

# Introduction to NGS

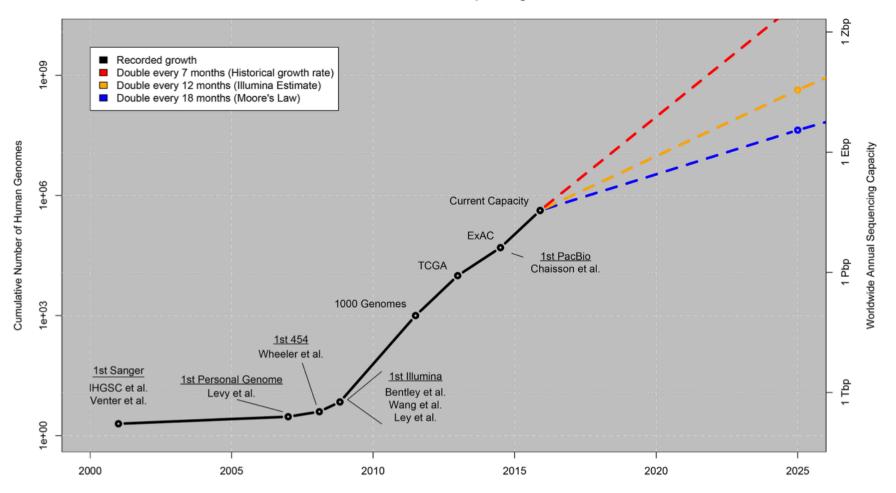
**Hubert Rehrauer** 





#### **NGS** Data Increase

#### **Growth of DNA Sequencing**



NGS data increases faster than computer speed

functional genomics center zurich

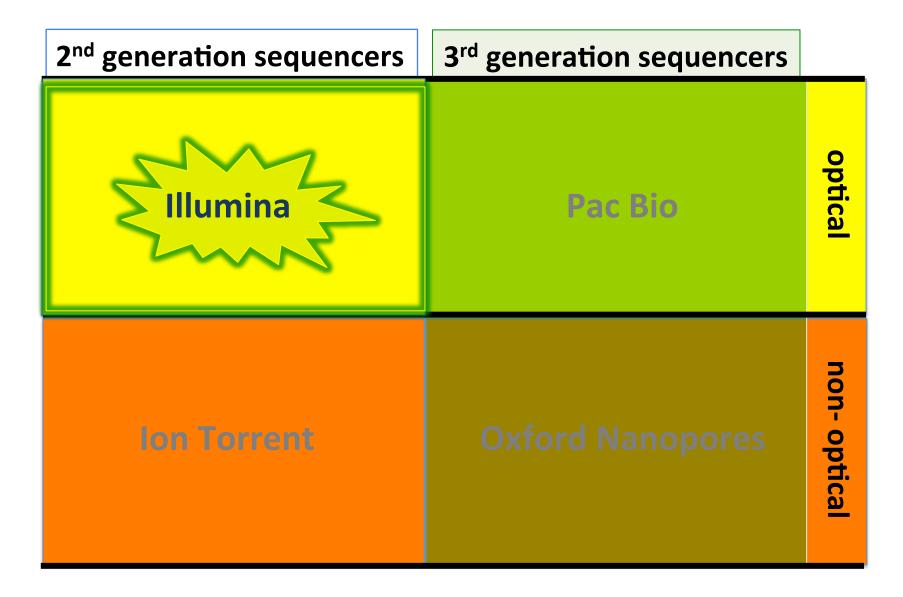
Q · C · Z ·



#### **Ingredients of the NGS success**

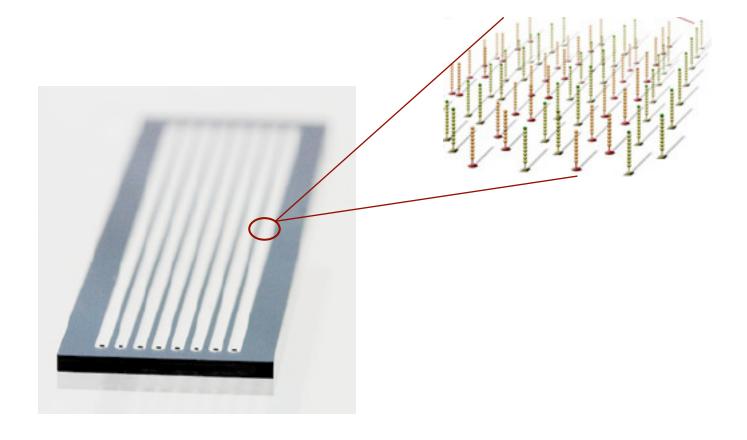
#### **Key ingredients**

- RNA/DNA have evolved over Millions of years into perfect carriers of information
  - they are designed for Read and Write operations
- parallelism
  - the 2018 technology allows in a single experiment the read operation of 10 billion molecules in parallel with high fidelity
- measurements are done by molecular machines
- single molecule manipulation
- Claim: NGS based experiments outperform other Omics technologies in terms of the amount you can learn with a fixed budget

