

# Greedy search

Torgeir R. Hvidsten

## This lecture

- Genome **rearrangements**
  - Sorting by reversals
  - Approximation algorithms
  - Breakpoints: a different face of greed
- Finding **regulatory motifs** in DNA Sequences
- Greedy search methods

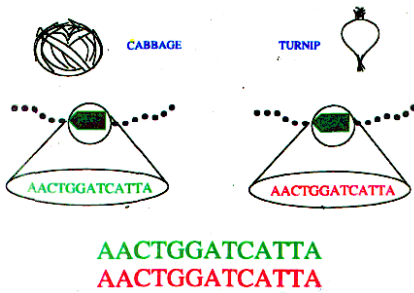
# Genome rearrangements

## Turnip vs cabbage: Look and taste different

Although cabbages and turnips share a recent common ancestor, they look and taste different



### Turnip vs cabbage: Comparing gene sequences yields no evolutionary information



### Turnip vs cabbage: Almost identical mtDNA gene sequences

- In 1980s Jeffrey Palmer studied evolution of plant organelles by comparing mitochondrial genomes of the cabbage and turnip
- 99% similarity between genes
- These surprisingly identical gene sequences **differed in gene order**
- This study helped pave the way to analyzing genome rearrangements in molecular evolution

### Turnip vs cabbage: Different mtDNA gene order

Gene order comparison:



### Turnip vs cabbage: Different mtDNA gene order

Gene order comparison:



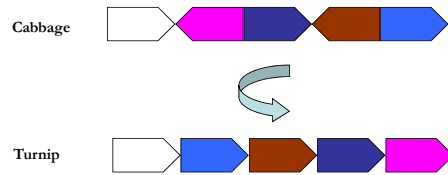
### Turnip vs cabbage: Different mtDNA gene order

Gene order comparison:



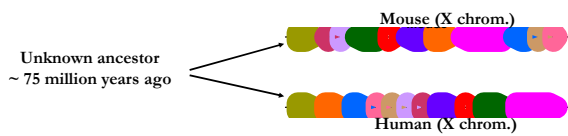
### Turnip vs cabbage: Different mtDNA gene order

Gene order comparison:



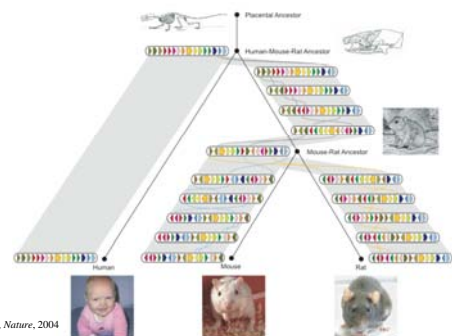
Evolution is manifested as the divergence in gene order

### Genome rearrangements

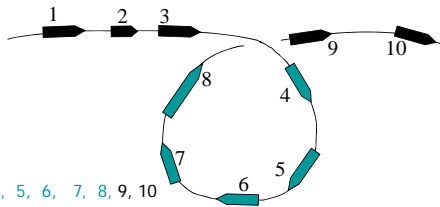


- What are the similarity blocks and how can they be found?
- What is the architecture of the ancestral genome?
- What is the evolutionary scenario for transforming one genome into the other?

### History of chromosome X



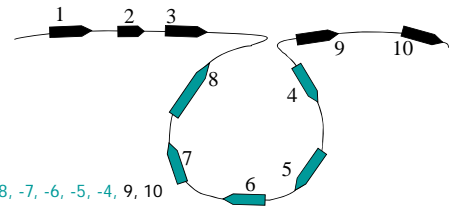
## Reversals



1, 2, 3, 4, 5, 6, 7, 8, 9, 10

Blocks represent conserved genes

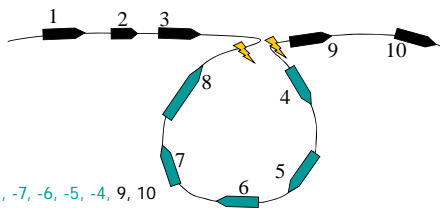
## Reversals



1, 2, 3, -8, -7, -6, -5, -4, 9, 10

In the course of evolution, blocks 1,...,10 could be misread as 1, 2, 3, -8, -7, -6, -5, -4, 9, 10

