

Algorithms in Comparative Genomics

Lecture:

Marília D. V. Braga
Thursdays, 10:15-11:45

Tutorial:

Leonard Bohnenkämper
Thursdays (or Fridays?), 8:30-10:00

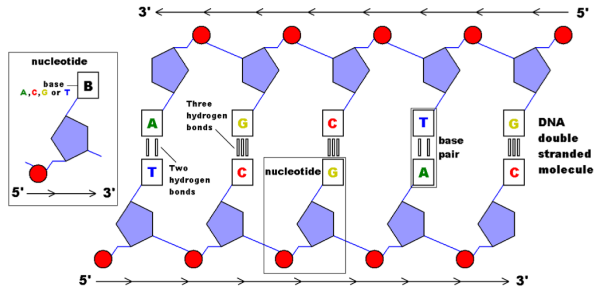
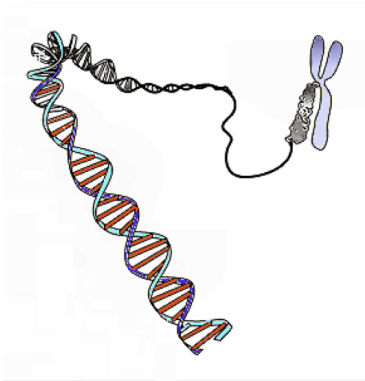
Topics:

1. Genomes as gene orders or list of adjacencies
2. Large-scale rearrangements
3. Types of genomes
4. Breakpoint distance, breakpoint double distance
5. Common intervals and gene clusters
6. Relational diagram of two genomes
7. Double-cut-and-join (DCJ) operation
 - 7.1 DCJ distance
 - 7.2 DCJ-indel distance
 - 7.3 Capping (circularizing genomes)
 - 7.4 Duplications and ILP
 - 7.5 Family-free setting
8. Inversion distance
9. Single-cut-or-join (SCJ) distance and median
10. Small parsimony

Topics of today - Introduction:

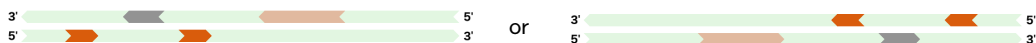
1. Genomes as gene orders or list of adjacencies
2. Large-scale rearrangements
3. Types of genomes
4. Breakpoint distance, breakpoint double distance

Genomes



Genomes as gene orders or list of adjacencies - linear chromosomes

Each gene is an oriented DNA fragment:
it lies on one of the two complementary anti-parallel DNA strands



↓ A chromosome is represented by its gene order ↓



[1 $\bar{2}$ 1 $\bar{3}$]

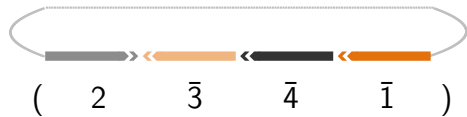
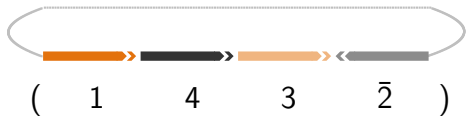
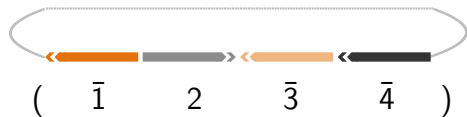
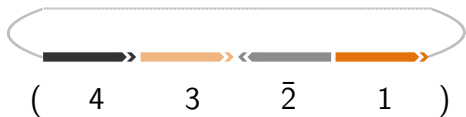
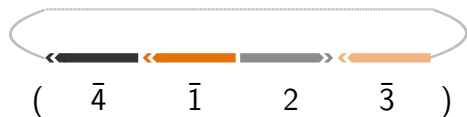
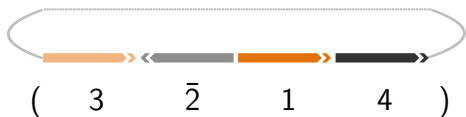
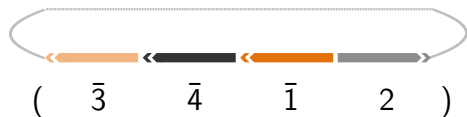
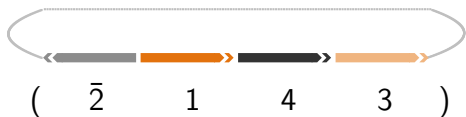
[3 $\bar{1}$ 2 $\bar{1}$]

