

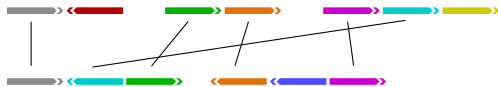
## Topics of today:

1. Recall concepts from lecture 01
2. Single-cut-or-join model, distance and double-distance
3. Formalizing the number of occurrences ( $\phi$ ) of families/adjacencies/telomeres
4. Other problems: median and halving

# Types of genome pairs/sets

## Pair/set of singular genomes:

each family occurs **at most once** in each genome



## Pair/set of balanced genomes:

each family occurs **the same number of times** in each genome



## Singular/duplicated canonical pair:

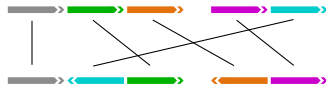
one genome is singular, the other is duplicated and the gene families of both genomes are the same



(whole genome duplication)



## Pair/set of canonical genomes: **singular** and **balanced**



## Definitions / notation (family-based setting)

Given genomes  $\mathbb{G}_1^f, \mathbb{G}_2^f, \dots, \mathbb{G}_k^f$ :

- ▶ Set of **common families** (occurring in each  $\mathbb{G}_i^f$ ):

$$\mathcal{F}_* = \mathcal{F}(\mathbb{G}_1^f) \cap \mathcal{F}(\mathbb{G}_2^f) \cap \dots \cap \mathcal{F}(\mathbb{G}_k^f)$$

- ▶ (Multi)set of annotated **common markers**:

$$\mathcal{G}_* = \mathcal{G}(\mathbb{G}_1^f) \cap \mathcal{G}(\mathbb{G}_2^f) \cap \dots \cap \mathcal{G}(\mathbb{G}_k^f)$$

$$|\mathcal{G}_*| = n$$

### Type

singular:  $\mathcal{F}_* = \mathcal{G}_*$  **(a)**

balanced:  $\mathcal{F}_* = \mathcal{F}(\mathbb{G}_1^f) = \mathcal{F}(\mathbb{G}_2^f) = \dots = \mathcal{F}(\mathbb{G}_k^f)$  **and**  $\mathcal{G}_* = \mathcal{G}(\mathbb{G}_1^f) = \mathcal{G}(\mathbb{G}_2^f) = \dots = \mathcal{G}(\mathbb{G}_k^f)$  **(b)**

canonical: **both (a) and (b)**

## Breakpoint distance

Given genomes  $\mathbb{A}^f$  and  $\mathbb{B}^f$ , let:

- ▶  $\Gamma_\star = \Gamma(\mathbb{A}^f) \cap \Gamma(\mathbb{B}^f)$  be the set of **common adjacencies**

$$|\Gamma_\star| = a$$

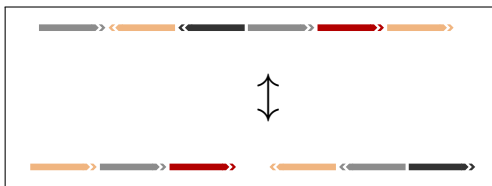
- ▶  $\Theta_\star = \Theta(\mathbb{A}^f) \cap \Theta(\mathbb{B}^f)$  be the set of **common telomeres**

$$|\Theta_\star| = t$$

The breakpoint distance of **canonical genomes**  $\mathbb{A}_\triangleright^f$  and  $\mathbb{B}_\triangleright^f$  is defined to be:

$$d_{\text{BP}}(\mathbb{A}_\triangleright^f, \mathbb{B}_\triangleright^f) = n - a - \frac{t}{2}$$

## Breakpoint distance of balanced genomes

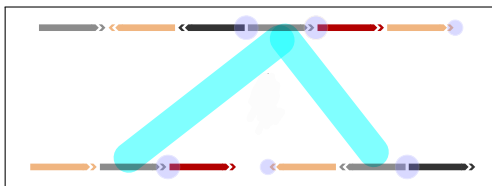


The breakpoint distance of **balanced genomes**  $\mathbb{A}^f$  and  $\mathbb{B}^f$  is:

$$d_{\text{BP}}(\mathbb{A}^f, \mathbb{B}^f) = \min_{f_m} d_{\text{BP}}(\mathbb{A}_{\triangleright}^{f_m}, \mathbb{B}_{\triangleright}^{f_m})$$

where  $f_m$  is any function that produces a maximal matching of  $f$ -families

## Breakpoint distance of balanced genomes



The breakpoint distance of **balanced genomes**  $\mathbb{A}^f$  and  $\mathbb{B}^f$  is:

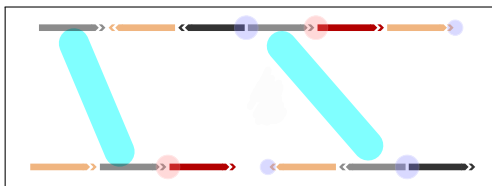
$$d_{BP}(\mathbb{A}^f, \mathbb{B}^f) = \min_{f_m} d_{BP}(\mathbb{A}_{\triangleright}^{f_m}, \mathbb{B}_{\triangleright}^{f_m})$$

where  $f_m$  is any function that produces a maximal matching of  $f$ -families

**Greedy approach:** take all common adjacencies/telomeres:  $|\mathcal{G}_\star| - a - \frac{t}{2} = 6 - 2 - \frac{1}{2} = 3.5$

may lead to **inconsistencies**

## Breakpoint distance of balanced genomes



The breakpoint distance of **balanced genomes**  $\mathbb{A}^f$  and  $\mathbb{B}^f$  is:

