

# NGS Data and Sequence Alignment

Manpreet S. Katari  
Aug 11, 2016

- ### Outline
- NGS Data
    - FastA
    - FastQ
    - SAM
    - BAM
    - GFF
  - Sequence Alignment
    - Global vs Local
    - Dynamic Programming
    - Burrow Wheeler's Algorithm.

### Important files types

FASTA		Sequence files
FASTQ		
SAM		Alignment files
BAM		
GFF		Annotation files

### Important file types: FASTA

A sequence in FASTA format begins with a single-line description, followed by lines of sequence data. The description line is distinguished from the sequence data by a greater-than (">") symbol in the first column. The word following the ">" symbol is the identifier of the sequence, and the rest of the line is the description (both are optional).

There should be no space between the ">" and the first letter of the identifier. It is recommended that all lines of text be shorter than 80 characters. The sequence ends if another line starting with a ">" appears; this indicates the start of another sequence.

### Important file types: FASTA

```
>chr1
CCACACACACCCACACACCCACACACCCACACACACACACACCCACACACCC
CACACACACACACCCACACACCCACACACCCACACACCCACACACCCACACCC
GCCACACCCGCTCCACACACCCACACACCCACACACCCACACACCCACACCC
CCGCGCCCAATCAACACACCCACACACCCACACACCCACACACCCACACCC
ACTACACCCACACCCACACCCACACACCCACACACCCACACACCCACACCC
CGACATACCCACACACCCACACACCCACACACCCACACACCCACACACCC
TGCTCTACCCACACATATTGAAAGCCACACAAATGATCGTAATAACA
CACACCGCTTACCGTACCCATTTATACACACACACACACATGCTACTCAC
CCCTACTTGTATACCTGATTTAGCTACGACACGGATGCTACGATATATA
CACTCCCAAACTACCTACCTCCAGATCCACCTTCCACTCCAGAGGCCAT
CCTCCTCAATCAGTACCAAAAGCCACTCACATCATATGACAGGCACTT
GCCCTACGGGCTTATACCCGCTTACCCATTTACCCATACCCGCTTATAT
CCACATTTGTATCTTATACCTCACTTGGGGGCTCCCAATATGATATAC
TGCCCTTAACATACCTTATACACCTTTTCCACATATACCTTACCTC
CAFTTATACACCTTATGCAATATACAGAAAATCCCAAAAATCA
CCTAAGCTAABAAATTTCTCTCTTCCACGATATACATACATATATG
GCTTGGGTAGCACACTATCATGCTTACTACTAACCTAATGTTCTTCA
TATTGCAATTTGCTTGAACGATGCTTATTCAGAAATTTTCTACTACA
CAGGCGTACATTTAGATATATGTCATCATCTCGTGGTAAACCTTTTA
GGST
```





