## Topics of today:

Canonical inversion distance and sorting:

1. Relational / Breakpoint diagram

2. Split / Neutral / Joining inversions

3. Good / bad components

4. Hurdles and fortress

# Canonical inversion model - circular chromosomes

(Unichromosomal genomes $\equiv$ chromosomes)

Given two canonical circular chromosomes $\mathbb{A}$ and $\mathbb{B}$,...

**Canonical Inversion Distance Problem:**    Compute the minimum number of inversions required to transform $\mathbb{A}$ into $\mathbb{B}$.

Denote by $d_{\text{INV}}(\mathbb{A}, \mathbb{B})$ the inversion distance of $\mathbb{A}$ and $\mathbb{B}$.

**Canonical Inversion Sorting Problem:**    Find a sequence of $d_{\text{INV}}(\mathbb{A}, \mathbb{B})$ inversions that transform $\mathbb{A}$ into $\mathbb{B}$.

# Relational diagram of canonical circular chromosomes

Given canonical circular chromosomes $\mathbb{A}$ and $\mathbb{B}$, their **relational diagram** $RD(\mathbb{A}, \mathbb{B}) = (V, E)$ is described as follows:

1. $V = V(\xi(\mathbb{A})) \cup V(\xi(\mathbb{B}))$ :    there is a vertex for each extremity of each gene in $\mathbb{A}$

   and a vertex for each extremity of each gene in $\mathbb{B}$

   The vertices corresponding to $\xi(\mathbb{A})$ are drawn in an upper line,
   while the vertices corresponding to $\xi(\mathbb{B})$ are drawn in a lower line.

   In each line, the vertices must follow the same (circular) order of the corresponding extremities in the respective chromosome, according to one of the two reading directions.

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

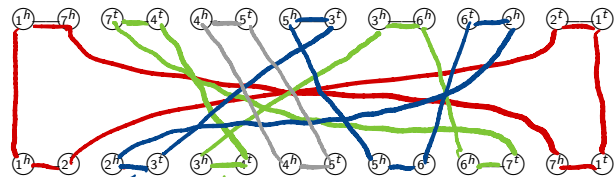2. $E = E_\alpha(\mathbb{A}) \cup E_\alpha(\mathbb{B}) \cup E_\xi$, where:

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

   ▶ **Extremity edges:** $E_\xi = \{uv : u \in V(\xi(\mathbb{A})) \text{ and } v \in V(\xi(\mathbb{B})) \text{ and } \ell(u) = \ell(v)\}$

Note that:

▶ Let $n = |\mathcal{G}_\star|$. The number of edges in $E_\alpha(\mathbb{A}) \cup E_\alpha(\mathbb{B})$ is $2n$ ($n$ adjacency edges per chromosome).

# Relational diagram of canonical circular chromosomes



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

$$\mathbb{B} = (1\;2\;3\;4\;5\;6\;7)$$

4-cycle

2 cycle

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

Every vertex has degree two:

$RD(\mathbb{A}, \mathbb{B})$ is a collection of (even) cycles (alternating edes in $E_\xi$ and in $E_\alpha(\mathbb{A}) \cup E_\alpha(\mathbb{B})$)

cycle with $k$ edges in $E_\alpha(\mathbb{A}) \cup E_\alpha(\mathbb{B})$: $k$-cycle

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

If $\mathbb{A} = \mathbb{B}$,
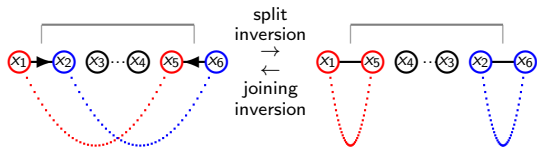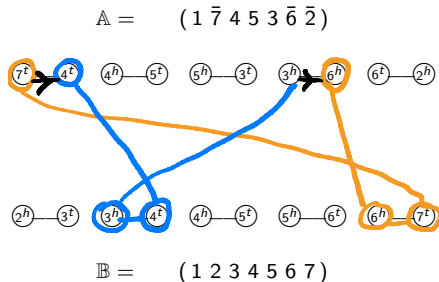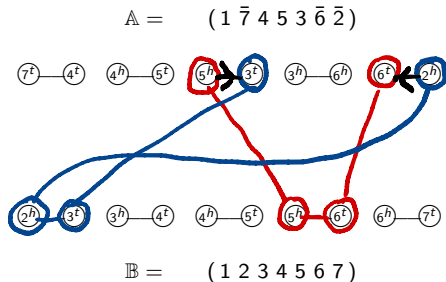$RG(\mathbb{A}, \mathbb{B})$ has only 2-cycles:

$$2n = 2|\mathcal{C}| \quad \Rightarrow \quad n = |\mathcal{C}|$$

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

$$n > |\mathcal{C}|$$

# Types of inversion and lower bound for the inversion distance



Assign one (arbitrary) direction to each cycle of $RD(\mathbb{A}, \mathbb{B})$

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

$\mathbb{B} = (1\ 2\ 3\ 4\ 5\ 6\ 7)$

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

$\mathbb{B} = (1\ 2\ 3\ 4\ 5\ 6\ 7)$

split inversion $\rightarrow$ $\leftarrow$ joining inversion

$\leftrightarrow$ neutral inversion

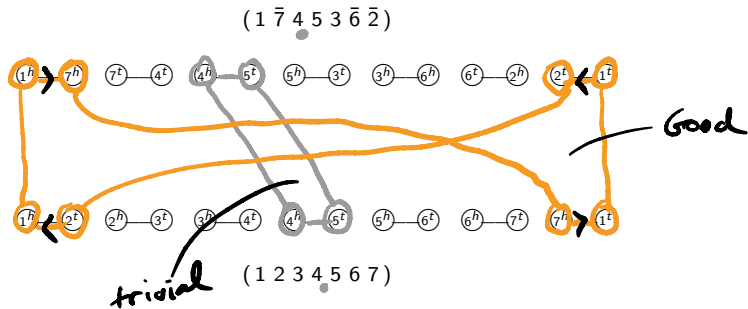Lower bound for the inversion distance: $\quad d_{\text{INV}}(\mathbb{A}, \mathbb{B}) \geq n - |\mathcal{C}|$

# Types of cycles

**Trivial cycle**: one adjacency in each chromosome

2-cycle (sorted)

**Good cycle**: $\begin{cases} \text{at least one pair of adjacencies with opposite directions in chromosome } \mathbb{A} \\ \text{at least one pair of adjacencies with opposite directions in chromosome } \mathbb{B} \end{cases}$
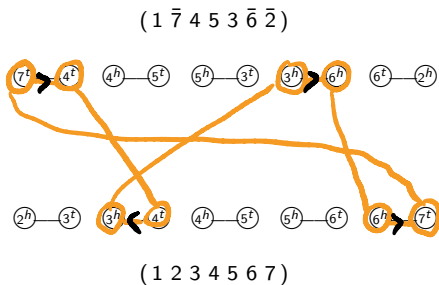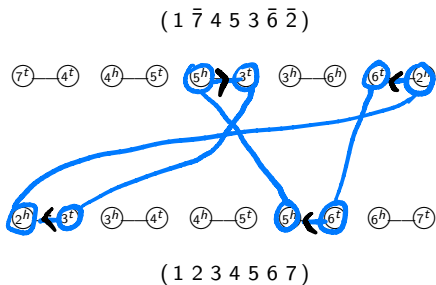
Can be split into two cycles by applying an inversion in $\mathbb{A}$ or in $\mathbb{B}$

# Types of cycles

**Semi-good cycle**: $\begin{cases} \text{at least one pair of adjacencies with opposite directions in one of the two chromosomes} \\ \text{all adjacencies have the same direction in the other chromosome} \end{cases}$
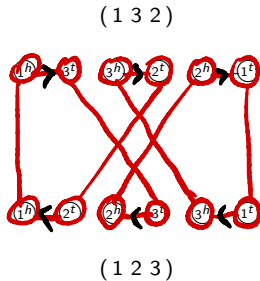
Can be split into two cycles by applying an inversion only in $\mathbb{A}$ or only in $\mathbb{B}$



