CHAPTER 3

DCJ model of canonical genomes

Let \mathbb{A} and \mathbb{B} be two annotated genomes and note that, if \mathbb{A} and \mathbb{B} form a canonical pair, then $\mathcal{F}(\mathbb{A}) = \mathcal{G}(\mathbb{A}) = \mathcal{G}(\mathbb{B}) = \mathcal{F}(\mathbb{B})$. We then denote by n_* the cardinality of all these sets: $n_* = |\mathcal{F}(\mathbb{A})| = |\mathcal{G}(\mathbb{A})| = |\mathcal{G}(\mathbb{B})| = |\mathcal{F}(\mathbb{B})|$. Recall that, in this case, only DCJ operations are used for sorting one genome into the other, and the corresponding DCJ distance is denoted by $d_{DCJ}(\mathbb{A}, \mathbb{B})$.

3.1 Relational graph of canonical genomes

Finding sorting DCJ operations and computing the DCJ distance between two canonical genomes \mathbb{A} and \mathbb{B} can be achieved with the help of the *relational graph* of \mathbb{A} and \mathbb{B} [18], denoted by $\mathsf{G}_{\mathrm{R}}(\mathbb{A},\mathbb{B}) = (V, E)$, whose sets of vertices and edges are defined as follows:

1. The set of vertices is $V = V(\mathbb{A}) \cup V(\mathbb{B})$, where

 $V(\mathbb{A})$ contains a vertex for each extremity of each marker in $\mathcal{M}(\mathbb{A})$ and $V(\mathbb{B})$ contains a vertex for each extremity of each marker in $\mathcal{M}(\mathbb{B})$.

Each vertex v has an identifier corresponding to the unannotated marker extremity it represents, and a label $\eta(v)$, corresponding to the annotated marker extremity it represents. Note that there are $4n_*$ vertices in $\mathsf{G}_{\mathsf{R}}(\mathbb{A},\mathbb{B})$, $2n_*$ per genome.

2. The set of edges is $E = E_{\mathsf{adj}}(\mathbb{A}) \cup E_{\mathsf{adj}}(\mathbb{B}) \cup E_{\mathsf{ext}}$, where the *adjacency edges* are sets

 $E_{\mathsf{adj}}(\mathbb{A}) = \{uv : u, v \in V(\mathbb{A}) \text{ and } uv \in \mathsf{adj}(\mathbb{A})\} \text{ and } E_{\mathsf{adj}}(\mathbb{B}) = \{uv : u, v \in V(\mathbb{B}) \text{ and } uv \in \mathsf{adj}(\mathbb{B})\},\$

and the set of *extremity edges* (whose cardinality is $2n_*$) is

$$E_{\mathsf{ext}} = \{ uv : u \in V(\mathbb{A}) \text{ and } v \in V(\mathbb{B}) \text{ and } \eta(u) = \eta(v) \}.$$

Since any vertex in $G_R(\mathbb{A}, \mathbb{B})$ has exactly one extremity edge and at most one adjacency edge, its degree is one or two. Therefore, $G_R(\mathbb{A}, \mathbb{B})$ is a collection of paths and cycles. A vertex (marker extremity) that has no adjacency edge corresponds to a telomere and is therefore also called *telomere*. Each cutpoint of each genome is represented in the graph either as an adjacency edge or as a telomere. Recall that the number of cutpoints in \mathbb{A} (respectively \mathbb{B}) is $n_* + \kappa(\mathbb{A})$ (respectively $n_* + \kappa(\mathbb{B})$).

Each connected component of the graph alternates between extremity edges and cutpoints, and we define the *length* of a component Γ to be the number of extremity edges in Γ . An *i-cycle* and an *i-path* denote respectively a cycle and a path of length *i*. Note that all cycles have even length, while paths start and end with extremity edges and can have even or odd lengths, called *even* and *odd paths* respectively.

A cycle can be simply denoted by \mathbb{O} . An odd path has one endpoint in a telomere from \mathbb{A} and the other endpoint in a telomere from \mathbb{B} and is called an $\mathbb{A}\mathbb{B}$ -path, simply denoted by $\mathbb{A}\mathbb{B}$. Even paths have either both endpoints in \mathbb{A} , being an $\mathbb{A}\mathbb{A}$ -path, simply denoted by $\mathbb{A}\mathbb{A}$, or both endpoints in \mathbb{B} , being a $\mathbb{B}\mathbb{B}$ -path, simply denoted by $\mathbb{B}\mathbb{B}$. Even paths can also be called *unbalanced paths*, while odd paths are also called *balanced paths*. Let the sets of cycles, $\mathbb{A}\mathbb{B}$ -, $\mathbb{A}\mathbb{A}$ - and $\mathbb{B}\mathbb{B}$ -paths be respectively denoted by \mathbb{C} , $\mathcal{P}_{\mathbb{A}\mathbb{B}}$, $\mathcal{P}_{\mathbb{A}\mathbb{A}}$ and $\mathcal{P}_{\mathbb{B}\mathbb{B}}$. Now let $\Upsilon(\Gamma)$ give the type of component Γ . For example, if Γ is a cycle, then $\Upsilon(\Gamma) = \mathbb{O}$. We can then explicitly write the above mentioned sets as:

$$\begin{split} \mathfrak{C} &= \{ \Gamma \mid \Gamma \in \mathsf{G}_{\mathrm{R}}(\mathbb{A},\mathbb{B}) \text{ and } \Upsilon(\Gamma) = \mathbb{O} \} \;, \\ \mathfrak{P}_{\mathbb{A}\mathbb{B}} &= \{ \Gamma \mid \Gamma \in \mathsf{G}_{\mathrm{R}}(\mathbb{A},\mathbb{B}) \text{ and } \Upsilon(\Gamma) = \mathbb{A}\mathbb{B} \} \;, \\ \mathfrak{P}_{\mathbb{A}\mathbb{A}} &= \{ \Gamma \mid \Gamma \in \mathsf{G}_{\mathrm{R}}(\mathbb{A},\mathbb{B}) \text{ and } \Upsilon(\Gamma) = \mathbb{A}\mathbb{A} \} \; \text{ and} \\ \mathfrak{P}_{\mathbb{B}\mathbb{B}} &= \{ \Gamma \mid \Gamma \in \mathsf{G}_{\mathrm{R}}(\mathbb{A},\mathbb{B}) \text{ and } \Upsilon(\Gamma) = \mathbb{B}\mathbb{B} \}. \end{split}$$

Recall that $\kappa(\mathbb{A})$ and $\kappa(\mathbb{B})$ are the numbers of linear chromosomes in genomes \mathbb{A} and \mathbb{B} . The endpoints of paths and chromosomes are the same telomeres, therefore we have $\kappa(\mathbb{A}) + \kappa(\mathbb{B}) = |\mathcal{P}_{\mathbb{A}\mathbb{B}}| + |\mathcal{P}_{\mathbb{A}\mathbb{A}}| + |\mathcal{P}_{\mathbb{B}\mathbb{B}}|$. Furthermore, the numbers of telomeres in each genome are even. Since each $\mathbb{A}\mathbb{A}$ - or $\mathbb{B}\mathbb{B}$ -path takes either zero or two telomeres per genome and each $\mathbb{A}\mathbb{B}$ -path takes one telomere per genome, the number of $\mathbb{A}\mathbb{B}$ -paths must be even.

Related graphs. As illustrated in Figure 3.1, the relational graph has the same properties of two simpler graphs that were proposed earlier:

1. The first is the so-called *breakpoint graph*, originally proposed in the seminal studies of the *inversion sorting and distance* [7]. It can be derived from the relational graph by contracting each extremity edge e of $G_R(\mathbb{A}, \mathbb{B}) = (V, E)$ and assigning to the resulting single vertex the common annotation of the vertices that are connected by e. In the breakpoint graph there are only adjacency edges. Furthermore, cycles also have even length, while \mathbb{AB} -paths are even and \mathbb{AA} - and \mathbb{BB} -paths are odd.



- **Figure 3.1:** For a canonical pair formed by multilinear genome $\mathbb{A} = \{ [1 5 3], [\overline{2} \overline{4}] \}$ and unilinear genome $\mathbb{B} = \{ [1 2 3 4 5] \}$, where $n_* = 5$, we represent the relational graph (in the middle) surrounded by the breakpoint graph (top) and by the adjacency graph (bottom). Note that the number of vertices in the breakpoint graph and the numbers of edges in both relational and adjacency graphs are equal to $2n_*$. In all graphs we have a (blue) 4-cycle, a (red) AA-path and two AB-paths.
- 2. The second is the so-called *adjacency graph*, which is bipartite and was originally proposed in the formalization of the *DCJ sorting and distance* [10]. It can be derived from the relational graph by contracting each adjacency edge *a* of $G_R(\mathbb{A}, \mathbb{B}) = (V, E)$, concatenating in the label of the resulting single vertex the annotations of the vertices that are connected by *a*. In other words, the vertices of the adjacency graph are the adjacencies and telomeres of \mathbb{A} and \mathbb{B} and all edges are extremity edges. Similarly to the relational graph, in the adjacency graph cycles have even length, \mathbb{AB} -paths are odd and \mathbb{AA} and \mathbb{BB} -paths are even.

