

## Are there Fragile Regions in the Human Genome? Combinatorial Algorithms

Phillip Compeau and Pavel Pevzner  
*Bioinformatics Algorithms: an Active Learning Approach*  
 ©2013 by Compeau and Pevzner. All rights reserved

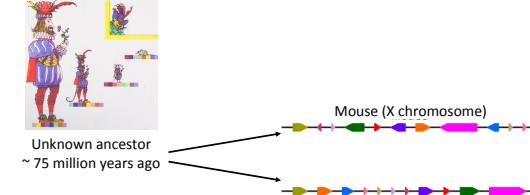
### Are There Fragile Regions in the Human Genome?

- **Transforming Men into Mice**

- Sorting by Reversals
- Breakpoint Theorem
- Rearrangements in Tumor Genomes
- 2-Break Distance Problem
- Breakpoint Graphs
- 2-Break Distance Theorem
- Rearrangement Hotspots in the Human Genome
- Synteny Block Construction

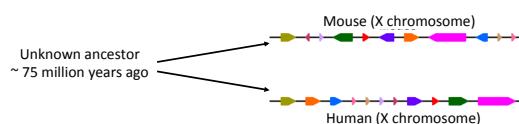


### Genome Rearrangements



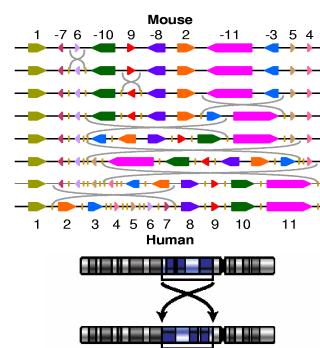
- What are the similarity blocks and how to find them?
- What is the evolutionary scenario for transforming one genome into the other?

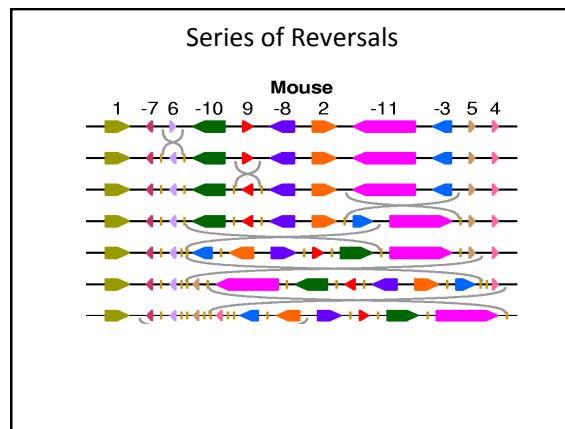
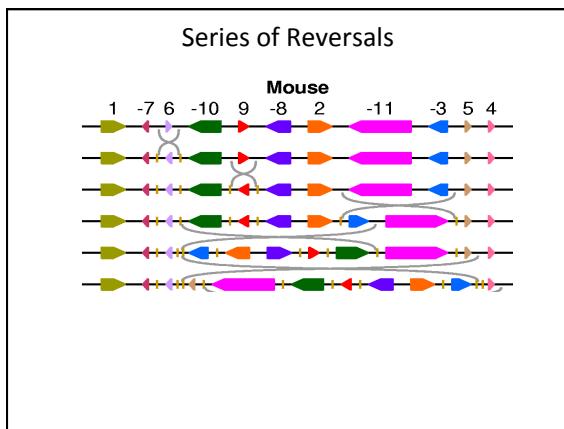
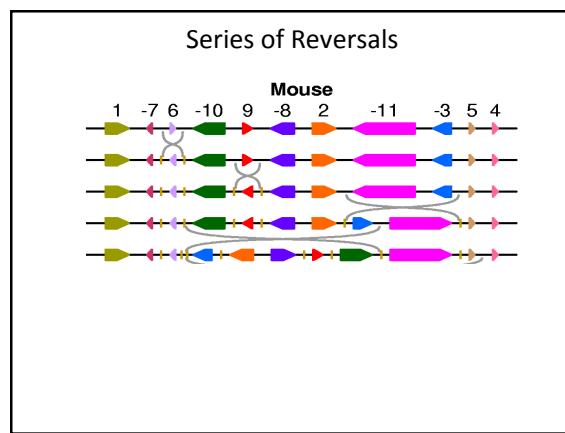
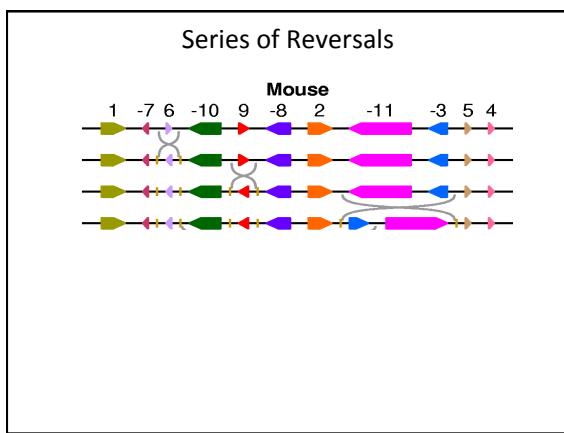
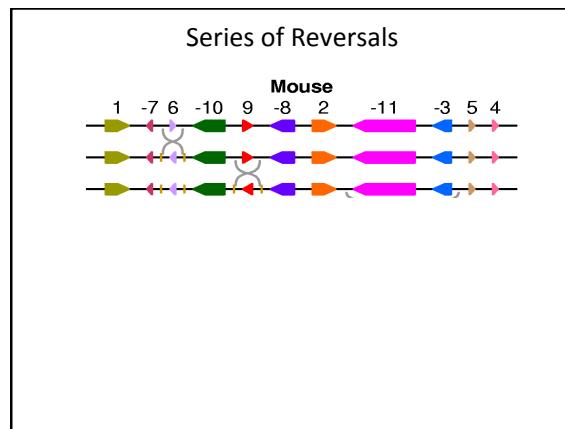
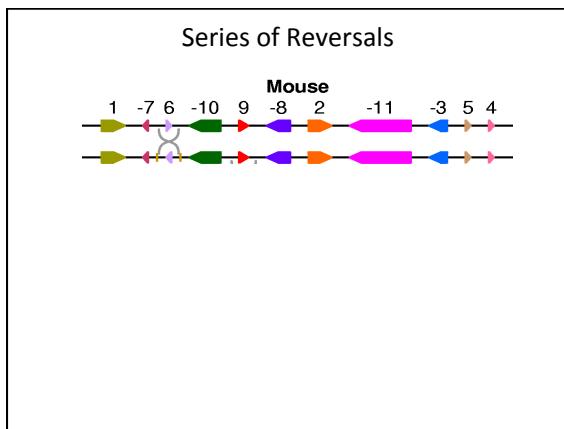
### Genome rearrangements

