

# Mapping Reads

## Purpose of Read Mapping

- Identify the origin of a read
- RNA-seq: Infer the gene that was expressed and that generated the read
- BUT: Is it always possible to find the read origin?

# Sequence Alignment

- Optimization problem: Find the alignment with the **highest score**
- Global alignment (end-to-end): Needleman-Wunsch
- Local alignment: Smith-Waterman

Score matrix:

	a	c	g	t
a	+2	-2	-1	-2
c	-2	+2	-2	-1
g	-1	-2	+2	-2
t	-2	-1	-2	+2

Gap score: -3

c	c	t	a	g	t	t	c	a	t
				x		x			
t	g	t	a	a	t	-	c	a	c
		+2	+2	-1	+2	-3	+2	+2	

Alignment score = 6

## Local alignments ( Smith-Waterman)

- Generally useful for local alignments, determining regions of similarity between two strings (here, DNA sequence)
- Dynamic programming based algorithm: gives optimal alignment, with respect to scoring system
- Need to set scores for **match**, penalties for **mismatch** and **gap** (typically, mismatch penalties are set according to evolutionary knowledge and error profiles)
- Match: +1; Mismatch: -1; Gap: -2

Reference Sequence

	.	A	A	T	G	T
.	0	0	0	0	0	0
A	0					
T	0					
G	0					
A	0					
C	0					

$$H(i, j) = \max \begin{cases} 0 \\ H(i-1, j-1) + w(a_i, b_j) \\ H(i-1, j) + w(a_i, -) \\ H(i, j-1) + w(-, b_j) \end{cases}, 1 \leq i \leq m, 1 \leq j \leq n$$

Match/Mismatch

Deletion

Insertion

## Local alignments ( Smith-Waterman)

- Generally useful for local alignments, determining regions of similarity between two strings (here, DNA sequence)
- Dynamic programming based algorithm: gives optimal alignment, with respect to scoring system
- Need to set scores for **match**, penalties for **mismatch** and **gap** (typically, mismatch penalties are set according to evolutionary knowledge and error profiles)
- Match: +1; Mismatch: -1; Gap: -2

	.	A	A	T	G	T
.	0	0	0	0	0	0
A	0	1	1			
T	0					
G	0					
A	0					
C	0					

$$H(i, j) = \max \begin{cases} 0 \\ H(i-1, j-1) + w(a_i, b_j) \\ H(i-1, j) + w(a_i, -) \\ H(i, j-1) + w(-, b_j) \end{cases}, 1 \leq i \leq m, 1 \leq j \leq n$$

Match/Mismatch

Deletion

Insertion

## Local alignments ( Smith-Waterman)

- Generally useful for local alignments, determining regions of similarity between two strings (here, DNA sequence)
- Dynamic programming based algorithm: gives optimal alignment, with respect to scoring system
- Need to set scores for **match**, penalties for **mismatch** and **gap** (typically, mismatch penalties are set according to evolutionary knowledge and error profiles)
- Match: +1; Mismatch: -1; Gap: -2

	.	A	A	T	G	T
.	0	0	0	0	0	0
A	0	1	1	0	0	0
T	0	0	0	2	0	1
G	0					
A	0					
C	0					

$$H(i, j) = \max \begin{cases} 0 \\ H(i-1, j-1) + w(a_i, b_j) & \text{Match/Mismatch} \\ H(i-1, j) + w(a_i, -) & \text{Deletion} \\ H(i, j-1) + w(-, b_j) & \text{Insertion} \end{cases}, 1 \leq i \leq m, 1 \leq j \leq n$$

## Local alignments ( Smith-Waterman)

- Generally useful for local alignments, determining regions of similarity between two strings (here, DNA sequence)
- Dynamic programming based algorithm: gives optimal alignment, with respect to scoring system
- Need to set scores for **match**, penalties for **mismatch** and **gap** (typically, mismatch penalties are set according to evolutionary knowledge and error profiles)
- Match: +1; Mismatch: -1; Gap: -2

	.	A	A	T	G	T
.	0	0	0	0	0	0
A	0	1	1	0	0	0
T	0	0	0	2	0	1
G	0	0	0	0	3	1
A	0	1	1	0	1	2
C	0	0	0	0	0	0

$$H(i, j) = \max \begin{cases} 0 \\ H(i-1, j-1) + w(a_i, b_j) \\ H(i-1, j) + w(a_i, -) \\ H(i, j-1) + w(-, b_j) \end{cases}, 1 \leq i \leq m, 1 \leq j \leq n$$

Match/Mismatch

Deletion

Insertion

## Local alignments ( Smith-Waterman)

- Traceback:
- Start with maximum score (here, 3)
- Follow path that gives maximum score

	.	A	A	T	G	T
.	0	0	0	0	0	0
A	0	1	1	0	0	0
T	0	0	0	2	0	1
G	0	0	0	0	3	1
A	0	1	1	0	1	2
C	0	0	0	0	0	0

$$H(i, j) = \max \begin{cases} 0 \\ H(i-1, j-1) + w(a_i, b_j) \\ H(i-1, j) + w(a_i, -) \\ H(i, j-1) + w(-, b_j) \end{cases}, 1 \leq i \leq m, 1 \leq j \leq n$$

Match/Mismatch

Deletion

Insertion



## Local alignments ( Smith-Waterman)

- Traceback:
- Start with maximum score (here, 3)
- Follow path that gives multiple score
- This gives:
- AATGT-
- -ATGAC

	.	A	A	T	G	T
.	0	0	0	0	0	0
A	0	1	1	0	0	0
T	0	0	0	2	0	1
G	0	0	0	0	3	1
A	0	1	1	0	1	2
C	0	0	0	0	0	0

$$H(i, j) = \max \begin{cases} 0 \\ H(i-1, j-1) + w(a_i, b_j) \\ H(i-1, j) + w(a_i, -) \\ H(i, j-1) + w(-, b_j) \end{cases}, 1 \leq i \leq m, 1 \leq j \leq n$$

Match/Mismatch

Deletion

Insertion

## Local alignments ( Smith-Waterman)

- Other considerations:
- Natural extension is to have a different *gap opening* and a *gap extension* penalty (former generally being larger)
- GO penalty=-2; GE penalty=-1; Mismatch=-1; Match=1
- With above penalties, which is the best scoring alignment?

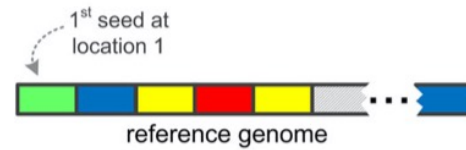
- |           |         |          |
|-----------|---------|----------|
| • AT-C-GT | ATC--GT | AT-C--GT |
| • ATTTTGT | ATTTTGT | ATT-TTGT |

## Alignment Tools

- Dynamic programming is slow
- Speed-up with heuristics
  - e.g. exactly align short subsequences and extend these alignments
    - e.g. BLAST / BLAT
  - no longer guaranteed to find the best alignments
  - exact matches are found by index lookup

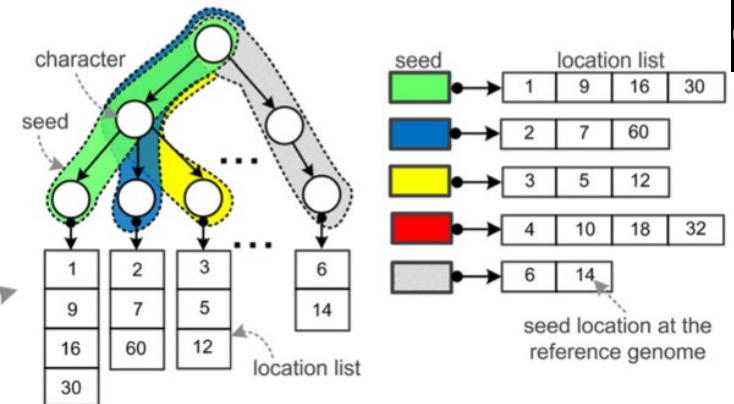
# Hashing based alignment

## a. Seed extraction from reference genome

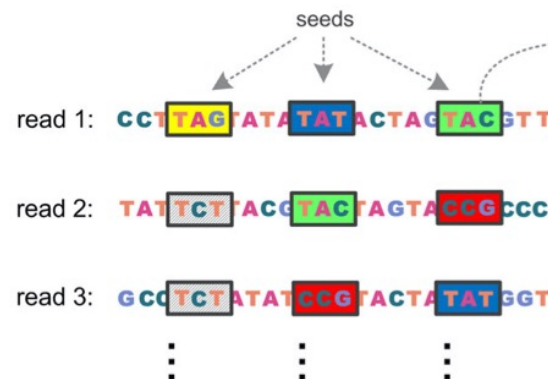


indexing

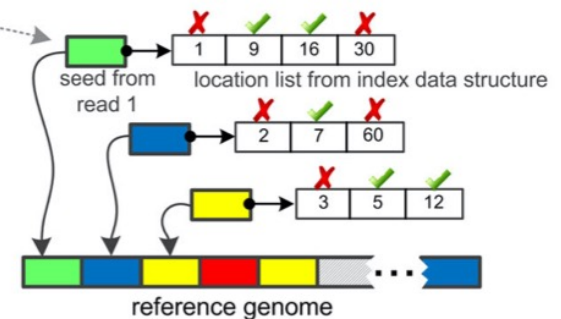
## b. Seed indexing using suffix tree or hash table