Relational graph of sorted and unsorted genomes. The smallest components that can occur in $G_{\mathbb{R}}(\mathbb{A}, \mathbb{B})$ are 2-cycles and $(\mathbb{A}\mathbb{B})$ 1-paths, denoted *short components*. A cycle whose length is greater than 2 or a path whose length is greater than 1 is called a *long component*. When canonical genomes \mathbb{A} and \mathbb{B} are identical (or *sorted*), their relational graph is a collection of short components: identical genomes have the same sets of adjacencies and telomeres,

and each common adjacency corresponds to a 2-cycle while each common telomere corresponds to a 1-path in $G_R(\mathbb{A}, \mathbb{B})$. Recall that the length of a component corresponds to its number of extremity edges and that $G_R(\mathbb{A}, \mathbb{B})$ has $2n_*$ extremity edges. Therefore, for sorted genomes we have $2n_* = 2|\mathbb{C}| + |\mathcal{P}_{\mathbb{A}\mathbb{B}}|$ and, consequently, $n_* = |\mathbb{C}| + \frac{|\mathcal{P}_{\mathbb{A}\mathbb{B}}|}{2}$. Otherwise, when canonical genomes \mathbb{A} and \mathbb{B} are distinct (or *unsorted*), their relational graph contains at least one long component. Therefore, in this case $n_* > |\mathbb{C}| + \frac{|\mathcal{P}_{\mathbb{A}\mathbb{B}}|}{2}$. With these observations we can already see that the DCJ operations that transform \mathbb{A} into \mathbb{B} must increase the numbers of cycles and/or of $\mathbb{A}\mathbb{B}$ -paths in $G_R(\mathbb{A}, \mathbb{B})$. In the following we will present the results from Bergeron *et al.* [10], explaining how this can be achieved.

3.2 Types of DCJ operation with respect to the relational graph

Note that, with respect to its effect on the relational graph, a DCJ ρ cuts one or two components, and rejoins the open ends to transform them into one or two new components. A DCJ operation ρ is said to be *internal* (INT) to a single component Γ , when ρ cuts only at cutpoint(s) that are in Γ . The result of an internal DCJ ρ can be one component (distinct from Γ) or two components. In contrast, A DCJ operation ρ is said to be a *recombination* (REC) when ρ cuts at cutpoints of two distinct components Γ and Γ' . The result of a recombination ρ can be either a single component or two components (distinct from Γ and Γ').

In order to describe all possible DCJ operations, we adopt the following notation to represent the types of components and their lengths:

- $\mathbb{O}^{\langle i \rangle}$: even cycle with length $i \in \{2, 4, 6, \ldots\}$;
- $\mathbb{AB}^{\langle i \rangle}$: balanced path with length $i \in \{1, 3, 5, \ldots\};$
- $\mathbb{AA}^{\langle i \rangle}$: unbalanced \mathbb{AA} -path with length $i \in \{2, 4, \ldots\}$;
- $\mathbb{BB}^{\langle i \rangle}$: unbalanced \mathbb{BB} -path with length $i \in \{2, 4, \ldots\}$;
- $\Gamma^{\langle i \rangle}$: component of any type with length $i \geq 1$.

The possible types of DCJ operation applied on cutpoints of genome \mathbb{A} are described in Table 3.1 [10]. Note that each DCJ operation must be of one of three types:

Gaining DCJ. Either increases $|\mathcal{C}|$ by one or increases $|\mathcal{P}_{AB}|$ by two.

Neutral DCJ. Does not change the cardinalities of the sets \mathcal{C} and \mathcal{P}_{AB} .

Losing DCJ. Either decreases $|\mathcal{C}|$ by one or decreases $|\mathcal{P}_{AB}|$ by two.