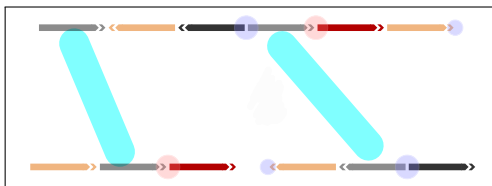
$$d_{BP}(\mathbb{A}^f, \mathbb{B}^f) = \min_{f_m} d_{BP}(\mathbb{A}_{\triangleright}^{f_m}, \mathbb{B}_{\triangleright}^{f_m})$$

where  $f_m$  is any function that produces a maximal matching of  $f$ -families

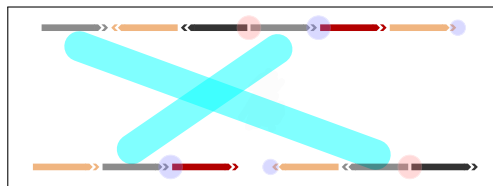
**Greedy approach:** take all common adjacencies/telomeres:  $|\mathcal{G}_\star| - a - \frac{t}{2} = 6 - 2 - \frac{1}{2} = 3.5$

may lead to **inconsistencies**

## Breakpoint distance of balanced genomes



or



The breakpoint distance of **balanced genomes**  $\mathbb{A}^f$  and  $\mathbb{B}^f$  is:

$$d_{BP}(\mathbb{A}^f, \mathbb{B}^f) = \min_{f_m} d_{BP}(\mathbb{A}_{\triangleright}^{f_m}, \mathbb{B}_{\triangleright}^{f_m})$$

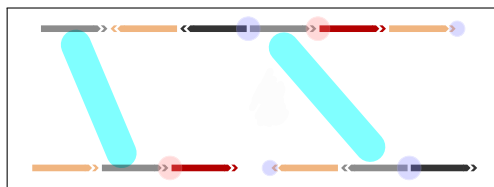
where  $f_m$  is any function that produces a maximal matching of  $f$ -families

**Greedy approach:** take all common adjacencies/telomeres:  $|\mathcal{G}_\star| - a - \frac{t}{2} = 6 - 2 - \frac{1}{2} = 3.5$

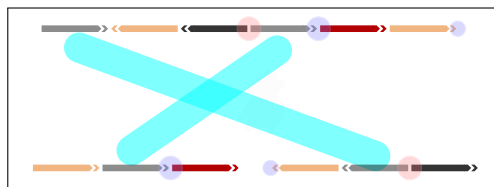
may lead to **inconsistencies**



## Breakpoint distance of balanced genomes



or



The breakpoint distance of **balanced genomes**  $\mathbb{A}^f$  and  $\mathbb{B}^f$  is:

$$d_{\text{BP}}(\mathbb{A}^f, \mathbb{B}^f) = \min_{f_m} d_{\text{BP}}(\mathbb{A}_{\triangleright}^{f_m}, \mathbb{B}_{\triangleright}^{f_m})$$

where  $f_m$  is any function that produces a maximal matching of  $f$ -families

**Greedy approach:** take all common adjacencies/telomeres:  $|\mathcal{G}_*| - a - \frac{t}{2} = 6 - 2 - \frac{1}{2} = 3.5$

may lead to **inconsistencies**

**Correct distance:**  $d_{\text{BP}}(\mathbb{A}^f, \mathbb{B}^f) = 6 - 1 - \frac{1}{2} = 4.5$

The breakpoint distance of balanced genomes is **NP-hard**

[Blin, Chauve and Fertin, 2004: The breakpoint distance for signed sequences]

# Breakpoint double distance

For a given singular genome  $\mathbb{S}_\triangleright^f$ , let  $2 \cdot \mathbb{S}_\triangleright^f$  be the set of doubled genomes derived from  $\mathbb{S}_\triangleright^f$ .

We define:

- ▶  $\mathcal{G}(2 \cdot \mathbb{S}_\triangleright^f) = \mathcal{G}(\mathbb{S}_\triangleright^f) \oplus \mathcal{G}(\mathbb{S}_\triangleright^f)$  : the multiset of **markers** in any doubled genome from the set  $2 \cdot \mathbb{S}_\triangleright^f$
- ▶  $\Gamma(2 \cdot \mathbb{S}_\triangleright^f) = \Gamma(\mathbb{S}_\triangleright^f) \oplus \Gamma(\mathbb{S}_\triangleright^f)$  : the multiset of **adjacencies** in any doubled genome from the set  $2 \cdot \mathbb{S}_\triangleright^f$
- ▶  $\Theta(2 \cdot \mathbb{S}_\triangleright^f) = \Theta(\mathbb{S}_\triangleright^f) \oplus \Theta(\mathbb{S}_\triangleright^f)$  : the multiset of **telomeres** in any doubled genome from the set  $2 \cdot \mathbb{S}_\triangleright^f$

**Breakpoint double distance:**

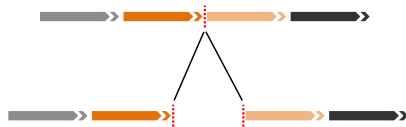
$$\begin{aligned}
 d_{\text{BP}}^2(\mathbb{S}_\triangleright^f, \mathbb{D}_\diamond^f) &= \min_{\mathbb{P}_\boxtimes^f \in 2 \cdot \mathbb{S}_\triangleright^f} d_{\text{BP}}(\mathbb{P}_\boxtimes^f, \mathbb{D}_\diamond^f) \Rightarrow \text{greedy approach} \\
 &\quad \text{is consistent} \\
 &= n' - |\Gamma(\mathbb{P}_\boxtimes^f) \cap \Gamma(\mathbb{D}_\diamond^f)| - \frac{|\Theta(\mathbb{P}_\boxtimes^f) \cap \Theta(\mathbb{D}_\diamond^f)|}{2} \\
 &= 2n - |\Gamma(2 \cdot \mathbb{S}_\triangleright^f) \cap \Gamma(\mathbb{D}_\diamond^f)| - \frac{|\Theta(2 \cdot \mathbb{S}_\triangleright^f) \cap \Theta(\mathbb{D}_\diamond^f)|}{2}
 \end{aligned}$$

$$n = |\mathcal{G}(\mathbb{S}_\triangleright^f)|$$

$$\begin{aligned}
 n' &= |\mathcal{G}(\mathbb{P}_\boxtimes^f) \cap \mathcal{G}(2 \cdot \mathbb{D}_\diamond^f)| \\
 &= |\mathcal{G}(2 \cdot \mathbb{S}_\triangleright^f) \cap \mathcal{G}(2 \cdot \mathbb{D}_\diamond^f)| \\
 &= |\mathcal{G}(2 \cdot \mathbb{S}_\triangleright^f)| \\
 &= 2|\mathcal{G}(\mathbb{S}_\triangleright^f)| \\
 &= 2n
 \end{aligned}$$