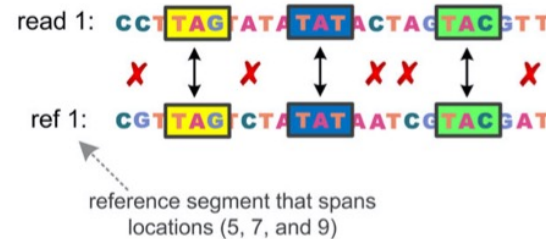
## c. Seed extraction from reads



## d. Seed querying and filtering



## e. Seed chaining and pre-alignment filtering



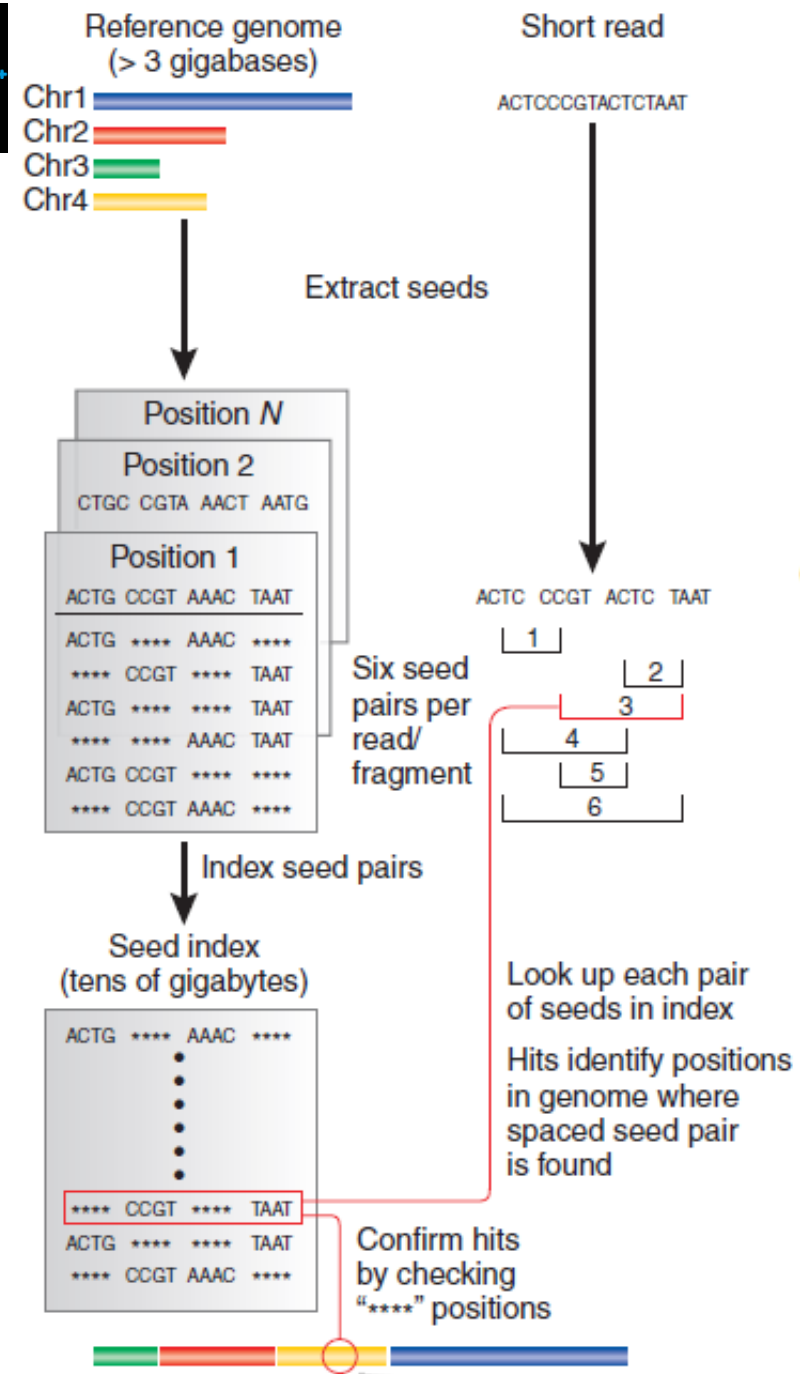
## f. Alignment verification

	C	G	T	T	A	G	T	C	T	A	...
C	0	0	0	0	0	0	0	0	0	0	...
G	0	2	2	2	2	2	2	2	2	2	...
T	0	2	3	3	3	3	3	3	4	4	...
A	0	2	3	5	5	5	5	5	5	6	...
T	0	2	3	5	7	7	7	7	7	7	...
A	0	3	3	5	7	9	9	9	9	9	...
G	0	2	4	5	7	9	11	11	11	11	...
T	0	2	4	6	7	9	11	13	13	13	...
A	0	2	4	6	7	9	11	13	14	14	...
T	0	2	4	6	8	9	11	13	14	16	...

.bam/.sam file contains necessary alignment information (e.g., type, location, and number of each edit)

# Index Genome: Spaced Seeds

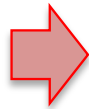
- Tags and tag-sized pieces of reference are cut into small “seeds.”
- Pairs of spaced seeds are stored in an index.
- Look up spaced seeds for each tag.
- For each “hit,” confirm the remaining positions.





# The Burrows-Wheeler transform (1994; 1983)

c a c a a c g \$



c	a	c	a	a	c	g	\$
a	c	a	a	c	g	\$	c
c	a	a	c	g	\$	c	a
a	a	c	g	\$	c	a	c
a	c	g	\$	c	a	c	a
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c
\$	c	a	c	a	a	c	g



\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c



g c c a a \$ a c

The “Last-First mapping” property

The relative ordering of a particular character (say **c**) in column 1 is the same as that in the last column

c<sub>1</sub> a c<sub>2</sub> a a c<sub>3</sub> g \$

\$	c <sub>1</sub>	a	c <sub>2</sub>	a	a	c <sub>3</sub>	g
a	a	c <sub>3</sub>	g	\$	c <sub>1</sub>	a	c <sub>2</sub>
a	c <sub>2</sub>	a	a	c <sub>3</sub>	g	\$	c <sub>1</sub>
a	c <sub>3</sub>	g	\$	c <sub>1</sub>	a	c <sub>2</sub>	a
c <sub>2</sub>	a	a	c <sub>3</sub>	g	\$	c <sub>1</sub>	a
c <sub>1</sub>	a	c <sub>2</sub>	a	a	c <sub>3</sub>	g	\$
c <sub>3</sub>	g	\$	c <sub>1</sub>	a	c <sub>2</sub>	a	a
g	\$	c <sub>1</sub>	a	c <sub>2</sub>	a	a	c <sub>3</sub>



## The “Last-First mapping” property

The relative ordering of a particular character (say **c**) in column 1 is the same as that in the last column

$c_1$  a  $c_2$  a a  $c_3$  g \$

\$	$c_1$	a	$c_2$	a	a	$c_3$	g
a	a	$c_3$	g	\$	$c_1$	a	$c_2$
a	$c_2$	a	a	$c_3$	g	\$	$c_1$
a	$c_3$	g	\$	$c_1$	a	$c_2$	a
$c_2$	a	a	$c_3$	g	\$	$c_1$	a
$c_1$	a	$c_2$	a	a	$c_3$	g	\$
$c_3$	g	\$	$c_1$	a	$c_2$	a	a
g	\$	$c_1$	a	$c_2$	a	a	$c_3$

Proof:

Suppose  $cX$  and  $cY$  are cyclic permutations of the input  $T$ .

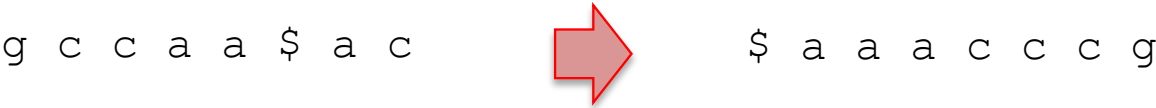
Suppose  $cX < cY$  (in lexicographical ordering)

Then  $Xc < Yc$  (in lexicographical ordering)

The LF-mapping property follows.

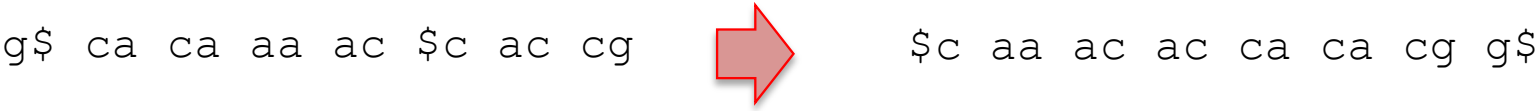
BWT is reversible

\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c



BWT is reversible

\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c



BWT is reversible

\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c

## Lookup AAC

\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c

Range  $\leftarrow$  range of last character in 1<sup>st</sup> column

While characters left (and nonzero range):

Lookup first and last match to preceding character in final column

Range  $\leftarrow$  LF-mapping of first and last match

## Lookup AAC

	\$	c	a	c	a	a	c	g
	a	a	c	g	\$	c	a	c
	a	c	a	a	c	g	\$	c
	a	c	g	\$	c	a	c	a
→	c	a	a	c	g	\$	c	a
	c	a	c	a	a	c	g	\$
→	c	g	\$	c	a	c	a	a
	g	\$	c	a	c	a	a	c


Range  $\leftarrow$  range of last character in 1<sup>st</sup> column

While characters left (and nonzero range):

Lookup first and last match to preceding character in final column

Range  $\leftarrow$  LF-mapping of first and last match

## Lookup AAC



\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c

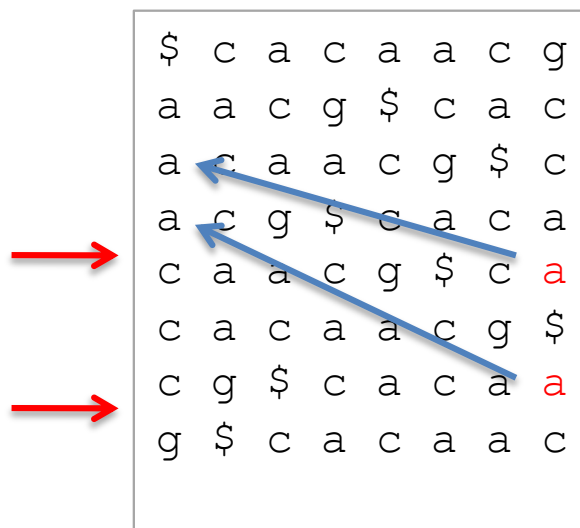
Range  $\leftarrow$  range of last character in 1<sup>st</sup> column

While characters left (and nonzero range):

Lookup first and last match to preceding character in final column

Range  $\leftarrow$  LF-mapping of first and last match

## Lookup AAC



Range  $\leftarrow$  range of last character in 1<sup>st</sup> column

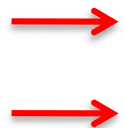
While characters left (and nonzero range):

Lookup first and last match to preceding character in final column

Range  $\leftarrow$  LF-mapping of first and last match



## Lookup A**AC**



\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c

Range  $\leftarrow$  range of last character in 1<sup>st</sup> column

While characters left (and nonzero range):

Lookup first and last match to preceding character in final column

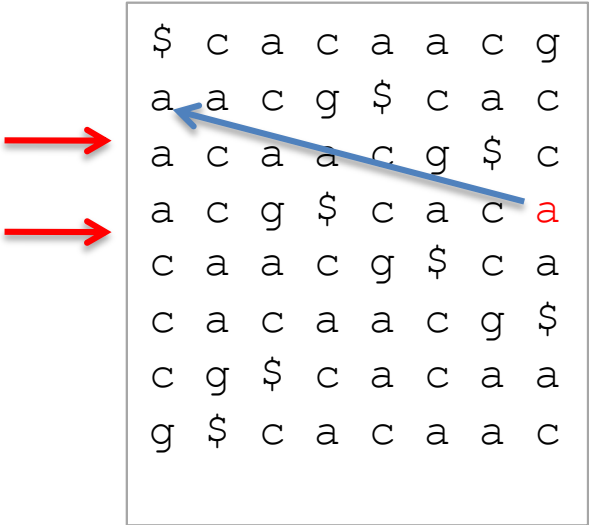
Range  $\leftarrow$  LF-mapping of first and last match

Lookup AAC

	\$	c	a	c	a	a	c	g
	a	a	c	g	\$	c	a	c
→	a	c	a	a	c	g	\$	c
→	a	c	g	\$	c	a	c	a
	c	a	a	c	g	\$	c	a
	c	a	c	a	a	c	g	\$
	c	g	\$	c	a	c	a	a
	g	\$	c	a	c	a	a	c

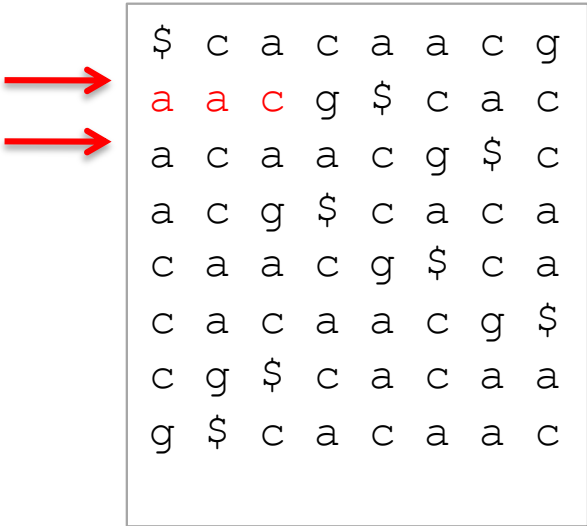
Range ← range of last character in 1<sup>st</sup> column  
While characters left (and nonzero range):  
    Lookup first and last match to preceding character in final column  
    Range ← LF-mapping of first and last match

# Lookup AAC



Range  $\leftarrow$  range of last character in 1<sup>st</sup> column  
While characters left (and nonzero range):  
    Lookup first and last match to preceding character in final column  
    Range  $\leftarrow$  LF-mapping of first and last match

Lookup AAC



\$	c	a	c	a	a	c	g
a	a	c	g	\$	c	a	c
a	c	a	a	c	g	\$	c
a	c	g	\$	c	a	c	a
c	a	a	c	g	\$	c	a
c	a	c	a	a	c	g	\$
c	g	\$	c	a	c	a	a
g	\$	c	a	c	a	a	c

