



10

01

101

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

01

1

010

01

..

011

10

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

..

Introduction to NGS

Hubert Rehrauer



University of
Zurich^{UZH}

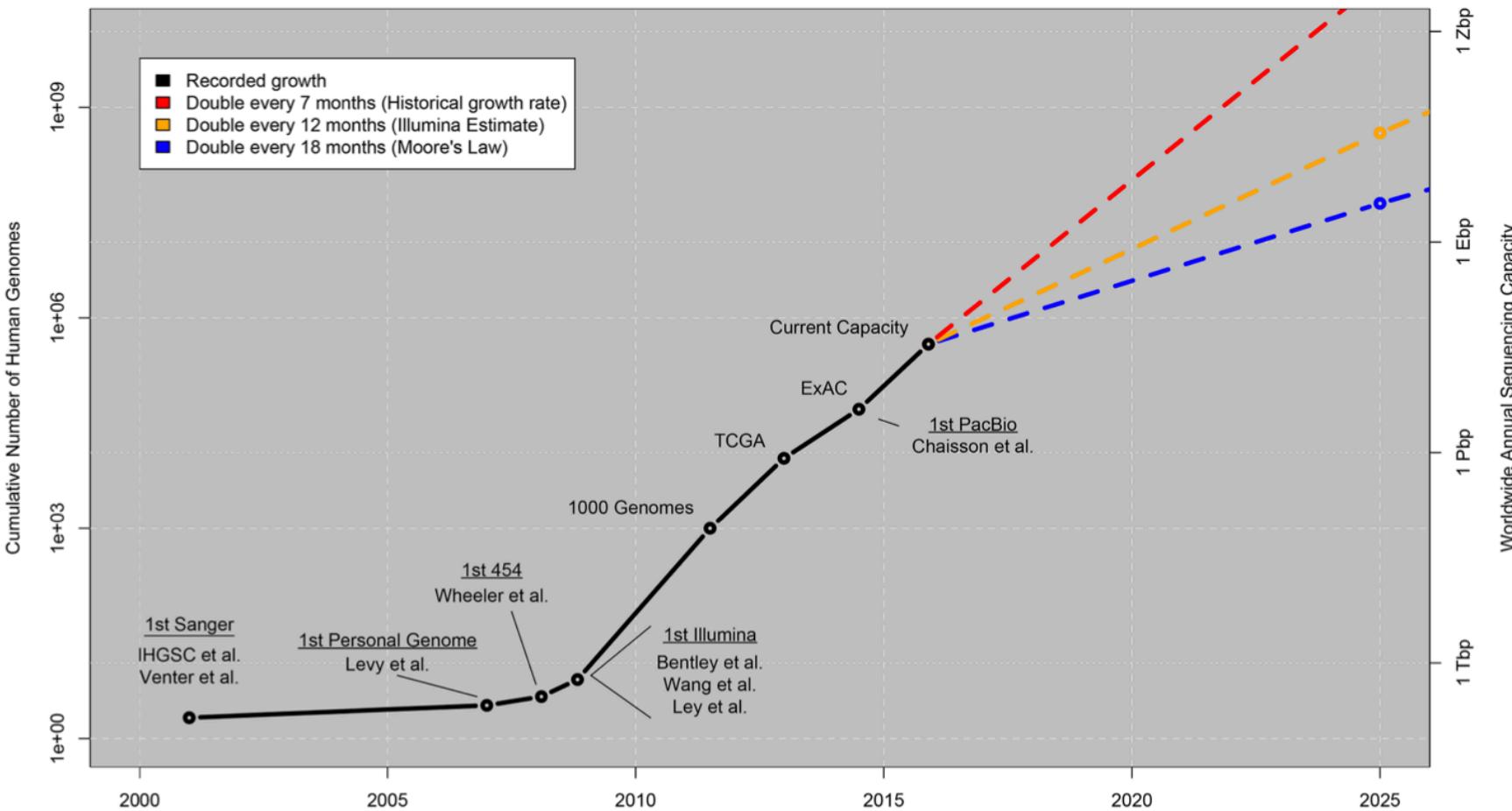
ETH

Eidgenössische Technische Hochschule Zürich
Swiss Federal Institute of Technology Zurich

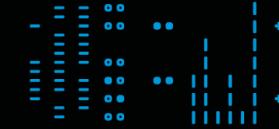


NGS Data Increase

Growth of DNA Sequencing



- NGS data increases faster than computer speed

10
01
101functional genomics center zurich
01 10
010 01
101 10
010 01
01 10
01 01
f g c z
10
01

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**)

