

# High-throughout sequencing and using short-read aligners

Simon Anders

# High-throughput sequencing (HTS)

Sequencing millions of short DNA fragments in parallel.

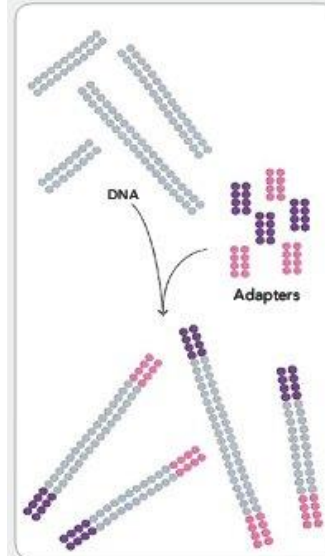
a.k.a.: next-generation sequencing (NGS)  
massively-parallel sequencing (MPS)

Market leader: Illumina (“HiSeq” instruments)

# Illumina's sequencing technology

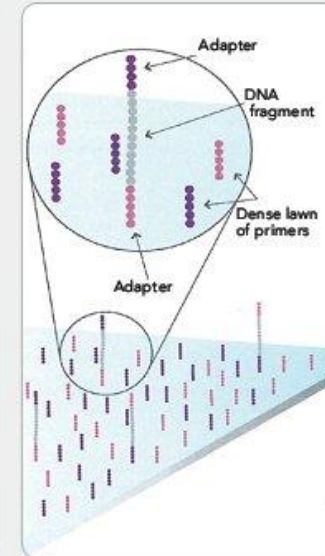
- Bridge PCR
- Sequencing-by-synthesis

1. PREPARE GENOMIC DNA SAMPLE



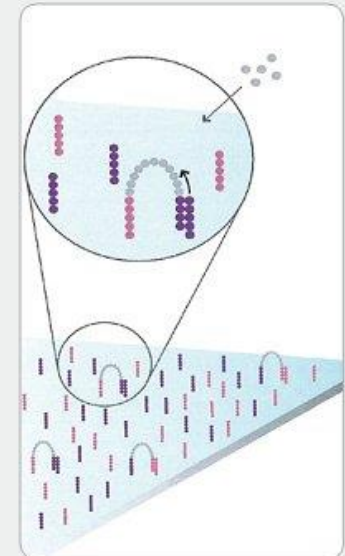
Randomly fragment genomic DNA and ligate adapters to both ends of the fragments.

2. ATTACH DNA TO SURFACE



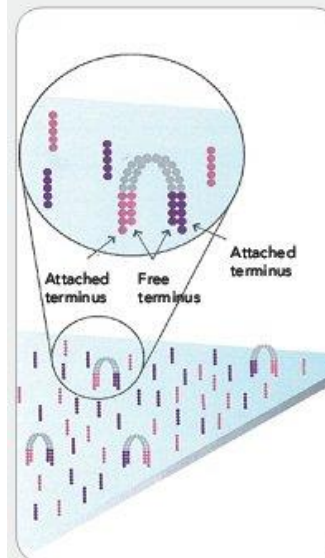
Bind single-stranded fragments randomly to the inside surface of the flow cell channels.

3. BRIDGE AMPLIFICATION



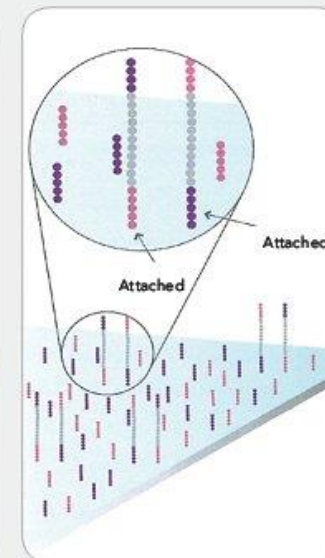
Add unlabeled nucleotides and enzyme to initiate solid-phase bridge amplification.

4. FRAGMENTS BECOME DOUBLE STRANDED



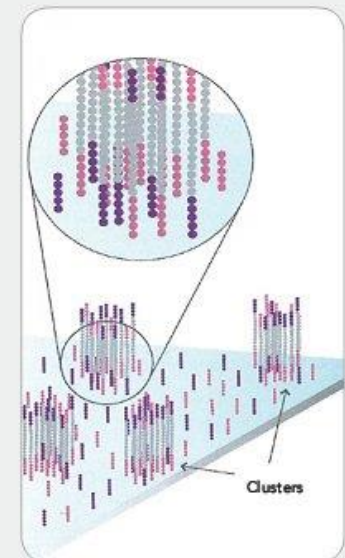
The enzyme incorporates nucleotides to build double-stranded bridges on the solid-phase substrate.

