

Dealing with 'raw reads'

Analysis of Next-Generation Sequencing Data

Friederike Dündar

Applied Bioinformatics Core

Slides at <https://bit.ly/2CUdS9z>¹

January 29, 2019

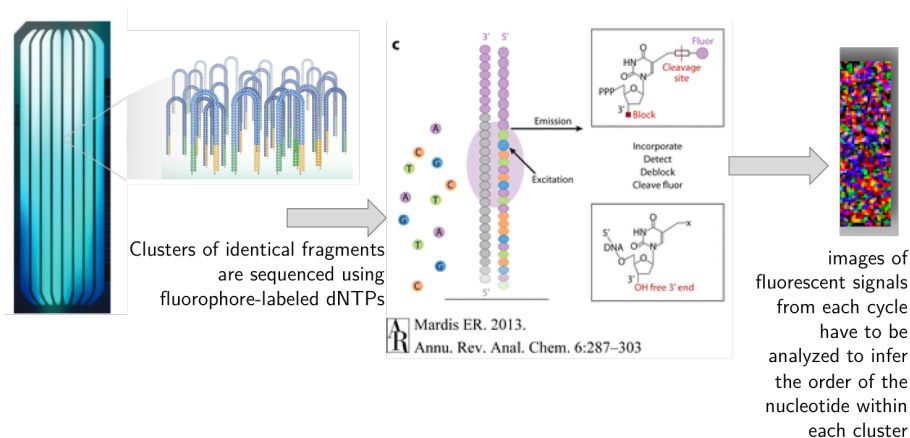


Weill Cornell Medicine

¹http://physiology.med.cornell.edu/faculty/skrabanek/lab/angsd/schedule_2018/

- 1 Paired-end reads
- 2 Illumina's "raw reads"
- 3 Quality control of sequencing reads
- 4 Sequence Read Archive
- 5 References

Re-cap: Sequencing by synthesis after library preparation



The number of **sequencing cycles**² determines the read **length**.

²(1) Incorporate fluor-dNTP, (2) detect, (3) deblock, (4) cleave fluor

Number of flowcell lanes determines the sequencing depth

Every read represents a cluster on the flowcell. The *lower* limit for the number of reads should follow ENCODE (<https://www.encodeproject.org/about/experiment-guidelines/>)

Application	Recommended seq. depth
differential gene expression	20 - 50 mio SR, 75 bp
variant calling	30-200x coverage
whole-genome bisulfite sequencing	30x coverage
ChIA-PET	200 mio PE

You may need more, longer, and possibly paired-end (PE) reads for:

- novel transcript identification, alternative splicing ³
- ChIP-seq for broad histone marks
- 3D chromatin structure assessment assays

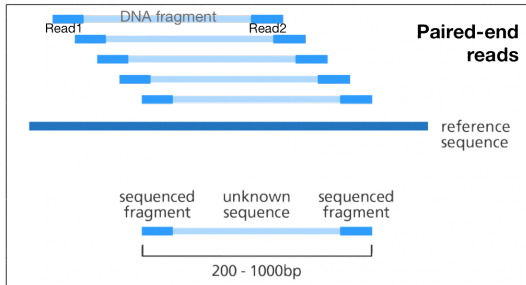
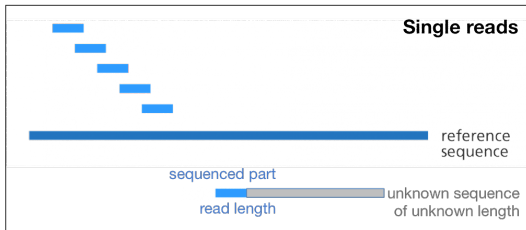
The addition of replicates may be more meaningful than increased sequencing depth!

³Most PIs that are serious about this will not use Illumina sequencing for this.

Paired-end reads

Types of reads

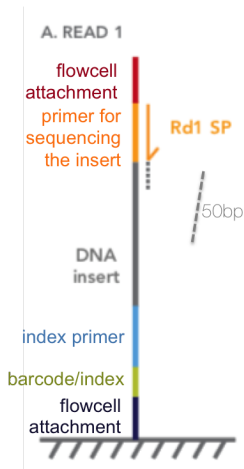
<https://www.yourgenome.org/facts/how-do-you-put-a-genome-back-together-after-sequencing>



