

NGS Characteristics

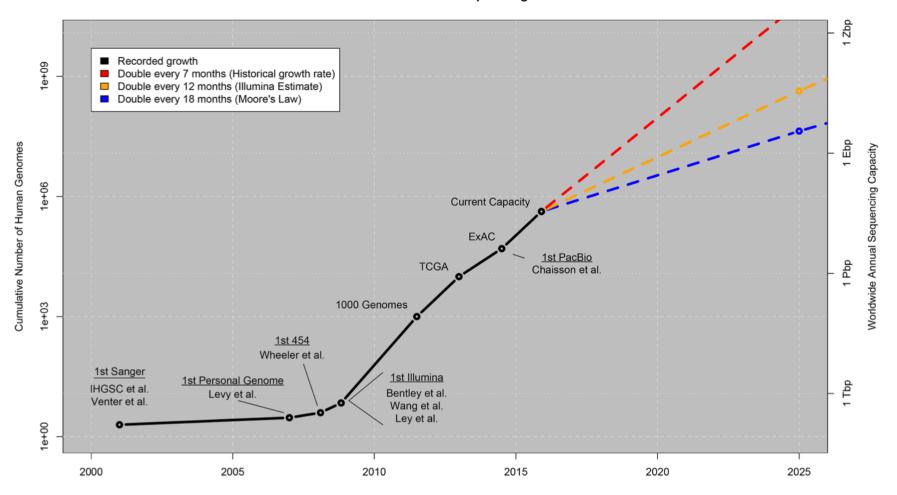
Hubert Rehrauer





NGS Data Increase

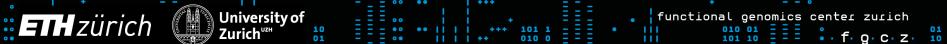
Growth of DNA Sequencing



NGS data increases faster than computer speed

functional genomics center zurich

Q · C · Z ·



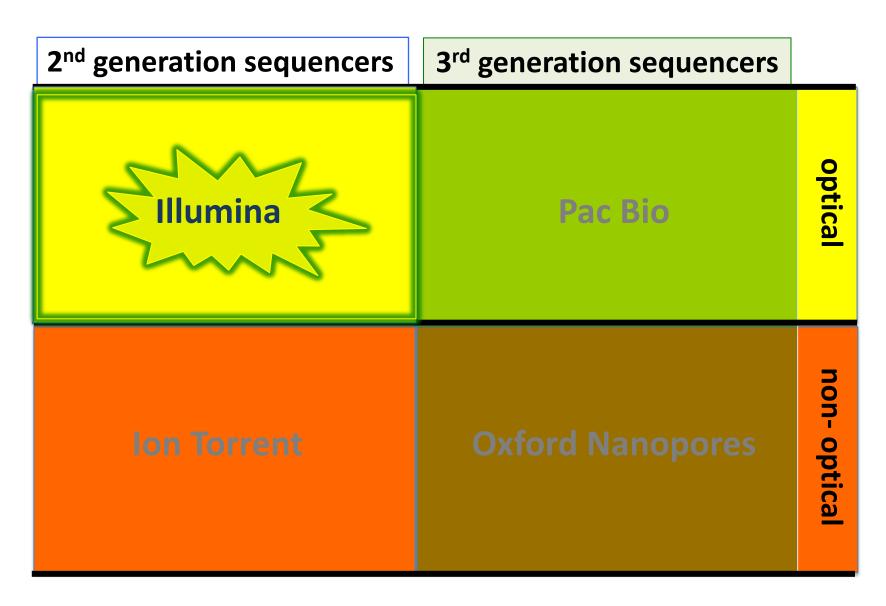
Ingredients for the success

- Evolution has yielded DNA and RNA molecules for information storage and transfer. They have good properties to be read (measured)
- NGS technologies rely on
 - massive parallelization
 - measurement process is done by individual molecules (cheap and fast)