10
01
101101 1
010 0
0101 1001 1
10
01 0

2nd generation sequencers



3rd generation sequencers

Pac Bio

Ion Torrent

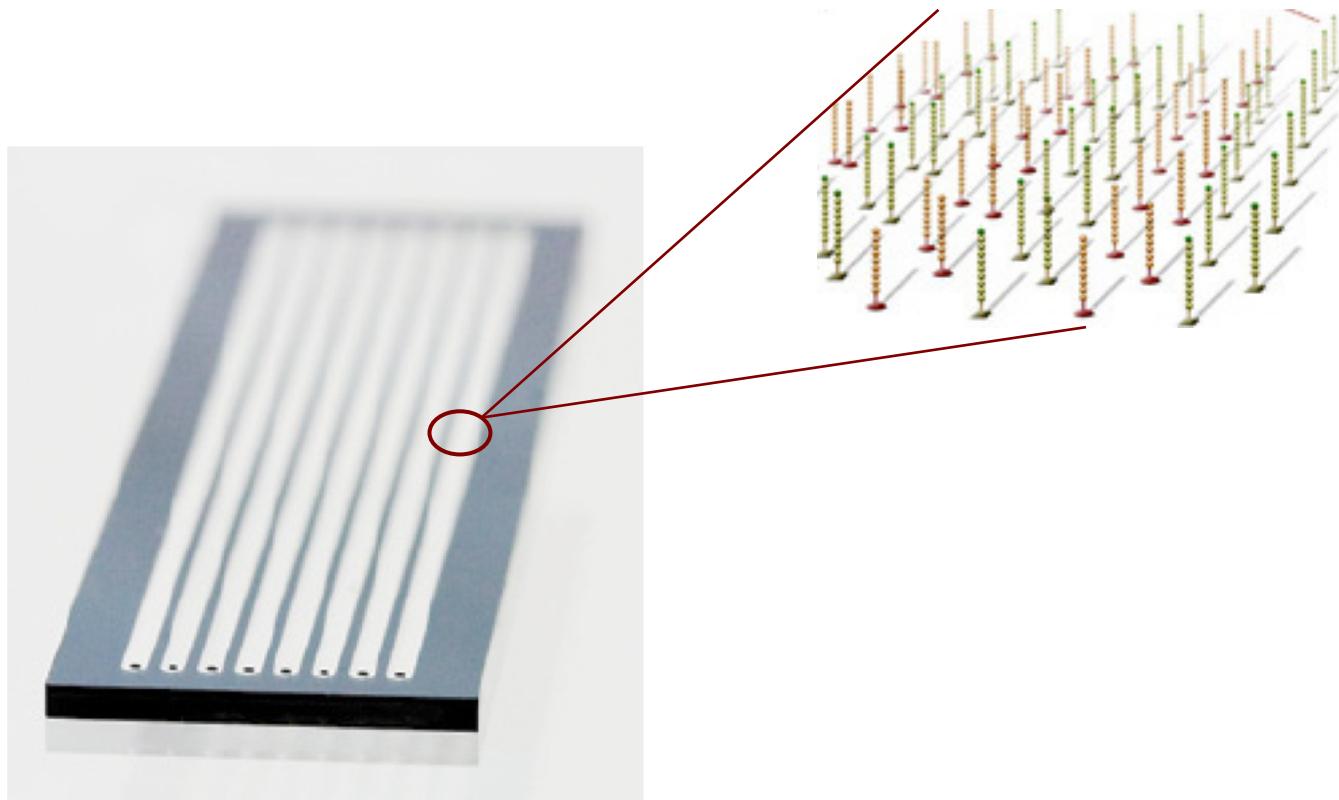
Oxford Nanopores

optical

non-optical

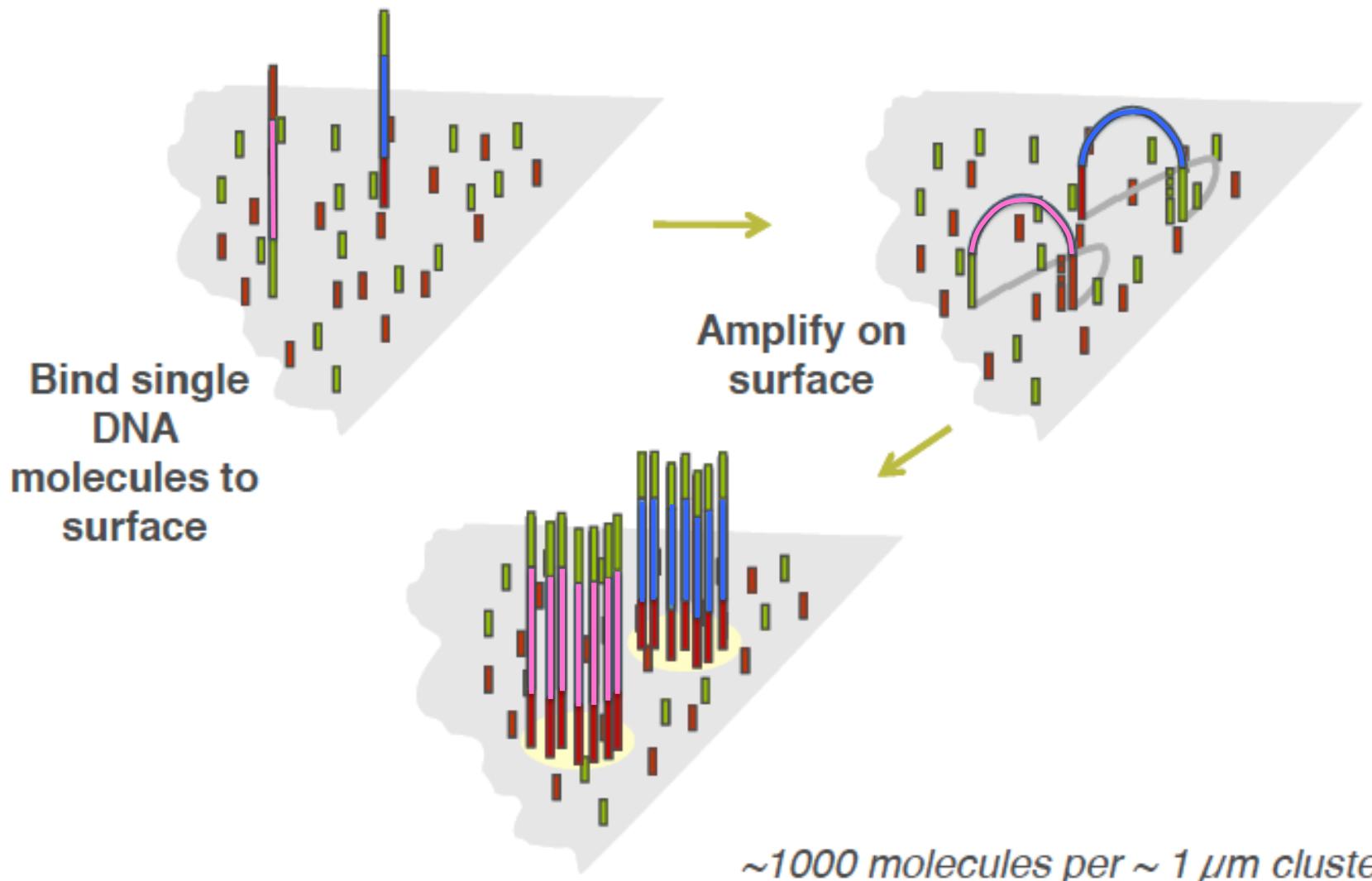
10
01
101101 1
010 0
0101 10... .
... .
... .| ++ +++
010 01. . f . g . c . z .
101 10
010 01.
10 01
01 00

