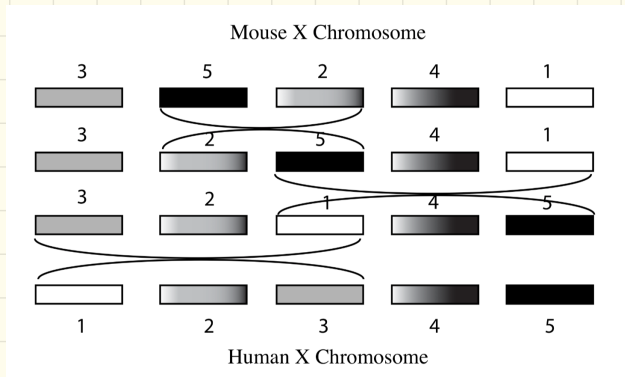# Algorithms in Bioinformatics

## Greedy Algorithms

# Recap

- Website should be updated
- HW - due next Tuesday
- Essays - hopefully Thursday?
- On Sep 27 - no class

## Today: Ch 5

# Problem Motivation: gene flipping



Mouse X Chromosome / Human X Chromosome

## Simple version:

Rearrangement events can be modeled by reversals.

Order of synteny blocks (genes) can be represented by a permutation: a reordering of $1 \ldots n$.

$$\pi = \pi_1 \pi_2 \ldots \pi_n$$

Reversal: $\pi = \pi_1 \ldots \pi_{i-1} \overline{\pi_i \ldots \pi_j} \pi_{j+1} \ldots \pi_n$

$$= \pi_1 \pi_2 \ldots \pi_i \pi_j \pi_{j-1} \ldots \pi_{i+1} \pi_i \ldots \pi_2$$

Ex: $\pi = 1\ 2\ 4\ 3\ 7\ 5\ 6 = 1\ 2, 5\ 7\ 3\ 4, 6$

# The problem:

---
**Reversal Distance Problem:**

*Given two permutations, find a shortest series of reversals that transforms one permutation into another.*

    **Input:** Permutations $\pi$ and $\sigma$.

    **Output:** A series of reversals $\rho_1, \rho_2, \ldots, \rho_t$ transforming $\pi$ into $\sigma$ (i.e., $\pi \cdot \rho_1 \cdot \rho_2 \cdots \rho_t = \sigma$), such that $t$ is minimum.

---

Then $t = d(\pi, \sigma)$ is the reversal distance between $\pi$ & $\sigma$.

In practice, one is the main one, & we reset $\sigma$ to be $1..n$. (& update $\pi$ accordingly)

    ↑ replace $\sigma_i$ (in $\pi$) with $i$

Then:

---
**Sorting by Reversals Problem:**

*Given a permutation, find a shortest series of reversals that transforms it into the identity permutation.*

    **Input:** Permutation $\pi$.

    **Output:** A series of reversals $\rho_1, \rho_2, \ldots, \rho_t$ transforming $\pi$ into the identity permutation such that $t$ is minimum.

---

Then $t = d(\pi)$ is the reversal distance of $\pi$.

prefix$(\pi) = 3$

Ex: $\pi = 123\,645$

$\uparrow P(4,5)$

What would you do?

$\pi \circ P(4,5)$

1 2 3 4 6 5 1

1 2 3 4 5 6

$d(\pi) \le 2$

Can it be $= 1$?

(lower bnd)
exhaustively, no

Define prefix($\pi$) = # of already sorted elements in $\pi$.

Try to increase this by 1 iteratively (and greedily):

SIMPLEREVERSALSORT($\pi$)
1   **for** $i \leftarrow 1$ **to** $n-1$
2       $j \leftarrow$ position 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**

1 2 3 $\underline{6\ 4}$ 5

$i = 4$
$j = 5$
$\rho(i,j) =$
swap from $i$ to $j$

Does this always work?

use induction:
After round $i$ of loop,
$i$ is in correct spot.

How many flips?
worst case: $n-1$
(not optimal)

Optimal?

$\overbrace{6\ 2}\ 1, 345$

Ex:  $\overset{\downarrow\ \downarrow}{\overbrace{6\ 1}}\ 2\ 3\ 4\ 5$

greedy!