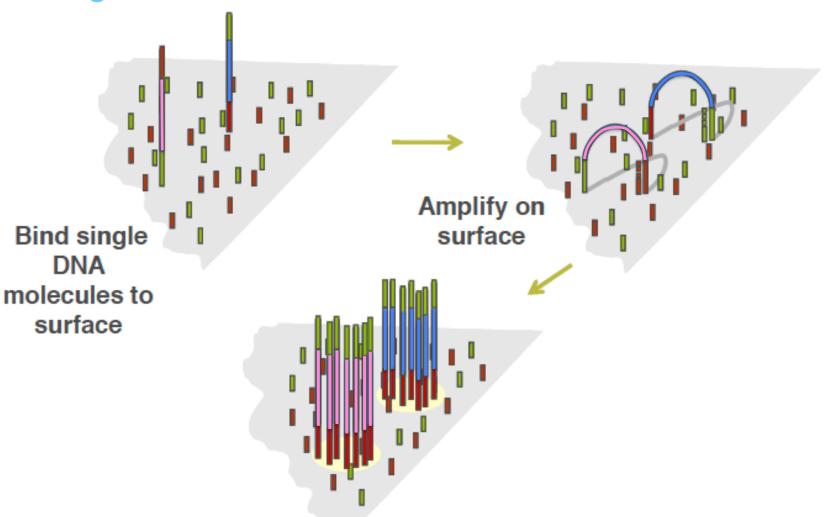




# **Illumina Flow cell**



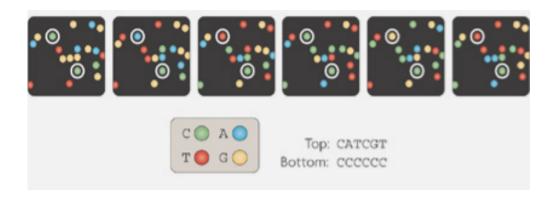
#### **Cluster generation overview**

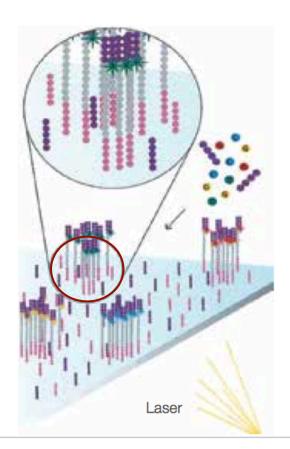




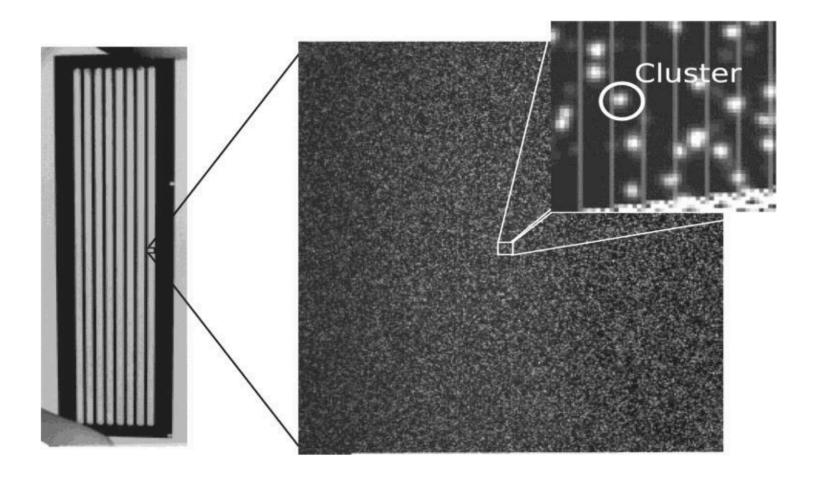


### **Illumina Sequencing**



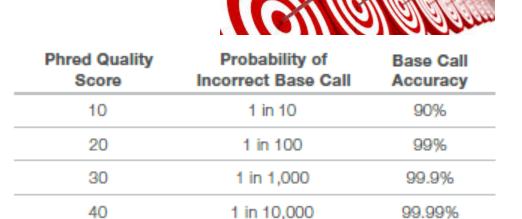


The first sequencing cycle begins by adding four labeled reversible terminators, primers, and DNA polymerase.



## Phred scores measure base call accuracy

- P
- error probability of a given base call
- Q
- -10log<sub>10</sub>P
- Assign to each base
- Range from 0-41



1 in 100,000

99.999%

Ewing B, Green P. 1998. Genome Res. 8(3):186-194.

http://en.wikipedia.org/wiki/Phred\_quality\_score

50

### Phred scores are stored with sequences

- FASTQ
  - 4 lines:
    - 1. Header line for Read (starts with "@" and the sequence ID)
    - 2. Sequence
    - 3. 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 Hx Oct Char	Dec Hx Oct Html Chr	Dec Hx Oct Html Chr Dec Hx Oct Html Chr
0 0 000 NUL (null)	32 20 040   Space	64 40 100 4#64; 0 96 60 140 4#96;
1 1 001 SOH (start of heading)	33 21 041 6#33;	65 41 101 a#65; A 97 61 141 a#97; a
2 2 002 STX (start of text)	34 22 042 @#34; "	66 42 102 a#66; B   98 62 142 a#98; b
3 3 003 ETX (end of text)	35 23 043 # #	67 43 103 a#67; C   99 63 143 a#99; C
4 4 004 EOT (end of transmission)	36 24 044 \$ 年	68 44 104 D D 100 64 144 d d
5 5 005 <b>ENQ</b> (enquiry)	37 25 045 4#37; %	69 45 105 6#69; E   101 65 145 6#101; e
6 6 006 ACK (acknowledge)	38 26 046 & &	70 46 106 F F   102 66 146 f f
7 7 007 BEL (bell)	39 27 047 ' '	71 47 107 6#71; G 103 67 147 6#103; g
8 8 010 <mark>BS</mark> (backspace)	40 28 050 @#40; (	72 48 110 6#72; H   104 68 150 6#104; h
9 9 011 TAB (horizontal tab)	41 29 051 @#41; )	73 49 111 6#73; I   105 69 151 6#105; i
