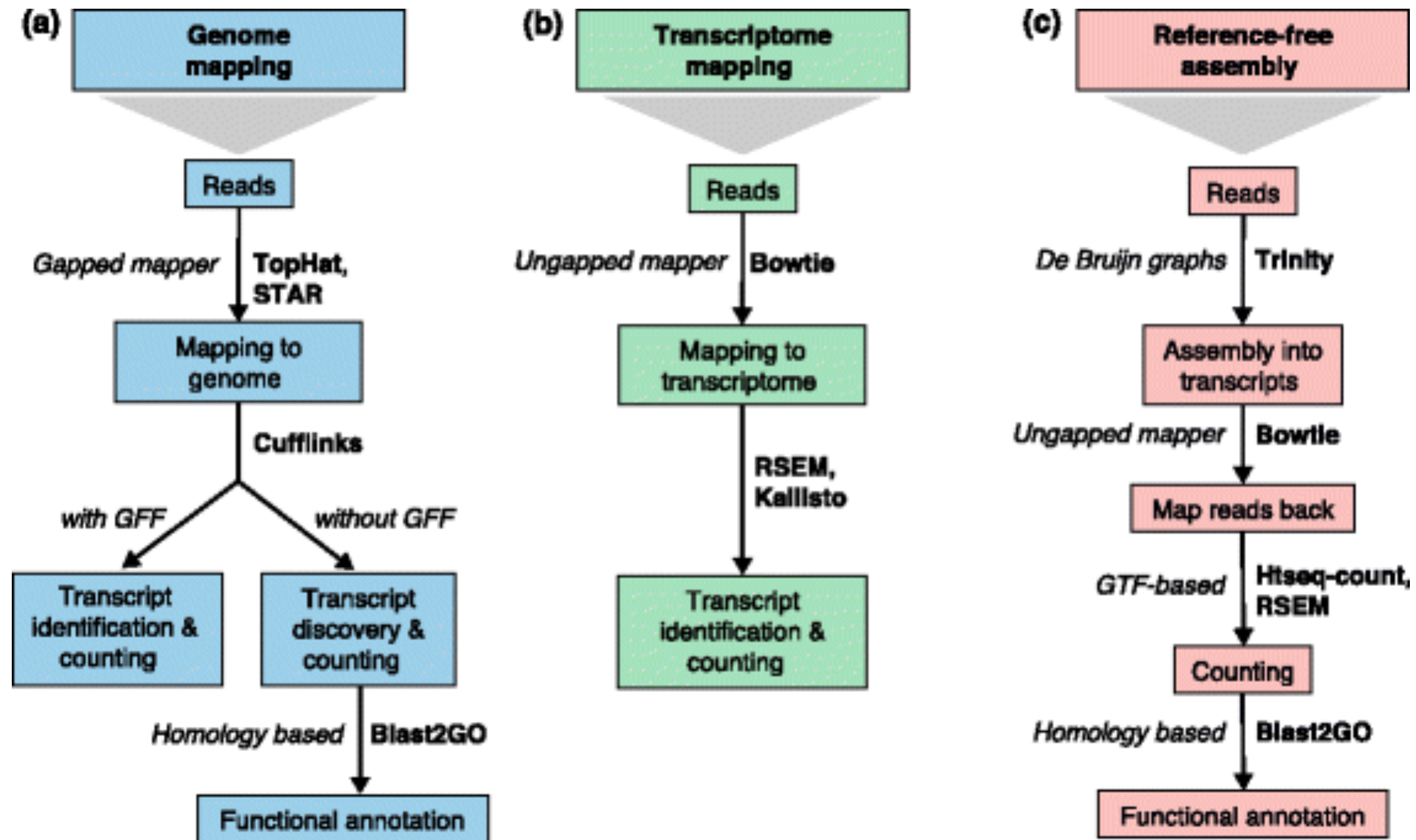
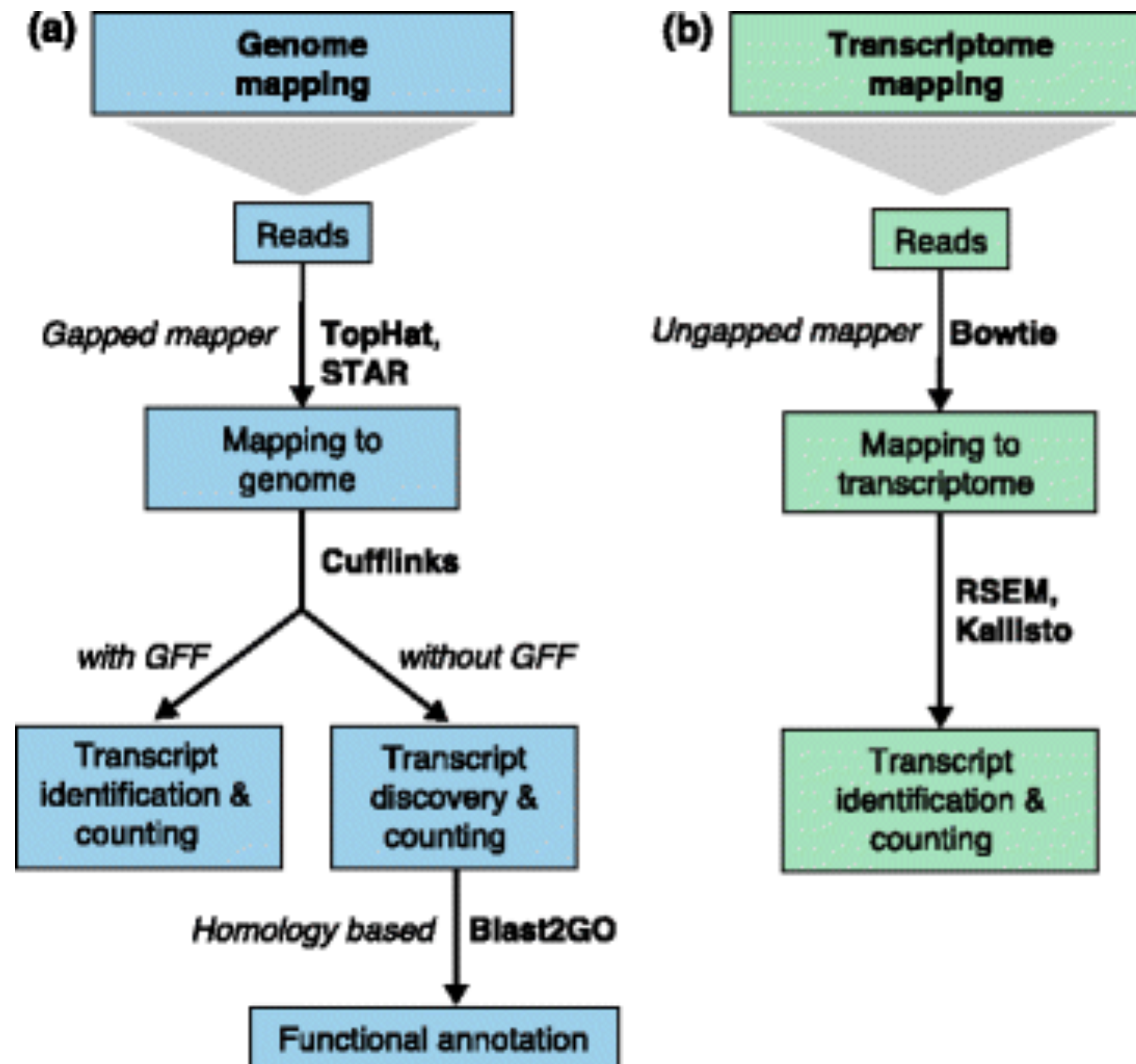