## Reversals and breakpoints



1, 2, 3, -8, -7, -6, -5, -4, 9, 10

The reversion introduced two **breakpoints** (disruptions in order)

## Reversals: Example



Break and Invert

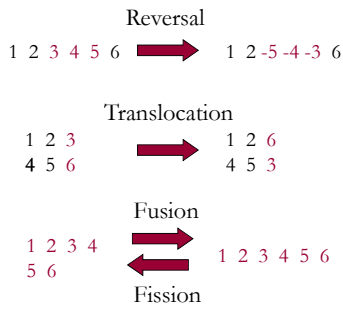
5' ATG**CCTG**ACTA 3'

3' TAC**GGAC**ATGAT 5'

5' ATGTAC**AGG**CCTA 3'

3' TACAT**GTCC**GAT 5'

## Types of rearrangements

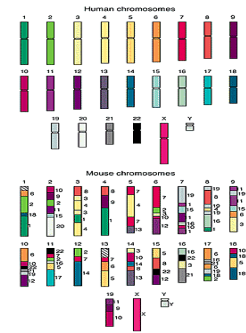


T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

17

## Comparative genomic architectures: Mouse vs human genome

- Humans and mice have similar genomes, but their genes are ordered differently
- ~245 rearrangements
  - Reversals
  - Fusions
  - Fissions
  - Translocation



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

18

## Waardenburg's syndrome: Mouse provides insight into human genetic disorder

- Waardenburg's syndrome is characterized by pigmentary dysphasia
- Gene implicated in the disease was linked to human chromosome 2 but it was not clear where exactly it is located on chromosome 2



T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

19

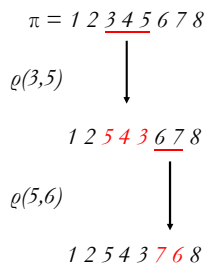
## Waardenburg's syndrome and splotch mice

- A breed of mice (with splotch gene) had similar symptoms caused by the same type of gene as in humans
- Scientists succeeded in identifying the location of the gene responsible for disorder in mice
- Finding the gene in mice gives clues to where the same gene is located in humans by analyzing the relative architecture of human and mouse genomes

T.R. Hvidsten: 1MB304: Discrete structures for bioinformatics II

20

## Reversals: Example



## Reversals and gene orders

- Gene order is represented by a permutation  $\pi$ :

$$\pi = \pi_1 \dots \pi_{i-1} \underline{\pi_i \pi_{i+1} \dots \pi_{j-1} \pi_j \pi_{j+1}} \dots \pi_n$$

$$\downarrow \rho(i,j)$$

$$\pi \cdot \rho(i,j) = \pi_1 \dots \pi_{i-1} \underline{\pi_j \pi_{j-1} \dots \pi_{i+1} \pi_i} \pi_{j+1} \dots \pi_n$$

- Reversal  $\rho(i,j)$  reverses (flips) the elements from  $i$  to  $j$  in  $\pi$

## Reversal distance problem

- **Goal:** Given two permutations, find the shortest series of reversals that transforms one into another
- **Input:** Permutations  $\pi$  and  $\sigma$
- **Output:** A series of reversals  $\rho_1, \dots, \rho_t$  transforming  $\pi$  into  $\sigma$ , such that  $t$  is minimum
- $t$  - reversal distance between  $\pi$  and  $\sigma$
- $d(\pi, \sigma)$  - smallest possible value of  $t$ , given  $\pi$  and  $\sigma$

## Sorting by reversals problem

- **Goal:** Given a permutation, find a shortest series of reversals that transforms it into the identity permutation  $(1\ 2\ \dots\ n)$
- **Input:** Permutation  $\pi$
- **Output:** A series of reversals  $\rho_1, \dots, \rho_t$  transforming  $\pi$  into the identity permutation such that  $t$  is minimum