10 A 012 LF (NL line feed, new line		74 4A 112 6#74; J 106 6A 152 6#106; j
ll B 013 VT (vertical tab)	43 2B 053 + +	75 4B 113 6#75; K 107 6B 153 6#107; k
12 C 014 FF (NP form feed, new page		76 4C 114 6#76; L 108 6C 154 6#108; L
13 D 015 CR (carriage return)	45 2D 055 - -	77 4D 115 6#77; M 109 6D 155 6#109; M
14 E 016 SO (shift out)	46 2E 056 . .	78 4E 116 6#78; N 110 6E 156 6#110; n
15 F 017 SI (shift in)	47 2F 057 @#47; /	79 4F 117 6#79; 0 111 6F 157 6#111; 0
16 10 020 DLE (data link escape)	48 30 060 4#48; 0	80 50 120 6#80; P   112 70 160 6#112; P
17 11 021 DC1 (device control 1)	49 31 061 449; 1	81 51 121 6#81; Q 113 71 161 6#113; q
18 12 022 DC2 (device control 2)	50 32 062 4#50; 2	82 52 122 6#82; R   114 72 162 6#114; r
19 13 023 DC3 (device control 3)	51 33 063 3 3	83 53 123 4#83; \$ 115 73 163 4#115; \$
20 14 024 DC4 (device control 4)	52 34 064 6#52; 4	84 54 124 T T   116 74 164 t t
21 15 025 NAK (negative acknowledge)	53 35 065 <b>6#53;</b> 5	85 55 125 6#85; U 117 75 165 6#117; u
22 16 026 SYN (synchronous idle)	54 36 066 @#5 <b>4; 6</b>	86 56 126 V ♥   118 76 166 v ♥
23 17 027 ETB (end of trans. block)	55 37 067 <b>6#55; 7</b>	87 57 127 4#87; ₩ 119 77 167 4#119; ₩
24 18 030 CAN (cancel)	56 38 070 <b>4#56;</b> 8	88 58 130 4#88; X 120 78 170 4#120; X
25 19 031 EM (end of medium)	57 39 071 9 9	89 59 131 6#89; Y 121 79 171 6#121; Y
26 1A 032 SUB (substitute)	58 3A 072 @#58; :	90 5A 132 6#90; Z 122 7A 172 6#122; Z
27 1B 033 <b>ESC</b> (escape)	59 3B 073 ; ;	91 5B 133 6#91; [   123 7B 173 6#123; {
28 1C 034 FS (file separator)	60 3C 074 @#60; <	92 5C 134 6#92; \ 124 7C 174 6#124;
29 1D 035 GS (group separator)	61 3D 075 = =	93 5D 135 6#93; ] 125 7D 175 6#125; }
30 1E 036 RS (record separator)	62 3E 076 > >	94 5E 136 6#94; ^ 126 7E 176 6#126; ~
31 1F 037 US (unit separator)	63 3F 077 ? ?	95 5F 137 6#95; _  127 7F 177 6#127; DEL

