# Topics of today:

Canonical inversion distance and sorting:

- 1. Breakpoint diagram
- 2. Split / Neutral / Joining inversions
- 3. Good / bad components
- 4. Safe inversions and overlap graph
- 5. Hurdles and fortress / component tree

### Canonical inversion model - circular chromosomes

(Unichromosomal genomes  $\equiv$  chromosomes)

Given two canonical circular chromosomes  $\mathbb A$  and  $\mathbb B,\ldots$ 

| Canonical Inversion Distance Problem: | Compute the minimum number of inversions required to transform $\mathbb{A}$ into $\mathbb{B}$ .                 |
|---------------------------------------|---|
|                                       | Denote by $d_{\rm INV}(\mathbb{A},\mathbb{B})$ the inversion distance of $\mathbb{A}$ and $\mathbb{B}.$         |
| Canonical Inversion Sorting Problem:  | Find a sequence of $d_{INV}(\mathbb{A}, \mathbb{B})$ inversions that transform $\mathbb{A}$ into $\mathbb{B}$ . |

#### Breakpoint diagram of canonical circular chromosomes

Let  $\mathbb{A}$  and  $\mathbb{B}$  be canonical circular chromosomes, with  $n = |\mathcal{G}_{\star}|$ . The **breakpoint diagram**  $BD(\mathbb{A}, \mathbb{B}) = (V, E)$  is described as follows:

1. 
$$V = \bigcup_{\mathbf{X} \in \mathcal{G}_{\star}} {\mathbf{X}^{h}, \mathbf{X}^{t}} \Rightarrow V = \xi(\mathbb{A}) = \xi(\mathbb{B}) ; |V| = 2n$$

there is a vertex for each extremity of each gene in  $\mathcal{G}_{\star}$ 

each vertex v has a label  $\ell(v)$ , that corresponds to the extremity it represents

The vertices are drawn in one line, next to each other.

The vertices must follow the same (circular) order of the corresponding extremities in chromosome  $\mathbb{A}$ , according to one of the two reading directions.

2.  $E = E_{\Gamma}(\mathbb{A}) \cup E_{\Gamma}(\mathbb{B})$ , where:

► Adjacency edges: 
$$\begin{cases} E_{\Gamma}(\mathbb{A}) = \{uv : u, v \in V(\xi(\mathbb{A})) \text{ and } \ell(u)\ell(v) \in \Gamma(\mathbb{A})\} \\ E_{\Gamma}(\mathbb{B}) = \{uv : u, v \in V(\xi(\mathbb{B})) \text{ and } \ell(u)\ell(v) \in \Gamma(\mathbb{B})\} \end{cases}$$

The number of edges is |E| = 2n (*n* adjacency edges per chromosome)

Two equivalent breakpoint diagrams

 $BD(\mathbb{A},\mathbb{B}) \cong BD(\mathbb{B},\mathbb{A})$ 



#### Properties of the breakpoint diagram

 $\mathbb{A} = (1\,\bar{7}\,4\,5\,3\,\bar{6}\,\bar{2})$ 



Every vertex has degree two:

 $BD(\mathbb{A}, \mathbb{B})$  is a collection of (even) cycles (alternating edes in  $E_{\Gamma}(\mathbb{A})$  and in  $E_{\Gamma}(\mathbb{B})$ ) cycle with *k* edges: *k*-cycle (always even)

 $\mathcal{C} = \mathsf{set} \mathsf{ of cycles in } BD(\mathbb{A},\mathbb{B})$ 

 $n = |\mathcal{G}_{\star}| = 7$ 

If  $\mathbb{A} = \mathbb{B}$ ,  $RG(\mathbb{A}, \mathbb{B})$  has only 2-cycles:  $2n = 2|\mathcal{C}| \implies n = |\mathcal{C}|$ 