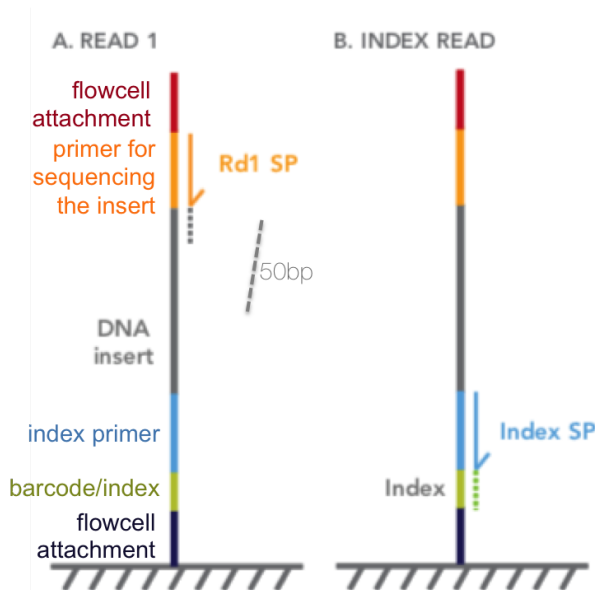
Paired-end (PE) reads are helpful for:

- **alignment** along repetitive regions
- chromosomal **rearrangements** and gene fusion detection
- *de novo* genome and transcriptome **assembly**
- precise information about the size of the original fragment (**insert size**)
- PCR duplicate identification

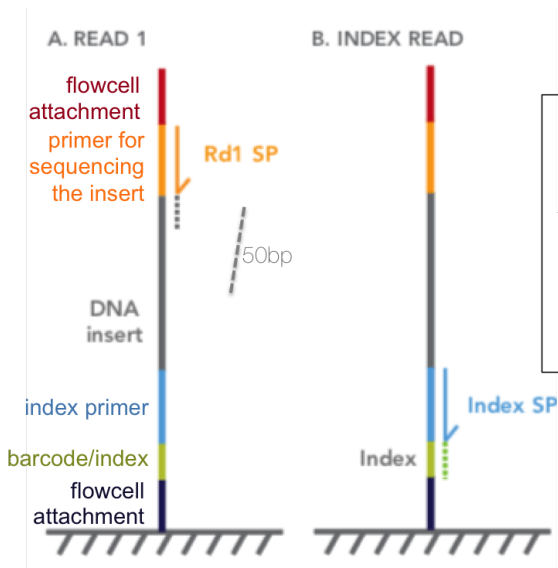
Paired-end read generation



Paired-end read generation



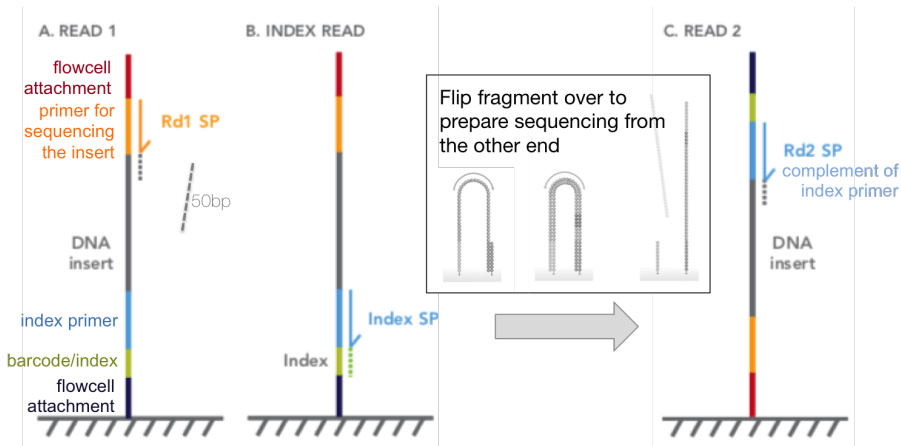
Paired-end read generation



Flip fragment over to prepare sequencing from the other end

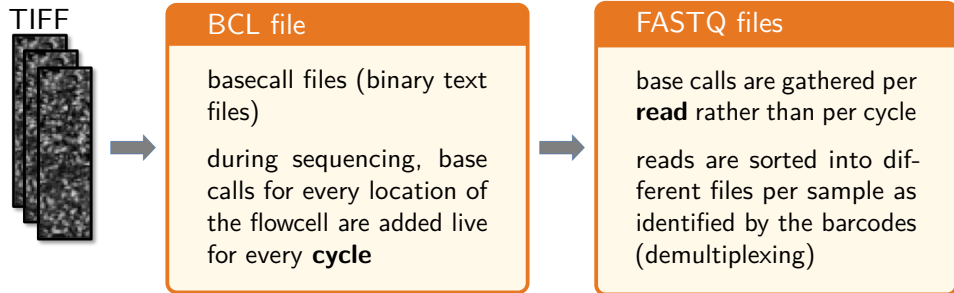


Paired-end read generation



Illumina's "raw reads"