Source: www.LookupTables.com



## Quality control is to know your reads

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

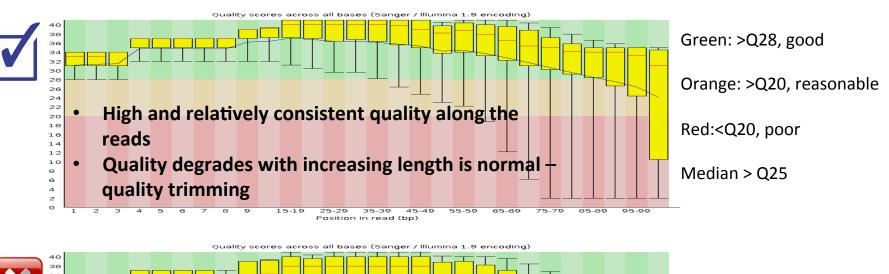
- Sequencing errors
  - Chemical, optical, computational

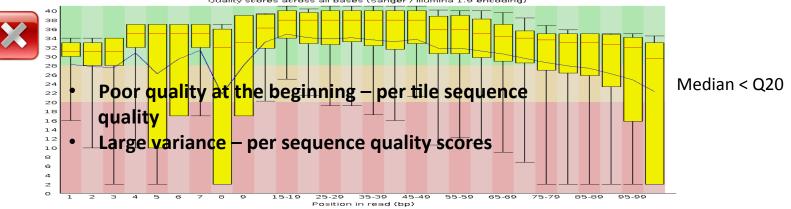
Platform	Primary error	Error rate (%)
Illumina	Substitution	0.1
PacBio	Indel	12
PGM	Indel	1
454	Indel	1
Oxford Nanopore	Indel	20-40



