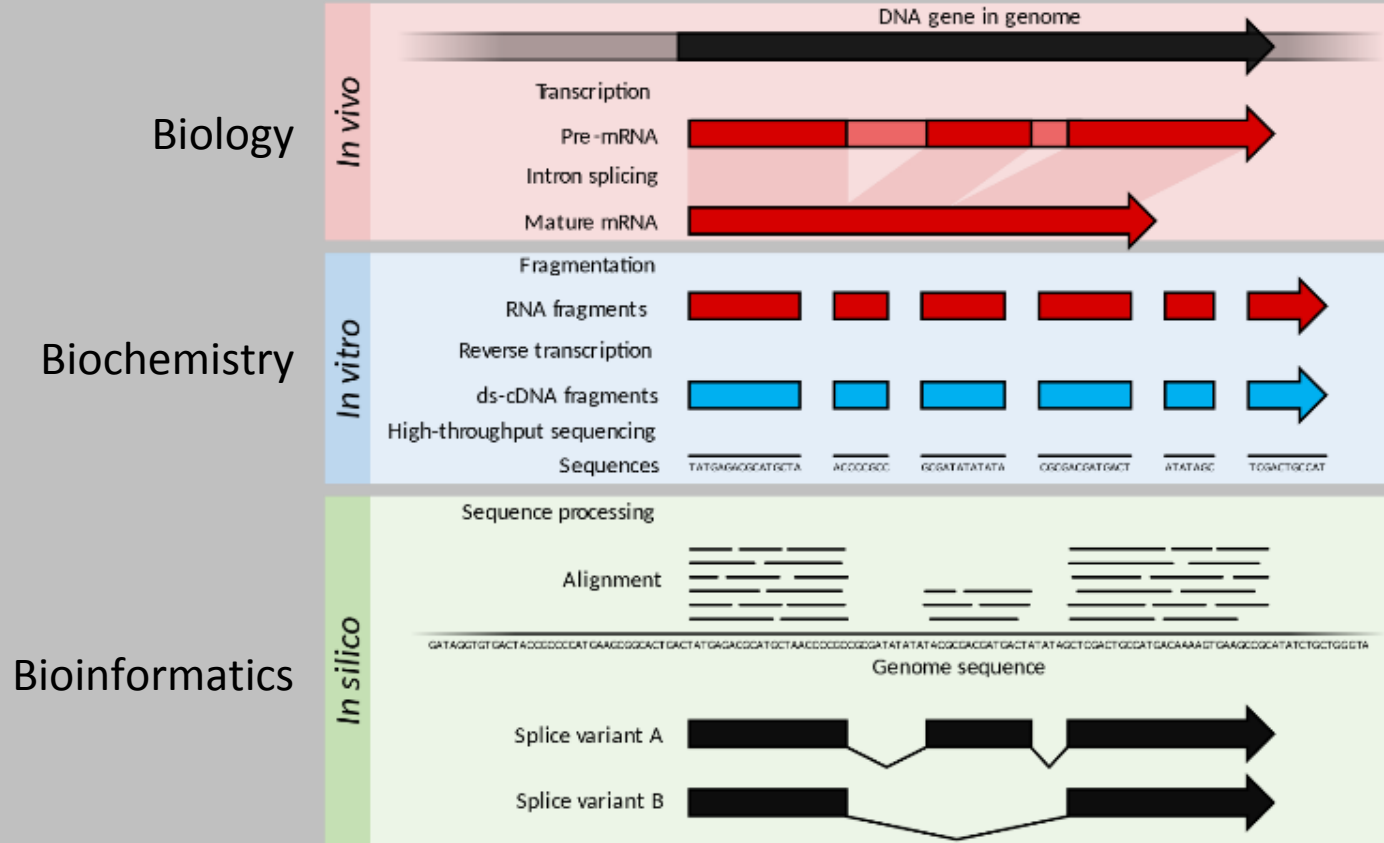


RNA Sequencing workflow

CANBIO

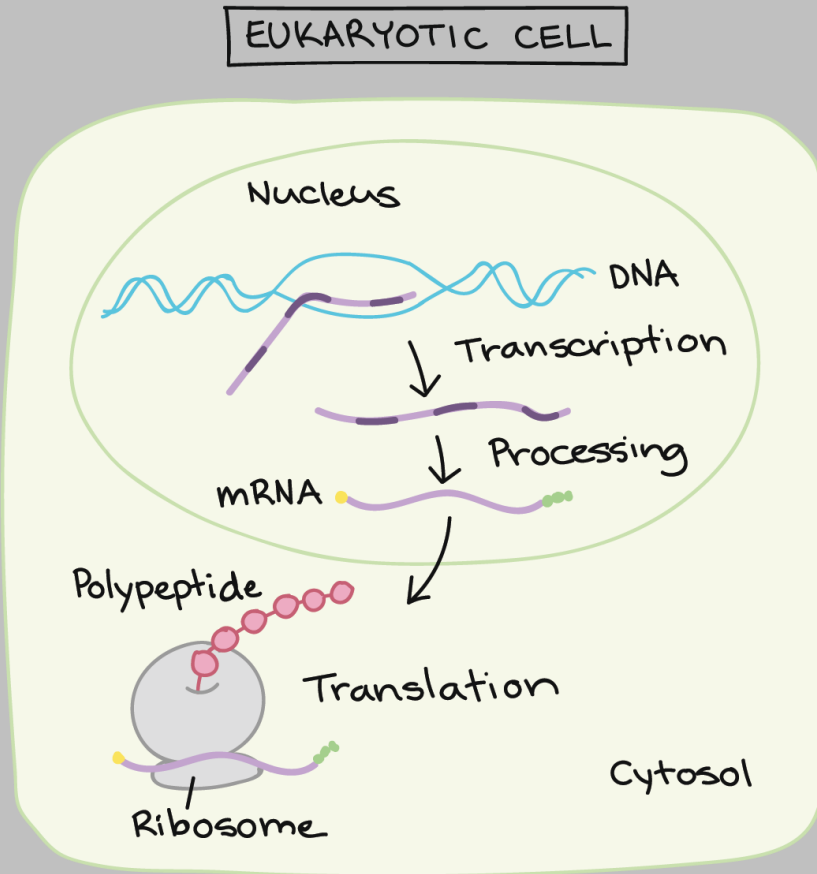
Arnaud Muller



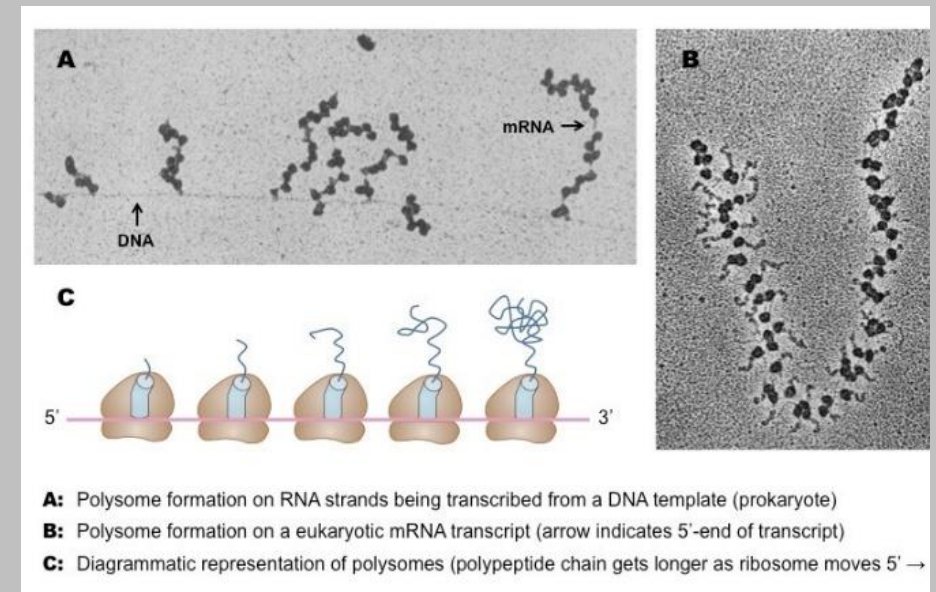
source: wikipedia

in vivo

mRNA processing

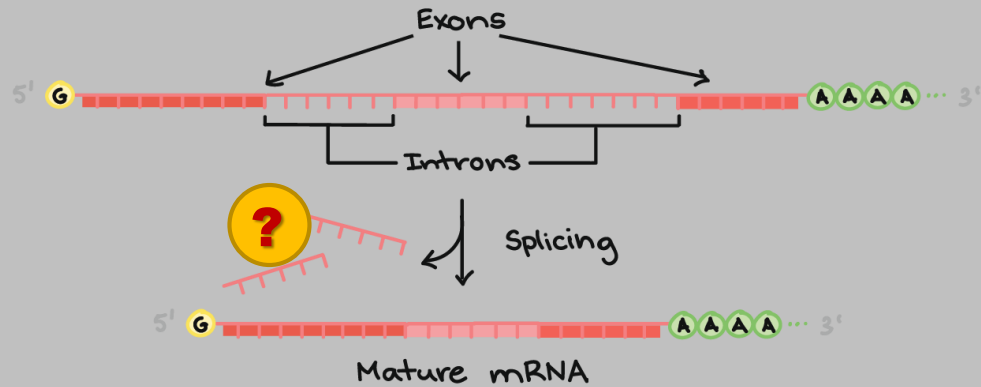


source: khanacademy



mRNA processing

- 5' capping: addition of a modified *guanine* 5'-end of the pre-mRNA
- 3' poly-A tail: a polyadenylation signal (AAUAAA) is recognized, followed by a cleavage and polyA tail synthesis (200 ntd)
- Splicing:



source: khanacademy

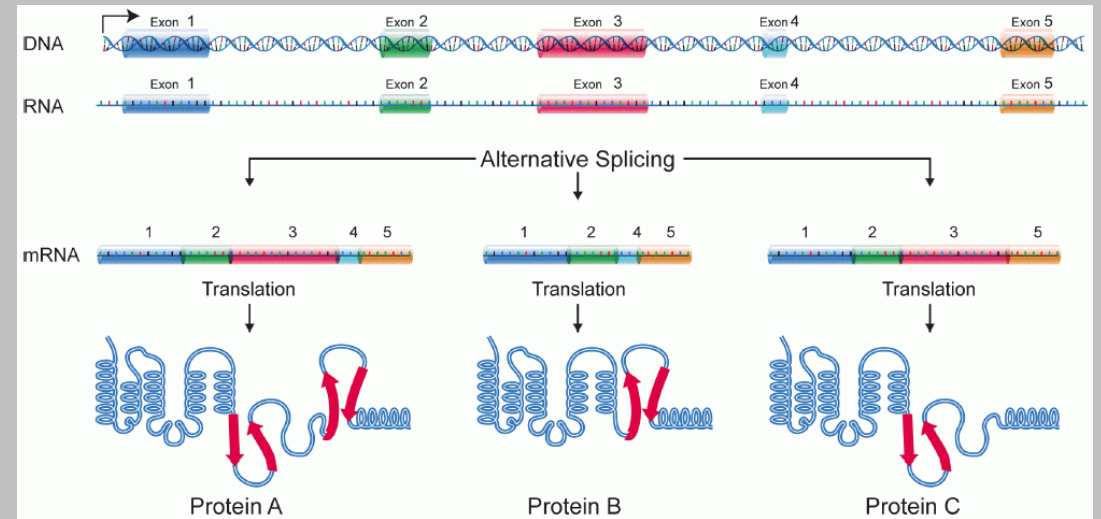
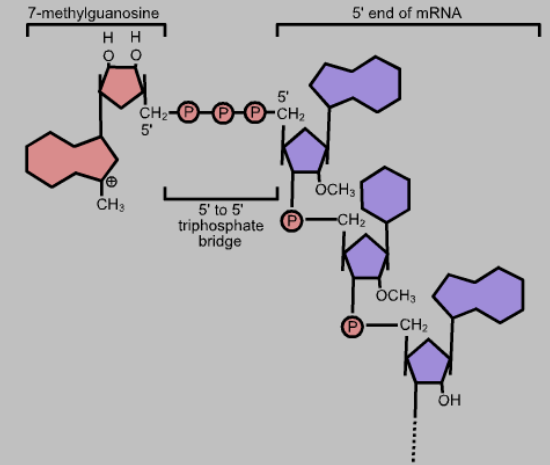


Image credit: "DNA, alternative splicing," by the National Human Genome Research Institute.