	$\operatorname{component}(s)$	DCJ	component(s)		affects
	[state]		[state]		
1.	$\frac{\Gamma^{\langle i+j\rangle}}{\left[{\rm \widehat{P}}^4 / {\rm \widehat{P}}^3 \right]}$	gaining (INT) losing (REC)	$ \begin{bmatrix} \mathbb{O}^{\langle i \rangle} \\ \Gamma^{\langle j \rangle} \end{bmatrix} \\ [\mathbb{O}^4 / \mathbb{O}^3] $	with $\begin{cases} \imath \in \{2,4,6,\ldots\} \\ \jmath \ge 1 \\ \Upsilon(\Gamma^{\langle \imath+\jmath\rangle}) = \Upsilon(\Gamma^{\langle \jmath\rangle}) \end{cases}$	$\Delta \mathcal{C} $ by ± 1
2.	$\frac{\Gamma^{\langle i+j\rangle}}{[{\bf \hat{r}}^4/{\bf \hat{r}}^3]}$	neutral (INT)	$\frac{\check{\Gamma}^{\langle i+j\rangle}}{[{\bf \hat{r}}^4/{\bf \hat{r}}^3]}$	with $\begin{cases} \imath \in \{2, 4, 6, \ldots\} \\ \jmath \ge 1 \\ \Gamma^{\langle \imath + \jmath \rangle}) \neq \check{\Gamma}^{\langle \imath + \jmath \rangle} \\ \Upsilon(\Gamma^{\langle \imath + \jmath \rangle}) = \Upsilon(\check{\Gamma}^{\langle \imath + \jmath \rangle}) \end{cases}$	
3.	$\frac{\mathbb{AA}^{\langle i+j\rangle}}{\mathbb{BB}^{\langle i'+j'\rangle}}$ $[\mathbf{\widehat{\Gamma}}^4/\mathbf{\widehat{\Gamma}}^3]$	gaining (INT)	$ \begin{array}{ c c c c c } \mathbb{AB}^{\langle i+i'\rangle} \\ \mathbb{AB}^{\langle j+j'\rangle} \\ \hline \hline & \hline &$	with $\begin{cases} \imath \in \{0, 2, 4, \ldots\} \\ \jmath \in \{2, 4, 6, \ldots\} \\ \imath', \jmath' \in \{1, 3, 5, \ldots\} \end{cases}$	$\Delta \mathcal{P}_{\mathbb{AB}} $ by ± 2
4.	$ \begin{array}{c} \mathbb{AB}^{\langle i+j'\rangle} \\ \mathbb{AB}^{\langle j+i'\rangle} \end{array} \\ \hline \left[\left(\mathbf{\hat{r}}^{4} / \left(\mathbf{\hat{r}} \right)^{3} \right] \end{array} \end{array} $	neutral (REC)	$ \begin{array}{ c c c c c } \mathbb{AB}^{\langle i+i'\rangle} \\ \mathbb{AB}^{\langle j+j'\rangle} \\ \hline [\mathbb{C}^4/\mathbb{C}^3] \end{array} \end{array} $	with $\begin{cases} \text{all } \mathbb{AB}\text{-paths distinct} \\ i \in \{0, 2, 4, \ldots\} \\ j \in \{2, 4, 6, \ldots\} \\ i', j' \in \{1, 3, 5, \ldots\} \end{cases}$	
5.	$ \begin{array}{c} \mathbb{A}\mathbb{A}^{\langle i+j\rangle} \\ \mathbb{A}\mathbb{B}^{\langle i'+j'\rangle} \end{array} \\ [\widehat{\mathbb{C}}^4 / \widehat{\mathbb{T}}^3] \end{array} $	neutral (REC)	$ \begin{array}{c} \mathbb{A}\mathbb{A}^{\langle i+i'\rangle} \\ \mathbb{A}\mathbb{B}^{\langle j+j'\rangle} \end{array} \\ [\widehat{\mathbf{\Gamma}}^4 / \widehat{\mathbf{\Gamma}}^3] \end{array} $	with $\begin{cases} \text{all paths distinct} \\ i, i' \in \{2, 4, 6, \ldots\} \\ j \in \{0, 2, 4, \ldots\} \\ j' \in \{1, 3, 5, \ldots\} \end{cases}$	
