#### Aligning reads to a genome Analysis of Next-Generation Sequencing Data

#### Luce Skrabanek

Applied Bioinformatics Core

Slides at https://bit.ly/2CUdS9z<sup>1</sup>

5 February, 2019

Weill Cornell Medicine

<sup>1</sup>http://physiology.med.cornell.edu/faculty/skrabanek/lab/angsd/schedule\_2018/

Luce Skrabanek (ABC, WCM)

Aligning reads to a genome



- 2 What do we align to?
- 3 How do we align?

#### 4 Output files



### Why do we align?

#### What do we learn?





**CNV** detection

SNP identification and frequency estimation







#### which genes are expressed, and how much

Luce Skrabanek (ABC, WCM)

Aligning reads to a genome

#### What do we align to?

#### What do we need?

- Reference sequence: the nucleotide sequence of the chromosomes of a species <sup>2</sup>
- Optional annotations: the gene/transcript models for a genome; includes the coordinates of the exons of a transcript on a reference genome, optionally the strand, gene name, coding portion of the transcript.

<sup>2</sup>see discussion on reference genomes in [Ballouz et al., 2019]

## Sources for reference genomes

#### Ensembl

- http://www.ensembl.org
- UCSC
  - https://genome.ucsc.edu/
- NCBI
  - https://www.ncbi.nlm.nih.gov/
- Gencode
  - https://www.gencodegenes.org/
- Organism-specific databases
  - (e.g., http://toxodb.org/toxo/)

Always note the source and version of your reference genome. Look out for chromosome naming conventions.

#### Annotations



RefSeq ncbi.nlm.nih.gov/refseq UCSC Known Genes genome.ucsc.edu Ensembl/Gencode gencodegenes.org

> 1/3 protein-coding genes > 17,000 non-coding RNAs > 15,000 pseudogenes

The chromosome names must match those in your reference genome; annotations must correspond to the same reference genome build as your reference genome fasta file.

#### Gene models can vary dramatically



## Which annotation should you use?

"More sensitive annotations, such as **Ensembl** (...) **should be preferred** over more specific annotations, such as RefSeq (...) if the aim is to obtain accurate expression estimates."

> Janes et al. (Briefings in Bioinformatics, 2015). doi: 10.1093/bib/bbv007

> > "We observe that **RefSeq Genes produces the most accurate fold-change measures** with respect to a ground truth of RT-qPCR gene expression estimates. "

> > > Wu et al. (BMC Bioinfo, 2013). doi: 10.1186/1471-2105-14-S11-S8

"In practice, there is **no simple answer to this question**, and it depends on the purpose of the analysis. (...) When choosing an annotation database, researchers should keep in mind that **no database is perfect** and **some gene annotations might be inaccurate or entirely wrong**."

Zhao & Zhang (BMC Genomics, 2015). doi:10.1186/s12864-015-1308-8

F 1 ··· .		1 71		 00101	
Luce Skrabanek	(ABC, WCM)	Aligning r	reads to a genome	5 February, 2019	10 / 29

#### Storing annotation information



- Represent genome coordinates and gene descriptions/names
- multiple formats: GFF2, GFF3, GTF<sup>3</sup>, BED, SAF...

 $^{3} http://genome.ucsc.edu/FAQ/FAQformat\#format4$ 

### How do we align?

#### Aligners

- Genomic aligners
  - BWA [Li and Durbin, 2009], Bowtie2
- Splice-aware aligners
  - STAR [Dobin et al., 2013], TopHat, HiSAT2
- Pseudo alignment
  - Salmon, kallisto, RSEM

#### Challenge

Mapping millions of reads accurately and in a reasonable amount of time, despite complications from sequencing errors, genomic variation and repetitive elements.

## Genomic aligner: BWA

BWA uses a canonical seed-and-extend paradigm. BWA is based on the Burrows-Wheeler Transform and uses the FM-index<sup>4</sup> to search for exact string matches.



This has a very small memory footprint.

<sup>4</sup>Full-text Minute-space, or Ferragina and Manzini [Ferragina and Manzini, 2010]

#### FM-index backwards search



#### BWA-MEM

BWA MEM [Li, 2013] is the next generation in the BWA family, and is one of the few that works well for both 70bp reads and long sequences up to a few megabases.

- allows long gaps
- 2 the allowable error rate adjusts with sequence length
- 3 can report multiple non-overlapping local hits
  - As for BWA, uses a canonical seed-and-extend paradigm, grouping seeds that are colinear and close to each other as a chain.
  - Each seed is extended using a banded affine-gap-penalty dynamic programming, stopping when the difference between the best and the current extension score is above some threshold, avoiding extension through poorly aligned regions
  - Keep track of the best extension score reaching the end of the query sequence. If the difference between the best score reaching the end and the best local alignment score is below a threshold, the local alignment will be rejected even if it has a higher score.

### Mapping to the transcriptome



- 1 Alignment of exon-exon spanning reads
- 2 Multiple isoforms
- 3 Identification of novel splice junctions

Is it better to map to the genome, or the transcriptome?