# Single-Cut-or-Join (SCJ) model

- ▶ A **cut** is an operation that breaks an adjacency of genome  $\mathbb{G}$  in two telomeres.
- ▶ A **join** is the reverse operation: joins two telomeres of  $\mathbb{G}$  into one adjacency.
- ▶ Any **single cut** or **single join** is a SCJ operation.



A canonical genome  $\mathbb{G}_{\triangleright}^f$  can be represented by its set of adjacencies  $\Gamma(\mathbb{G}_{\triangleright}^f)$   
(the set of telomeres  $\Theta(\mathbb{G}_{\triangleright}^f)$  can be derived from  $\Gamma(\mathbb{G}_{\triangleright}^f)$ )

Then, SCJ operations can be seen as set operations:

- ▶ A cut of an adjacency  $xy$ :  $\Gamma(\mathbb{G}_{\triangleright}^f) \setminus \{xy\}$ .
- ▶ A join of an adjacency  $xy$ :  $\Gamma(\mathbb{G}_{\triangleright}^f) \cup \{xy\}$ .

# SCJ distance and sorting of canonical genomes

The SCJ distance  $d_{\text{SCJ}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f)$  is the minimum number of SCJs that transform  $\Gamma(\mathbb{A}_{\triangleright}^f)$  into  $\Gamma(\mathbb{B}_{\triangleright}^f)$

The only allowed operations are to remove an element from and to include an element in a set

↘ A lower bound is derived from the simple difference between the two given sets:

$$d_{\text{SCJ}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) \geq |\Gamma(\mathbb{A}_{\triangleright}^f) \setminus \Gamma(\mathbb{B}_{\triangleright}^f)| + |\Gamma(\mathbb{B}_{\triangleright}^f) \setminus \Gamma(\mathbb{A}_{\triangleright}^f)|$$

We can achieve this lower bound by ensuring that all adjacencies that must be included are available (the corresponding involved extremities are “free”):

1. First, remove all elements of  $\Gamma(\mathbb{A}_{\triangleright}^f)$  that are not present in  $\Gamma(\mathbb{B}_{\triangleright}^f)$ :

$$\# \text{ of single cut operations} = |\Gamma(\mathbb{A}_{\triangleright}^f) \setminus \Gamma(\mathbb{B}_{\triangleright}^f)|$$

2. Then, include in  $\Gamma(\mathbb{A}_{\triangleright}^f)$  all elements of  $\Gamma(\mathbb{B}_{\triangleright}^f)$  that are not already present in  $\Gamma(\mathbb{A}_{\triangleright}^f)$ :

$$\# \text{ of single join operations} = |\Gamma(\mathbb{B}_{\triangleright}^f) \setminus \Gamma(\mathbb{A}_{\triangleright}^f)|$$

## SCJ distance

$$d_{\text{SCJ}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) = |\Gamma(\mathbb{A}_{\triangleright}^f) \setminus \Gamma(\mathbb{B}_{\triangleright}^f)| + |\Gamma(\mathbb{B}_{\triangleright}^f) \setminus \Gamma(\mathbb{A}_{\triangleright}^f)|$$

# SCJ sorting of $\mathbb{A}_{\triangleright}^f$ into $\mathbb{B}_{\triangleright}^f$

$\Gamma(\mathbb{A}_{\triangleright}^f) =$

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



$\Gamma(\mathbb{B}_{\triangleright}^f) =$

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



# SCJ sorting of $\mathbb{A}_{\triangleright}^f$ into $\mathbb{B}_{\triangleright}^f$

$\Gamma(\mathbb{A}_{\triangleright}^f) =$	$\{1^h 3^h, 3^t 2^h, 2^t 4^t\}$		
$\Gamma(\mathbb{I}'_{\triangleright}^f) =$	$\Gamma(\mathbb{A}_{\triangleright}^f) \setminus \{1^h 3^h\} =$	$\{3^t 2^h, 2^t 4^t\}$	
$\Gamma(\mathbb{I}''_{\triangleright}^f) =$	$\Gamma(\mathbb{I}'_{\triangleright}^f) \setminus \{2^t 4^t\} =$	$\{3^t 2^h\}$	
$\Gamma(\mathbb{I}'''_{\triangleright}^f) =$	$\Gamma(\mathbb{I}''_{\triangleright}^f) \cup \{1^h 2^t\} =$	$\{1^h 2^t, 2^h 3^t\}$	
$\Gamma(\mathbb{B}_{\triangleright}^f) =$	$\Gamma(\mathbb{I}'''_{\triangleright}^f) \cup \{3^h 4^t\} =$	$\{1^h 2^t, 2^h 3^t, 3^h 4^t\}$	

