

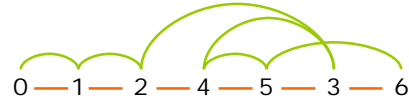
## How good is simple reversal sort?

- ρ Not so good actually
- ρ It has to do at most  $n-1$  reversals with permutation of length  $n$
- ρ The algorithm can return a distance that is as large as  $(n-1)/2$  times the correct result  $d(\Pi)$ 
  - For example, if  $n = 1001$ , result can be as bad as  $500 \times d(\Pi)$

311

## Estimating reversal distance by cycle decomposition

- ρ We can estimate  $d(\Pi)$  by *cycle decomposition*
- ρ Lets represent permutation  $\Pi = 1\ 2\ 4\ 5\ 3$  with the following graph

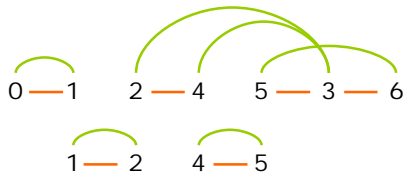


where edges correspond to adjacencies (identity, permutation  $F$ )

312

## Estimating reversal distance by cycle decomposition

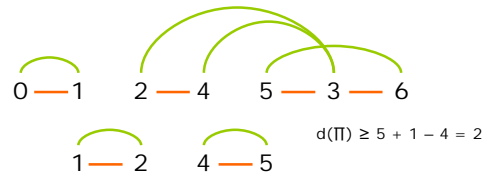
- ρ Cycle decomposition: a set of cycles that
  - have edges with alternating colors
  - do not share edges with other cycles (=cycles are edge disjoint)



313

## Cycle decompositions

- ρ Let  $c(\Pi)$  the maximum number of alternating, edge-disjoint cycles in the graph representation of  $\Pi$
- ρ The following formula allows estimation of  $d(\Pi)$ 
  - $d(\Pi) \geq n + 1 - c(\Pi)$ , where  $n$  is the permutation length



Claim in Deonier: equality holds for "most of the usual and interesting biological systems."

314

## Cycle decompositions

- ρ Cycle decomposition is NP-complete
  - We cannot solve the general problem exactly for large instances
- ρ However, with signed data the problem becomes easy
  - Before going into signed data, lets discuss another algorithm for the general case

315

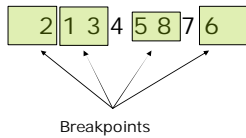
## Computing reversals with breakpoints

- ρ Lets investigate a better way to compute reversal distance
- ρ First, some concepts related to permutation  $\Pi_1 \Pi_2, \dots, \Pi_{n-1} \Pi_n$ 
  - Breakpoint: two elements  $\Pi_i$  and  $\Pi_{i+1}$  are a *breakpoint*, if they are not consecutive numbers
  - Adjacency: if  $\Pi_i$  and  $\Pi_{i+1}$  are consecutive, they are called *adjacency*

316

## Breakpoints and adjacencies

This permutation contains four breakpoints *begin-2*, 13, 58, 6-*end* and five adjacencies 21, 34, 45, 87, 76



317

## Breakpoints

- Each breakpoint in permutation needs to be removed to get to the identity permutation (=our target)
  - Identity permutation does not contain any breakpoints

$$[2, 1, 3, 4, 5, 8, 7, 6] \quad b(\pi) = 4$$

- First and last positions special cases
- Note that each reversal can remove *at most* two breakpoints
- Denote the number of breakpoints by  $b(\pi)$

318

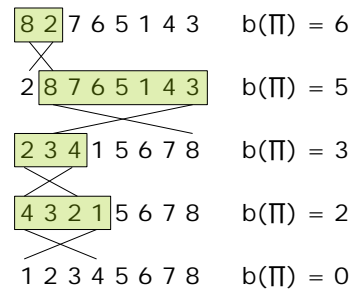
## Breakpoint reversal sort

Idea: try to remove as many breakpoints as possible (max 2) in every step

- While  $b(\pi) > 0$
- Choose reversal  $p$  that removes most breakpoints
- Perform reversal  $p$  to  $\pi$
- Output  $\pi$
- return

319

## Breakpoint removal: example



320

## Breakpoint removal

- The previous algorithm needs refinement to be correct
- Consider the following permutation:

1 5 6 7 2 3 4 8

- There is no reversal that decreases the number of breakpoints!
- See Jones & Pevzner for detailed description on this

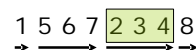
321

Strip: maximal segment without breakpoints

## Breakpoint removal

→ Increasing strip  
← Decreasing strip

- Reversal can only decrease breakpoint count if permutation contains *decreasing strips*



322

## Improved breakpoint reversal sort

1. While  $b(\Pi) > 0$
2.   If  $\Pi$  has a decreasing strip
3.     Do reversal  $p$  that removes most BPs
4.   Else
5.     Reverse an increasing strip
6.   Output  $\Pi$
7. return

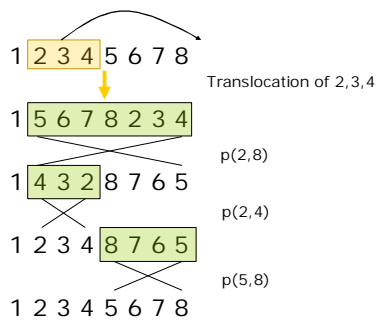
323

## Is Improved BP removal enough?

- p The algorithm works pretty well:
  - n It produces a result that is at most four times worse than the optimal result
  - n ...is this good?
- p We considered only reversals
- p What about translocations & duplications?

