J CATGCATGCA1	FGCATGCAAT			
JCATGCATGCATGC	ATGCATGCATG			
.GCTGATGCATGCATGC	ATGCATGCATE			
TGCATGCATGCATGCAT	SCAGTTGCATG			
TGCATGCATGCATGCA	TGCATGCATGC			
GCATGCATGCATGCA	TGCATGCATGCA			
CATGCATGCATGCAT	GCATGCATGCA			
TGCATGCAT	GCATGGCATGCA			
"GCATGC	ATGCATGCATGC			
CATGCA	TGCATGCATGCA			
TGCATO	SCACTGCATGCAT6	- 5		
TCCAT	GCCATGCAATGCAT. AT	GCATE AT		
T	GCATGCATGCATGCATGC	ATGCATE		
	SCATGCAGGTTGCATGCA	TECATE		
	'GCATGCATGCATGCATG	CATECAL ATECA	TECATEC	
	ATGCATGCATGCATGCG	CATGCATGCATCGCATGC	ATCGCA	
	IGCATGCATGCATGCATG	CATGEATGEATGEATGEATGEA	TGCATG	
	(CCATECAATECATECATE	CATECATECATECATEC	ATGCATG	
THE	ΑΤΑΓΑΤΑΓΑΤΑΓΑΤΑΓΑΤΑ	CATECATECATECATEC	ATCCATC	
110C/	CATECATECATECATECAT	CATCCATCCATCCATCC	ATCCAT	
			CATCCAT	
	COMPENSION	ATCCATCCATCCA	ATCCAT	
CA.		CATCCAT		
ATC C	JCP			
CCC .			AA/ 77	
A 1 6 1		GCATGCA		
All		G (/	r	

GD Genetic Diversity Centre Zurich

Genetic Diversity: Andlysis

STI S











- 1 Sample Quality Control
- 2 Library Quality Control
- 3 Run Quality Control
- 4 Sequencing Quality Control
- 5 Data Quality Control (e.g. bias, outlier)



- 1 Sample Quality Control
- 2 Library Quality Control
- 3 Run Quality Control
- 4 Sequencing Quality Control
- 5 Data Quality Control







GD Genetic Diversity Centre Zurich



Illumina Sequencing Analysis Viewer 1.8.37 - 151204_M03930_0010_00000000-AK6GL

Sequencing Analysis Viewer

Run Folder: D:\Illumina\MiSeqAnalysis\151204_M03930_0010_00000000-AK6GL Refresh Browse Analysis Imaging Summary Tile Status TruSeq Controls Indexing Cycle All Surface Both Swath All Lane All 🔾 А 🔵 С 🚫 С 🔘 Т Ŧ 7 💣 1 1 🗎 📴 Index Cycle Surface Swath Time P90 A Lane Tile Section 363 📤 Тор 12/04/201. Top 12/04/201 12/04/201 ... Top Тор 12/04/201 Тор 12/04/201 ... 12/04/201 ... Top 12/04/201 ... Top Top 12/04/201 Top 12/04/201 12/04/201 Тор Top 12/04/201 ... <u>n</u> Ó Top 12/04/201 ... Тор 12/04/201 Top 12/04/201 12/04/201 ... Top 12/04/201 Top Top 12/04/201 -× Rows=14504 Disp=14504 Sel=1 Filter

GD Genetic Diversity Centre Zurich

GDC Genetic Diversity Centre Zurich

Cluster Density: 1017 K/mm2 (optimal 865-965 k/mm²) Reads Total: 27.69 M (goal 30 M) Reads PF: 21.60 M (78%) PhiX Conc: 2.03 % (loaded 2%) %>=Q30: Total 63.06% (should be at least 70%)

- The **density** of clusters for each tile (in thousands per mm²) and the number of **clusters** for each tile (in millions).
- Total **yield** is the number of bases generated in the run.
- The calculated **error rate**, as determined by a spiked in PhiX control sample if available and it refers to the percentage of bases called incorrectly at any one cycle.
- The total fraction of passing filter reads (**PF**) assigned to an index.
- % Q-score >= Q30 (percentage of bases that have a Q-score above or equal to 30; Q30 is a probability of incorrect base calling of 1 in 1000).
- The **signal to noise ratio** is calculated as mean called intensity divided by standard deviation of non-called intensities. Not calculated for NextSeq two-channel sequencing or HiSeq X.
- The percentage of molecules in a cluster for which sequencing falls behind (**phasing**) or jumps ahead (**prephasing**) the current cycle within a read.