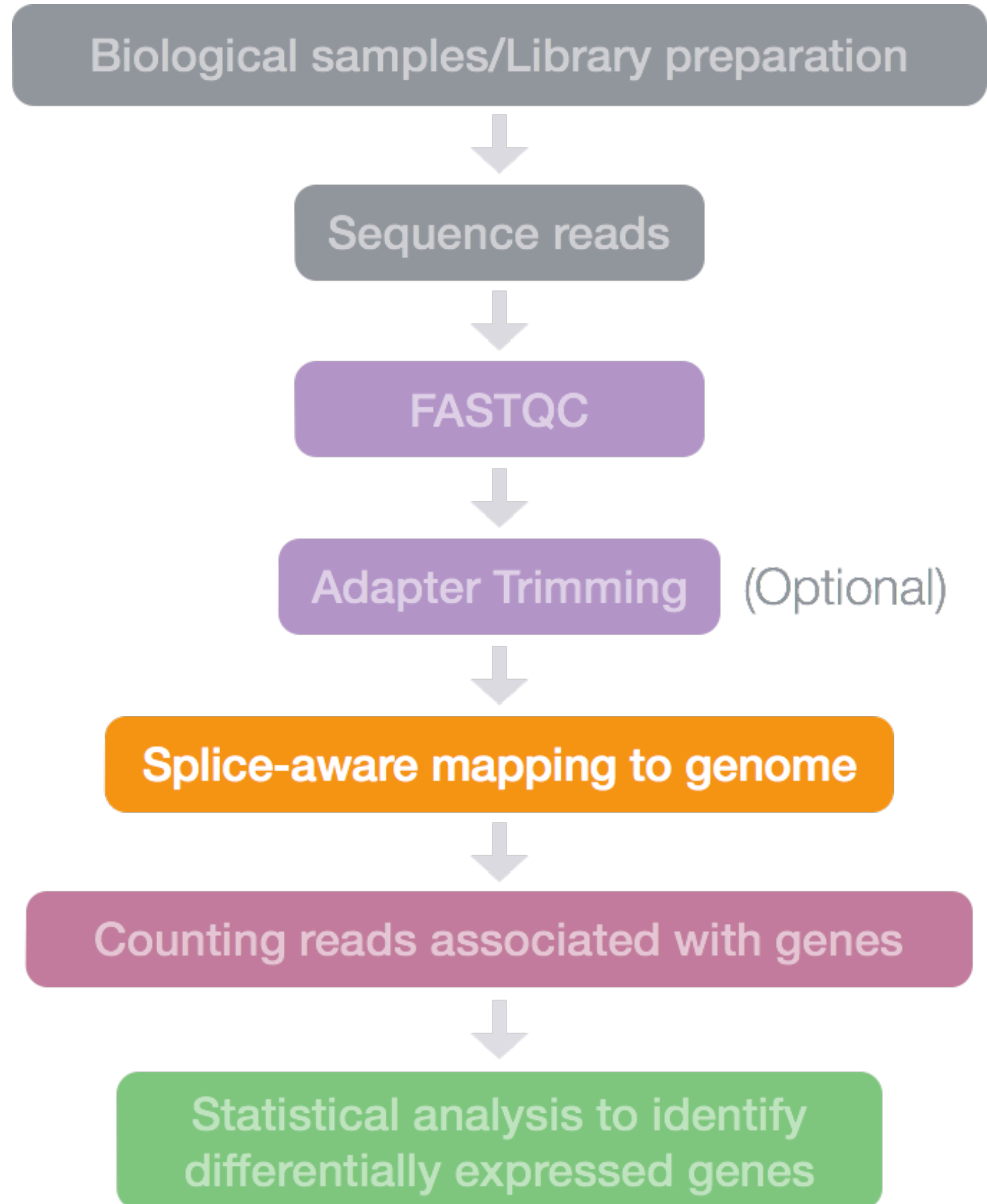
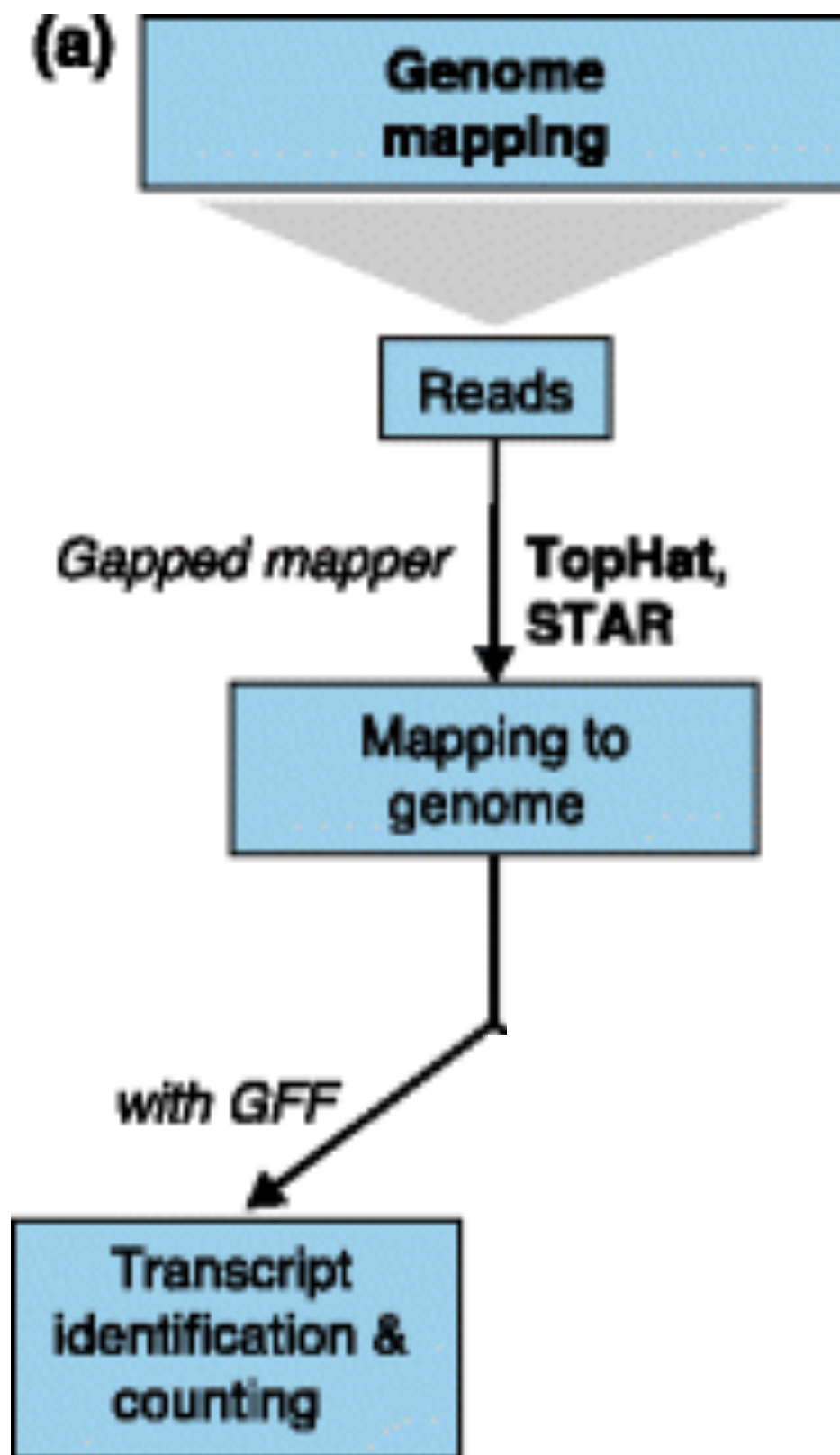
Aligning reads: tools and theory