Otherwise, if  $\mathbb{A} \neq \mathbb{B}$ :

n > |C|

#### Types of inversion and lower bound for the inversion distance



Lower bound for the inversion distance:  $\mathsf{d}_{\scriptscriptstyle\mathrm{INV}}(\mathbb{A},\mathbb{B})\geq n-|\mathcal{C}|$ 

# Types of cycles

Trivial cycle: one adjacency in each chromosome

2-cycle (sorted)

 $\ensuremath{\textbf{Good}}\xspace$  cycle: at least one pair of adjacencies with opposite directions

Can be split into two cycles by applying an inversion



(1234567)

# Types of cycles

Bad cycle: all adjacencies have the same direction

Cannot be split into two cycles



(123)

# (Interleaving) components

#### Breakpoint diagram:



Interleaving sequence of cycles:

 $c_1, c_2, ..., c_k$  such that  $c_i$  and  $c_{i+1}$  are interleaving for all  $1 \le i \le k-1$ 

Interleaving component or simply component K:

either a cycle *c* that does not interleave with any other cycle

or  $\begin{cases}
\text{for each pair of cycles } c, c' \in K \text{ there is an interleaving sequence from } c \text{ to } c' \\
K \text{ is maximal}
\end{cases}$ 

# Types of (interleaving) components

Trivial component: only one trivial 2-cycle

Good component: at least one good cycle



# Types of (interleaving) components

Bad component: only bad cycles



# Overlap graph of a component



# Another example



(123456)

Overlap graph:



Effects on the overlap graph by inverting a bad adjacency

Three overlapping adjaconcies:



Effects on the overlap graph by inverting a bad adjacency

Three overlapping adjaconcies:



V







#### Sorting a bad component with a neutral inversion



Sorting a bad component with a neutral inversion

### Sorting bad components with a joining inversion



By joining with an inversion two cycles  $c_1$  and  $c_2$ , that belong to two distinct components  $K_1$  and  $K_2$  respectively, we merge not only the components  $K_1$  and  $K_2$ , but also all components that separate  $K_1$  and  $K_2$ , into a single **good** component K.

Hurdle: a bad component that does not separate two bud components

Sorting bad components with a joining inversion

By joining with an inversion two cycles  $c_1$  and  $c_2$ , that belong to two distinct components  $K_1$  and  $K_2$ respectively, we merge not only the components  $K_1$  and  $K_2$ , but also all components that separate  $K_1$  and  $K_2$ , into a single **good** component K.

bad component ale bad component, won hurdle Another

# Quiz 1

- 1 Which of the following statements about the breakpoint diagram are true?
  - $\mathbf{X}$  A cycle can always be split into two cycles with an inversion.
  - X A neutral inversion cannot be optimal.
  - $\checkmark$  A joining inversion cannot be optimal.
  - D It is always possible to split a good cycle into two.
  - E A bad cycle cannot be split by an inversion.



Effects on the overlap graph by inverting a good adjacency

Three overlapping adjaconcies:







Effects on the overlap graph by inverting a good adjacency



Effects on the overlap graph by inverting a good adjacency



Sorting a good component - finding safe split inversions

score 
$$(xy)$$
: # good adjacencies in the diagram after fixing  $xy$   
score  $(xy)$ :  $G + b(xy) - g(xy) - 1$ 

Sorting a good component - finding safe split inversions

Sorting a good component - finding safe split inversions

### Sorting bad components with a joining inversion



By joining with an inversion two cycles  $c_1$  and  $c_2$ , that belong to two distinct components  $K_1$  and  $K_2$  respectively, we merge not only the components  $K_1$  and  $K_2$ , but also all components that separate  $K_1$  and  $K_2$ , into a single **good** component K.