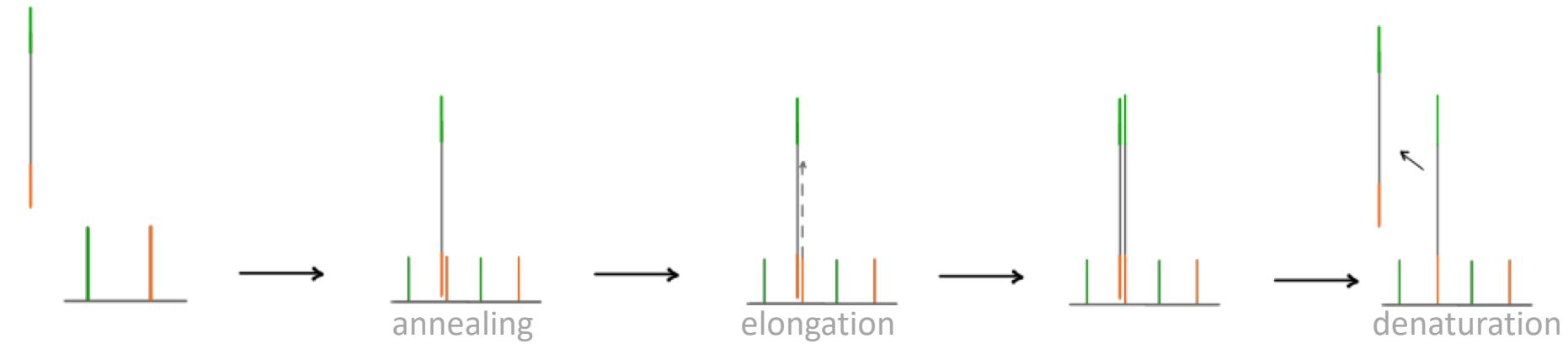
Illumina Flow cell



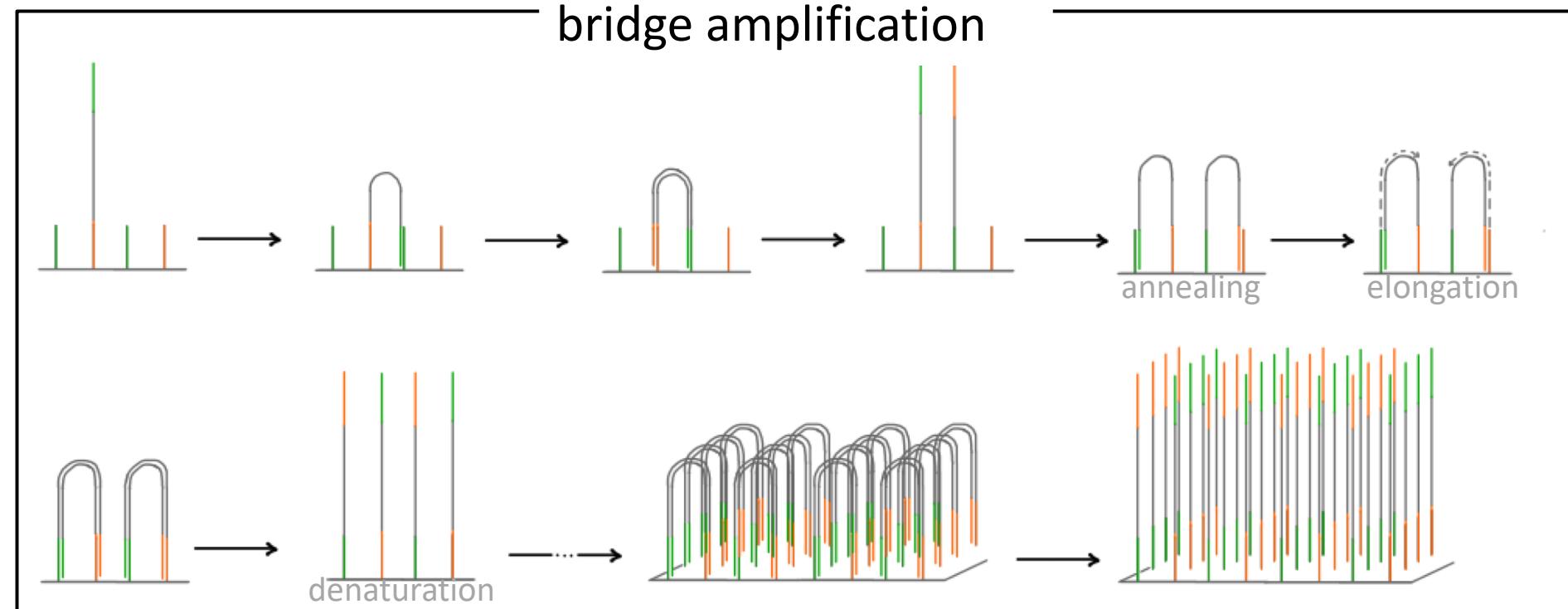
10
01
101101 1
010 0
0101 10010 01
101 10
010 01
01 01

Cluster generation overview

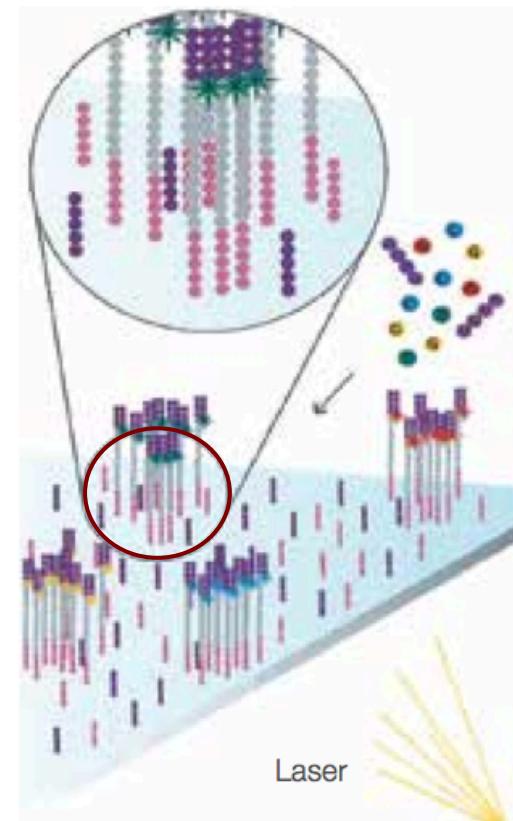
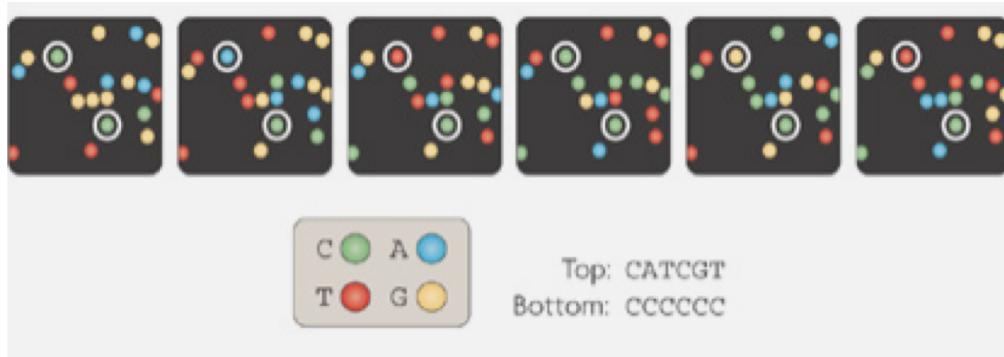