1_1^t $1_1^h 2^h$ $2^t 1_2^t$ $1_2^h 3^h$ 3^t

3^t $3^h 1_2^h$ $1_2^t 2^t$ $2^h 1_1^h$ 1_1^t

(genes of the same color/number belong to the same **family**)

Genomes as gene orders or list of adjacencies - circular chromosomes

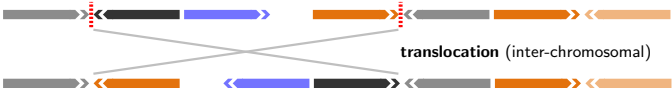


$1^h 4^t$ $4^h 3^t$ $3^h 2^h$ $2^t 1^t$

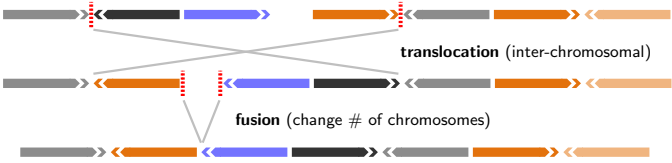
Large-scale genome rearrangements



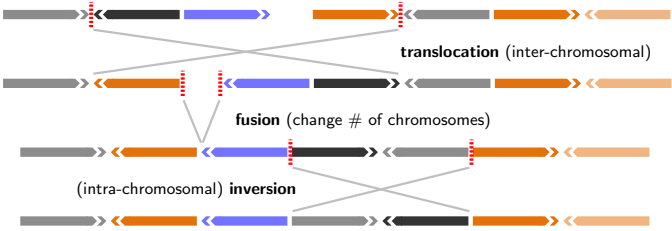
Large-scale genome rearrangements



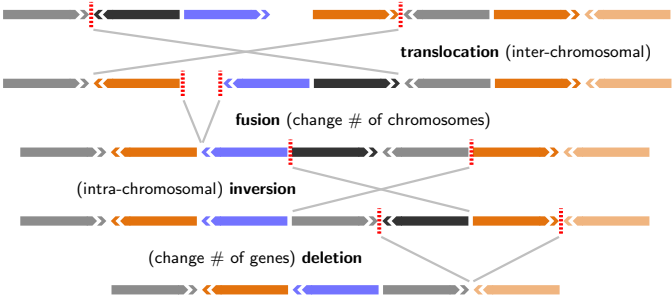
Large-scale genome rearrangements



Large-scale genome rearrangements



Large-scale genome rearrangements



Types of genomes

- ▶ Unichromosomal × multichromosomal
- ▶ Linear, circular, mixed
- ▶ Concerning the gene content:

1. **Singular genome:** each family occurs **exactly once**



2. **Duplicated genome:** each family occurs **exactly twice**



3. **Perfectly duplicated or doubled genome:** each adjacency or telomere occurs **exactly twice**



4. **Natural genome:** **no restriction** on the number of occurrences of families



Comparison of genomes



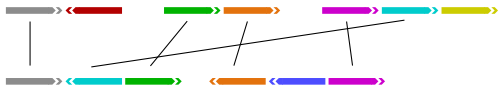
↕ find a distance/similarity measure
↕ between two gene orders



Types of genome pairs

Pair of singular genomes:

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



Pair of balanced genomes:

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



Pair of natural genomes:

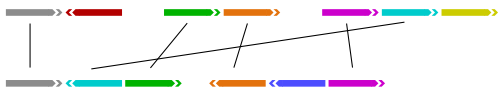
no restriction on the number of occurrences of families



Types of genome pairs

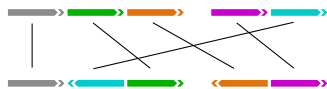
Pair of singular genomes:

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



Pair of canonical genomes:

singular and balanced



Pair of balanced genomes:

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



Pair of natural genomes:

no restriction on the number of occurrences of families



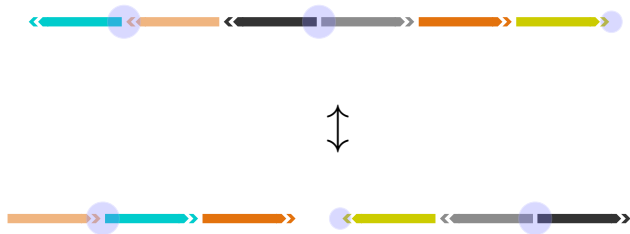
Breakpoint distance of canonical genomes