Hurdle: a bad component that does not separate two bud components

#### Sorting bad components - simple hurdles and super hurdles

*h* : number of hurdles in  $BD(\mathbb{A}, \mathbb{B})$ 

hordle: bad component that does not separate 2 bad components Super hurdle K: fixing K by a neutral inversion creater a new hurdle On the previous pyce, both green and blue are super hurdles fixing only the green or only the blue component with a new trail invorsion would turn the red component into a hurdle

Sorting bad components - simple hurdles and super hurdles

fortress { \* there is an odd number of hurdles Sorting bad components - fortress A,B,C,D,E,F are bed components sopor hurdles are B separates  $f: \begin{cases} 0 & BD(\mathbb{A}, \mathbb{B}) \text{ is not a fortress} \\ 1 & BD(\mathbb{A}, \mathbb{B}) \text{ is a fortress} \end{cases}$ (A and D A and C A and F A and E 18 2 4 6 8 7 7 9 5 10 / 12 14 / 13 15 11 16 3 17 1 WW NUSPA C syarates A, D and F are E separatos ) and A super-hurdles, | F and A D and B but joining my ⇒ only A, I and F are hurdles F and B 0 and E pair among A,D,F '0 and F F and C creates a new hurdle F and D

# Canonical inversion distance of circular chromosomes

$$\mathsf{d}_{ ext{INV}}(\mathbb{A},\mathbb{B})=n-|\mathcal{C}|+h+f$$

1 Which of the following statements about the inversion model are true?

X The inversion distance depends only on the number of cycles in the breakpoint diagram.

X Every bad component in the diagram is a hurdle.

X A split inversion is always optimal.

A good component can always be sorted with (safe) split inversions.

X A super hurdle can be optimally sorted with a neutral inversion.

F A diagram with an even number of bad components can be a fortress.

#### Chained and nested components on the breakpoint diagram

Alternative to component separation: chaining and nesting relationships between components

 $\begin{array}{l} \textbf{Chain:} & \begin{cases} \text{sequence of components } K_1, K_2, ..., K_\ell \\ & \\ \text{the rightmost adjacency-edge of } K_i \text{ is succeeded by} \\ & \\ \text{the leftmost adjacency-edge of } K_{i+1}, \text{ for } 1 \leq i \leq \ell \end{cases} \end{array}$ 

Maximal chain: cannot be extended to the left nor to the right.

A maximal chain H is **nested** in a component K when the leftmost adjacency-edge of H is preceded by an adjacency-edge of K and the rightmost adjacency-edge of H is succeeded by an adjacency-edge of K.



 $K_1 = \{c_1\} \qquad K_2 = \{c_2\} \qquad \qquad K_3 = \{c_3\} \qquad \qquad K_4 = \{c_4\} \qquad \qquad K_5 = \{c_5, c_6\} \qquad \qquad K_6 = \{c_7\}$ 

 $\begin{array}{l} \mbox{Maximal chains:} \\ \begin{cases} H_1 = K_1 \Join K_5 \ , \\ H_2 = K_2 \Join K_4 \ (\mbox{nested in component} \ K_1) \ , \\ H_3 = K_3 \ (\mbox{nested in component} \ K_2) \ , \\ H_4 = K_6 \ (\mbox{nested in component} \ K_5) \end{cases}$ 

# Chained component tree $\Upsilon_{\blacksquare}$ (rooted)



$$\begin{split} & \mathcal{K}_{1} = \{c_{1}\} \ \ \mathcal{K}_{2} = \{c_{2}\} & \mathcal{K}_{3} = \{c_{3}\} & \mathcal{K}_{4} = \{c_{4}\} & \mathcal{K}_{5} = \{c_{5}, c_{6}\} & \mathcal{K}_{6} = \{c_{7}\} \\ & \mathsf{Maximal chains:} & \begin{cases} \mathcal{H}_{1} = \mathcal{K}_{1} \bowtie \mathcal{K}_{5} \ , \\ \mathcal{H}_{2} = \mathcal{K}_{2} \bowtie \mathcal{K}_{4} \ (\mathsf{nested in component} \ \mathcal{K}_{1}) \ , \\ \mathcal{H}_{3} = \mathcal{K}_{3} \ (\mathsf{nested in component} \ \mathcal{K}_{5}) \end{cases} \\ & \mathcal{K}_{6} = \{c_{7}\} \end{cases}$$

- 1. One round node per component  $K_i$ :  $\begin{cases}
  bad node (\circ): K_i \text{ is a bad component;} \\
  good node (\bullet): K_i \text{ is a trivial or a good component.}
  \end{cases}$
- One square (■) node per maximal chain H<sub>i</sub>, whose children are the round nodes corresponding to the components of H<sub>i</sub>. A square node is either the root or a child of the component in which H<sub>1</sub> is nested.