Fasta

- Fastq (Fasta with Quality)
- Bam (PacBio)
- ► POD5/FAST5 (ONT)



Sequence Data Format: Fasta (>)



 $N_{rows} \ge 2$ per sequences read

Sequence (nucleotide or protein)

File extension: sequence(s).fa, sequence(s).fasta

Special cases: sequences.**mfa** (multiple - aligned - sequences) sequences.**afa** (aligned sequences)

(Illumina) Sequence Data Format: Fastq (@)



ASCII encoded quality scores per base

rich



Illumina Fastq Header Format (version > 1.8)



a. unique instrument name

- b. run id
- c. flowcell id
- d. flowcell lane
- e. tile number within the flowcell lane
- f. x-coordinate of the cluster within the tile
- g. y-coordinate of the cluster within the tile

h. the member of a pair, 1 or 2 (paired-end or mate-pair reads only)

- i. Y if the read fails filter (read is bad), N otherwise (read passed filter)
- j. 0 when no control bits are on
- k. index sequence

Older Illumina Fastq Header Format (version < 1.8)

- a. unique instrument name
- b. flowcell lane
- c. tile number within the flowcell lane
- d. x-coordinate of the cluster within the tile
- e. y-coordinate of the cluster within the tile
- f. index number for a multiplexed sample (0 for no indexing)
- g. the member of a pair, /1 or /2 (paired-end or mate-pair reads only)

GDG Genetic Diversity Centre Zurich

PacBio (Revio) Fastq Header Format

BAM → pbtk::bam2fastx → FASTQ

+

PacBio (Revio) Fastq Header Format

The QNAME for by-strand CCS reads includes a suffix **fwd** or **rev** to indicate strand relative to the other by-strand read for the ZMW. Strand assignment by CCS is arbitrary and does not imply the strand that may be assigned during mapping.

GADE Genetic Diversity Centre Zurich

ONT Fastq Header Format

+

#\$\$#\$\$%'('*/001:<;?

<<;920000==::;;;QHIHHKLFKSSNSKNJMF?;78888??;BILGKJIKRMHHHSGSJKKSHLPSISLSSMKNLHJIPHSHHPSPMSHNSSMQSSPM NKJHGJSSSNSHIPLJHIKSOHKLJHHGJSGOLILSMMSNIL9999;ISNKMS99999SKSHSSISJSIHIMKSHJLSSNMSKHLSSMKKKIIHKSPMJG KSSJSSJSJSIKSIGIJNJEFJ755)

The headers in FASTQ files generated by Oxford Nanopore Technologies (ONT) basecallers are defined by the specific basecaller used (such as Dorado, Guppy, Albacore, Bonito, etc.) and have changed many times over the years.

GD Genetic Diversity Centre Zurich

ONT Fastq Header Format

```
@c66cf9d5-8ceb-49a2-87ba-573bdf744476
runid=da30afcc85dd427c946e129369e698f701722998
read=19
ch=21
start_time=2023-12-13T17:44:43.197441+01:00
flow_cell_id=FAY08440
protocol_group_id=p1023_run240512_Meta
sample_id=Shower-123sw
parent_read_id=c66cf9d5-8ceb-49a2-87ba-573bdf744476
basecall_model_version_id=dna_r10.4.1_e8.2_400bps_sup@v4.2.0
```

The FASTQ headers generated by Dorado, like those from other ONT basecallers, are designed to provide comprehensive metadata about each read. Understanding the specific format used by the version of Dorado you are working with is critical to effective data management and analysis.