5. DENATURE THE DOUBLE-STRANDED MOLECULES



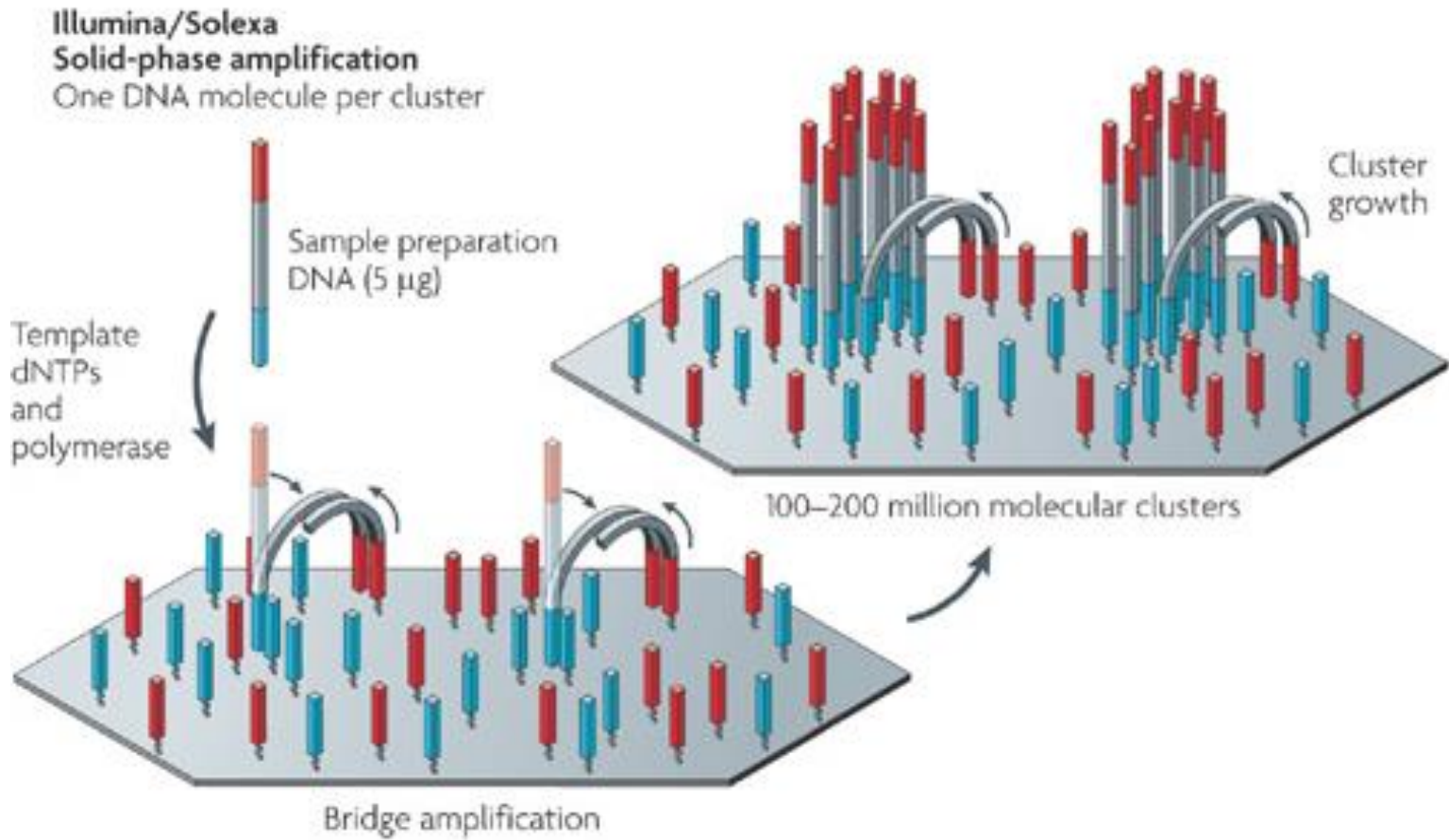
Denaturation leaves single-stranded templates anchored to the substrate.

6. COMPLETE AMPLIFICATION



Several million dense clusters of double-stranded DNA are generated in each channel of the flow cell.

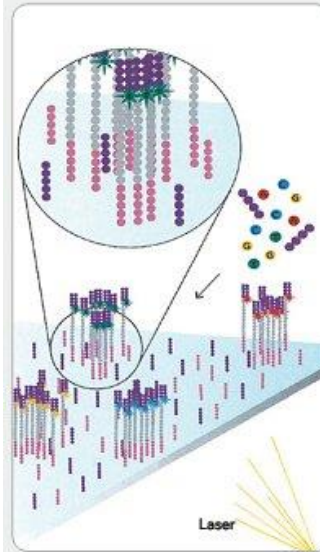
# Bridge PCR



# Illumina's sequencing technology

- Bridge PCR
- Sequencing-by-synthesis

7. DETERMINE FIRST BASE



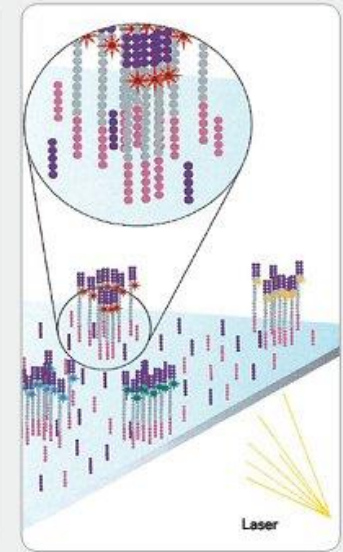
First chemistry cycle: to initiate the first sequencing cycle, add all four labeled reversible terminators, primers and DNA polymerase enzyme to the flow cell.

8. IMAGE FIRST BASE



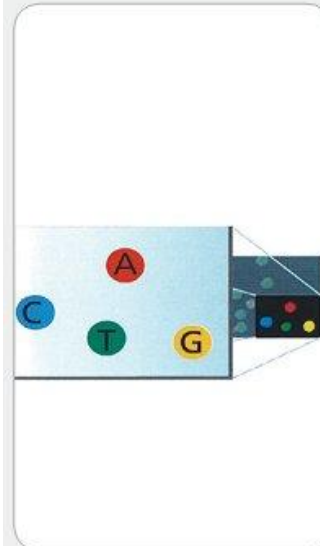
After laser excitation, capture the image of emitted fluorescence from each cluster on the flow cell. Record the identity of the first base for each cluster.

9. DETERMINE SECOND BASE



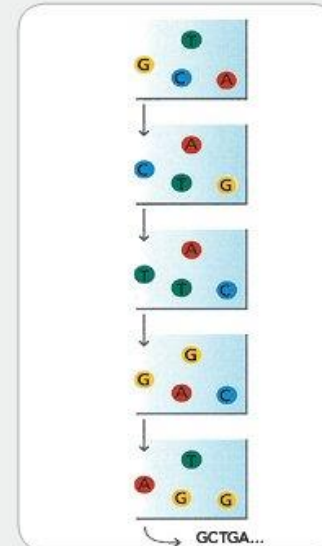
Second chemistry cycle: to initiate the next sequencing cycle, add all four labeled reversible terminators and enzyme to the flow cell.

10. IMAGE SECOND CHEMISTRY CYCLE



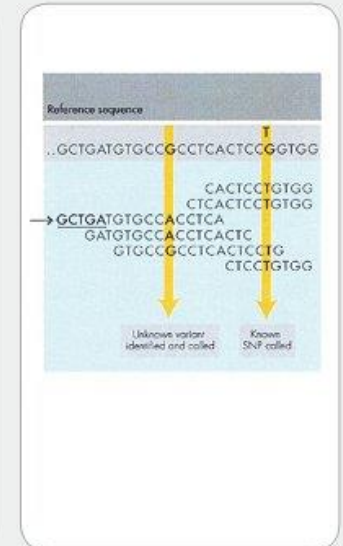
After laser excitation, collect the image data as before. Record the identity of the second base for each cluster.

11. SEQUENCE READS OVER MULTIPLE CHEMISTRY CYCLES



Repeat cycles of sequencing to determine the sequence of bases in a given fragment a single base at a time.

12. ALIGN DATA



Align data, compare to a reference, and identify sequence differences.