## Alternative formula for the SCJ distance of canonical genomes

$$\begin{aligned}
 d_{\text{SCJ}}(\mathbb{A}_\triangleright^f, \mathbb{B}_\triangleright^f) &= |\Gamma(\mathbb{A}_\triangleright^f) \setminus \Gamma(\mathbb{B}_\triangleright^f)| + |\Gamma(\mathbb{B}_\triangleright^f) \setminus \Gamma(\mathbb{A}_\triangleright^f)| \\
 &= \underbrace{|\Gamma(\mathbb{A}_\triangleright^f)| - |\Gamma(\mathbb{A}_\triangleright^f) \cap \Gamma(\mathbb{B}_\triangleright^f)|}_{\text{blue}} + \underbrace{|\Gamma(\mathbb{B}_\triangleright^f)| - |\Gamma(\mathbb{A}_\triangleright^f) \cap \Gamma(\mathbb{B}_\triangleright^f)|}_{\text{red}} \\
 &= |\Gamma(\mathbb{A}_\triangleright^f)| + |\Gamma(\mathbb{B}_\triangleright^f)| - 2|\Gamma(\mathbb{A}_\triangleright^f) \cap \Gamma(\mathbb{B}_\triangleright^f)| \\
 &= |\Gamma(\mathbb{A}_\triangleright^f)| + |\Gamma(\mathbb{B}_\triangleright^f)| - 2|\Gamma_\star|
 \end{aligned}$$

Note that:  $|\Theta(\mathbb{A}_\triangleright^f)| = 2(n - |\Gamma(\mathbb{A}_\triangleright^f)|) \Rightarrow |\Gamma(\mathbb{A}_\triangleright^f)| = n - \frac{|\Theta(\mathbb{A}_\triangleright^f)|}{2}$ , where  $n = |\mathcal{G}(\mathbb{A}_\triangleright^f)|$

$$\begin{aligned}
 d_{\text{SCJ}}(\mathbb{A}_\triangleright^f, \mathbb{B}_\triangleright^f) &= n - \frac{|\Theta(\mathbb{A}_\triangleright^f)|}{2} + n - \frac{|\Theta(\mathbb{B}_\triangleright^f)|}{2} - 2|\Gamma_\star| \\
 &= 2n - 2a - \frac{|\Theta(\mathbb{A}_\triangleright^f)| + |\Theta(\mathbb{B}_\triangleright^f)|}{2} \\
 &= 2n - 2a - \kappa(\mathbb{A}) - \kappa(\mathbb{B})
 \end{aligned}$$

where  $n = |\mathcal{G}_\star| = |\mathcal{G}(\mathbb{A}_\triangleright^f)| = |\mathcal{G}(\mathbb{B}_\triangleright^f)|$ ,  $a = |\Gamma_\star|$  and  $\kappa(\cdot)$  is the number of linear chromosomes in the respective genome

## Breakpoint distance $\times$ SCJ distance

$$d_{\text{BP}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) = n - a - \frac{t}{2}$$

$$\begin{aligned}d_{\text{SCJ}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) &= 2n - 2a - \kappa(\mathbb{A}) - \kappa(\mathbb{B}) \\&= 2n - 2a - \kappa(\mathbb{A}) - \kappa(\mathbb{B}) - t + t \\&= 2n - 2a - t - \kappa(\mathbb{A}) - \kappa(\mathbb{B}) + t \\&= 2\left(n - a - \frac{t}{2}\right) - \kappa(\mathbb{A}) - \kappa(\mathbb{B}) + t \\&= 2d_{\text{BP}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) - \kappa(\mathbb{A}) - \kappa(\mathbb{B}) + t\end{aligned}$$

Note that:  $t \leq \kappa(\mathbb{A}) + \kappa(\mathbb{B})$

For circular genomes:

$$d_{\text{SCJ}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) = 2d_{\text{BP}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f)$$

In general:

$$d_{\text{BP}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) \leq d_{\text{SCJ}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f) \leq 2d_{\text{BP}}(\mathbb{A}_{\triangleright}^f, \mathbb{B}_{\triangleright}^f)$$



## SCJ double distance

The SCJ distance of **balanced genomes**  $A^f$  and  $B^f$  is:

$$d_{\text{SCJ}}(A^f, B^f) = \min_{f_m} d_{\text{SCJ}}(A_{\triangleright}^{f_m}, B_{\triangleright}^{f_m})$$

where  $f_m$  is any function that produces a maximal matching of the families defined by  $f$

### SCJ double distance:

$$\begin{aligned} d_{\text{SCJ}}^2(\mathbb{S}_{\triangleright}^f, \mathbb{D}_{\diamond}^f) &= \min_{\mathbb{P}_{\bowtie}^f \in 2 \cdot \mathbb{S}_{\triangleright}^f} d_{\text{SCJ}}(\mathbb{P}_{\bowtie}^f, \mathbb{D}_{\diamond}^f) = \text{greedy approach} \\ &= |\Gamma(\mathbb{P}_{\bowtie}^f) \setminus \Gamma(\mathbb{D}_{\diamond}^f)| + |\Gamma(\mathbb{D}_{\diamond}^f) \setminus \Gamma(\mathbb{P}_{\bowtie}^f)| \\ &= |\Gamma(2 \cdot \mathbb{S}_{\triangleright}^f) \setminus \Gamma(\mathbb{D}_{\diamond}^f)| + |\Gamma(\mathbb{D}_{\diamond}^f) \setminus \Gamma(2 \cdot \mathbb{S}_{\triangleright}^f)| \end{aligned}$$

Ex:  $\mathbb{S} = [\bar{2} \ 1 \ \bar{3}]$  and  $\mathbb{D} = [3\bar{1} \ \bar{2} \ 3 \ \bar{1} \ 2]$

# Quiz 1

Given genomes  $\mathbb{G}_1^f = [\bar{2} \bar{1}] [\bar{4} \bar{3} 5]$ ,  $\mathbb{G}_2^f = [1 2 3 4 5]$  and  $\mathbb{G}_3^f = (1 2 3 4) [1 \bar{5} \bar{4} 5 \bar{3} \bar{2}]$ :

1 What is the SCJ distance of  $\mathbb{G}_1^f$  and  $\mathbb{G}_2^f$ ?

A 2

B 2.5

C 3

D 4

2 What is the SCJ double distance of  $\mathbb{G}_2^f$  and  $\mathbb{G}_3^f$ ?

A 6

B 7

C 7.5

D 8

# Occurrences of families

Given a family  $\mathbf{X}$  and a genome  $\mathbb{G}^f$ , let  $\phi(\mathbf{X}, \mathbb{G}^f)$  be the number of occurrences of  $\mathbf{X}$  in  $\mathcal{G}(\mathbb{G}^f)$ .