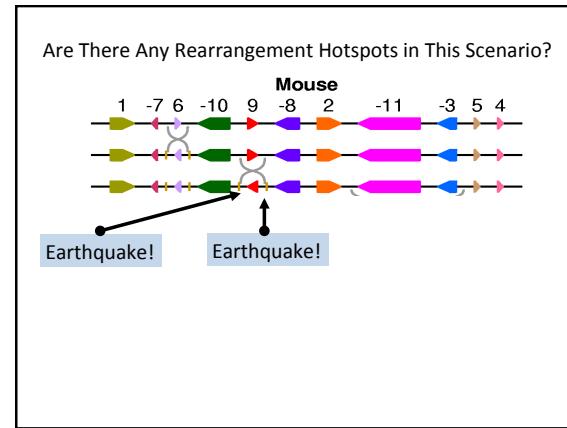
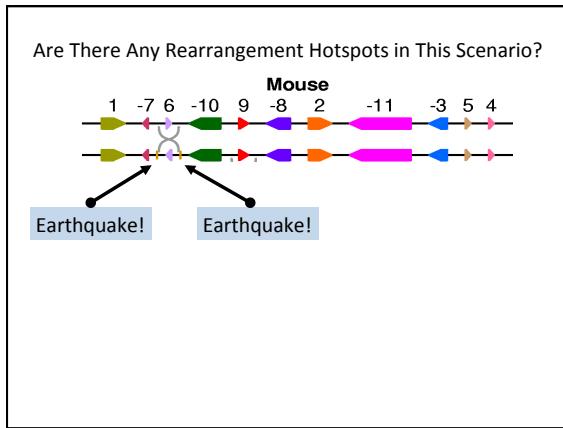
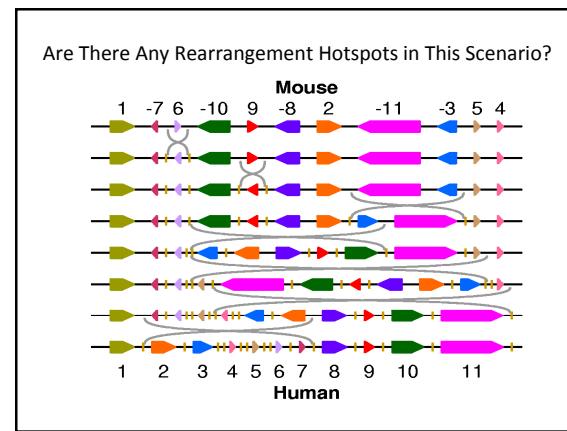
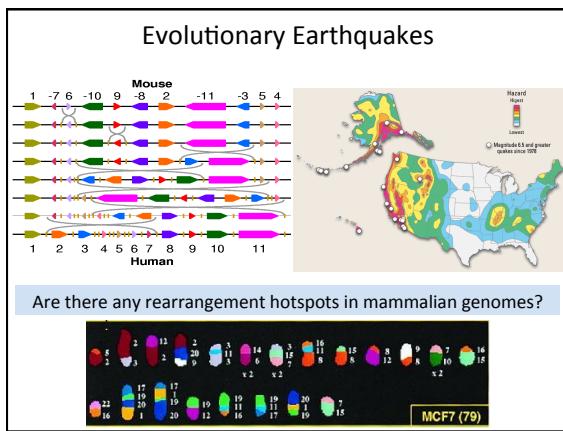
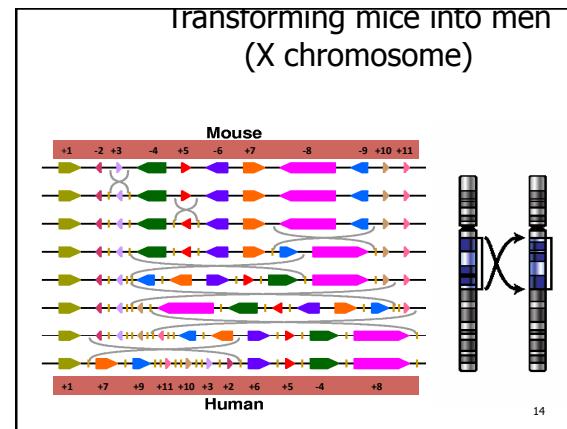
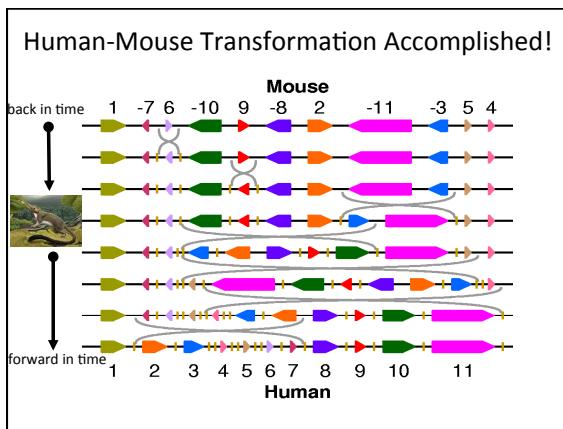


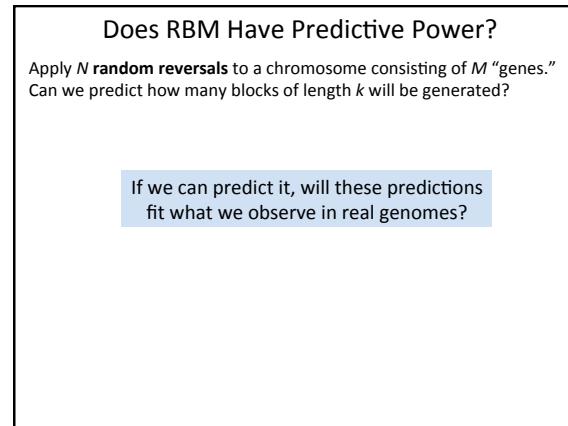
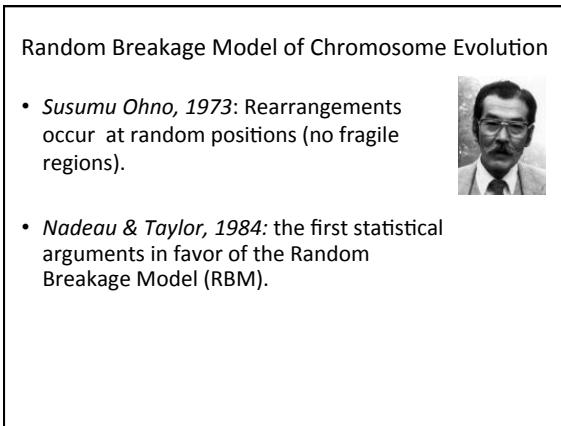
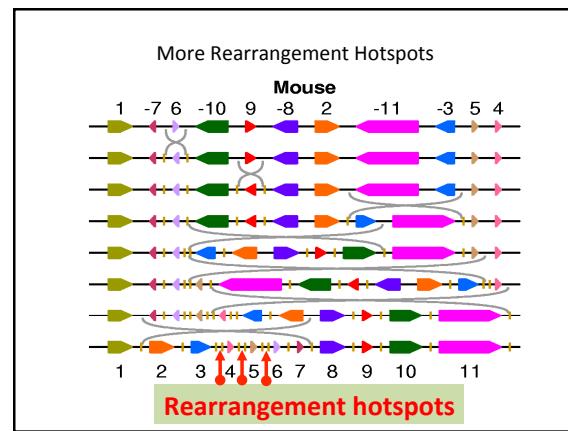
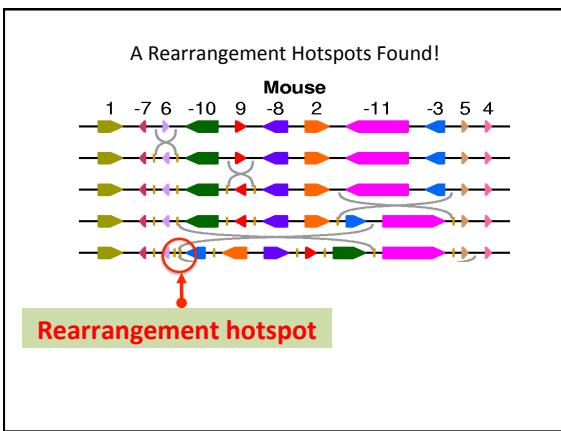
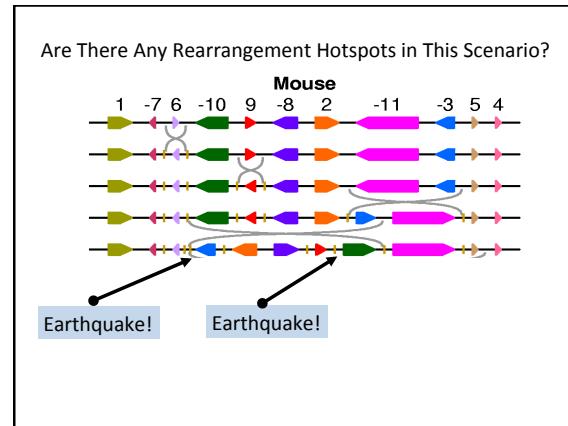
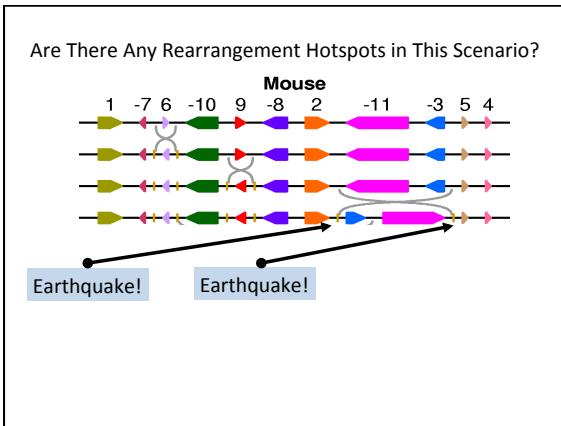
- What are the similarity blocks and how to find them?
- What is the evolutionary scenario for transforming one genome into the other?

### Transforming Mice into Men



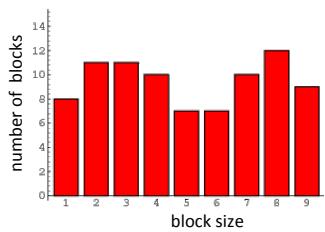






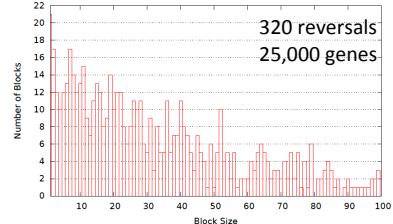
### Does RBM Have Predictive Power?

Apply  $N$  random reversals to a chromosome consisting of  $M$  "genes." Can we predict how many blocks of length  $k$  will be generated?