STAR generates the SA from both genomic sequence, as well as the sequence spanning known exon-exon boundaries (transcriptome) to generate the SA.

STAR can also identify novel junctions, if it finds enough reads as support. Users can define how many reads must span a novel junction, and how many bases must be covered on either side of the junction.

# Splice-aware aligner: STAR [Spliced Transcripts Alignment to a Reference]



Luce Skrabanek (ABC, WCM)

Aligning reads to a genome

5 February, 2019 18 / 29

## Running STAR



STAR has many parameters (familiarize yourself with the manual)! See [Ballouz et al., 2018] for a discussion of how parameter selection affects mapping (e.g., handling of multi-mapped reads, intron sizes).

Output files

#### Output files

#### SAM files



Each line of the optional header section starts with Q. and includes information such as chromosomes names (SN) and their lengths (LN). The vast majority of lines within a SAM file are compact representations of the read alignments where each read is described by the 11 mandatory entries and a variable number of optional fields [Li et al., 2009].

### SAM FLAG field

	- 2	3	- 4	0	0		8	9	10		>11	
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT	

2<sup>nd</sup> field: binary FLAG

Binary (Decimal)	Hex	Description
00000000001 (1)	0x1	Is the read paired?
0000000010 (2)	0x2	Are both reads in a pair mapped "properly" (i.e., in the correct orientation with respect to one another)?
00000000100 (4)	0x4	Is the read itself unmapped?
00000001000 (8)	0x8	Is the mate read unmapped?
00000010000 (16)	0x10	Has the read been mapped to the reverse strand?
00000100000 (32)	0x20	Has the mate read been mapped to the reverse strand?
00001000000 (64)	0x40	Is the read the first read in a pair?
00010000000 (128)	0x80	Is the read the second read in a pair?
00100000000 (256)	0x100	Is the alignment not primary? (A read with split matches may have multiple primary alignment records.)
01000000000 (512)	0x200	Does the read fail platform/vendor quality checks?
10000000000 (1024)	0x400	Is the read a PCR or optical duplicate?

The FLAG field includes information about the mapping of the individual read. It is a bitwise flag, compactly storing answers to multiple binary Yes/No questions as a short series of bits where each of the single bits can be addressed separately.

See https://broadinstitute.github.io/picard/explain-flags.html to interpret bit flag values.

# CIGAR [Concise Idiosyncratic Gapped Alignment Report string]

1	2	3	4	5	6	7	8	9	10	11	>11
QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT

6<sup>th</sup> field: CIGAR string – which hoops did the aligner have to jump through to align the read to the <u>genome</u> locus that it thought was the best fit?

- M alignment (match or mismatch!!)
- I (N)
   insertion (large insertions)
   spliced out introns = sequences are missing in the read, i.e., they need to be inserted in order to align the read to the genome

   D
   clipping
   align the read to the genome

									F	Refe	ere	nc	e s	eq	ue	nc	e v	vitł	n a	ılig	ne	d r	ea	ds							CIGAR string	Explanation
(	2	Т	G	С	А	Т	G	Т	Т	А	G	А	Т	А	А	*	*	G	А	Т	А	G	С	Т	G	Т	G	С	Т	Α		
															А	Α	G	G	А	т	А	*	С	т	G						1M <mark>2I</mark> 4M1D3M	<b>Insertion &amp; Deletion</b>
2											G	А	т	А	А	*	G	G	А	т	А					_					5M1P1I4M	Padding & Insertion
ğ						т	G	т	т	А																т	G	С	т	А	5M15N5M	Spliced read
-	а	а	а	С	А	т	G	т	т	А	G																				388M	Soft clipping
4	¥.	A	A	С	A	т	G	т	т	А	G																				3H8M	Hard clipping