6.	$\frac{\mathbb{BB}^{\langle i+j\rangle}}{\mathbb{AB}^{\langle i'+j'\rangle}}$ $[\mathbb{C}^4/\mathbb{C}^3]$	neutral (REC)	$ \begin{array}{c} \mathbb{BB}^{\langle i+j'\rangle} \\ \mathbb{AB}^{\langle i'+j\rangle} \end{array} \\ \hline \left(\mathbb{C}^4 / \mathbb{C}^3 \right] \end{array} $	with $\begin{cases} \text{all paths distinct} \\ \imath, \jmath, \jmath' \in \{1, 3, 5, \ldots\} \\ \imath' \in \{0, 2, 4, \ldots\} \end{cases}$	
7.	$\frac{\mathbb{AA}^{\langle i+j'\rangle}}{\mathbb{AA}^{\langle j+i'\rangle}}$ $\boxed{\left[\mathbb{C}^{4}/\mathbb{C}^{3}\right]}$	neutral (REC)	$ \begin{array}{c} \mathbb{A}\mathbb{A}^{\langle i+i'\rangle} \\ \mathbb{A}\mathbb{A}^{\langle j+j'\rangle} \end{array} \\ [\mathbb{D}^4 / \mathbb{D}^3] \end{array} $	with $\begin{cases} \text{all } \mathbb{A}\text{-paths distinct} \\ i \in \{0, 2, 4, \ldots\} \\ j, i', j' \in \{2, 4, 6, \ldots\} \end{cases}$	
8.	$ \begin{array}{c} \mathbb{BB}^{\langle i+j'\rangle} \\ \mathbb{BB}^{\langle j+i'\rangle} \end{array} \\ \hline [\mathbb{C}^4] \end{array} $	neutral (REC)	$ \begin{array}{c} \mathbb{BB}^{\langle i+i'\rangle} \\ \mathbb{BB}^{\langle j+j'\rangle} \end{array} \\ \hline [(\mathbf{\hat{r}}^4] \end{array} $	with $\begin{cases} \text{all } \mathbb{B}\mathbb{B}\text{-paths distinct} \\ \imath, \jmath, \imath', \jmath' \in \{1, 3, 5, \ldots\} \end{cases}$	
9.		gaining (INT)	$\boxed{\mathbb{O}^{\langle i \rangle}}{[\mathbb{f})^{2 1}]}$	with $i \in \{2, 4, 6,\}$	$\Delta \mathcal{C} $ by ± 1
10.	$\boxed{\mathbb{BB}^{\langle i+j\rangle}}$ $[\mathbb{t}^{2 1}]$	gaining (REC)	$ \begin{bmatrix} \mathbb{AB}^{\langle i \rangle} \\ \mathbb{AB}^{\langle j \rangle} \end{bmatrix} $ [P ^{2 2}]	with $i, j \in \{1, 3, 5,\}$	$\Delta \mathcal{P}_{\mathbb{AB}} $ by ± 2
11.	$\boxed{\mathbb{A}\mathbb{A}^{\langle i+j\rangle}}$ $[\mathbb{P}^{2 1}]$	neutral (INT)	$ \begin{bmatrix} \mathbb{A} \mathbb{A}^{\langle i \rangle} \\ \mathbb{A} \mathbb{A}^{\langle j \rangle} \end{bmatrix} \\ [(\widehat{\mathbf{I}})^{2 2}] $	with $i, j \in \{2, 4, 6,\}$	
12.	$\mathbb{AB}^{\langle i+j\rangle}$ $[\mathbb{T}^{2 1}]$	neutral (INT)	$ \begin{bmatrix} \mathbb{A}\mathbb{A}^{\langle i \rangle} \\ \mathbb{A}\mathbb{B}^{\langle j \rangle} \end{bmatrix} \\ [\mathbb{P}^{2 2}] $	with $\begin{cases} i \in \{2, 4, 6 \dots\} \\ j \in \{1, 3, 5, \dots\} \end{cases}$	