Illumina's read output: turning images into text files



All steps here are performed by Illumina's proprietary CASAVA software. The file name usually includes some information about the sample:
`<sample name>_<barcode sequence>_<L(lane)>_<R(read number)>_<set number>.fastq.gz`, e.g. `MyExperiment_AGCTTGTTTC_L001_R1_001.fastq.gz`

The FASTQ format: FASTA + quality score

1 read = 4 lines

```
1 @ERR459145.1 DHKW5DQ1:219:D0PT7ACXX:2:1101:1590:2149/1
2 GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGATCGGAAGAGCGGTTCAGC
3 +
4 @7<DBADDDDBH?DHHI@DH>HHHEGHIIIGGIFFGIBFAAGAFHA '5?B@D
```

- ① @Read ID and sequencing run information
- ② sequence
- ③ + (additional description possible; usually an empty line)
- ④ quality scores

The read ID line is standardized by Casava 1.8

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
```

EAS139	the unique instrument name
136	the run id
FC706VJ	the flowcell id
2	flowcell lane
2104	tile number within the flowcell lane
15343	'x'-coordinate of the cluster within the tile
197393	'y'-coordinate of the cluster within the tile
1	the member of a pair, 1 or 2 (<i>paired-end or mate-pair reads only</i>)
Y	Y if the read is filtered, N otherwise
18	0 when none of the control bits are on, otherwise it is an even number
ATCACG	index sequence

CAUTION

This will only be true if you receive FASTQ files fresh off the sequencer. If you download FASTQ files from public repositories, the read ID might have been changed significantly.

see https://en.wikipedia.org/wiki/FASTQ_format

The quality scores: summarizing numerical scores into single-character representations

```
@ERR459145.1 DHKW5DQ1:219:D0PT7ACXX:2:1101:1590:2149/1
GATCGGAAGAGCGGTTCAGCAGGAATGCCGAGATCGGAAGAGCGGTTCAGC
+
@7<DBADDDDBH?DHHI@DH>HHHEGHHIIIGGIFFGIBFAAGAFHA'5?B@D
```



Illumina's CASAVA pipeline:

Base calls are immediately recorded with an error probability⁴ (BCL files), which are translated into ASCII symbols in the FASTQ files.

⁴See the QC section for reasons for base call uncertainties.

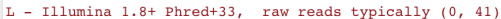
ASCII symbols

DEC	OCT	HEX	BIN	Symbol
32	040	20	00100000	
33	041	21	00100001	!
34	042	22	00100010	"
35	043	23	00100011	#
36	044	24	00100100	\$
37	045	25	00100101	%
38	046	26	00100110	&
39	047	27	00100111	'
40	050	28	00101000	(
41	051	29	00101001)
42	052	2A	00101010	*
43	053	2B	00101011	+
65	101	41	01000001	A
66	102	42	01000010	B
67	103	43	01000011	C
68	104	44	01000100	D
69	105	45	01000101	E
70	106	46	01000110	F
71	107	47	01000111	G
72	110	48	01001000	H

www.ascii-code.com

ASCII encodes 128 specified characters into seven-bit integers, which is useful for digital communication. The first 33

characters represent unprintable control codes (e.g. "Start of Text"), therefore the Phred scores were originally encoded by using an **offset of +33**.



F. Dünder (ABC, WCM)

Different offsets have been used by different Casava versions

```
SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS.....
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
.....IIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIIII.....
.....JJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJJ.....
LLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLLL.....
!"#$%&'()*+,-./0123456789:;<=?@ABCDEFGHIJKLMNOPQRSTUVWXYZ[\]^_`abcdefghijklmnopqrstuvwxyz{|}~
|                                     |               |
33                                59    64          73                               104                      126
0.....26...31.....40
      -5....0.....9.....40
              0.....9.....40
                  3.....9.....40
0.2.....26...31.....41