Strategies for read mapping with RNA-seq



Strategies for read mapping with RNA-seq



Transcriptome quantification

Genome

chrX: 52139280 152139290 152139300 152139310 152139320 152139330
--->CGCCGTCCCTCAGAAATGGAAACCTCGCTTCTCTCTGCCCCACAATGCGCAAGTCAG

Sequence read

CGTCCCTCAGAAATGGAAACCTCGCTT

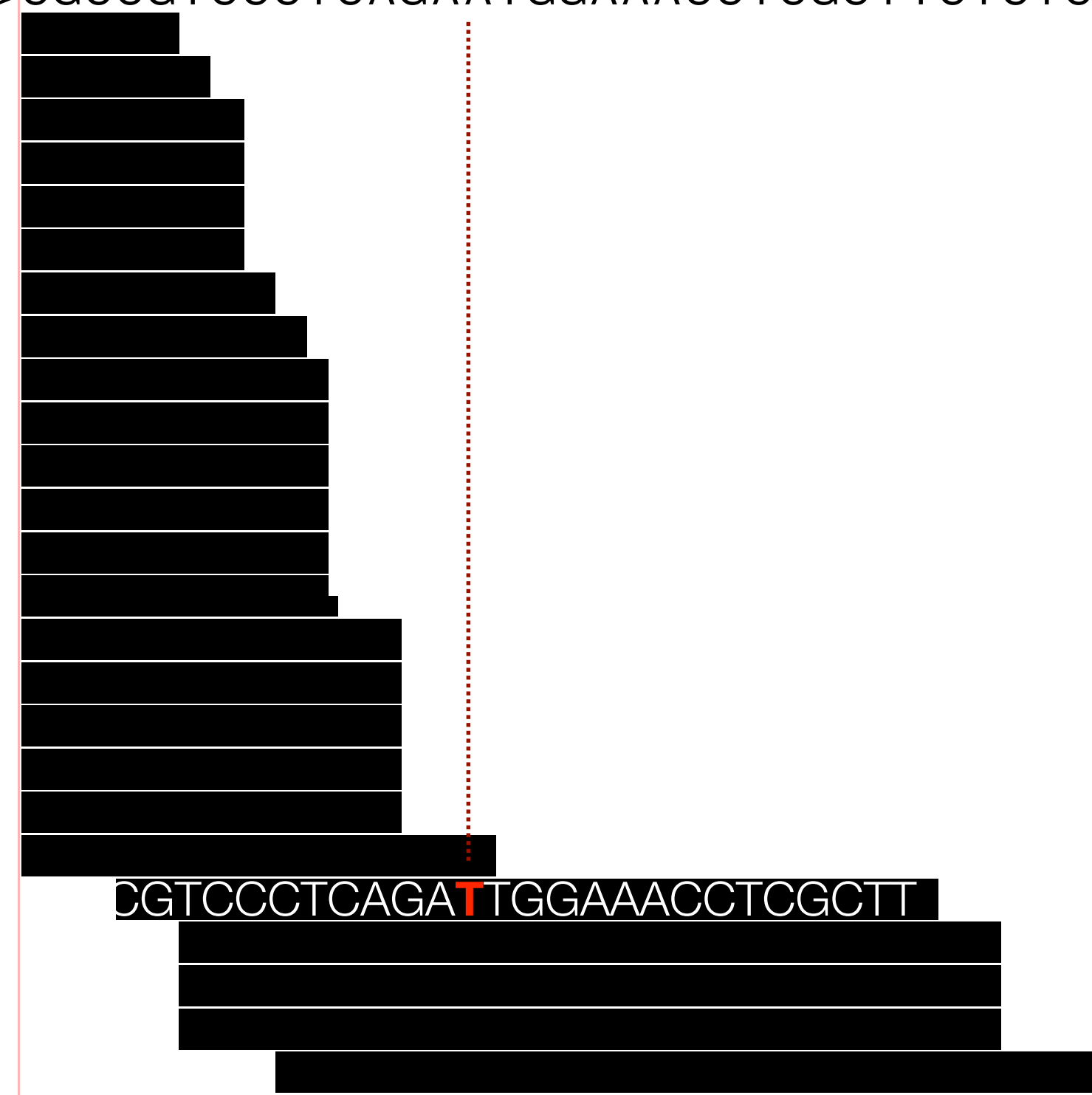


A simple case of string matching

Genome

chrX: 52139280 152139290 152139300 152139310 152139320 152139330
--->CGCCGTCCCTCAGAAATGGAAACCTCGCTTCTCTCTGCCCCACAATGCGCAAGTCAG

Sequence reads



Difficult in practice

- Volume of data: ~3 Gbp
- ~50% of genome is repeat regions that cannot be covered by reads
 - Simple repeats, tandem, interspersed
 - Transposons
 - Segmental duplications where mapping is unclear
- Gap or unfinished regions
 - peri-centromere, sub-telomere
 - ~5Mb unique to ethnic groups (e.g., African, Asian)
- Finishing errors(1/10,000bp), miscalled base incorporated

Challenges:

Human genome is large and complex ⁷

- Short reads: 50-150 bp (versus a very long reference)
 - Non-unique alignment
 - Sensitive to sequencing errors
- Massive amount of short reads: one lane produces ≥ 150 million 100 nucleotide reads
- Small insert size: 200-500 bp libraries

Challenges: short read NGS data

Reference ATCTCCATAGGACTAGGAAGTAG
Substitution ATCTCCATAGCACTAGGAAGTAG
Deletion ATCTCCATAGGAC-AGAAGTAG
Insertion ATCTCCATAGGACTAGGAAGTAG
3bp deletion ATCTC--AGGACTAGGAAGTAG