Range ← range of last character in 1<sup>st</sup> column  
While characters left (and nonzero range):  
    Lookup first and last match to preceding character in final column  
    Range ← LF-mapping of first and last match

## Comparison

### Spaced seeds

- Depending on choices up to 50Gb of memory for humang genome
- More straightforward to program.
- Example:
  - Subread
- More tolerant to
  - sequence variations
  - sequencing errors

### Burrows-Wheeler

- The Burrows-Wheeler transform can be compressed to obtain a compact presentation of the index. This is then called **FM-Index**, ('Full-text index in Minute space', developed by Ferragina and Manzini, 2000
- Compressed index requires <2Gb of memory for human genome.
- can be very fast
- More complicated to program.
- Examples:
  - bowtie2
  - BWA

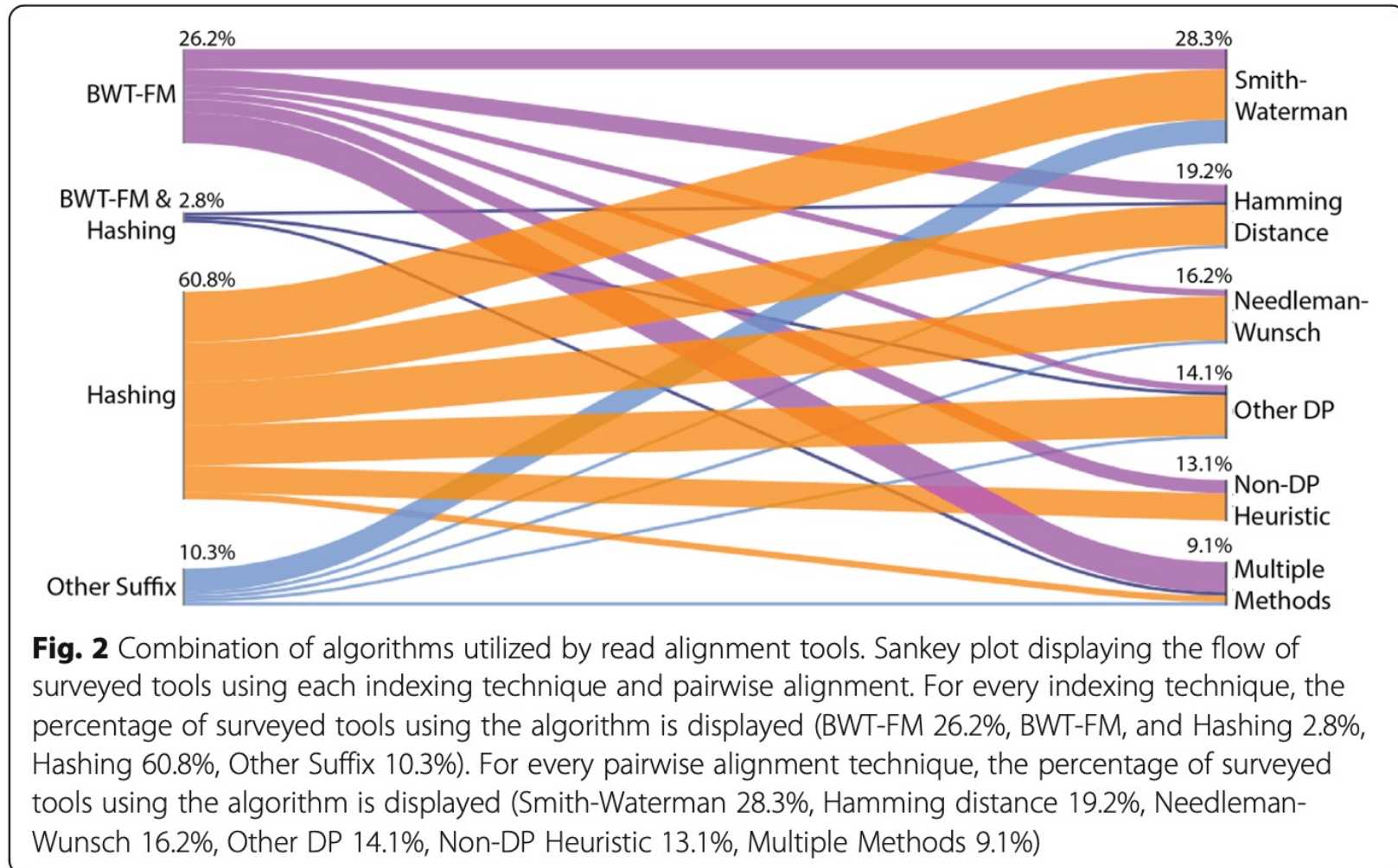
## Characteristics of algorithms

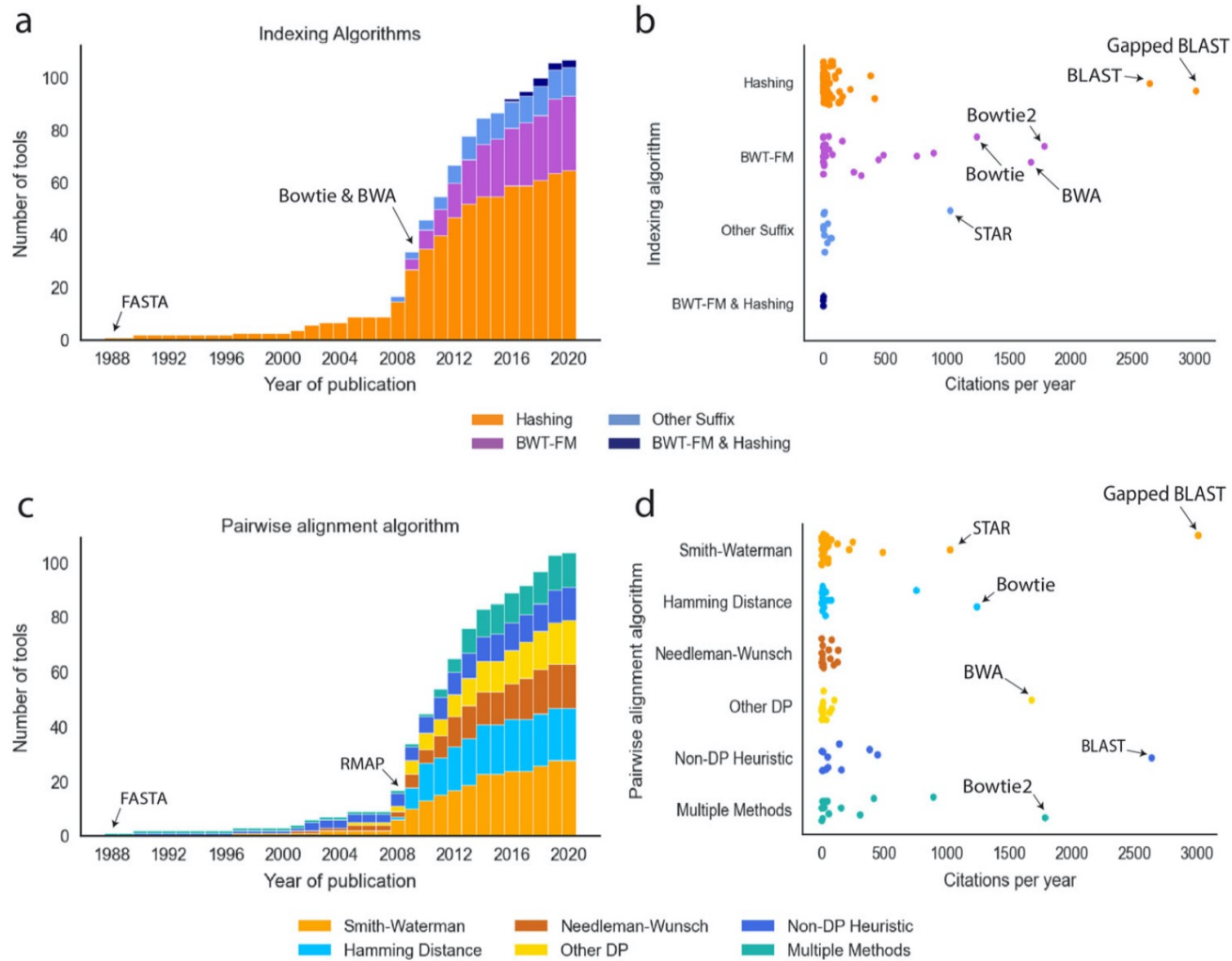
	Hashing	Suffix tree and BWT-FM
Easy to implement	Yes	No
Search for exact/inexact match	Exact	Exact and inexact
Index size	Large	Compressed (small)
Indexing time	Small	Large
Seed query speed	O(1), fast	Slow
Seed length	Fixed length per index	Can be fixed or variable

Alser et al.

<https://genomebiology.biomedcentral.com/articles/10.1186/s13059-021-02443-7>

## Combining index lookup and pairwise alignment







## Alignment with Mismatches

- Mismatches can occur because of
  - sequencing error (error rate  $\sim 1/1000$ )
  - mutation; (human mutation rate  $\sim 1/1e4$ )
  - $\rightarrow$  if reads are long ( $> 100nt$ ) reads with mismatches will not be rare; more than  $\sim 10\%$  of the reads may have a mismatch
- If there is a sequence mismatch the index lookup fails! How to find nevertheless an alignment?

## BWT Alignment with Mismatches

- Strategies:
  - If the BWT lookup fails at position k, try all different bases at position k  
→ drawback: computing effort grows exponentially with the number of mismatches  
→ implemented e.g. in bowtie
  - Chop reads in segments (seeds) and align those mismatch-free and stitch seed alignments together  
→ implemented e.g. in bowtie

## Consideration: Sequencing Errors – PHRED quality

- Each called base is given a quality score Q
- $Q = 10 * \log_{10}(\text{estimated probability that call is wrong})$

10 prob = 0.1

20 prob = 0.01

30 prob = 0.001

[Q30 often used as a threshold for useful sequence data]

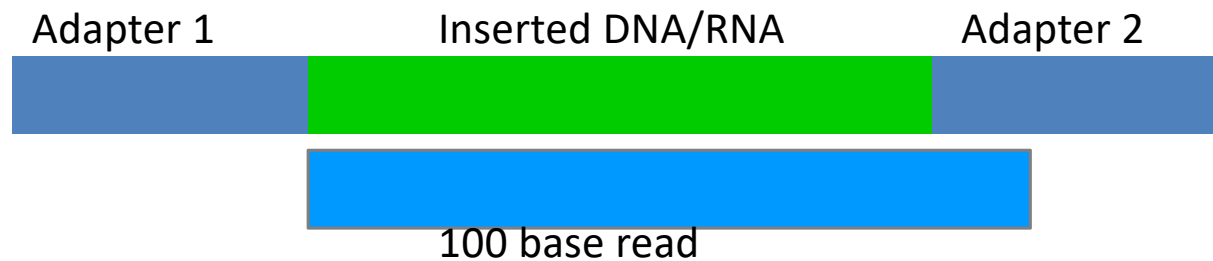
→ down-weight low-quality bases when computing the alignment score

- These quality scores are called PHRED scores.
- PHRED scores are determined by the sequencer that directly rates how reliable the measured signal is

## Considerations: Read trimming

- Sequencers have systematic errors
- Illumina sequencers have a higher error rate at the first few bases (1-5)
- Basically, all sequencers have increasing error rate towards the end of the read
- Hard trimming: trim a fixed number of bases from the beginning and/or end
- Quality trimming: cut the end of the read as soon as the base quality drops below a threshold

## Considerations: Read trimming



- If the inserted DNA/RNA fragment is too short, the read will contain part of the adapter
- Adapter trimming can be challenging if
  - if the insert is 90 – 100 bases
  - if the read has many sequencing errors