If genome  $\mathbb{S}_{\triangleright}^f$  is singular, then  $\phi(\mathbf{X}, \mathbb{S}_{\triangleright}^f) = 1$  for each  $\mathbf{X} \in \mathcal{F}(\mathbb{S}_{\triangleright}^f)$ .

If genome  $\mathbb{D}_{\diamond}^f$  is duplicated, then  $\phi(\mathbf{X}, \mathbb{D}_{\diamond}^f) = 2$  for each  $\mathbf{X} \in \mathcal{F}(\mathbb{D}_{\diamond}^f)$ .

If genomes  $\mathbb{S}_{\triangleright}^f$  and  $\mathbb{S}'_{\triangleright}^f$  are canonical, then

$$\mathcal{F}_{\star} = \mathcal{F}(\mathbb{S}_{\triangleright}^f) = \mathcal{F}(\mathbb{S}'_{\triangleright}^f) \text{ and } \phi(\mathbf{X}, \mathbb{S}_{\triangleright}^f) = \phi(\mathbf{X}, \mathbb{S}'_{\triangleright}^f) = 1 \text{ for each } \mathbf{X} \in \mathcal{F}_{\star}.$$

If genomes  $\mathbb{B}_1^f$  and  $\mathbb{B}_2^f$  are balanced, then

$$\mathcal{F}_{\star} = \mathcal{F}(\mathbb{B}_1^f) = \mathcal{F}(\mathbb{B}_2^f) \text{ and } \phi(\mathbf{X}, \mathbb{B}_1^f) = \phi(\mathbf{X}, \mathbb{B}_2^f) \text{ for each } \mathbf{X} \in \mathcal{F}_{\star}.$$

A maximal matching of the genes of two genomes  $\mathbb{A}_1^f$  and  $\mathbb{A}_2^f$  has size:

$$\sum_{\mathbf{X} \in \mathcal{F}(\mathbb{A}_1^f) \cup \mathcal{F}(\mathbb{A}_2^f)} \min\{\phi(\mathbf{X}, \mathbb{A}_1^f), \phi(\mathbf{X}, \mathbb{A}_2^f)\}$$

# Occurrences of adjacencies

Given an adjacency  $xy$  and a genome  $\mathbb{G}^f$ , let  $\phi(xy, \mathbb{G}^f)$  be the number of occurrences of  $xy$  in  $\Gamma(\mathbb{G}^f)$ .

If genome  $\mathbb{S}_{\triangleright}^f$  is singular, then  $\phi(xy, \mathbb{S}_{\triangleright}^f) = \begin{cases} 1, & xy \in \Gamma(\mathbb{S}_{\triangleright}^f), \\ 0, & xy \notin \Gamma(\mathbb{S}_{\triangleright}^f). \end{cases}$

If genome  $\mathbb{D}_{\diamond}^f$  is duplicated, then  $\phi(xy, \mathbb{D}) \in \{0, 1, 2\}$ .

Given an adjacency  $xy$  and a set of  $k$  genomes  $\mathcal{A}^f = \{\mathbb{A}_1^f, \mathbb{A}_2^f, \dots, \mathbb{A}_k^f\}$ , we define:

$$\phi(xy, \mathcal{A}^f) = \phi(xy, \mathbb{A}_{1..k}^f) = \sum_{i=1}^k \phi(xy, \mathbb{A}_i^f)$$

# Occurrences of telomeres

Given a telomere  $x$  and a genome  $\mathbb{G}^f$ , let  $\phi(x, \mathbb{G}^f)$  be the number of occurrences of  $x$  in  $\Theta(\mathbb{G}^f)$ .

If genome  $\mathbb{S}_{\triangleright}^f$  is singular, then  $\phi(x, \mathbb{S}_{\triangleright}^f) = \begin{cases} 1, & x \in \Theta(\mathbb{S}_{\triangleright}^f), \\ 0, & x \notin \Theta(\mathbb{S}_{\triangleright}^f). \end{cases}$

If genome  $\mathbb{D}_{\diamond}^f$  is duplicated, then  $\phi(x, \mathbb{D}) \in \{0, 1, 2\}$ .

Given a telomere  $x$  and a set of  $k$  genomes  $\mathcal{A}^f = \{\mathbb{A}_1^f, \mathbb{A}_2^f, \dots, \mathbb{A}_k^f\}$ , we define:

$$\phi(x, \mathcal{A}^f) = \phi(x, \mathbb{A}_{1..k}^f) = \sum_{i=1}^k \phi(x, \mathbb{A}_i^f)$$

## Quiz 2

1 Let  $\mathbb{D}_\diamond^f = (1234) [1\bar{5}\bar{4}5\bar{3}\bar{2}]$ . Give, respectively, the values of  $\phi(3^h 5^t, \mathbb{D}_\diamond^f)$ ,  $\phi(2^h 3^t, \mathbb{D}_\diamond^f)$ ,  $\phi(4^h 1^t, \mathbb{D}_\diamond^f)$ ,  $\phi(1^t, \mathbb{D}_\diamond^f)$ :

