
s_{\mathrm{BP}}(\mathbb{M}_\triangleright^f) &= d_{\mathrm{BP}}(\mathbb{M}_\triangleright^f, \mathbb{C}_1^f) + d_{\mathrm{BP}}(\mathbb{M}_\triangleright^f, \mathbb{C}_2^f) + d_{\mathrm{BP}}(\mathbb{M}_\triangleright^f, \mathbb{C}_k^f) \\
&= n - \textstyle\sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_1^f) - \sum_{x \in \Theta(\mathbb{M}_\triangleright^f)} \frac{\phi(x, \mathbb{C}_1^f)}{2} + n - \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_2^f) \\
&\quad - \textstyle\sum_{x \in \Theta(\mathbb{M}_\triangleright^f)} \frac{\phi(x, \mathbb{C}_2^f)}{2} + n - \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_3^f) - \sum_{x \in \Theta(\mathbb{M}_\triangleright^f)} \frac{\phi(x, \mathbb{C}_3^f)}{2} \\
&= 3n - \textstyle\sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \phi(xy, \mathbb{C}_{1..3}^f) - \sum_{x \in \Theta(\mathbb{M}_\triangleright^f)} \frac{\phi(x, \mathbb{C}_{1..3}^f)}{2}
\end{aligned}
$$

Here $\omega'(\mathbb{M}_\triangleright^f) = \sum_{xy \in \Gamma(\mathbb{M}_\triangleright^f)} \omega'(xy) + \sum_{x \in \Theta(\mathbb{M}_\triangleright^f)} \omega'(x)$, where $\omega'(xy) = \phi(xy, \mathbb{C}_{1..3}^f)$ and $\omega'(x) = \frac{\phi(x, \mathbb{C}_{1..3}^f)}{2}$.

1. Build the complete graph $\mathfrak{G}$
2. Assign weights to each edge $xy$ of $\mathfrak{G}$: $\omega'(xy) = \phi(xy, \mathbb{C}_{1..k}^f)$.
3. Build the complete graph $\mathfrak{G}_t$ with vertices $V(\mathfrak{G}_t) = \bigcup_{\chi \in \mathcal{G}_\star} \{t_{\chi^h}, t_{\chi^t}\}$
4. Assign weight 0 to each edge of $\mathfrak{G}_t$
5. Add one edge connecting each vertex $x$ in $\mathfrak{G}$ to the corresponding vertex $t_x$ in $\mathfrak{G}_t$, with weight $\omega'(xt_x) = \frac{\phi(x, \mathbb{C}_{1..k}^f)}{2}$

Perfect matching $M$ in $\mathfrak{G} + \mathfrak{G}_t \Leftrightarrow$ Genome $\mathbb{M}^f$ ; with $\omega'(M) = \omega'(\mathbb{M}^f)$

A matching $M_{\mathrm{MIN}}$ with **maximum weight** gives a BP median $\mathbb{M}_{min}^f$ with **maximum weight**

$$L_1 = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$$

$$L_2 = \begin{bmatrix} 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 4 \end{bmatrix}$$

$$L_3 = \begin{bmatrix} \bar{1} & 2 & \bar{4} & \bar{3} \end{bmatrix}$$

$$M = \begin{bmatrix} \bar{1} & 2 & 3 & 4 \end{bmatrix}$$

$$w(M) = 6$$

$L_1 = \begin{bmatrix} 1 & 2 & 3 & 4 \end{bmatrix}$

$L_2 = \begin{bmatrix} 3 & 2 & 1 \end{bmatrix} \begin{bmatrix} 4 \end{bmatrix}$

$L_3 = \begin{bmatrix} \bar{1} & 2 & \bar{4} & \bar{3} \end{bmatrix}$

$M_2 = \begin{bmatrix} 3 & 4 & \bar{2} & 1 \end{bmatrix}$

$\omega(M_2) = 6$

The breakpoint halving can

be computed in a similar way

# Quiz 1

1 Which of the following statements are true?

    A The breakpoint median can only be computed for circular genomes.

    B The circular SCJ median is equivalent to the circular breakpoint median of three canonical circular genomes.

    C The problem of computing a circular breakpoint halving of a circular duplicated genome is polynomial.

