Alignment of Long Sequences

BMI/CS 776 www.biostat.wisc.edu/bmi776/ Spring 2011 Mark Craven craven@biostat.wisc.edu

Goals for Lecture

the key concepts to understand are the following

- how large-scale alignment differs from the simple case
- the canonical three step approach of large-scale aligners
- using suffix trees to find MUMs (alignment seeds)
- using tries and threaded tries to find alignment seeds
- constrained dynamic programming to align between/ around anchors
- using sparse DP to find a chain of local alignments

Pairwise Large-Scale Alignment: Task Definition

Given

- a pair of large-scale sequences (e.g. chromosomes)
- a method for scoring the similarity of a pair of characters

Do

 – construct global alignment: identify matches between sequences as well as various non-match features

Large Scale Alignment Example: Mouse Chr6 vs. Human Chr12



250000





General Strategy



 perform pattern matching to find seeds for global alignment



2. find a good chain of anchors

Figure from: Brudno et al. Genome Research, 2003



 fill in remainder with standard but constrained alignment method

Comparison of Large-Scale Alignment Methods

| Method | Pattern matching | Chaining |
|--------|--------------------------------------|---------------------------|
| MUMmer | suffix tree - MUMs | LIS variant |
| AVID | suffix tree - exact & wobble matches | Smith-Waterman variant |
| LAGAN | <i>k</i> -mer trie, inexact matches | sparse DP |

The MUMmer System

Delcher et al., Nucleic Acids Research, 1999

Given: genomes A and B

- 1. find all maximal, unique, matching subsequences (MUMs)
- 2. extract the longest possible set of matches that occur in the same order in both genomes
- 3. close the gaps



Suffix Trees

- substring problem:
 - given text *S* of length *m*
 - preprocess *S* in *O*(*m*) time
 - such that, given query string Q of length n, find occurrence (if any) of Q in S in O(n) time
- suffix trees solve this problem, and others

Suffix Tree Definition

- a suffix tree *T* for a string *S* of length *m* is tree with the following properties:
 - rooted and directed
 - -m leaves, labeled 1 to m

key property



- each edge labeled by a substring of S
 - concatenation of edge labels on path from root
 to leaf *i* is suffix *i* of *S* (we will denote this by *S_{i...m}*)
- each internal non-root node has at least two children
- edges out of a node must begin with different characters

Suffixes

S = "banana\$" suffixes of S \$ a\$ na\$ ana\$ nana\$ anana\$

banana\$

Suffix Tree Example S ="banana\$" add '\$' to end so that suffix tree exists (no suffix is a prefix of another suffix) a n a n а а 6

Solving the Substring Problem

1

5

- assume we have suffix tree T
- FindMatch(Q, T): •
 - follow (unique) path down from root of Taccording to characters in Q
 - if all of Q is found to be a prefix of such a path return label of some leaf below this path
 - else, return no match found _



MUMs and Generalized Suffix Trees

- build one suffix tree for both genomes A and B
- label each leaf node with genome it represents









- sort MUMs according to position in genome A
- solve variation of *Longest Increasing Subsequence* (LIS) problem to find sequences in ascending order in both genomes



Figure from: Delcher et al., Nucleic Acids Research 27, 1999

Finding Longest Subsequence

- unlike ordinary LIS problems, MUMmer takes into account
 - lengths of sequences represented by MUMs
 - overlaps
- requires O(k log k) time where k is number of MUMs

Types of Gaps in a MUMmer Alignment

1. SNP: exactly one base (indicated by ^) differs between the two sequences. It is surrounded by exact-match sequence.

Genome A: cgtcatgggcgttcgtcgttg Genome B: cgtcatgggcattcgtcgttg

2. Insertion: a sequence that occurs in one genome but not the other.

| Genome A: | cggggtaaccgccctggtcggg |
|-----------|--|
| Genome B: | cggggtaaccgcgttgctcggggtaaccgccctggtcggg |
| | ~~~~~~~~~~~~~~~~ |

3. Highly polymorphic region: many mutations in a short region.

Genome A: ccgcctcgcctgg.gctggcgcccgctc Genome B: ccgcctcgccagttgaccgcgcccgctc

4. Repeat sequence: the repeat is shown in uppercase. Note that the first copy of the repeat in Genome B is imperfect, containing one mismatch to the other three identical copies.