Challenges: non-exact matching

Local alignment vs Global alignment

- ▶ **Local alignment** matches the query with a *substring* (k-mer) of the reference
 - ▶ Tailored towards finding *regions of highly similar sequence* and aligning around those by working outwards to align the rest

Local Alignment

```
5' ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA 3'
      |||| | |||| | |||| |||| |||| |||| |||| |||| ||||
5' TACTCACGGATGAGGTACTTTAGAGGC 3'
```

Global Alignment

```
5' ACTACTAGATTACTTACGGATCAGGTACTTTAGAGGCTTGCAACCA 3'
||| |||| |||| |||| |||| |||| |||| |||| |||| ||||
5' ACTACTAGATT----ACGGATC--GTACTTTAGAGGCTAGCAACCA 3'
```

- ▶ A **global alignment** performs end-to-end alignment between the query and the reference

Reference ATCTCCATAGGACTAGGAAGTAG
Substitution ATCTCCATAG**C**ACTAGGAAGTAG
Deletion ATCTCCATAGGAC-AGAAGTAG
Insertion ATCTCCATAGGACTAGGAAGT**T**AG
3bp deletion ATCTC--AGGACTAGGAAGTAG

General concepts: edit distance

Reference CGTCCCTCAGATTGGAA—CCTCGCTT
Read TCCCTCAGAATGGAAACCTCGCT

Edit distance =3

General concepts: edit distance

Building an index

- ▶ For each read we need to scan the entire corpus as fast as possible
- ▶ Having an index of the reference genome provides an efficient way to search
- ▶ Once index is built, it can be queried any number of times
- ▶ Indexes are genome and tool-specific



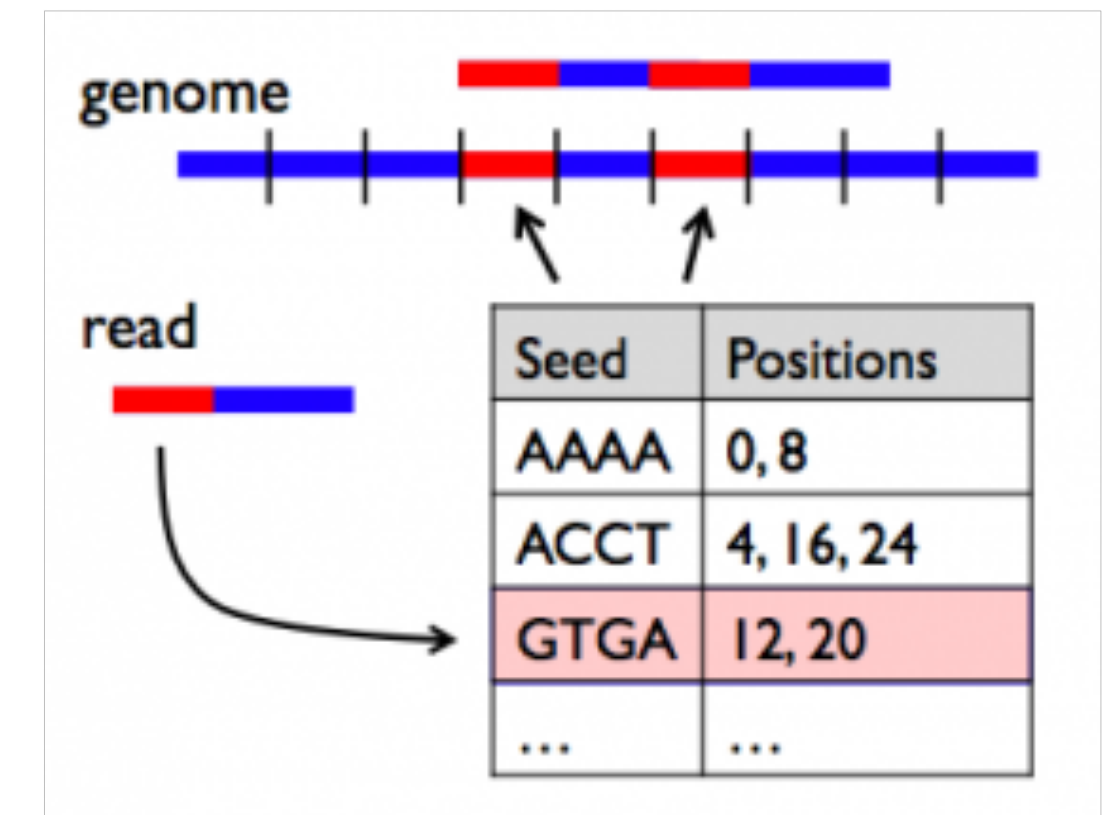
Alignment tools can be grouped based on indexing method

- ▶ Some examples include:
 - ▶ Hash-based
 - ▶ Suffix arrays
 - ▶ Burrows-Wheeler Transform

Hash-based alignment (circa 1990)



- ▶ Pick k-mer size, build lookup of every k-mer in the reference mapped to its positions (the index)
- ▶ Break the query into k-mers
- ▶ Seed-and-extend strategy
- ▶ For BLAST, 100% match the query k-mer to reference then extend until score drops below 50%
- ▶ 0.1 - 1 sec per query; not feasible for NGS data



Hash-based alignment (present day)

- ▶ Need to make some concessions on sensitivity by making adaptations for use on NGS data:
 - ▶ allow for mismatches and/or gaps (ELAND, MAQ, SOAP)
 - ▶ using multiple seeds (BLAT, ELAND2)
- ▶ Memory intensive and slower (~16GB RAM required for hg19)
- ▶ Simpler in design but more sensitive

Suffix arrays

- ▶ A sorted table of all suffixes (substrings) of a given string
- ▶ A suffix array will contain integers that represent the starting indexes of the all the suffixes of a given string, after the aforementioned suffixes are sorted
- ▶ Requires large amount of memory to load the suffix array and genome sequence prior to alignment
- ▶ Popular Tools:
STAR (2012)

Let the given string be “mississippi”

Suffixes	ID	Sorted Suffixes	Suffix Array
mississippi\$	1	\$	12
issippi\$	2	i\$	11
ssissippi\$	3	ippi\$	8
sissippi\$	4	issippi\$	5
issippi\$	5	issippi\$	2
ssippi\$	6	mississippi\$	1
sippi\$	7	pi\$	10
ippi\$	8	ppi\$	9
ppi\$	9	sippi\$	7
pi\$	10	sissippi\$	4
i\$	11	ssippi\$	6
\$	12	ssissippi\$	3