## Sorting by reversals: Example

➤  $t = d(\pi)$  - reversal distance of  $\pi$

➤ Example :

$$\pi = \begin{array}{cccccccc} \underline{3} & \underline{4} & 2 & 1 & 5 & 6 & 7 & 10 & 9 & 8 \\ 4 & 3 & 2 & 1 & 5 & 6 & 7 & \underline{10} & \underline{9} & 8 \\ \underline{4} & \underline{3} & \underline{2} & 1 & 5 & 6 & 7 & 8 & 9 & 10 \\ 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 \end{array}$$

So  $d(\pi) = 3$

## Sorting by reversals: A greedy algorithm

➤ When sorting permutation  $\pi = 1 \ 2 \ 3 \ 6 \ 4 \ 5$ , the first three elements are already in order so it makes no sense to break them

➤ The length of the already sorted prefix of  $\pi$  is denoted  $prefix(\pi)$

-  $prefix(\pi) = 3$

➤ This results in an idea for a greedy algorithm: increase  $prefix(\pi)$  at every step

## Greedy algorithm: An example

➤ Sorting  $\pi$ :

$$\begin{array}{cccc} 1 & 2 & 3 & \underline{6} & 4 & 5 \\ & & & \downarrow & & \\ 1 & 2 & 3 & 4 & \underline{6} & 5 \\ & & & & \downarrow & \\ 1 & 2 & 3 & 4 & 5 & 6 \end{array}$$

➤ Number of steps to sort permutation of length  $n$  is at most  $(n - 1)$

## Greedy Algorithm: Pseudo-code

SimpleReversalSort( $\pi$ )

```

1 for  $i \leftarrow 1$  to  $n - 1$ 
2    $j \leftarrow$  Index of element  $i$  in  $\pi$  (i.e.,  $\pi_j = i$ )
3   if  $j \neq i$ 
4      $\pi \leftarrow \pi \cdot \rho(i, j)$ 
5   output  $\pi$ 
6 if  $\pi$  is the identity permutation
7   return
    
```

## Analyzing SimpleReversalSort (I)

- SimpleReversalSort does not guarantee the smallest number of reversals
- It takes five steps on  $\pi = 612345$ :
  - Step 1:  $162345$
  - Step 2:  $126345$
  - Step 3:  $123645$
  - Step 4:  $123465$
  - Step 5:  $123456$

## Analyzing SimpleReversalSort (II)

- But it can be sorted in two steps:  $\pi = 612345$ 
  - Step 1:  $543216$
  - Step 2:  $123456$
- So, SimpleReversalSort( $\pi$ ) is not optimal
- SimpleReversalSort is not **correct** (according to the definition in lecture 1)
- Optimal/correct algorithms are unknown for many problems; approximation algorithms are used

## Approximation algorithms

- These algorithms find **approximate solutions** rather than **optimal solutions**
- The **approximation ratio** of an algorithm A on input  $\pi$  is:

$$A(\pi) / \text{OPT}(\pi)$$

where

$A(\pi)$  - solution produced by algorithm A  
 $\text{OPT}(\pi)$  - optimal solution of the problem

## Performance guarantee

- Approximation ratio (**performance guarantee**) of algorithm A is the maximal approximation ratio of all inputs of size  $n$
- For algorithm A that minimizes the objective function (minimization algorithm):
  - $\max_{|\pi|=n} A(\pi) / \text{OPT}(\pi)$
- For maximization algorithms
  - $\min_{|\pi|=n} A(\pi) / \text{OPT}(\pi)$



## Performance guarantee: Example

- SimpleReversalSort is a minimization algorithm
- Performance guarantee:
  - $\max_{|\pi|=n} \Lambda(\pi) / \text{OPT}(\pi)$
  - We have seen that  $\Lambda(\pi) = n-1$  for a problem that could be solved in two steps
  - So, the approximation ratio is at least  $(n-1)/2$