S - Sanger           Phred+33, raw reads typically (0, 40)
G - Solexa           Solexa+64, raw reads typically (-5, 40)
I - Illumina 1.3+   Phred+64, raw reads typically (0, 40)
J - Illumina 1.5+   Phred+64, raw reads typically (3, 40)
                    with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)
                    (Note: See discussion above).
L - Illumina 1.8+   Phred+33, raw reads typically (0, 41)
```

image from https://en.wikipedia.org/wiki/FASTQ_format

Different offsets have been used by different Casava versions

Both the **range of the base call score** as well as its translation via the ASCII code (**offset**) are somewhat arbitrary and have undergone numerous changes.

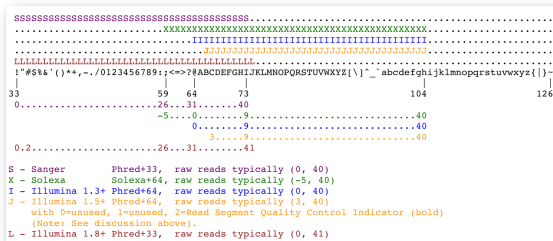


image from https://en.wikipedia.org/wiki/FASTQ_format

Today's standard:

- min. score: 0
- max. score: 41
- ASCII offset: 33

Make sure you know which version you're dealing with.

Quality control of sequencing reads

Two basic QC questions

- ① Did our **library prep** generate a **faithful representation** of the DNA/RNA molecules our our samples?
 - ▶ ideally, the entire universe of nucleotides was captured (diverse library)
 - ▶ no contaminations
 - ▶ no degradation
 - ▶ no bias towards fragments of certain GC contents and/or sizes
- ② How successful was the actual **sequencing**?
 - ▶ consistently high base call confidence
 - ▶ uniform nucleotide frequencies

Biases

QC should help identify **systematic distortions** of data and their possible sources.

FastQC

<http://www.bioinformatics.babraham.ac.uk/projects/fastqc>

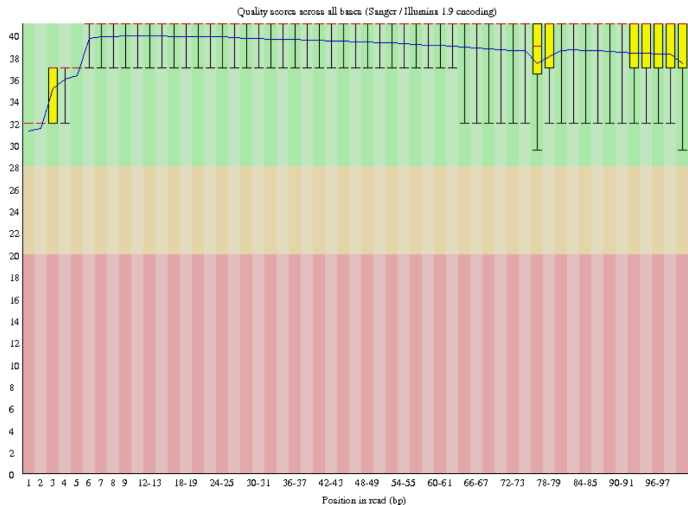
- unpublished, but most widely used QC tool
- supports all NGS technologies
- continuously developed and maintained by long-time bioinformatics experts
- will only use the first 200K reads for the diagnosis!

Function	A quality control tool for high throughput sequence data.
Language	Java
Requirements	A suitable Java Runtime Environment The Picard BAM/SAM Libraries (included in download)
Code Maturity	Stable. Mature code, but feedback is appreciated.
Code Released	Yes, under GPL v3 or later .
Initial Contact	Simon Andrews

Sequencing quality

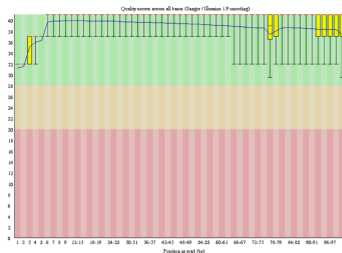
Based on ASCII-encoded Phred scores within the fastq file.

✔ Per base sequence quality

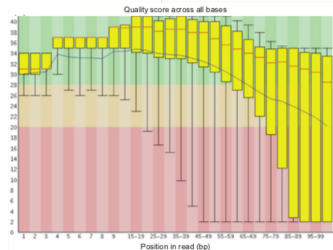


Sequencing quality

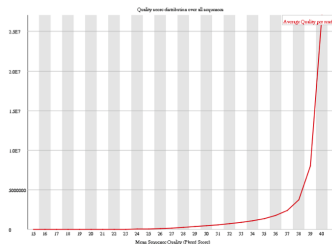