THEDOGRAMAPQANANDAYAPTQMTETHEHAT

Are you able to visually splice that sequence?

THEDOGRAMAPQANANDAYAPTQMTETHEHAT

THEDOGRAMAPQANANDAYAPTQMTETHEHAT

Here are the introns.

Splicing exercise

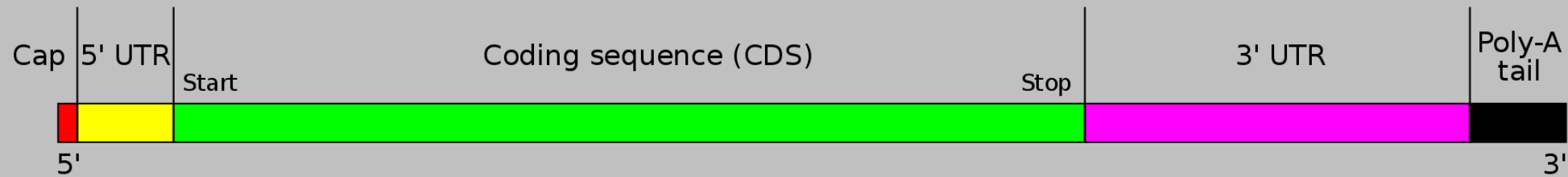
THE DOGRAMAPQANANDAYAPTQMTETHEHAT

THE DOGRAMAPQANANDAYAPTQMTETHEHAT

THE DOG R-----AN AND A-----TE THE HAT

mature mRNA structure

The structure of a typical human protein coding mRNA including the untranslated regions (UTRs)



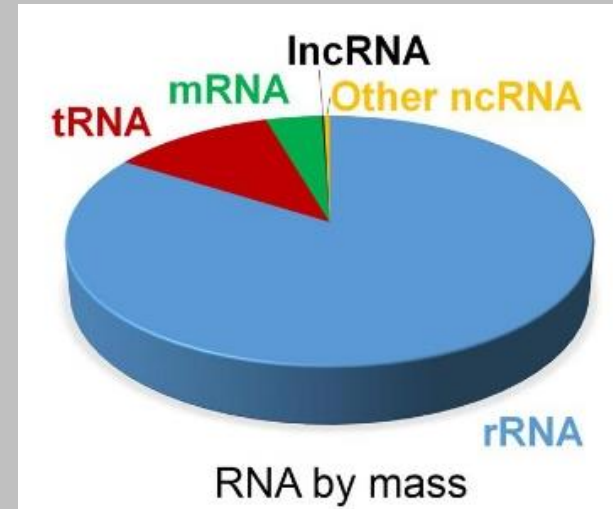
This dogmatic view is INcomplete.

RNA biotypes abundance

Type	Percent of total RNA by mass	Molecules per cell	Average size (kb)	Total weight picograms/cell	Notes	Reference
rRNAs	80 to 90	3–10 × 10 ⁶ (ribosomes)	6.9	10 to 30		Blobel and Potter (1967), Wolf and Schlessinger (1977), Duncan and Hershey (1983)
tRNA	10 to 15	3–10 × 10 ⁷	<0.1	1.5 to 5	About 10 tRNA molecules /ribosome	Waldron and Lacroute (1975)
mRNA	3 to 7	3–10 × 10 ⁵	1.7	0.25 to 0.9		Hastie and Bishop (1976), Carter et al. (2005)
hnRNA (pre-mRNA)	0.06 to 0.2	1–10 × 10 ³	10*	0.004 to 0.03	Estimated at 2–4% of mRNA by weight	Mortazavi et al. (2008), Menet et al. (2012)
Circular RNA	0.002 to 0.03	3–20 × 10 ³	~0.5	0.0007 to 0.005	Estimated at 0.1–0.2% of mRNA**	Salzman et al. (2012), Guo et al. (2014)
snRNA	0.02 to 0.3	1–5 × 10 ⁵	0.1–0.2	0.008 to 0.04		Kiss and Filipowicz (1992), Castle et al. (2010)
snoRNA	0.04 to 0.2	2–3 × 10 ⁵	0.2	0.02 to 0.03		Kiss and Filipowicz (1992), Cooper (2000), Castle et al. (2010)
miRNA	0.003 to 0.02	1–3 × 10 ⁵	0.02	0.001 to 0.003	About 10 ⁵ molecules per 10 pg total RNA	Bissels et al. (2009)
7SL	0.01 to 0.2	3–20 × 10 ⁴	0.3	0.005 to 0.03	About 1–2 SRP molecules/100 ribosomes	Raue et al. (2007), Castle et al. (2010)
Xist	0.0003 to 0.02	0.1–2 × 10 ³	2.8	0.0001 to 0.003		Buzin et al. (1994), Castle et al. (2010)
Other lncRNA	0.03 to 0.2	3–50 × 10 ³	1	0.002 to 0.03	Estimated at 1–4% of mRNA by weight	Mortazavi et al. (2008), Ramsköld et al. (2009), Menet et al. (2012)

*The size for the average unspliced pre-mRNA is 17 kb; however, most pre-mRNAs are partially spliced at any given time, and the average size of hnRNA is estimated at 10 kb (Salditt-Georgieff et al., 1976).

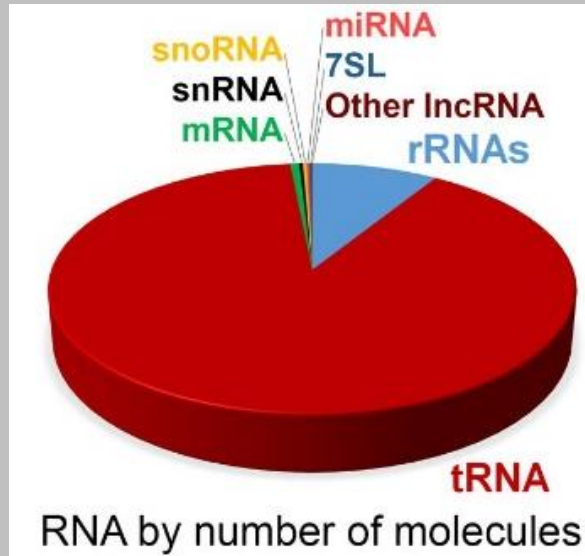
**Based on the finding that 1–2% of all mRNA species generate circular RNA, which is present at 10% of the level of the parental mRNA.



Alexander F. Palazzo and Eliza S. Lee, Front. Genet., 2015

- ribosomal RNA represents up to 90% of the total RNA mass, the second most represented biotype is
- tRNA, followed by
- mRNA (3% to 7%) and
- pre-mRNA (<0.2%).
- ... only later comes the lncRNA, microRNA, circRNA ...

Types/abundance of RNA-Seq



Alexander F. Palazzo and Eliza S. Lee, Front. Genet., 2015

Number of molecules per cell:

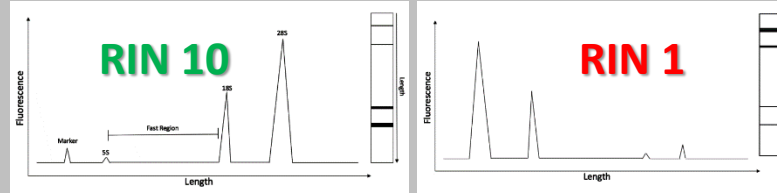
- tRNA: 100E06
- rRNA: 10E06
- mRNA: 1E06
- microRNA: 0.1E06

in vitro

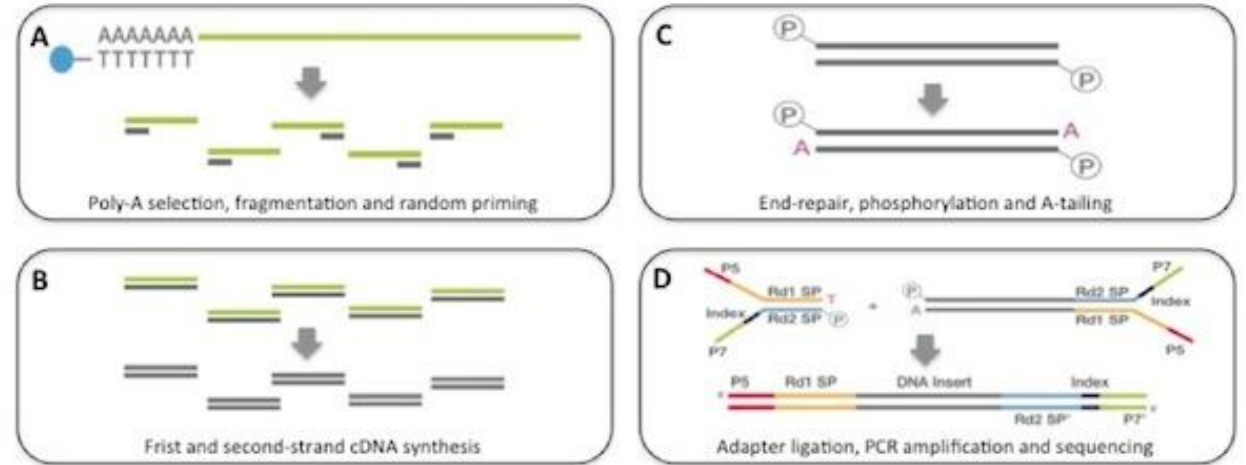
Library preparation

Prepare a complementary DNA (cDNA) library ready for sequencing, with the following main steps:

- RNA extraction (RNA Integrity Number)
- RNA selection or depletion:
 - polyA selection
 - ribosomal RNA depletion
 - targeted RNA capture
 - small RNA
 - ... many library preparation methods ...
- Fragmentation
- double strand cDNA synthesis (random priming)
- Addition of adapters, indexes and sequencing primers
- Amplification



Illumina Tru-Seq RNA-seq protocol

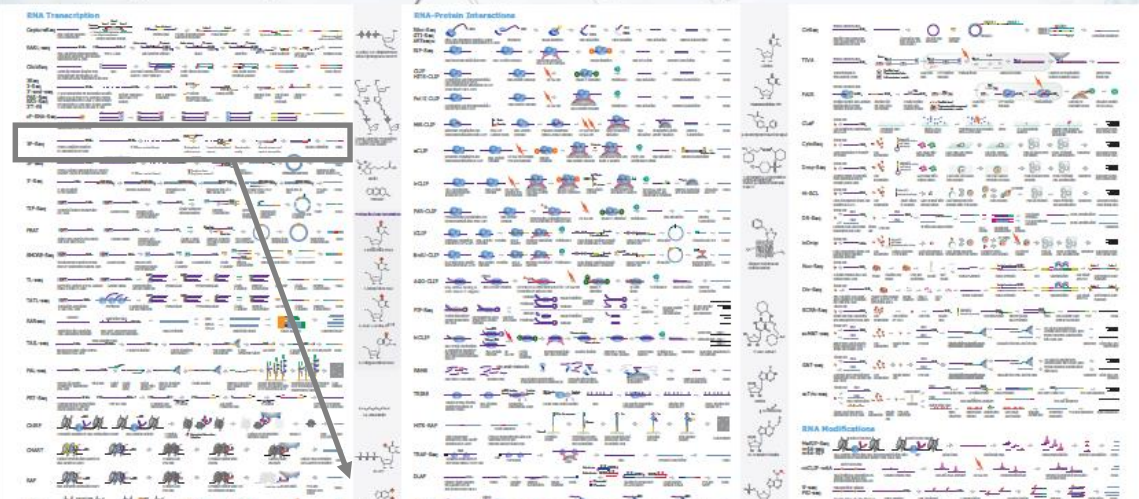


Library prep begins from 100ng-1ug of Total RNA which is poly-A selected (A) with magnetic beads. Double-stranded cDNA (B) is phosphorylated and A-tailed (C) ready for adapter ligation. The library is PCR amplified (D) ready for clustering and sequencing.

Library preparation

RNA

For all you seq...