## Adjacencies

$$\pi = \pi_1 \pi_2 \pi_3 \dots \pi_{n-1} \pi_n$$

- A pair of elements  $\pi_i$  and  $\pi_{i+1}$  are **adjacent** if

$$\pi_{i+1} = \pi_i \pm 1$$

- For example:

$$\pi = 1 \ 9 \ \underline{3} \ \underline{4} \ \underline{7} \ \underline{8} \ 2 \ \underline{6} \ 5$$

- (3, 4), (7, 8) and (6, 5) are adjacent pairs

## Breakpoints: An example

- There is a **breakpoint** between any neighboring elements that are non-consecutive (**not adjacent**):

$$\pi = 1 \mid 9 \mid 3 \ 4 \mid 7 \ 8 \mid 2 \mid 6 \ 5$$

- Pairs (1,9), (9,3), (4,7), (8,2) and (2,6) form breakpoints of permutation  $\pi$
- $b(\pi)$  - number of breakpoints in permutation  $\pi$

## Adjacencies and breakpoints

- An **adjacency** - a pair of neighboring elements that are **consecutive**
- A **breakpoint** - a pair of neighboring elements that are **not consecutive**

$$\pi = 5 \ 6 \ 2 \ 1 \ 3 \ 4 \longrightarrow \text{Extend } \pi \text{ with } \pi_0 = 0 \text{ and } \pi_7 = 7$$



## Reversal distance and breakpoints

Each reversal eliminates at most 2 breakpoints i.e. reversal distance  $\geq b(\pi) / 2$

$\pi = 2\ 3\ 1\ 4\ 6\ 5$

$0\  \underline{2\ 3}\  \underline{1}\  \underline{4}\  \underline{6\ 5}\  \underline{7}$	$b(\pi) = 5$
$0\ 1\  \underline{3\ 2}\  \underline{4}\  \underline{6\ 5}\  \underline{7}$	$b(\pi) = 4$
$0\ 1\ 2\ 3\ 4\  \underline{6\ 5}\  \underline{7}$	$b(\pi) = 2$
$0\ 1\ 2\ 3\ 4\ 5\ 6\ 7$	$b(\pi) = 0$

## Sorting by reversals: A better greedy algorithm

BreakPointReversalSort( $\pi$ )

```

1 while  $b(\pi) > 0$ 
2   Among all possible reversals, choose reversal
    $\rho$  minimizing  $b(\pi \cdot \rho)$ 
3    $\pi \leftarrow \pi \cdot \rho(i, j)$ 
4   output  $\pi$ 
5 return
    
```

**Problem:** It is not obvious that this algorithm will terminate!

## Strips

**Strip:** an interval between two consecutive breakpoints in a permutation

- **Decreasing strip:** strip of elements in decreasing order (e.g.  $6\ 5$  and  $3\ 2$ ).
- **Increasing strip:** strip of elements in increasing order (e.g.  $7\ 8$ )

$0\ \underline{1}\ \underline{2}\ \underline{4}\ \underline{3}\ \underline{7}\ \underline{8}\ \underline{2}\ \underline{5}\ \underline{6}\ \underline{10}$

- A single-element strip can be declared either increasing or decreasing. We will choose to declare them as decreasing with exception of the strips with  $0$  and  $n+1$

## Reducing the Number of Breakpoints

Theorem 1:

If permutation  $\pi$  contains at least one decreasing strip, then there exists a reversal  $\rho$  which decreases the number of breakpoints (i.e.  $b(\pi \cdot \rho) < b(\pi)$ )

## “Proof”

For  $\pi = 1\ 4\ 6\ 5\ 7\ 8\ 3\ 2\ 9$

$0\ 1\ |4\ 6\ 5|\ 7\ 8|\ 3\ 2|\ 9$        $b(\pi) = 5$

- Choose the decreasing strip with the smallest element  $k$  in  $\pi$  ( $k = 2$  in this case)
- Find  $k - 1$  in the permutation
- Reverse the segment between  $k$  and  $k-1$ :