324

## Translocations via reversals



325

## Genome rearrangements with reversals

- p With *unsigned* data, the problem of finding minimum reversal distances is *NP-complete*
  - n Why is this so if sorting is easy?
- p An algorithm has been developed that achieves 1.375-approximation
- p However, reversal distance in *signed data* can be computed quickly!
  - n It takes linear time w.r.t. the length of permutation (Bader, Moret, Yan, 2001)

326

## Cycle decomposition with signed data

- p Consider the following two permutations that include *orientation* of markers
  - n  $J$ : +1 +5 -2 +3 +4
  - n  $K$ : +1 -3 +2 +4 -5
- p We modify this representation a bit to include both endpoints of each marker:
  - n  $J'$ : 0 1a 1b 5a 5b 2b 2a 3a 3b 4a 4b 6
  - n  $K'$ : 0 1a 1b 3b 3a 2a 2b 4a 4b 5b 5a 6

327

## Graph representation of $J'$ and $K'$

- p Drawn online in lecture!

328

## Multiple chromosomes

- ρ In unichromosomal genomes, inversion (reversal) is the most common operation
- ρ In multichromosomal genomes, inversions, translocations, *fissions* and *fusions* are most common

329

## Multiple chromosomes

- ρ Lets represent multichromosomal genome as a set of permutations, with \$ denoting the boundary of a chromosome:

```

5 9 $           Chr 1
1 3 2 8 $      Chr 2
7 6 4 $         Chr 3
    
```

This notation is frequently used in software used to analyse genome rearrangements.

330

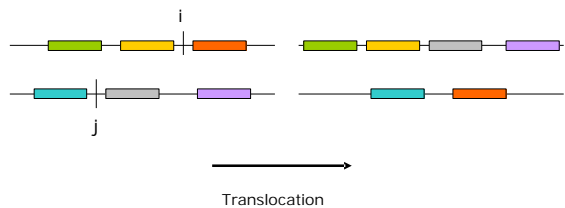
## Multiple chromosomes

- ρ Note that when dealing with multiple chromosomes, you need to specify numbering for elements on both genomes

331

## Reversals & translocations

- ρ Reversal  $p(\Pi, i, j)$
- ρ Translocation  $p(\Pi, \sigma, i, j)$



332

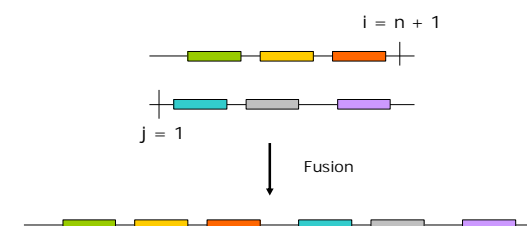
## Fusions & fissions

- ρ Fusion: merging of two chromosomes
- ρ Fission: chromosome is split into two chromosomes
- ρ Both events can be represented with a translocation

333

## Fusion

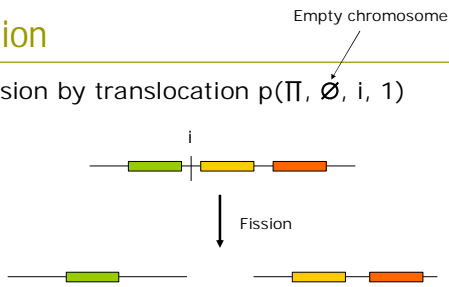
- ρ Fusion by translocation  $p(\Pi, \sigma, n+1, 1)$



334

## Fission

- ρ Fission by translocation  $p(\Pi, \emptyset, i, 1)$



335

## Algorithms for general genomic distance problem

- ρ Hannenhalli, Pevzner: Transforming Men into Mice (polynomial algorithm for genomic distance problem), *36th Annual IEEE Symposium on Foundations of Computer Science*, 1995

336

## Human & mouse revisited

- ρ Human and mouse are separated by about 75-83 million years of evolutionary history
- ρ Only a few hundred rearrangements have happened after speciation from the common ancestry
- ρ Pevzner & Tesler identified in 2003 for 281 synteny blocks a rearrangement from mouse to human with
  - η 149 inversions
  - η 93 translocations
  - η 9 fissions

337

## Discussion

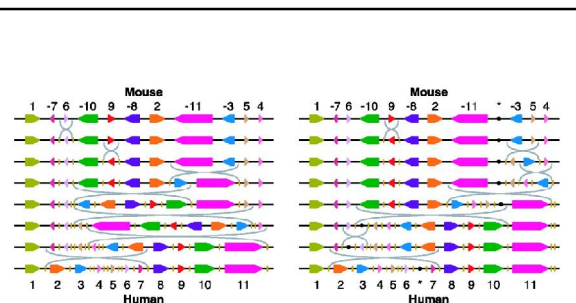
- ρ Genome rearrangement events are very rare compared to, e.g., point mutations
  - η We can study rearrangement events further back in the evolutionary history
- ρ Rearrangements are easier to detect in comparison to many other genomic events
- ρ We cannot detect homologs 100% correctly so the input permutation can contain errors

338

## Discussion

- ρ Genome rearrangement is to some degree constrained by the number and size of repeats in a genome
  - η Notice how the importance of genomic repeats pops up once again
- ρ Sequencing gives us (usually) signed data so we can utilize faster algorithms
- ρ What if there are more than one optimal solution?

339



Two different genome rearrangement scenarios giving the same result.

340