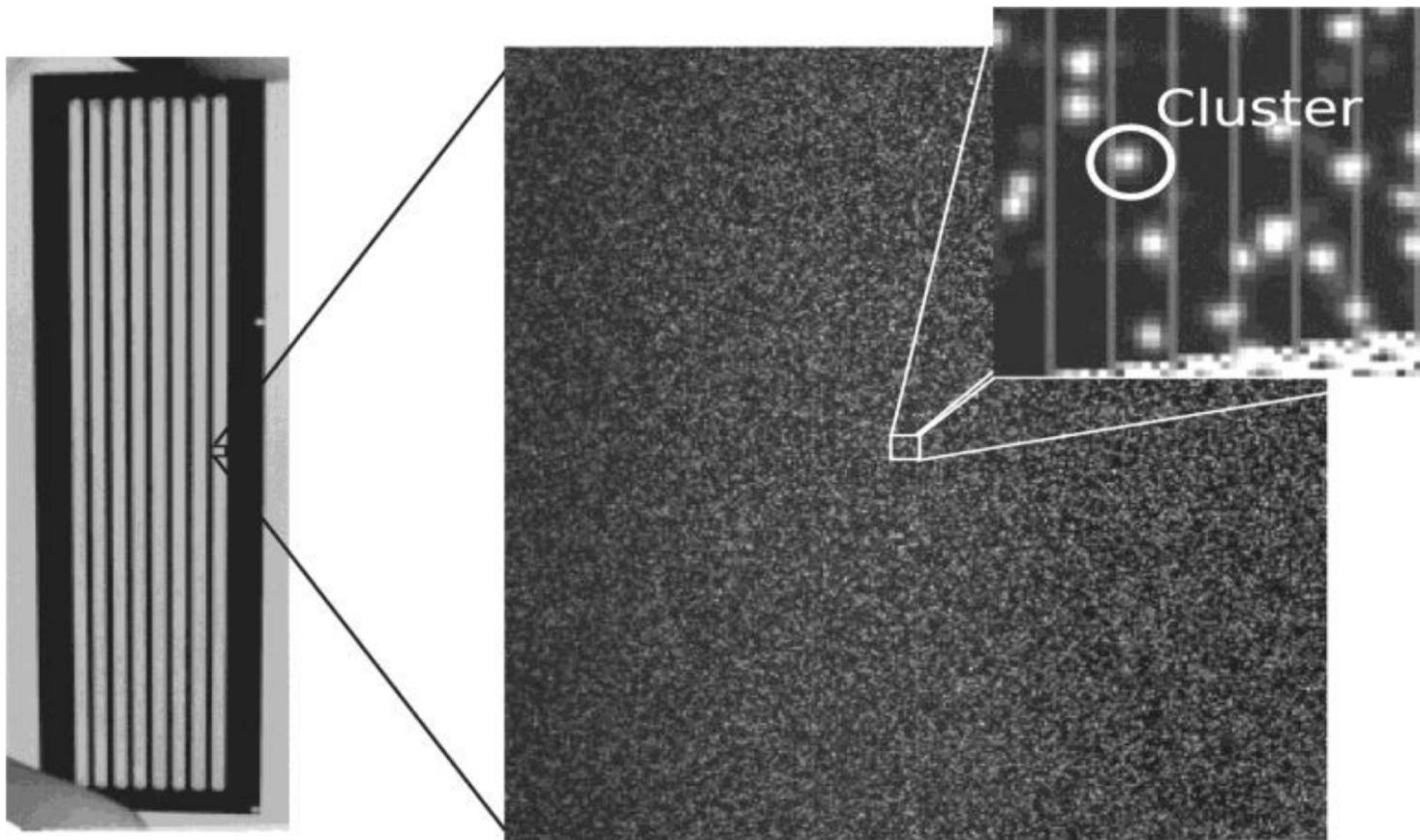
bridge amplification



Illumina Sequencing



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



10
01
101101
1
010
0
0101
10010
01
101
10
010
0101
0
01
10
01

Phred scores measure base call accuracy

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



Phred Quality Score	Probability of Incorrect Base Call	Base Call Accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1,000	99.9%
40	1 in 10,000	99.99%
50	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

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)

@HWI-ST1034:40:C08PJACXX:2:1101:20681:1994 1:N:0:ATCACG
CTCGNAGACTGGCAACTTGTCTGGTTACTGCACCTCTTTAAAGGCAGAAAGGC
+
CCCF#2ADHHHGHJJJJHIIIIJHIIJJHJJJJJJJBGIIJJJJJJJJJJJJJJ

10
01
1010101
0100
0101 10

... .

| ++

++

++

++

++

++

++

++

++

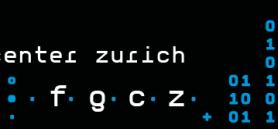
++

++

++

++

++



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



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	@	Ø	96	60	140	`	`
1	1	001	SOH (start of heading)	33	21	041	!	!	65	41	101	A	A	97	61	141	a	a
2	2	002	STX (start of text)	34	22	042	"	"	66	42	102	B	B	98	62	142	b	b
3	3	003	ETX (end of text)	35	23	043	#	#	67	43	103	C	C	99	63	143	c	c
4	4	004	EOT (end of transmission)	36	24	044	$	\$	68	44	104	D	D	100	64	144	d	d
5	5	005	ENQ (enquiry)	37	25	045	%	%	69	45	105	E	E	101	65	145	e	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	G	G	103	67	147	g	g
8	8	010	BS (backspace)	40	28	050	((72	48	110	H	H	104	68	150	h	h
9	9	011	TAB (horizontal tab)	41	29	051))	73	49	111	I	I	105	69	151	i	i
10	A	012	LF (NL line feed, new line)	42	2A	052	*	*	74	4A	112	J	J	106	6A	152	j	j
11	B	013	VT (vertical tab)	43	2B	053	+	+	75	4B	113	K	K	107	6B	153	k	k
12	C	014	FF (NP form feed, new page)	44	2C	054	,	,	76	4C	114	L	L	108	6C	154	l	l
13	D	015	CR (carriage return)	45	2D	055	-	-	77	4D	115	M	M	109	6D	155	m	m
14	E	016	SO (shift out)	46	2E	056	.	.	78	4E	116	N	N	110	6E	156	n	n
15	F	017	SI (shift in)	47	2F	057	/	/	79	4F	117	O	O	111	6F	157	o	o
16	10	020	DLE (data link escape)	48	30	060	0	0	80	50	120	P	P	112	70	160	p	p
17	11	021	DC1 (device control 1)	49	31	061	1	1	81	51	121	Q	Q	113	71	161	q	q
18	12	022	DC2 (device control 2)	50	32	062	2	2	82	52	122	R	R	114	72	162	r	r
19	13	023	DC3 (device control 3)	51	33	063	3	3	83	53	123	S	S	115	73	163	s	s
20	14	024	DC4 (device control 4)	52	34	064	4	4	84	54	124	T	T	116	74	164	t	t
21	15	025	NAK (negative acknowledge)	53	35	065	5	5	85	55	125	U	U	117	75	165	u	u
22	16	026	SYN (synchronous idle)	54	36	066	6	6	86	56	126	V	V	118	76	166	v	v
23	17	027	ETB (end of trans. block)	55	37	067	7	7	87	57	127	W	W	119	77	167	w	w
24	18	030	CAN (cancel)	56	38	070	8	8	88	58	130	X	X	120	78	170	x	x
25	19	031	EM (end of medium)	57	39	071	9	9	89	59	131	Y	Y	121	79	171	y	y
26	1A	032	SUB (substitute)	58	3A	072	:	:	90	5A	132	Z	Z	122	7A	172	z	z
27	1B	033	ESC (escape)	59	3B	073	;	:	91	5B	133	[[123	7B	173	{	{
28	1C	034	FS (file separator)	60	3C	074	<	<	92	5C	134	\	\	124	7C	174	|	
29	1D	035	GS (group separator)	61	3D	075	=	=	93	5D	135]]	125	7D	175	}	}
30	1E	036	RS (record separator)	62	3E	076	>	>	94	5E	136	^	^	126	7E	176	~	~
31	1F	037	US (unit separator)	63	3F	077	?	?	95	5F	137	_	_	127	7F	177		DEL