# Illunia HiSeq instrument

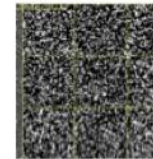
## Illumina HiSeq2000



- 2 flow cells (can be run independently)
- Up to 320Gb mapped sequence per FC
- 64Gb sequence per day (2 flow cells)



**c-Bot**  
(automated cluster generation)



Cluster density  
750-850/mm<sup>2</sup>



**HiSeq Flow Cells**

# Illumina HiSeq -- typical numbers

- 1 or 2 flow cells, each with 8 lanes
- up to ~ 200 M clusters per lane
- max 2x 150 bp read length
- 350 GB per flow cell and run
- 3 d per run (23 h in “rapid mode”)
  
- price per instrument: ~ €500k
- reagent costs per run: ~ €10k



# HTS: Application

- de-novo sequencing
- resequencing, variant calling  
(whole genome or targeted)
- RNA-Seq
- ChIP-Seq and other enrichment assays
- barcodes (tagged mutation collections)
- and many more ...



# Short reads and longer reads



Illumina:

maximum fragment length: 1 .. 2 kb

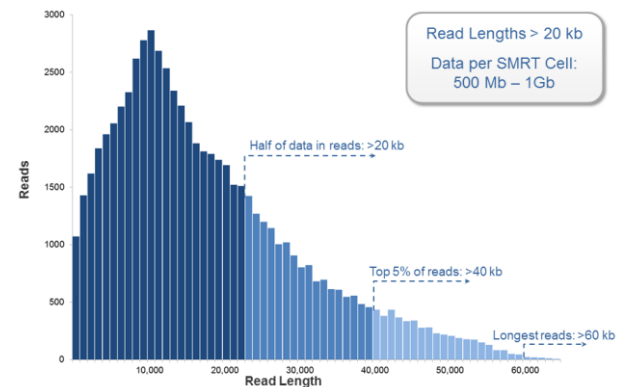
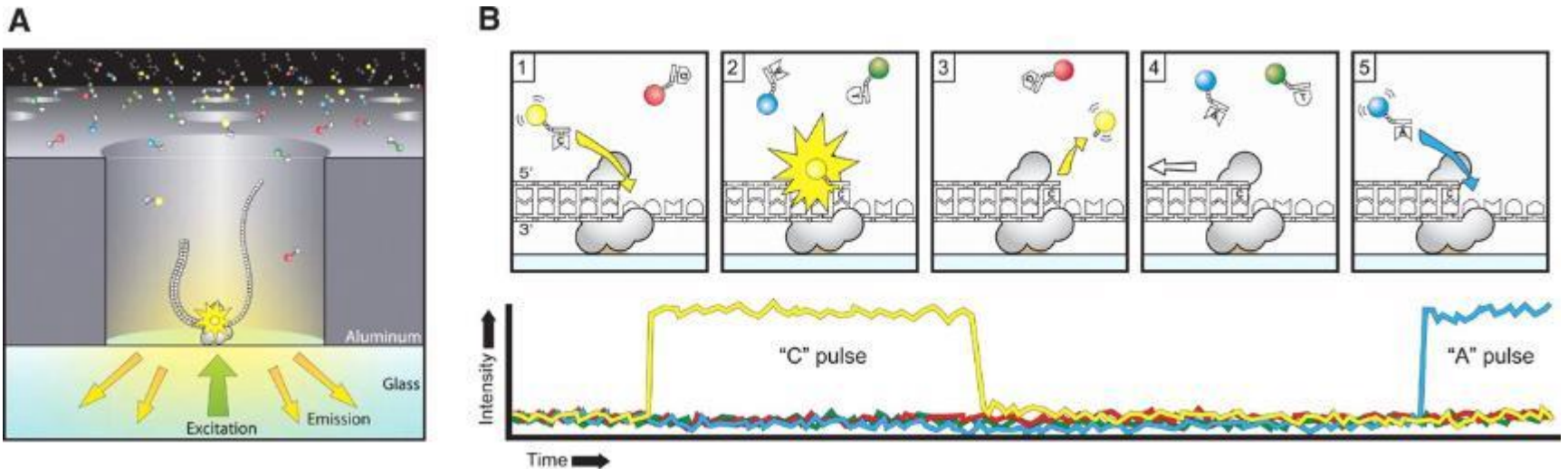
maximum read length: 150 .. 200 bp (x2)