✓ Per base sequence quality



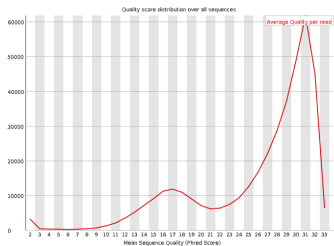
✗ Per base sequence quality



✓ Per sequence quality scores



✓ Per sequence quality scores

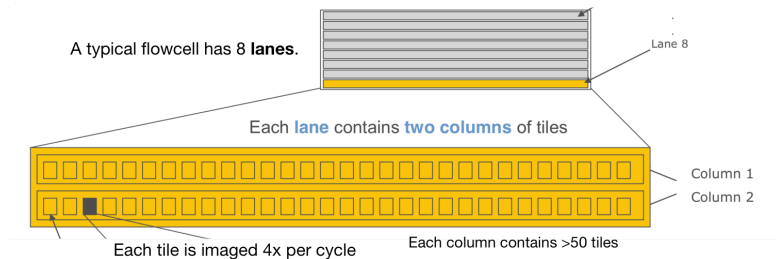


Sequencing quality: reasons for sequencing noise

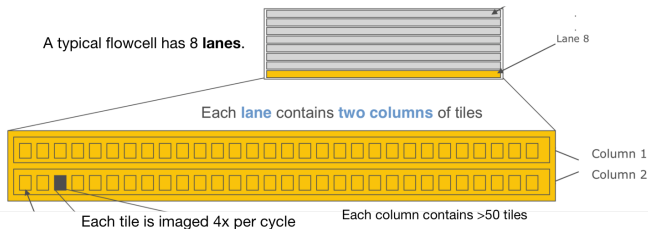
Noise = fluorophore intensity signal is not as strong and clear as expected.

- laser not well **calibrated**
- **interfering signals** from neighbouring clusters or bases with similar emission spectra
- **unsynchronized fragments** in each cluster:
 - ▶ *phasing*: small fraction of fragments in each cluster fails to incorporate any base
 - ▶ *prephasing*: more than one base is incorporated
- **decaying** chemicals (runs often last several days to a week!)
- **extraneous objects** on the flow cell (e.g. dust, air bubbles)

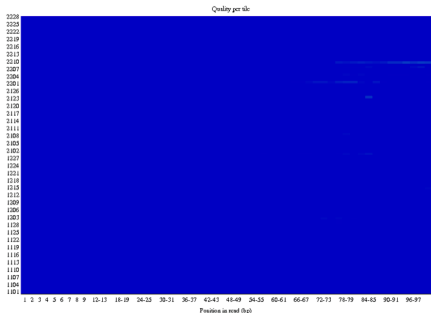
Physically localized error rates: tiles vs. time



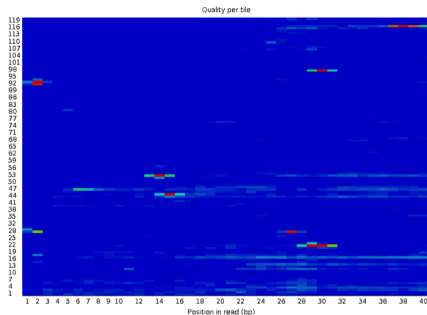
Physically localized error rates: tiles vs. time



✓ Per tile sequence quality



✗ Per tile sequence quality



Contaminations: threats to the full representation of our original fragment pool

- Sources:

- ▶ **primer** contamination
- ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
- ▶ DNA from other species/libraries

- Consequences:

- ▶ noise
- ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Contaminations: threats to the full representation of our original fragment pool

- Sources:

- ▶ **primer** contamination
- ▶ **adapter** contamination

- sequence read length larger than the fragment size (3' contamination)
- adapter dimers without insert

- ▶ DNA from other species/libraries

- Consequences:

- ▶ noise
- ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Contaminations: threats to the full representation of our original fragment pool

- Sources:

- ▶ **primer** contamination
- ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
- ▶ DNA from other species/libraries

- Consequences:

- ▶ noise
- ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Contaminations: threats to the full representation of our original fragment pool

- Sources:
 - ▶ **primer** contamination
 - ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
 - ▶ DNA from other species/libraries
- Consequences:
 - ▶ noise
 - ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Contaminations: threats to the full representation of our original fragment pool