The suffix array will be:
{12, 11, 8, 5, 2, 1, 10, 9, 7, 4, 6, 3}

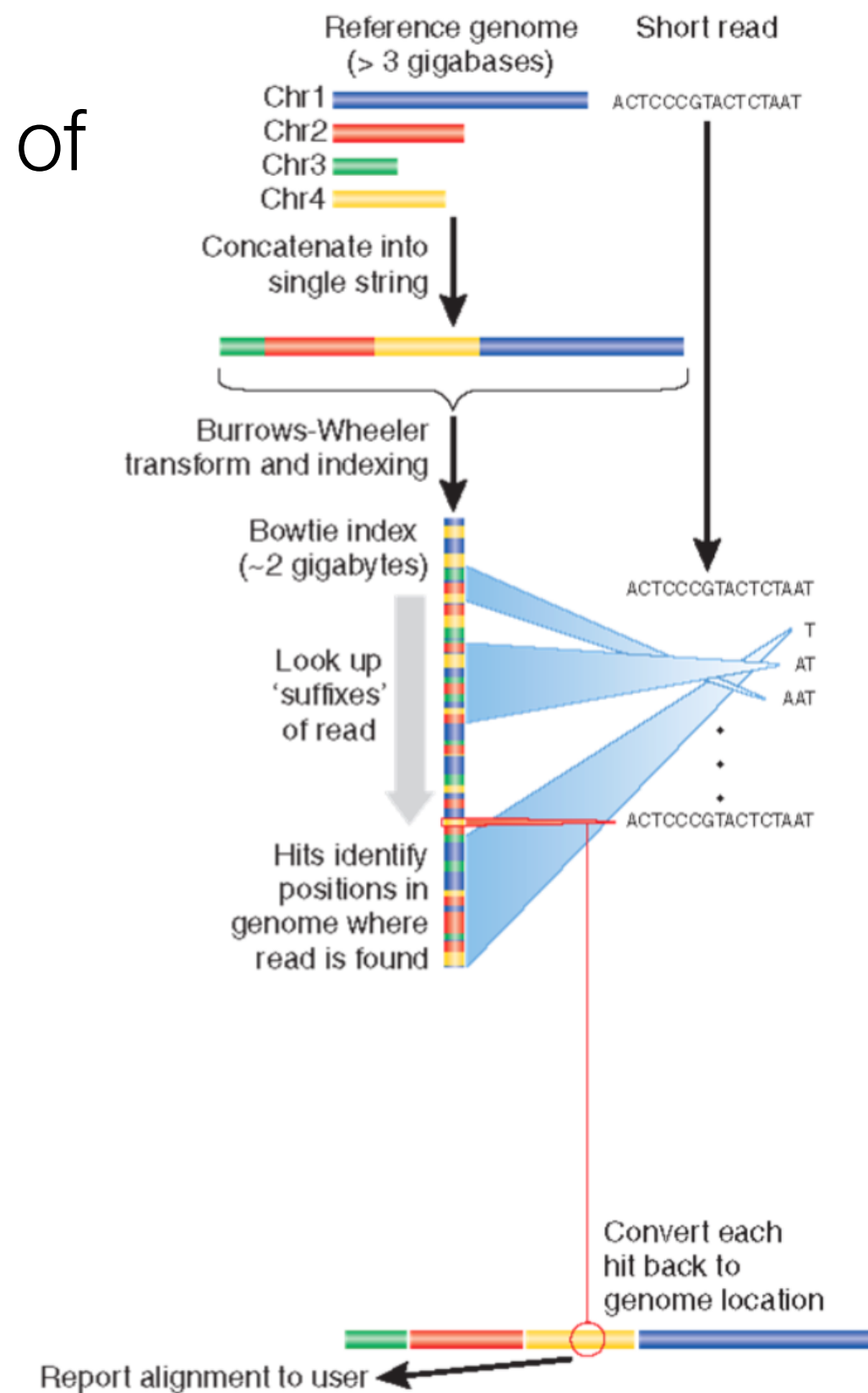
Burrows-Wheeler transform

- ▶ A compressed form of suffix arrays
- ▶ Tends to put runs of the same character together rather than alphabetically, which makes the compression work well

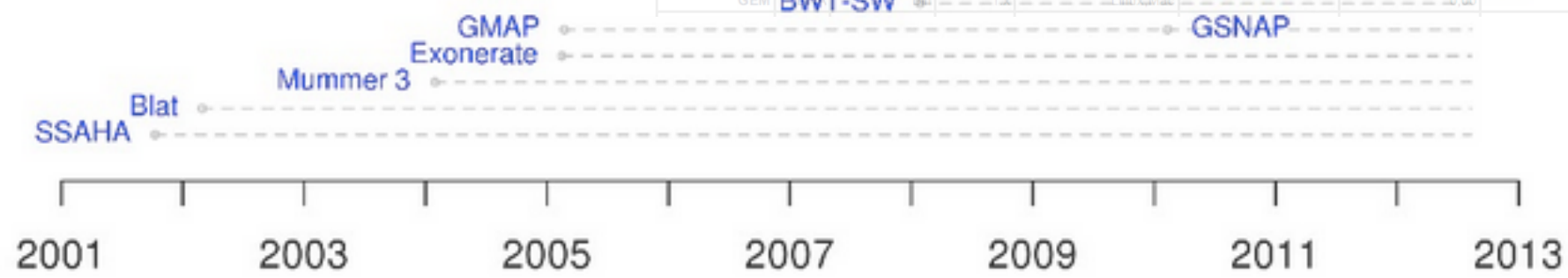
Suffixes	ID	Sorted Suffixes	Suffix Array	Sorted Rotations (A_s matrix)	BWT Output (L)
mississippi\$	1	\$	12	\$mississippi	i
ississippi\$	2	i\$	11	i\$mississipp	p
ssissippi\$	3	ippi\$	8	ippi\$mississ	s
sissippi\$	4	issippi\$	5	issippi\$miss	s
issippi\$	5	ississippi\$	2	ississippi\$m	m
ssippi\$	6	mississippi\$	1	mississippi\$	\$
sippi\$	7	pi\$	10	pi\$mississip	p
ippi\$	8	ppi\$	9	ppi\$mississi	i
ppi\$	9	sippi\$	7	sippi\$missis	s
pi\$	10	sissippi\$	4	sissippi\$mis	s
i\$	11	ssippi\$	6	ssippi\$missi	i
\$	12	ssissippi\$	3	ssissippi\$mi	i

Burrows-Wheeler transform

- ▶ Much less memory because of compression; ~1.5 GB of RAM required for hg19 index
- ▶ But compression results in diminished efficiency of the string search operations
- ▶ Popular Tools:
 - Bowtie2 (2012)
 - SOAP2
 - BWA-MEM (2013)