# NP-hardness of unichromosomal breakpoint median

A unichromosomal circular genome $\mathbb{C}$ can be represented as a simple directed cycle graph:

Ex:  $\mathbb{C} = (1\,\bar{2}\,3)$



OR

Assume that the genes in three canonical circular genomes $\mathbb{C}_1^f$, $\mathbb{C}_2^f$ and $\mathbb{C}_3^f$ have the same relative orientation and represent these three genomes in the same directed cycle graph:

Ex:  $\mathbb{C}_1^f = (1\,2\,3\,4)$ ,  $\mathbb{C}_2^f = (2\,4\,1\,3)$ ,  $\mathbb{C}_3^f = (2\,3\,1\,4)$
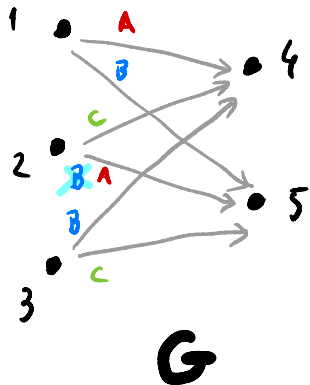


Every vertex has indegree = outdegree = 3

$M = (1\;2\;3\;4)$

# NP-hardness of unichromosomal breakpoint median

The Problem of determining whether a directed graph $G$ has a hamiltonian cycle is NP-complete, even if $G$ has maximum indegree and maximum outdegree equal to 3.
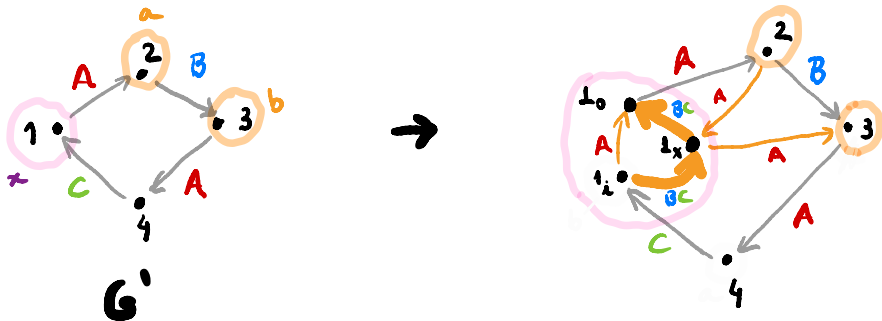
Reduction of this problem to the problem of computing a breakpoint median of three canonical circular genomes **A**, **B** and **C** that have the same relative orientation:

We need to transform $G$ into another directed graph $G''$, such that $G''$ is the union of three hamiltonian cycles (each one representing one input genome of the median problem)

# NP-hardness of unichromosomal breakpoint median

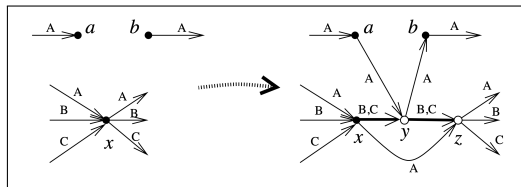Build a modified directed graph $G''$, such that $G''$ is the union of three hamiltonian cycles (each one representing one genome among $\mathbf{A}$, $\mathbf{B}$ and $\mathbf{C}$)



$G''$ has only adjacencies that occur in one or in two genomes

Let $\mathbb{M}$ be a solution to the circular breakpoint median of $\mathbf{A}$, $\mathbf{B}$ and $\mathbf{C}$:

$\mathbb{M}$ contains all adjacencies common to two input genomes and no "new" adjacency

$\updownarrow$

Initial graph $G$ has an hamiltonian cycle

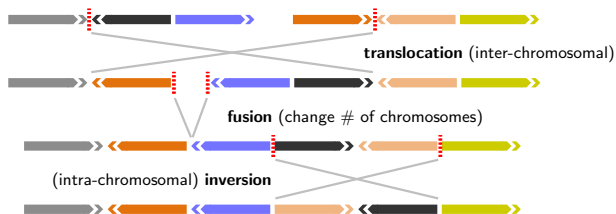# NP-hardness of unichromosomal breakpoint median

# Quiz 2

1 Which of the following statements are true?

A There is a polynomial time algorithm for solving the unichromosomal breakpoint median.