SAM format

```

@SQ   SN:YBR27M   LN:3580
@SQ   SN:YBR27L   LN:3533
@SQ   SN:YBR27R   LN:3537
@SQ   SN:YBR28   LN:3861
@SQ   SN:YBR29   LN:4953
@SQ   SN:YBR31   LN:4687
@SQ   SN:YBR32   LN:3397
@SQ   SN:YBR33C   LN:3390
@SQ   SN:YBR33L   LN:4390
@SQ   SN:YBR34   LN:4643
@SQ   SN:YBR35   LN:3588
@SQ   SN:YBR36   LN:3288
@SQ   SN:YBR37   LN:3685
@SQ   SN:YBR38   LN:3685
@SQ   SN:YBR39   LN:3685
@SQ   SN:YBR40   LN:3685
@SQ   SN:YBR41   LN:3685
@SQ   SN:YBR42   LN:3685
@SQ   SN:YBR43   LN:3685
@SQ   SN:YBR44   LN:3685
@SQ   SN:YBR45   LN:3685
@SQ   SN:YBR46   LN:3685
@SQ   SN:YBR47   LN:3685
@SQ   SN:YBR48   LN:3685
@SQ   SN:YBR49   LN:3685
@SQ   SN:YBR50   LN:3685
@SQ   SN:YBR51   LN:3685
@SQ   SN:YBR52   LN:3685
@SQ   SN:YBR53   LN:3685
@SQ   SN:YBR54   LN:3685
@SQ   SN:YBR55   LN:3685
@SQ   SN:YBR56   LN:3685
@SQ   SN:YBR57   LN:3685
@SQ   SN:YBR58   LN:3685
@SQ   SN:YBR59   LN:3685
@SQ   SN:YBR60   LN:3685
@SQ   SN:YBR61   LN:3685
@SQ   SN:YBR62   LN:3685
@SQ   SN:YBR63   LN:3685
@SQ   SN:YBR64   LN:3685
@SQ   SN:YBR65   LN:3685
@SQ   SN:YBR66   LN:3685
@SQ   SN:YBR67   LN:3685
@SQ   SN:YBR68   LN:3685
@SQ   SN:YBR69   LN:3685
@SQ   SN:YBR70   LN:3685
@SQ   SN:YBR71   LN:3685
@SQ   SN:YBR72   LN:3685
@SQ   SN:YBR73   LN:3685
@SQ   SN:YBR74   LN:3685
@SQ   SN:YBR75   LN:3685
@SQ   SN:YBR76   LN:3685
@SQ   SN:YBR77   LN:3685
@SQ   SN:YBR78   LN:3685
@SQ   SN:YBR79   LN:3685
@SQ   SN:YBR80   LN:3685
@SQ   SN:YBR81   LN:3685
@SQ   SN:YBR82   LN:3685
@SQ   SN:YBR83   LN:3685
@SQ   SN:YBR84   LN:3685
@SQ   SN:YBR85   LN:3685
@SQ   SN:YBR86   LN:3685
@SQ   SN:YBR87   LN:3685
@SQ   SN:YBR88   LN:3685
@SQ   SN:YBR89   LN:3685
@SQ   SN:YBR90   LN:3685
@SQ   SN:YBR91   LN:3685
@SQ   SN:YBR92   LN:3685
@SQ   SN:YBR93   LN:3685
@SQ   SN:YBR94   LN:3685
@SQ   SN:YBR95   LN:3685
@SQ   SN:YBR96   LN:3685
@SQ   SN:YBR97   LN:3685
@SQ   SN:YBR98   LN:3685
@SQ   SN:YBR99   LN:3685
@SQ   SN:YBR100   LN:3685

```

SAM format gets converted to a BAM file



Annotation Formats

- Mostly tab delimited files that describe the location of genome features (i.e., genes, etc.)
- Also used for displaying annotations on standard genome browsers
- Important for associating alignments with specific genome features
- Descriptions

GFF3 Format

GFF3 format is a flat tab-delimited file. The first line of the file is a comment that identifies the file format and version. This is followed by a series of data lines, each one of which corresponds to an annotation. Here is a miniature GFF3 file:

```

##gff-version 3
ctg123 . exon 1300 1500 . + . ID=exon0001
ctg123 . exon 1050 1500 . + . ID=exon0002
ctg123 . exon 3000 3902 . + . ID=exon0003
ctg123 . exon 5000 5500 . + . ID=exon0004
ctg123 . exon 7000 9000 . + . ID=exon0005

```

The `##gff-version 3` line is required and *must* be the first line of the file. It introduces the annotation section of the file.

Columns:  
Seqid, Source, Type, Start, End, Score, Strand, Phase, Attribute(Identifier)

GFF format

```

##gff-version 3
##sequence-region 1 1 1000000
seqid source type start end score strand phase attribute
chr1 chr1 gene 1000 2000 0.0 + . ID=gene1
chr1 chr1 exon 1000 1200 0.0 + . ID=exon1
chr1 chr1 exon 1200 1500 0.0 + . ID=exon2
chr1 chr1 exon 1500 1800 0.0 + . ID=exon3
chr1 chr1 exon 1800 2000 0.0 + . ID=exon4
chr1 chr1 intron 1200 1500 0.0 - . ID=intron1
chr1 chr1 intron 1500 1800 0.0 - . ID=intron2
chr1 chr1 intron 1800 2000 0.0 - . ID=intron3
chr1 chr1 UTR 5' 1000 1000 0.0 + . ID=UTR5
chr1 chr1 UTR 3' 2000 2000 0.0 + . ID=UTR3
chr1 chr1 CDS 1000 1000 0.0 + . ID=CDS1
chr1 chr1 CDS 1000 1200 0.0 + . ID=CDS2
chr1 chr1 CDS 1200 1500 0.0 + . ID=CDS3
chr1 chr1 CDS 1500 1800 0.0 + . ID=CDS4
chr1 chr1 CDS 1800 2000 0.0 + . ID=CDS5
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS6
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS7
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS8
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS9
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS10
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS11
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS12
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS13
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS14
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS15
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS16
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS17
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS18
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS19
chr1 chr1 CDS 2000 2000 0.0 + . ID=CDS20

```

Global and local approaches to aligning sequences

**GLOBAL:** Attempt to "match" and assess similarity between two entire sequences

**LOCAL:** Find subsequences of high similarity

... and then possibly "stitch" (chain, net, thread) together local alignments to obtain an overall comparison of the original sequences.

