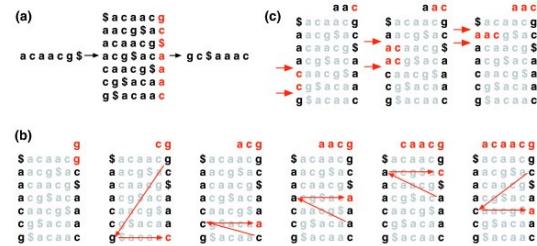


SRA: Burrows-Wheeler Transform

- BWT is an efficient data indexing technique that maintains a relatively small memory footprint when searching through a given data block. BWT was extended by Ferragina and Manzini to a newer data structure, named FM-index, to support exact matching. By transforming the genome into an FM-index, the lookup performance of the algorithm improves for the cases where a single read matches multiple locations in the genome. However, the improved performance comes with a significantly large index build up time compared to hash tables.



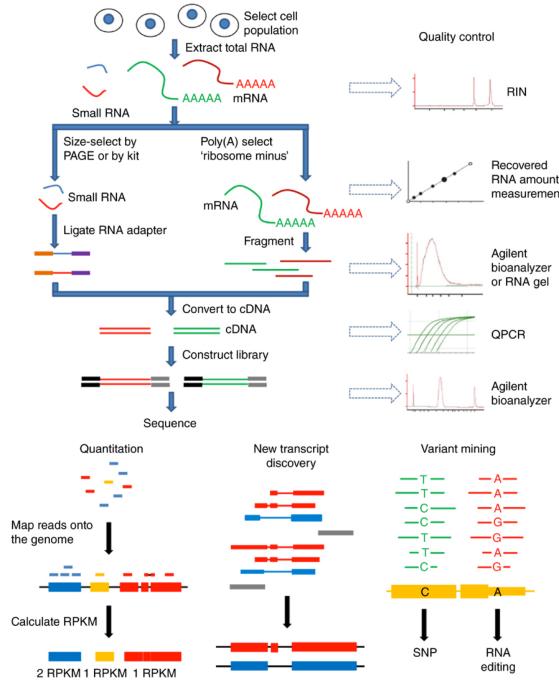
SRA: BWA

- BWA is a BWT based tool. The BWA tool uses the Ferragina and Manzini matching algorithm to find exact matches. To find inexact matches, the authors provided a new backtracking algorithm that searches for matches between substring of the reference genome and the query within a certain defined distance.
- BWA is fast, and can do gapped alignments. When run without seeding, it will find all hits within a given edit distance. Long read aligner is also fast, and can perform well for 454, Ion Torrent, Sanger, and PacBio reads. BWA is actively maintained and has a strong user community.

SRA: Bowtie

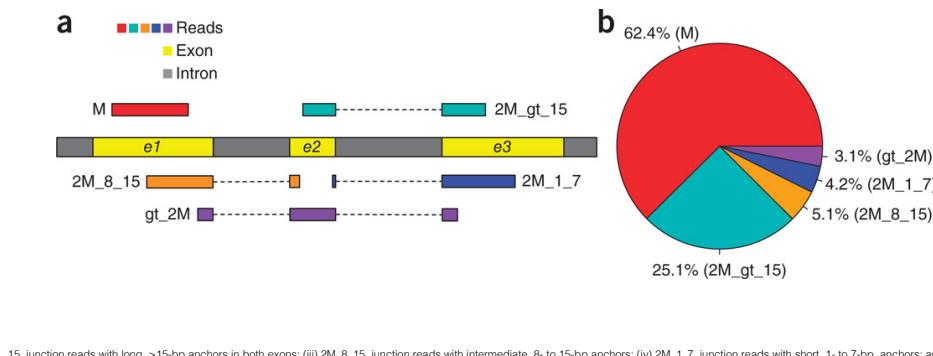
- Bowtie starts by building an FM-index for the reference genome and then uses the modified Ferragina and Manzini [39] matching algorithm to find the mapping location. There are two main versions of Bowtie namely Bowtie and Bowtie 2. Bowtie 2 is mainly designed to handle reads longer than 50 bps. Additionally, Bowtie 2 supports features not handled by Bowtie.
- Bowtie2 is faster than BWA for some types of alignment, but it takes a hit in sensitivity and specificity in some applications.