## Multiple Alignments

- A read may have multiple valid alignments with identical or similarly good alignment scores
- Aligners allow to choose different reporting strategies:
  - Randomly select one alignment from the top-scoring alignments
  - Report all alignments that are within *delta* of the top-scoring alignment; clip if more than *Nmax* alignments are found
  - Report only alignments if they are unique (no other alignment within *delta* of the alignment score)
  - Do not report anything if more than *Nmax* valid alignments are found
  - ...
- Whether a read has a **unique alignment** depends on
  - the read sequence and the sequence homology of the organism
  - the search algorithm of the aligner
  - the reporting strategy of the aligner

## Read Alignment Summary (I)

### How to map billions of reads?

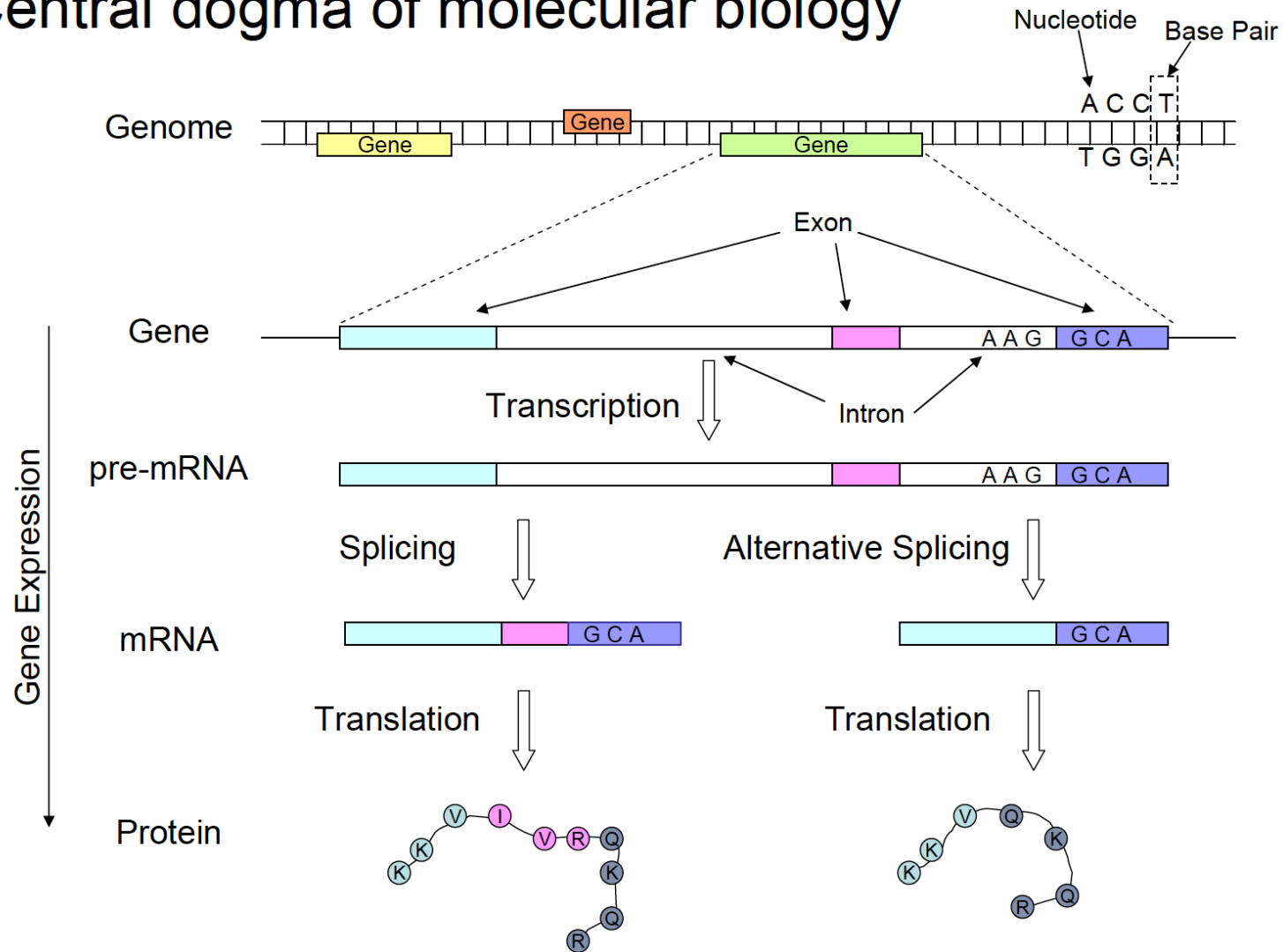
- The major aligners use the Burrows-Wheeler-Transform to create an index of the reference. Even for the human genome the index fits into 3GB RAM.
- Reads are aligned by **index lookup** not by sequence comparison
- The lookup of a perfect match read is faster than loading the read and writing the alignment coordinates to the output file!
- If the lookup fails because of SNPs or sequencing errors, an actual sequence comparison is done –and this can take time!

## Read Alignment Summary (II)

- Alignment is a score optimization problem
- Aligners find the alignment with the **highest alignment score**
- Aligners do not run an extensive search, they use **index lookup and heuristics** to find the alignment with highest score; these heuristics are not guaranteed to find the alignment with the maximum score
- If the heuristics are well chosen and the scores are well defined, then the alignment will be very often the highest-scoring alignment and that will be the correct alignment



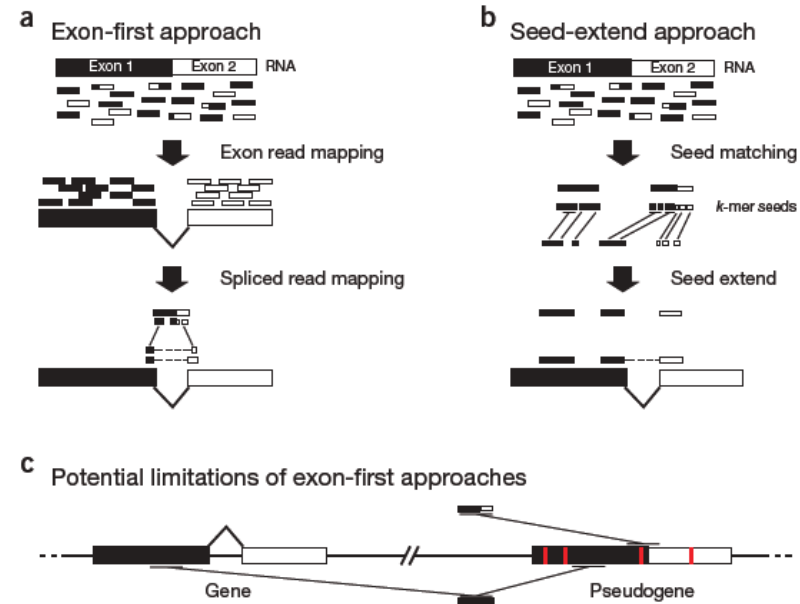
# Central dogma of molecular biology



92–94% of human genes undergo alternative splicing,

## RNA-seq Mapping

- Mapping targets:
  - transcriptome
  - splice junction library
  - genome
- Mouse retina 60+60bp reads:
  - 41 Mio of 91 Mio map to junctions

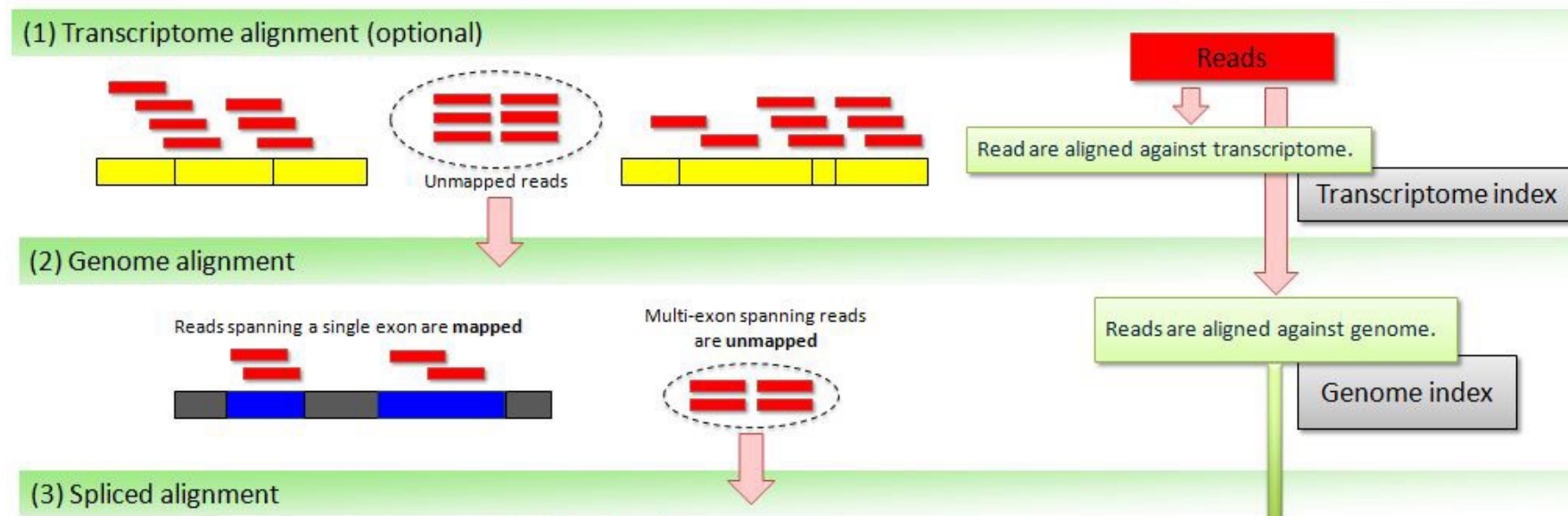


### Strategies:

- Exon-first: fast
- Seed-extend:
  - good with polymorphisms
  - simultaneous spliced/unspliced mapping

## Tophat2 pipeline – illustrative of the challenge of mapping RNA-seq reads

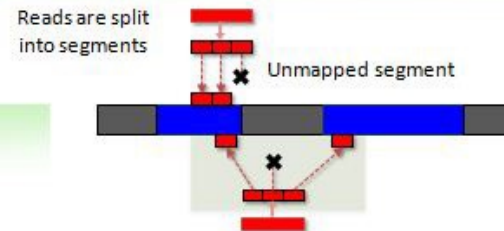
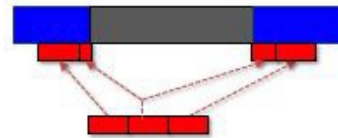
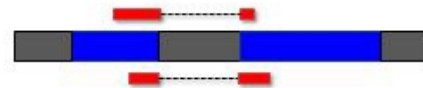
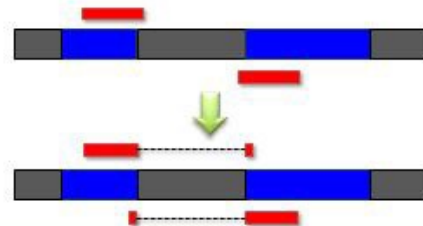
- Note: Tophat2 is outdated, use STAR, HISAT, Subread, ...



## Tophat2 pipeline

## (3) Spliced alignment

(3-1) Segment alignment to genome

(3-2) Identification of splice sites  
(including indels and fusion break points)(3-3) Segments aligned to junction  
flanking sequences(3-4) Segment alignments stitched  
together to form whole read alignments(3-5) Re-alignment of reads minimally  
overlapping intronsReads are split into smaller segments  
which are then aligned to the genome.