- Sources:
 - ▶ **primer** contamination
 - ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
 - ▶ DNA from other species/libraries
- Consequences:
 - ▶ noise
 - ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Contaminations: threats to the full representation of our original fragment pool

- Sources:
 - ▶ **primer** contamination
 - ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
 - ▶ DNA from other species/libraries
- Consequences:
 - ▶ noise
 - ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Contaminations: threats to the full representation of our original fragment pool

- Sources:
 - ▶ **primer** contamination
 - ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
 - ▶ DNA from other species/libraries
- Consequences:
 - ▶ noise
 - ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Contaminations: threats to the full representation of our original fragment pool

- Sources:
 - ▶ **primer** contamination
 - ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
 - ▶ DNA from other species/libraries
- Consequences:
 - ▶ noise
 - ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

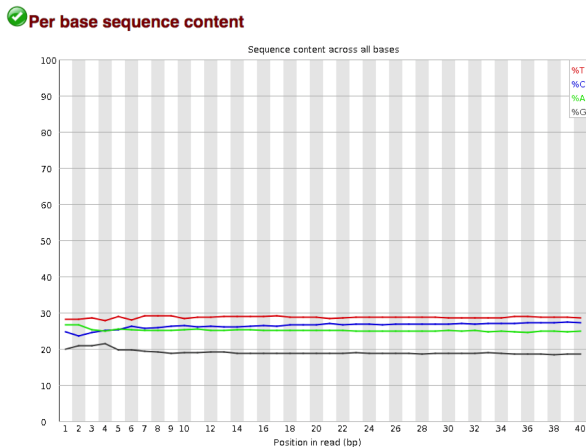
Contaminations: threats to the full representation of our original fragment pool

- Sources:
 - ▶ **primer** contamination
 - ▶ **adapter** contamination
 - sequence read length larger than the fragment size (3' contamination)
 - adapter dimers without insert
 - ▶ DNA from other species/libraries
- Consequences:
 - ▶ noise
 - ▶ reduced alignment rates

Can be identified by examining **sequence composition** and **overrepresented sequences/k-mers**.

Detecting contaminations

Per Base Sequence Content

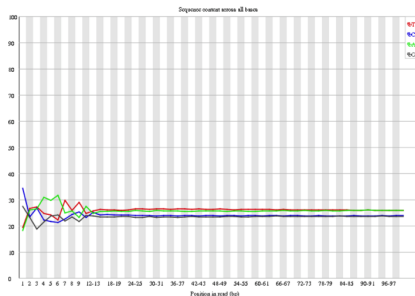


If the fragments represent a random and diverse representation of the entire genome, there should be a uniform distribution of all four bases across all cycles.

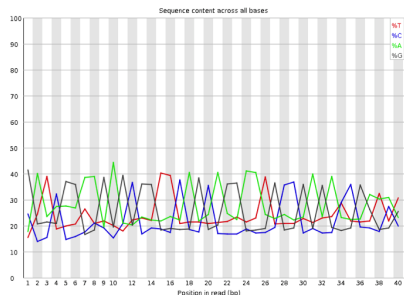
Detecting contaminations

Per Base Sequence Content – more examples

✓ Per base sequence content



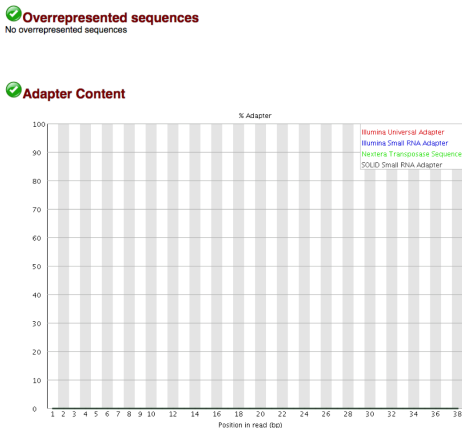
✗ Per base sequence content



- irregularities in the first ca. 8 bp are often seen for RNA-seq and ATAC-seq and indicate a bias for certain sequences at the fragment beginning
- more severe deviations from uniformity often indicate contaminations and/or lack of library diversity