- Introduction to Mapping
- Short Read Aligners
- DNA vs RNA
- Alignment Quality
- Pitfalls and Improvements



RNASeq

HISAT2

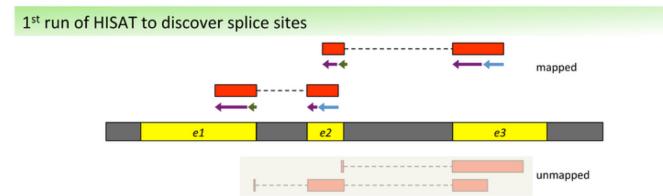


15 junction reads with long, >15-bp anchors in both exons; (iii) 2M, 8, 15 junction reads with intermediate, 8- to 15-bp anchors; (iv) 2M, 1, 7 junction reads with short, 1- to 7-bp anchors; and

Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory requirements. Nat Methods. 2015 Apr;12(4):357-60. doi: 10.1038/nmeth.3317. Epub 2015 Mar 9. PubMed PMID: 25751142; PubMed Central PMCID: PMC4655817.

SRA Features: Splice-Aware

- Splicing refers to the process of cutting the RNA to remove the non-coding part (introns) and keeping only the coding part (exons) and joining them together. Therefore, when sequencing the RNA, a read might be located across exon-exon junctions. The process of mapping such reads back to the genome is hard due to the variability of the intron length. For instance, the intron length ranges between 250 and 65,130 nt in eukaryotic model organisms [37].



HISAT2

Sensitivity and precision of leading spliced aligners

	no. of splice sites reported	no. of true splice sites reported	sensitivity (%)	Precision (%)
Program	reported	sites reported	(%)	(%)
HISATx1	91,904	85,546	97.3	93.1
HISATx2	90,331	85,603	97.3	94.8
HISAT	90,300	85,587	97.3	94.8
STAR	95,892	84,678	96.3	88.3
STARx2	92,254	84,734	96.3	91.8
GSNAP	92,547	85,598	97.3	92.5
oLego	86,779	82,879	94.2	95.5
TopHat2	96,474	79,705	90.6	82.6

Sensitivity and precision of leading spliced aligners for 87,944 true splice sites contained in 20 million simulated reads from the human genome, with a mismatch rate of 0.5%. Sensitivity is the percentage of true splice sites found out of the total that were present. Precision (or positive predictive value) is the percentage of reported splice sites that are correct.

Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory requirements. *Nat Methods*. 2015 Apr;12(4):357-60. doi: 10.1038/nmeth.3317. Epub 2015 Mar 9. PubMed PMID: 25751142; PubMed Central PMCID: PMC4655817.

- Introduction to Mapping
- Short Read Aligners
- DNA vs RNA
- Alignment Quality
- Pitfalls and Improvements

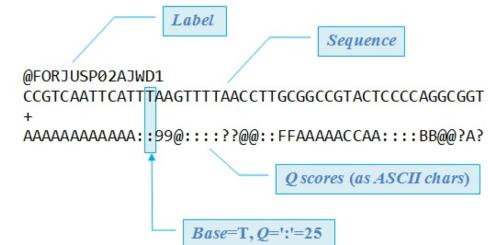
Alignment Metrics

- Alignment Rate (Mapping Rate)
- Paired Alignments
 - Properly Paired Mapping
 - Comparing the rate of read pairs mapped within a certain proximity
 - Average Insert Size
 - Distance between adapter sequences
- Duplication Rate
 - Same fragment size can be duplicated during library preparation (PCR) or during sequencing (colony formation)

Mapping Quality

- Probability that a read is mapped incorrectly
- Factors include:
 - uniqueness (one best scoring alignment)
 - number of mismatches
 - number of gaps
 - quality of the bases (Phred)

File Formats: FastQ



File Formats: SAM

Field	Regular expression	Range	Description
QNAME	[^ \t\n\r]+		Query pair NAME if paired; or Query NAME if unpaired ²
FLAG	[0-9]+	[0,2 ¹⁶ -1]	bitwise FLAG (Section 2.2.2)
RNAME	[^ \t\n\r@=]+		Reference sequence NAME ³
POS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost POSition/coordinate of the clipped sequence
MAPQ	[0-9]+	[0,2 ⁸ -1]	MAppling Quality (phred-scaled posterior probability that the mapping position of this read is incorrect) ⁴
CIGAR	(([0-9]+[MIDNSHP]) *)		extended CIGAR string
MRNM	[^ \t\n\r@=]+		Mate Reference sequence NaMe; "=" if the same as <RNAME> ³
MPOS	[0-9]+	[0,2 ²⁹ -1]	1-based leftmost Mate POSition of the clipped sequence
ISIZE	-?[0-9]+	[-2 ²⁹ ,2 ²⁹]	inferred Insert SIZE ⁵
SEQ	[acgtnACGTN.=] \ *		query SEQuence; "=" for a match to the reference; n/N. for ambiguity; cases are not maintained ^{6,7}
QUAL	[!~-]+ \ *	[0,93]	query QUALity; ASCII-33 gives the Phred base quality ^{6,7}
TAG	[A-Z][A-Z0-9]		TAG
VTYPE	[AifZH]		Value TYPE
VALUE	[^ \t\n\r]+		match <VTYPE> (space allowed)