Table 3.1: Effects of DCJ operations (applied on cutpoints of \mathbb{A}) on the relational graph

3.3 DCJ distance formula

Gaining DCJ operations lead to a "sorted" graph. In a sorting procedure it is always possible to find a gaining DCJ at each step, therefore any optimal operation is gaining [10]. This can be verified simply by looking at Table 3.1: it is obvious that an unsorted graph will have at least one of the situations leading to the gaining operations shown in lines 1, 3, 9 and 10.

Since gaining operations are those that produce the maximum possible increase of cycles or AB-paths, the following theorem holds.

Theorem 1 ([10]) Given a canonical pair of genomes \mathbb{A} and \mathbb{B} , their DCJ distance is

$$\mathrm{d}_{\mathrm{DCJ}}(\mathbb{A},\mathbb{B}) = n_* - \left(|\mathcal{C}| + \frac{|\mathcal{P}_{\mathbb{A}\mathbb{B}}|}{2} \right).$$

3.4 Triangular inequality

Given any three canonical genomes \mathbb{A} , \mathbb{B} and \mathbb{C} , consider without loss of generality that $d_{DCJ}(\mathbb{A}, \mathbb{B}) \geq d_{DCJ}(\mathbb{A}, \mathbb{C})$ and $d_{DCJ}(\mathbb{A}, \mathbb{B}) \geq d_{DCJ}(\mathbb{B}, \mathbb{C})$. Then the triangular inequality is the property that guarantees that $d_{DCJ}(\mathbb{A}, \mathbb{B}) \leq d_{DCJ}(\mathbb{A}, \mathbb{C}) + d_{DCJ}(\mathbb{B}, \mathbb{C})$. It obviously holds for the DCJ distance: by combining a sorting scenario from \mathbb{A} to \mathbb{C} with $d_{DCJ}(\mathbb{A}, \mathbb{C})$ steps and a sorting scenario from \mathbb{C} to \mathbb{B} with $d_{DCJ}(\mathbb{B}, \mathbb{C})$ steps we trivially get a sorting scenario from \mathbb{A} to \mathbb{B} with $d_{DCJ}(\mathbb{A}, \mathbb{C}) + d_{DCJ}(\mathbb{B}, \mathbb{C})$ steps. Therefore, it is clear that $d_{DCJ}(\mathbb{A}, \mathbb{B})$ cannot be greater than $d_{DCJ}(\mathbb{A}, \mathbb{C}) + d_{DCJ}(\mathbb{B}, \mathbb{C})$, otherwise it would contradict the fact that it corresponds to the length of a most parsimonious sorting scenario.

3.5 On DCJ sorting

Proposition 1 For any DCJ-state of type $(\mathbf{r})^4$ or $(\mathbf{r})^3$ or $(\mathbf{r})^{2|2}$ whose cutpoints are in genome \mathbb{A} and belong to the same long component Γ of $G_R(\mathbb{A}, \mathbb{B})$, there is one, and only one (internal) gaining DCJ operation.

Proof: For each pair of A-cutpoints, such that at most one is a telomere, we know that there are two different DCJ operations. When the cutpoints belong to the same component Γ , one of the two operations simply inverts a fragment, not changing the structure of Γ , and is therefore neutral (Table 3.1, line 2.). The second operation is a gaining DCJ that splits Γ into a cycle and a smaller component of the same type as Γ , increasing the number of cycles (Table 3.1, line 1.). This includes all pairs of A-cutpoints in cycles, AB-paths and BB-paths, and all pairs of A-cutpoints in AA-paths excluding the case where the two cutpoints are telomeres. For the particular case where the two cutpoints are telomeres of an AA-path forming a DCJ-state of type $(\hat{\Gamma}^{2|2})$, there is only one (internal) DCJ operation, and this operation is gaining (Table 3.1, line 9).

Proposition 1 guarantees that there is a gaining DCJ operation for each DCJ-state of type $(\hat{r}^4 \text{ or } (\hat{r})^3 \text{ or } (\hat{r})^{2|2}$ internal to any long component of $G_R(\mathbb{A}, \mathbb{B})$. In particular, this DCJ-state can be a pair of cutpoints directly connected to an adjacency α in genome \mathbb{B} . If the result

of this particular cut is to join such that the adjacency α is created in \mathbb{A} , a new 2-cycle appears in the relational graph, meaning that a DCJ operation that creates an adjacency of genome \mathbb{B} in genome \mathbb{A} is gaining.

Corollary 1 Let α be an adjacency of \mathbb{B} that is not present in \mathbb{A} . The DCJ operation that reconstructs α in \mathbb{A} is gaining.

Once all adjacencies of \mathbb{B} are reconstructed in \mathbb{A} , the only long components that can exist in the relational graph are $\mathbb{B}\mathbb{B}$ -paths of length 2. Any $\mathbb{B}\mathbb{B}$ -path of length 2 can be split into two $\mathbb{A}\mathbb{B}$ -paths with a gaining DCJ (Table 3.1 line 10):

Corollary 2 For any DCJ-state of type $(\hat{r})^{2|1}$ whose cutpoint is an A-adjacency of a BB-path there is only one (internal) DCJ operation, and this operation is gaining.

The results above give a simple greedy algorithm to find one optimal sequence of DCJ operations (exclusively composed of gaining DCJs) to sort \mathbb{A} into \mathbb{B} consisting of simply reconstructing each adjacency of \mathbb{B} that is not in \mathbb{A} , followed by reconstructing each telomere of \mathbb{B} that is not in \mathbb{A} [10].

3.6 Complexity of DCJ distance and greedy sorting

The cutpoints of genome A with n_* markers can be obtained by reading the chromosomes of A once and storing the cutpoints in a vector of length $n_* + \kappa(A)$. Simultaneously, we can store in another vector of length $2n_*$ a pointer to the cutpoint to which each marker extremity belongs. The cutpoints and pointers of genome B can be obtained similarly. Since $\kappa(A) \leq n_*$ and $\kappa(B) \leq n_*$, this procedure takes $O(n_*)$ time and space.

These four vectors are an implict representation of the relational graph: by navigating on them we can easily obtain the cycles and paths of the graph, still in $O(n_*)$.