3P-Seq

Poly(A)-position profiling
by sequencing (3P-Seq)

T1 RNase partial digest

TT(T)_n Biotin

Biotinylated
splint primer

AA(A)_n
TT(T)_n
Ligate biotinylated
primer

Streptavidin
purify

Anneal primer and
reverse-transcribe

RNase H digestion

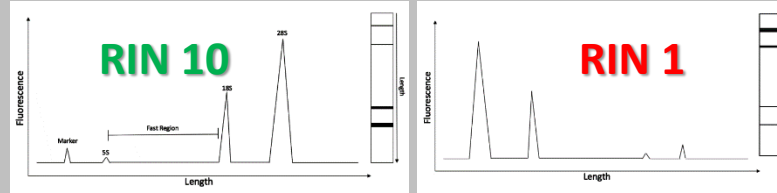
cDNA

illumina

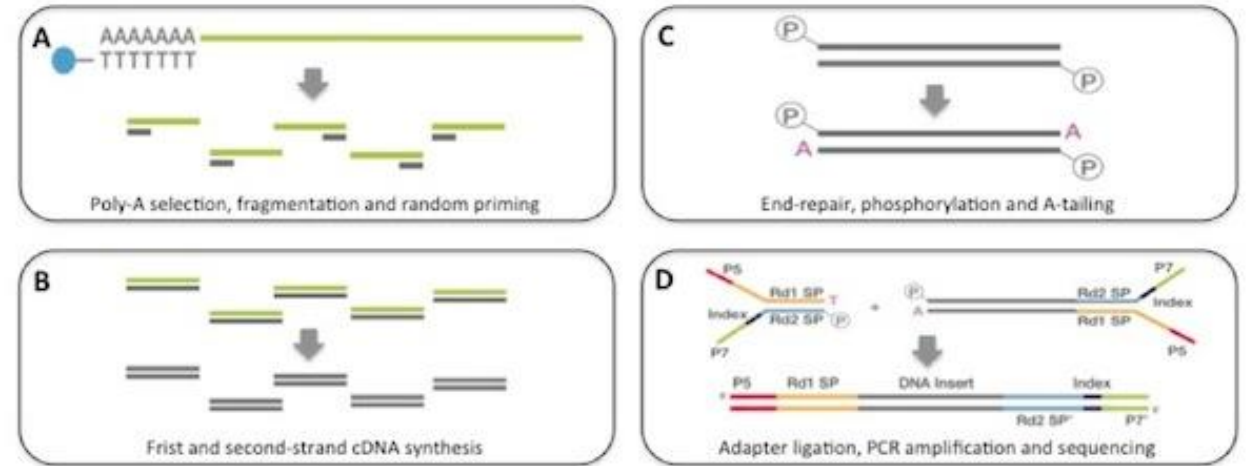
Library preparation

Prepare a complementary DNA (cDNA) library ready for sequencing, with the following main steps:

- RNA extraction (RNA Integrity Number)
- RNA selection or depletion:
 - polyA selection
 - ribosomal RNA depletion
 - targeted RNA capture
 - small RNA
 - ... many library preparation methods ...
- Fragmentation
- double strand cDNA synthesis (random priming)
- Addition of adapters, indexes and sequencing primers
- Amplification



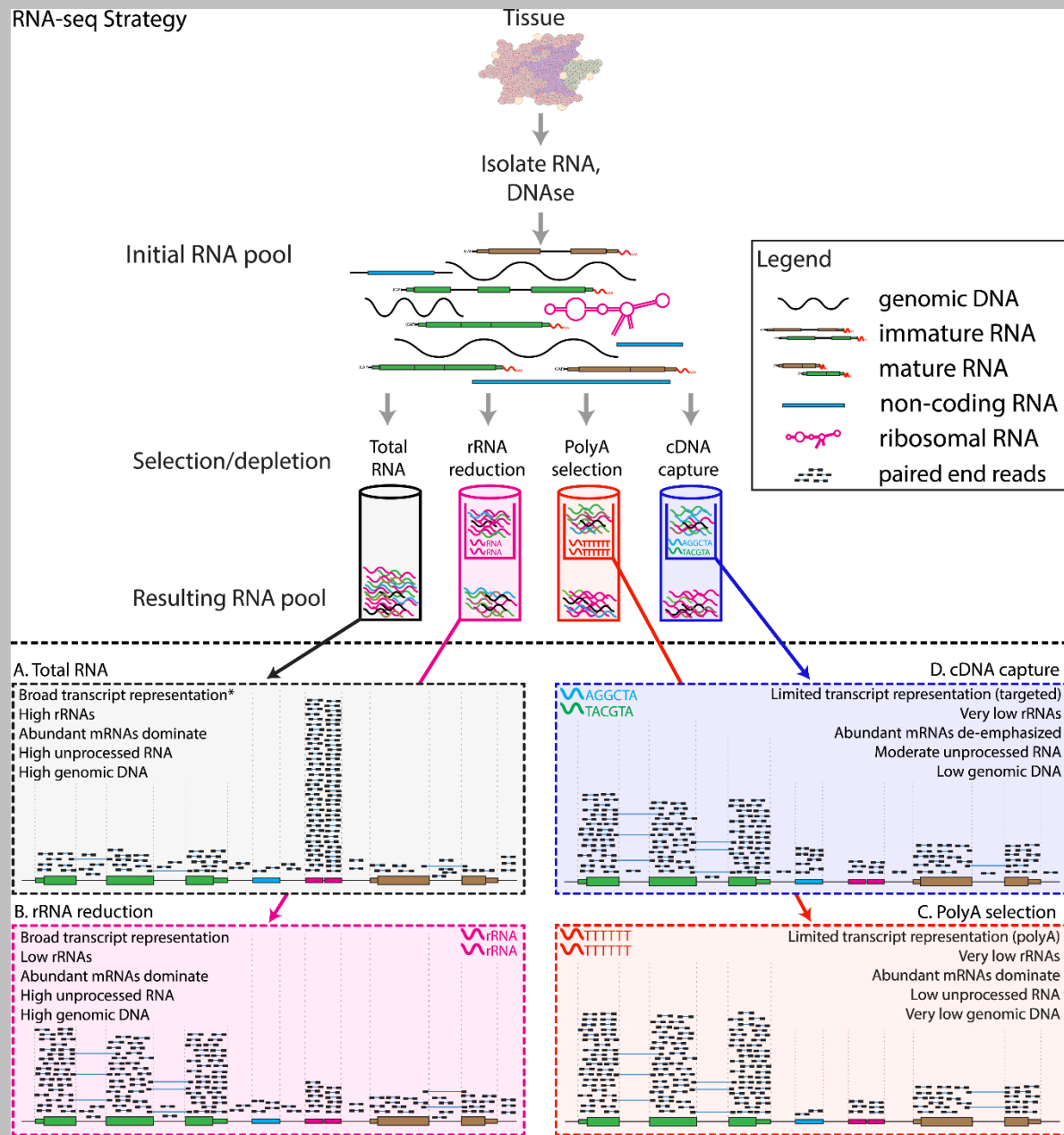
Illumina Tru-Seq RNA-seq protocol



Library prep begins from 100ng-1ug of Total RNA which is poly-A selected (A) with magnetic beads. Double-stranded cDNA (B) is phosphorylated and A-tailed (C) ready for adapter ligation. The library is PCR amplified (D) ready for clustering and sequencing.

Library preparation

RNA-seq Strategy



Expected Alignments

Informatics for RNA Sequencing: A Web Resource for Analysis on the Cloud Malachi Griffith

Library preparation

Since the library is prepared, it's now ready for sequencing:

- cluster generation
- Sequencing By Synthesis

HiSeq 2000
New flow cell design

LARGER, DUAL-SURFACE ENABLED
>5x increase in imaging area
Retains 8 lane format

Compatible with cBot

Cluster density
750-850/mm²

HiSeq Flow Cells

A

flowcell ID
flow in
flow out
barcode

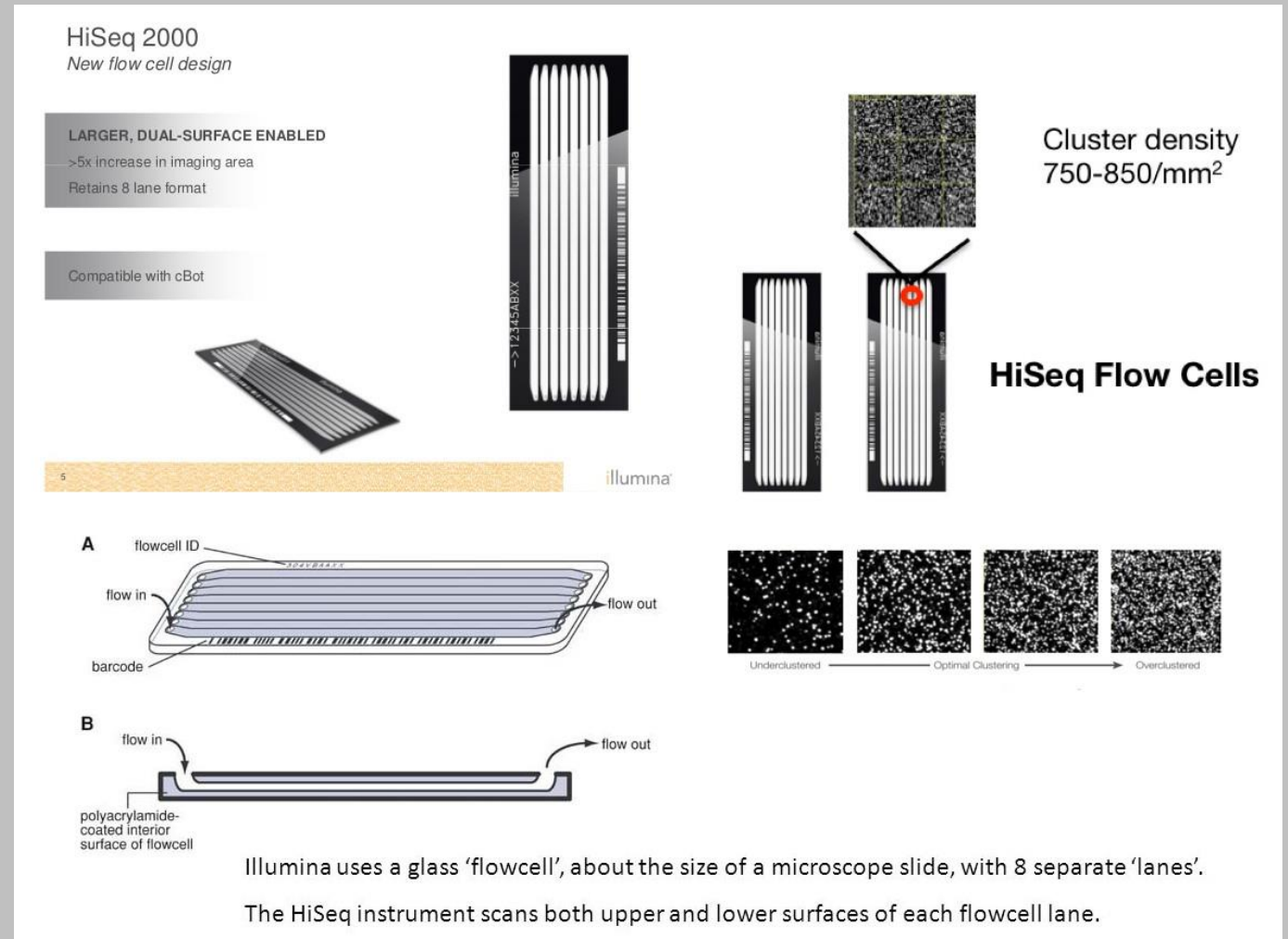
B

flow in
flow out
polyacrylamide-coated interior surface of flowcell

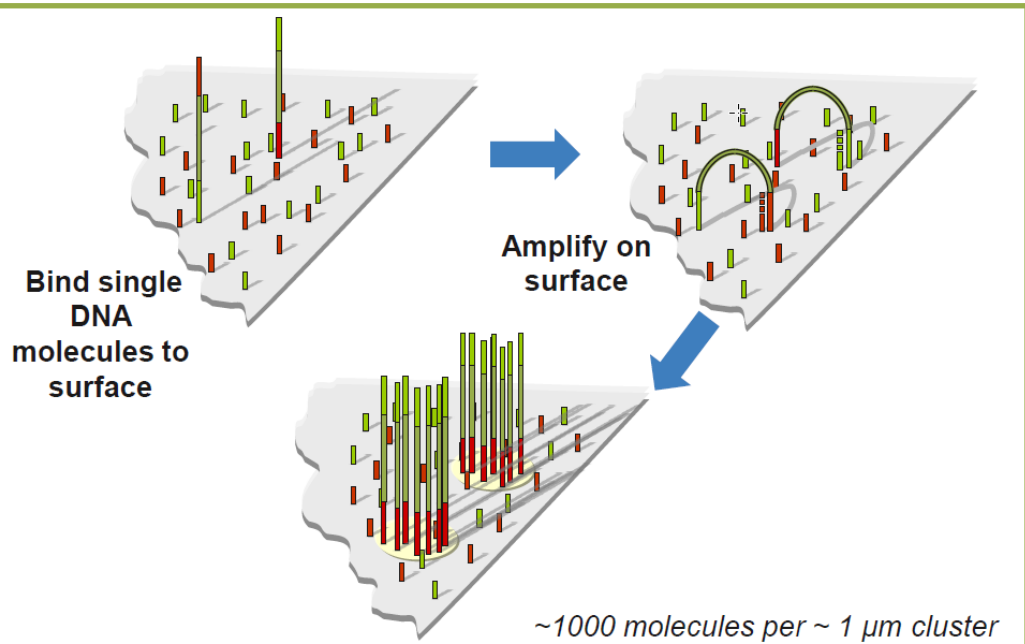
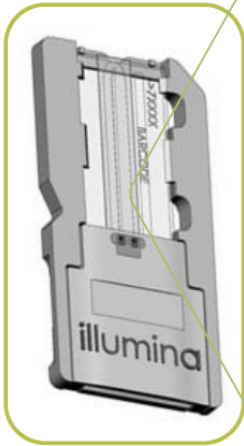
Underclustered Optimal Clustering Overclustered

illumina

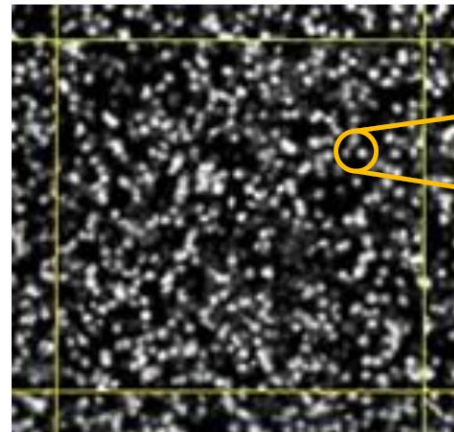
illumina uses a glass 'flowcell', about the size of a microscope slide, with 8 separate 'lanes'.
The HiSeq instrument scans both upper and lower surfaces of each flowcell lane.



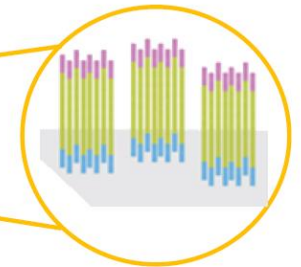
Cluster generation



Clusters are bright spots on an image

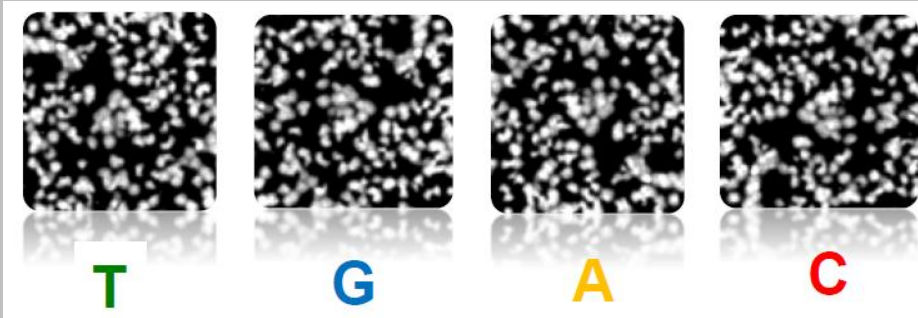


Each cluster represents thousands of copies of the same DNA strand in a 1–2 micron spot



Sequencing By Synthesis

1. Addition of A C G T which contains a fluorophore and a terminator + polymerase
2. Integration of the nucleoside
3. Laser excitation + image capture (4 times!)
4. Removal of fluorescent dye and terminator

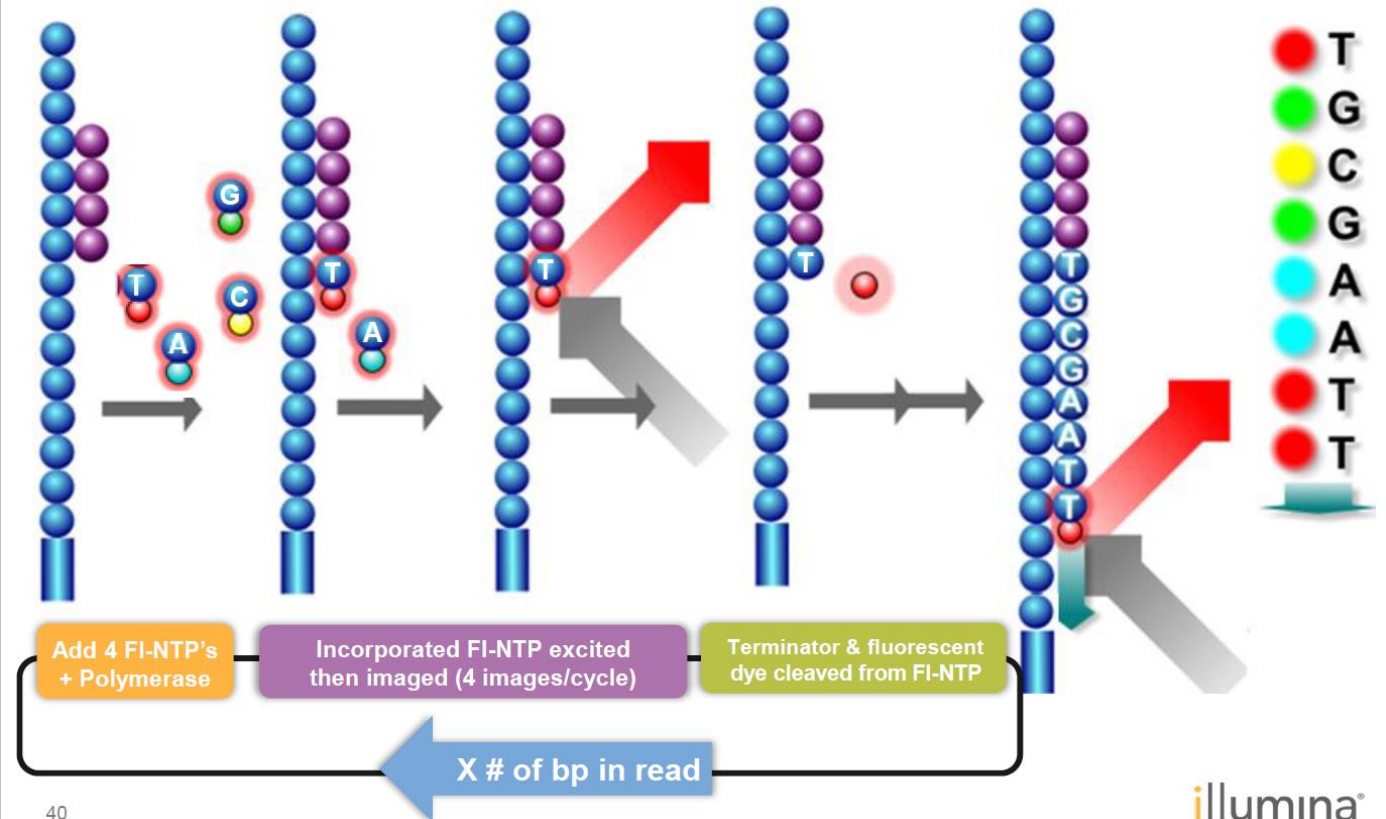


These 4 steps are repeated as many times as the read length.

- > computational/storage cost: made on the instrument, on the fly
- > biochemical cost
- > microfluidics
- > optics
- > time: 10-30 h

A Closer Look At 4-Dye Chemistry

Four channel chemistry

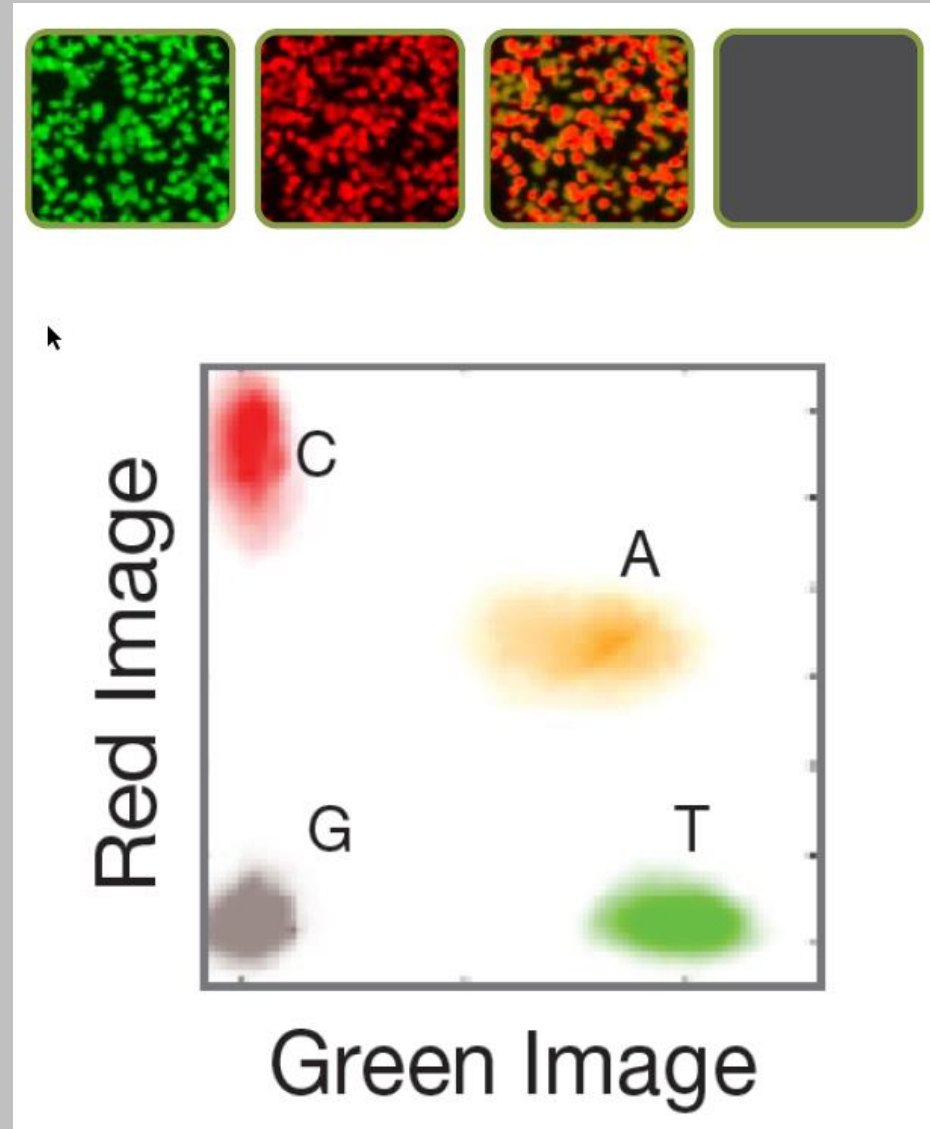


Sequencing By Synthesis

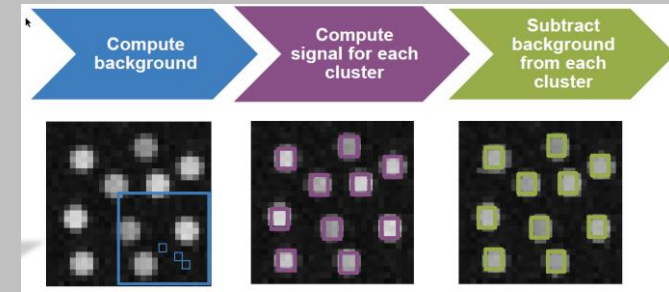
Two channel SBS uses 2 images

- Builds template over 5 cycles
- Clusters appearing in green only are T
- Clusters appearing in red only are C
- Clusters appearing in both images are A
- Clusters not present in either green nor red are G
- Cluster intensities are plotted and bases are called accordingly

-> time cost reduced by 50%



- The Sequencing By Synthesis steps generates BCL files.
- Binary Base Call (BCL) files contains the base and the confidence in the call.
- This file is generated on the fly.
- The bcl generation workflow is the following:
 - Template generation (cluster map)
 - extract intensities
 - intensity norm
 - phasing estimate
 - base call filtering
 - quality score (phred score): $Q40 \Leftrightarrow \text{Error probability} < 0.0001$
 - quality score is a combination of metrics, such as: intensity, S/N ratio, phasing..
- Because of the size of the optical sensor, the flow cell is split into multiple tiles, generating several bcl files.
- After the experiment, bcl file is converted into fastq file.