Genome index

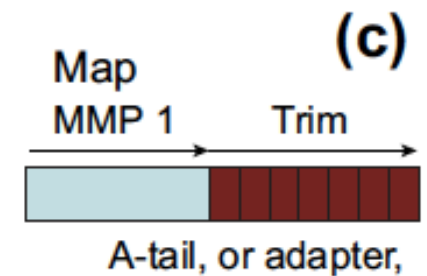
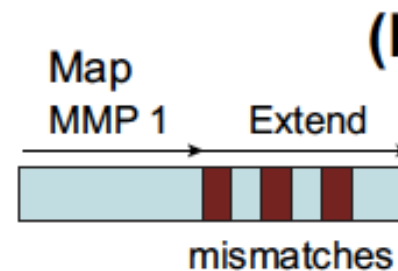
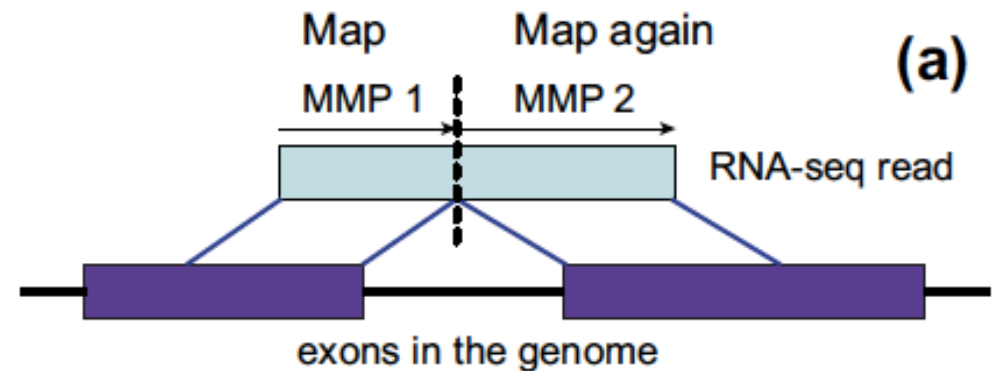
Segment mappings are used to find potential splice sites  
usually when the distance between the mapped positions  
of the left and the right segments are longer than the  
length of the middle part of a read.Sequences flanking a splice site are concatenated  
and segments are aligned to them.

Junction flanking index

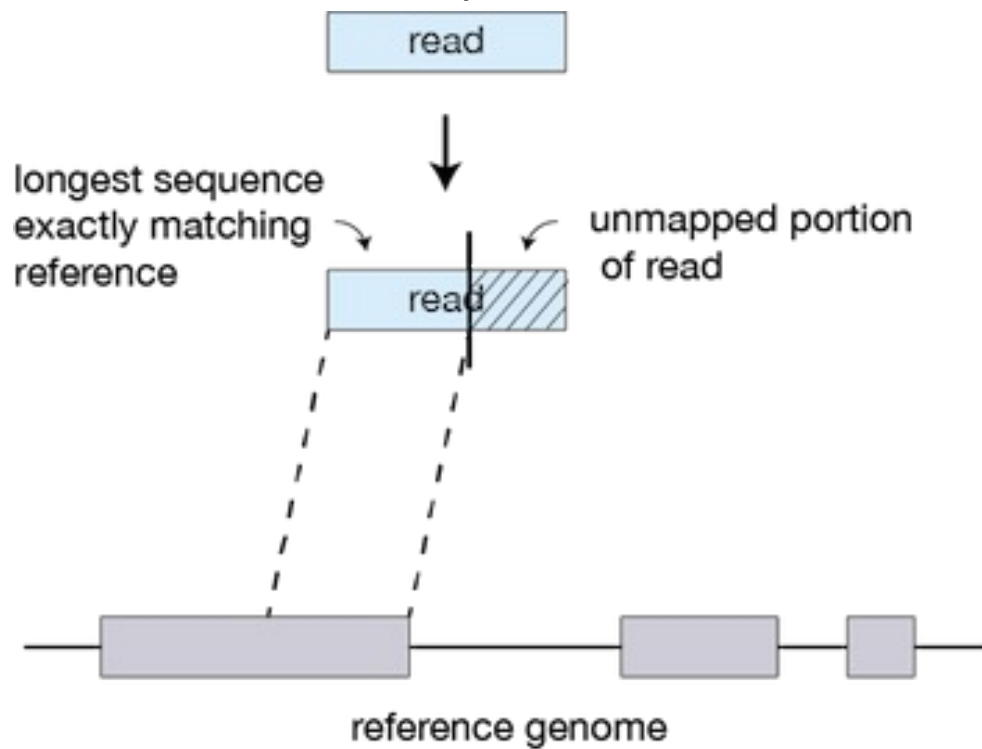
Mapped segments against either genome or flanking  
sequences are gathered to produce whole read alignments.Genome mapped reads with alignments extending a few  
bases into introns are re-aligned to exons instead.

## STAR: universal RNA-seq aligner

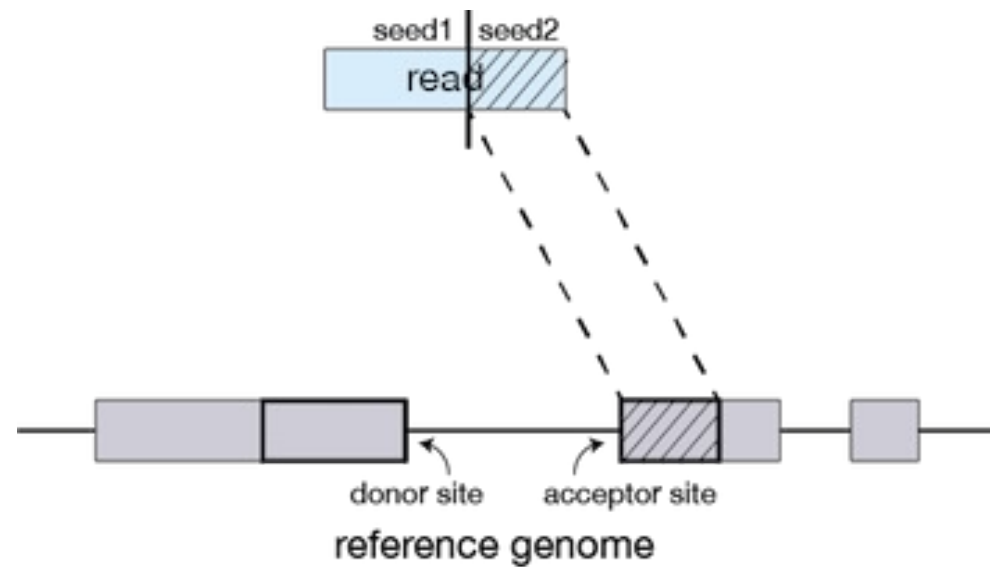
- Designed to align the non-contiguous sequences directly to the reference genome
- Steps:
  - Search Maximal Mappable Prefix (MMP)
  - clustering/stitching/scoring