Not really enough for

- calling structural variants
- assembling genomes
- phasing SNPs

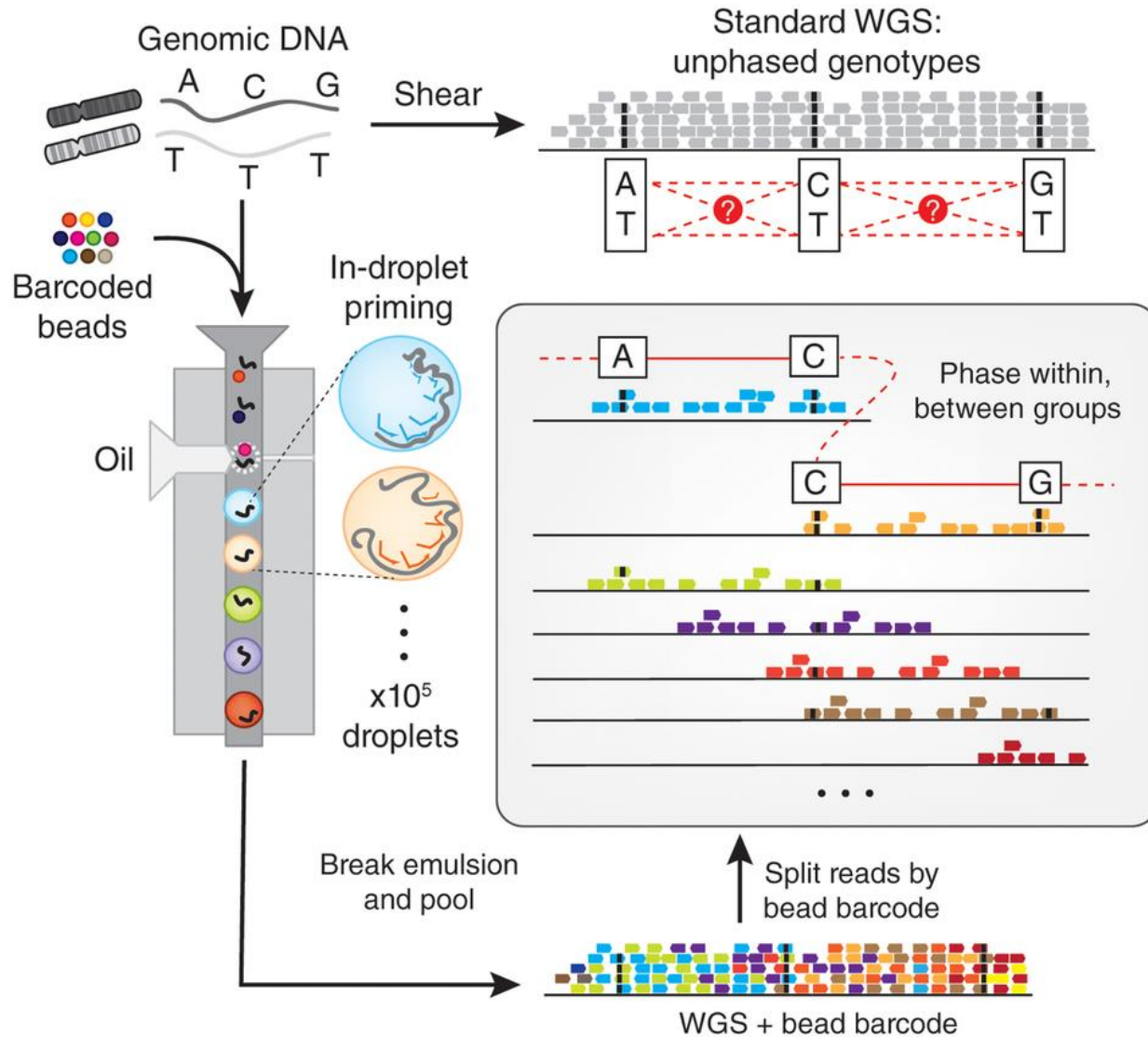
# Short reads and longer reads

## Pacific Biosystems: SMRT sequencing



# Droplet-based sequencing

10X Genomics



Next hot topic:

## Single-cell sequencing

Each read gets barcoded to indicate which cell it came from,  
by clever use of microfluidics.

Practical part:

Getting started with raw sequencing data







# FASTQ format

Each read is represented by four lines:

- '@', followed by read ID
- sequence
- '+', optionally followed by repeated read ID
- quality string:
  - same length as sequence
  - each character encodes the base-call quality of one base

# FASTQ quality scores

quality score $Q$	error prob. $p$	characters
0 .. 9	1 .. 0.13	!"#\$%&'()*
10 .. 19	0.1 .. 0.013	+,-./01234
20 .. 29	0.01 .. 0.0013	56789:;<=>
30 .. 39	0.001 .. 0.00013	?@ABCDEFGH
40	0.0001	I

$$Q = -10 \log_{10} p$$

# FASTQ files for paired end

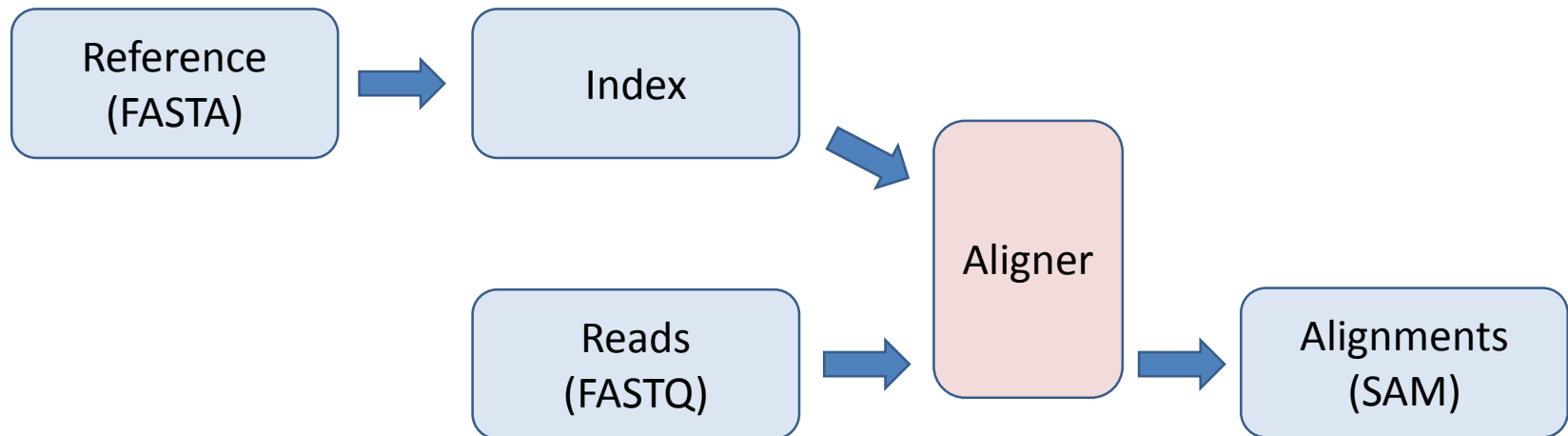
Simple: two FASTQ files,

- one with all first reads, one with all second reads
- Read names and ordering must match.

# Alignment

Given a reference genome, *where* did each read come from?

A *short-read aligner* finds for each read all possible matches in the reference genome.



# Popular aligners

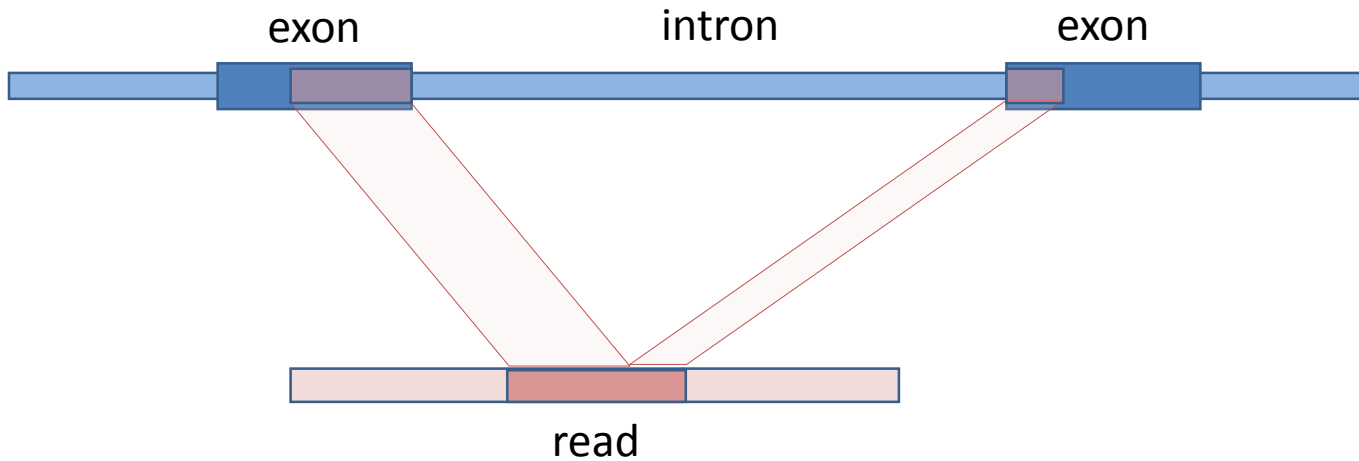
For genomic DNA

- Bowtie2
- BWA
- ...

