C G T A C G T A

A C G T A C G T

Whole Genome Alignment

Adam M. Phillippy

CMSC701: April 18, 2019





The Forefront of Genomics[®]

Goal of whole genome alignment

• For genomes A and B, find a mapping from each position in A to its corresponding position in B



CCGGTAGGCTATTAAACGGGGGTGAGGAGCGTTGGCATAGCA



Not so fast...

• Genome A may have insertions, deletions, translocations, inversions, or duplications with respect to B





Global vs. local alignments

- Global pairwise alignment
 - . . AAGCTTGGCTTAGCTGCTAGGGTAGGCTTGGG . . .
 - ... AAGCTGGGCTTAGTTGCTAG.. TAGGCTTTGG...
 - A
 AA
 A

- Whole genome alignment
 - A "collection" of local alignments





Global alignment visualization



NIH

Whole genome alignment visualization

- How can we visualize *whole* genome alignments?
- With an alignment dot plot
 - N x M matrix
 - Let *i* = position in genome *A*
 - Let j = position in genome B
 - Fill cell (*i*,*j*) if A_i shows similarity to B_j



• A perfect alignment between A and B would completely fill the positive diagonal







Drosophila genome evolution



NIH NHGRI

Homology map of human chr7



- Homology
 - Shared ancestry

- Showing:
 - >85% identity
 - Blue: 1–5 kbp
 - Red: 5–10 kbp
 - Black: >10 kbp



Homologs: orthologs vs. paralogs





Paralogous genes





Goal of whole genome alignment

- For genomes A and B, find a mapping from each position in A to its <u>orthologous</u> position in B
- Requires an evolutionary model
 - Compared to nucleotide substitution models genome rearrangement models are not well defined
 - Typically based on "synteny" arguments
 - Genomic regions tend to be inherited in blocks



Dotplot quiz





http://mummer.sf.net/manual/AlignmentTypes.pdf

NIH

NHGRI

MUMmer

- Maximal Unique Matcher (MUM)
 - match
 - exact match of a minimum length
 - maximal
 - cannot be extended in either direction without a mismatch
 - unique
 - occurs only once in both sequences (MUM)
 - occurs only once in the reference sequence (MAM)
 - occurs one or more times in either sequence (MEM)



MEMs and MUMs



NHGRI

How do we turn MUMs into alignments?

- Nucmer approach
 - Find all exact matches
 - Chain exact matches
 - Gapped alignment
- With which algorithms?





Seed, chain, and extend

FIND all MUMs CHAIN consistent MUMs EXTEND alignments





Seed

NIH

Suffix indexes for MUM finding

	Suffix tree	Suffix array	FM Index
Time: Does P occur?	O(n)	O(n log m)	O(n)
Time: Count k occurrences of P	O(n + k)	O(n log m)	O(n)
Time: Report k locations of P	O(n + k)	O(n log m + k)	O(n + k)
Space	O(m)	O(m)	O(m)
Needs T?	yes	yes	no
Bytes per input character	>15	~4	~0.5

m = |T|, n = |P|, k = # occurrences of P in T



Chain

A



Global: Longest Increasing Subsequence

- A subsequence of a permutation is a collection of elements of the permutation in the order that they appear. For example, (5, 3, 4) is a subsequence of (5, 1, 3, 4, 2).
- A subsequence is **increasing** if the elements of the subsequence increase. For example, given the permutation (8, 2, 1, 6, 5, 7, 4, 3, 9), an increasing subsequence is (2, 6, 7, 9).
 - Given: A a permutation π of length n.
 - Return: A longest increasing subsequence of π .

From rosalind.info

NHGRI

Applied LIS example





```
LIS pseudo code
for i in [1...n]
   for j in [1..i)
      score = max[j] + len(i) - olap(j,i)
      if score > max[i]
         max[i] = score
         pre[i] = j
```

NIH Can also define arbitrary gap penalties, and speed up with early termination heuristics: see minimap2

Local: Agglomerative Clustering

cluster length = $\sum m_i$

gap distance = A

diagonal difference = |B - A| or |B - A| / B





Agglomerative pseudo code

for i in [1...n)

- for j in [i+1..n]
 - if gap(i,j) < g and diagdiff(i,j) < d
 union(find(i),find(j))</pre>



Extend

A



Banded alignment extension



Alignment extension



break point = B

break length = A



Filter



Which are orthologs?



NIH Mark a segment as "best" if it has the highest score of all segments intersecting the sweep line

NHGRI

Best for *A*





Best for *B*





Best for A and B (orthologs only)





Best for A or B (paralogs too)





SuperMap



Multiple whole-genome alignments without a reference organism. Dubchak *et al.* 2009