Common adjacency \times breakpoint



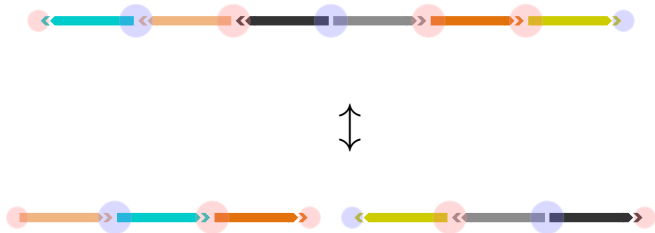
Breakpoint distance of canonical genomes

Common adjacency \times breakpoint



Breakpoint distance of canonical genomes

Common adjacency \times breakpoint



$$d_{BP}(A, B) = n - a - \frac{t}{2}$$

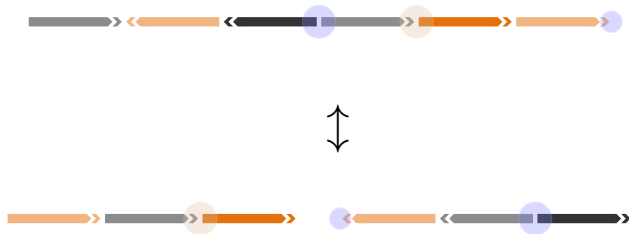
Breakpoint distance of balanced genomes

Common adjacency \times breakpoint



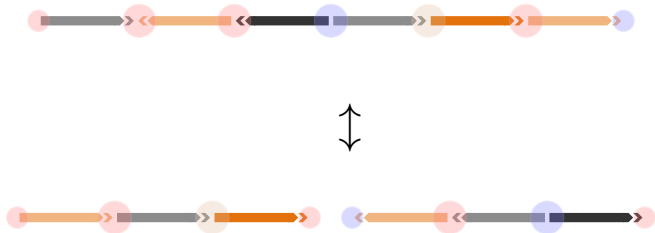
Breakpoint distance of balanced genomes

Common adjacency \times breakpoint



Breakpoint distance of balanced genomes

Common adjacency \times breakpoint



$$d_{BP}(A, B) = n - a - \frac{t}{2}$$

Breakpoint double distance

Given a singular genome S , let $S \oplus S$ be a doubled genome (obtained by duplicating each adjacency and each telomere of S).

Examples:

$$S_1 = [\bar{2}13] \text{ and } S_1 \oplus S_1 = [\bar{2}13] [\bar{2}13]$$
$$S_2 = (\bar{2}13) \text{ and } S_2 \oplus S_2 = (\bar{2}13) (\bar{2}13) \text{ or } S_2 \oplus S_2 = (\bar{2}13\bar{2}13)$$

Given a duplicated genome D , the breakpoint double distance is defined as:

$$d_{\text{BP}}^2(D, S) = d_{\text{BP}}(D, S \oplus S)$$

Ex: $D = [\bar{1}2\bar{3}132]$ and $S = [\bar{2}13]$



Reference

Multichromosomal median and halving problems under different genomic distances

(Eric Tannier, Chunfang Zheng and David Sankoff)

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

Quiz

Given genomes $A = (1\ 2\ 3\ 4)\ [1\bar{5}\bar{4}\ 5\bar{3}\bar{2}]$, $B = [1\ 2\ 3\ 4\ 5]$ and $C = [\bar{2}\ \bar{1}]\ [\bar{4}\ \bar{3}\ 5]$.

1 Which of the following statements are true?

A Genome A is linear.

B Genome A is multichromosomal.

C Genome A is duplicated.

D Genome A is doubled.

2 What is the breakpoint distance of B and C ?

A 1,5

B 2

C 2,5

D 3

2 How many families occur in genome A ?

A 4

B 5

C 5,5

D 6

4 What is the breakpoint double distance of A and B ?

A 4

B 4,2

C 4,5

D 5

Breakpoint halving

Given a duplicated genome D ,
find a singular genome S that minimizes
the breakpoint double distance:

$$d_{\text{BP}}^2(D, S) = d_{\text{BP}}(D, S \oplus S)$$