For RNA

- Tophat2
- GSNAP
- STAR
- ...

# Spliced alignment



For RNA-Seq data, we need a **splice-aware** aligner.

# Criteria to choose an aligner

- Reputation
- Publication year
- Still developed and maintained?
  
- Speed and memory requirements
- Sensitivity
  
- Special functions  
(e.g., bisulfite mapping, fusion gene detection, multiple references, etc.)



# Example: Aligning RNA-Seq data with STAR

## STAR: ultrafast universal RNA-seq aligner

Alexander Dobin<sup>1,\*</sup>, Carrie A. Davis<sup>1</sup>, Felix Schlesinger<sup>1</sup>, Jorg Drenkow<sup>1</sup>, Chris Zaleski<sup>1</sup>, Sonali Jha<sup>1</sup>, Philippe Batut<sup>1</sup>, Mark Chaisson<sup>2</sup> and Thomas R. Gingeras<sup>1</sup>

<sup>1</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA and <sup>2</sup>Pacific Biosciences, Menlo Park, CA, USA

Associate Editor: Inanc Birol

### ABSTRACT

**Motivation:** Accurate alignment of high-throughput RNA-seq data is a challenging and yet unsolved problem because of the non-contiguous transcript structure, relatively short read lengths and constantly increasing throughput of the sequencing technologies. Currently available RNA-seq aligners suffer from high mapping error rates, low mapping speed, read length limitation and mapping biases.

**Results:** To align our large (>80 billion reads) ENCODE Transcriptome RNA-seq dataset, we developed the Spliced Transcripts Alignment to a Reference (STAR) software based on a previously undescribed RNA-seq alignment algorithm that uses sequential maximum mappable seed search in uncompressed suffix arrays followed by seed clustering and stitching procedure. STAR outperforms other aligners by a factor of >50 in mapping speed, aligning to the human genome 550 million 2 × 76 bp paired-end reads per hour on a modest 12-core server, while at the same time improving alignment sensitivity and precision. In addition to unbiased *de novo* detection of canonical junctions, STAR can discover non-canonical splices and chimeric (fusion) transcripts, and is also capable of mapping full-length RNA se-

unique challenges to detection and characterization of spliced transcripts. Two key tasks make these analyses computationally intensive. The first task is an accurate alignment of reads that contain mismatches, insertions and deletions caused by genomic variations and sequencing errors. The second task involves mapping sequences derived from non-contiguous genomic regions comprising spliced sequence modules that are joined together to form spliced RNAs. Although the first task is shared with DNA resequencing efforts, the second task is specific and crucial to the RNA-seq, as it provides the connectivity information needed to reconstruct the full extent of spliced RNA molecules. These alignment challenges are further compounded by the presence of multiple copies of identical or related genomic sequences that are themselves transcribed, making precise mapping difficult.

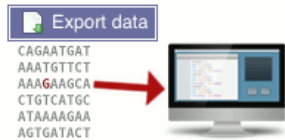
Various sequence alignment algorithms have been recently developed to tackle these challenges (Au *et al.*, 2010; De Bona, *et al.*, 2008; Grant *et al.*, 2011; Han *et al.*, 2011; Trapnell *et al.*, 2009; Wang *et al.*, 2010; Wu and Nacu, 2010; Zhang *et al.*, 2012). However, application of these algorithms invokes compromises

# Using STAR: Creating the index

We need

- a FASTA file with the human reference genome sequence
  - e.g., from Ensembl
- a GTF file with gene models
  - e.g., from Ensembl
- the STAR software:
  - from [github.com/alexdobin/STAR](https://github.com/alexdobin/STAR)

### Download a sequence or region

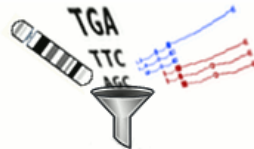


Click on the 'Export data' button in the lefthand menu of most pages to export:

- FASTA sequence
- GTF or GFF features

...and more!

### Customise your download



Custom datasets can be retrieved using the BioMart data-mining tool.

You may find exploring this web-based query tool easier than extracting information direct from our databases.

### Fetch data programmatically



Write your own Perl scripts to retrieve small-to-medium datasets. All our data, as well as added functionality, is available through the Ensembl Perl API.

Use the API to retrieve gene and transcript sets, fetch alignments between sequences, compare allele frequencies and much more!

You can also use our [REST API](#) to retrieve data to process in the programming language of your choice.

### Download **databases & software**

































All of our data and software, including pipelines and web code, is available free.

- [Download data via FTP](#)
- [Ensembl pipeline on GitHub](#)
- [Set up your own Ensembl website](#)

# Ensembl FTP download server

★	Species	DNA (FASTA)	cDNA (FASTA)	CDS (FASTA)	ncRNA (FASTA)	Protein sequence (FASTA)	Annotated sequence (EMBL)	Annotated sequence (GenBank)	Gene sets	Whole databases	Variation (GVF)	Variation (VCF)	Variation (VCF)
Y	<a href="#">Human</a> <i>Homo sapiens</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	<a href="#">GVF</a>	<a href="#">VCF</a>	<a href="#">VCF</a>
Y	<a href="#">Mouse</a> <i>Mus musculus</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	<a href="#">GVF</a>	<a href="#">VCF</a>	<a href="#">VCF</a>
Y	<a href="#">Zebrafish</a> <i>Danio rerio</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	<a href="#">GVF</a>	<a href="#">VCF</a>	<a href="#">VCF</a>
	<a href="#">Alpaca</a> <i>Vicugna pacos</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	-	-	<a href="#">VCF</a>
	<a href="#">Amazon molly</a> <i>Poecilia formosa</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	-	-	<a href="#">VCF</a>
	<a href="#">Anole lizard</a> <i>Anolis carolinensis</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	-	-	<a href="#">VCF</a>
	<a href="#">Armadillo</a> <i>Dasyus novemcinctus</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	-	-	<a href="#">VCF</a>
	<a href="#">Bushbaby</a> <i>Otolemur garnettii</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	-	-	<a href="#">VCF</a>
	<a href="#">C.intestinalis</a> <i>Ciona intestinalis</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	-	-	<a href="#">VCF</a>
	<a href="#">C.savignyi</a> <i>Ciona</i>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">FASTA</a>	<a href="#">EMBL</a>	<a href="#">GenBank</a>	<a href="#">GTF</a> <a href="#">GFF3</a>	<a href="#">MySQL</a>	-	-	<a href="#">VCF</a>