↳ $1\ \overset{\downarrow}{6}\ \overset{\downarrow}{2}\ 3\ 4\ 5$

↳ $1\ 2\ \overbrace{6\ 3}\ 4\ 5$

5 swaps

↳ $1\ 2\ 3\ \overbrace{6\ 4}\ 5$

↳ $1\ 2\ 3\ 4\ \overbrace{6\ 5}$

↳ $1\ 2\ 3\ 4\ 5\ 6$

best?

no  2

$\underbrace{\overbrace{6\ 1\ 2\ 3\ 4\ 5}}_{}$

$\underbrace{5\ 4\ 3\ 2\ 1,\ 6}$

Worse:  $\pi = n\ 1\ 2\ \dots\ (n\text{-}2)\ (n\text{-}1)$

# Similar CS problem:
## Pancake flipping

Stack of pancakes (out of order)

input:



goal:



All you have is a spatula.
(Basically, permutations will all be
prefixes.)

Ex: <u>1 2 3 6 4 5</u>

⤷ 6 3 2 1 4 5

(The minimum # of needed
flips is unknown.)

OK, so greedy doesn't always work!
Let's talk approximation:

Performance guarantee:

$A(\pi) =$ solution produced
by algorithm A on
input $\pi$

$OPT(\pi) =$ best solution on input $\pi$

Approximation ratio (for minimization
problem) =

$$\max_{|\pi|=n} \frac{A(\pi)}{OPT(\pi)}$$

(Worst case scenerio)

Ex: We showed approx ratio
of SIMPLE REVERSE SORT
is $\geq \dfrac{n-1}{2}$

If $n=1001$, alg could return
$500 \cdot OPT$ permutations.

# Better: a different notion of greed

## Breakpoints:

To make it easier, extend $\Pi$:

$$\Pi = \Pi_0 \ \Pi_1 \Pi_2 \cdots \Pi_n \ \Pi_{n+1}$$

For any $i \leq n$,

- if $\Pi_i \Pi_{i+1}$ are consecutive, call it an adjacency.

- if not, a <u>break point</u> (& let $b(\Pi) = \#$ breakpoints)

<u>Ex</u> $\Pi = 0, 2 \ 1, 3 \ 4 \ 5, 8 \ 7 \ 6, 9$
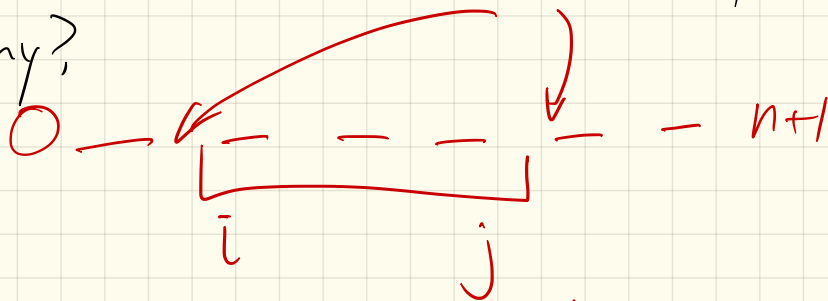
$$b(\Pi) = 4$$

$$0 \ n \ n\text{-}1 \cdots 1 \ n+1$$

<u>Note</u>: only 1 permutation has 0 break points.

<u>Observation</u>: Every reversal can
eliminate at most 2 breakpoints

Why?

$$0 - - - \overbrace{\underbrace{\phantom{- - - - -}}_{}}^{} - - - n+1$$

$i$        $j$

So: $d(\pi) \geq \dfrac{b(\pi)}{2}$ — invent some parameter to measure hardness for lower bnds

<u>Goal</u>: try to choose flips that lower $b(\pi)$

But — can removing a breakpoint increase new ones later? Or could we get stuck?