$(1\ \bar{7}\ 4\ 5\ 3\ \bar{6}\ \bar{2})$              $(1\ \bar{7}\ 4\ 5\ 3\ \bar{6}\ \bar{2})$

$(1\ 2\ 3\ 4\ 5\ 6\ 7)$              $(1\ 2\ 3\ 4\ 5\ 6\ 7)$
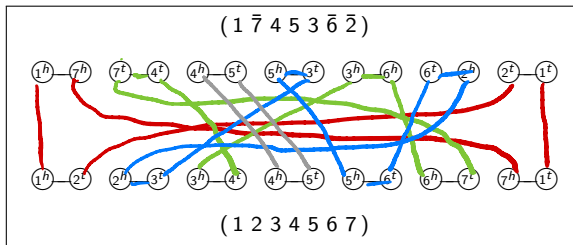
# Types of cycles

**Bad cycle**: $\begin{cases} \text{all adjacencies in chromosome } \mathbb{A} \text{ have the same direction} \\ \text{all adjacencies in chromosome } \mathbb{B} \text{ have the same direction} \end{cases}$

Cannot be split into two cycles

# Relational diagram $\cong$ Breakpoint diagram

**Relational diagram:**



$(1\ \bar{7}\ 4\ 5\ 3\ \bar{6}\ \bar{2})$

$(1\ 2\ 3\ 4\ 5\ 6\ 7)$

Looking either at the top line or
at the bottom line of the diagram:

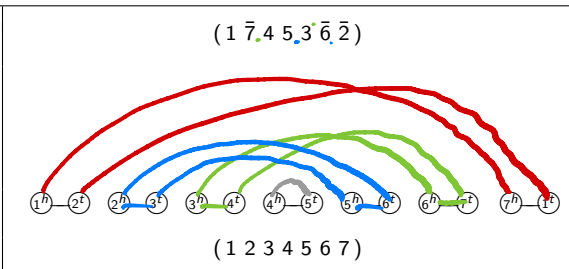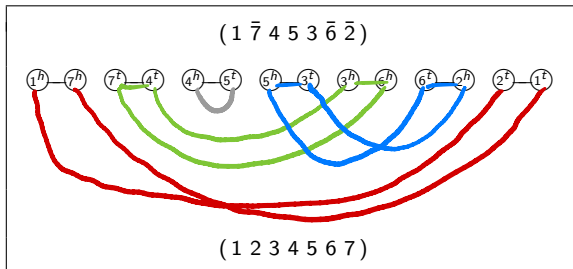**Two interleaving cycles:** $c \ldots c' \ldots c \ldots c'$

**Interleaving sequence of cycles:** $c_1, c_2, ..., c_k$ such
that $c_i$ and $c_{i+1}$ are interleaving for all $1 \leq i \leq k-1$

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



**Breakpoint diagrams:**



$(1\ \bar{7}\ 4\ 5\ 3\ \bar{6}\ \bar{2})$

$(1\ 2\ 3\ 4\ 5\ 6\ 7)$

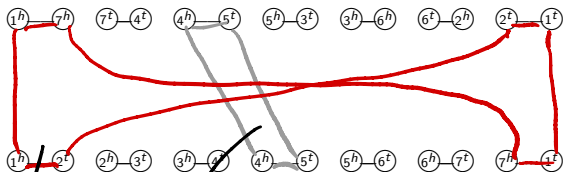$(1\ \bar{7}\ 4\ 5\ 3\ \bar{6}\ \bar{2})$