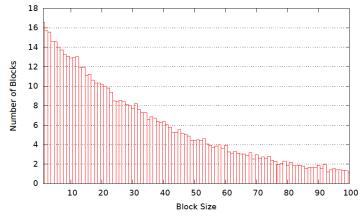
### Does RBM Have Predictive Power?

Apply  $N$  random reversals to a chromosome consisting of  $M$  "genes." Can we predict how many blocks of length  $k$  will be generated?

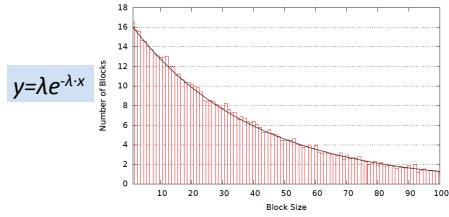


Despite the fact that reversals occur at random positions, we can predict (roughly) how many blocks of each length will be generated!

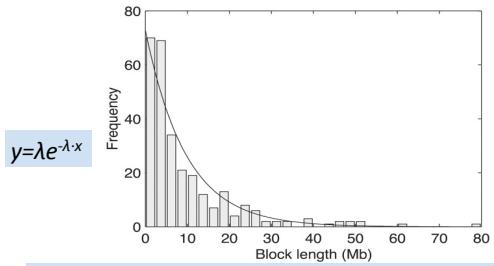
### Histogram of Synteny Block Lengths (averaged over 100 simulations)



### Exponential Distribution (simulated data)



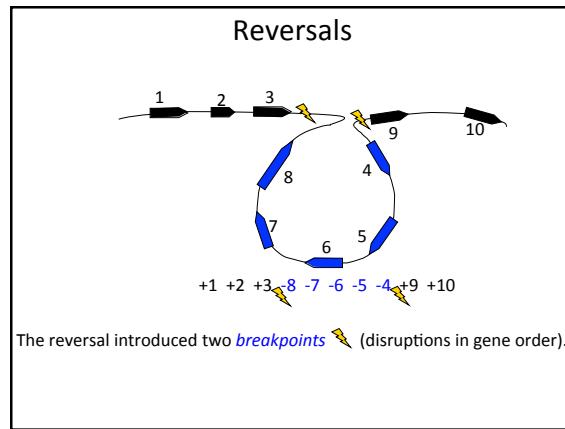
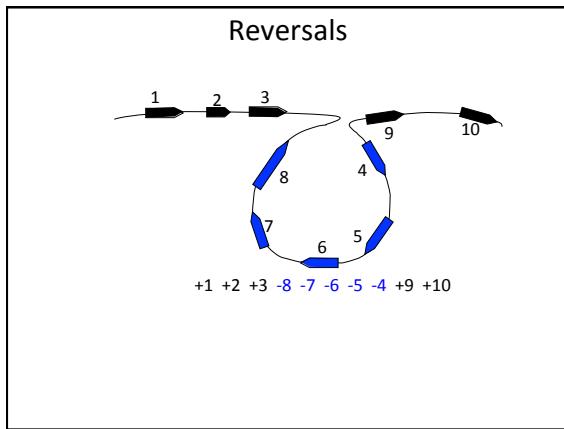
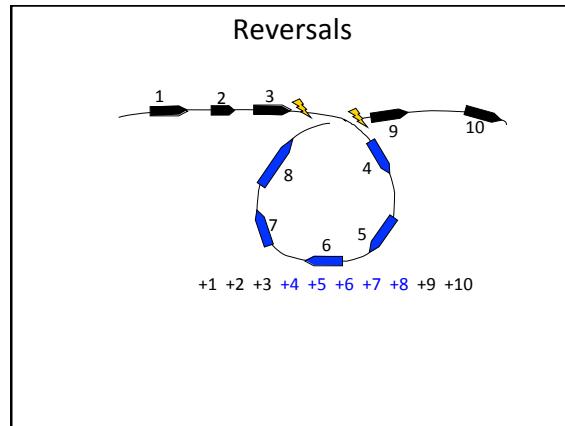
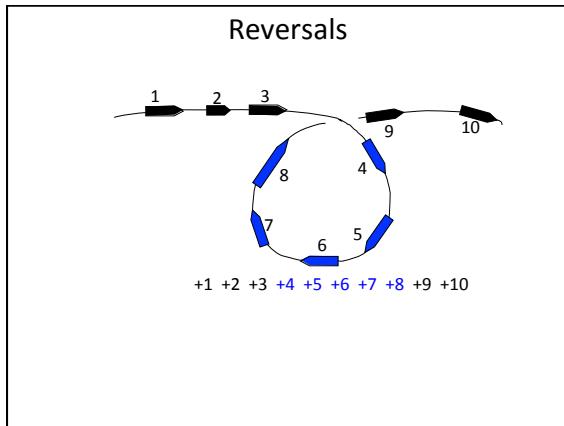
### Lengths of Real Human-Mouse Synteny Blocks



1990s: RBM was embraced by biologists and has become *de facto* theory of chromosome evolution.

### Are There Fragile Regions in the Human Genome?

- Transforming Men into Mice
- **Sorting by Reversals**
- Breakpoint Theorem
- Rearrangements in Tumor Genomes
- 2-Breaks
- Breakpoint Graphs
- 2-Break Distance Theorem
- Rearrangement Hotspots in the Human Genome
- Synteny Block Construction



**Rearrangement Scenario with 5 Reversals**

<b>Step 0:</b>	2	-4	-3	5	-8	-7	-6	1
<b>Step 1:</b>	2	3	4	5	-8	-7	-6	1
<b>Step 2:</b>	2	3	4	5	6	7	8	1
<b>Step 3:</b>	2	3	4	5	6	7	8	-1
<b>Step 4:</b>	-8	-7	-6	-5	-4	-3	-2	-1
<b>Step 5:</b>	1	2	3	4	5	6	7	8

**Rearrangement Scenario with 4 Reversals**

<b>Step 0:</b>	2	-4	-3	5	-8	-7	-6	1
<b>Step 1:</b>	2	3	4	5	-8	-7	-6	1
<b>Step 2:</b>	-5	-4	-3	-2	-8	-7	-6	1
<b>Step 3:</b>	-5	-4	-3	-2	-1	6	7	8
<b>Step 4:</b>	1	2	3	4	5	6	7	8

**Reversal distance:** the minimum number of reversals to transform one permutation into another.

### Sorting by 4 Reversals

**Step 0:**  $2 \ -4 \ -3 \ 5 \ -8 \ -7 \ -6 \ 1$   
**Step 1:**  $\underline{2 \ 3 \ 4} \ 5 \ -8 \ -7 \ -6 \ 1$   
**Step 2:**  $-5 \ -4 \ -3 \ -2 \ \underline{-8 \ -7 \ -6 \ 1}$   
**Step 3:**  $-5 \ -4 \ -3 \ -2 \ \underline{-1 \ 6 \ 7 \ 8}$   
**Step 4:**  $1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8$

**Sorting by Reversals Problem:** Calculate the reversal distance between a permutation and the identity permutation ( $+1 +2 \dots +n$ ).

- **Input:** A permutation  $P$ .
- **Output:** The reversal distance between  $P$  and the identity permutation.

### Greedy Sorting by Reversals

 (+1 -7 +6 -10 +9 -8 +2 -11 -3 +5 +4)

### Greedy Sorting by Reversals

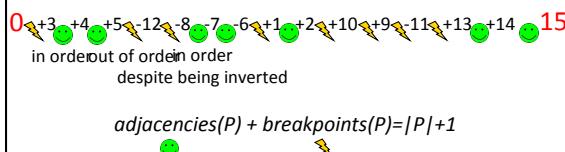
 (+1 -7 +6 -10 +9 -8 +2 -11 -3 +5 +4)  
(+1 -2 +8 -9 +10 -6 +7 -11 -3 +5 +4)  
(+1 +2 +8 -9 +10 -6 +7 -11 -3 +5 +4)  
(+1 +2 +3 +11 -7 +6 -10 +9 -8 +5 +4)  
(+1 +2 +3 -4 -5 +8 -9 +10 -6 +7 -11)  
(+1 +2 +3 +4 -5 +8 -9 +10 -6 +7 -11)  
(+1 +2 +3 +4 +5 +8 -9 +10 -6 +7 -11)  
(+1 +2 +3 +4 +5 +6 -10 +9 -8 +7 -11)  
(+1 +2 +3 +4 +5 +6 -7 +8 -9 +10 -11)  
(+1 +2 +3 +4 +5 +6 +7 +8 -9 +10 -11)  
(+1 +2 +3 +4 +5 +6 +7 +8 +9 +10 -11)  
(+1 +2 +3 +4 +5 +6 +7 +8 +9 +10 +11)



### Are there Fragile Regions in the Human Genome?

- Transforming Men into Mice
- Sorting by Reversals
- **Breakpoint Theorem**
- Rearrangements in Tumor Genomes
- 2-Breaks
- Breakpoint Graphs
- 2-Break Distance Theorem
- Rearrangement Hotspots in the Human Genome
- Synteny Block Construction

### Adjacencies and Breakpoints



What is the number of breakpoints in the identity permutation ( $+1 +2 \dots +n$ )?