**The second approach is more meaningful**  
(especially for long sequences, of different lengths, like whole genomes)

Two protein or DNA sequences are unlikely to present a straightforward overall "match", even if they are closely related.

**Why?** Substitutions are not the only process by which they diverge: insertions, deletions and rearrangements

### Dynamic programming: Pairwise sequence alignment

Create a matrix table for comparing sequences with one sequence along each axis (size  $m+1, n+1$ )

Fill in partial alignment scores until score for entire sequence has been calculated:

- Assign score for each position  $ij$ , progressing from top left to bottom right
- Score rule: take maximum of the three choices
  1. Take value from left, assign gap penalty (gap along left axis)
  2. Take value from top, assign gap penalty (gap along top axis)
  3. Take value from diagonal above left, assign match/mismatch score
- Repeat until table is completed.

Use trace-back to obtain full alignment:

```

AC--TCG
ACAGTAG
    
```

### BLAST: heuristic database search using local alignment

#### Basic Local Alignment Search Tool

... T L S R D Q H A W R L S ... query sequence

**QW (query word)**, size  $W=3$ .

**A cartoon for BLAST (parameters)**

For each QW, use the scoring matrix to form a neighborhood NB = (all words of size  $W$  with a score  $\geq T=13$ )

Find match(es) to words belonging to NB in the subject (target) sequence

For each match, use the scoring matrix and gap penalties to produce an HSP (High Scoring Segment Pair), then extend alignment on both sides, until

- drop is  $> X$ , or
- score goes below  $S_{min}$  (minimum hit score)

### About Blat (from genome.ucsc.edu)

- "BLAT on DNA is designed to quickly find sequences of 95% and greater similarity of length 25 bases or more."
- "It may miss more divergent or shorter sequence alignments. It will find perfect sequence matches of 20 bases."
- "BLAT is not BLAST."
- "DNA BLAT works by keeping an index of the entire genome in memory. The index consists of all overlapping 11-mers stepping by 5 except for those heavily involved in repeats."
- "The index takes up about 2 gigabytes of RAM. The genome itself is not kept in memory, allowing BLAT to deliver high performance on a reasonably priced Linux box."
- "The index is used to find areas of probable homology, which are then loaded into memory for a detailed alignment."

### Hash Table (BLAT)

NATURE METHODS SUPPLEMENT | VOL.6 NO.134 | NOVEMBER 2009

### Short Read Applications

```

...CCATAG      TATGCGCCC      CCGAAATT      GGTATAC...
..CCAT      CTATATGG      TCGGAAAT      CGGTATAC
..CCAT      GGCCTATG      CTATGCGAA      GCGGTATAC
..CCA      AGGCTATAT      CCTATCGGA      TTGCGGTA      C...
..CCA      AGGCTATAT      GGCCTATG      AAATTTCG      ATAC...
..CC      AGGCTATAT      GGCCTATG      AAATTTCG      ATAC...
..CC      TGGCTATAT      GGCCTATG      AAATTTCG      ATAC...
...CCATAGGCTATATGCGCCTATCGGCAATTTCGCGTATAC...
...
...CCATAGGCTATATGCGCCTATCGGCAATTTCGCGTATAC...
    
```

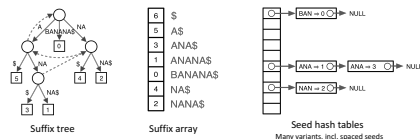
Finding the alignments is typically the performance bottleneck

### New alignment algorithms must address the requirements and characteristics of NGS reads

- Millions of reads per run (30x of genome coverage)
- Short Reads (as short as 36bp)
- Different types of reads (single-end, paired-end, mate-pair, etc.)
- Base-calling quality factors
- Sequencing errors (~ 1%)
- Repetitive regions
- Sequencing organism vs. reference genome
- Must adjust to evolving sequencing technologies and data formats

### Indexing

- Genomes and reads are too large for direct approaches like dynamic programming
- **Indexing is required**