Sequencers

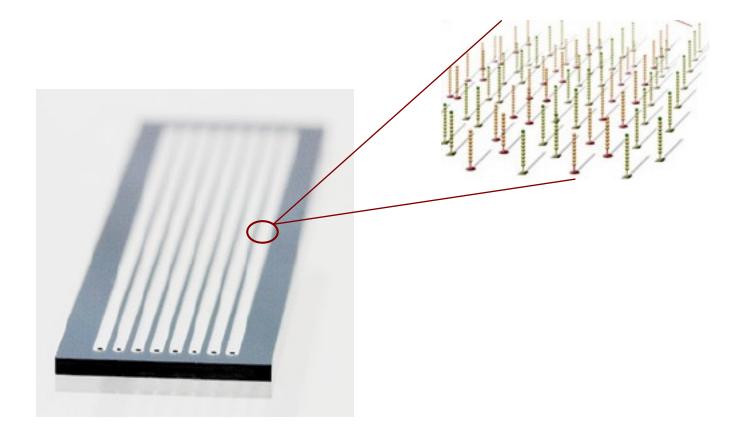
- short-read sequencers (up to ~ 600bp)
 - illumina
- long-read, single molecule technologies (500bp megabases)
 - PacBio
 - Oxford NanoPore
- example overview:
 - https://genohub.com/ngs-instrument-guide/