B There cannot be a polynomial time algorithm for solving the unichromosomal breakpoint median.

C The unichromosomal breakpoint median is NP-hard because it can be reduced to the hamiltonian cycle problem.

D The unichromosomal breakpoint median is NP-hard because the hamiltonian cycle problem can be reduced to it.
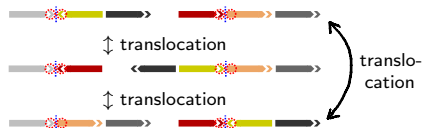
# Double-cut-and-join (DCJ) model

**Double-cut-and-join (DCJ) operation:** two cuts + two joins

- ▶ Cuts the genome twice and rejoins loose ends in a different way.
- ▶ Represents most large-scale genome rearrangements (inversions, translocations, fusions, fissions... )



**translocation** (inter-chromosomal)

**fusion** (change # of chromosomes)

(intra-chromosomal) **inversion**

# The double-cut-and-join (DCJ) operation

Cuts the genome in (at most) 2 positions and rejoins the open ends in a distinct way

# The double-cut-and-join (DCJ) operation

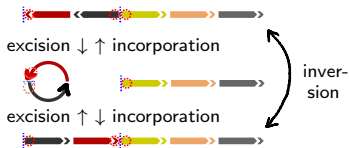Cuts the genome in (at most) 2 positions and rejoins the open ends in a distinct way
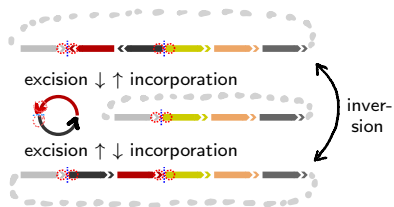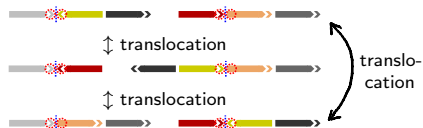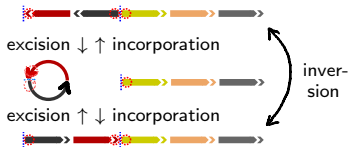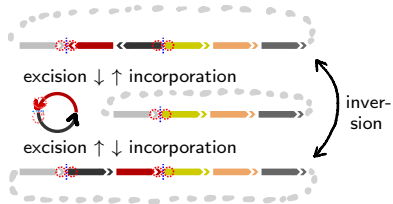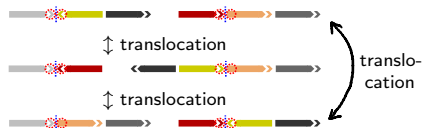
# The double-cut-and-join (DCJ) operation

Cuts the genome in (at most) 2 positions and rejoins the open ends in a distinct way

# DCJ model



**DCJ operation involving two adjacencies**

$$xv + wz$$

$$xz + wv \qquad xw + vz$$

**two possibilities** of rejoining in a different way

**Cases:**

**A.** Each adjacency is in a distinct linear chromosome:

$$[1 \; {}^x_v \; 2 \; 3] \quad [4 \; {}^w_z \; 5 \; 6]$$

reciprocal translocation $\triangle$ reciprocal translocation

$$[1 \; {}^x_z \; 5 \; 6] \quad [4 \; {}^w_v \; 2 \; 3] \quad \text{reciprocal translocation} \quad [1 \; {}^x_w \; \bar{4}] \quad [\bar{3} \; \bar{2} \; {}^v_z \; 5 \; 6]$$

**B.** Both adjacencies are in the same chromosome, or one is in a circular chromosome:

$$([1 \; {}^x_v \; 2 \; 3 \; 4 \; {}^z_w \; 5 \; 6])$$

inversion $\triangle$ excision/integration

$$([1 \; {}^x_z \; \bar{4} \; \bar{3} \; \bar{2} \; {}^v_w \; 5 \; 6]) \quad \text{excision/integration} \quad ([1 \; {}^x_w \; 5 \; 6]) \quad (3 \; 4 \; {}^z_v \; 2)$$