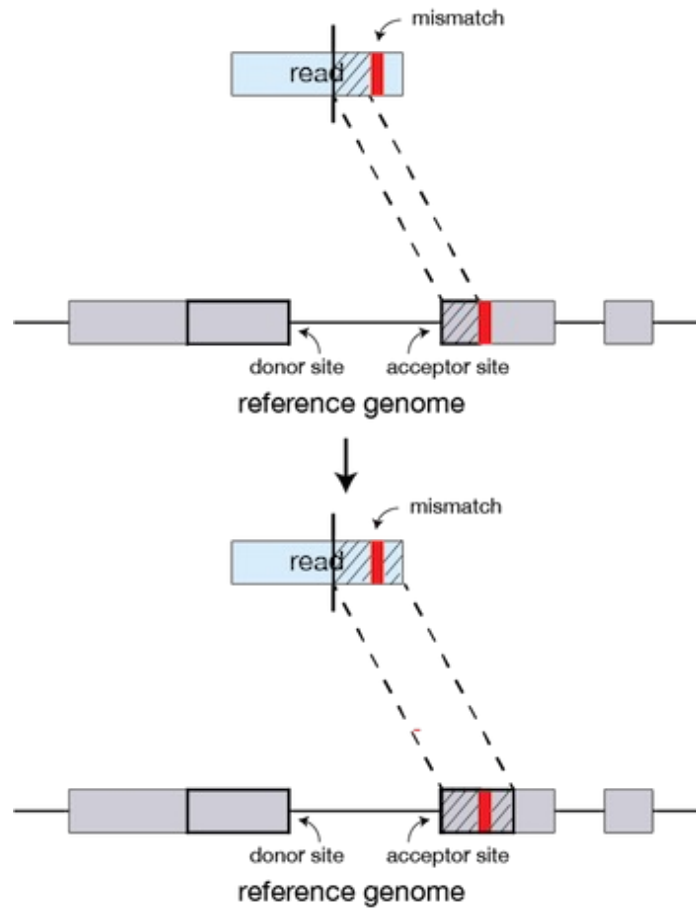
## Seed lookup



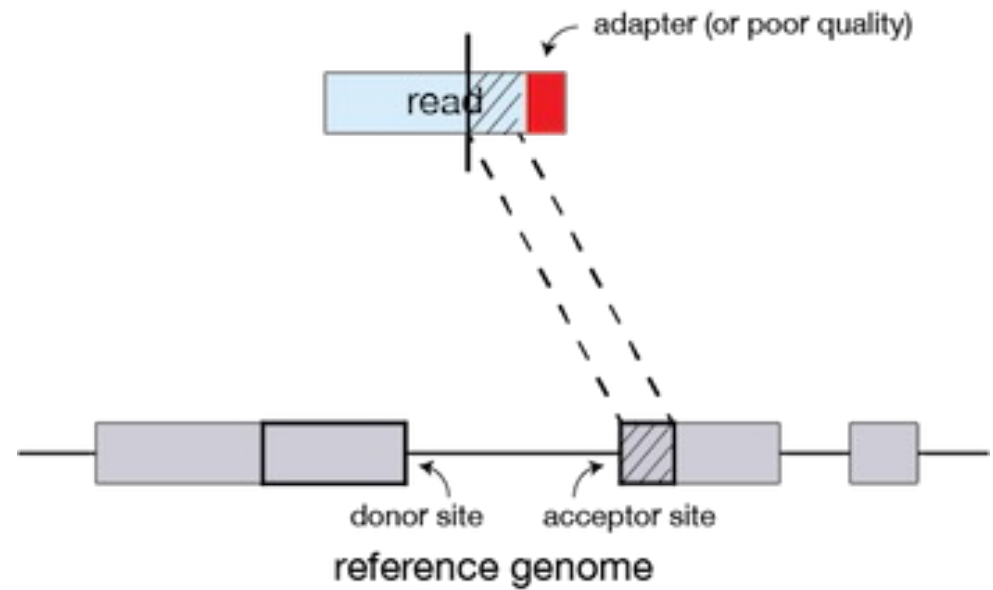
## Reseed



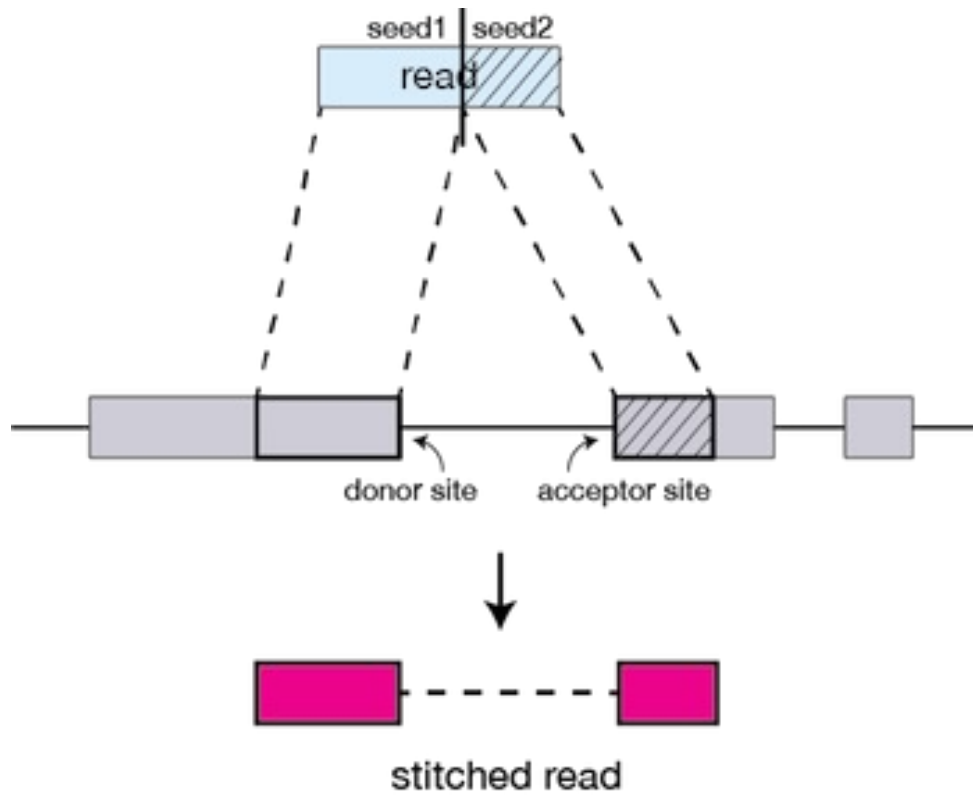
## Handle mismatches



## Soft clip ends



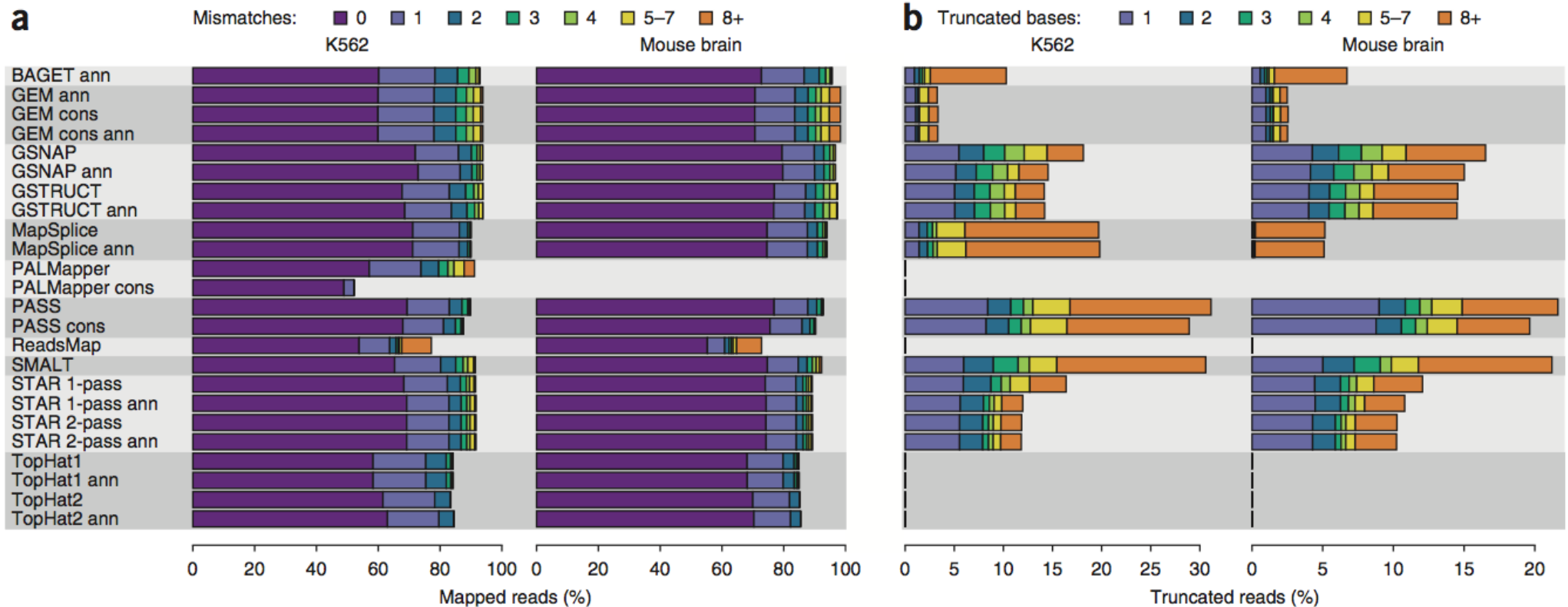
## Stitch spliced







## Performance comparison of RNA-seq mappers



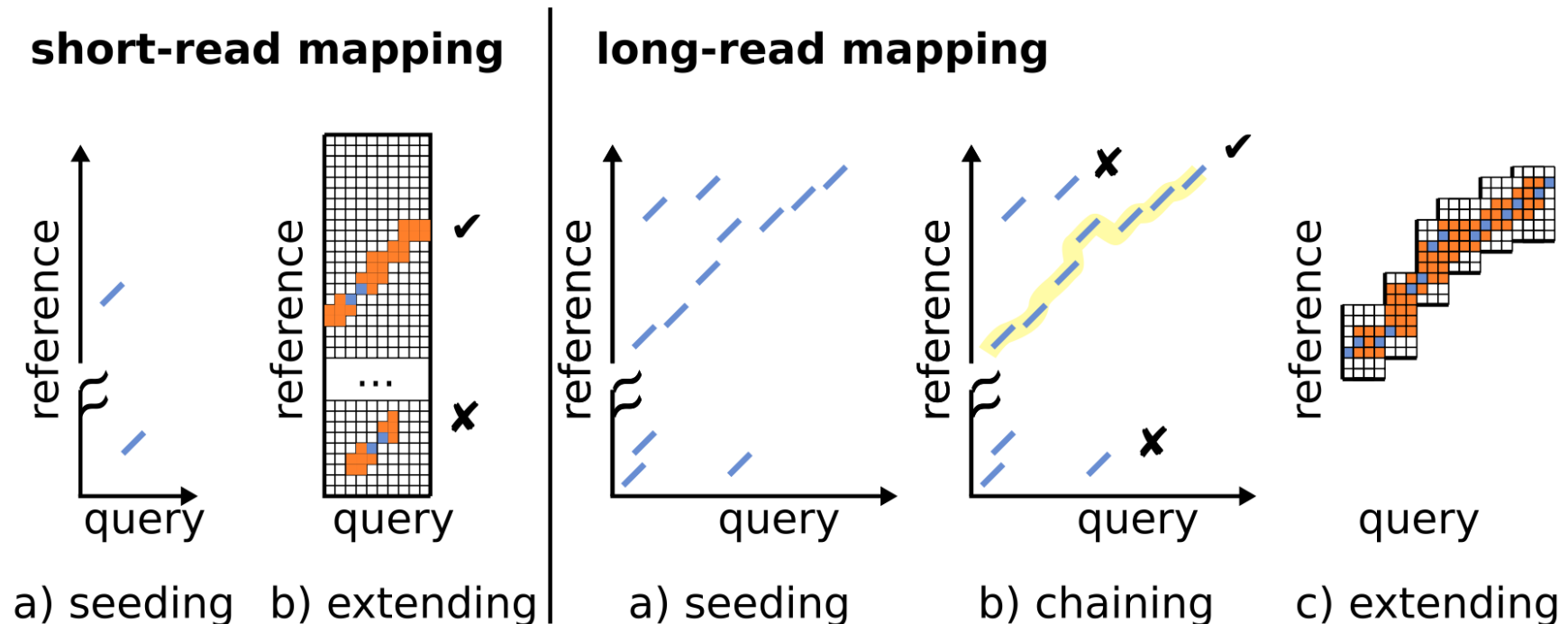
**Figure 2** | Mismatch and truncation frequencies. **(a)** Percentage of sequenced reads mapped with the indicated number of mismatches. **(b)** Percentage of sequenced reads truncated at either or both ends. Bar colors indicate the number of bases removed.

Engström et al. Systematic evaluation of spliced alignment programs for RNA-seq data. Nat Methods. 2013.

## STAR 2-pass mode

- goal: improve alignment around novel splice junctions
- does not discover more splice junctions but aligns more reads to discovered splice junctions
- initial mapping might fail to align reads where only a short part of the read spans the splice junction

# Mapping Long Reads



Mainly, short-read approaches extend (orange parts) from a single anchor (in blue) on the whole read length while long-read approaches gather multiple anchors, and chain (yellow line) them in for a candidate extending procedure that is done between pairs of anchors

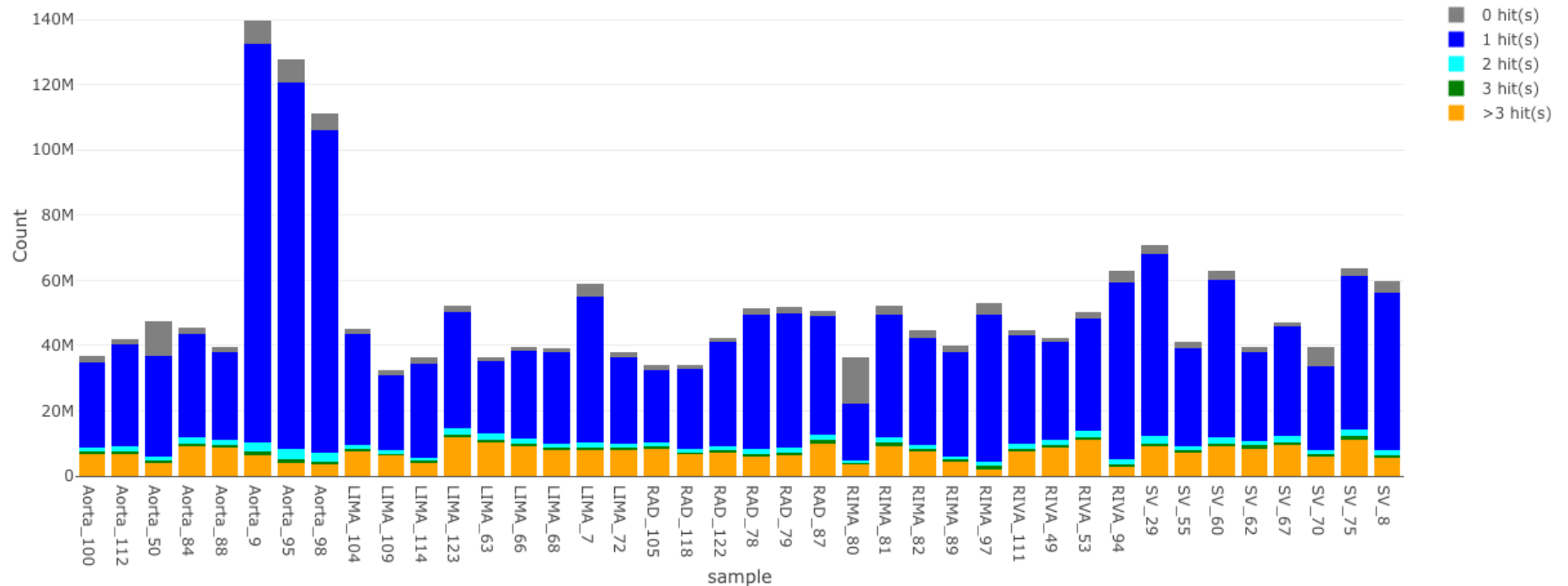
# How long should a read be so that it can be uniquely aligned??