A 1, 1, 2, 0

C 0, 2, 1, 1

B 0, 2, 0, 2

D 1, 2, 0, 2

2 Let  $\mathbb{C}_1^f = [12345]$  and  $\mathbb{C}_2^f = [\bar{2}\bar{1}] [\bar{4}\bar{3}5]$ . Give, respectively, the values of  $\phi(3^h 5^t, \{\mathbb{C}_1^f, \mathbb{C}_2^f\})$ ,  $\phi(2^h 3^t, \{\mathbb{C}_1^f, \mathbb{C}_2^f\})$ ,  $\phi(1^h 2^t, \{\mathbb{C}_1^f, \mathbb{C}_2^f\})$ ,  $\phi(1^t, \{\mathbb{C}_1^f, \mathbb{C}_2^f\})$ :

A 0, 1, 1, 2

C 1, 1, 2, 0

B 0, 1, 2, 2

D 1, 2, 0, 2

## SCJ model - expressing the double distance via adjacency occurrences

$$\begin{aligned}d_{\text{SCJ}}^2(\mathbb{S}_{\triangleright}^f, \mathbb{D}_{\diamond}^f) &= |\Gamma(2 \cdot \mathbb{S}_{\triangleright}^f) \setminus \Gamma(\mathbb{D}_{\diamond}^f)| + |\Gamma(\mathbb{D}_{\diamond}^f) \setminus \Gamma(2 \cdot \mathbb{S}_{\triangleright}^f)| \\&= \sum_{xy \in \Gamma(\mathbb{H}_{\triangleright}^f)} (2 - \phi(xy, \mathbb{D}_{\diamond}^f)) + \sum_{xy \notin \Gamma(\mathbb{H}_{\triangleright}^f)} \phi(xy, \mathbb{D}_{\diamond}^f) \\&= \sum_{xy} [ \phi(xy, \mathbb{H}_{\triangleright}^f) \cdot (2 - \phi(xy, \mathbb{D}_{\diamond}^f)) + (1 - \phi(xy, \mathbb{H}_{\triangleright}^f)) \cdot \phi(xy, \mathbb{D}_{\diamond}^f) ] \\&= \sum_{xy} [ 2 \cdot \phi(xy, \mathbb{H}_{\triangleright}^f) - \phi(xy, \mathbb{H}_{\triangleright}^f) \cdot \phi(xy, \mathbb{D}_{\diamond}^f) + \phi(xy, \mathbb{D}_{\diamond}^f) - \phi(xy, \mathbb{H}_{\triangleright}^f) \cdot \phi(xy, \mathbb{D}_{\diamond}^f) ] \\&= |\Gamma(\mathbb{D}_{\diamond}^f)| + \sum_{xy} [ \phi(xy, \mathbb{H}_{\diamond}^f) (2 - 2 \cdot \phi(xy, \mathbb{D}_{\diamond}^f)) ] \\&= |\Gamma(\mathbb{D}_{\diamond}^f)| + \sum_{xy \in \Gamma(\mathbb{H}_{\diamond}^f)} (2 - 2 \cdot \phi(xy, \mathbb{D}_{\diamond}^f))\end{aligned}$$

# SCJ halving of a duplicated genome

Given a duplicated genome  $\mathbb{D}_\diamond^f$ , find a singular genome  $\mathbb{H}_\triangleright^f$  that minimizes the SCJ double distance:

$$\begin{aligned}d_{\text{SCJ}}^2(\mathbb{H}_\triangleright^f, \mathbb{D}_\diamond^f) &= d_{\text{SCJ}}(2 \cdot \mathbb{H}_\triangleright^f, \mathbb{D}_\diamond^f) = |\Gamma(\mathbb{D}_\diamond^f)| + \sum_{xy \in \Gamma(\mathbb{H}_\triangleright^f)} (2 - 2 \cdot \phi(xy, \mathbb{D}_\diamond^f)) \\ &= |\Gamma(\mathbb{D}_\diamond^f)| + \omega(\mathbb{H}_\triangleright^f)\end{aligned}$$

Since  $|\Gamma(\mathbb{D}_\diamond^f)|$  is given (does not depend on  $\mathbb{H}_\triangleright^f$ ), for minimizing  $d_{\text{SCJ}}^2(\mathbb{H}_\triangleright^f, \mathbb{D}_\diamond^f)$  we need to minimize:

$$\omega(\mathbb{H}_\triangleright^f) = \sum_{xy \in \Gamma(\mathbb{H}_\triangleright^f)} \omega(xy) = \sum_{xy \in \Gamma(\mathbb{H}_\triangleright^f)} (2 - 2 \cdot \phi(xy, \mathbb{D}_\diamond^f)), \text{ where } \omega(xy) = 2 - 2 \cdot \phi(xy, \mathbb{D}_\diamond^f) \in \{-2, 0, +2\}$$

For minimizing  $\omega(\mathbb{H}_\triangleright^f)$ :

- ▶ Do not add to  $\mathbb{H}_\triangleright^f$  any adjacency  $xz$  that have  $\omega(xz) > 0$ :  
this happens when  $\phi(xz, \mathbb{D}_\diamond^f) = 0$  ( $xz$  does not occur in  $\mathbb{D}_\diamond^f$ ).
- ▶ Add to  $\mathbb{H}_\triangleright^f$  any adjacency  $xy$  that have  $\omega(xy) < 0$ :  
this happens when  $\phi(xy, \mathbb{D}_\diamond^f) = 2$  ( $xy$  occurs twice in  $\mathbb{D}_\diamond^f$ ).
- ▶ For  $z \neq y$ :  $\omega(xz) > 0 \Leftrightarrow \omega(xy) < 0$ .
- ▶ Any adjacency  $xy$  with  $\omega(xy) = 0$  (occurs once in  $\mathbb{D}_\diamond^f$ ) is optional (can be added to  $\mathbb{H}_\triangleright^f$  or not).