File Formats: SAM Flags

Flag	Chr	Description
0x0001	p	the read is paired in sequencing
0x0002	P	the read is mapped in a proper pair
0x0004	u	the query sequence itself is unmapped
0x0008	U	the mate is unmapped
0x0010	r	strand of the query (1 for reverse)
0x0020	R	strand of the mate
0x0040	1	the read is the first read in a pair
0x0080	2	the read is the second read in a pair
0x0100	s	the alignment is not primary
0x0200	f	the read fails platform/vendor quality checks
0x0400	d	the read is either a PCR or an optical duplicate

Bitwise:

```

00000000001 => 0x0001 => 2^0 = 1    => PAIRED
00000000010 => 0x0002 => 2^1 = 2    => PAIR MAPPED
00000000100 => 0x0004 => 2^2 = 4    => READ UNMAPPED
00000001000 => 0x0008 => 2^3 = 8    => MATE UNMAPPED
00000010000 => 0x001 => 2^4 = 16   => READ REVERSE
00000100000 => 0x0002 => 2^5 = 32   => MATE REVERSE
00001000000 => 0x0004 => 2^6 = 64   => FIRST IN PAIR
00010000000 => 0x0008 => 2^7 = 128  => SECOND IN PAIR
00100000000 => 0x0001 => 2^8 = 256  => ALIGN NOT PRIM.
01000000000 => 0x0002 => 2^9 = 512  => QUALITY FAILS
10000000000 => 0x0004 => 2^10 = 1024 => PCR DUPLICATE

```

File Formats: SAM

```

@HD VN:1.4
@SQ SN:insert LN:599
@SQ SN:ref1 LN:45
@SQ SN:ref2 LN:48
@SQ SN:ref3 LN:4
@RG ID:fish PG:donkey
@RG ID:fish1 PG:donkey
@RG ID:fish2 PG:donkey
@RG ID:mosse PG:mosse
@RG ID:cow PG:cow
r000 99 insert 59 39 19M = 88 39 ATTTTACCTAC AAAAAAAA RG:Z:cow PG:Z:bull
r000 211 insert 89 39 19M = 59 -38 CCATCAATT AAAAAAAA RG:Z:cow PG:Z:bull
r001 163 ref1 7 39 5M14M23M * 37 39 TTAGATAAACGCGATCTG * * XX:B:S,12561,2,20,112 YY:t:108 XB:t:-18 RG:Z:fish PG:Z:cool
r002 8 ref1 9 38 13216M11P11P114H21 * 0 0 AGCTAA * RG:Z:cow XA:Z:cabc
r003 8 ref1 9 38 5M6M * 0 0 AGCTAA * RG:Z:cool PG:Z:cool
r004 8 ref1 16 38 6M14M11M * 0 0 ATAGCTCTACGC * RG:Z:cool PG:Z:cool
r005 16 ref1 29 38 6M5M * 0 0 TAGGC * RG:Z:cool PG:Z:cool
r006 83 ref1 37 38 9M = 7 -39 CAGCCCAT * * RG:Z:fish PG:Z:cool
r007 8 ref2 1 38 20M * 0 0 AGTTTAAACAGCAAA * RG:Z:cool PG:Z:bull
r008 6 38 21M * 0 0 GGTTTTAACAGCAAAATT ?????????????????? RG:Z:cool PG:Z:bull
r009 6 38 9M113M * 0 0 TTATTAACAAATTAAGTCACAG * ?????????????????? RG:Z:fish PG:Z:cool
r010 8 ref2 18 38 25M * 0 0 CAATAATTAGCTCACAGGAGACM ?????????????????? RG:Z:fish PG:Z:bull
r011 8 ref2 12 38 24M * 0 0 AATAATTAGCTCACAGGAGCACT ?????????????????? RG:Z:fish PG:Z:bull
r012 8 ref2 14 38 23M * 0 0 TAATTAGCTCACAGGACACTA ?????????????????? RG:Z:cow PG:Z:bull
u1 4 * 0 38 23M * 0 0 TAATTAGCTCACAGGAAAAAA ???????????????????

```

@SQ = Contigs/Chromosomes

@RG = Read Group

@PG = Program Info

- Introduction to Mapping
- Short Read Aligners

- DNA vs RNA
- Alignment Quality
- Pitfalls and Improvements

Pitfalls

- All of the down stream analysis is based on the alignment
- Improper alignment can lead to false positive variate calls
- Inaccurate abundance calculations

Misalignment

- The Human Genome has many duplications
- Misalignment can result from sequence repetition
- Misalignment can be improved with:
 - Increased read lengths
 - Finding multiple 35bp matches is more likely than finding multiple 100bp matches
 - Paired Sequencing
 - A mate pair can “anchor” another when there is multiple mapping of the other mate pair.