### Sorting by Reversals as Breakpoint Elimination

$breakpoints(P)$

**Step 0:**  $2 \ -4 \ -3 \ 5 \ -8 \ -7 \ -6 \ 1 \ 6$   
**Step 1:**  $\underline{2 \ 3 \ 4} \ 5 \ -8 \ -7 \ -6 \ 1 \ 4$   
**Step 2:**  $-5 \ -4 \ -3 \ -2 \ \underline{8 \ -7 \ -6 \ 1} \ 4$   
**Step 3:**  $-5 \ -4 \ -3 \ -2 \ -1 \ 6 \ 7 \ 8 \ 2$   
**Step 4:**  $1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8 \ 0$

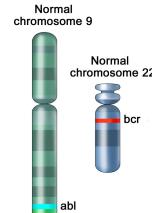
How many breakpoints can be eliminated by a single reversal?

**Breakpoint Theorem**  
**Hint:**   $\geq breakpoints(P)/2$

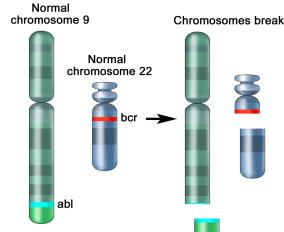
### Are There Fragile Regions in the Human Genome?

- Transforming Men into Mice
- Sorting by Reversals
- Breakpoint Theorem
- **Rearrangements in Tumor Genomes**
- 2-Breaks
- Breakpoint Graphs
- 2-Break Distance Theorem
- Rearrangement Hotspots in the Human Genome
- Synteny Block Construction

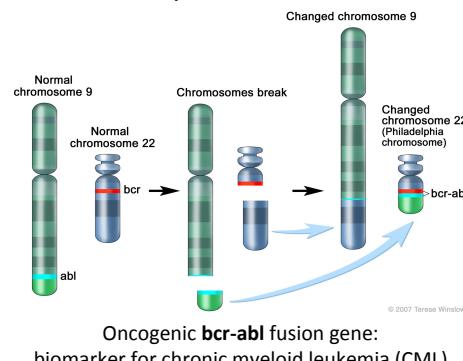
### Philadelphia Chromosome



### Philadelphia Chromosome

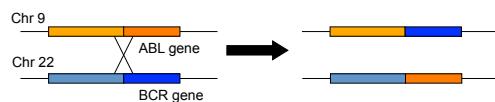


### Philadelphia Chromosome



### Translocations

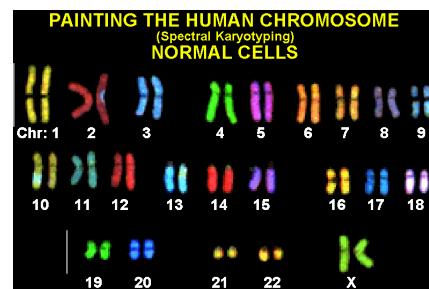
Translocation in CML yields Philadelphia chromosome:



Thousands of rearrangements known for different tumors.

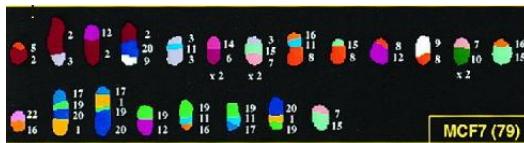


### Painting Chromosomes (Normal Cells)



### Painting Chromosomes (Tumor Cells)

- MCF7 is human breast cancer cell line:

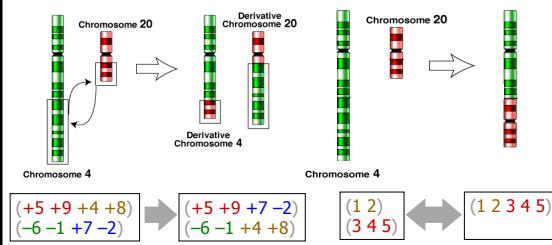


What sequence of rearrangements has produced MCF7?

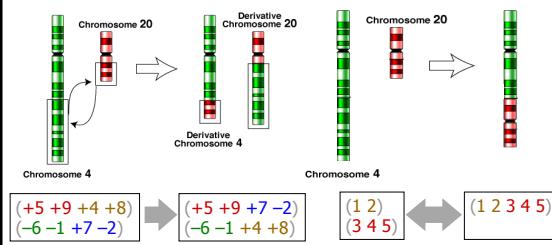
### Rearrangements in Multichromosomal Genomes

#### translocations

Before translocation      After translocation



#### fusions and fissions



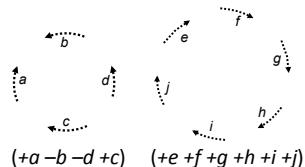
### Are There Fragile Regions in the Human Genome?

- Transforming Men into Mice
- Sorting by Reversals
- Breakpoint Theorem
- Rearrangements in Tumor Genomes
- 2-Break Distance Problem**
- Breakpoint Graphs
- 2-Break Distance Theorem
- Rearrangement Hotspots in the Human Genome
- Synteny Block Construction

### From Linear to Circular Chromosomes

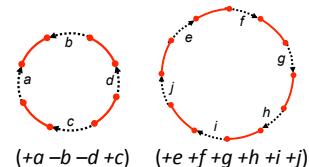
$$(+a -b -d +c) \quad (+e +f +g +h +i +j)$$

### From Linear to Circular Chromosomes

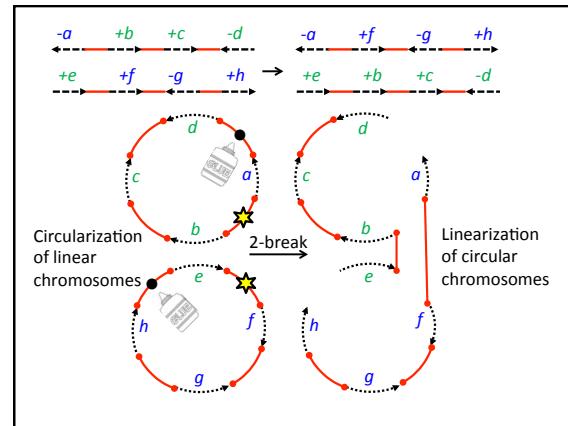
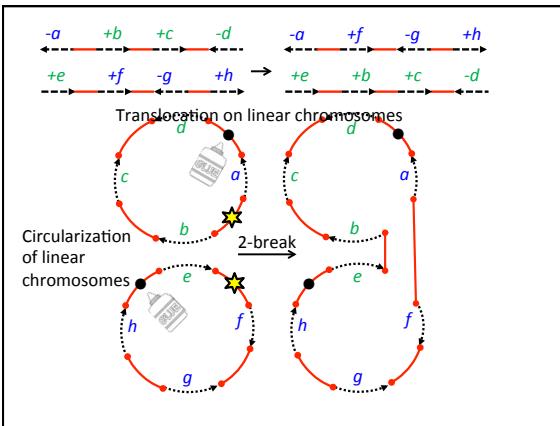
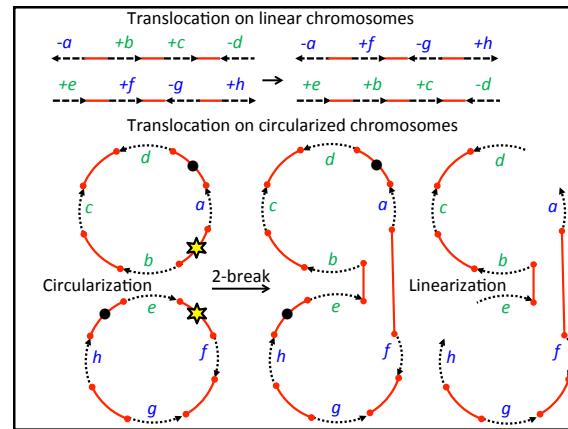
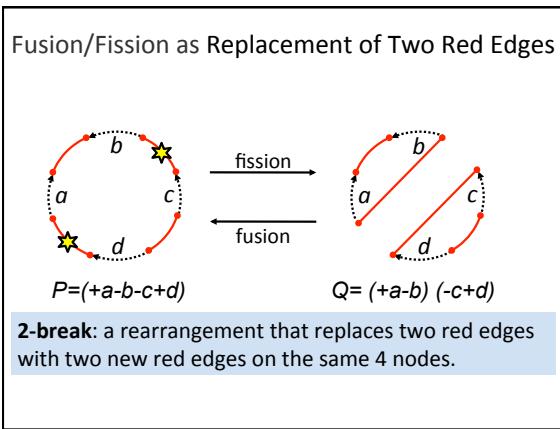
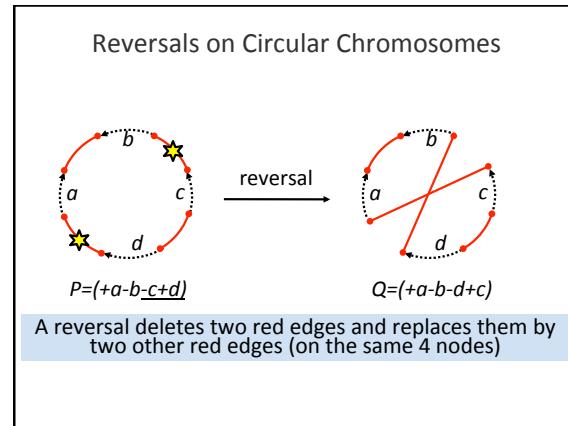
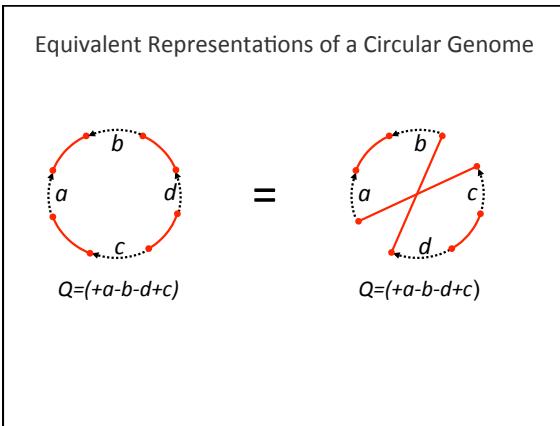