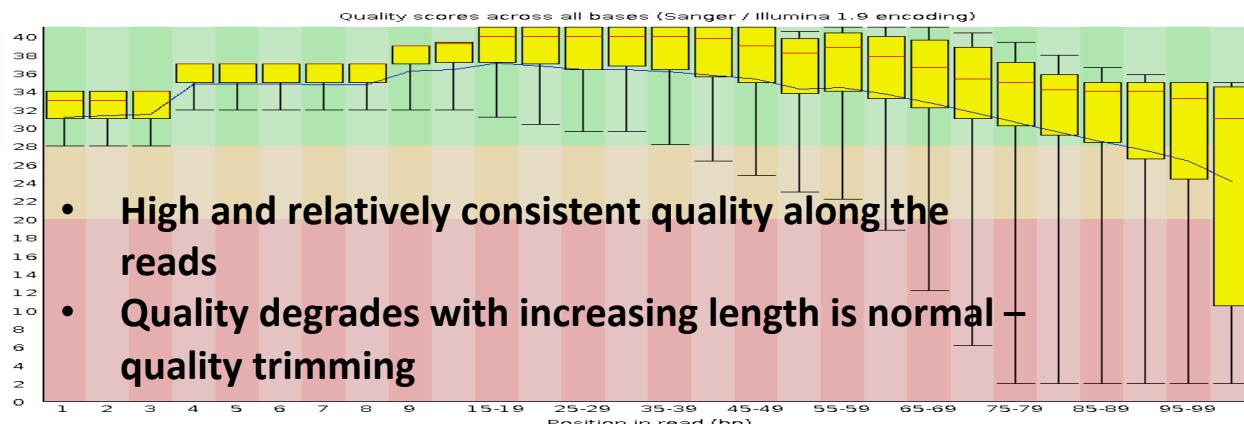
Read Quality Control

- 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	Indel	3 - 20
Nanopore		

Per base sequence quality - FastQC

- Range of quality values across all bases at each position

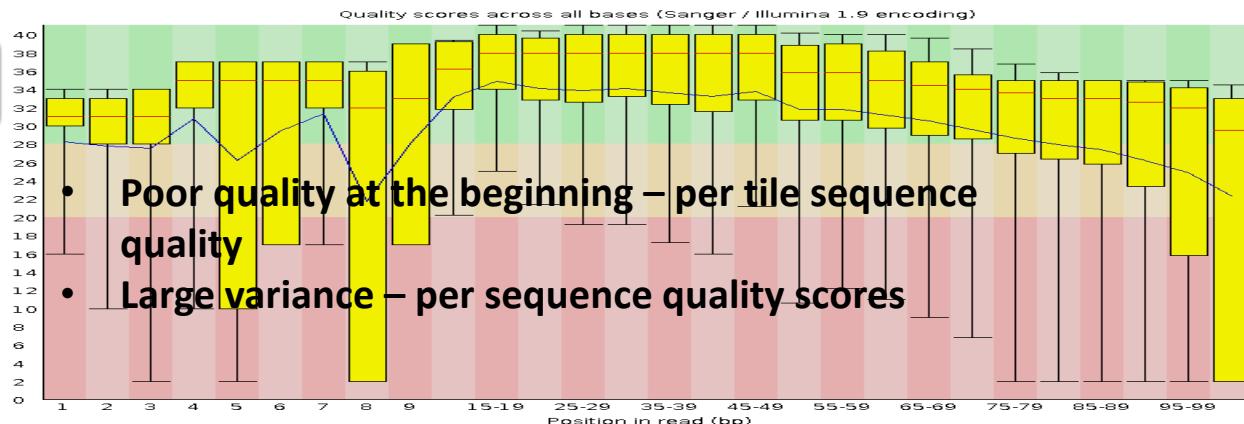


Green: >Q28, good

Orange: >Q20, reasonable

Red:<Q20, poor

Median > Q25

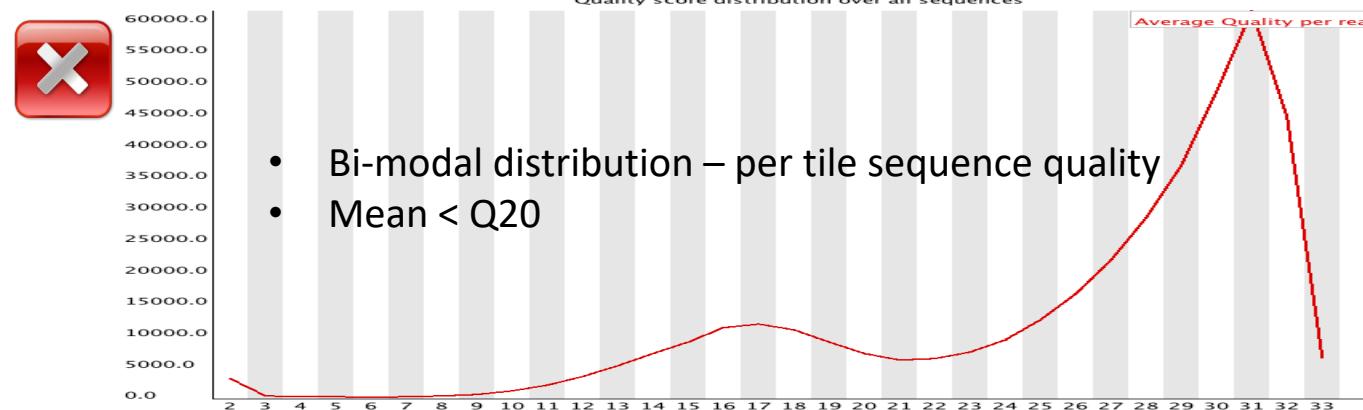


Median < Q20

10
01
101functional genomics center zurich
01 1
01 01
101 10
010 01
101 10
01 01
01 10
01 01

Per sequence quality scores - FastQC

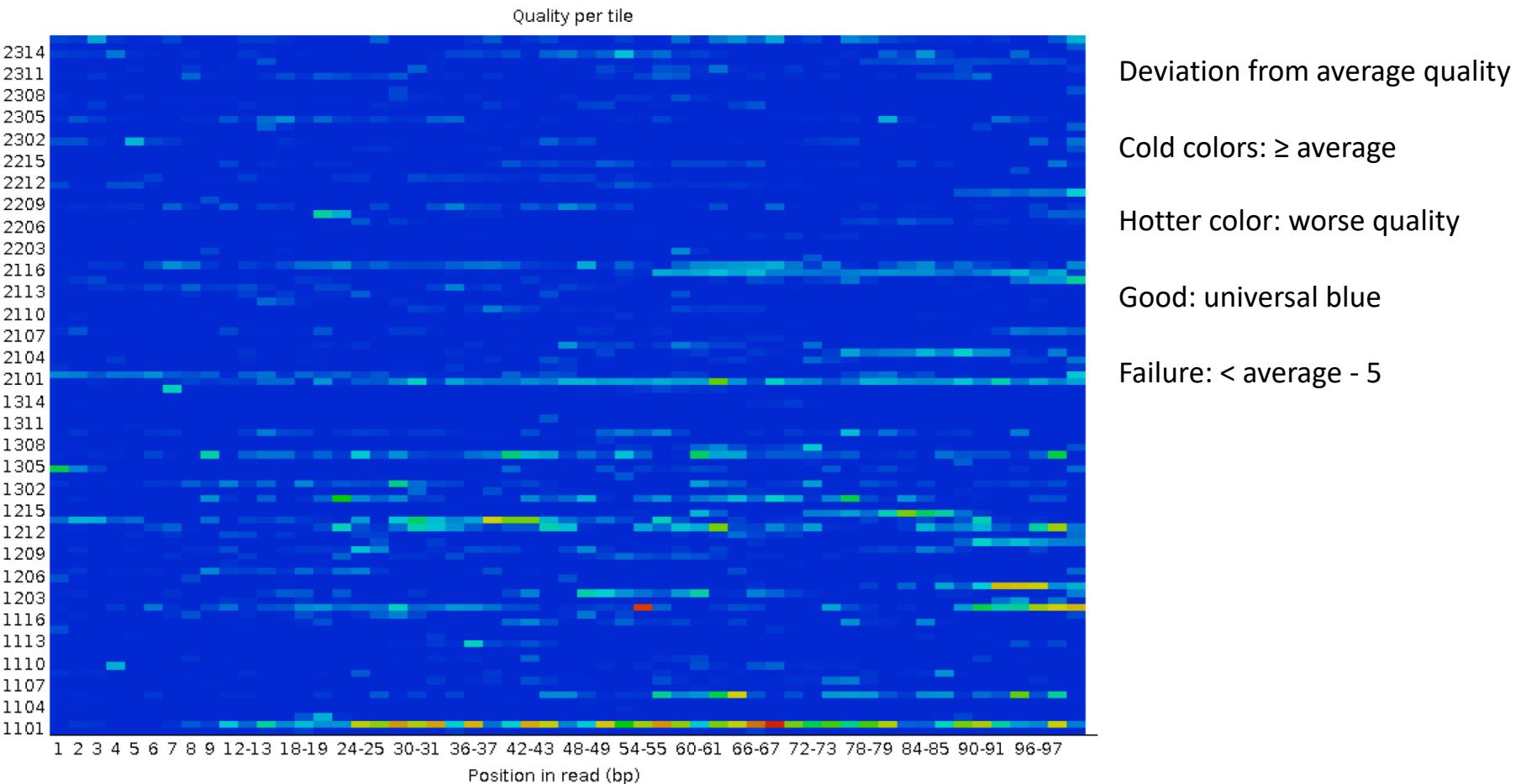
- Subset of sequences with universally low quality values