$0\ 1\ \underline{4\ 6\ 5\ 7\ 8\ 3\ 2}\ 9$        $b(\pi) = 5$

↓

$0\ 1\ 2\ 3\ 8\ 7\ 5\ 6\ 4\ 9$        $b(\pi) = 4$

## Reducing the number of breakpoints again

- If there is no decreasing strip, there may be no reversal  $\rho$  that reduces the number of breakpoints
- By reversing an increasing strip (the number of breakpoints stay unchanged), we will create a decreasing strip at the next step. Then the number of breakpoints will be reduced in the next step (Theorem 1)

## ImprovedBreakpointReversalSort

ImprovedBreakpointReversalSort( $\pi$ )

```

1 while  $b(\pi) > 0$ 
2   if  $\pi$  has a decreasing strip
3     Among all possible reversals, choose reversal  $\rho$ 
       that minimizes  $b(\pi \cdot \rho)$ 
4   else
5     Choose a reversal  $\rho$  that flips an increasing strip in  $\pi$ 
6      $\pi \leftarrow \pi \cdot \rho$ 
7   output  $\pi$ 
8 return
```

## ImprovedBreakpointReversalSort: Performance guarantee

ImprovedBreakpointReversalSort is an approximation algorithm with a performance guarantee of at most 4

- It eliminates at least one breakpoint in every two steps; at most  $2b(\pi)$  steps
- Approximation ratio:  $2b(\pi) / d(\pi)$
- Optimal algorithm eliminates at most 2 breakpoints in every step:  $d(\pi) \geq b(\pi) / 2$
- Performance guarantee:

$$(2b(\pi) / d(\pi)) \leq [2b(\pi) / (b(\pi) / 2)] = 4$$

## Finding regulatory motifs in DNA sequences

## Implanting motif AAAAAAGGGGGG with four random mutations

```

atgccgggactactgtagAAgAAGGtGGGggcgtacacattagataaacgtatgaagtcgttagactcggcgcgcg
accctatTTTTgagcagatttagtgaacctggaaaaaatttagtacaaaactttccgaatacAAATAAACGGcGGa
tgatataccctgggtagcttAAAAAATGGcGTGGtctctcccgattttgaaatgtaggactatccagggtccga
gctgagaattggatgCAAAAAAGGGGtTtGccacgcaatcgcgaaccaacgpgaccacaaggaagcgaataaggaga
tccctttgggtaetgtgcccgggctgggttcgtgagggcccttaecggacttaATAAAGGaaGGctttag
gtcaatcgttcttgtgaatgattAAcAAAGGGcTGGgacccttggcaccacaattcagtgaggcagcaca
cggittggccctgttagagcccccgtATAAACAGGcGGcCaattetgagagcctaactctatcggctggttcaat
aacttagtAAAAAATAGGGcCccctggggcacatacaagaggctcttcttatacgttaagctgatacactatgta
ttggccattggctaaagcccaactgcaaatggagaatgatcttgcataCTAAAAAGGcGcGGaccgaagggaag
ctggtgagcaacgacagattcttactgcttagctcgtctccgggatctaatagcacgaagcttCAAAAAAGGcGGcGaa
    
```