$(1\ 2\ 3\ 4\ 5\ 6\ 7)$

# Types of (interleaving) components

**Trivial component**: only one trivial 2-cycle

**Good component**: at least one good or semi-good cycle

# Types of (interleaving) components

**Bad component**: only bad cycles
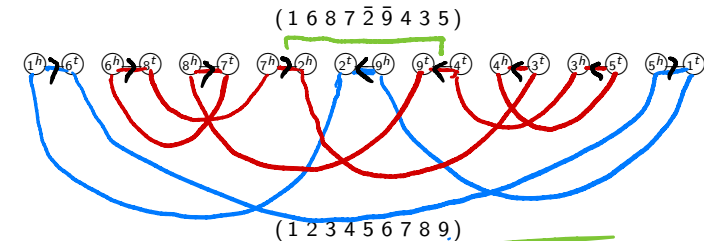
( 1 4 3 2 )



( 1 2 3 4 )

Bad

# Quiz 1

1 Which of the following statements about the relational diagram are true?

✗ A cycle can always be split into two cycles with an inversion.

✗ A joining inversion cannot be optimal.

✗ A split inversion is always optimal.

Ⓓ It is always possible to split a good or a semi-good cycle into two.

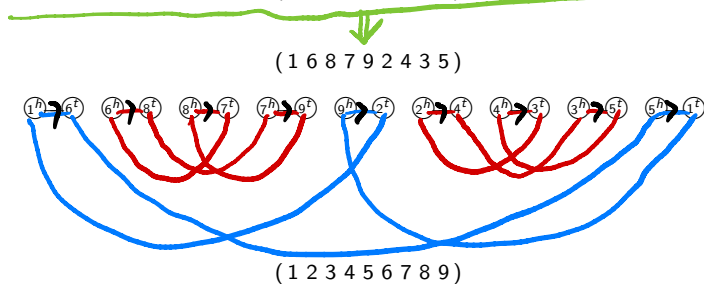Ⓔ A bad cycle cannot be split by an inversion.

# Unsafe inversions

*good*

A split inversion applied to a cycle of a good component can create bad components



$( 1\ 6\ 8\ 7\ \overline{2}\ \overline{9}\ 4\ 3\ 5 )$

$( 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9 )$

$( 1\ 6\ 8\ 7\ 9\ 2\ 4\ 3\ 5 )$

$( 1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9 )$

\- One good component

\- 3 bad components

# Sorting a good component - finding safe split inversions

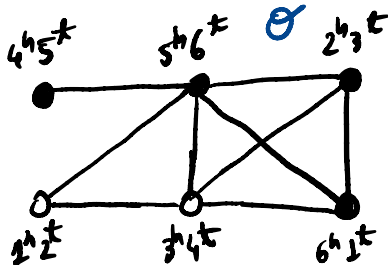**Target adjacency:** $\begin{cases} \text{good} \\ \text{bad} \end{cases}$

**Overlapping** target adjacencies

**Overlap graph** of a good component

Good:   Bad: 

Overlapping: 

Non overlapping:  OR 

$( \; 1 \quad 5 \quad \bar{4} \quad 2 \quad \bar{6} \quad \bar{3} \; )$



$( \; 1 \quad 2 \quad 3 \quad 4 \quad 5 \quad 6 \; )$

$4^h 5^t \qquad 5^h 6^t \qquad 2^h 3^t$

$1^h 2^t \qquad 3^h 4^t \qquad 6^h 1^t$



good adjacency
↓
black vertex

bad adjacency
↓
white vertex

$( 1 \; 5 \; \bar{4} \; 2 \; \bar{6} \; \bar{3} )$

$1^h \cdots \cdots 1^t$

$( 1 \; 2 \; 3 \; 4 \; 5 \; 6 )$

$4^h5^t \quad 5^h6^t \quad \mathcal{O} \quad 2^h3^t$

$1^h2^t \quad 3^h4^t \quad 6^h1^t$

$\mathcal{O}_{xy}$: subgraph of $\mathcal{O}$ composed of $xy$ and its adjacent vertices
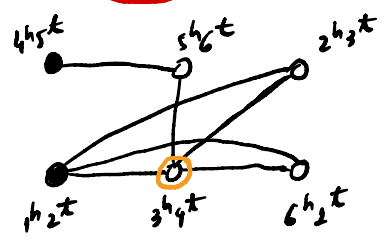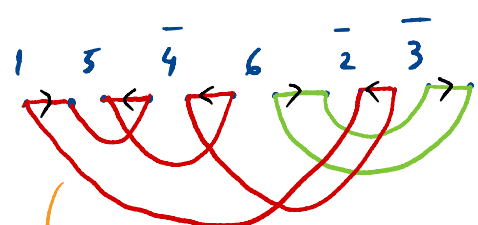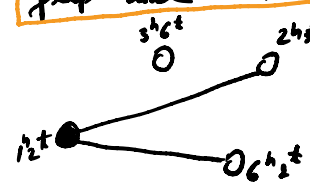
Effect of an inversion on the overlap graph $\mathcal{O}$:

if $xy$ is bad: 
1) flip the types of the vertices of $\mathcal{O}_{xy} - \{xy\}$
2) complement the edges of $\mathcal{O}_{xy} - \{xy\}$

Ex: $xy = 3^h4^t$
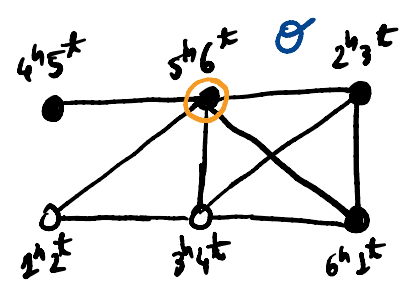
$\mathcal{O}_{3^h4^t}$:

$3^h6^t \quad 2^h3^t$
$1^h2^t \quad 3^h4^t \quad 6^h1^t$

$\mathcal{O}_{3^h4^t} - \{3^h4^t\}$

$5^h6^t \quad 2^h3^t$
$1^h2^t \quad 6^h1^t$

$1 \quad 5 \quad \bar{4} \quad 6 \quad \bar{2} \quad \bar{3}$

flip and complement

$3^h6^t \quad 2^h3^t$
$1^h2^t \quad 6^h1^t$

$4^h5^t \quad 5^h6^t \quad 2^h3^t$
$1^h2^t \quad 3^h4^t \quad 6^h1^t$

if $xy$ is good:
1) flip the types of all vertices of $\mathcal{O}_{xy}$
2) complement the edges of $\mathcal{O}_{xy}$

