Introduction to Bioinformatics

Lecture 4: Genome rearrangements

# Why study genome rearrangements?

Provide insight into evolution of speciesFun algorithmic problem!

p Structure of this lecture:

- n The biological phenomenon
- n How to computationally model it?
- n How to compute interesting things?
- Studying the phenomenon using existing tools (continued in exercises)

# Genome rearrangements as an algorithmic problem



# Background

p Genome sequencing enables us to compare genomes of two or more different species

n -> Comparative genomics

p Basic observation:

- n Closely related species (such as human and mouse) can be almost identical in terms of genome contents...
- n ...but the order of genomic segments can be very different between species

### Synteny blocks and segments

- p Synteny derived from Greek 'on the same ribbon' – means genomic segments located on the same chromosome
  - n Genes, markers (any sequence)
- p Synteny block (or syntenic block)
  - n A set of genes or markers that co-occur together in two species
- p Synteny segment (or syntenic segment)
  - n Syntenic block where the *order* of genes or markers is preserved

# Synteny blocks and segments



Chromosome j, species C

# Observations from sequencing

- Large chromosome *inversions* and *translocations* (we'll get to these shortly) are common
  - n ... Even between closely related species
- 2. Chromosome inversions are usually symmetric around the *origin of DNA replication*
- 3. Inversions are less common *within species*...

### What causes rearrangements?

- P RecA, Recombinase A, is a protein used to repair chromosomal damage
- P It uses a duplicate copy of the damaged sequence as template
- P Template is usually a homologous sequence on a sister chromosome





# What effects does RecA have on genome?

- P Repeated sequences cause RecA to fail to choose correct recombination start position
- p This leads to
  - n Tandem duplications \_
  - n Translocations
  - n Inversions







X, Y, Z and W are repeats of the same sequence.

a, b, c and d are sequences on genome bounded by repeats.

In a tandem duplication example, RecA recombines a sequence that starts from Y instead of Z after Z.

This leads to duplication of segment Y-Z.

Diarmaid Hughes: Evaluating genome dynamics: the constraints on rearrangements within bacterial genomes, *Genome Biology* 2000, 1



Recombination of two repeat sequences in the *same* chromosome can lead to a fragment translocation

Here sequence d is translocated



Diarmaid Hughes: Evaluating genome dynamics: the constraints on rearrangements within bacterial genomes, *Genome Biology* 2000, 1



Diarmaid Hughes: Evaluating genome dynamics: the constraints on rearrangements within bacterial genomes, *Genome Biology* 2000, 1

### Example: human vs mouse genome

- P Human and mouse genomes share thousands of homologous genes, but they are
  - n Arranged in different order
  - n Located in different chromosomes

#### p Examples

- n Human chromosome 6 contains elements from six different mouse chromosomes
- n Analysis of X chromosome indicates that rearrangements have happened primarily *within* chromosome



Fig. 5.1. Syntenic blocks conserved between human chromosome Hsa6 and mouse chromosomes. Broken lines indicate regions that appear in inverted orders in the two organisms. Reprinted, with permission, from Gregory SG et al. (2002) *Nature* 418:743–750. Copyright 2002 Nature Publishing Group.



Fig. 5.3. Synteny blocks shared by human and mouse X chromosomes. The arrowhead for each block indicates the direction of increasing coordinate values for the human X chromosome. Reprinted, with permission, from Pevzner P and Tesler G (2003) Genome Research 13:37–45. Copyright 2003 Cold Spring Harbor Laboratory Press.

#### Representing genome rearrangments

- P When comparing two genomes, we can find homologous sequences in both using BLAST, for example
- p This gives us a map between sequences in both genomes



Fig. 5.1. Syntenic blocks conserved between human chromosome Hsa6 and mouse chromosomes. Broken lines indicate regions that appear in inverted orders in the two organisms. Reprinted, with permission, from Gregory SG et al. (2002) Nature 418:743–750. Copyright 2002 Nature Publishing Group.

# Representing genome rearrangments

- P We assign numbers 1,...,n to the found homologous sequences
- P By convention, we number the sequences in the first genome by their order of appearance in chromosomes
- p If the homolog of *i* is in reverse orientation, it receives number –*i* (signed data)
- p For example, consider human vs mouse gene numbering on the right

Human		Μοι	Mouse	
1	(gnat2)	12	(inpp1)	
2	(nras)	13	(cd28)	
3	(ngfb)	14	(fn1)	
4	(gba)	15	(pax3)	
5	(pklr)	-9	(il10)	
6	(at3)	-8	(pdc)	
7	(lamc1)	-7	(lamc1)	
8	(pdc)	-6	(at3)	
9	(il10)	/		

List order corresponds to *physical order* on chromosomes!

#### Permutations

- p The basic data structure in the study of genome rearrangements is *permutation*
- P A permutation of a sequence of n numbers is a reordering of the sequence
- p For example, 4 1 3 2 5 is a permutation of 1 2 3 4 5

# Genome rearrangement problem

- p Given two genomes (set of markers), how many
  - n duplications,
  - n inversions and
  - n translocations
  - do we need to do to transform the first genome to the second?

Minimum number of operations? What operations? Which order?

# Genome rearrangement problem









Keep in mind, that the two genomes have been evolved from a common ancestor genome! Genome rearrangements using reversals (=inversions) only

- p Lets consider a simpler problem where we just study reversals with unsigned data
- P A reversal p(i, j) reverses the order of the segment  $\Pi_{i} \Pi_{i+1} \dots \Pi_{j-1} \Pi_{j}$  (indexing starts from 1)
- p For example, given permutation
  - 6 1 2 3 4 5 and reversal p(3, 5) we get permutation 6 1 4 3 2 5



... note that we do not care about exact positions on the genome

#### Reversal distance problem

- p Find the shortest series of reversals that, given a permutation ∏, transforms it to the *identity* permutation (1, 2, ..., n)
- p This quantity is denoted by  $d(\Pi)$
- p Reversal distance for a pair of chromosomes:
  - n Find synteny blocks in both
  - n Number blocks in the first chromosome to identity
  - n Set ∏ to correspond matching of second chromosome's blocks against the first
  - n Find reversal distance

#### Reversal distance problem: discussion

p If we can find the minimal series of reversals for some pair of genomes
 n Is that what happened during evolution?
 n If not, is it the correct number of reversals?

p In any case, reversal distance gives us a measure of evolutionary distance between the two genomes and species

# Solving the problem by sorting

- p Our first approach to solve the reversal distance problem:
  - n Examine each position i of the permutation
  - n At each position, if  $\prod_i \neq i$ , do a reversal such that  $\prod_i = i$
- p This is a greedy approach: we try to choose the best option at each step

#### Simple reversal sort: example

Reversal series: p(1,2), p(2,3), p(3,4), p(5,6)

Is d(6 1 2 3 4 5) then 4?

$$D(6 1 2 3 4 5) = 2$$

# Pancake flipping problem

- No pancake made by the chef is of the same size
- Pancakes need to be rearranged before delivery
- Flipping operation:
  take some from the
  top and flip them over
- p This corresponds to always reversing the sequence prefix



123456