Finally, for greedily sorting we only need to visit the adjacencies of genome \mathbb{B} followed by the telomeres of genome \mathbb{B} and reconstruct one by one in genome \mathbb{A} . Since each of these reconstructions requires constant time access to the vetcors of \mathbb{A} , the greedy sorting can be easily done in $O(n_*)$ time and space [10].

3.7 Capping the canonical relational graph optimally

The paths of the relational graph can be converted into cycles with a technique called *capping* [55]. It consists of modifying the graph by adding *artificial* extremities, called *cap extremities*, that link all paths into cycles. Each telomere must be connected to a distinct cap extremity by an additional *semi-artificial adjacency edge*. Furthermore, two cap extremities that are in the same genome and not connected to telomeres can be connected to each other by an *artificial adjacency edge*. Then, each cap extremity in genome \mathbb{A} must be connected to a distinct cap extremity in genome \mathbb{B} by a cap extremity edge. Therefore, the capping always connect a telomere in genome \mathbb{A} to a telomere in genome \mathbb{B} via two cap extremities.

Since the number of telomeres in each genome is even, the number of cap extremities is also even. Each pair of cap extremities compose an "artificial" marker called *cap*.

There many ways of performing a capping. A procedure for obtaining an *optimal capping*, that preserves the DCJ distance, is given in the following.

Optimal capping of \mathbb{AB} -**paths.** For each \mathbb{AB} -path whose telomeres are $\gamma_{\mathbb{A}}$ and $\gamma_{\mathbb{B}}$, add cap extremity vertices $\circ_{\mathbb{A}}$ and $\circ_{\mathbb{B}}$ and connect with (semi-artificial) adjacency edges $\circ_{\mathbb{A}}$ to $\gamma_{\mathbb{A}}$ and $\circ_{\mathbb{B}}$ to $\gamma_{\mathbb{B}}$. Furthermore, connect with a (cap) extremity edge the vertices $\circ_{\mathbb{A}}$ to $\circ_{\mathbb{B}}$. Note that this removes one \mathbb{AB} -path, adds one cycle, but also adds half a cap to the set of markers. Since the number of \mathbb{AB} -paths is even, each pair of capped \mathbb{AB} -paths removes two \mathbb{AB} -paths, adds two cycles, but also adds one cap to the set of markers. We denote the capping of path \mathbb{AB} by $\zeta(\mathbb{AB})$.

Optimal capping of AA- and BB-paths. Let an AA-path have telomeres $\gamma_{\mathbb{A}}^1$ and $\gamma_{\mathbb{A}}^2$ and a BB-path have telomeres $\gamma_{\mathbb{B}}^1$ and $\gamma_{\mathbb{B}}^2$. These paths can be optimally linked together into a single cycle as follows. Add cap extremity vertices $\circ_{\mathbb{A}}^1$, $\circ_{\mathbb{A}}^2$, $\circ_{\mathbb{B}}^1$ and $\circ_{\mathbb{B}}^2$ and connect with (semiartificial) adjacency edges $\circ_{\mathbb{A}}^1$ to $\gamma_{\mathbb{A}}^1$, $\circ_{\mathbb{A}}^2$ to $\gamma_{\mathbb{B}}^2$, $\circ_{\mathbb{B}}^1$ to $\gamma_{\mathbb{B}}^1$, and $\circ_{\mathbb{B}}^2$ for $\gamma_{\mathbb{B}}^2$. Furthermore, connect with (cap) extremity edges $\circ_{\mathbb{A}}^1$ to $\circ_{\mathbb{B}}^1$ and $\circ_{\mathbb{A}}^2$ to $\circ_{\mathbb{B}}^2$. Note that this adds one cycle, but also adds one *cap* to the set of markers. We denote the capping of paths AA and BB together by $\zeta(AA, BB)$. Alternatively, each AA- or BB-path can also be separately capped into a cycle, and this alternative is also optimal. Again, let an AA-path have telomeres $\gamma_{\mathbb{A}}^1$ and $\gamma_{\mathbb{A}}^2$. Add cap extremity vertices $\circ_{\mathbb{A}}^1$, $\circ_{\mathbb{A}}^2$, $\circ_{\mathbb{A}}^2$ and $\circ_{\mathbb{B}}^2$ and connect with (semi-artificial) adjacency edges $\circ_{\mathbb{A}}^1$ to $\gamma_{\mathbb{A}}^1$ and $\circ_{\mathbb{A}}^2$ to $\gamma_{\mathbb{A}}^2$. Furthermore, connect with an (artificial) adjacency edge the two cap vertices $\circ_{\mathbb{B}}^1$ and $\circ_{\mathbb{A}}^2$ and $\circ_{\mathbb{B}}^2$ and connect with an $\circ_{\mathbb{A}}^2$ to $\circ_{\mathbb{B}}^2$. Note that this adds one cycle, but also adds one *cap* to the set of markers. Capping a BB-path separately can be done analogously. We denote the cappings of path AA (respectively BB) by $\zeta(AA, \mathring{\alpha}_{\mathbb{B}})$ (respectively by $\zeta(BB, \mathring{\alpha}_{\mathbb{A}})$).