$( 1 \; 5 \; \bar{4} \; 2 \; \bar{6} \; \bar{3} )$

$1^h \cdots \cdots 1^t$

$( 1 \; 2 \; 3 \; 4 \; 5 \; 6 )$

$4^h5^t \quad 5^h6^t \quad \mathcal{O} \quad 2^h3^t$

$1^h2^t \quad 3^h4^t \quad 6^h1^t$

flip and complement

Ex: $xy = 5^h6^t$

$\mathcal{O}_{5^h6^t} = \mathcal{O}$

$( 1 \; 5 \; 6 \; \bar{2} \; 4 \; \bar{3} )$

$1^t$

$( 1 \; 2 \; 3 \; 4 \; 5 \; 6 )$

$4^h5^t \quad 5^h6^t \quad 2^h3^t$
$1^h2^t \quad 3^h4^t \quad 6^h1^t$

# Sorting a good component - finding safe split inversions

$G$: # of good adjacencies in $RD(A, B)$

$g(xy)$: # of good adjacencies overlapping $xy$ in $RD(A, B)$

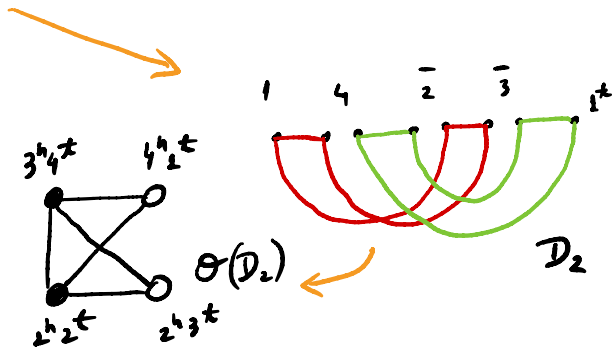$b(xy)$: # of bad adjacencies overlapping $xy$ in $RD(A, B)$
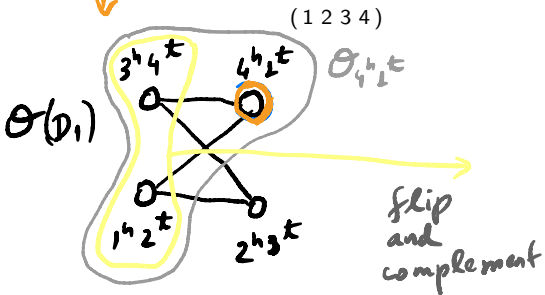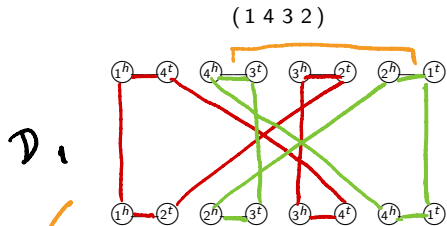
score $(xy)$: # good adjacencies in the diagram after fixing $xy$

$$\text{score}(xy) = G + b(xy) - g(xy) - 1$$

An inversion that fixes a good target adjacency with maximal score is SAFE. (does not create new bad components)

# Sorting a bad component with a neutral inversion

**Overlap graph** of a bad component



( 1 4 3 2 )

$D_1$

$\Theta(D_1)$

( 1 2 3 4 )

$3^h 4^t$   $4^h 2^t$   $\Theta_{4^h 2^t}$

$1^h 2^t$   $2^h 3^t$

flip and complement

Any neutral inversion applied to a bad adjacency of a bad component $k$ turns $k$ into a good component

$3^h 4^t$   $4^h 2^t$

$2^h 2^t$   $2^h 3^t$

$\Theta(D_2)$

$D_2$

# Sorting bad components - hurdles

$K_1$, $K_2$ and $K_3$ are three distinct components in $RD(\mathbb{A}, \mathbb{B})$ so that $K_3 \ldots K_1 \ldots K_1 \ldots K_3 \ldots K_2 \ldots K_2$

$\Rightarrow$ $K_3$ **separates** $K_1$ and $K_2$

$$(\,1\ 6\ 8\ 7\ 9\ 2\ 4\ 3\ 5\,)$$

$1^h$—$6^t$  $6^h$—$8^t$  $8^h$—$7^t$  $7^h$—$9^t$  $9^h$—$2^t$  $2^h$—$4^t$  $4^h$—$3^t$  $3^h$—$5^t$  $5^h$—$1^t$



$$(\,1\ 2\ 3\ 4\ 5\ 6\ 7\ 8\ 9\,)$$

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 $K$.

# Sorting bad components - simple hurdles and super hurdles

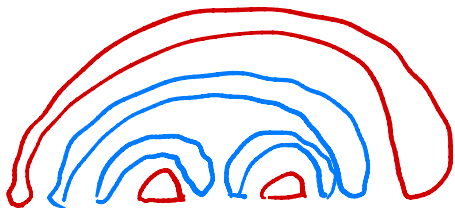$h$ : number of hurdles in $RD(\mathbb{A}, \mathbb{B})$

hurdle : bad component that does
not separate 2 bad components

super hurdle $K$ : fixing $K$ by a neutral
inversion creates a new
hurdle

# Sorting bad components - fortress

$$f : \begin{cases} 0 & RD(\mathbb{A}, \mathbb{B}) \text{ is not a fortress} \\ 1 & RD(\mathbb{A}, \mathbb{B}) \text{ is a fortress} \end{cases}$$



fortress: a) odd # of hurdles

0) all hurdles are super hurdles

# Canonical inversion distance of circular chromosomes

$$d_{\mathrm{INV}}(\mathbb{A}, \mathbb{B}) = n - |\mathcal{C}| + h + f$$

# Quiz 2

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

    A The inversion distance depends only on the number of cycles in the relational diagram.

    B Every bad component in the diagram is a hurdle.

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

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

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

# 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, No. 1, pages. 1–27 (1999)

The Inversion Distance Problem

(Anne Bergeron, Julia Mixtacki and Jens Stoye)

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