Base quality score

Q-score	Base Call Accuracy	Probability of Incorrect Based Call	ASCII Quality Score
Q10	90%	0.1	+
Q20	99%	0.01	5
Q30	99.9%	0.001	?
Q40	99.99%	0.0001	



in silico

coordinates of the cluster Y: cluster did not pass the filter

@<instrument>:<run number>:<flowcellID>:<lane>:<tile>: **<x-pos>:<y-pos>** <read> **<is filtered>** :
 <control number> : <sample number>

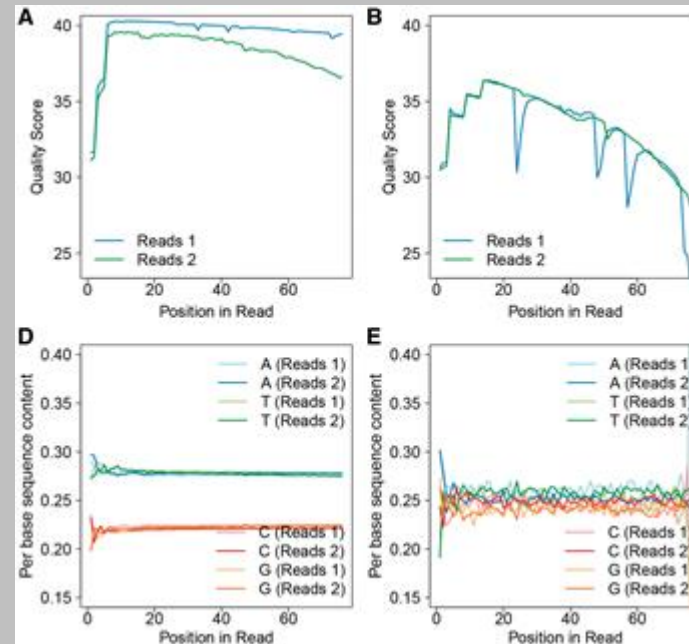
unique ID

```

$ gunzip -c 2T210_S1_R1_001.fastq.gz | head -4
@NB551409:68:H7CN7BGXB:1:11101:16989:1046 1:N:0:ATCCACTG+ACGCACCT
CTCTATACCANTGGTCCAATGGGCTTAAAAAAGAGCAAATATTACCAAATGGATATGCTCTGAAGTTGTCGTTAAT
+
AAAAAEEEEEE#EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE
  
```

- Line 1 begins with a '@' character and is followed by a sequence identifier and an *optional* description (like a [FASTA](#) title line).
- Line 2 is the raw sequence letters.
- Line 3 begins with a '+' character, this is historical and mysterious.
- Line 4 encodes the quality values for the sequence in Line 2, and must contain the same number of symbols as letters in the sequence.

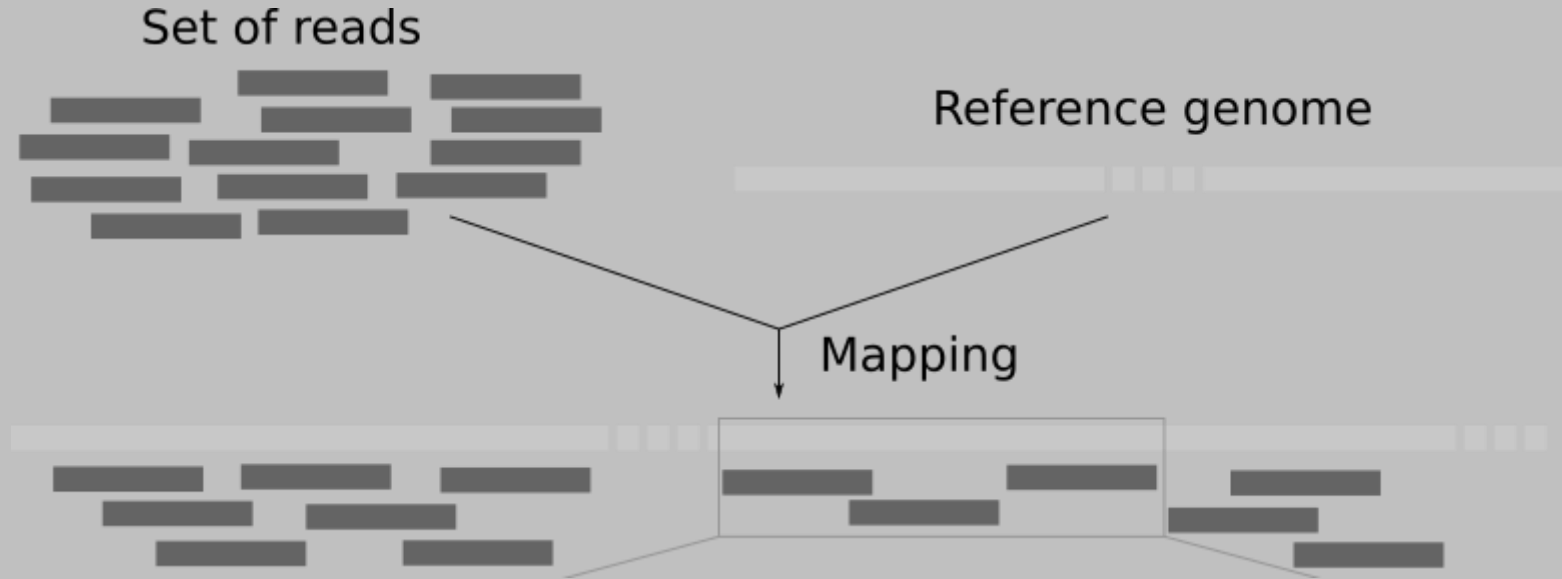
Quality Controls



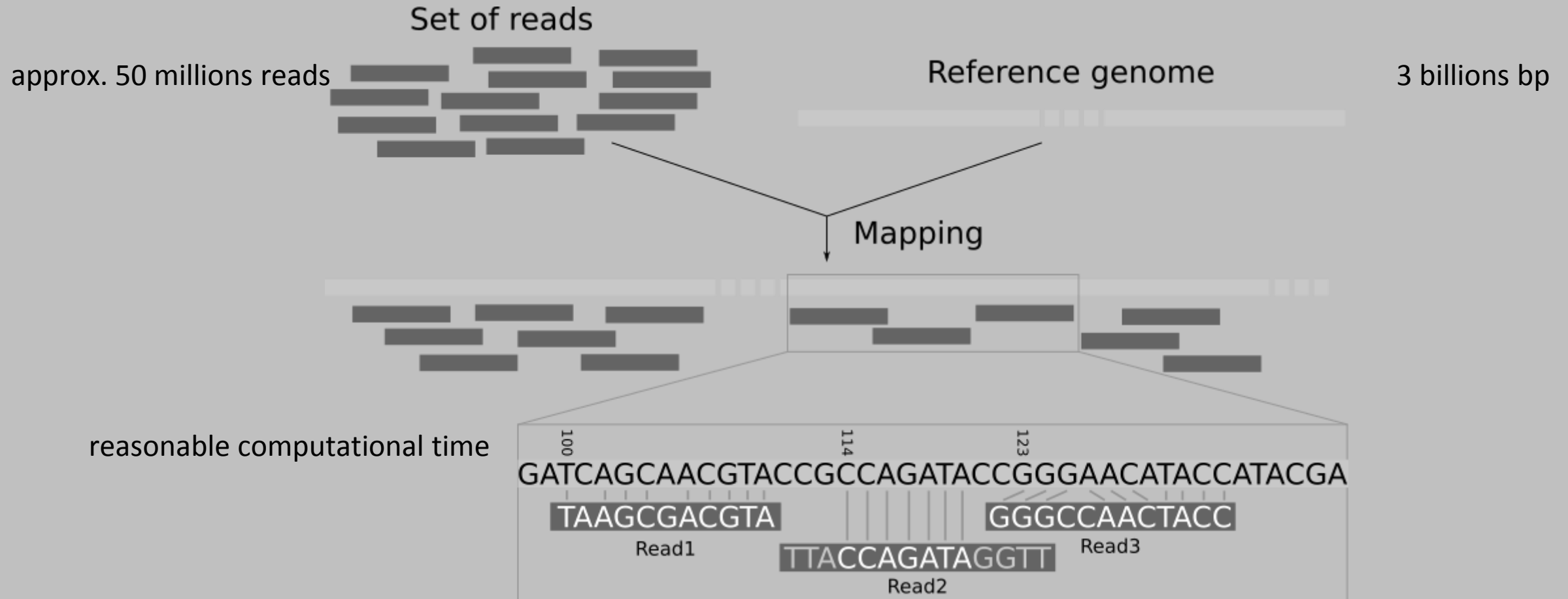
Quanhui Sheng et al., 2017, *Briefings in Functional Genomics*

- (A)** Example of a long RNA-seq sample with expected base quality score. Read 2 tends to have a slightly lower median base score than read 1, but it is not usually a quality concern.
- (B)** Example of a long RNA-seq sample with potential base quality problem, as denoted by the sudden drops of median base quality in read 2 of pair-end read sequencing.
- (D)** Example of a long RNA-seq sample with expected nucleotide distribution, as denoted by the stable nucleotide distribution across the samples.
- (E)** Example of a long RNA-seq sample with a potential nucleotide distribution issue, as denoted by the unstable distribution across the cycles.