<u>Dfn</u>: Strips: an interval between 2 breakpoints

<u>Ex</u>: 0, 2 1, 3 4 5, 8 7 6, 9

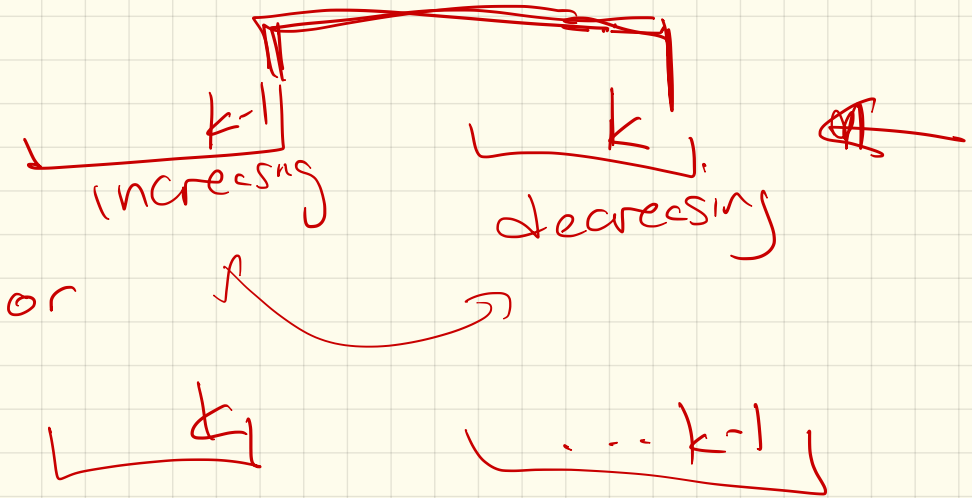Either increasing or decreasing
(or size 1, in which case both)

<u>Thm</u>: If a permutation $\pi$ contains a decreasing strip or a strip of size 1, then there is a flip that decreases the # of breakpoints.

<u>Pf</u>: Choose decreasing (or single) strip with smallest k.

012 765 8 43 9

Note k-1 can't be in a decreasing strip,
& if in fact ends some increasing strip

So

$k-1$     increasing

$k$     decreasing

$\cancel{k}$

or

$\cancel{k}$ 1

$\ldots k-1$

reverse between
$k$ & $k-1$

So — breakpts
decrease by 1
(maybe more)
$k-1$ & $k$ are now adjacent

So: if have a decreasing strip,
    can    reduce.

If not : all strips are increasing.
What to do ?

flip some
   increasing strip

# Algorithm:

IMPROVEDBREAKPOINTREVERSALSORT($\pi$)
1  **while** $b(\pi) > 0$
2      **if** $\pi$ has a decreasing strip
3          Among all reversals, choose reversal $\rho$ minimizing $b(\pi \cdot \rho)$
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**

← this reduces $b(\pi)$ by $\geq 1$

↑ for every 2 iterations, do the if iterations

Thm: This algorithm is a 4-approximation.
  (ie it's # of swaps is $\leq 4 \cdot OPT$)

pf: # steps in alg
      $\leq 2 \, b(\pi)$

$\dfrac{A(\pi)}{OPT(\pi)}$

since $b(\pi)$ goes down by (at least) 1 for every 2 iterations

and observation 3 slides ago:

$$OPT = d(\pi) \geq \frac{b(\pi)}{2}$$

↳ rearrange $b(\pi) \leq 2d(\pi)$

# steps in alg $\leq 2b(\pi) \leq 2\,(2d(\pi))$
$\leq 4\, d(\pi)$

# Another example: Motifs again

Motif finding: $O(\ell n^t)$ or $O(4^\ell nt)$

Greedy approach:
- Scan each DNA sequence <u>once</u>
- Decide which $\ell$-mers have best contribution as we go

More detail:

- Scan sequence $1 + 2$ & find 2 closest $\ell$-mers (in Hamming distance)

  Time: $\ell(n-\ell+1)$

- Then as $j = 3$ to $t$ scan $j$-th sequence & select the best $\ell$-mer (max Score $(s,i)$)

  Time: $\ell(n-\ell+))$ each time

# Pseudocode:

```
GREEDYMOTIFSEARCH(DNA, t, n, l)
 1  bestMotif ← (1, 1, ..., 1)
 2  s ← (1, 1, ..., 1)
 3  for s₁ ← 1 to n − l + 1
 4      for s₂ ← 1 to n − l + 1
 5          if Score(s, 2, DNA) > Score(bestMotif, 2, DNA)
 6              BestMotif₁ ← s₁
 7              BestMotif₂ ← s₂
 8      s₁ ← BestMotif₁
 9      s₂ ← BestMotif₂
10  for i ← 3 to t
11      for sᵢ ← 1 to n − l + 1
12          if Score(s, i, DNA) > Score(bestMotif, i, DNA)
13              bestMotifᵢ ← sᵢ
14      sᵢ ← bestMotifᵢ
15  return bestMotif
```

Time:

This tool is called CONSENSUS.

No approximation ratio!
So optimal can be missed
entirely, & can be bad.

<u>Note</u>: — usually do more than
2 rows

— usually run several
times in random order

## Next time:

Greedy approach for Shortest common superstring.

mention in 8.4

(my notes are based on another source — I'll post a link)