Black directed edges represent synteny blocks.

### From Linear to Circular Chromosomes



Black directed edges represent synteny blocks.  
Red undirected edges connect adjacent synteny blocks.



## 2-Break Distance

### 2-Break distance $d(P, Q)$ :

minimum number of 2-breaks transforming genome  $P$  into genome  $Q$

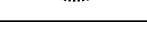
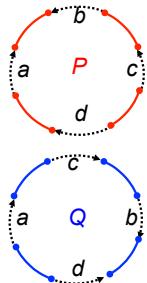
**2-Break Distance Problem.** Find the 2-break distance between two genomes.

- **Input.** Two genomes on the same set of synteny blocks.
- **Output.** The 2-break distance between these genomes.

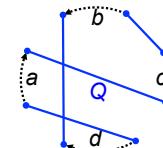
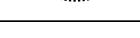
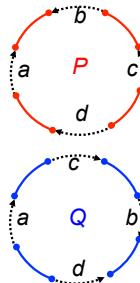
## Are There Fragile Regions in the Human Genome?

- Transforming Men into Mice
- Sorting by Reversals
- Breakpoint Theorem
- Rearrangements in Tumor Genomes
- 2-Break Distance Problem
- **Breakpoint Graphs**
- 2-Break Distance Theorem
- Rearrangement Hotspots in the Human Genome
- Synteny Block Construction

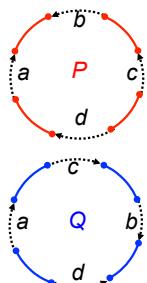
## Comparing Genomes $P$ and $Q$



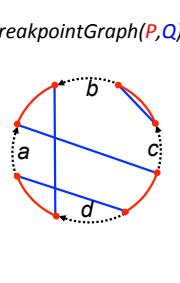
## Different Drawing of $Q$



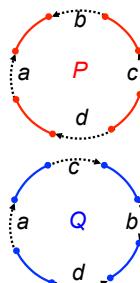
## Superimposing $P$ and $Q$



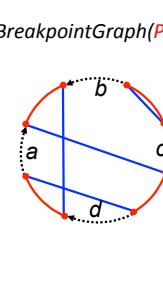
*BreakpointGraph( $P, Q$ )*

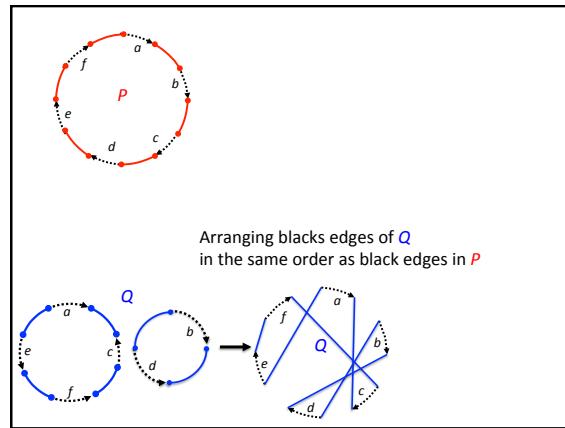
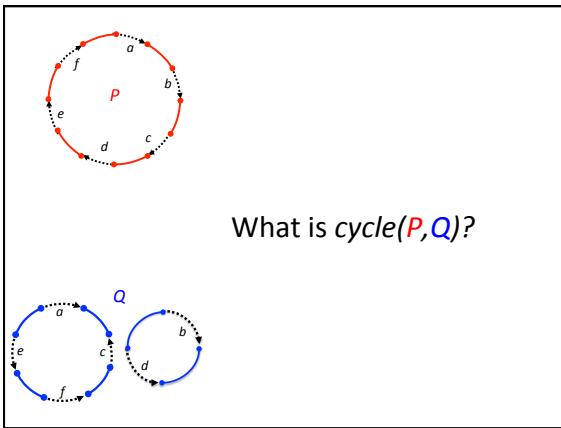
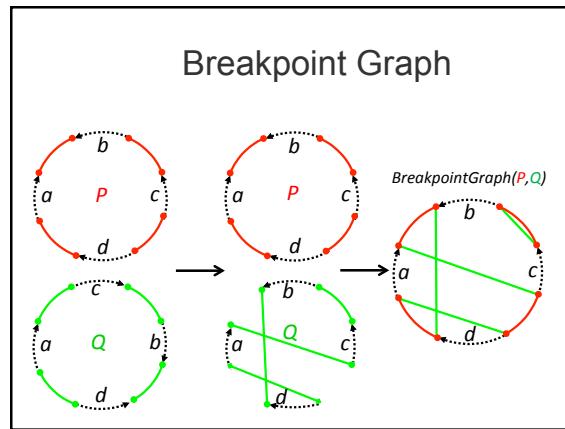
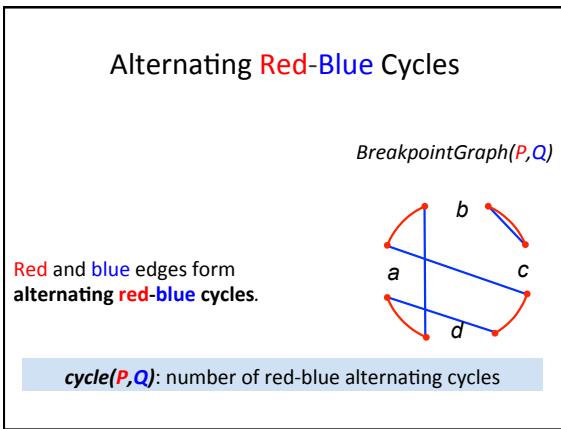
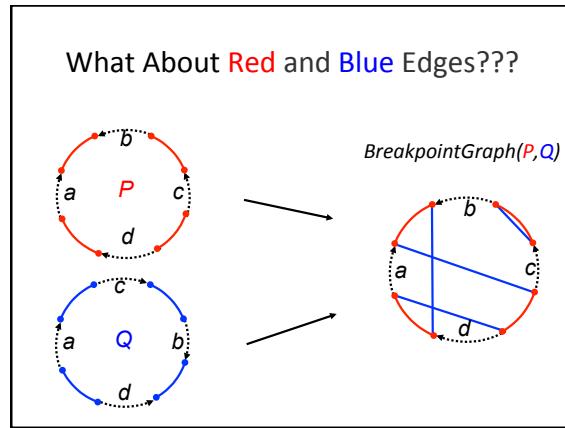
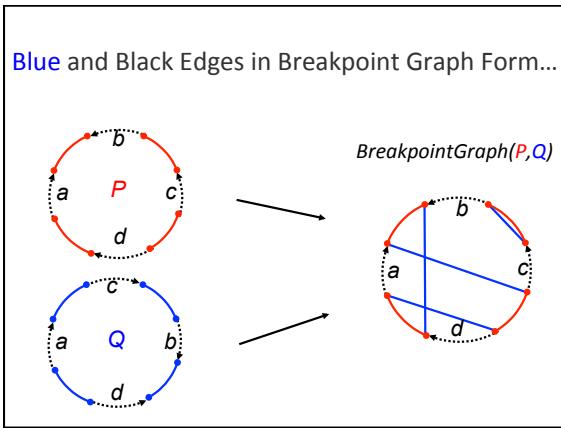