10
01
10101 10
01 01
101 10
010 01
101 10
01 10
01 01

Per tile sequence quality - FastQC

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

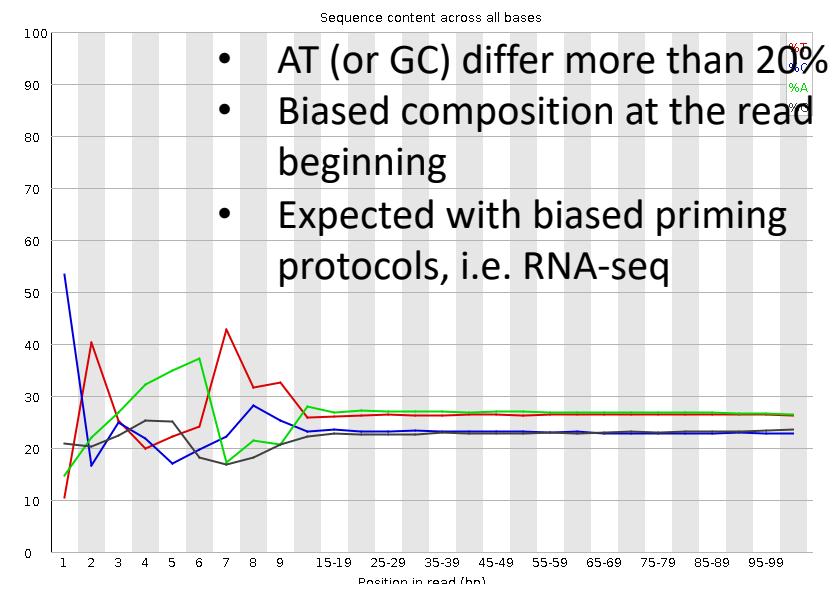
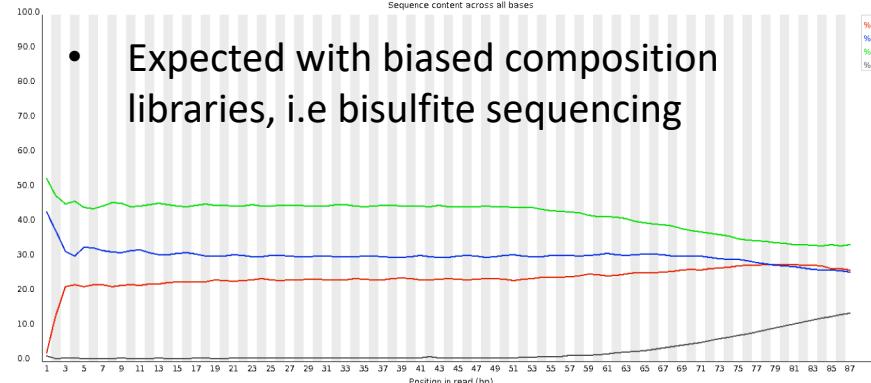
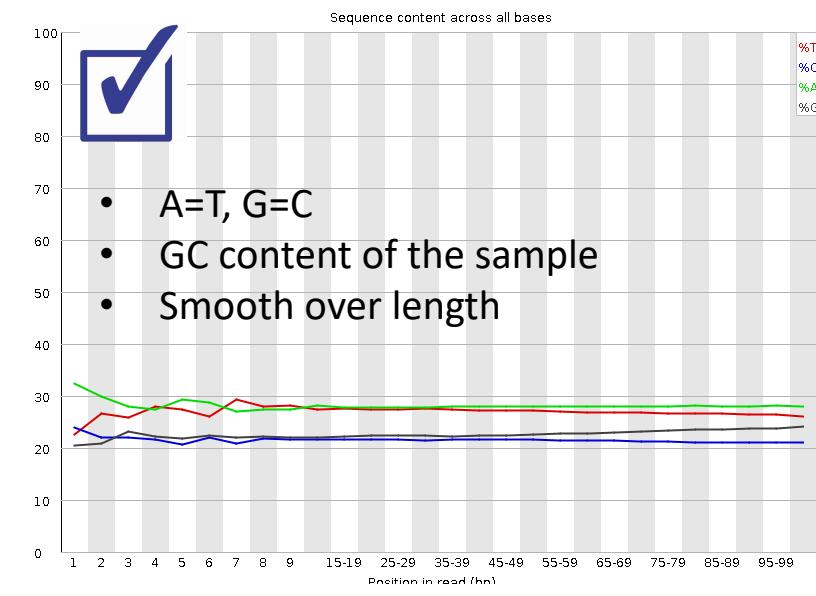


10
01
101101
010
010
0101
10010 01
101 10
010 01
010 01
01 10
01 01

f g c z

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

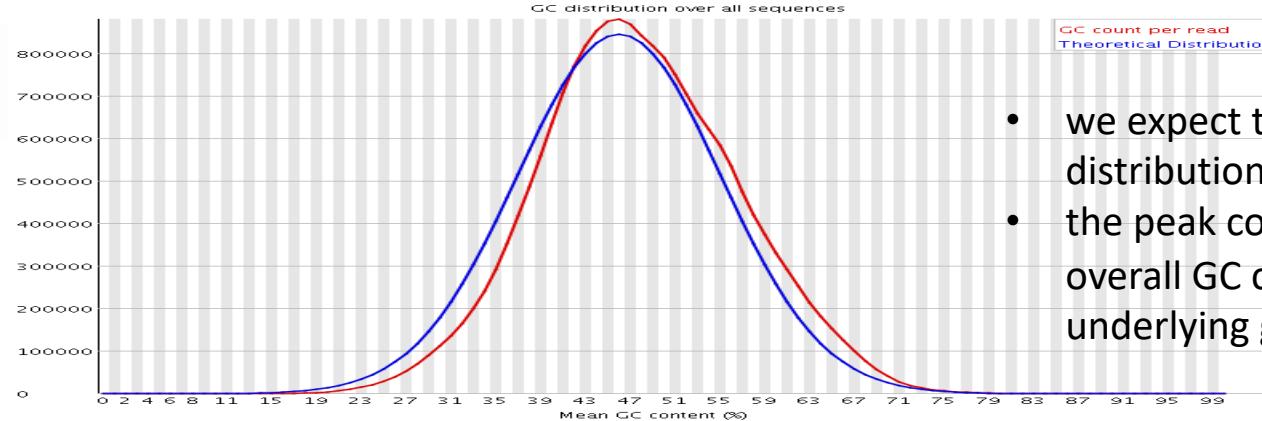
Kasper D. Hansen^{1,*}, Steven E. Brenner² and Sandrine Dudoit^{1,3}

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.