Genome A: cTGGGTGGGACAACGTaaaaaaaaaTGGGTGGGACAACGTc Genome B: aTGGGTGGGGCgACGTgggggggggGGGGTGGGACAACGTa

Figure from: Delcher et al., Nucleic Acids Research 27, 1999

Step 3: Close the Gaps

- SNPs:
 - between MUMs: trivial to detect
 - otherwise: handle like repeats
- inserts
 - transpositions (subsequences that were deleted from one location and inserted elsewhere): look for out-of-sequence MUMs
 - simple insertions: trivial to detect



- polymorphic regions
 - short ones: align them with dynamic programming method
 - long ones: call MUMmer recursively w/ reduced min MUM length
- repeats
 - detected by overlapping MUMs



Figure from: Delcher et al. Nucleic Acids Research 27, 1999

The LAGAN Method

Brudno et al., Genome Research, 2003

```
Given: genomes A and B

anchors = find_anchors(A, B)

step 3: finish global alignment with DP constrained by anchors

find_anchors(A, B)

step 1: find local alignments by matching, chaining k-mer seeds

step 2: anchors = highest-weight sequence of local alignments

for each pair of adjacent anchors a_1, a_2 in anchors

if a_1, a_2 are more than d bases apart

A', B' = sequences between a_1, a_2

sub-anchors = find_anchors(A', B')

insert sub-anchors between a_1, a_2 in anchors

return anchors
```

Step 1a: Finding Seeds in LAGAN

- *degenerate k-mers*: matching *k*-long sequences with a small number mismatches allowed
- by default, LAGAN uses 10-mers and allows 1 mismatch

cacg cgcg<mark>c</mark>tacat acct acta cgcggtacat cgta







Step 2: Chaining in LAGAN

• use *sparse dynamic programming* to chain local alignments



The Problem: Find a Chain of Local Alignments





(x,y) → (x',y') requires x < x' y < y'

Each local alignment has a weight

FIND the chain with highest total weight

Sparse DP for rectangle chaining

1,..., N: rectangles
(h_j, l_j): y-coordinates of rectangle j
w(j): weight of rectangle j
V(j): optimal score of chain ending in j
L: list of triplets (l_j, V(j), j)
L is sorted by l_j: smallest (North) to largest (South) value
L is implemented as a balanced binary tree

Slide from Serafim Batzoglou, Stanford University

Sparse DP for rectangle chaining



Main idea:

- Sweep through xcoordinates
- To the right of b, anything chainable to a is chainable to b
- Therefore, if V(b) > V(a), rectangle a is "useless" for subsequent chaining
- In L, keep rectangles j sorted with increasing l_jcoordinates ⇒ sorted with increasing V(j) score

Slide from Serafim Batzoglou, Stanford University

Sparse DP for rectangle chaining

Go through rectangle x-coordinates, from lowest to highest:

- 1. When on the leftmost end of rectangle i:
 - a. j: rectangle in L, with largest $l_i < h_i$
 - b. V(i) = w(i) + V(j)
- 2. When on the rightmost end of i:
 - a. k: rectangle in L, with largest $I_k \leq I_i$
 - b. If V(i) > V(k):
 - i. INSERT (I_i, V(i), i) in L
 - ii. **REMOVE** all $(I_j, V(j), j)$ with $V(j) \le V(i) \& I_j \ge I_j$

Slide from Serafim Batzoglou, Stanford University



V(b)

V(a)



Time Analysis

- 1. Sorting the x-coords takes O(N log N)
- 2. Going through x-coords: N steps
- 3. Each of N steps requires O(log N) time:
 - Searching L takes log N
 - Inserting to L takes log N
 - All deletions are consecutive, so log N per deletion
 - Each element is deleted at most once: N log N for all deletions
 - Recall that INSERT, DELETE, SUCCESSOR, take O(log N) time in a balanced binary search tree

Slide from Serafim Batzoglou, Stanford University



Step 3: Computing the Global Alignment in LAGAN

- given an anchor that starts at (i, j) and ends at (i', j'), LAGAN limits the DP to the unshaded regions
- thus anchors are somewhat flexible



Figure from: Brudno et al. Genome Research, 2003