Solution with the minimum number of adjacencies:  $\Gamma(\mathbb{H}_\triangleright^f) = \{xy : \phi(xy, \mathbb{D}_\diamond^f) = 2\}$

Solution with the maximum number of adjacencies:  $\Gamma(\mathbb{H}_\triangleright^f) = \{xy : \phi(xy, \mathbb{D}_\diamond^f) \geq 1\}$



# SCJ median of three canonical genomes

Given three canonical genomes  $C_1^f$ ,  $C_2^f$  and  $C_3^f$ , find another genome  $M_{\triangleright}^f$  such that:

1.  $M_{\triangleright}^f$  is canonical with  $C_1^f$ ,  $C_2^f$  and  $C_3^f$ ,
2.  $M_{\triangleright}^f$  minimizes the sum  $s_{\text{SCJ}}(M_{\triangleright}^f) = d_{\text{SCJ}}(M_{\triangleright}^f, C_1^f) + d_{\text{SCJ}}(M_{\triangleright}^f, C_2^f) + d_{\text{SCJ}}(M_{\triangleright}^f, C_3^f)$ .

Note that:

$$\begin{aligned} d_{\text{SCJ}}(M_{\triangleright}^f, C_i^f) &= |\Gamma(M_{\triangleright}^f) \setminus \Gamma(C_i^f)| + |\Gamma(C_i^f) \setminus \Gamma(M_{\triangleright}^f)| \\ &= \sum_{xy \in \Gamma(M_{\triangleright}^f)} (1 - \phi(xy, C_i^f)) + \sum_{xy \notin \Gamma(M_{\triangleright}^f)} \phi(xy, C_i^f) \end{aligned}$$

Therefore:

$$\begin{aligned} s_{\text{SCJ}}(M_{\triangleright}^f) &= \sum_{xy \in \Gamma(M_{\triangleright}^f)} [1 - \phi(xy, C_1^f) + (1 - \phi(xy, C_2^f)) + (1 - \phi(xy, C_3^f))] \\ &\quad + \sum_{xy \notin \Gamma(M_{\triangleright}^f)} [\phi(xy, C_1^f) + \phi(xy, C_2^f) + \phi(xy, C_3^f)] \\ &= \sum_{xy \in \Gamma(M_{\triangleright}^f)} (3 - \phi(xy, C_{1..3}^f)) + \sum_{xy \notin \Gamma(M_{\triangleright}^f)} \phi(xy, C_{1..3}^f) \\ &= \sum_{xy} [\phi(xy, M_{\triangleright}^f) \cdot (3 - \phi(xy, C_{1..3}^f)) + (1 - \phi(xy, M_{\triangleright}^f)) \cdot \phi(xy, C_{1..3}^f)] \\ &= \sum_{xy} [3 \cdot \phi(xy, M_{\triangleright}^f) - \phi(xy, M_{\triangleright}^f) \cdot \phi(xy, C_{1..3}^f) + \phi(xy, C_{1..3}^f) - \phi(xy, M_{\triangleright}^f) \cdot \phi(xy, C_{1..3}^f)] \\ &= |\Gamma(C_1^f)| + |\Gamma(C_2^f)| + |\Gamma(C_3^f)| + \sum_{xy} [\phi(xy, M_{\triangleright}^f)(3 - 2 \cdot \phi(xy, C_{1..3}^f))] \\ &= |\Gamma(C_1^f)| + |\Gamma(C_2^f)| + |\Gamma(C_3^f)| + \sum_{xy \in \Gamma(M_{\triangleright}^f)} (3 - 2 \cdot \phi(xy, C_{1..3}^f)) \end{aligned}$$

# SCJ median of three canonical genomes

$$\begin{aligned} s_{\text{SCJ}}(\mathbb{M}_D^f) &= |\Gamma(\mathbb{C}_1^f)| + |\Gamma(\mathbb{C}_2^f)| + |\Gamma(\mathbb{C}_3^f)| + \sum_{xy \in \Gamma(\mathbb{M}_D^f)} (3 - 2 \cdot \phi(xy, \mathbb{C}_{1..3}^f)) \\ &= |\Gamma(\mathbb{C}_1^f)| + |\Gamma(\mathbb{C}_2^f)| + |\Gamma(\mathbb{C}_3^f)| + \omega(\mathbb{M}_D^f) \end{aligned}$$

Since  $|\Gamma(\mathbb{C}_1^f)| + |\Gamma(\mathbb{C}_2^f)| + |\Gamma(\mathbb{C}_3^f)|$  is given (does not depend on  $\mathbb{M}_D^f$ ), for minimizing  $s_{\text{SCJ}}(\mathbb{M}_D^f)$  we need to minimize:

$$\omega(\mathbb{M}_D^f) = \sum_{xy \in \Gamma(\mathbb{M}_D^f)} \omega(xy) = \sum_{xy \in \Gamma(\mathbb{M}_D^f)} (3 - 2 \cdot \phi(xy, \mathbb{C}_{1..3}^f))$$

where  $\omega(xy) = 3 - 2 \cdot \phi(xy, \mathbb{C}_{1..3}^f) \in \{-3, -1, +1, +3\}$ .

For minimizing  $\omega(\mathbb{M}_D^f)$ :

- ▶ Do not add to  $\mathbb{M}_D^f$  any adjacency  $xz$  that have  $\omega(xz) > 0$ :  
this happens when  $\phi(xz, \mathbb{C}_{1..3}^f) \leq 1$  ( $xz$  occurs in at most one genome among  $\mathbb{C}_1^f$ ,  $\mathbb{C}_2^f$  and  $\mathbb{C}_3^f$ ).
- ▶ Add to  $\mathbb{M}_D^f$  any adjacency  $xy$  that have  $\omega(xy) < 0$ :  
this happens when  $\phi(xy, \mathbb{C}_{1..3}^f) \geq 2$  ( $xy$  occurs in at least two genomes among  $\mathbb{C}_1^f$ ,  $\mathbb{C}_2^f$  and  $\mathbb{C}_3^f$ ).
- ▶ For  $z \neq y$ :  $\omega(xz) > 0 \Leftrightarrow \omega(xy) < 0$ .