## Current practice

use  $\geq 75$ nt for aligning RNA-reads to mammals

use  $\geq 100+100$ nt for aligning DNA reads for variant calling

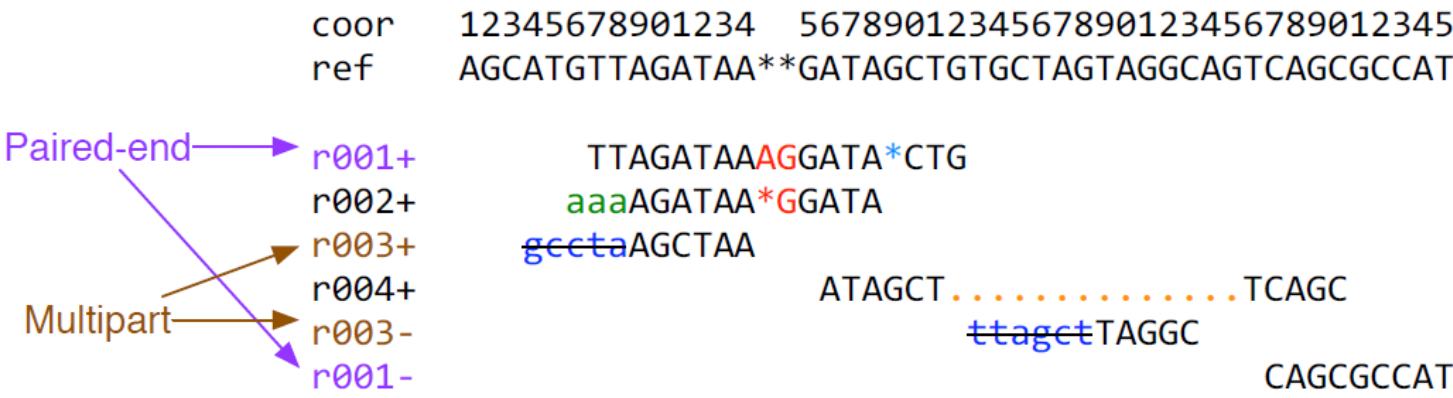
use 21nt for aligning micro-RNA reads to miRbase



## Sequence / Alignment (SAM) files

- **SAM (Sequence Alignment/Map)**
  - Single unified format for storing read alignments to a reference genome
  - Large plain text file
- **BAM (Binary Alignment/Map)**
  - Binary equivalent of SAM
  - Compressed data plus index (bai)
  - Developed for fast processing/indexing

# An example of read mapping







# CIGAR string - compact representation of an alignment

- M - match or mismatch
- I – insertion
- D – deletion
- S - soft clip
  - Clipped sequences stored in SAM
- H - hard clip
  - Clipped sequences not stored in SAM
- N – skipped reference bases, splicing

Match/mismatch, indels

Ref: ACGCA**T**GT**--**GT

Read: ATGCA**-**TG**C**AGT

Cigar: 5M**1**D2M**2**I2M

Soft clipping

REF: **A**TCGTGTAAC**C**TGACTAGT**TAA**

READ: **g**ggGTGTAAC**C**-GACTAG**g**ggg

Cigar: **3**S8M**1**D5M**4**S

Hard clipping

REF: **A**TCGTGTAAC**C**TGACTAGT**TAA**

READ: **g**ggGTGTAAC**C**-GACTAG**g**ggg

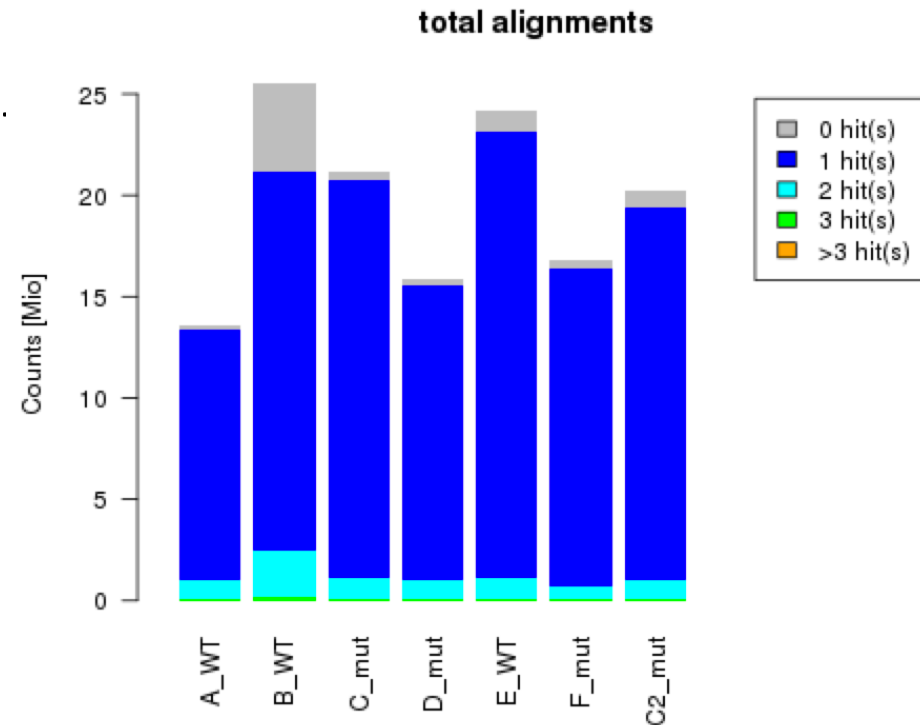
Cigar: **3**H8M**1**D5M**4**H

## Mapping quality – the MAPQ tag

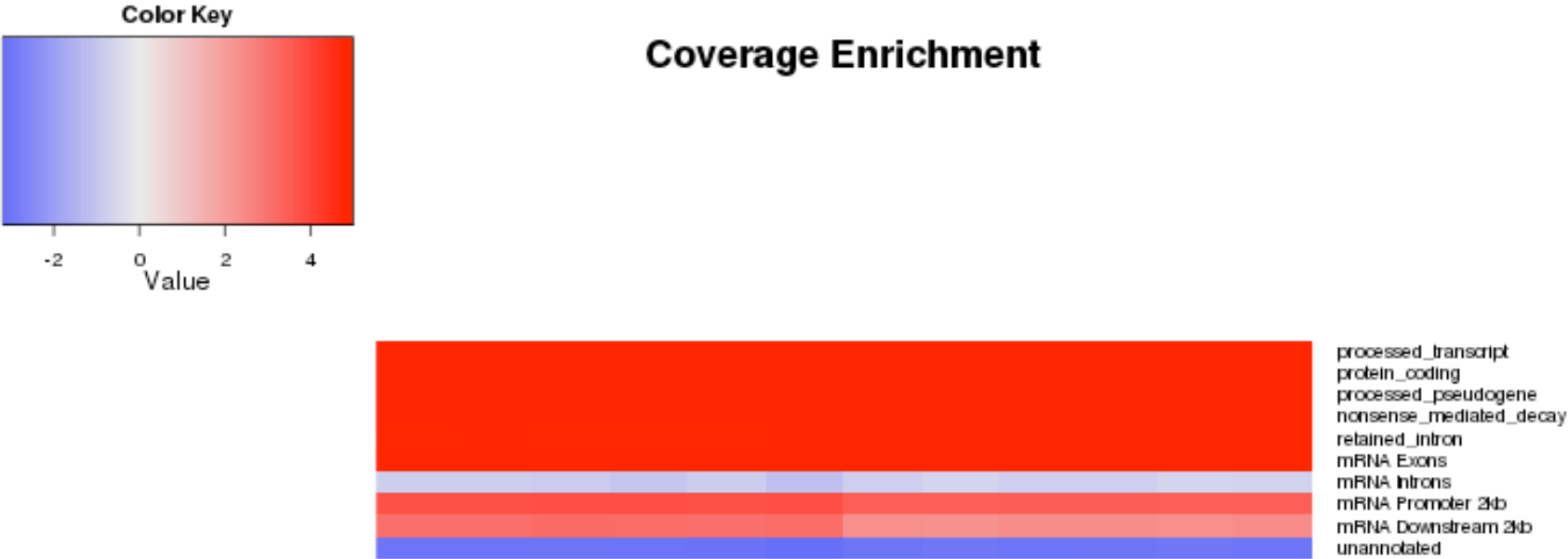
- The MAPQ tag does indicate the quality of an alignment
- Equivalent to the PHRED score for base
- However – the error probability can not be as accurately determined as the base quality
- The different aligners have different methods for assigning the quality

## Mapping QC: Mapping statistics

- How well did my sequence library align to my reference?
- Summary statistics (per read library)
  - % reads with unique alignment
  - % reads with multiple alignment
  - % reads with no alignment
  - % reads properly paired (for paired end libraries)



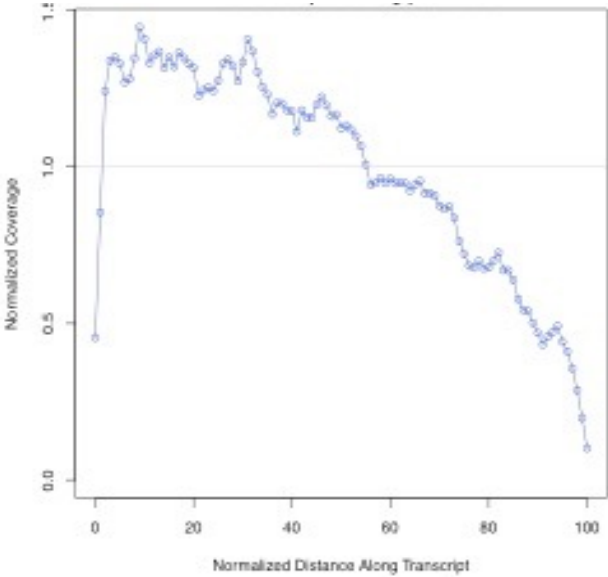
Mapping QC: Profiling efficiency by transcript annotation



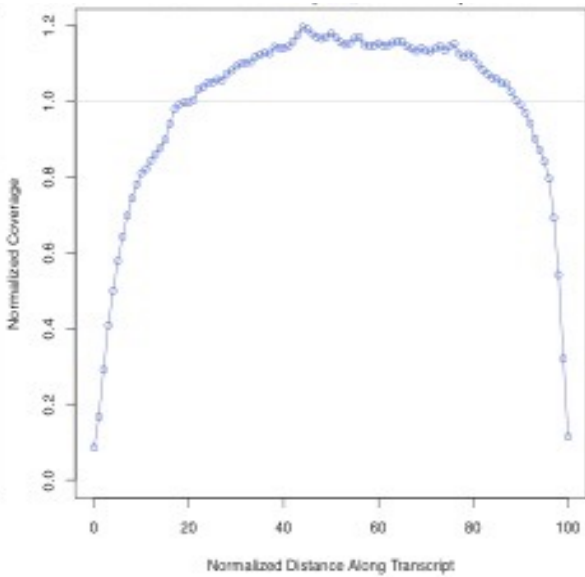
- Transcript annotation: intron, exon, up/down stream, unannotated
- Expression profiling efficiency