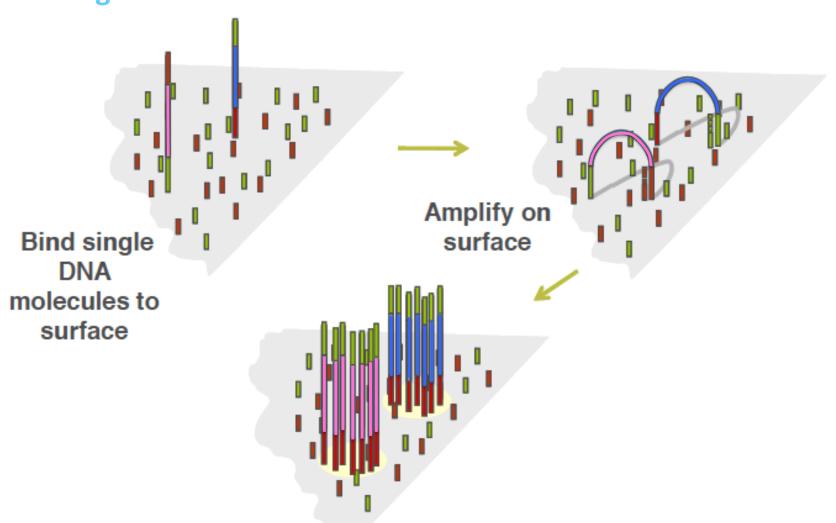
Illumina Flow cell

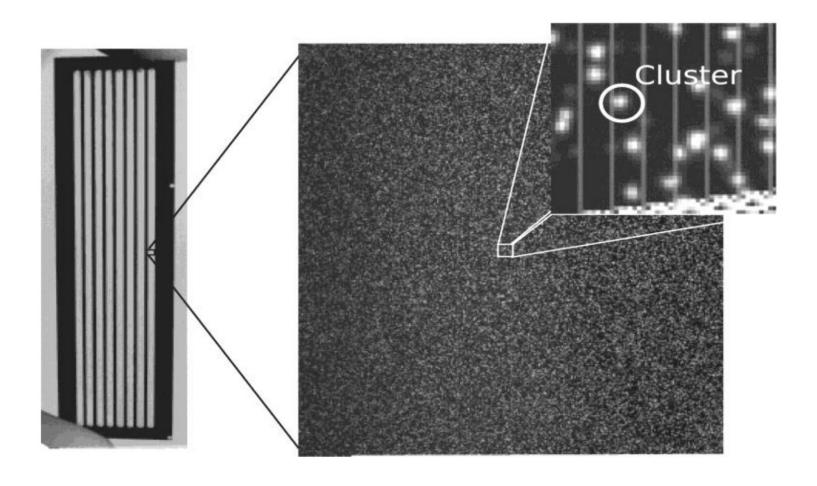




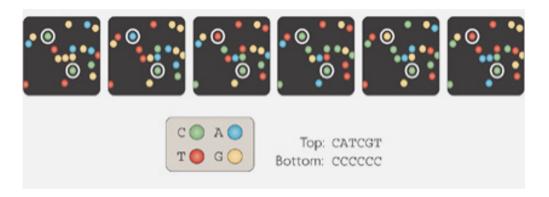


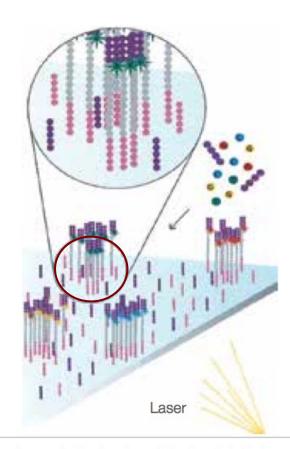
Cluster generation overview





Illumina Sequencing

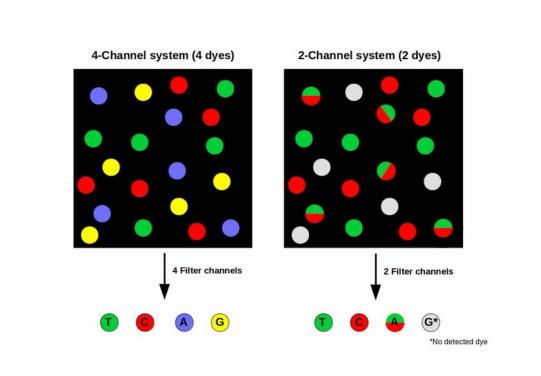


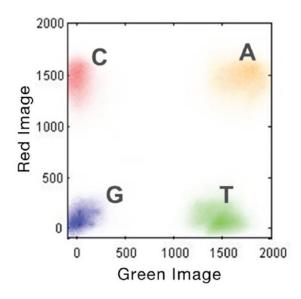


The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase.



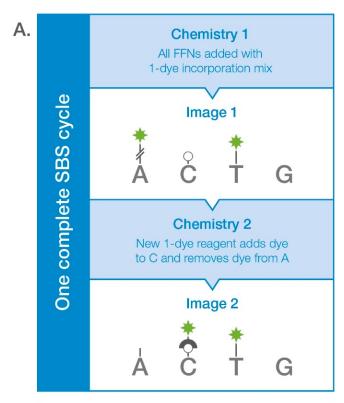
Color coding of bases







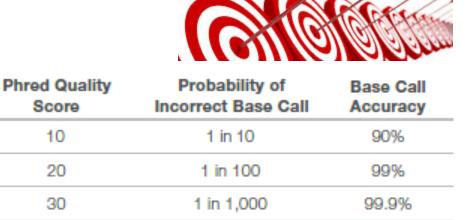
Color coding of illumina iSeq



B.	Image 1	Image 2	Result
	ON	OFF	А
	OFF	ON	С
	ON	ON	Т
	OFF	OFF	G

Phred scores measure base call accuracy

- P
- error probability of a given base call
- Q
- $-10log_{10}P$
- Assign to each base
- Range from 0-41



99.99%

99.999%

1 in 10,000

1 in 100,000

Ewing B, Green P. 1998. Genome Res. 8(3):186-194. http://en.wikipedia.org/wiki/Phred_quality_score

Score

10

20

30

40

50

Phred scores are stored with sequences

- FASTQ
 - 4 lines:
 - 1. Header line for Read (starts with "@" and the sequence ID)
 - 2. Sequence
 - Header line for Qualities (starts with "+")
 - 4. Quality score (represented in ASCII format)

Phred scores can be ASCII encoded

- Add an offset and convert the sum to ASCII
- Current format
 - Illumina 1.9 (i.e. Sanger format)
 - Phred scoring: 0-41;
 - Offset: 33
 - 41+33=74 (J)
 - All current sequencers