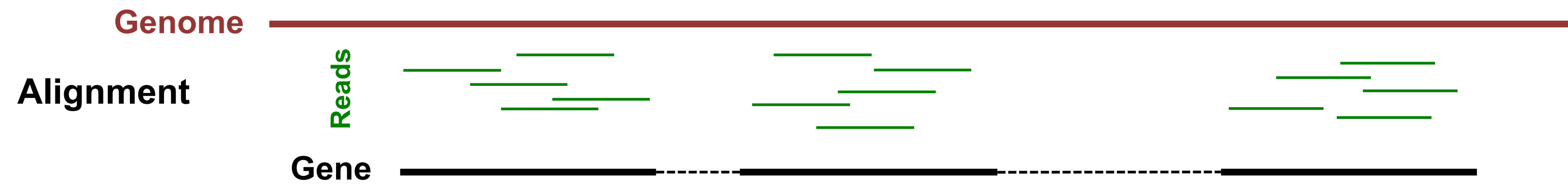
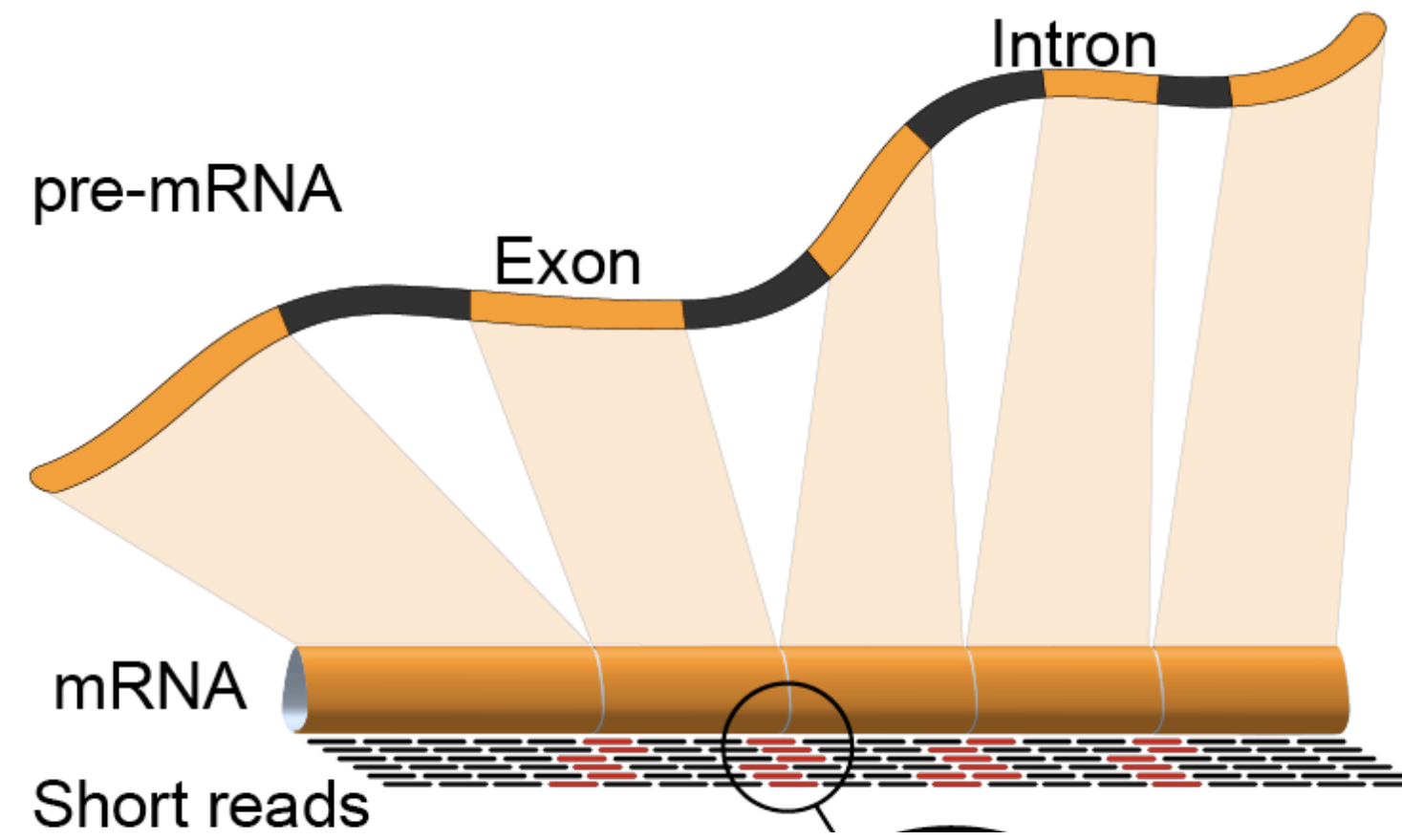