#### Per base sequence quality - FastQC

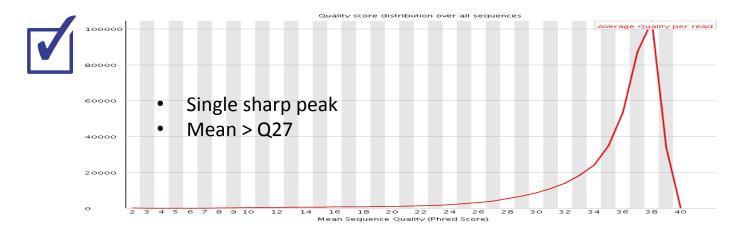
Range of quality values across all bases at each position

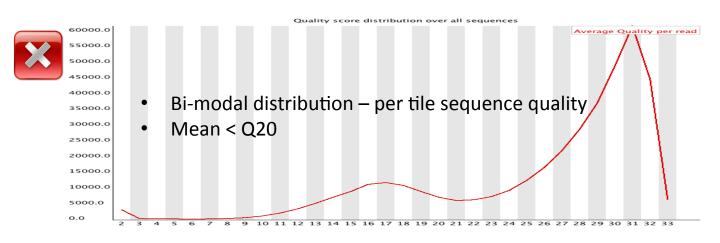




#### Per sequence quality scores - FastQC

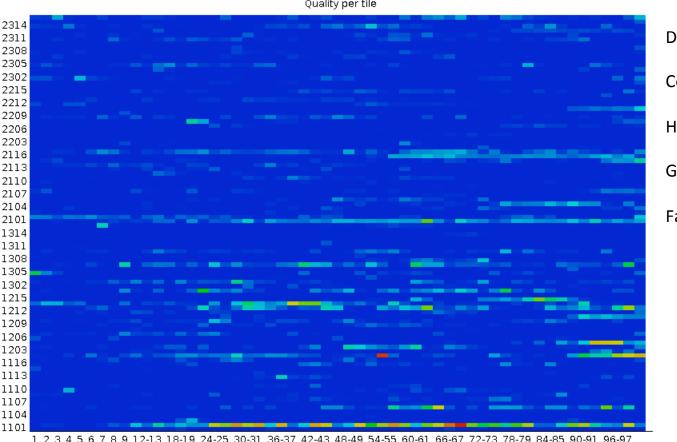
Subset of sequences with universally low quality values