# Ensembl genome reference sequence

 Homo_sapiens.GRCh38.dna.chromosome.21.fa.gz	11.2 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.22.fa.gz	10.9 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.2.fa.gz	69.4 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.3.fa.gz	57.0 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.4.fa.gz	54.7 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.5.fa.gz	52.0 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.6.fa.gz	49.1 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.7.fa.gz	45.2 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.8.fa.gz	41.6 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.9.fa.gz	34.8 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.MT.fa.gz	5.4 kB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.X.fa.gz	44.0 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.chromosome.Y.fa.gz	6.8 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.nonchromosomal.fa.gz	2.9 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna.primary_assembly.fa.gz	3.10 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.alt.fa.gz	142 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.10.fa.gz	21.7 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.11.fa.gz	20.9 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.12.fa.gz	20.6 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.13.fa.gz	16.1 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.14.fa.gz	14.4 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.15.fa.gz	13.5 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.16.fa.gz	12.9 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.17.fa.gz	13.0 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.18.fa.gz	12.7 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.19.fa.gz	7.8 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.1.fa.gz	36.5 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.20.fa.gz	9.8 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.21.fa.gz	6.3 MB	25/02/2016, 04:27:00
 Homo_sapiens.GRCh38.dna_rm.chromosome.22.fa.gz	6.0 MB	25/02/2016, 04:27:00



# STAR: Generate genome index

```
Terminal
anderss@lx5-fimm8:~/work/STAR_demo$ ls download
Homo_sapiens.GRCh38.84.gtf Homo_sapiens.GRCh38.dna.primary_assembly.fa
anderss@lx5-fimm8:~/work/STAR_demo$ mkdir index_GRCh38.84
anderss@lx5-fimm8:~/work/STAR_demo$ STAR --runMode genomeGenerate --runThreadN
2 --genomeDir index_GRCh38.84 --genomeFastaFiles download/Homo_sapiens.GRCh38.d
na.primary_assembly.fa --sjdbGTFfile download/Homo_sapiens.GRCh38.84.gtf --sjdb
Overhang 99 --limitGenomeGenerateRAM 2500000000
Jul 06 17:03:26 ..... started STAR run
Jul 06 17:03:26 ... starting to generate Genome files
```

# STAR: Generate genome index

```
STAR \
  --runMode genomeGenerate \
  --runThreadN 2 \
  --genomeDir index_GRCh38.84 \
  --genomeFastaFiles download/Homo_sapiens.GRCh38.\
    dna.primary_assembly.fa \
  --sjdbGTFfile download/Homo_sapiens.GRCh38.84.gtf
```



# STAR: Align the reads

```
STAR \  
  --runThreadN 2 \  
  --genomeDir index_GRCh38.84 \  
  --readFilesIn sample1_1.fastq \  
                 sample1_2.fastq \  
  --outFileNamePrefix sample1_ \  
  --outSAMtype BAM SortedByCoordinate
```

**and maybe:**

```
--quantMode GeneCounts
```

# STAR: output

## Created files:

<u>size</u>	<u>file name</u>
<b>34M</b>	sample1_Aligned.sortedByCoord.out.bam
4.0K	sample1_Log.final.out
16K	sample1_Log.out
4.0K	sample1_Log.progress.out
20K	sample1_ReadsPerGene.out.tab
72K	sample1_SJ.out.tab

# Alignment output: SAM files

Part of a SAM file:

- Header
  - (lines start with @)
  - command line
  - chromosome names
  - etc
- Alignments
  - tab-separated table, at least 11 columns
  - one line per alignment

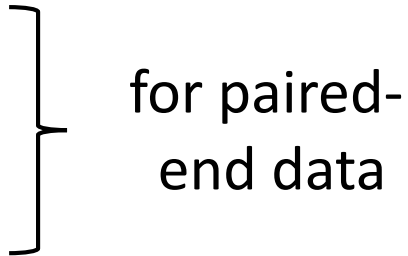
# SAM alignment section

[...]

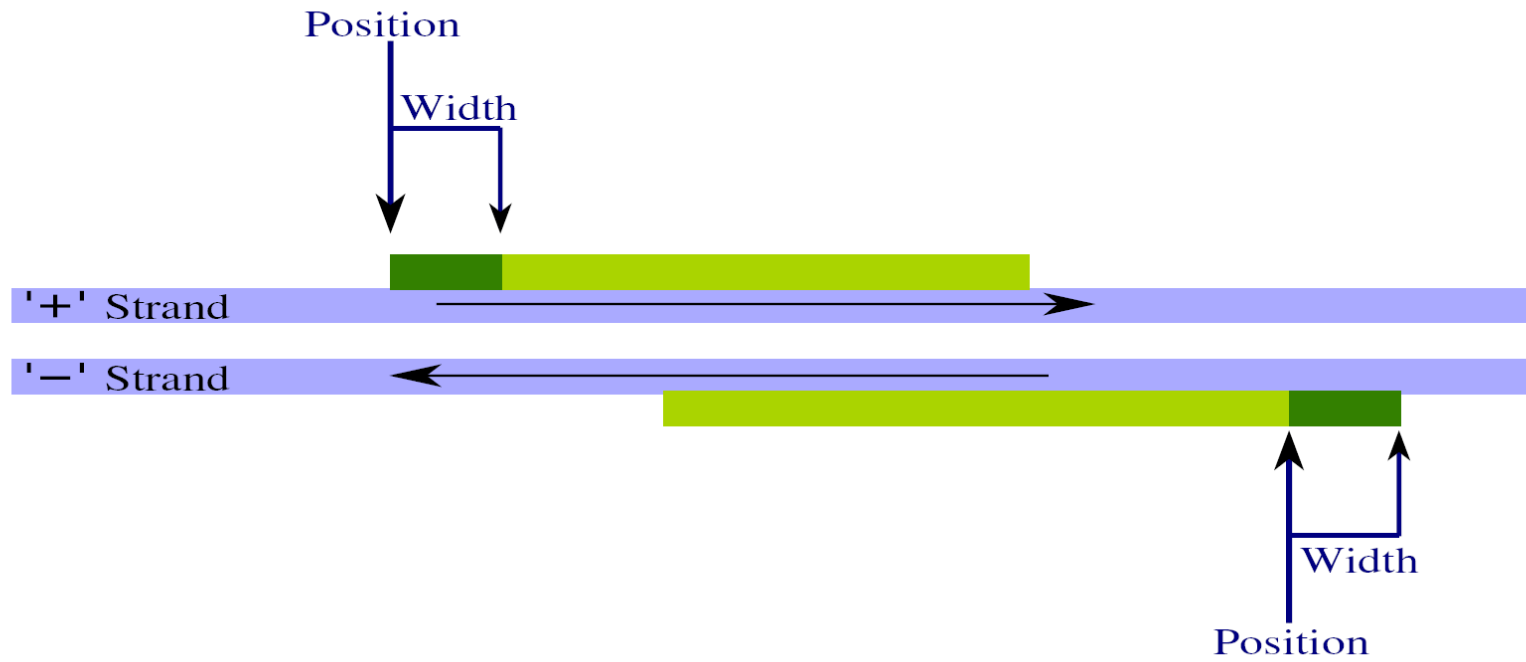
HWI-EAS225_309MTAAXX:5:1:689:1485	0	XIII	863564	25	36M	*
0	0	GAAATATATACGTTTTTATCTATGTTACGTTATATA				
CCCCCCCCCCCCCCCCCCCC4CCCB4CA?AAA<		NM:i:0	X0:i:1	MD:Z:36		
HWI-EAS225_309MTAAXX:5:1:689:1485	16	XIII	863766	25	36M	*
0	0	CTACAATTTTGCACATCAAAAAAGACCTCCAACACTAC				
=8A=AA784A9AA5AAAAAAAAAAAA=AAAAAAAAAA		NM:i:0	X0:i:1	MD:Z:36		
HWI-EAS225_309MTAAXX:5:1:394:1171	0	XII	525532	25	36M	*
0	0	GTTTACGGCGTTGCAAGAGGCCTACACGGGCTCATT				
CCCCCCCCCCCCCCCCCCCC?CCACCACA7?<???		NM:i:0	X0:i:1	MD:Z:36		
HWI-EAS225_309MTAAXX:5:1:394:1171	16	XII	525689	25	36M	*
0	0	GCTGTTATTTCTCCACAGTCTGGCAAAAAAAGAAA				
7AAAAAA?AA<AA?AAAAA5AAA<AAAAAAAAAAAA		NM:i:0	X0:i:1	MD:Z:36		
HWI-EAS225_309MTAAXX:5:1:393:671	0	XV	440012	25	36M	*
0	0	TTTGGTGATTTTCCCGTCTTTATAATCTCGGATAAA				
AAAAAAAAAAAAAAAA<AAAAAAAA<AAAA5<AAAA3		NM:i:0	X0:i:1	MD:Z:36		
HWI-EAS225_309MTAAXX:5:1:393:671	16	XV	440188	25	36M	*
0	0	TCATAGATTCCATATGAGTATAGTTACCCCATAGCC				
?9A?A?CC?<ACCCCCCCCCCCCCCCCCACCCCCC		NM:i:0	X0:i:1	MD:Z:36		