- Choice of index is key to performance

### Suffix Array

Find "ctat" in the reference

[00] ggglaaagctataactatgacaggcgt	[04] aaagctataactatgacaggcgt
[01] gglaaagctataactatgacaggcgt	[12] aactatgacaggcgt
[02] gaaagctataactatgacaggcgt	[05] aagctataactatgacaggcgt
[03] laaagctataactatgacaggcgt	[13] actatgacaggcgt
[04] aaagctataactatgacaggcgt	[06] agctataactatgacaggcgt
[05] aagctataactatgacaggcgt	[23] agcgt
[06] agctataactatgacaggcgt	[10] alaaactatgacaggcgt
[07] gctataactatgacaggcgt	[20] atcaggcgt
[08] ctataactatgacaggcgt	[16] atgacaggcgt
[09] tataactatgacaggcgt	[22] caggcgt
[10] ataactatgacaggcgt	[26] cgt
[11] laactatgacaggcgt	[08] ctataactatgacaggcgt
[12] aactatgacaggcgt	[14] ctatgacaggcgt
[13] actatgacaggcgt	[19] gatcaggcgt
[14] ctatgacaggcgt	[25] gcgt
[15] tatgacaggcgt	[07] gctataactatgacaggcgt
[16] atgacaggcgt	[24] ggcgt
[17] tgaacaggcgt	[00] ggglaaagctataactatgacaggcgt
[18] tgaacaggcgt	[01] gglaaagctataactatgacaggcgt
[19] gatcaggcgt	[02] gaaagctataactatgacaggcgt
[20] atcaggcgt	[27] gtt
[21] tccaggcgt	[29] t
[22] caggcgt	[03] laaagctataactatgacaggcgt
[23] agcgt	[11] laactatgacaggcgt
[24] ggcgt	[09] tataactatgacaggcgt
[25] gcgt	[15] tatgacaggcgt
[26] cgt	[21] tccaggcgt
[27] gtt	[18] tgaacaggcgt
[28] tt	[28] tt
[29] t	[17] tgaacaggcgt

### NGS Read Alignment Burrows Wheeler Transformation (BWT)

- Invented by David Wheeler in 1983 (Bell Labs). Published in 1994. "A Block-Sorting Lossless Data Compression Algorithm" Systems Research Center Technical Report No 124. Palo Alto, CA: Digital Equipment Corporation, Burrows M., Wheeler DJ. 1994
- Originally developed for compressing large files (bzp2, etc.)
- Lossless, Fully Reversible
- Alignment Tools based on BWT: *bowtie*, *BWA*, *SOAP2*, etc.
- Approach:
  - Align reads on the *transformed* reference genome, using an efficient index (FM index)
  - Solve the simple problem first (align one character) and then build on that solution to solve a slightly harder problem (two characters) etc.
- Results in great speed and efficiency gains (a few GigaByte of RAM for the entire H. Genome). Other approaches require tens of GigaBytes of memory and are much slower.

### Burrows Wheeler Transformation

Text = c t g a a a c t g g t \$

1 2 3 4 5 6 7 8 9 10 11 0

➤ Introduce \$ at the end and construct all cyclic permutations of Text

```

c t g a a a c t g g t $ $
t g a a a c t g g t $ c
g a a a c t g g t $ c t
a a a c t g g t $ c t g
a a c t g g t $ c t g a
a c t g g t $ c t g a a
c t g g t $ c t g a a a
t g g t $ c t g a a a c
g g t $ c t g a a a c t
g t $ c t g a a a c t g
t $ c t g a a a c t g g
t g a a a c t g g t $ c
t g g t $ c t g a a a
    
```

### Burrows Wheeler Transformation

- Sort rows alphabetically, keeping of which row went where

```

$ c t g a a a c t g g t t
a a a c t g g t $ c t g
a a a c t g g t $ c t g a
a c t g g t $ c t g a a
c t g a a a c t g g t $
c t g a c t $ c t g a a a
g a a a c t g g t $ c t
g g t $ c t g a a a c t
g t $ c t g a a a c t g
t $ c t g a a a c t g g
t g a a a c t g g t $ c
t g g t $ c t g a a a
    
```

Burrows Wheeler Matrix

BWT (Text) = t g a a \$ a t t g g c c

11 3 4 5 0 6 2 8 9 10 1 7

### Exact Matching with FM Index

• In progressive rounds, **top** & **bot** delimit the range of rows beginning with progressively longer suffixes of Q