10 01 101

Dec	Нх	Oct	Cha	r	Dec	Нх	Oct	Html	Chr	Dec	Нх	Oct	Html	Chr	Dec	Нх	Oct	Html CI	hr_
0	0	000	NUL	(null)	32	20	040	@#32 ;	Space	64	40	100	a#64;	0	96	60	140	& # 96;	8
1	1	001	SOH	(start of heading)	33	21	041	@#33;	!	65	41	101	@#65;	A	97	61	141	a	a
2	2	002	STX	(start of text)	34	22	042	@#3 4 ;	rr	66	42	102	B	В	98	62	142	& # 98;	b
3	3	003	ETX	(end of text)	35	23	043	#	#	67	43	103	C	С	99	63	143	c	C
4	4	004	EOT	(end of transmission)	36	24	044	@#36;	ş	68	44	104	4#68;	D	100	64	144	d	d
5	5	005	ENQ	(enquiry)	37	25	045	@#37;	*				4#69;		101	65	145	e	е
6	6	006	ACK	(acknowledge)				&		70	46	106	a#70;	F	102	66	146	f	f
7	7	007	BEL	(bell)	39	27	047	'	1	71	47	107	G	G		_		g	
8	8	010	BS	(backspace)	40	28	050	&# 4 0;	(72	48	110	6#72;	H	104	68	150	h	h
9	9	011	TAB	(horizontal tab))					6#73;					i	
10	A	012	LF	(NL line feed, new line)				@# 4 2;					6#74;					j	
11	В	013	VT	(vertical tab)				&#43;</td><td></td><td></td><td></td><td></td><td><u>475;</u></td><td></td><td></td><td></td><td></td><td>k</td><td></td></tr><tr><td>12</td><td>С</td><td>014</td><td>FF</td><td>(NP form feed, new page)</td><td></td><td></td><td></td><td>,</td><td></td><td></td><td></td><td></td><td>L</td><td></td><td></td><td></td><td></td><td>l</td><td></td></tr><tr><td>13</td><td></td><td>015</td><td></td><td>(carriage return)</td><td></td><td></td><td></td><td>&#45;</td><td></td><td></td><td></td><td></td><td>a#77;</td><td></td><td></td><td></td><td></td><td>m</td><td></td></tr><tr><td>14</td><td></td><td>016</td><td></td><td>(shift out)</td><td></td><td></td><td></td><td>&#46;</td><td></td><td></td><td></td><td></td><td>a#78;</td><td></td><td></td><td></td><td></td><td>n</td><td></td></tr><tr><td>15</td><td></td><td>017</td><td></td><td>(shift in)</td><td></td><td></td><td></td><td>6#47;</td><td>-</td><td></td><td></td><td></td><td>6#79;</td><td></td><td></td><td></td><td></td><td>o</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(data link escape)</td><td></td><td></td><td></td><td>a#48;</td><td></td><td></td><td></td><td></td><td>%#80;</td><td></td><td></td><td></td><td></td><td>p</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(device control 1)</td><td></td><td></td><td></td><td>&#49;</td><td></td><td>ı</td><td></td><td></td><td>Q</td><td></td><td></td><td></td><td></td><td>q</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(device control 2)</td><td></td><td></td><td></td><td>2</td><td></td><td></td><td></td><td></td><td>R</td><td></td><td></td><td></td><td></td><td>r</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(device control 3)</td><td></td><td></td><td></td><td>3</td><td></td><td></td><td></td><td></td><td>S</td><td></td><td></td><td></td><td></td><td>s</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(device control 4)</td><td></td><td></td><td></td><td>4</td><td></td><td> </td><td></td><td></td><td>4;</td><td></td><td></td><td></td><td></td><td>t</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(negative acknowledge)</td><td></td><td></td><td></td><td>6#53;</td><td></td><td></td><td></td><td></td><td>6#85;</td><td></td><td></td><td></td><td></td><td>u</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(synchronous idle)</td><td>I</td><td></td><td></td><td><u>@</u>#54;</td><td></td><td></td><td></td><td></td><td>486;</td><td></td><td></td><td></td><td></td><td>v</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(end of trans. block)</td><td></td><td></td><td></td><td>a#55;</td><td></td><td></td><td></td><td></td><td>a#87;</td><td></td><td></td><td></td><td></td><td>w</td><td></td></tr><tr><td></td><td></td><td></td><td></td><td>(cancel)</td><td></td><td></td><td></td><td>8</td><td></td><td></td><td></td><td></td><td>4#88;</td><td></td><td></td><td>-</td><td></td><td>x</td><td></td></tr><tr><td></td><td></td><td>031</td><td></td><td>(end of medium)</td><td></td><td></td><td></td><td>9;</td><td></td><td></td><td></td><td></td><td>4#89;</td><td></td><td></td><td></td><td></td><td>y</td><td></td></tr><tr><td></td><td></td><td>032</td><td></td><td>(substitute)</td><td></td><td></td><td></td><td>:</td><td></td><td></td><td></td><td></td><td>6#90;</td><td></td><td></td><td></td><td></td><td>z</td><td></td></tr><tr><td></td><td></td><td></td><td>ESC</td><td>(escape)</td><td></td><td></td><td></td><td><u>4</u>,59;</td><td></td><td></td><td></td><td></td><td>6#91;</td><td>_</td><td></td><td></td><td></td><td>{</td><td></td></tr><tr><td></td><td></td><td>034</td><td></td><td>(file separator)</td><td></td><td></td><td></td><td><u>@#60;</u></td><td></td><td></td><td></td><td></td><td>6#92;</td><td></td><td></td><td></td><td></td><td>4;</td><td></td></tr><tr><td></td><td></td><td>035</td><td></td><td>(group separator)</td><td></td><td></td><td></td><td>=</td><td></td><td></td><td></td><td></td><td>a#93;</td><td>-</td><td></td><td></td><td></td><td>}</td><td></td></tr><tr><td></td><td></td><td>036</td><td></td><td>(record separator)</td><td></td><td></td><td></td><td>></td><td></td><td></td><td></td><td></td><td>a#94;</td><td></td><td></td><td></td><td></td><td>a#126;</td><td></td></tr><tr><td>31</td><td>lF</td><td>037</td><td>US</td><td>(unit separator)</td><td>63</td><td>3F</td><td>077</td><td>?</td><td>2</td><td>95</td><td>5F</td><td>137</td><td>a#95;</td><td>_</td><td>127</td><td>7F</td><td>177</td><td></td><td>DEL</td></tr></tbody></table>											