# DCJ model



**DCJ operation involving one adjacency and one telomere**

$$x \;+\; wz$$
$$xz \;+\; w \qquad xw \;+\; z$$

**two possibilities** of rejoining in a different way

**Cases:**

**A.** The adjacency and the telomere are in distinct linear chromosomes:

$$[\,1\ 2\ 3^{x}\,]\quad[\,4^{wz}\ 5\ 6\,]$$

translocation $\qquad$ translocation

$$[\,1\ 2\ 3^{xz}\ 5\ 6\,]\quad[\,4^{w}\,]\quad\text{translocation}\quad[\,1\ 2\ 3^{xw}\ \bar{4}\,]\quad[\,{}^{z}\ 5\ 6\,]$$

**B.** The adjacency is in the same linear chromosome, or in a circular chromosome:

$$[\,1\ 2\ 3\ 4^{zw}\ 5\ 6^{x}\,]$$

inversion $\qquad$ excision/ integration

$$[\,1\ 2\ 3\ 4^{zx}\ \bar{6}\ \bar{5}^{w}\,]\quad\text{excision/}\atop\text{integration}\quad[\,1\ 2\ 3\ 4^{z}\,]\quad(\,6^{xw}\ 5\,)$$

# DCJ model

$$\boxed{x \ + \ z}$$

**DCJ operation
involving one adjacency
or two telomeres**

$\updownarrow$

**one possibility**
of rejoining
in a different way

$$\boxed{xz}$$

**Cases:**

**A.** The adjacency is in a linear chromosome / the telomeres are in two distinct chromosomes:

$$[\,1\ 2\ 3\,{}^x_{\blacktriangledown}{}_{\blacktriangledown}\,]\quad[\,{}_{\blacktriangledown}{}^z_{\blacktriangledown}\,4\ 5\,]$$

fusion $\downarrow\uparrow$ fission

$$[\,1\ 2\ 3\,{}^x_{\blacktriangledown}{}^z_{\blacktriangledown}\,4\ 5\,]\quad[\,{}_{\blacktriangledown\blacktriangledown}\,]$$

**B.** The adjacency is in a circular chromosome / the telomeres are in the same chromosome:

$$[\,{}_{\blacktriangledown}{}^x_{\blacktriangledown}\,1\ 2\ 3\ 4\ 5\,{}^z_{\blacktriangledown}{}_{\blacktriangledown}\,]$$

circularization $\downarrow\uparrow$ linearization

$$(\,2\ 3\ 4\ 5\,{}^z_{\blacktriangledown}{}^x_{\blacktriangledown}\,1\,)\quad[\,{}_{\blacktriangledown\blacktriangledown}\,]$$

# Quiz 3

1 Which transformations can be done with a single DCJ operation?

A $[1\,2\,3]$ $[4\,5]$ $\leftrightarrow$ $[1\,2\,4\,5\,3]$

B $[1\,2\,3]$ $[4\,5]$ $\leftrightarrow$ $[1\,2\,3\,\bar{5}\,\bar{4}]$

C $[1\,2\,3]$ $[4\,5]$ $\leftrightarrow$ $[1\,2\,5]$ $[4\,3]$

D $[1\,2\,3\,4\,5]$ $\leftrightarrow$ $[1\,\bar{4}\,3\,\bar{2}\,5]$

E $[1\,2\,3\,4\,5]$ $\leftrightarrow$ $[1\,2\,\bar{5}\,\bar{4}\,\bar{3}]$

F $[1\,2\,3]$ $(4\,5)$ $\leftrightarrow$ $[1\,2\,4\,5\,3]$

G $[1\,2\,3]$ $(4\,5)$ $\leftrightarrow$ $[1\,2\,5\,4\,3]$

H $(1\,2\,3\,4\,5)$ $\leftrightarrow$ $[3\,4\,5\,1\,2]$

# References

Multichromosomal median and halving problems under different genomic distances

(Eric Tannier, Chunfang Zheng and David Sankoff)

BMC Bioinformatics volume 10, Article number: 120 (2009)


SCJ: A Breakpoint-Like Distance that Simplifies Several Rearrangement Problems

(Pedro Feijão and João Meidanis)

TCBB volume 8 Number: 5 (2011)


The complexity of the breakpoint median problem

(David Bryant)

Tech. Rep. CRM-2579, Centre de recherches mathématiques, Université de Montréal, 1998