### 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

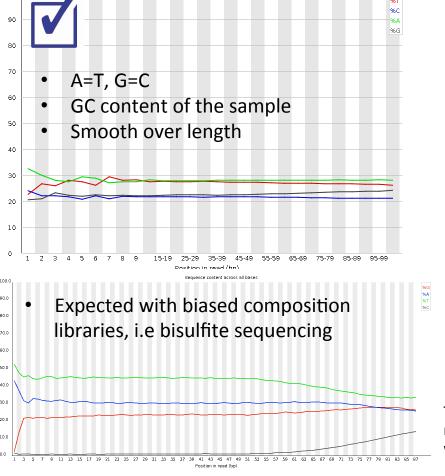
Hotter color: worse quality

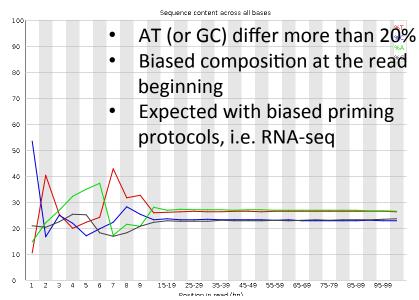
Good: universal blue

Failure: < average - 5

#### Per base sequence content - FastQC

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





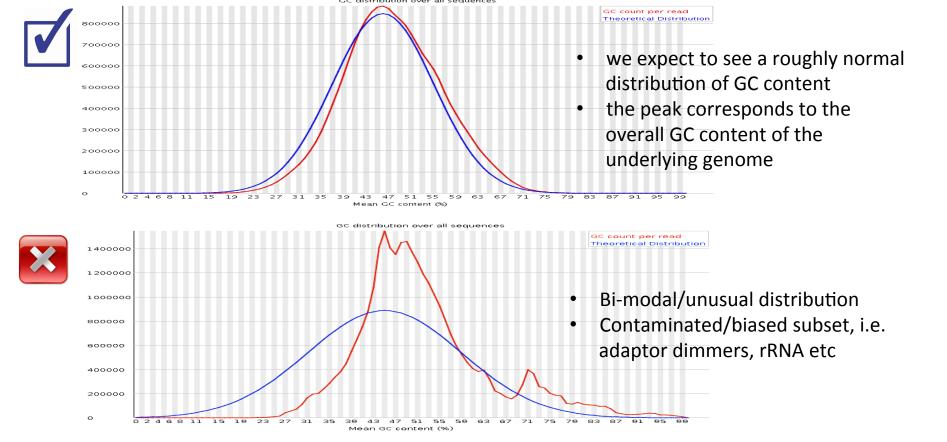
Biases in Illumina transcriptome sequencing caused by random hexamer priming

Kasper D. Hansen<sup>1,\*</sup>, Steven E. Brenner<sup>2</sup> and Sandrine Dudoit<sup>1,3</sup>

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.

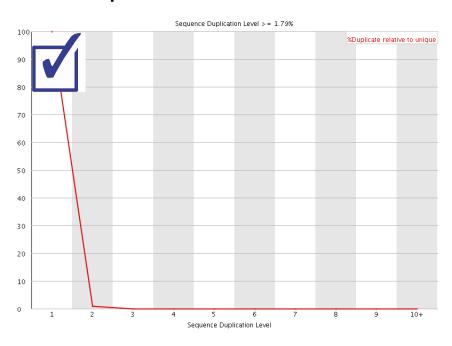
### Per sequence GC content - FastQC

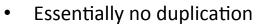
Distribution of average GC in all reads

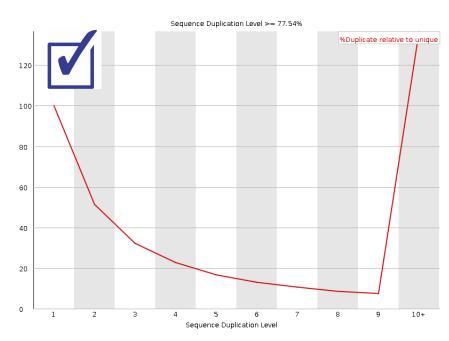


#### **Sequence duplication - FastQC**

 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)

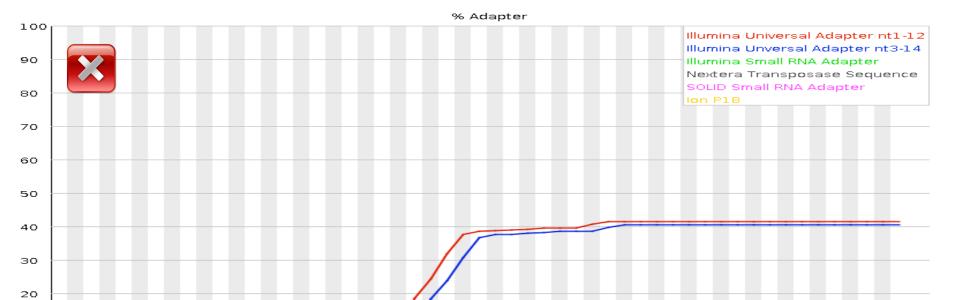
- Can be normal and biologically meaningful
  - highly expressed transcripts
  - high copy number repeats
  - Less diverse library (amplicons)