Mapping - the ingredients

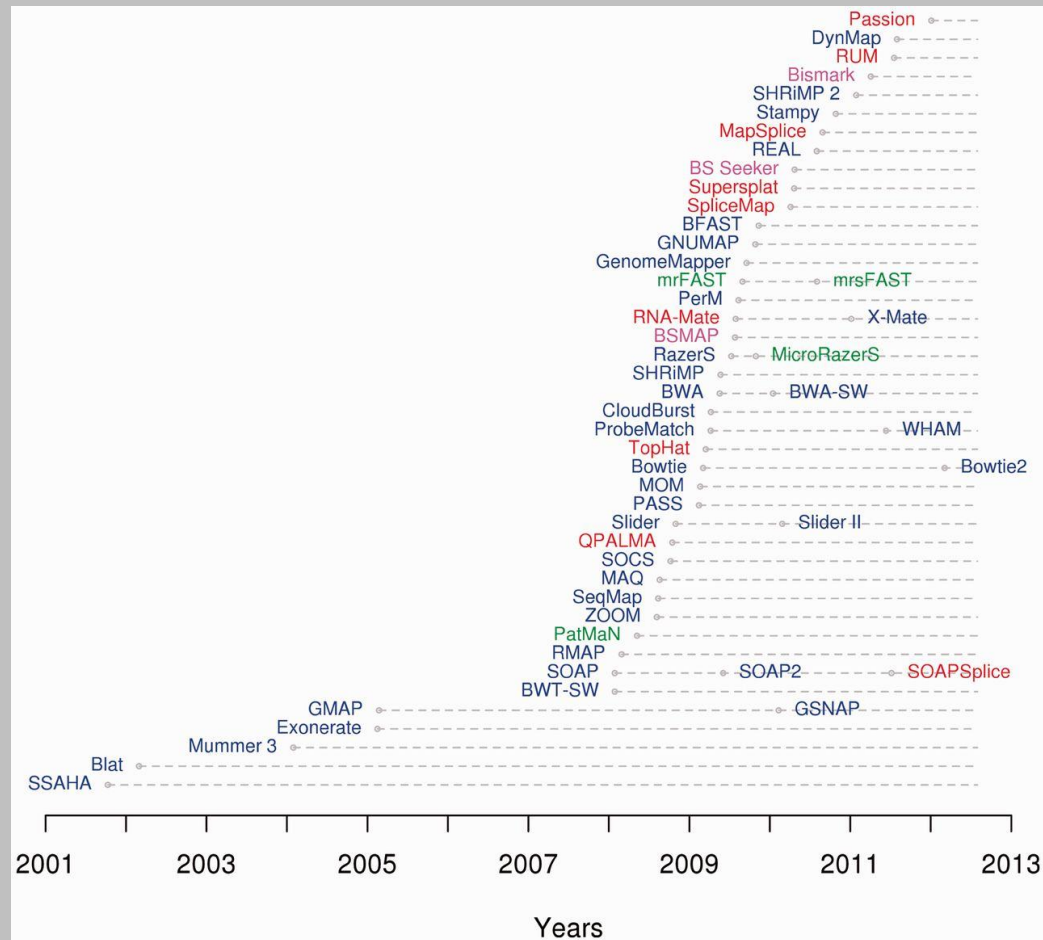


Mapping - the challenge



Mapping tools

Mappers timeline (since 2001). DNA mappers are plotted in blue, RNA mappers in red, miRNA mappers in green.



Alignment scoring

Sequence One : GGCTGG						
Sequence Two : GAGG						
	G	G	C	T	G	G
	G	A	-	-	G	G
	10	-5	-5	-1	10	10
	10	5	0	-1	9	19
Your cumulative score						

- **Reward** for a match (e.g. +10), **penalty** for a mismatch (e.g. -5)
- **Penalty** for gaps
 - *Linear*: every gap same penalty (e.g. -5)
 - *Affine*: gap open vs gap extend (e.g. -5 and -1)
- Different tools use different scoring values (and give different results)

Alignment scoring

- Reference: AAA CAGTGA GAA
- Observed: AAA TCTCT GAA

Alignment:

```
AAA-CAGTGAGAA
|||-|--|::|||
AAATC-TCTGAA
```

```
AAACAGTGAGAA
|||-::|::|||
AAA-TCTCTGAA
```

```
AAACAGTGAGAA
|||:-:|::|||
AAAT-CTCTGAA
```

```
AAACAGTCA-----GAA
|||-----|||
AAA-----TCTCTGAA
```


Alignment scoring

- Reference: AAA CAGTGA GAA
- Observed: AAA TCTCT GAA

Alignment:

```
AAA-CAGTGAGAA
|||---|::|||
AAATC-TCTGAA
```

Novoalign

```
AAACAGTGAGAA
|||---::|::|||
AAA-TCTCTGAA
```

Ssaha2

```
AAACAGTGAGAA
|||:-:|::|||
AAAT-CTCTGAA
```

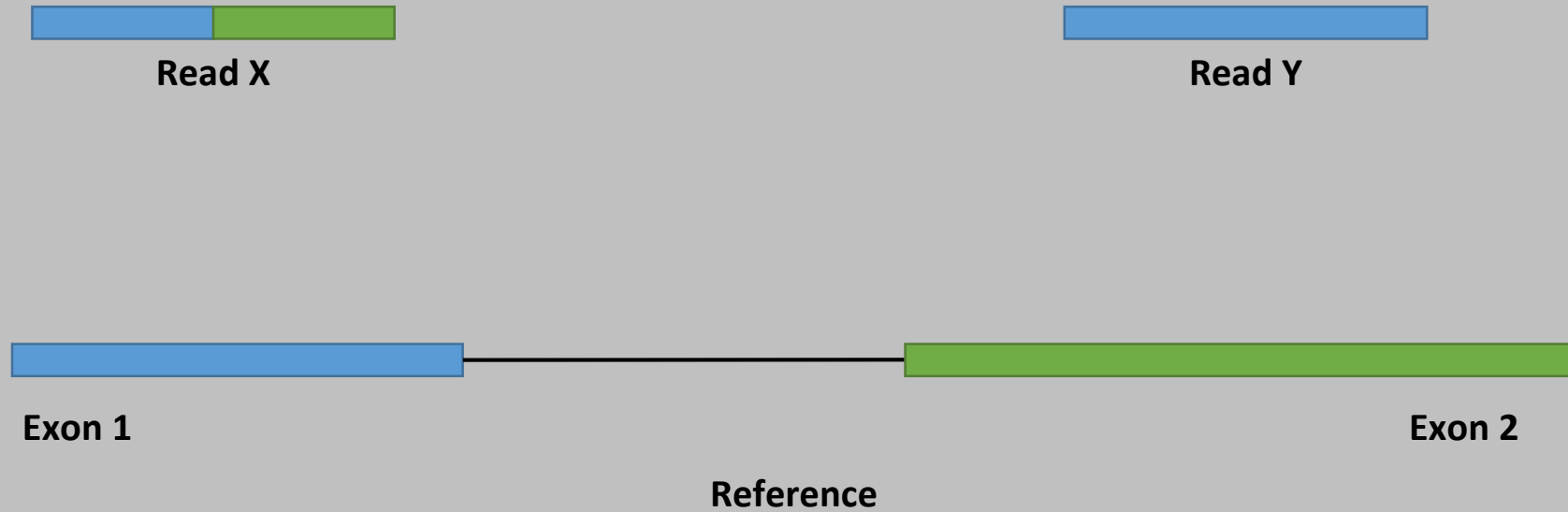
BWA

```
AAACAGTCA-----GAA
|||-----|||
AAA-----TCTCTGAA
```

Complete Genomics

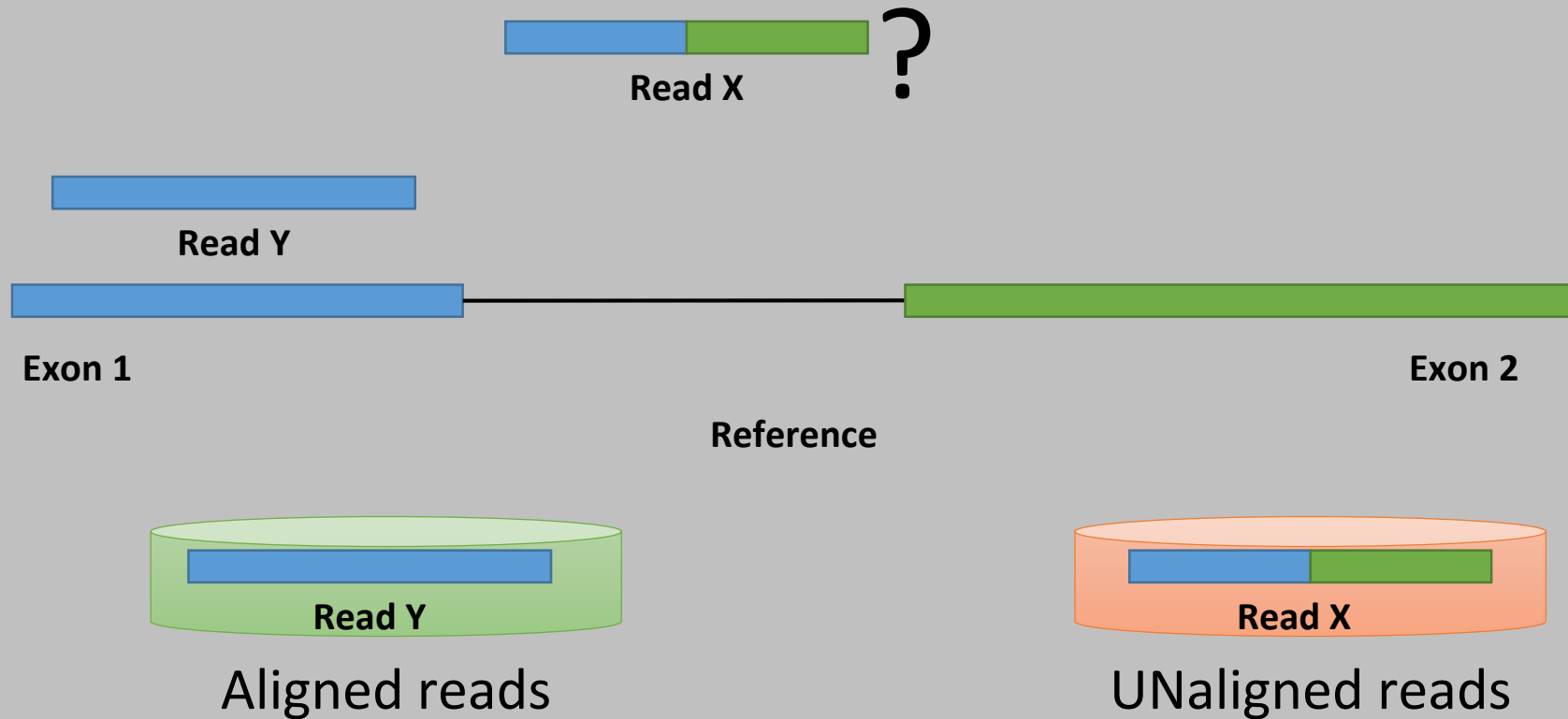
Mapping - bowtie

Bowtie / TopHat



Mapping - bowtie

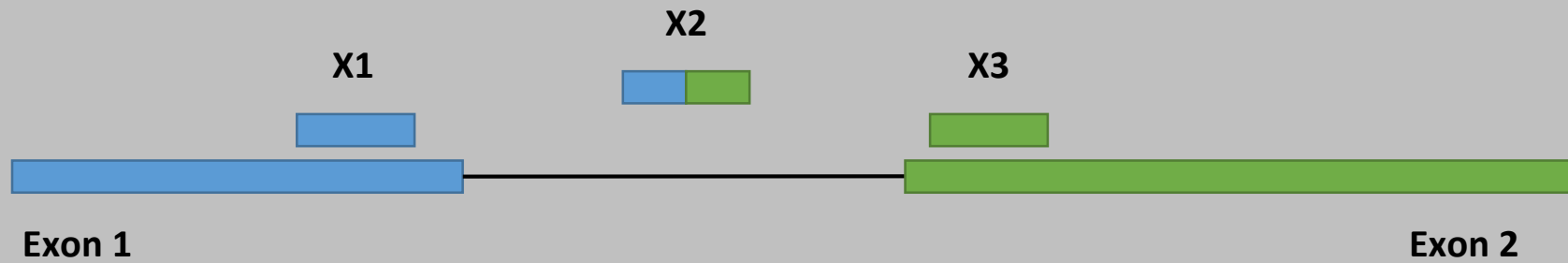
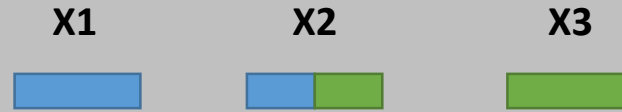
Bowtie / TopHat