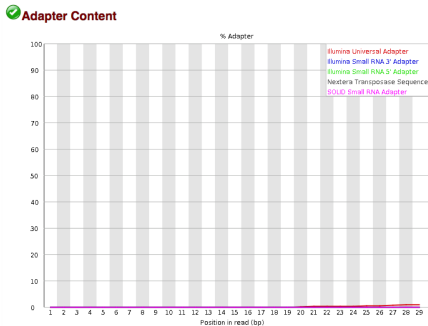
Detecting contaminations

Overrepresented sequences & adapter sequence frequencies



Overrepresented sequences

Sequence	Count	Percentage	Possible Source
AGAGTTTATGCTTCATGACGZGAGCTTACACTTTC	2045	8.5224039181558763	No hit
GATTGGCTATCCAACTGACAGATTTTATGCTTCATG	2047	8.5178502762542754	No hit
ATTGGCTATCCAACTGACAGATTTTATGCTTCATG	2014	8.5095019327480071	No hit
CGATAAAATGATTGCTATCCAACTGACAGATTTTAT	1913	8.483950420979134	No hit
GTTATCAACTGACAGATTTTATGCTTCATGACGAGA	1879	8.4703496165660066	No hit
AAAAATGATTGGCTATCCAACTGACAGATTTTATGCT	1846	8.4678012750197325	No hit



Trimming contaminations & low-quality bases

- Can be done before alignment or, if contaminations/low-quality bases are low in number, might be left to the “soft-clipping” function⁵ of read aligners.
- There are numerous tools out there to do the job, e.g. Cutadapt and TrimGalore.
- For *de novo* assemblies, it is probably more meaningful to perform some error-correction based on overlapping reads rather than trimming the reads [Salzberg et al., 2012],[Yang et al., 2013]

⁵ignoring mis-matched bases at the beginning/end of a read

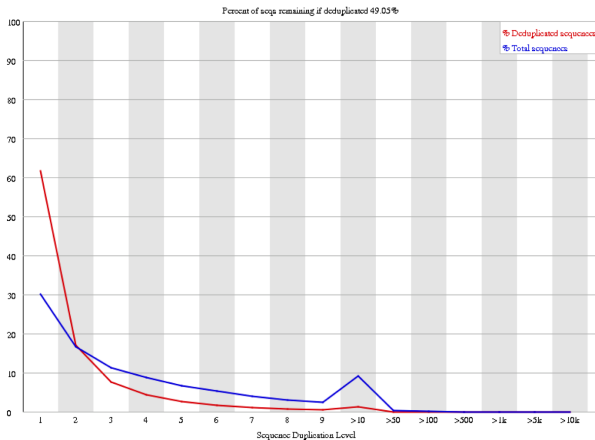


- 33 / 44

Duplicate reads: FastQC assessment

Proportion of reads (y-axis) that contain sequences in each of the different duplication level bins (x-axis).

Sequence Duplication Levels

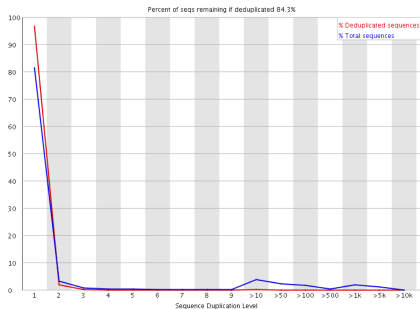


Blue line: all reads (= first 100K!) – how many times are individual sequences found?

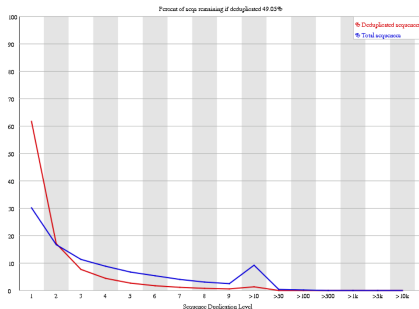
Red line: sequences after de-duplication – how many *different* sequences were found to be duplicated?

Duplicate reads: FastQC assessment

✓ Sequence Duplication Levels



✗ Sequence Duplication Levels



Check that the red line is flat and that the number of remaining reads after de-duplication is acceptable.

Two basic QC questions

- ① Did our **library prep** generate a **faithful representation** of the DNA/RNA molecules our our samples?
 - ▶ ideally, the entire universe of nucleotides was captured (diverse library)
 - ▶ no contaminations
 - ▶ no bias towards fragments of certain GC contents and/or sizes
 - ▶ no degradation
- ② How successful was the actual **sequencing**?
 - ▶ consistently high base call confidence
 - ▶ uniform nucleotide frequencies

QC summary

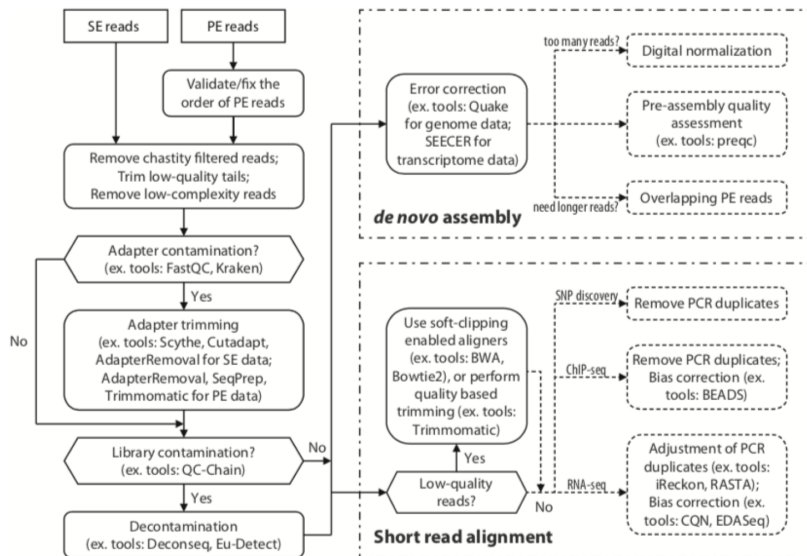