Source: www.LookupTables.com

Read Quality Control

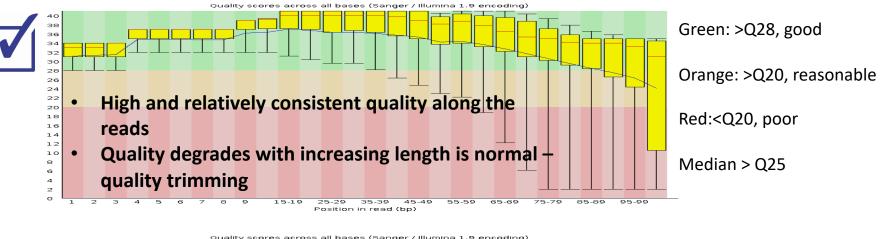
- Library construction could introduce bias
 - Fragmentation, ligation, amplification
 - GC bias
 - Over-amplification
 - Contamination

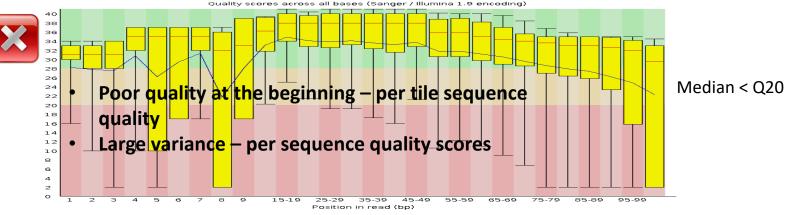
- Sequencing errors
 - Chemical, optical, computational

Platform	Primary error	Error rate (%)
Illumina	Substitution	0.1
PacBio	Indel	12 (consensus: 1)
Oxford Nanopore	Indel	3 - 20

Per base sequence quality - FastQC

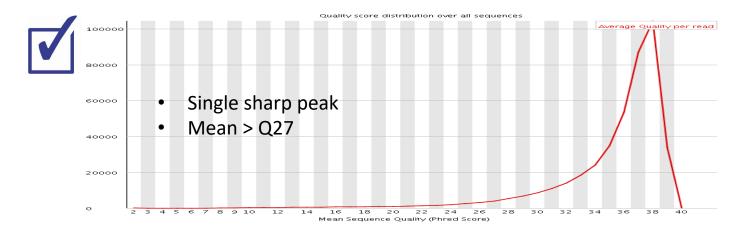
Range of quality values across all bases at each position

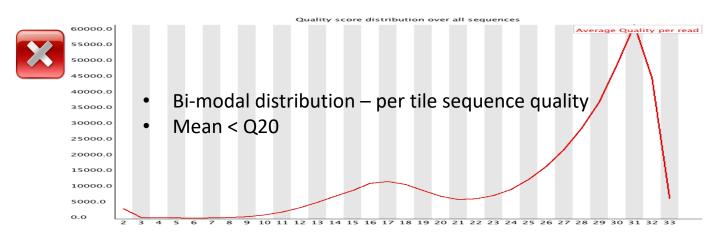




Per sequence quality scores - FastQC

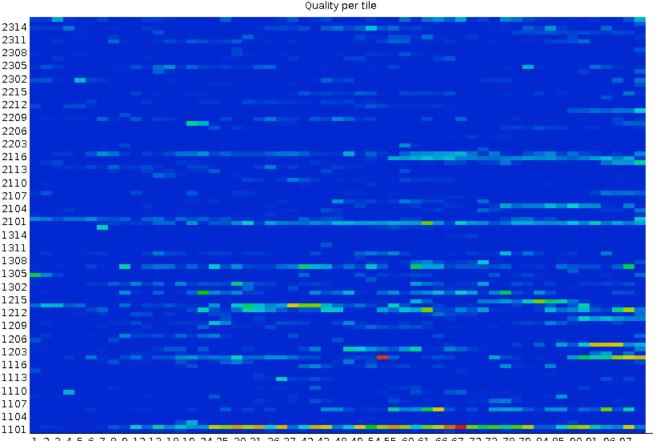
Subset of sequences with universally low quality values