Table 3.2 summarizes the effect of capping the paths of the graph as described above. Each line of the table has neutral effect on the distance, being therefore optimal. Moreover, it is easy to see that it is not possible to obtain a "better" capping, with more cycles and/or less caps, that would even decrease the distance. An example of a graph with optimally capped paths is given in Figure 3.2.

Table 3.2: Linking paths from $G_R(\mathbb{A}, \mathbb{B})$ of canonical genomes. The symbol $\mathring{\alpha}_{\mathbb{A}}$ represents an artificial adjacency in \mathbb{A} and the symbol $\mathring{\alpha}_{\mathbb{B}}$ represents an artificial adjacency in \mathbb{B} . Observe that $\Delta d_{DCJ} = \Delta n_* - (\Delta(|\mathcal{C}|) + \Delta(2|\mathcal{P}_{\mathbb{A}\mathbb{B}}|)).$

	paths	linking cycle	Δn_*	$\Delta(\mathcal{C})$	$\Delta(2 \mathcal{P}_{\mathbb{AB}})$	$\mathbf{\Delta} \mathrm{d}_{\mathrm{DCJ}}$
1.	$\mathbb{AB}^{\langle i \rangle}$	$\mathbb{O}^{\langle i+1 \rangle} = \zeta(\mathbb{AB}^{\langle i \rangle})$	+0.5	+1	-0.5	0
2.	$\mathbb{A}\mathbb{A}^{\langle i\rangle} + \mathbb{B}\mathbb{B}^{\langle j\rangle}$	$\mathbb{O}^{\langle i+j+2\rangle} = \zeta(\mathbb{A}\mathbb{A}^{\langle i\rangle}, \mathbb{B}\mathbb{B}^{\langle j\rangle})$	+1	$^{+1}$	0	0
3.	$\mathbb{AA}^{\langle i \rangle}$	$\mathbb{O}^{\langle i+2 \rangle} = \zeta(\mathbb{A}\mathbb{A}^{\langle i \rangle}, \mathring{\alpha}_{\mathbb{B}})$	$^{+1}$	$^{+1}$	0	0
4.	$\mathbb{BB}^{\langle j \rangle}$	$\mathbb{O}^{\langle j+2 \rangle} = \zeta(\mathbb{BB}^{\langle j \rangle}, \mathring{\alpha}_{\mathbb{A}})$	+1	$^{+1}$	0	0

Recall that $\kappa(\mathbb{A}) + \kappa(\mathbb{B}) = |\mathcal{P}_{\mathbb{A}\mathbb{B}}| + |\mathcal{P}_{\mathbb{A}\mathbb{A}}| + |\mathcal{P}_{\mathbb{B}\mathbb{B}}|$. Now let $a_* = |\kappa(\mathbb{A}) - \kappa(\mathbb{B})| = ||\mathcal{P}_{\mathbb{A}\mathbb{A}}| - |\mathcal{P}_{\mathbb{B}\mathbb{B}}||$, $p_* = \max\{\kappa(\mathbb{A}), \kappa(\mathbb{B})\} = \frac{|\mathcal{P}_{\mathbb{A}\mathbb{B}}|}{2} + \max\{|\mathcal{P}_{\mathbb{A}\mathbb{A}}|, |\mathcal{P}_{\mathbb{B}\mathbb{B}}|\}$ and $u = \min\{|\mathcal{P}_{\mathbb{A}\mathbb{A}}|, |\mathcal{P}_{\mathbb{B}\mathbb{B}}|\}$. An optimal capping of all paths with the minimum number of caps maximizes the use of line (2.) of Table 3.2 and has p_* caps and a_* artificial adjacencies, while an optimal capping with the maximum number of caps uses only lines (1.) , (3.) and (4.) of Table 3.2 and has $p_* + u$ caps and $a_* + 2u$ artificial adjacencies.



Figure 3.2: Optimal capping of the paths of a relational graph. Cap vertices are drawn in gray.