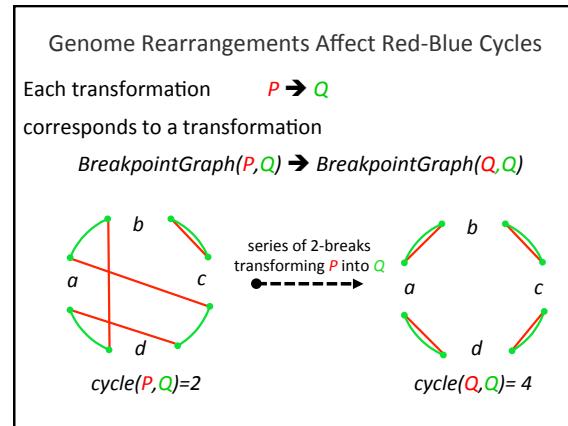
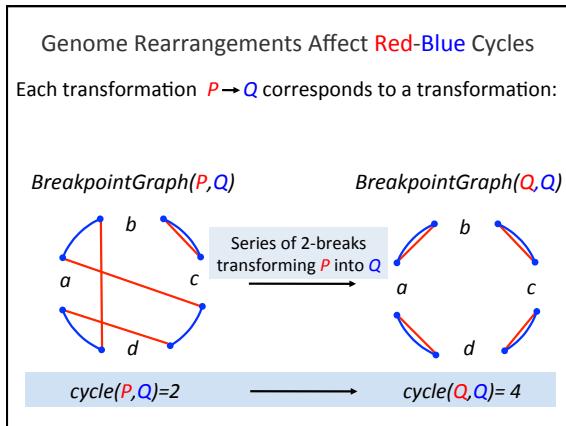
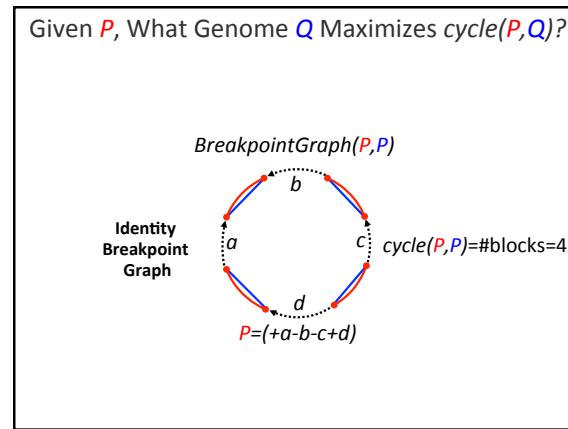
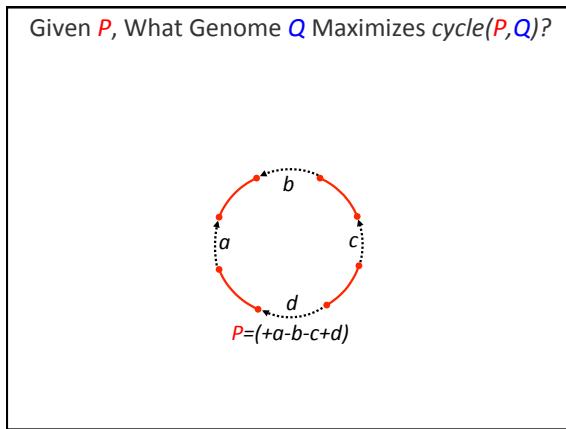
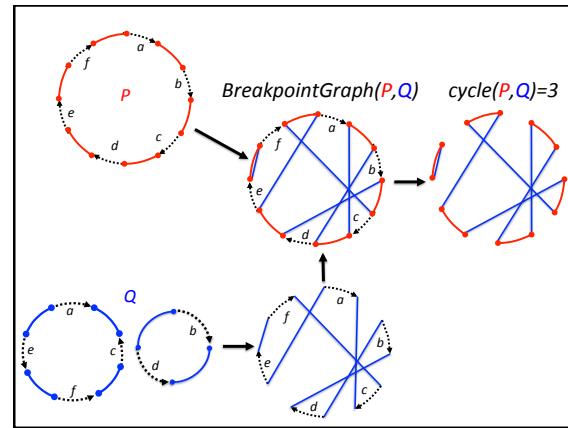
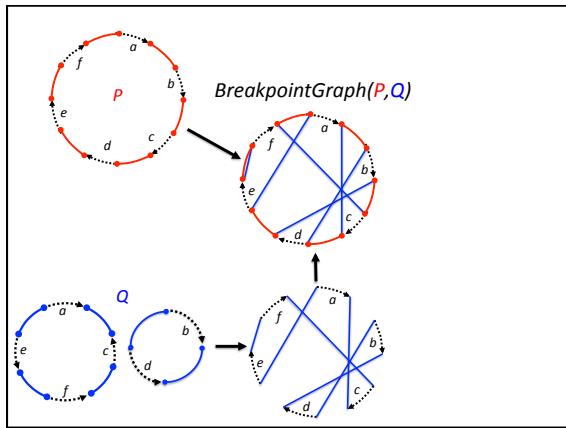
## Red and Black Edges in Breakpoint Graph Form..

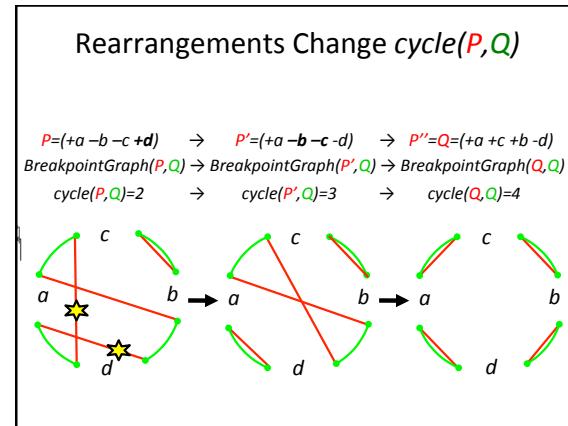
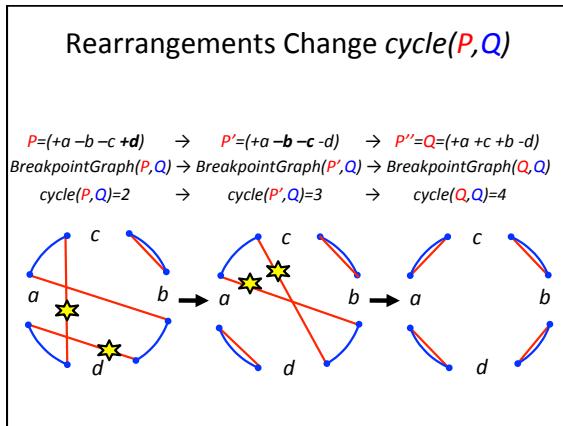


*BreakpointGraph( $P, Q$ )*









- Are There Fragile Regions in the Human Genome?**
- Transforming Men into Mice
  - Sorting by Reversals
  - Breakpoint Theorem
  - Rearrangements in Tumor Genomes
  - 2-Break Distance Problem
  - Breakpoint Graphs
  - **2-Break Distance Theorem**
  - Rearrangement Hotspots in the Human Genome
  - Synteny Block Construction

**Sorting by 2-Breaks**

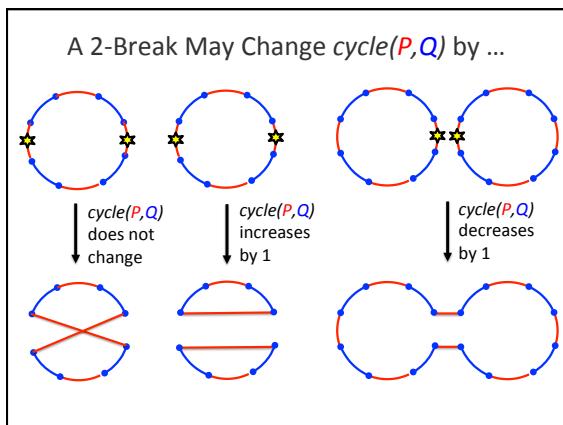
2-breaks  
 $P \rightarrow \dots \rightarrow Q$

$\text{BreakpointGraph}(P, Q) \rightarrow \dots \rightarrow \text{BreakpointGraph}(Q, Q)$

$\text{cycle}(P, Q) \rightarrow \dots \rightarrow \text{cycle}(Q, Q) = \text{blocks}(Q, Q)$

# of red-blue cycles increases by  $\text{blocks}(P, Q) - \text{cycle}(P, Q)$

How much each 2-break can contribute to this increase?