Mapping - bowtie



UNaligned reads



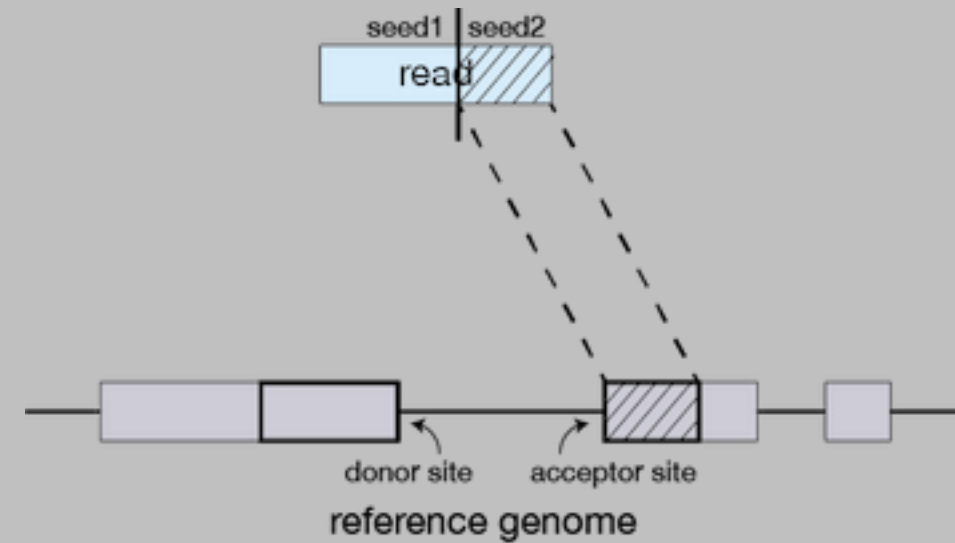
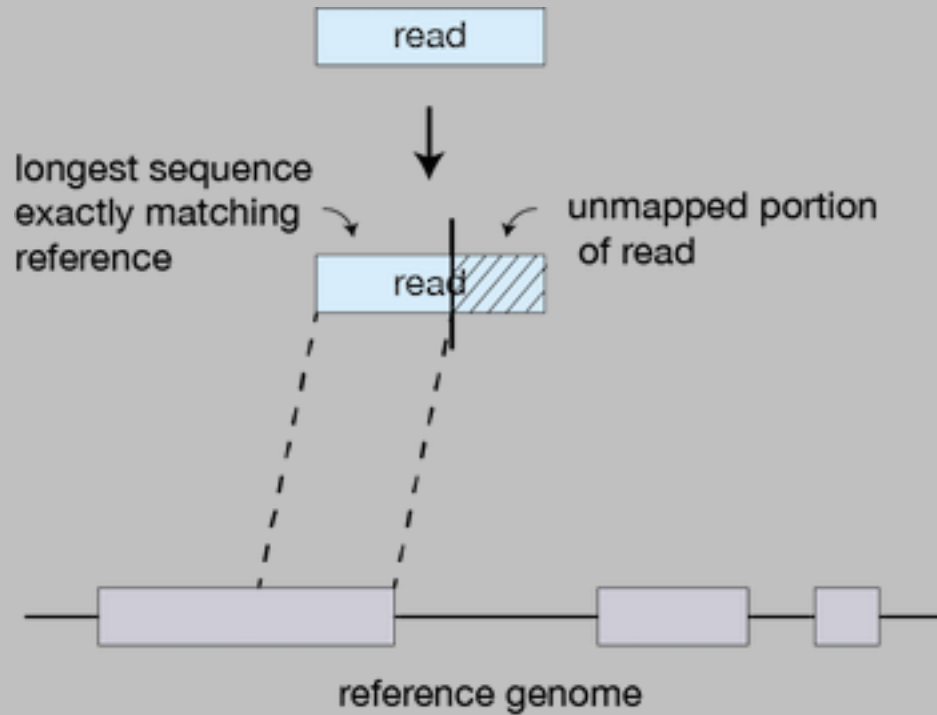
Reference

Collect Mapping Information for X1 and X3



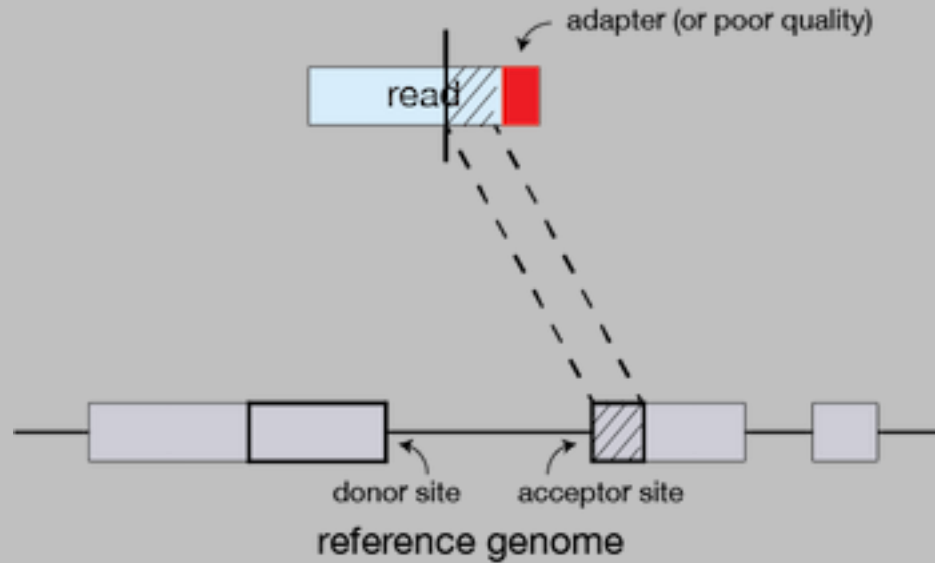
Construct a Splice Library

Mapping - STAR



1. Find the longest perfect match seed
2. Independent search for the unmapped portion

Mapping - STAR



If extension does not give a good alignment, then the poor quality or adapter sequence (or other contaminating sequence) will be **soft clipped**.

The separate seeds are **stitched** together to create a complete read by first clustering the seeds together based on proximity to a set of 'anchor' seeds, or seeds that are not multi-mapping.

Then the seeds are stitched together based on the best alignment for the read (scoring based on mismatches, indels, gaps, etc.).

Genome versions

- Completion of the *Human Genome Project* in 2003.
- Continuation by the *Genome Reference Consortium* (GRC):
 - NCBI
 - Wellcome Trust Sanger Institute
 - European Bioinformatics Institute
 - Genome Institute at Washington University
- Now: closing gaps, centromeres/telomeres representation, error correction

	Release Name	UCSC Version	Release Date
Human	NCBI Build 34	hg16	July 2003
	NCBI Build 35	hg17	May 2004
	NCBI Build 36.1	hg18	March 2006
	GRCh37	hg19	February 2009
	GRCh38	hg38	December 2013

Genome sequence file

```
head ~/Tools/Genome/Human/hg19/chrM.fa
>chrM
GATCACAGGTCTATCACCTATTAACCACTCACGGGAGCTCTCCATGCAT
TTGGTATTTTCGTCTGGGGGGTGTGCACGCGATAGCATTGCGAGACGCTG
GAGCCGGAGCACCTATGTCGCAGTATCTGTCTTTGATTCCTGCCTCATT
CTATTATTTATCGCACCTACGTTCAATATTACAGGCGAACATACCTACTA
AAGTGTGTTAATTAATTAATGCTTGTAGGACATAATAACAATTGAAT
GTCTG
```

```
head ~/Tools/Genome/Human/hg19/chrM.fa
>chrM
GATCACAGGTCTATCACCTATTAACCACTCACGGGAGCTCTCCATGCAT
TTGGTATTTTCGTCTGGGGGGTGTGCACGCGATAGCATTGCGAGACGCTG
GAGCCGGAGCACCTATGTCGCAGTATCTGTCTTTGATTCCTGCCTCATT
CTATTATTTATCGCACCTACGTTCAATATTACAGGCGAACATACCTACTA
AAGTGTGTTAATTAATTAATGCTTGTAGGACATAATAACAATTGAAT
GTCTGCACAGCCGCTTCCACACAGACATCATAACAAAAAATTTCCACCA
```


Genome sequence file

```
$ wc -l `ls |egrep 'chr[0-9XYM]+.fa'`
 2710696 chr10.fa
 2700132 chr11.fa
 2677039 chr12.fa
 2303399 chr13.fa
 2146992 chr14.fa
 2050629 chr15.fa
 1807097 chr16.fa
 1623906 chr17.fa
 1561546 chr18.fa
 1182581 chr19.fa
 4985014 chr1.fa
 1260512 chr20.fa
 962599 chr21.fa
 1026093 chr22.fa
 4863989 chr2.fa
 3960450 chr3.fa
 3823087 chr4.fa
 3618307 chr5.fa
 3422303 chr6.fa
 3182775 chr7.fa
 2927282 chr8.fa
 2824270 chr9.fa
    333 chrM.fa
 3105413 chrX.fa
 1187473 chrY.fa
```

61913917 total

=> leading to 3 095 695 850 nucleotides.

Sequence Alignment Map - SAM

```
$ more S1AR005ACAGTG.sam
@HD      VN:1.0   SO:unsorted
@SQ      SN:1     LN:248956422
@SQ      SN:2     LN:242193529
@PG      ID:bowtie2      PN:bowtie2      VN:2.3.4.1      CL:"/mnt/pcpnfs/homedirs/luxgen/Tools/Bowtie/bowtie2-
2.3.4.1-linux-x86_64/bowtie2-align-1 --wrapper basic-0 -p20 -t -x
/mnt/pcpnfs/homedirs/luxgen/Tools/Genome/ensembl/Homo_sapiens_GRCh38 --met-file S1AR005
ACAGTG.metrics --passthrough -1 ../fastq/S1AR005ACAGTG_S1_R1_001.fastq.gz -2
../fastq/S1AR005ACAGTG_S1_R2_001.fastq.gz"

NB551409:37:H5FHLBGXB:1:11101:3778:1083 83      16      177285 1      75M      =      177055 -305
CCTAAGCCACTGCCTGCTGGTGACCCTGGCCGCCACCTCCCNNGAGTTACCCNTGCGGTGCACGCCTCCCG
EEE/AEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE#E#EAEEEEEEEEE#EEEEEEEEEEEEEEEEAAAAA AS:i:-9 XS:i:-9 XN:i:0 XM:i:5
XO:i:0 XG:i:0 NM:i:5 MD:Z:42C0G1C10C17T0 YS:i:-3 YT:Z:CP

NB551409:37:H5FHLBGXB:1:11101:3778:1083 163      16      177055 1      76M      =      177285 305
GGACGACATGCCCAACGCGCTGTCCGCCCTGAGCGCCCTGCACGCGCACAAGCTTCGGGTGGACCCGGTCAACTTC
AA/AA/E/EEEEEA/EEEEEEEEEEEEEEEEEEEEEEEE/EEEEEEEEEEEEEEEEEEEEEEEE6E<EEE/EE/EEEE AS:i:-3 XS:i:-3 XN:i:0 XM:i:1
XO:i:0 XG:i:0 NM:i:1 MD:Z:35A40 YS:i:-9 YT:Z:CP

NB551409:37:H5FHLBGXB:1:11101:8321:1083 99      10      1051875 42      76M      =      1051950 150
AGGAAACTTTTTTTTTTGANACAGAGTCTCNTNNTTTCGCCCAGGCTGAAGTGCAGTGGTGCATCTTGGCTCACTG
AAAAAEEEEEEEEEEEEEEEE#EEEEEEEEEEEE#E##EEEEEEEEEEEEEEEEEEEEEEEEEEEEEEEE<AEEE AS:i:-4 XN:i:0 XM:i:4 XO:i:0
XG:i:0 NM:i:4 MD:Z:18G10G1T0C43 YS:i:0 YT:Z:CP
```

Sequence Alignment Map - SAM

```
$ more S1AR005ACAGTG.sam
@HD      VN:1.0   SO:unsorted
@SQ      SN:1     LN:248956422
@SQ      SN:2     LN:242193529
@PG      ID:bowtie2      PN:bowtie2      VN:2.3.4.1      CL:"/mnt/pcpnfs/homedirs/luxgen/Tools/Bowtie/bowtie2-
2.3.4.1-linux-x86_64/bowtie2-align-1 --wrapper basic-0 -p20 -t -x
/mnt/pcpnfs/homedirs/luxgen/Tools/Genome/ensembl/Homo_sapiens_GRCh38 --met-file S1AR005
ACAGTG.metrics --passthrough -1 ../fastq/S1AR005ACAGTG_S1_R1_001.fastq.gz -2
../fastq/S1AR005ACAGTG_S1_R2_001.fastq.gz"
```

Header:

- version and sorting status
- reference sequence name (chromosome) + length
- program used for read processing

More information can be added: <https://samtools.github.io/hts-specs/SAMv1.pdf>

@<instrument>:<run number>:<flowcellID>:<lane>:<tile>: <x-pos>:<y-pos>

CIGAR

read sequence

read quality score

- The purpose of the counting is to estimate the number of reads (counts) associated with each feature of interest (gene, exon, transcript).
- Input required: BAM/SAM files + General Transfer Format (**GTF**) file.
- Output generated: a gene expression matrix with genes as rows and samples as column.
- Top popular tools *htseq-count* (more cited, same developers as *DESeq*) and *featureCounts* (a way faster, Unix and R package, a bit more liberal).