Per tile sequence quality - FastQC

 Quality scores from each tile across all bases - loss in quality associated with only one part of the flowcell



Deviation from average quality

Cold colors: ≥ average

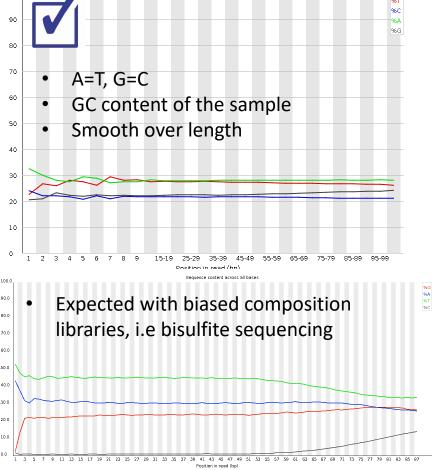
Hotter color: worse quality

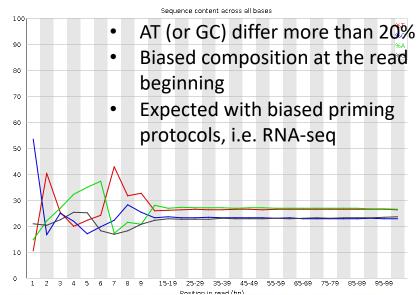
Good: universal blue

Failure: < average - 5

Per base sequence content - FastQC

The portion of A, T, G, and C at each position





Biases in Illumina transcriptome sequencing caused by random hexamer priming

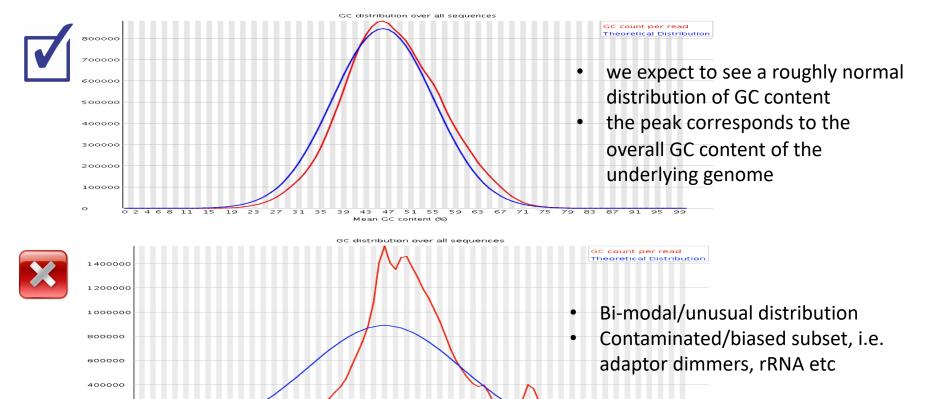
Kasper D. Hansen1,*, Steven E. Brenner2 and Sandrine Dudoit1,3

Treatment of DNA with bisulfite converts cytosine to uracil, but leaves methylated cytosine unaffected. Therefore, DNA that has been treated with bisulfite retains only methylated cytosines.

Per sequence GC content - FastQC

200000

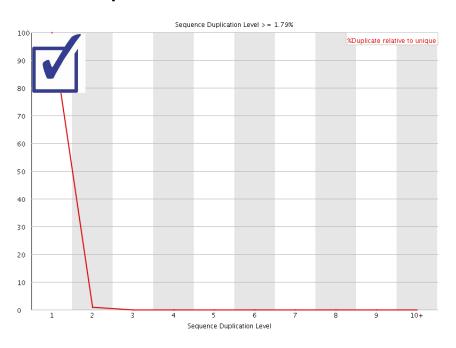
Distribution of average GC in all reads

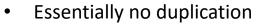


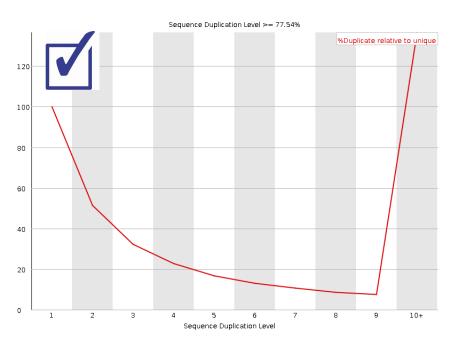
Mean GC content (%)