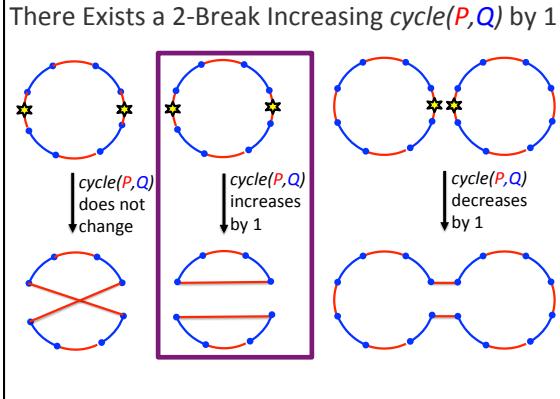
**Each 2-Break Increases #Cycles by at Most 1**

A 2-break:

adds 2 new red edges and thus creates at most 2 new cycles (containing two new red edges)

removes 2 red edges and thus destroys at least 1 old cycle (containing two old edges)

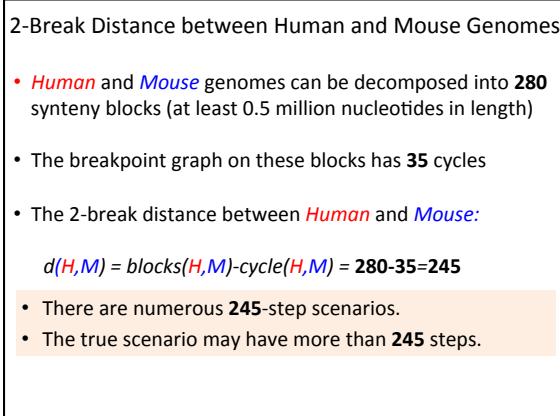
change in the number of cycles  $\leq 2-1=1$ .



### 2-Break Distance Theorem

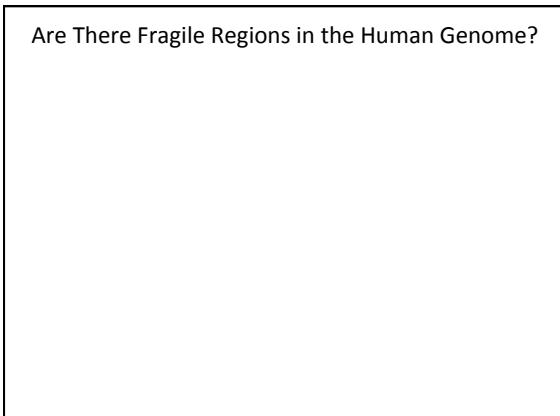
- A 2-break increases #cycles by at most 1.
- There exists a 2-break increasing #cycles by 1.
- Every sorting by 2-breaks must increase #cycles by  $blocks(P, Q) - cycle(P, Q)$
- 2-break distance between genomes  $P$  and  $Q$ :

$$d(P, Q) = blocks(P, Q) - cycle(P, Q)$$



### Are There Fragile Regions in the Human Genome?

- Transforming Men into Mice
- Sorting by Reversals
- Breakpoint Theorem
- Rearrangements in Tumor Genomes
- 2-Break Distance Problem
- Breakpoint Graphs
- 2-Break Distance Theorem
- **Rearrangement Hotspots in the Human Genome**
- Synteny Block Construction



### Are There Fragile Regions in the Human Genome?

**Rearrangement Hotspots Theorem:** Yes!

### Are There Fragile Regions in the Human Genome?

**Rearrangement Hotspots Theorem:** Yes!

**Proof:** If the Random Breakage Model is correct, then  $N$  rearrangements applied to circular chromosomes will produce approximately  $2N$  synteny blocks.

- Since there are 280 human-mouse synteny blocks, there must have been approximately  $280/2 = 140$  2-breaks on the human-mouse evolutionary path.
- However, the 2-Break Distance Theorem implies that there are at least 245 2-breaks on this path.

**STOP and Think.** Is  $245 \approx 140$ ?

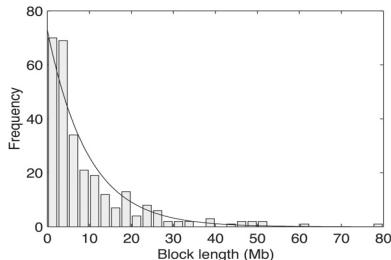
### A Contradiction!

Since 245 is much larger than 140, we arrived at a contradiction implying that **one of our assumptions is incorrect!** Which one?

But the only assumption we made in this proof was:

**"If the Random Breakage Model is correct..."**

### If RBM Is Wrong, How Would You Explain the Exponential Distribution?



### Computational Tests vs. Biological Models

- Why have biologists embraced the Random Breakage Model?

Test

### Computational Tests vs. Biological Models

- Why have biologists embraced the Random Breakage Model?
  - A logical fallacy: RBM is not the only model that complies with the "exponential distribution" test.

Test Model	Exponential distribution
<b>RBM</b>	<b>YES</b>

### Computational Tests vs. Biological Models

- Why have biologists embraced the Random Breakage Model?
  - A logical fallacy: RBM is not the only model that complies with the "exponential distribution" test.
- Why was RBM refuted?
  - It does not comply with the observed "breakpoint reuse."

Test Model	Exponential distribution	Breakpoint reuse
<b>RBM</b>	<b>YES</b>	<b>NO</b>

## A New Model Needed

- Why have biologists embraced the Random Breakage Model?
  - A logical fallacy: RBM is not the only model that complies with the “exponential distribution” test.
- Why was RBM refuted?
  - It does not comply with the “breakpoint reuse” observed in genomes.
- Is there a model that complies with both the “exponential distribution” and the “breakpoint reuse” tests?

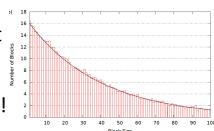
Test Model	Exponential distribution	Breakpoint reuse
RBM	YES	NO
???	YES	YES

## Fragile Breakage Model

- Genome is a mosaic of:
  - fragile regions* with high propensity for rearrangements and
  - solid regions* with low propensity for rearrangements.
- Fragile regions (regions between consecutive synteny blocks) are small, accounting for less than ~ 5% of genome.

## Does FBM Explain BOTH Exponential Distribution and Rearrangement Hotspots?

- A small number of short fragile regions explain rearrangement hotspots.
- If the fragile regions are somewhat randomly distributed throughout the genome, the synteny blocks follow the exponential distribution!

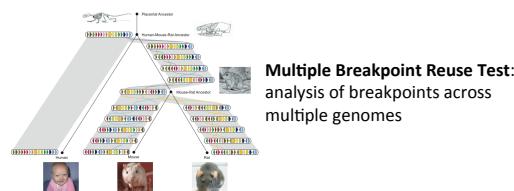


Test Model	Exponential distribution	Breakpoint reuse	A New Test?
RBM	YES	NO	NO
FBM	YES	YES	NO

## But Is There a Test that **BOTH** RBM and FBM Fail?

Test Model	Exponential distribution	Breakpoint reuse	A New Test?
RBM	YES	NO	NO
FBM	YES	YES	NO

## Information About Multiple Genomes Enables a New Test



Test Model	Exponential distribution	Breakpoint reuse	A New Test?
RBM	YES	NO	NO
FBM	YES	YES	NO

## Birth and Death of Fragile Regions

- Recent studies revealed evidence for the “*birth and death*” of the fragile regions, implying that they move to different locations in different lineages.
- This discovery resulted in the **Turnover Fragile Breakage Model (TFBM)** that complies with a new Multiple Breakpoint Reuse (MBR) Test.
- TFBM points to locations of the *currently* fragile regions.

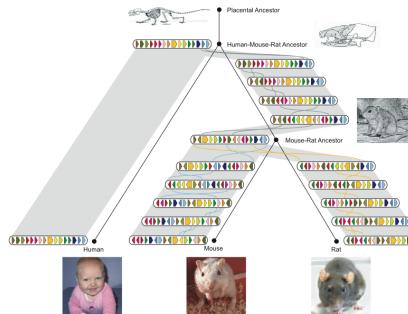


### Birth and Death of Fragile Regions

- Recent studies revealed evidence for the “***birth and death***” of the fragile regions, implying that they move to different locations in different lineages.
- This discovery resulted in the ***Turnover Fragile Breakage Model (TFBM)*** that complies with a new Multiple Breakpoint Reuse (**MBR**) Test.
- TFBM points to locations of the ***currently*** fragile regions.

Test Model \	Exponential distribution	Breakpoint reuse	MBR
RBM	YES	NO	NO
FBM	YES	YES	NO
TFBM	YES	YES	YES

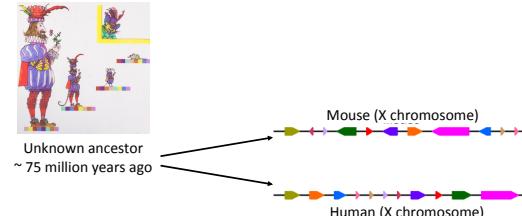
### Where Are the Fragile Regions Located? What Causes Fragility?



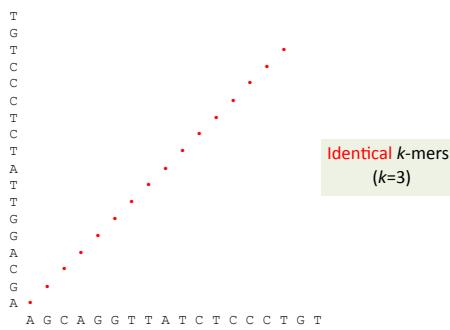
### Are There Fragile Regions in the Human Genome?