General Transfer Format - GTF

```
$ head ../../ensembl/Homo_sapiens.GRCh38.93.gtf|cut -f 1-8
#!genome-build GRCh38.p12
#!genome-version GRCh38
#!genome-date 2013-12
#!genome-build-accession NCBI:GCA_000001405.27
#!genebuild-last-updated 2018-01
1      havana  gene      11869      14409      .          +          .
1      havana  transcript      11869      14409      .          +          .
1      havana  exon       11869      12227      .          +          .
1      havana  exon       12613      12721      .          +          .
1      havana  exon       13221      14409      .          +          .
```

chr id	gene	feature	end location	strand	
source		start location	score	reading frame	attributes

```
$ head -8 ../../ensembl/Homo_sapiens.GRCh38.93.gtf|cut -f 9
#!genome-build GRCh38.p12
#!genome-version GRCh38
#!genome-date 2013-12
#!genome-build-accession NCBI:GCA_000001405.27
#!genebuild-last-updated 2018-01
gene_id "ENSG00000223972"; gene_version "5"; gene_name "DDX11L1"; gene_source "havana"; gene_biotype "transcribed_unprocessed_pseudogene";
gene_id "ENSG00000223972"; gene_version "5"; transcript_id "ENST00000456328"; transcript_version "2"; gene_name "DDX11L1"; gene_source "havana"; gene_biotype
"transcribed_unprocessed_pseudogene"; transcript_name "DDX11L1-202"; transcript_source "havana"; transcript_biotype "processed_transcript"; tag "basic";
transcript_support_level "1";
gene_id "ENSG00000223972"; gene_version "5"; transcript_id "ENST00000456328"; transcript_version "2"; exon_number "1"; gene_name "DDX11L1"; gene_source "havana";
```

- The purpose of the counting is to estimate the number of reads (counts) associated with each feature of interest (gene, exon, transcript).
- Input required: BAM/SAM files + General Transfer Format (GTF) file.
- Output generated: a gene expression matrix with genes as rows and samples as column.
- Top popular tools *htseq-count* (more cited, same developers as *DESeq*) and *featureCounts* (a way faster, Unix and R package, a bit more liberal).

aligned read:

start: 113217600 end: 113217650

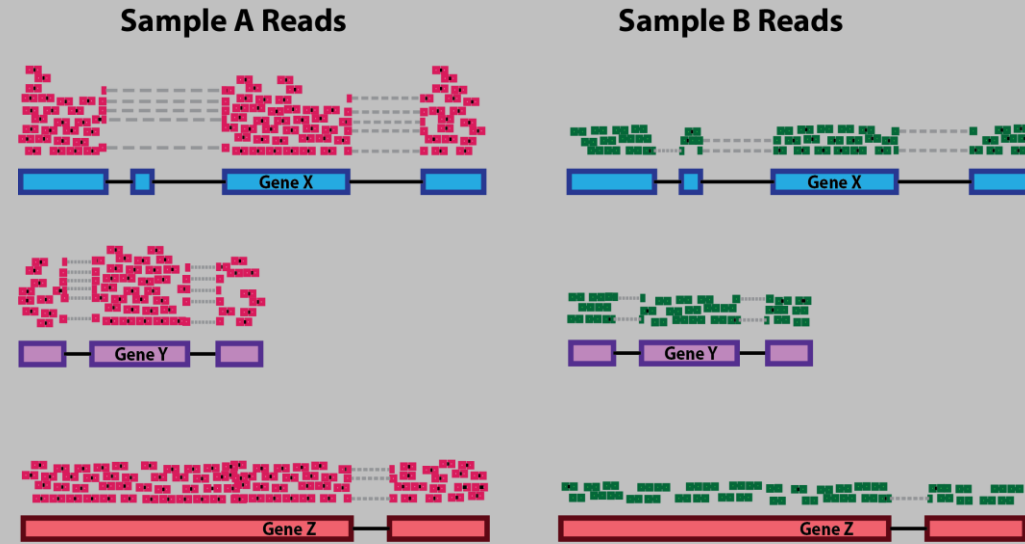


- Input: BAM file (alignment output, indexed) + genome sequences + genome annotation.



Counting issue: let's consider a gene A with 100 counts, a gene B with 100 counts as well but gene A is TWICE longer than gene B...

"Are these 2 genes expressed at the same level?"



2 biases to correct:

- library size or sequencing depth
- gene length

Gene name	Rep1 Count	Rep2 Count	Rep3 Count
A (2kb)	10	12	30
B (4kb)	20	25	60
C (1 kb)	5	8	15
D (10 kb)	0	0	1

Normalisation with TPM

Transcript Per Million

Gene name	Rep1 Count	Rep2 Count	Rep3 Count
A (2kb)	10	12	30
B (4kb)	20	25	60
C (1 kb)	5	8	15
D (10 kb)	0	0	1

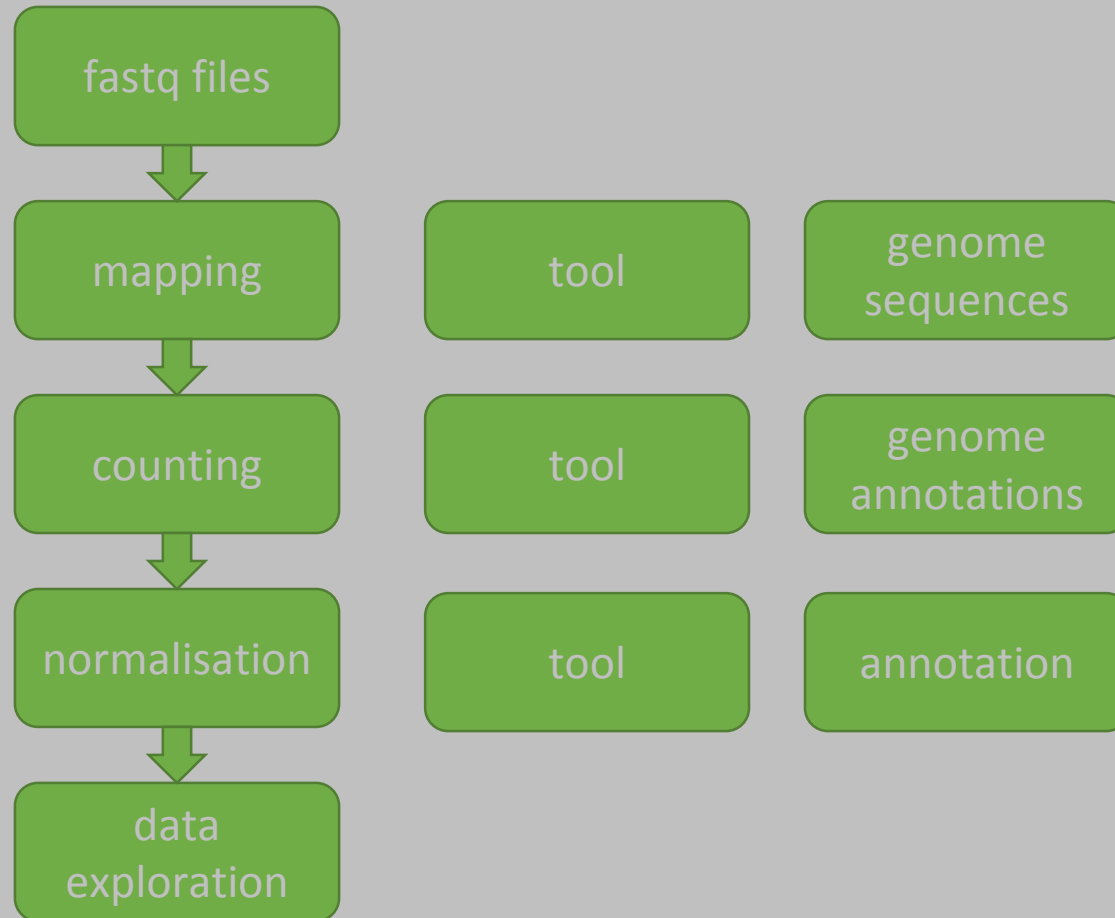
Step 1:
Normalisation in
gene length (divide
by gene length).

Gene name	Rep1 RPK	Rep2 RPK	Rep3 RPK
A (2kb)	5	6	15
B (4kb)	5	6.25	15
C (1 kb)	5	8	15
D (10 kb)	0	0	0.1

Total RPK:	15	20.25	45.1
Tens of RPK:	1.5	2.025	4.51

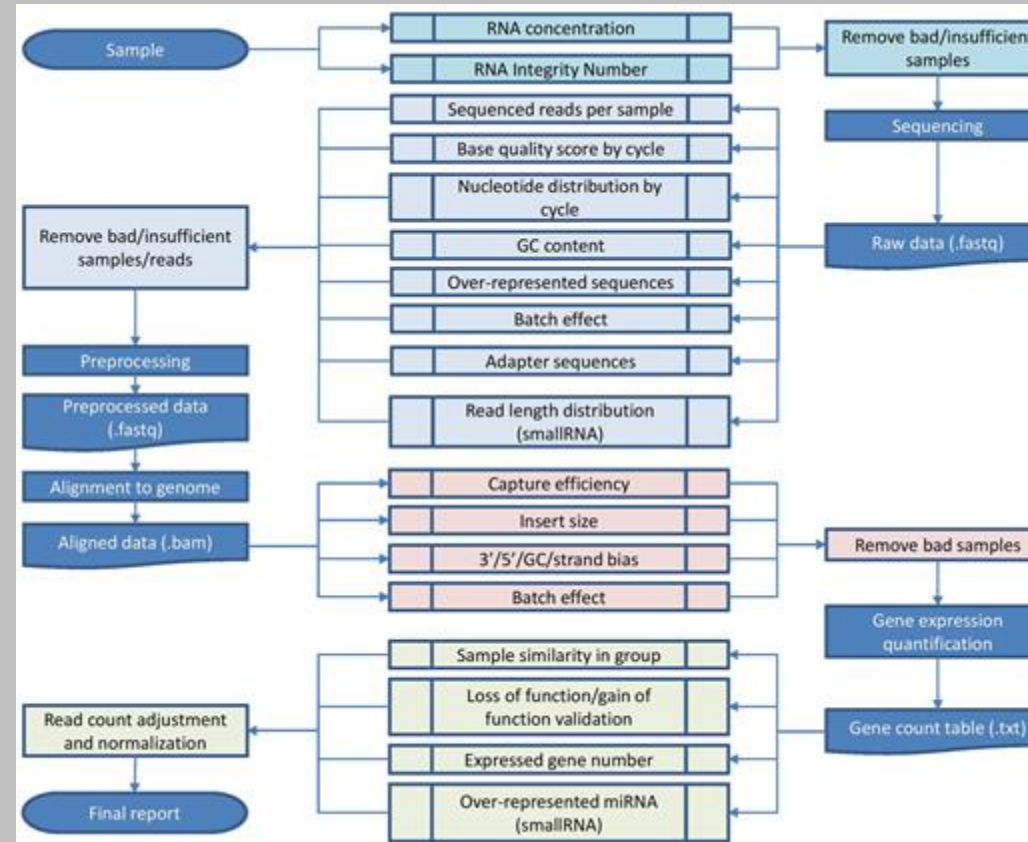
Step 2:
Normalisation in
sequencing depth
(divide by **scaling
factor**).

Gene name	Rep1 TPM	Rep2 TPM	Rep3 TPM
A (2kb)	3.33	2.96	3.326
B (4kb)	3.33	3.09	3.326
C (1 kb)	3.33	3.95	3.326
D (10 kb)	0	0	0.02



**Thank you for your
attention**

Quality Controls



Third gen sequencing

Principle of nanopore and single-molecule real-time (SMRT) sequencing