10

#### Adapter Content - FastQC





1 2 3 4 5 6 7 8 9 12-13 18-19 24-25 30-31 36-37 42-43 48-49 54-55 60-61 66-67 72-73 78-79 84-85 90-91 Position in read (bp)

functional genomics center zurich

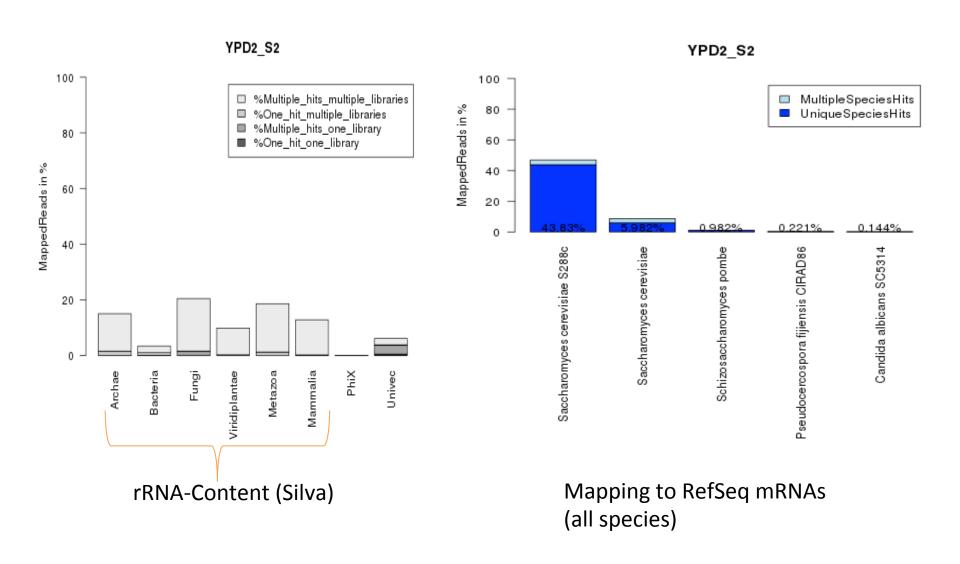
f. g. c. z.

# Millions of reads with base resolution

@HWI-ST1034:40:C08PJACXX:2:1101:20681:1994 1:N:0:ATCACG  ${\tt CTCGNAGACTGGCAACTTGTTCTGGTTTACTGCACCTTCTTTTAAAGGCAGAAAGGCTTTTTGATAAAGAAGTTGTGAAAAGGCTACATGAGCTGCTT$ @HWI-ST1034:40:C08PJACXX:2:1101:1907:2005 1:N:0:ATCACG  $\mathtt{CTCACCTTCAACTGTATTCACGCTTGGACCACAGATCTTGGCCTTAGTGCGATATAGGACAC$ FACGTGCGCATGCTCGCTCAGCTCTTCCCAGACTACCCAATCCTTGCCAAGTGC @HWI-ST1034:40:CO8PJACXX:2:1101:2463:2168 1:N:0:ATCACG  ${\tt CGTTCATATGCAAAAGAAGCTTCTCAGTCTGCTTTACCACCTCTTAAAGGGGGATCAAATGTTGAAGAACATCTTTTTTGAGGTAAAGAACAAATTTGATAT$ MHWI-ST1034:40:C08PJACXX:2:1101:2378:2207 1:N:0:ATCACG BCCFFFFDHHHHHJJJJJJJJJJJJJJJJJJJJJJJIJIGIJHEHBD8>6?:ABCDDDCCDDCCBCCCDEBBBBBACCEEEECCBCDDDCDBBB?BBBDDDDC



### **Contamination Check - FastqScreen**



## **Data preprocessing common tasks**

- 1. Trimming: remove bad bases from (end(s) of) reads
  - Adaptor 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</li>



# **Data preprocessing software**

**University of** 

- PRINSEQ
  - http://prinseq.sourceforge.net/
  - Quality/hard trimming, quality filtering, reformat, ...
- Trimmomatic
  - http://www.usadellab.org/ cms/?page=trimmomatic
  - Adaptor trimming, quality trimming &filtering, ...
- FlexBar (FAR)
  - http://sourceforge.net/ projects/theflexibleadap/
  - Flexible barcode detection and adapter removal

#### FASTX

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

#### **Recommendations**

- Always generate quality plots for all libraries
- Trim and/or filter data if needed
  - always trim and filter for de-novo transcriptome assembly