# Contraction of $\Upsilon_{\blacksquare}$ into unrooted component tree $\Upsilon_{\circ}$



#### Topology and paths in the component tree $\Upsilon_{\!\circ}$

All leaves in  $\Upsilon_{\circ}$  are bad nodes ( $\equiv$  hurdles)

 $\mathcal{L} {:}~\#$  of leaves in  $\Upsilon_{\circ}$ 

Traversal: path connecting two leaves of  $\Upsilon_{\circ}$ 



**Branching node** of  $\Upsilon_{\circ}$ : any node whose degree is  $\geq 3$  **Leaf-branch** of  $\Upsilon_{\circ}$ :  $\begin{cases}
\text{if } \mathcal{L} \leq 2: \text{ the complete tree } \Upsilon_{\circ} \\
\text{if } \mathcal{L} \geq 3: \text{ maximal path } u_1, u_2, ..., u_k, \text{ such that } u_1 \text{ is a leaf of } \Upsilon_{\circ} \text{ and,} \\
\text{ for } i = 2, ..., k, \text{ the degree of internal node } u_i \text{ in T is two}
\end{cases}$ 

A leaf-branch may be a path of length 1 (a leaf directly connected to a branching node of  $\Upsilon_{\circ}$ )

#### Cost of covering the component tree $\Upsilon_{\! \circ}$

path *P* in  $\Upsilon_0$ :  $\begin{cases}
\text{short: contains a single bad node} &\Rightarrow \text{ if } P \text{ is a branch it corresponds to a simple hurdle} \\
\text{long: contains at least two bad nodes} &\Rightarrow \text{ if } P \text{ is a branch it corresponds to a super hurdle} \\
\text{cost of path } P: \tau(P) \begin{cases}
P \text{ is short: } \tau(P) = 1 \text{ (cut a bad component)} \\
P \text{ is long: } \tau(P) = 2 \text{ (merge two or more bad components)}
\end{cases}$ 

Cover of  $\Upsilon_{\circ}$ : set of paths  $\widehat{\mathcal{P}}$  such that each bad node of  $\Upsilon_{\circ}$  is contained in at least one path  $P \in \widehat{\mathcal{P}}$ Cost of cover  $\widehat{\mathcal{P}}$ :  $\tau(\widehat{\mathcal{P}}) = \sum_{P \in \widehat{\mathcal{P}}} \tau(P)$ Cost of an optimal cover of  $\Upsilon_{\circ}$ :  $\tau(\Upsilon_{\circ}) = \min_{\widehat{\mathcal{P}} \text{ is a cover of } \Upsilon_{\circ}} \tau(\widehat{\mathcal{P}})$ 



# Covering the component tree $\Upsilon_{\! \circ}$

 $\mathcal{L}: \ \# \ \text{of leaves in} \ \Upsilon_{\circ} \ ; \quad \textbf{Branching node of} \ \Upsilon_{\circ}: \ \text{any node whose degree is} \geq 3$ 

 $\begin{array}{l} \mbox{Leaf-branch of } \Upsilon_{\circ} \colon \begin{cases} \mbox{if } \mathcal{L} \leq 2 \colon \mbox{the complete tree } \Upsilon_{\circ} \\ \mbox{if } \mathcal{L} \geq 3 \colon \mbox{maximal path } u_1, u_2, ..., u_k, \mbox{ such that } u_1 \mbox{ is a leaf of } \Upsilon_{\circ} \mbox{ and}, \\ \mbox{ for } i = 2, ..., k, \mbox{ the degree of internal node } u_i \mbox{ in } \mathsf{T} \mbox{ is two} \end{cases}$ 