[...]

# Columns in a SAM file

- QNAME ID of the read (“query”)
  - FLAG alignment flags
  - RNAME ID of the reference (typically: chromosome name)
  - POS Position in reference (1-based, left side)
  - MAPQ Mapping quality (as Phred score)
  - CIGAR Alignment description (gaps etc.) in CIGAR format
  - RNEXT Mate reference sequence name
  - PNEXT Mate position
  - ISIZE inferred insert size
  - SEQ sequence of the read
  - QUAL quality string of the read
  - extra fields
- 
- for paired-end data

# Reads and fragments



# SAM format: FLAG field

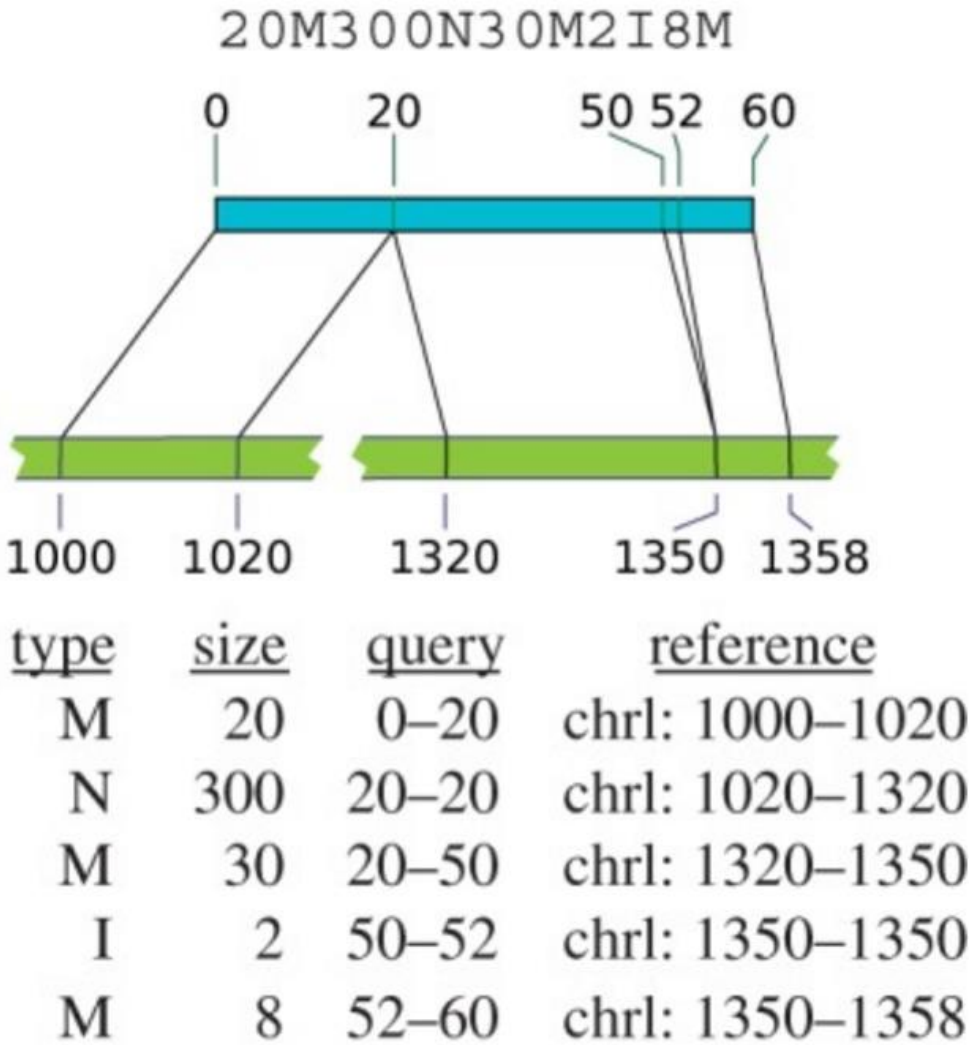
bit	hex	decimal	
0	00 01	1	read is a paired-end read
1	00 02	2	read pair is properly matched
2	00 04	4	read has not been mapped
3	00 08	8	mate has not been mapped
4	00 10	16	read has been mapped to "-" strand
5	00 20	32	mate has been mapped to "-" strand
6	00 40	64	read is the first read in a pair
7	00 80	128	read is the second read in a pair
8	01 00	256	alignment is secondary
9	02 00	512	read did had not passed quality check
10	04 00	1024	read is a PCR or optical duplicate

# SAM format: Optional fields

- last column
- triples of the format TAG : VTYPE : VALUE
- some important tag types:
  - NH: number of reported alignments
  - NM: number of mismatches
  - MD: positions of mismatches



# SAM format: CIGAR string



## SAM format: Paired-end and multiple alignments

- Each line is one alignment for one read.
- Multiple alignments for a read take several lines
  - The Read ID groups them.
- Paired-end alignments take two lines.
- All these lines are not necessarily in adjacent lines!