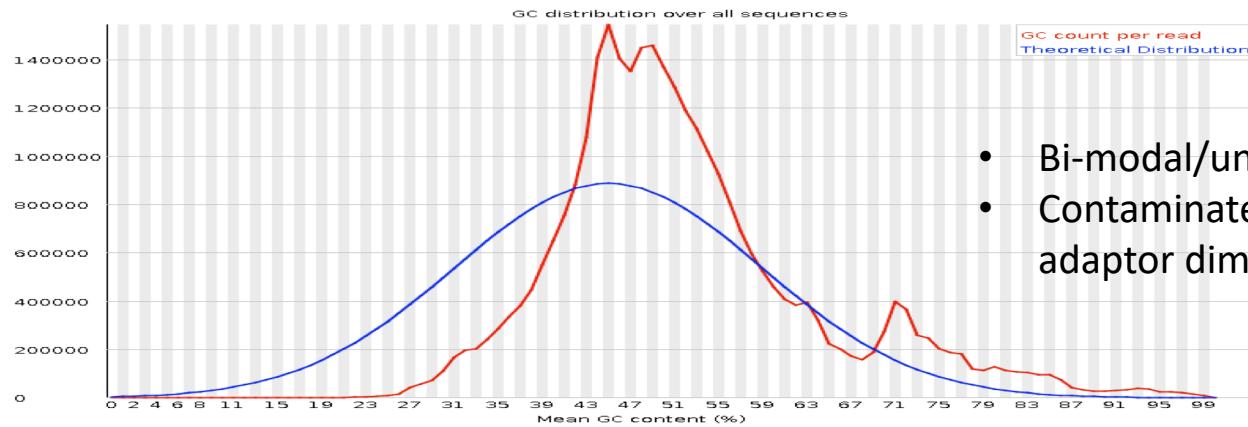
10
01
10100
01
010..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01..
00
01

Per sequence GC content - FastQC

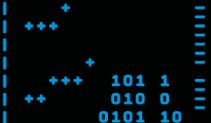
- Distribution of average GC in all reads



- we expect to see a roughly normal distribution of GC content
- the peak corresponds to the overall GC content of the underlying genome

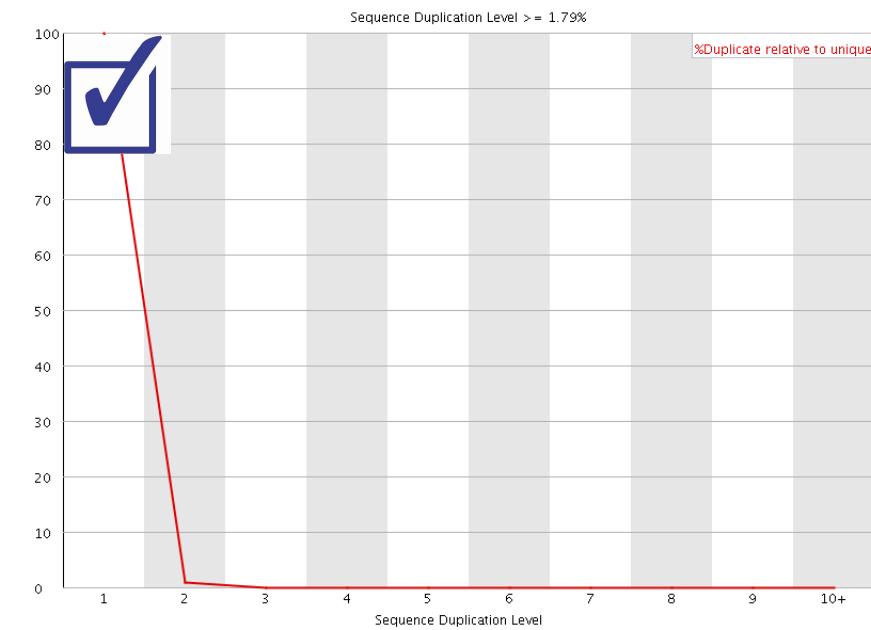


- Bi-modal/unusual distribution
- Contaminated/biased subset, i.e. adaptor dimmers, rRNA etc

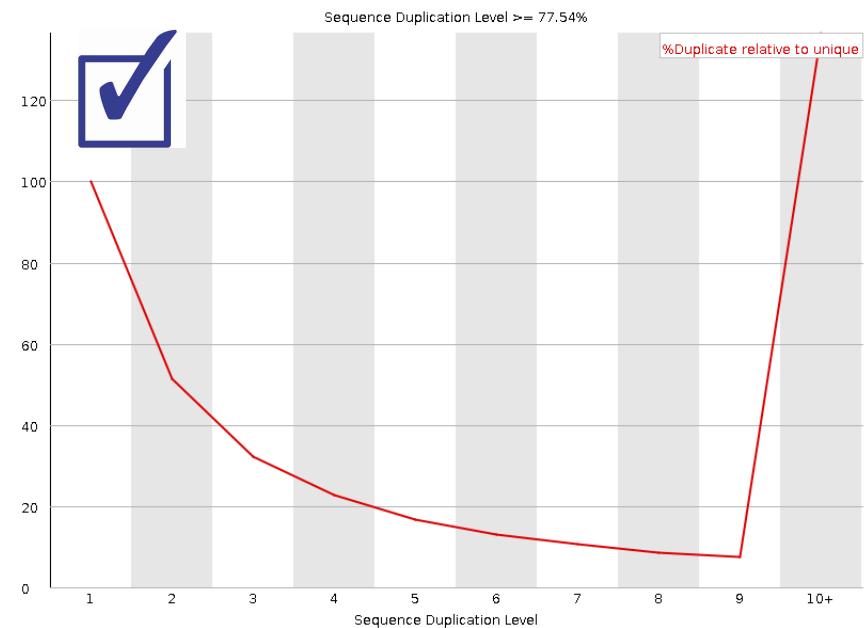
10
01
101functional genomics center zurich
01 10
01 01
101 10
010 01
101 10
010 01
01 10
01 01

Sequence duplication - FastQC

- Relative number of sequences with different degrees of duplication



- Essentially no 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)

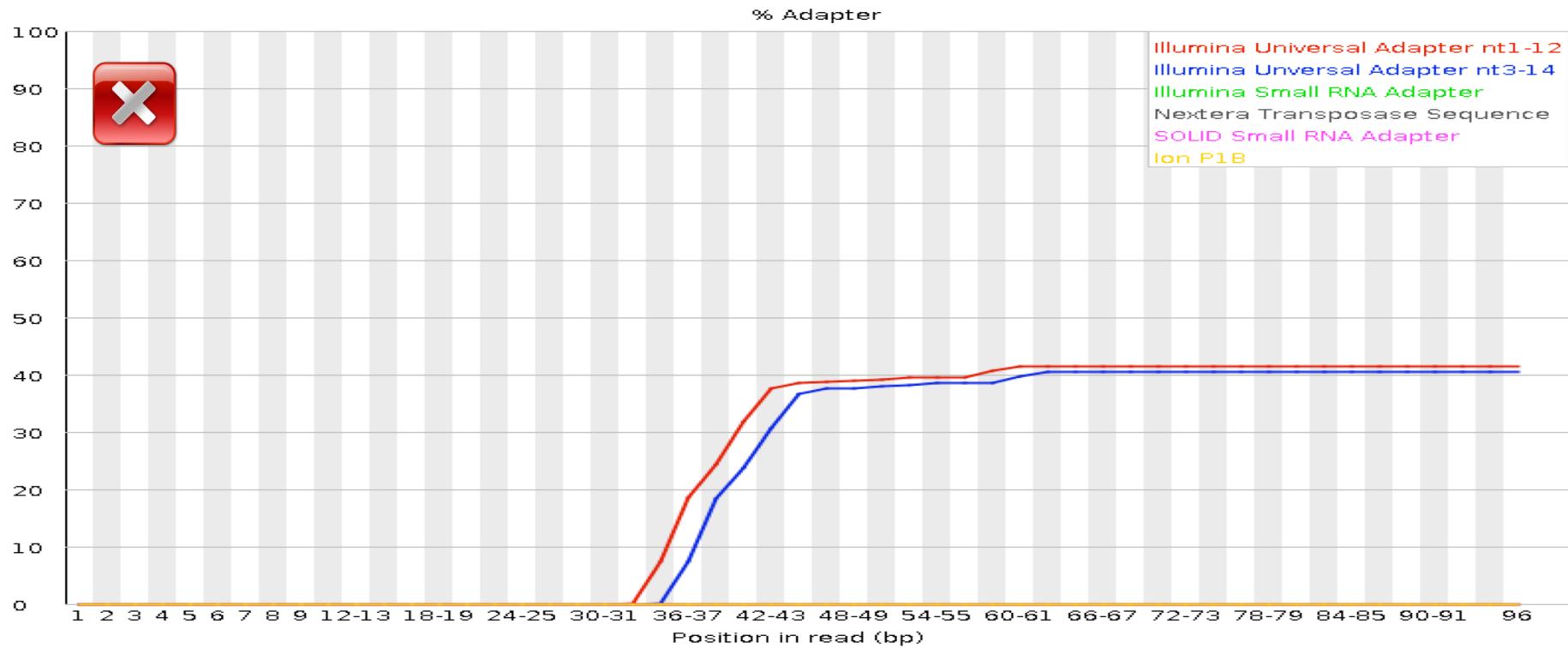
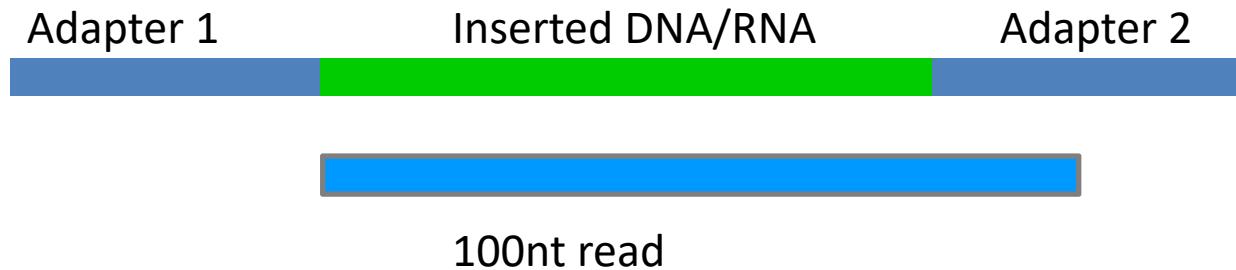