There is no adjacency  $xy$  with  $\omega(xy) = 0$ . Therefore, the SCJ median problem has a unique solution:

$$\Gamma(\mathbb{M}_D^f) = \{xy : \phi(xy, \mathbb{C}_{1..3}^f) \geq 2\}$$

## SCJ median of three canonical genomes - intuition

Let  $\mathcal{F}_* = \mathcal{G}_* = \{1, 2, 3, \dots, n\}$

and start with  $\mathbb{M}_\triangleright^f = [1] [2] \dots [n]$

$$\Gamma(\mathbb{M}_\triangleright^f) = \emptyset \quad \text{and} \quad s_{\text{SCJ}}(\mathbb{M}_\triangleright^f) = |\Gamma(\mathbb{C}_1^f)| + |\Gamma(\mathbb{C}_2^f)| + |\Gamma(\mathbb{C}_3^f)|$$

Effect of adding an adjacency  $xy$  to  $\mathbb{M}_\triangleright^f$ :

1. If  $xy$  is not present in any genome among  $\{\mathbb{C}_1^f, \mathbb{C}_2^f, \mathbb{C}_3^f\}$ , then  $\Delta s_{\text{SCJ}} = +3$ .
2. If  $xy$  is present in exactly one genome among  $\{\mathbb{C}_1^f, \mathbb{C}_2^f, \mathbb{C}_3^f\}$ , then  $\Delta s_{\text{SCJ}} = +1$ .  
( $\Delta d_{\text{SCJ}}(\mathbb{M}_\triangleright^f, \mathbb{C}_i^f) = -1$ , but  $2 \times \Delta d_{\text{SCJ}}(\mathbb{M}_\triangleright^f, \mathbb{C}_i^f) = +1$ )
3. If  $xy$  is present in exactly two genomes among  $\{\mathbb{C}_1^f, \mathbb{C}_2^f, \mathbb{C}_3^f\}$ , then  $\Delta s_{\text{SCJ}} = -1$ .  
( $2 \times \Delta d_{\text{SCJ}}(\mathbb{M}_\triangleright^f, \mathbb{C}_i^f) = -1$ , but  $\Delta d_{\text{SCJ}}(\mathbb{M}_\triangleright^f, \mathbb{C}_i^f) = +1$ )
4. If  $xy$  is present in all three genomes  $\{\mathbb{C}_1^f, \mathbb{C}_2^f, \mathbb{C}_3^f\}$ , then  $\Delta s_{\text{SCJ}} = -3$ .

# SCJ median of $k$ canonical genomes

Given  $k$  canonical genomes  $C_1^f, C_2^f, \dots, C_k^f$ , find another canonical genome  $M_{\triangleright}^f$  that minimizes the sum:

$$\begin{aligned} s_{\text{SCJ}}(M_{\triangleright}^f) &= d_{\text{SCJ}}(M_{\triangleright}^f, C_1^f) + d_{\text{SCJ}}(M_{\triangleright}^f, C_2^f) + \dots + d_{\text{SCJ}}(M_{\triangleright}^f, C_k^f) \\ &= |\Gamma(C_1^f)| + |\Gamma(C_2^f)| + \dots + |\Gamma(C_k^f)| + \omega(M_{\triangleright}^f) \end{aligned}$$

Analogously to the median of three genomes, we need to minimize:

$$\omega(M_{\triangleright}^f) = \sum_{xy \in \Gamma(M_{\triangleright}^f)} \omega(xy)$$

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

For minimizing  $\omega(M_{\triangleright}^f)$ :

- ▶ Do not add to  $M_{\triangleright}^f$  any adjacency  $xz$  that have  $\omega(xz) > 0$ :  
this happens when  $\phi(xz, C_{1..k}^f) < \frac{k}{2}$  ( $xz$  occurs in less than half of the genomes among  $C_1^f, C_2^f, \dots, C_k^f$ ).
- ▶ Add to  $M_{\triangleright}^f$  any adjacency  $xy$  that have  $\omega(xy) < 0$ :  
this happens when  $\phi(xy, C_{1..k}^f) > \frac{k}{2}$  ( $xy$  occurs in more than half of the genomes among  $C_1^f, C_2^f, \dots, C_k^f$ ).
- ▶ For  $z \neq y$ :  $\omega(xz) > 0 \Leftrightarrow \omega(xy) < 0$ .
- ▶ Any adjacency  $xy$  with  $\omega(xy) = 0$  is optional (can be added to the median or not). If there is no such an adjacency (e.g., if  $k$  is odd), the SCJ median problem has a unique solution.

In general, the following set of adjacencies define a SCJ median of  $k$  genomes:

$$\Gamma(M_{\triangleright}^f) = \left\{ xy : \phi(xy, C_{1..k}^f) > \frac{k}{2} \right\}$$

# SCJ linear median of $k$ canonical linear genomes

1. Compute the general SCJ median  $\mathbb{M}_{\triangleright}^f$  as described above.
2. For each circular chromosome in  $\mathbb{M}_{\triangleright}^f$ , remove one adjacency  $xy$  with smallest weight  $\omega(xy)$ .

# Quiz 3

1 Which of the following statements are true?

- A The SCJ halving is always satisfied by a unique singular genome.
- B The SCJ halving cannot be satisfied by a unique singular genome.
- C The SCJ median of four canonical genomes is always unique.
- D The SCJ median of four canonical genomes cannot be unique.
- E The SCJ median of three canonical genomes is always unique.
- F The SCJ linear median of three canonical linear genomes is always unique.

# 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)