

Sequence Mapping/Alignment

BCB 504: Applied Bioinformatics

Matt Settles

University of Idaho
Bioinformatics and Computational Biology Program

April 4, 2012

- 1 Introduction
- 2 Alignment Algorithms
 - BLAST
 - BLAT
 - Improvements
- 3 Illumina Data
 - Hash Based
 - Burrows-Wheeler Transform
- 4 SAM/BAM output

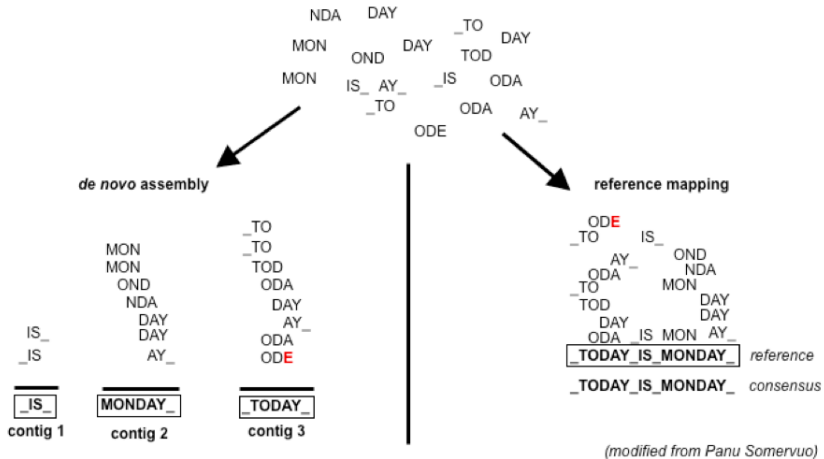
Mapping/Alignment

Given sequence data,

Assembly seeks to put together the puzzle without knowing what the picture is

Mapping tries to put together the puzzle pieces directly onto an image of the picture

In mapping the question is more, given a small chunk of sequence, where in the genome did this piece most likely come from.



Basic Local Alignment Search Tool (BLAST)

Some say the first bioinformatics tool, developed at NIH and published in 1990.

Problem:

- Exact algorithms like Smith-Waterman and Needleman-Wunsch (dynamic programming) are slow, when the search space becomes large.
- With the advent of automated DNA sequencing technology, the database of possible matches was becoming increasingly larger.

the BLAST algorithm emphasizes speed over sensitivity, and does not guarantee an optimal alignment.

BLAST is a few to many algorithm - performs gapped alignment

BLAST Like Alignment Tool (BLAT)

Blat (Jim Kent, UCSC, 2002) was designed to solve the problem of performing comparisons between large genomes and was one of the first algorithms to efficiently search many query sequences against a large database (a genome). Blat also performs a gapped-alignment for searching RNA sequences against a genome and handling splice junctions.

gapped-alignment alignment allowing for insertions and deletions greater than a few base pairs. Gapped alignment are less efficient, but more accurate.

BLAT is a many to many algorithm - performs gapped alignments

Improving Algorithms

Many additional algorithms have been developed since BLAST and BLAT, mainly improving on either speed or accuracy, or both.

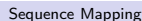
and then came Illumina data -
New Problem:

- We still have a large search space (aka genome)
- Very small pieces, many possible close matches
- Millions or tens of millions of query sequences

Types

Hash based First generation of alignment algorithms relied on hashes (Eland [Illumina], RMAP, MAQ, SHRiMP, SOAP)

Burrows-Wheeler Second generation algorithms with a reduced memory footprint (BWA, SOAP2, Bowtie)



hash based example: MAQ

- Index reference genome (or sequence reads) => creates hash index
 - Big file: >50GB
 - takes a long time (hours or overnight), but only need to do once
- Divide each read into segments (seeds) and look up in table
 - Search stage finds regions in the genome that can potentially be homologous to the read.
 - Alignment stage verifies these regions to check if they are indeed homologous. More computationally intensive

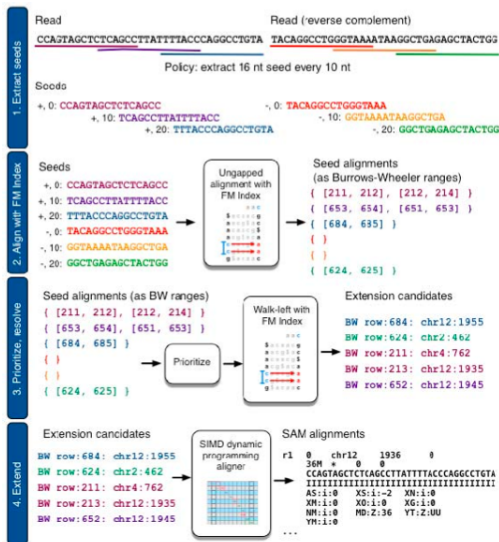
burrows-wheeler example: Bowtie2

Used in data compression (e.g. bzip) => index: much smaller than hash-based index (<2GB)

- Alignment speed: 30x faster than MAQ

Steps:

- Create BWT index of genome
- Align read 1 character at a time to BWT-transformed genome



considerations

- placing reads in regions that do not exist in the reference genome
- sequencing errors and variations: alignment between read and true source in genome may have more differences than alignment with some other copy of repeat.
- What if many nucleotide differences with closest fully sequenced genome? (3% is a common alignment capability)
- placing reads in repetitive regions: Some algorithms only return 1 mapping; If multiple: map quality = 0
- algorithms that use paired-end information => might prefer correct distance over correct alignment.

SAM/BAM format

SAM (Sequence Alignment/Map) format = unified format
for storing read alignments to a reference genome

BAM = binary version of SAM for fast querying

```
7172283 163 chr9 139389330 60 90M = 139389482 242 TAGGAGG... EHHHHHH...
7705896 83 chr9 139389513 60 90M = 139389512 -91 GCTGGGG... EBCHHFC...
7705896 163 chr9 139389512 60 90M = 139389513 91 AGCTGGG... HHHHHHH...
```

1	QNAME	query template name
2	FLAG	bitwise flag
3	RNAME	reference sequence name
4	POS	1-based leftmost mapping position
5	MAPQ	mapping quality
6	CIGAR	CIGAR string
7	RNEXT	reference name of mate
8	PNEXT	position of mate
9	TLEN	observed template length
10	SEQ	sequence
11	QUAL	ASCII of Phred-scaled base quality