# SAM and BAM

- Text SAM files (".sam"): human readable
- Binary SAM files (".bam"): smaller, quicker to read

Use *samtools* for conversion.

# SAMtools

A small program to

- convert between .bam and .sam
- sort and merge SAM files
- index SAM and FASTA files for fast access
- calculate tallies (“flagstat”)
- ...

The SAMtools C API is a library for tool development.

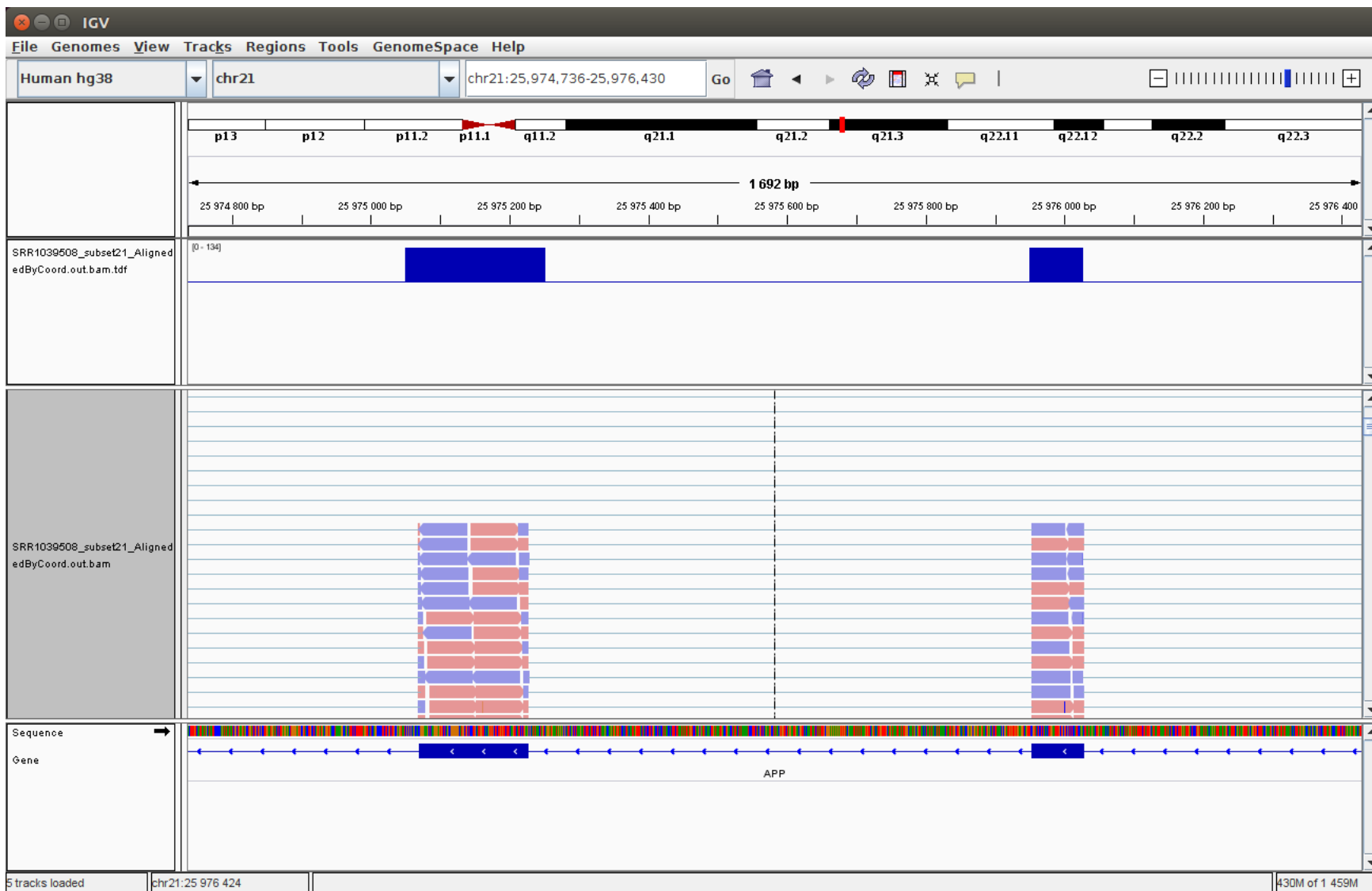
# Visually inspecting SAM files

*Look at all your big files!*

Learn to use text views *and* graphical tools.

For SAM, use e.g. the *Integrated Genome Viewer* (IGV).

# IGV



# Aligning multiple files

Sample table:

	<b>Library_ID</b>	<b>Cell_line</b>	<b>Treatment</b>
1	SRR1039508	N61311	none
2	SRR1039509	N61311	dexa
3	SRR1039512	N052611	none
4	SRR1039513	N052611	dexa
5	SRR1039516	N080611	none
6	SRR1039517	N080611	dexa
7	SRR1039520	N061011	none
8	SRR1039521	N061011	dexa

# Generating multiple alignment commands with R

```
sample_table <- read.csv( "reads/sample_table.csv" )

command <- paste0(
  "STAR --runThreadN 2 --genomeDir GRCh38.84 \\n",
  "  --readFilesIn reads/###_1.fastq \\n",
  "    reads/###_2.fastq \\n",
  "  --outSAMtype BAM SortedByCoordinate \\n",
  "  --outFileNamePrefix star_out/###_ \\n",
  "  --quantMode GeneCounts\\n",
  "\\n" )

f <- file( "aligner_commands.sh", open="w" )
for( sampleID in sample_table$Library_ID ) {
  cat( gsub( "###", sampleID, command, fixed=TRUE ),
    file=f )
}
close( f )
```



# Next step: Using the alignments

## RNA-Seq:

- What transcripts are there?
- Which genes change expression due to treatment?
  - Mike's talk and this afternoon's lab

## ChIP-Seq:

- Peak finding and comparing
  - Thursday lab

## In general:

- The real work starts only *now*.  
So far, this was just preparation.

\*