Mapper	Data	Availability	Version	O.S.	Number Citations	Citations/Years	Seq.Plat.	Input	Output	Min. RL	Max. RL	Mismatches	Indels	Gaps	Align. Reported	Alignment	Parallel	QA	PE	Splicing	
Bowtie	DNA	OS	0.12.7	Linux,Mac,Windows	1168	335.04	1,So,4,Sa,P	(C)FAST(A/Q)	SAM TSV	4	1K	Score	Score	N	A,B,R,S	G L	SM	Y	Y	N	
Blat	DNA	OS	34	Linux,Mac	2844	268.37	N	FASTA	TSV BLAST	11	5000K	Score	Score	Y	B	L	N	N	N	De novo	
MAQ	DNA	OS	0.7.1	Linux,Mac	957	237.27	1,So	(C)FAST(A/Q)	TSV	8	63	Y	Y	N			N	Y	Y	N	
BWA	DNA	OS	0.6.2	Linux,Mac,Windows	738	225.15	1,So,4,Sa,P	FASTA/Q	SAM	4	200	Y	8	Y	R,S	G	SM	Y	Y	N	
TopHat	RNA	OS	1.4.1	Linux,Mac	389	112.66	1	FASTA/Q, GFF	BAM	-	-	2	0	N	B,S	-	SM	Y	Y	De novo	
SOAP	DNA	OS	1.11	Linux,Mac	451	98.04	1	FASTA/Q	TSV	7	60	5	3	N	B,R,S		SM	N	Y	N	
SOAPdenovo	DNA	OS	2.04	Linux,Mac	284	98.04	1	FASTA/Q	SAM TSV	27	1K	2	0	Y	A,B,R	L	SM	N	Y	N	
Mummer 3	DNA	OS	3.23	Linux,Mac	653	78.93	N	FASTA	TSV	10	*	Y	Y	Y	A,B	G	N	N	N	N	
BWA-SW	DNA	OS	0.6.2	Linux,Mac,Windows	189	61.41	1,So,4,Sa,P	FASTA/Q	SAM	4	1000K	0.1	0.1	Y	R,S	L	SM	Y	N	N	
mrFAST	miRNA	OS	2.1.0.4	Linux	159	59.56	1	FASTA/Q	SAM	25	300	Score	6	N	A,B	G	N	N	Y	N	
SHRIMP	DNA	OS	1.3.2	Linux,Mac	647	47.46	1,So,4,Hel	(C)FAST(A/Q)	TSV	14	1K	Score	Score	Y	B,S	G	SM	N	Y	N	
SSAHA	DNA	OS	3.1	Linux,Mac	43	43.24	N	FASTA/Q	TSV	15	*	Y	Y	Y	B,S	G L	N	N	N	N	
CloudBurst	DNA	OS	1.1	Linux,Mac,Windows	13	43.08	1	FASTA	TSV		1K	Y	Y	Y	A,B	G	Cloud	N	N	N	
RMAP	DNA	OS	2.05	Linux,Mac	13	35.89	1,So,4	(C)FAST(A/Q)	BED	11	10K	Y	0	N	B,S		N	Y	Y	N	
SeqMap	DNA	OS	1.013	Linux,Mac	12	35.04	1	FASTA	ELAND	15	500	5	3	N	A		SM	N	N	N	
BFAST	DNA	OS	0.7.0	Linux,Mac	38	33.74	1,So,4,Hel	(C)FAST(A/Q)	SAM TSV		*	Y	Y	Y	B,R,U	G	SM	N	Y	N	
Exonerate	DNA	OS	2.2	Linux,Mac,Windows	259	33.55	N	FASTA	TSV	20	*	Score	Score	Y	B,S	G L	N	N	N	De novo	
GMAP	DNA	OS	2012-04-27	Linux,Unix,Mac,Windows	21	33.06	1,4,So,Hel,Ion,P	FASTA/Q	SAM, GFF	8	*	Y	Y	Y	B	G L	SM	N	N	De novo	
GSNAP	DNA	OS	2012-04-27	Linux,Unix,Mac,Windows	17	32.92	1,4,So,Hel,Ion,P	FASTA/Q	SAM	8	250	Y	Y	Y	A,B,U,S	G L	SM	N	Y	Lib and de novo	
ZOOM	DNA	Com	1.5	Linux,Mac	12	28.34	1,So,4	(C)FAST(A/Q)	SAM BED GFF	12	240	Y	Y	N	B,U,S	G	SM/DM	Y	Y	N	
SpliceMap	RNA	OS	3.3.5.2	Linux,Mac	61	29.43	1	FASTA/Q	SAM BED	-	-	0.1		Y	A	-	SM	N	Y	Lib and/or de novo	
MapSplice	RNA	OS	1.15.2	Linux,Mac	50	25.25	1	FASTA/Q	SAM BED	-	-	3		Y	B	-	SM	N	Y	De novo	
QPALMA	RNA	OS	0.9.2	Linux,Mac	75	19.35	1,4	Specific	TSV	-	-	Y	Y	Y	B	L	N	Y	N	Lib and de novo	
RazerS	RNA	OS	1.1	Linux,Mac	58	19.19	1	FASTQ	TSV ELAND	11	*	Score	Score	Y	A,B,S	G	N	N	N	N	
mrFAST	miRNA	OS	2.3.0	Linux,Mac	52	18.75	1,So	FASTA/Q	SAM	25	200	Y	0	N	A	G	N	N	Y	N	
Stampy	DNA	Bin	1.0.16	Linux,Mac	25	17.75	1	FASTA/Q	SAM TSV	4	4K	0.15	30	N	B,R,S	G	N	Y	Y	N	
PASS	DNA	Bin	1.0	Linux,Mac	14	17.22	1,So,4	(C)FAST(A/Q)	SAM GFF3 BLAST	23	1K	Y	Y	Y	A,B	G	SM	Y	Y	De novo	
SOCS	DNA	OS	1.0	Linux,Mac	42	16.67	So	(C)FAST(A/Q)	TSV		64	Y	0	N	A,B		SM	Y	N	N	
GenomeMapper	DNA	OS	1.0	Linux,Mac	12	16.57	1	FASTA/Q	BED TSV	12	2K	10	10	Y	A,B,R	G	SM	N	N	N	
Slider	DNA	OS	0.6	Linux,Mac,Windows	32	16.17	1	FASTA/Q	TSV		62	3	0	N	B,S		N	Y	Y	N	
BSMAP	Bisulfite	OS	1.0	Linux,Mac	12	16.05	1	FASTA/Q	SAM TSV	8	144	15	0	N	B,U,S		SM	N	Y	N	
PerM	DNA	OS	1.0	Linux,Mac	7	9.87	1,So	(C)FAST(A/Q)	SAM TSV	20	128	9	0	Y	A,U	G	DM	Y	Y	N	
BWT-SW	DNA	OS	2007	Linux,Mac	25	9.78	N	FASTA	TSV		1K	Score	Score	Y	A		N	N	N	N	
SHRIMP 2	DNA	OS	1.0	Linux,Mac	12	9.74	1,So,4	FASTA/Q	SAM	30	1K	Y	Score	N	B,U,S	G	SM	Y	Y	N	
RNA-Mate	RNA	OS	1.0	Linux,Mac	25	9.71	So	CFASTA	BED Counts	-	-	Y	0	N	S	-	DM	Y	N	Lib	
Supersplat	RNA	OS	1.0	Linux,Mac	9	9.66	N	FASTA	TSV			0	0	Y	A,U	G	N	N	N	De novo	
PatMaN	miRNA	OS	1.0	Linux,Mac	14	9.54	N	FASTA	TSV	1	*	Y	Y	N	A	G	N	N	N	N	
BS Seeker	Bisulfite	OS	1.0	Linux,Mac	16	9.19	1	FASTA/Q	SAM	-	-	3	0	N	U	-	SM	Y	N	N	
Slider II	DNA	OS	1.0	Linux,Mac,Windows	18	9.13	1	FASTA/Q	TSV		93	Y		N	B,S		N	N	Y	N	
GNUMAP	DNA	OS	1.0	Linux,Mac	15	9.10	1	FASTA/Q Illumina	SAM TSV	16	1K	Score	Score	Y	B	G	SM/DM	Y	N	N	
MOM	DNA	Bin	1.0	Linux,Mac,Windows	16	9.10	1,4	FASTA	TSV			Y	0	N	A	L	SM	N	Y	N	
Bismark	Bisulfite	OS	1.0	Linux,Mac	27	9.10	1	FASTA/Q	SAM	16	10K	Score	Score	N	U	-	SM	Y	Y	N	
BRAT	Bisulfite	OS	1.0	Linux,Mac	10	8.74	1	FASTA/Q	TSV			Y	0	N			N	N	Y	N	
SOAPSplice	RNA	OS	1.0	Linux,Mac	13	8.70	1,4	FASTA/Q	TSV	13	3K	5	2	Y	U	-	SM	Y	Y	De novo	
WHAM	DNA	OS	1.0	Linux,Mac	12	8.64	N	FASTQ	SAM	5	128	5	3	N	A,B,R,U,S	G	N	Y	Y	De novo	
MicroRazerS	miRNA	OS	1.0	Linux,Mac	4	8.49	N	FASTA	SAM TSV	10	*	Score	0	N	S	G	N	N	N	N	
RUM	RNA	OS	1.1	Linux,Mac	2	1.88	1,4	FASTA/Q	SAM TSV BED	-	-	Y	Y	Y	B	-	SM	N	Y	De novo	
ProbeMatch	DNA	OS	1.0	Linux,Mac	6	1.77	1,4,So	FASTA	ELAND	38	50	3	Y	N	A,B		N	N	N	N	
X-Mate	DNA	OS	1.0	Linux,Mac	1	0.92	1,So,4	(C)FAST(A/Q)	SAM BED Counts	-	-	Y	0	N	S	-	DM	Y	N	Lib	
SSAHA2	DNA	OS	2.3.5	Linux,Mac	1	0.90	1,4,So	FASTA/Q	SAM	15	48K	Score	Score	N	B,S	L	N	N	Y	N	
Novoalign	DNA	OS	2.8.5	Linux,Mac	1	0.90	1,So,4,Ion,P	(C)FAST(A/Q)	Illumina SAM TSV	30	300	8	2	N	A,B,R,U,S	G SM/DM/Cloud	Y	Y	Lib		
VMATCH	DNA	OS	1.0	Linux,Mac	1	0.90	N	FASTA	TSV			Score	Score	Y	A,B,S	G L	N	N	N	N	
ELAND	DNA	OS	1.0	Linux,Mac	1	0.90	1	FASTA	TSV			32	2	0	N	B		N	N	N	N
GEM	DNA	OS	1.0	Linux,Mac	1	0.90	1,So	FASTA/Q	SAM, Counts		*	Y	Y	Y	A,S	G	SM	Y	Y	Lib and de novo	