A leaf-branch may be a path of length 1 (a leaf directly connected to a branching node of  $\Upsilon_{\circ}$ )

Traversal: path connecting two leaves of  $\Upsilon_{\!\scriptscriptstyle O}$ 

Suppose  $\mathcal{L} = 2, 4, 6, ...$ :

 $\widehat{\mathcal{P}}_{_{\mathrm{T}}}(\Upsilon_{\circ})$ : smallest set of traversals covering all nodes of  $\Upsilon_{\circ}$ :  $|\widehat{\mathcal{P}}_{_{\mathrm{T}}}(\Upsilon_{\circ})| = \frac{\mathcal{L}}{2}$ 

```
COVERTREEWITHTRAVERSALS

Input: unrooted tree \Upsilon_{\circ} with \mathcal{L} = 2n leaves

Output: set \widehat{\mathcal{P}}_{r} of n traversals covering all nodes of \Upsilon_{\circ}

Based on any planar view of \Upsilon_{\circ}, enumerate the leaves from 1 to 2n in circular order;

\widehat{\mathcal{P}}_{r} = \emptyset;

for i = 1 to n do

\widehat{\mathcal{P}}_{r} = \widehat{\mathcal{P}}_{r} \cup \{\text{traversal connecting leaves } i \text{ and } i + n\};

Return \widehat{\mathcal{P}}_{r};
```

L = h

# Computing $\tau(\Upsilon_{\circ})$

Lower bound for the cost of an optimal cover of  $\Upsilon_{\circ}$ :  $\tau(\Upsilon_{\circ}) \geq \mathcal{L}$ 

Each traversal T has cost  $\tau(T) = 2$ 

If  $\mathcal{L}$  is even,  $\widehat{\mathcal{P}}_{T}(\Upsilon_{\circ})$  is an optimal cover:  $\Rightarrow \tau(\Upsilon_{\circ}) = \tau\left(\widehat{\mathcal{P}}_{T}(\Upsilon_{\circ})\right) = 2\frac{\mathcal{L}}{2} = \mathcal{L}$ 

If  $\mathcal{L}$  is odd and  $\Upsilon_{\circ}$  has a short leaf-branch s ( $\tau(s) = 1$ ):  $\Rightarrow \tau(\Upsilon_{\circ}) = \tau\left(\widehat{\mathcal{P}}_{_{\mathrm{T}}}(\Upsilon_{\circ} \setminus s)\right) + \tau(s) = 2\frac{\mathcal{L}-1}{2} + 1 = \mathcal{L}$ 

If  $\mathcal{L}$  is odd and  $\Upsilon_{\circ}$  has no short leaf-branch ("fortress"); let  $\ell$  be any long leaf-branch of  $\Upsilon_{\circ}$  ( $\tau(\ell) = 2$ ):  $\Rightarrow \tau(\Upsilon_{\circ}) = \tau\left(\widehat{\mathcal{P}}_{\mathbb{T}}(\Upsilon_{\circ} \setminus \ell)\right) + \tau(\ell) = 2\frac{\mathcal{L}-1}{2} + 2 = \mathcal{L} + 1$ 

The cost of any optimal cover of  $\Upsilon_{\circ}$  is:

$$\tau(\Upsilon_{\circ}) = \begin{cases} \mathcal{L} + 1 & \text{if } \mathcal{L} \text{ is odd and all leaf-branches are long ("fortress"),} \\ \mathcal{L} & \text{otherwise.} \end{cases}$$

# Canonical inversion distance

$$\mathsf{d}_{\scriptscriptstyle\mathrm{INV}}(\mathbb{A},\mathbb{B})=n-|\mathcal{C}|+ au_*$$

where

 $au_* = au(\Upsilon_{\circ}(\mathbb{A},\mathbb{B})) = h + f$ 

#### Components are framed conserved intervals

Assuming that  $\mathbb{B} = (1 \ 2 \ 3 \ \dots \ 16)$ , let us identify its framed conserved intervals with respect to