## Where is the motif?

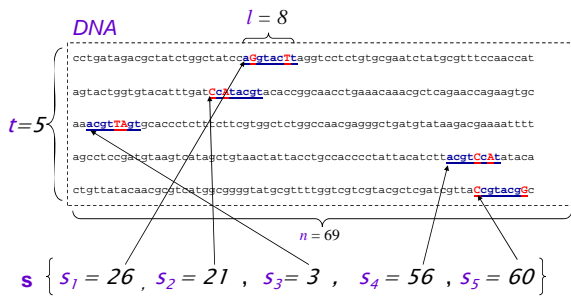
```

atgccgggactactgtagaagaaggTgggggcgtacacattagataaacgtatgaagtcgttagactcggcgcgcg
accctatTTTTgagcagatttagtgaacctggaaaaaatttagtacaaaactttccgaatacAAATAAACGGcGGa
tgatataccctgggtagcttAAaataaggatgggtgctctcccgattttgaaatgtaggactatccagggtccga
gctgagaattggatcaaaaaaggattgtccacgcaatcgcgaaccaacgpgaccacaaggaagcgaataaggaga
tccctttgggtaetgtgcccgggctgggttcgtgagggcccttaecggacttaataaataaaggaggctttag
gtcaatcgttcttgtgaatgatttaacaataaggcTgggacccttggcaccacaattcagtgaggcagcaca
cpgitttggccctgttagagcccccgtataacaaggaggccaattatgagagactaactctatcgtgctgttcaat
aacttgagttaaaaaatagggcccctggggcacatacaagaggcttcccttatacgttaagctgatacactatgta
ttggccattggctaaagcccaactgcaaatggagaatgatacttgcatactaaaggagcggaccgaagggaag
ctggtgagcaacgacagattcttactgcttagctcgtctccgggatctaatagcacgaagcttcaaaaaaggagcga
    
```

## Definitions

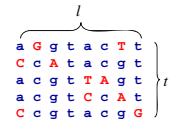
- $t$  number of sample DNA sequences
- $n$  length of each DNA sequence
- DNA** sample of DNA sequences ( $t \times n$  array)
- $l$  length of the motif ( $l$ -mer)
- $s_i$  starting position of an  $l$ -mer in sequence  $i$
- $S = (s_1, s_2, \dots, s_i)$  array of motif starting positions

## Parameters



## Scoring motifs: consensus score

➤ Given  $s = (s_1, \dots, s_l)$  and  $DNA$ :



➤  $\text{Score}(s, DNA) =$

$$\sum_{i=1}^l \max_{k \in \{A, T, C, G\}} \text{count}(k, i)$$

A	3	0	1	0	3	1	1	0
C	2	4	0	0	1	4	0	0
G	0	1	4	0	0	0	3	1
T	0	0	0	5	1	0	1	4

Consensus a c g t a c g t

Score 3+4+4+5+3+4+3+4=30

## Greedy motif finding

➤ Partial score:  $\text{Score}(s, i, DNA)$

- The consensus score for the first  $i$  sequences

➤ Algorithm:

- Find the optimal motif for the two first sequences
- Scan the remaining sequences only once, and choose the motif with the best contribution to the partial score

## Greedy motif finding

GreedyMotifSearch( $DNA, l, n, l$ )

```

1   $s \leftarrow (l, l, \dots, l)$ 
2   $bestMotif \leftarrow s$ 
3  for  $s_1 \leftarrow 1$  to  $n - l + 1$ 
4      for  $s_2 \leftarrow 1$  to  $n - l + 1$ 
5          if  $\text{Score}(s, 2, DNA) > \text{Score}(bestMotif, 2, DNA)$ 
6               $bestMotif_1 \leftarrow s_1$ 
7               $bestMotif_2 \leftarrow s_2$ 
8   $s_1 \leftarrow bestMotif_1$ 
9   $s_2 \leftarrow bestMotif_2$ 
10 for  $i \leftarrow 3$  to  $l$ 
11     for  $s_i \leftarrow 1$  to  $n - l + 1$ 
12         if  $\text{Score}(s, i, DNA) > \text{Score}(bestMotif, i, DNA)$ 
13              $bestMotif_i \leftarrow s_i$ 
14          $s_i \leftarrow bestMotif_i$ 
15 return  $bestMotif$ 
    
```

## Running time

- Optimal motif for the two first sequences
  - $(n - l + 1)^2$  operations
- The remaining  $t-2$  sequence
  - $(t-2)(n - l + 1)$  operations
- Running time
  - $O(ln^2 + tn)$  or  $O(ln^2)$  if  $n \gg t$
- Vastly better than
  - BruteForceMotifSearch:  $(n - l + 1)^t$
  - BruteForceMedianStringSearch:  $4^l$