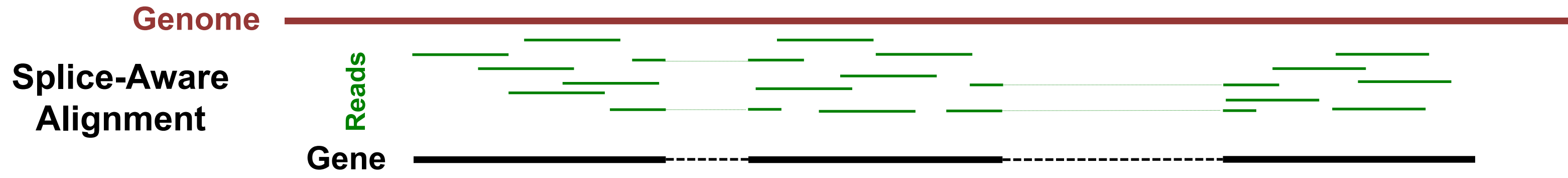


http://wwwdev.ebi.ac.uk/fg/hts_mappers/

Short-read aligners: choices



Versus



Splice-aware alignment

Splice-aware alignment tools:

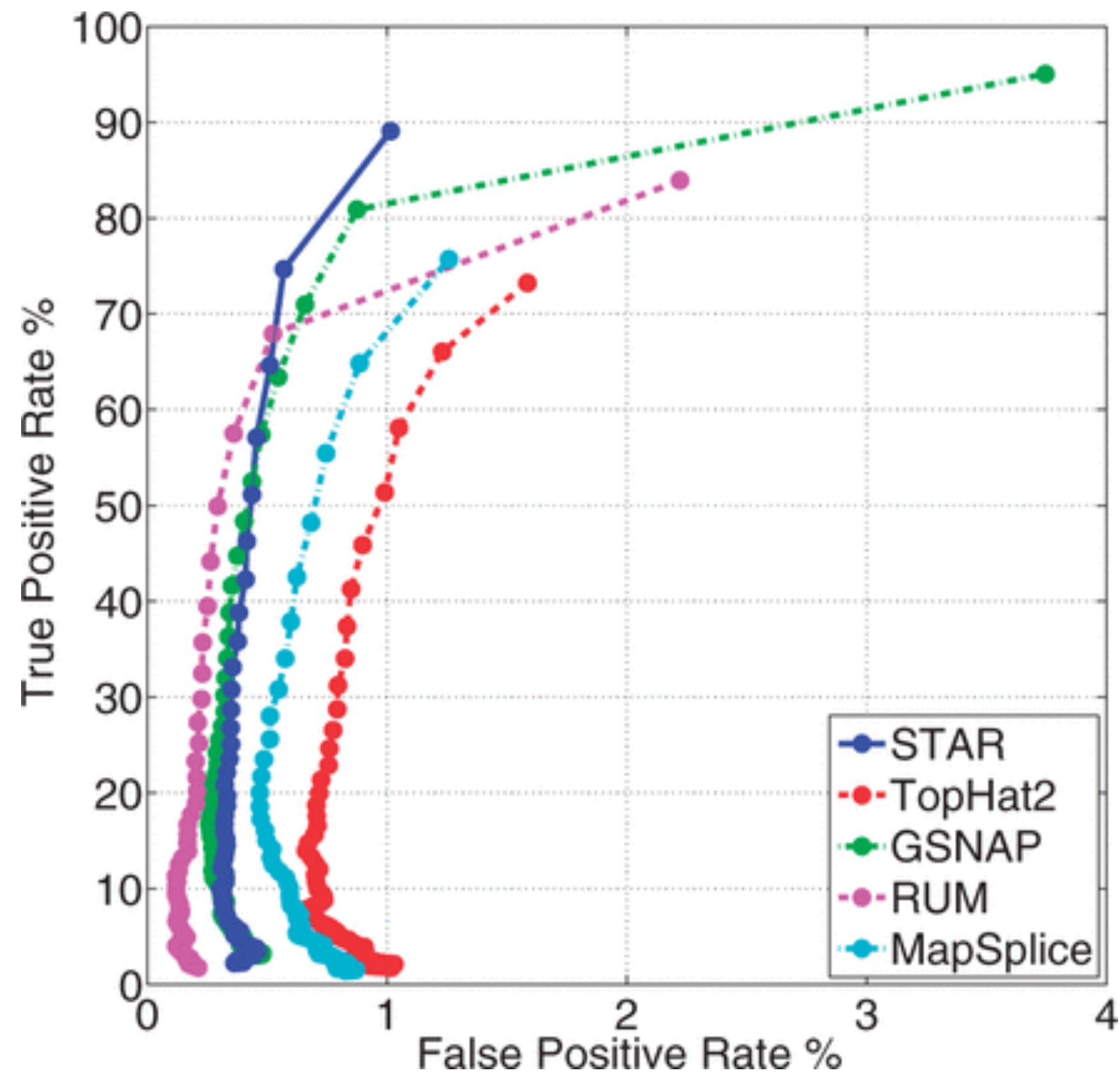
HISAT2, STAR, MapSplice, SOAPSplice, Passion, SpliceMap,
RUM, ABMapper, CRAC, GSNAP, HMMSplicer, Olego, BLAT

There are excellent aligners available that are not splice-aware. These are useful for aligning directly to genes.

However, you will lose isoform information.

Bowtie2, BWA, Novoalign (not free), SOAPaligner

Splice-aware alignment



Aligner	Mapping speed: million read pairs/hour		Peak physical RAM, GB	
	6 threads	12 threads	6 threads	12 threads
STAR	309.2	549.9	27.0	28.4
STAR sparse	227.6	423.1	15.6	16.0
TopHat2	8.0	10.1	4.1	11.3
RUM	5.1	7.6	26.9	53.8
MapSplice	3.0	3.1	3.3	3.3
GSNAP	1.8	2.8	25.9	27.0

Bioinformatics (2013) 29 (1): 15-21

The RNA-Seq specific tools

Alignment for RNA-Seq

- ▶ Use the strategy that is most relevant based on the quality of your genome and GTF
- ▶ Choose an aligner that can allow for a read to be “split” across distant regions to account for splice events
- ▶ Evaluate your computational resources and use an aligner that would work best within the confines of the available memory and CPU