Sequence duplication - FastQC

 Relative number of sequences with different degrees of duplication







High duplication levels:

- DNA-seq: PCR over amplification, too little input material
- Normal in RNA-seq: high expression

Overrepresented sequences - FastQC

- Sequences make up >0.1 % of the total
- Compare those with a contamination database for finding contamination (i.e. adaptor dimmers)

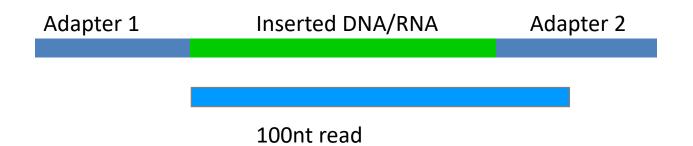
Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GGAAGAGCACACGTCTGAACTCCAGTCACCAGATCATCTCGTATGCCGTC	75874	1.5613887498682963	TruSeq Adapter, Index 7 (100% over 50bp)
GGAAGAGCACACGTCTGAACTCCAGTCACCGATGTATCTCGTATGCCGTC	7636	0.15713900010536297	TruSeq Adapter, Index 2 (100% over 50bp)
GGAAGAGCACACGTCTGAACTCCAGTCACACAGTGATCTCGTATGCCGTC	7539	0.1551428656095248	TruSeq Adapter, Index 5 (100% over 50bp)
GGAAGAGCACACGTCTGAACTCCAGTCACGCCAATATCTCGTATGCCGTC	5117	0.10530123933199874	TruSeq Adapter, Index 6 (100% over 50bp)

- Can be normal and biologically meaningful
 - highly expressed transcripts
 - high copy number repeats
 - Less diverse library (amplicons)

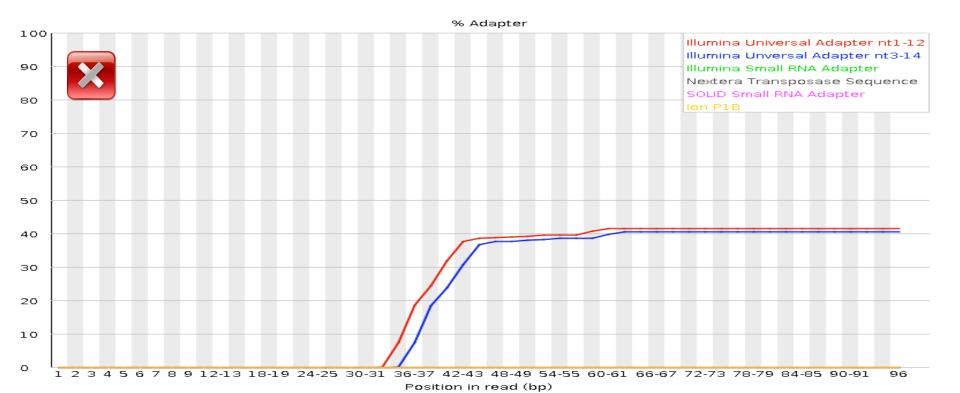


Adapter Content - FastQC



functional genomics center zurich

f. g. c. z.