Overrepresented sequences

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



Millions of reads with base resolution

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



10

01

101

00

00

001

00

00

001

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

00

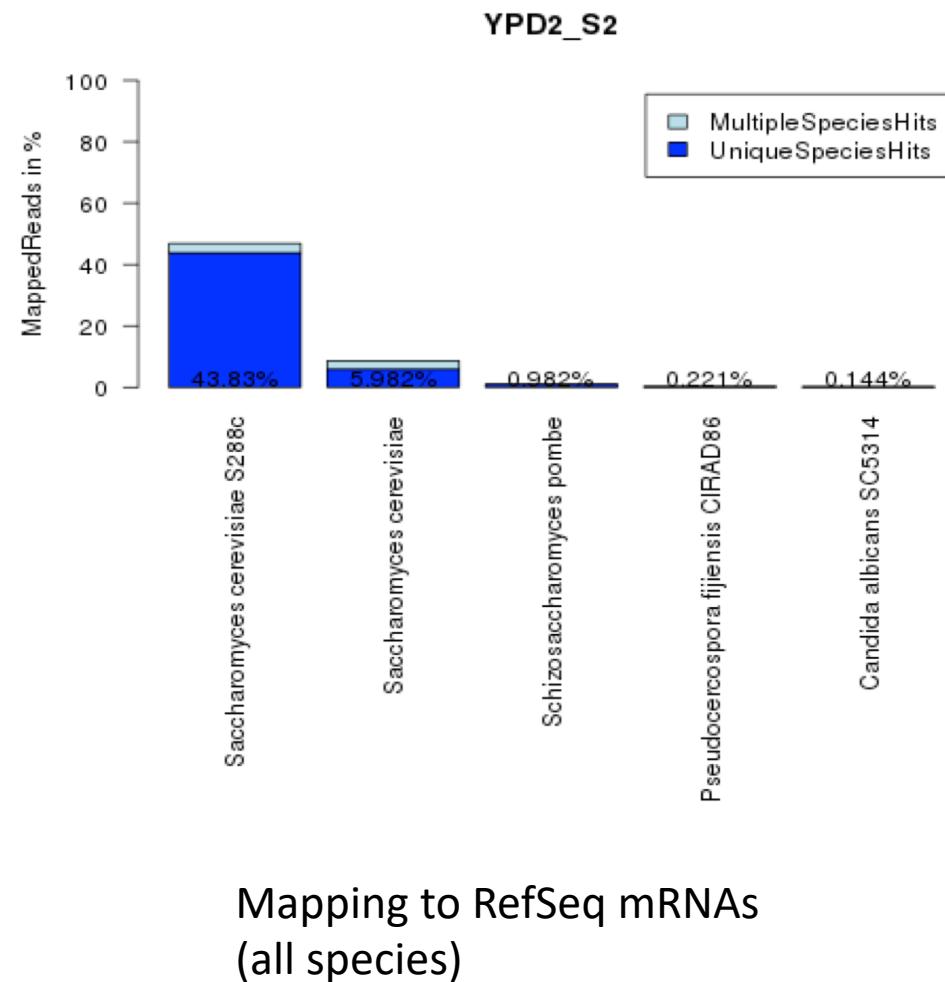
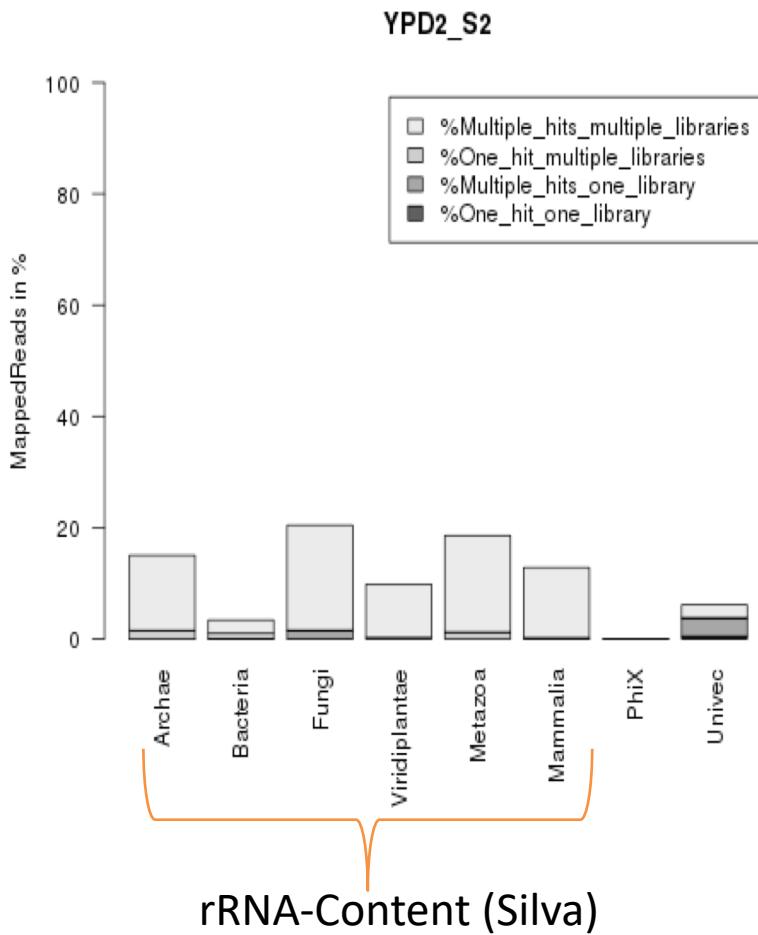
00

00

00

00

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

Data preprocessing software

- 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)
 - <https://github.com/seqan/flexbar>
 - 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
- Applications where erroneous reads are detrimental
 - de novo assembly
 - low coverage variant calling