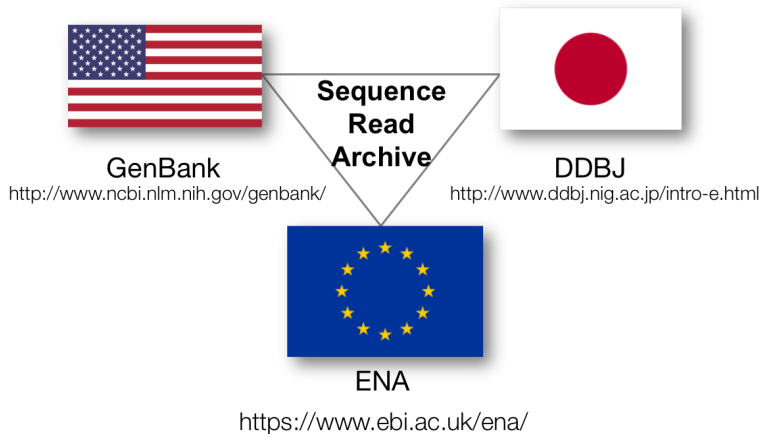


Figure from Zhou and Rokas [2014] (highly recommended reading!)

Sequence Read Archive

Where are all the reads?

The SRA is the main repository for publicly available DNA and RNA sequencing data of which three instances are maintained world-wide. GEO (<https://www.ncbi.nlm.nih.gov/geo/>) can be used to find SRA data, too.



See O'Sullivan et al. [2017] for many more details.

References

[Mardis, 2013, Illumina Inc, 2015, 2008, 2013, Andrews]

References

- Simon Andrews. FastQC Analysis Module Documentation. URL <https://www.bioinformatics.babraham.ac.uk/projects/fastqc/Help/3%20Analysis%20Modules/>.
- Illumina Inc. Multiplexed Sequencing with the Illumina Genome Analyzer System. In *Illumina Sequencing*. 2008. URL https://www.illumina.com/documents/products/datasheets/datasheet_sequencing_multiplex.pdf.
- Illumina Inc. bcl2fastq Conversion User Guide. *Illumina*, Version 1.(March), 2013. URL https://support.illumina.com/downloads/bcl2fastq_conversion_software_184.html.
- Illumina Inc. Patterned Flow Cell Technology. In *Technical Spotlight: Sequencing*, pages 1–2. 2015. URL <https://www.illumina.com/content/dam/illumina-marketing/documents/products/technotes/patterned-flow-cell-technology-technical-note-770-2015-010.pdf>.
- Elaine R. Mardis. Next-Generation Sequencing Platforms. *Annual Review of Analytical Chemistry*, 2013. doi: 10.1146/annurev-anchem-062012-092628.

- Christopher O'Sullivan, Benjamin Busby, and Ilene Karsch Mizrahi. In Jonathan M. Keith, editor, *Bioinformatics: Volume I: Data, Sequence Analysis, and Evolution*, chapter Managing Sequence Data. Humana Press, 2017. doi: 10.1007/978-1-4939-6622-6_4.
- Steven L. Salzberg, Adam M. Phillippy, Aleksey Zimin, Daniela Puiu, Tanja Magoc, Sergey Koren, Todd J. Treangen, Michael C. Schatz, Arthur L. Delcher, Michael Roberts, Guillaume Marcxais, Mihai Pop, and James A. Yorke. GAGE: A critical evaluation of genome assemblies and assembly algorithms. *Genome Research*, 2012. doi: 10.1101/gr.131383.111.
- Xiao Yang, Sriram P. Chockalingam, and Srinivas Aluru. A survey of error-correction methods for next-generation sequencing. *Briefings in Bioinformatics*, 2013. doi: 10.1093/bib/bbs015.
- Xiaofan Zhou and Antonis Rokas. Prevention, diagnosis and treatment of high-throughput sequencing data pathologies. *Molecular Ecology*, 23(7): 1679–1700, 2014. doi: 10.1111/mec.12680.