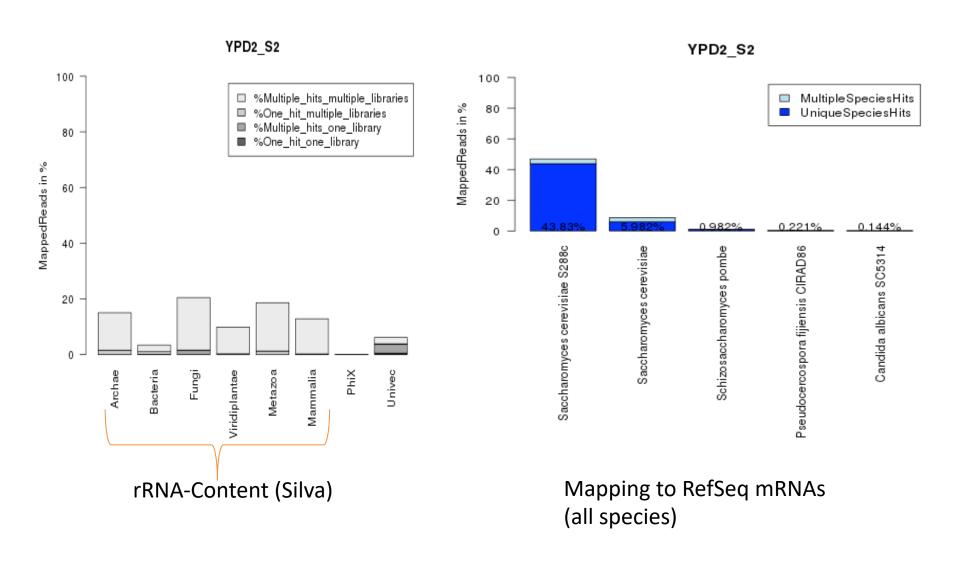
Millions of reads with base resolution

@HWI-ST1034:40:CO8PJACXX:2:1101:20681:1994 1:N:0:ATCACG ${\tt CTCGNAGACTGGCAACTTGTTCTGGTTTACTGCACCTTCTTTTAAAGGCAGAAAGGCTTTTTGATAAAGAAGTTGTGAAAAGGCTACATGAGCTGCTTTTA$ @HWI-ST1034:40:CO8PJACXX:2:1101:1907:2005 1:N:0:ATCACG @HWI-ST1034:40:CO8PJACXX:2:1101:2155:2031 1:N:0:ATCACG ${\tt CAATCAATTAACAATATTAGTTACATAAGCACTTCCTTAACCACCCTCTCAAAGTTGGCAAATGAAGAACCCCCTTTCTCAATAGCTTTAACCGCGCTCTC$ @HWI-ST1034:40:C08PJACXX:2:1101:2220:2057 1:N:0:ATCACG @HWI-ST1034:40:CO8PJACXX:2:1101:2460:2116 1:N:0:ATCACG @HWI-ST1034:40:CO8PJACXX:2:1101:2463:2168 1:N:0:ATCACG ${\tt CGTTCATATGCAAAAGAAGCTTCTCAGTCTGCTTTACCACCTCTTAAAAGGGGATCAAATGTTGAAGAACATCTTTTTTTGAGGTAAAGAACAAATTTGATAT$ @HWI-ST1034:40:CO8PJACXX:2:1101:2378:2207 1:N:0:ATCACG ${\tt CACGCGGTGTGGAAAACCCCTTCACATCCATCAATGGCGGCTCGGAGGGGATTCAAAATCAAGCATATCCGCTTTGTACAGCACAAGACGATCCGATGCTCC$

- How accurate was the sequencing → Fastqc
- Are these reads the intented ones → FastqScreen



Contamination Check - FastqScreen



Data preprocessing common tasks

- 1. Trimming: remove bad bases from (end(s) of) reads
 - Adapter sequence
 - Low quality bases
- 2. Filtering: remove bad reads
 - Low quality reads
 - Contaminating sequences
 - Low complexity reads (repeats)
 - Short (<20bp) reads they slow down mapping software



Data preprocessing software

- fastp
 - https://github.com/OpenG ene/fastp
 - Adapter trimming, quality trimming &filtering, ...
- Trimmomatic
 - https://github.com/usadellab/T rimmomatic
 - Adapter trimming, quality trimming &filtering, ...
- FlexBar (FAR)
 - https://github.com/seqan/flexb
 ar
 - Flexible barcode detection and adapter removal

FASTX

- http://hannonlab.cshl.edu/fast x toolkit/
- Reformat, stats, collapse duplicated reads, trim, filter, reverse compliment
- TagCleaner
 - http://tagcleaner.sourceforge.n
 et
 - Trim MIDs or adaptors, demultiplexing
- DeconSeq
 - http://deconseq.sourceforge. net
 - Remove potential contaminants

Recommendations

- Always generate quality control plots visualizing key characteristics for all libraries
- Trim and/or filter data if needed
- Applications where erroneous reads are of concern:
 - de novo assembly
 - low coverage variant calling
- Applications that are more tolerant to low quality bases
 - RNA-seq