 $\mathbb{A} = (1 \quad \overline{4} \quad 2 \quad 3 \quad 5 \quad 7 \quad 6 \quad 8 \quad \overline{16} \quad \overline{14} \quad \overline{15} \quad \overline{13} \quad \overline{11} \quad \overline{12} \quad \overline{10} \quad 9)$ 

For given  $i \ge 1$  and  $j \ge 1$  such that  $i+j \le n+1$ :

**Conserved interval:** interval of A composed of values i, i+1, ..., i+j (assuming  $n+1 \equiv 1$ )

**Framed** conserved interval  $\begin{cases} \text{direct: first element is } i \text{ and last element is } i+j; \text{ or } \\ \text{reverse: first element is } \overline{i+j} \text{ and last element is } \overline{i} \end{cases}$ 

Direct: [1..5]; [2..3]; [5..8]; [8..17] Reverse: [16..13]; [13..10], [16..10]

Component: framed conserved interval that is not a union of framed conserved intervals

Direct: [1..5]; [2..3]; [5..8]; [8..17] Reverse: [16..13]; [13..10]



### Components are framed conserved intervals





# Complexity of inversion distance and sorting

The inversion distance can be computed in linear time, by efficiently identifying chains of framed conserved intervals (Bergeron *et al.*, 2002: Common intervals and sorting by reversals: a marriage of necessity)

An optimal inversion sorting scenario can be computed in subquadratic time. (Tannier and Sagot, 2004: Sorting by reversals in subquadratic time)

#### Canonical inversion distance of linear chromosomes

Given canonical linear chromosomes  $\mathbb A$  and  $\mathbb B :$ 

Add one new family (e.g. 0) and circularize chromosome  $\mathbb B$  into  $\mathbb B' = (0 \ \mathbb B)$ 

$$d_{\text{INV}}(\mathbb{A},\mathbb{B}) = \min \begin{cases} d_{\text{INV}}((0 \ \mathbb{A}),\mathbb{B}') \\ d_{\text{INV}}((\overline{0} \ \mathbb{A}),\mathbb{B}') \end{cases}$$

#### Example:

$$\begin{split} \mathbb{A} &= [\bar{5} \ 1 \ 2 \ \bar{3} \ 4] \quad \text{and} \quad \mathbb{B} = [1 \ 2 \ 3 \ 4 \ 5] \\ \mathbb{B}' &= (0 \ 1 \ 2 \ 3 \ 4 \ 5) \\ d_{\rm INV}((0 \ \bar{5} \ 1 \ 2 \ \bar{3} \ 4), \mathbb{B}') &= 3 \\ d_{\rm INV}((\bar{0} \ \bar{5} \ 1 \ 2 \ \bar{3} \ 4), \mathbb{B}') &= 2 \\ d_{\rm INV}((\bar{0} \ \bar{5} \ 1 \ 2 \ \bar{3} \ 4), \mathbb{B}') &= 2 \end{split}$$

1 What is the bottleneck of the running time of inversion sorting?

- A Finding inversions that fix bad components.
- B Finding split inversions.
- C Finding safe split inversions.
- D Finding inversions that merge bad components.

#### References

Transforming cabbage into turnip: polynomial algorithm for sorting signed permutations by reversals Sridhar Hannenhalli and Pavel A. Pevzner Journal of the ACM, vol 46, issue 1, pages 1–27 (1999)

Reversal Distance without Hurdles and Fortresses (Anne Bergeron, Julia Mixtacki and Jens Stoye) LNCS, volume 3109, pages 388-399 (2004)

The Inversion Distance Problem

(Anne Bergeron, Julia Mixtacki and Jens Stoye)

In: Mathematics of Evolution and Phylogeny. Gascuel O (Ed); (2005)