- Transforming Men into Mice
- Sorting by Reversals
- Breakpoint Theorem
- Rearrangements in Tumor Genomes
- 2-Break Distance Problem
- Breakpoint Graphs
- 2-Break Distance Theorem
- Rearrangement Hotspots in the Human Genome
- Synteny Block Construction**

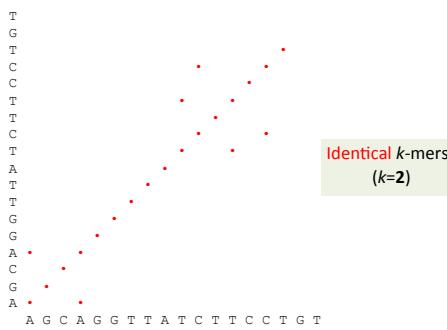
### How Do We Construct the Synteny Blocks?

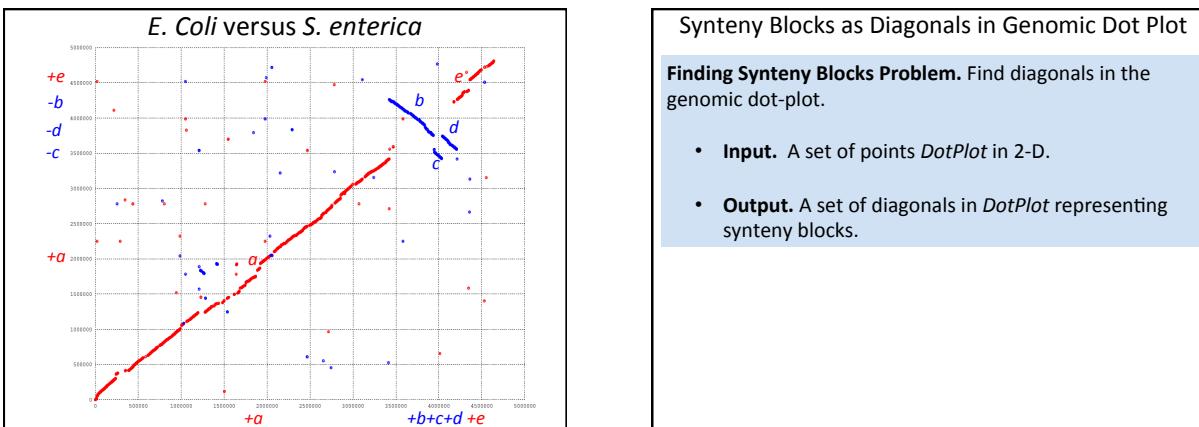
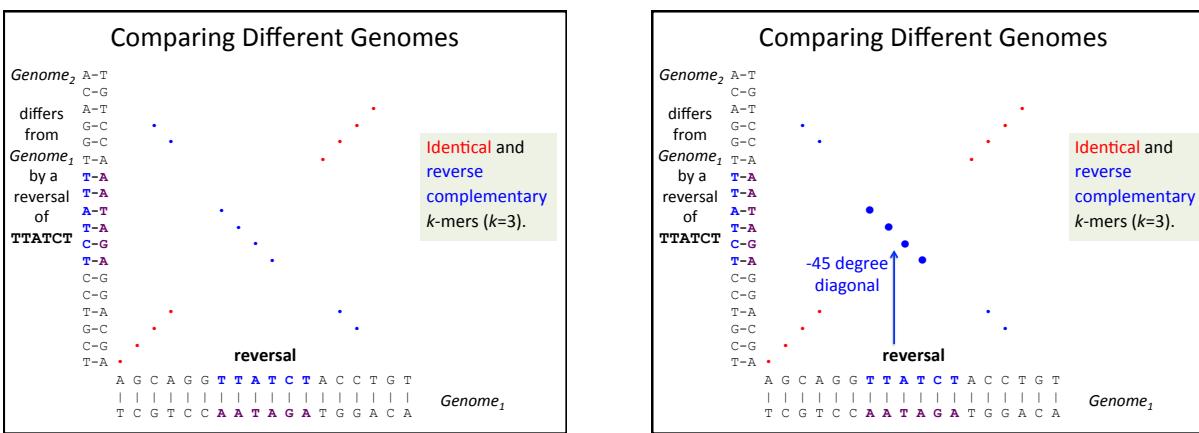
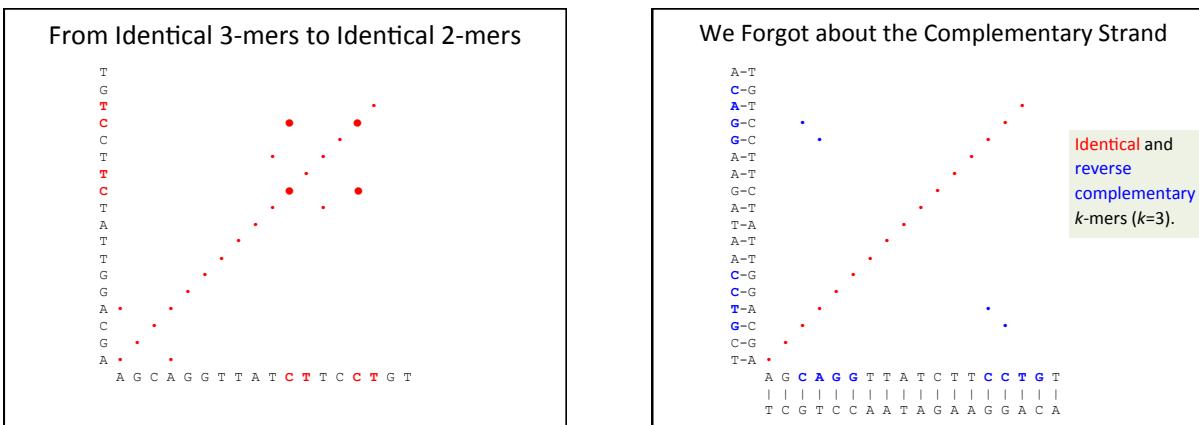


### Comparing Genome with Itself (in 2-D)



### From Identical 3-mers to Identical 2-mers





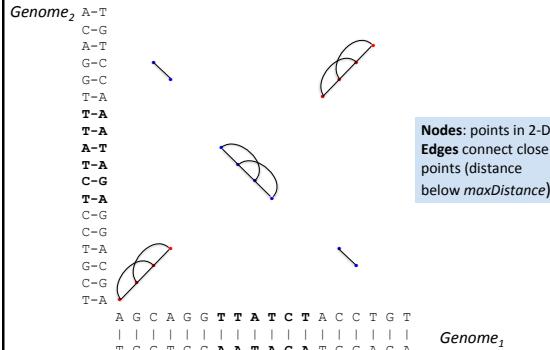
### Synteny Blocks as Diagonals in Genomic Dot Plot

**Finding Synteny Blocks Problem.** Find diagonals in the genomic dot-plot.

- **Input.** A set of points *DotPlot* in 2-D.
- **Output.** A set of diagonals in *DotPlot* representing synteny blocks.



### Connecting Closely Located Points in Genomic Dot-Plot



### Synteny Block Generation Algorithm

**Synteny**(*DotPlot, maxDistance, minSize*)

*maxDistance*: gap size

*minSize*: minimum synteny block size

- Form a graph whose node set is the set of points in *DotPlot*
- Connect two nodes by an edge if the 2-D distance between them is < *maxDistance*. The connected components in the resulting graph define synteny blocks
- Delete small synteny blocks (length < *minSize*)

### Another Synteny Block Generation Algorithm

**Amalgamate**(*DotPlot, maxDistance, minSize*)

*maxDistance*: gap size

*minSize*: minimum synteny block size

- Define each point in *DotPlot* as a separate block and iteratively amalgamate the resulting blocks
- Amalgamate two blocks if they contain two points that are separated by < *maxDistance* in another genome.
- Delete small synteny blocks (length < *minSize*)

### Two Synteny Block Generation Algorithms: Which One is Better?

**Synteny**(*DotPlot, maxDistance, minSize*)

- Form a graph whose node set is the set of points in *DotPlot*
- Connect two nodes by an edge if the 2-D distance between them is < *maxDistance*. The connected components in the resulting graph define synteny blocks
- Delete small synteny blocks (length < *minSize*)

**Amalgamate**(*DotPlot, maxDistance, minSize*)

- Define each point in *DotPlot* as a separate block and iteratively amalgamate the resulting blocks
- Amalgamate two blocks if they contain two points that are separated by < *maxDistance* in another genome.
- Delete small synteny blocks (length < *minSize*)

The algorithms look very similar but do they produce similar results?

### Genomic Dot Plot (Human vs. Mouse X Chromosome)