#### Output files

## SAM OPT field

	1						7			175			
	QNAME	FLAG	RNAME	POS	MAPQ	CIGAR	RNEXT	PNEXT	TLEN	SEQ	QUAL	OPT	
afte	r 11 <sup>th</sup> fiel	d: OPT	ΓIONAL	infor	matior	ı							
AS:i BC:Z HI:i NH:i NM:i MD:Z RG:Z	Alignme Barcodo Query i Number Edit dis String t Read g	ent scor e seque: s <i>i</i> -th h r of rep stance of that cor coup (sl	re nce nit stored orted alig of the que ntains the hould ma	in the gnmen ery to e exac tch th	e file its for t the ref t position ie entry	he quer erence ons of m after II	y sequen nismatch D if @RG	ce es (shou is preser	<b><tag< b=""> tags a ld com</tag<></b>	>: <t are no pleme e head</t 	YPE> ot star nt the der.	: <b><val< b=""> ndardiz CIGAR s</val<></b>	UE> ed!
4 T W A	.2.2 SAM the SAM at thich accept S nM jM ji	<b>A attri</b> tributes t a list <b>I XS</b> . Bj	butes. can be sj of 2-chara y default,	pecifie acter S STAF	d by the SAM att t output	e user us ributes. ts NH HI	ingou The imp AS nM a	tSAMatt: plemente uttribute	ributes d attrib	s <b>A1 A</b> butes a	2 A3 are: NH	opt I HI NM	ion MD
	NH HI	NM MD	have stan	dard n	neaning	as defin	ed in the	SAM fo	rmat sp	ecifica	tions.		
	AS id t	he local	alignmer	t scor	e (paire	d for pai	red-end	reads).					
	nM is the num	he num mber of	ber of mis mismatcl	match nes in	ies per each ma	(paired) ate.	alignmer	nt, not to	be cor	nfused	with N	M, which	ı is
	jM:B:c,M1,M2, intron motifs for all junctions (i.e. N in CIGAR): 0: non-canonical; 1: GT/AG, 2: CT/AC, 3: GC/AG, 4: CT/GC, 5: AT/AC, 6: GT/AT. If splice junctions database is used, and a junction is annotated, 20 is added to its motif value.												1: ase
	jI:B:I	,Start	1,End1,S	tart2	,End2,	Star	t and Er	nd of intr	ons for	all ju	nctions	(1-base	d).
	jM jI some d	attribut ownstre	es require am tools	e samt such a	tools 0.1 as Cuffli	.18 or la nks.	ater, and	were rej	ported t	to be i	ncomp	atible w	ith

The number of optional SAM/BAM fields, their value types and the information stored within them depends on the alignment program and can vary substantially. Luce Skrabanek (ABC, WCM) Aligning reads to a genome 5 February, 2019 24 / 29

## Exploring SAM/BAM files

The most widely used tool to explore and manipulate SAM/BAM files is samtools.

There are many options to subset reads based on SAM fields such as chromosomal location, or FLAG value, or mapping quality.

samtools view <in.bam>

Use egrep to subset reads based on the optional tags.

Most downstream applications also require the BAM file to be indexed by reference sequence position, to allow the efficient retrieval of all reads aligning to a locus.

samtools index <in.bam>

References

## References

- Sara Ballouz, Alexander Dobin, Thomas R Gingeras, and Jesse Gillis. The fractured landscape of RNA-seq alignment: the default in our STARs. *Nucleic Acids Research*, 46(10):5125–5138, 05 2018. ISSN 0305-1048. doi: 10.1093/nar/gky325. URL https://dx.doi.org/10.1093/nar/gky325.
- Sara Ballouz, Alexander Dobin, and Jesse Gillis. Is it time to change the reference genome? *bioRxiv*, 2019. doi: 10.1101/533166. URL https://www.biorxiv.org/content/early/2019/01/29/533166.
- Alexander Dobin, Carrie A. Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha, Philippe Batut, Mark Chaisson, and Thomas R. Gingeras. STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics*, 29 (1):15–21, 2013. doi: 10.1093/bioinformatics/bts635.
- Paolo Ferragina and Giovanni Manzini. Opportunistic Data Structures with Applications. Technical report, 2010.

Jürgen Jänes, Fengyuan Hu, Alexandra Lewin, and Ernest Turro. A comparative study of RNA-seq analysis strategies. *Briefings in Bioinformatics*, (January):1–9, 2015. ISSN 1467-5463. doi: 10.1093/bib/bbv007.

- Heng Li. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv e-prints*, art. arXiv:1303.3997, March 2013.
- Heng Li and Richard Durbin. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*, 25(14):1754–1760, 05 2009. ISSN 1367-4803. doi: 10.1093/bioinformatics/btp324. URL https://dx.doi.org/10.1093/bioinformatics/btp324.
- Heng Li, Bob Handsaker, Alec Wysoker, Tim Fennell, Jue Ruan, Nils Homer, Gabor Marth, Goncalo Abecasis, and Richard Durbin. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16):2078–9, August 2009. ISSN 1367-4811. doi: 10.1093/bioinformatics/btp352.
- S.P. Pfeifer. From next-generation resequencing reads to a high-quality variant data set. *Heredity*, 118(2):111–124, 2017. doi: 10.1038/hdy.2016.102.

- Knut Reinert, Ben Langmead, David Weese, and Dirk J. Evers. Alignment of next-generation sequencing reads. Annual Review of Genomics and Human Genetics, 16:133–151, 8 2015. ISSN 1527-8204. doi: 10.1146/annurev-genom-090413-025358.
- Po-Yen Wu, John H. Phan, and May D. Wang. Assessing the impact of human genome annotation choice on rna-seq expression estimates. *BMC Bioinformatics*, 14(11):S8, Nov 2013. ISSN 1471-2105. doi: 10.1186/1471-2105-14-S11-S8. URL https://doi.org/10.1186/1471-2105-14-S11-S8.
- Shanrong Zhao and Baohong Zhang. A comprehensive evaluation of ensembl, refseq, and ucsc annotations in the context of rna-seq read mapping and gene quantification. *BMC Genomics*, 16(1):97, Feb 2015. ISSN 1471-2164. doi: 10.1186/s12864-015-1308-8. URL https://doi.org/10.1186/s12864-015-1308-8.