Mapping QC: Coverage bias

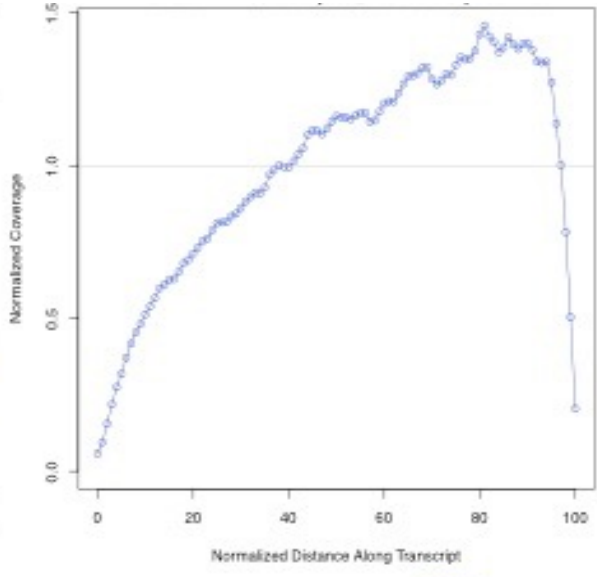
3' bias



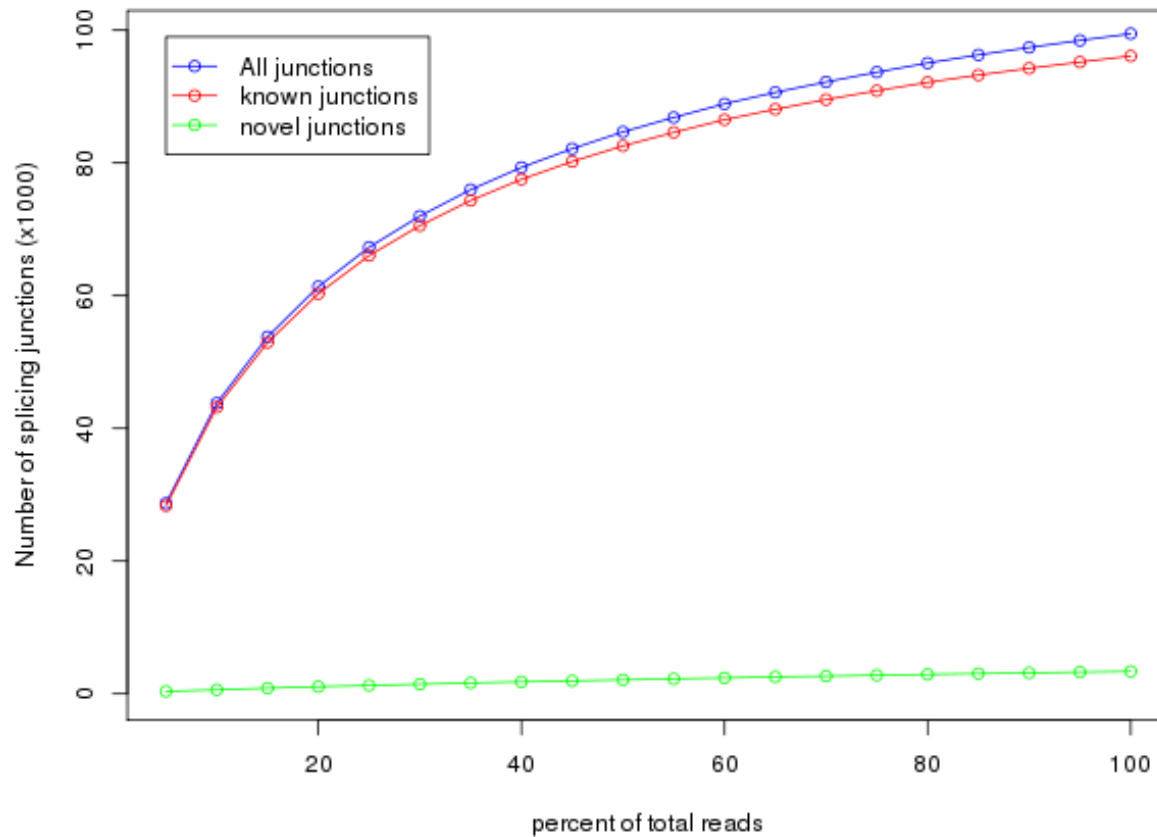
No bias



5' bias

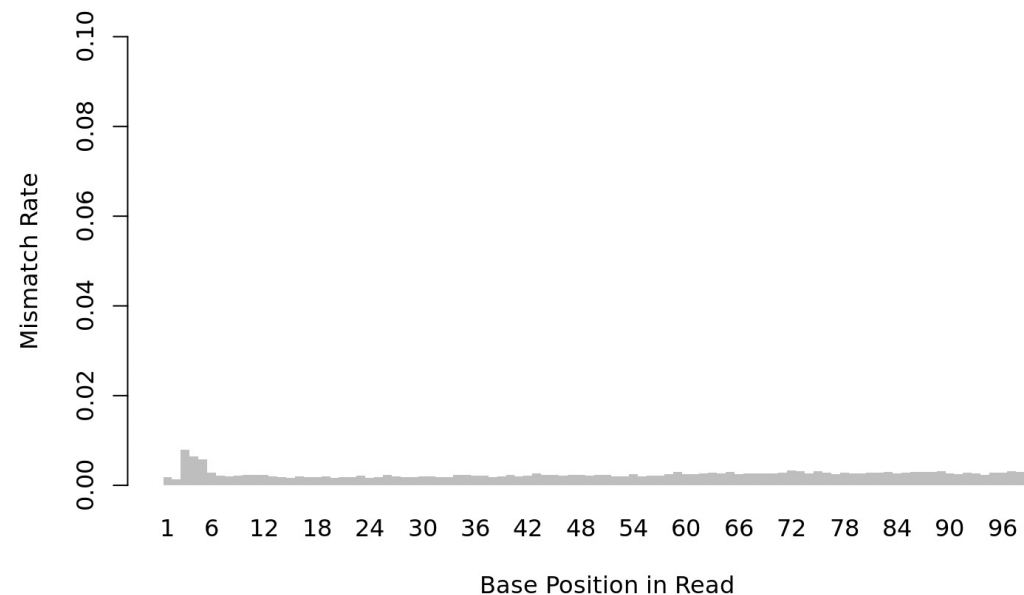
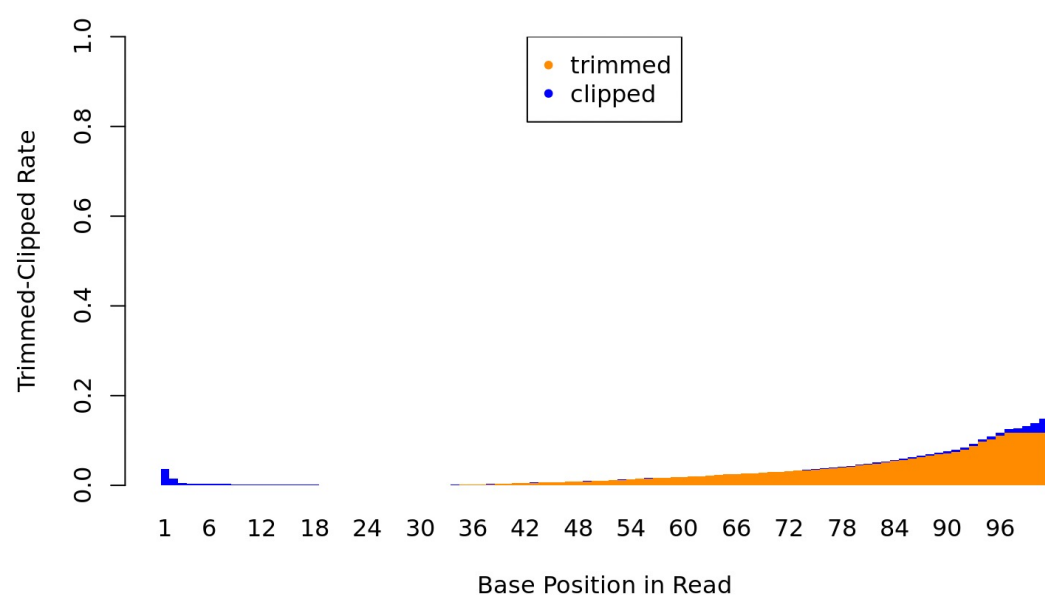


## Mapping QC: Junction saturation



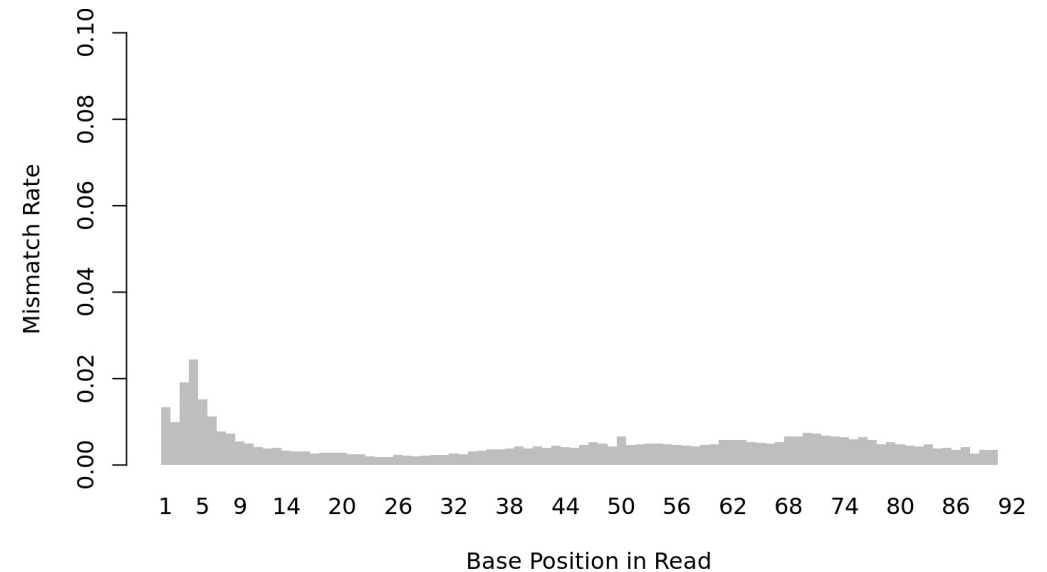
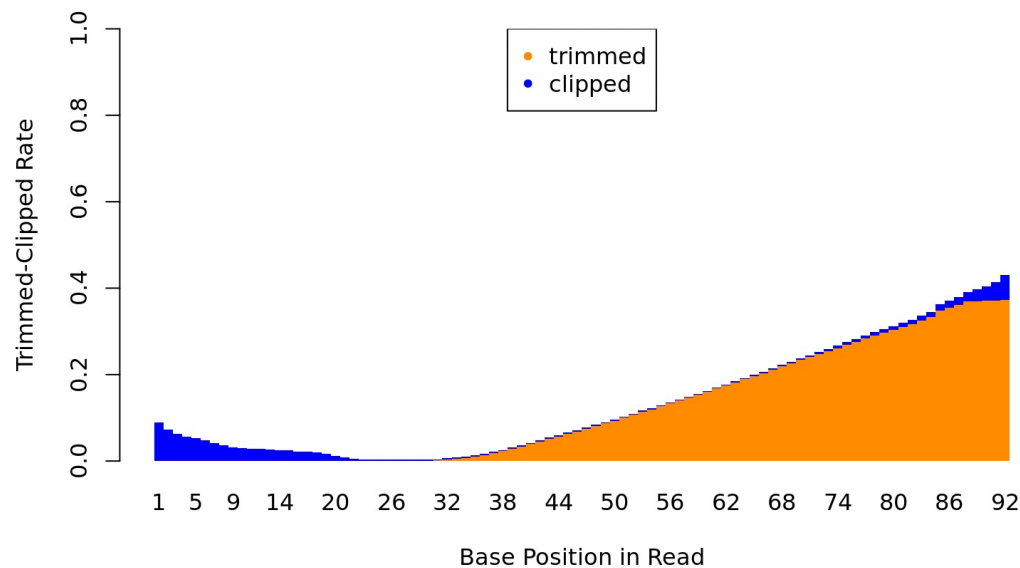
- Depth needed for alternative splicing analysis
- All annotated splice junctions are detected - a saturated RNA-seq dataset

## Statistics on Trimming, Soft Clipping and Mismatches



- Mismatches depend on position in the read
- If bases 1 or 2 do not match, the STAR aligner soft-clips them

## Statistics on Trimming, Soft Clipping and Mismatches



- Example from a low quality sample (low input <1ng, slight degradation)
- Short inserts, increased mismatch rate



## Where mappers have issues

- paralog genes that exist on X and Y chromosomes
- MHC genes (part of the adaptive immune response)
  - poly-morphic: several alleles (versions) exist on that gene's locus

# RNA-seq mapping QC tools

- RNA-SeQC
  - <https://github.com/getzlab/rnaseqc>
- RSeQC
  - <http://rseqc.sourceforge.net/>