Genetic

iversity

Zurich

entre



Fastq Sequence Read

- 1 @HWI-ST486:166:C06K9ACXX:7:1101:1443:1995 1:N:0:ACAGTG
- 2 ACTGAGCGTGGGCGAGCCGCACGGCACCATCCTCTGGCACACCCTCTCCTC
- 3 +

 分 ASCII **encoded** quality scores per base



Sequencing (phred) quality scores (Q) measures the (error) probability (P) that a base is called (in)correctly.

position 1 2 3 4 ... nucleotide A C G T ... quality score (Q) 20 20 22 21 ...

https://www.phrap.com/phred/



Sequencing (<u>phred</u>) **quality scores (Q)** measure the (error) **probability (P)** that a base is called incorrectly.

Base-Calling Error Probability

$$P = 10^{\frac{-Q}{10}}$$

Phred Quality Score

$$Q = -10\log_{10} P$$



Sequencing (<u>phred</u>) **quality scores (Q)** measures the (error) **probability (P)** that a base is called incorrectly.

position	1	2	3	4	• • •		
nucleotide	Α	С	G	Τ	• • •		
quality score (Q)	20	20	22	21	• • •		
$P = 10^{\frac{-Q}{10}} = 10^{-2} = 0.01 \rightarrow 1\% \ (\rightarrow 99\%)$							



Sequencing (<u>phred</u>) **quality scores (Q**) measure the (error) **probability (P)** that a base is called incorrectly.

position	1	2	3	4	•••
nucleotide	Α	С	G	Τ	• • •
quality score (Q)	20	20	22	21	• • •
probability (P)	0.01	0.01	0.006	0.008	•••
accuracy (1-P)	0.99	0.99	0.994	0.992	•••



One character encoding!

1234	1234	1234
ACGT	or ACGT	or ACGT
a a c b	+ + * "	• • • •

 $20 \rightarrow a$ $20 \rightarrow +$ $20 \rightarrow \odot$

 $21 \rightarrow b$ $21 \rightarrow "$ $21 \rightarrow \bigcirc$

 $22 \rightarrow c$ $22 \rightarrow *$ $22 \rightarrow \Rightarrow$

1 2 3 4 A C G T 20 20 22 21

GD Genetic Diversity Centre Zurich

ASCII TABLE

Decimal	Hex	Char	Decimal	Hex	Char	JDecimal	Hex	Char	Decimal	Hex	Char
0	0	[NULL]	32	20	[SPACE]	64	40	0	96	60	`
1	1	[START OF HEADING]	33	21	1	65	41	Α	97	61	а
2	2	[START OF TEXT]	34	22		66	42	В	98	62	b
3	3	[END OF TEXT]	35	23	#	67	43	С	99	63	с
4	4	[END OF TRANSMISSION]	36	24	\$	68	44	D	100	64	d
5	5	[ENQUIRY]	37	25	%	69	45	E	101	65	е
6	6	[ACKNOWLEDGE]	38	26	&	70	46	F	102	66	f
7	7	[BELL]	39	27	1	71	47	G	103	67	g
8	8	[BACKSPACE]	40	28	(72	48	н	104	68	h
9	9	[HORIZONTAL TAB]	41	29)	73	49		105	69	i
10	А	[LINE FEED]	42	2A	*	74	4A	J	106	6A	j
11	В	[VERTICAL TAB]	43	2B	+	75	4B	Κ	107	6B	k
12	С	[FORM FEED]	44	2C	,	76	4C	L	108	6C	1
13	D	[CARRIAGE RETURN]	45	2D	-	77	4D	M	109	6D	m
14	E	[SHIFT OUT]	46	2E		78	4E	Ν	110	6E	n
15	F	[SHIFT IN]	47	2F	1	79	4F	0	111	6F	0
16	10	[DATA LINK ESCAPE]	48	30	0	80	50	Р	112	70	р
17	11	[DEVICE CONTROL 1]	49	31	1	81	51	Q	113	71	q
18	12	[DEVICE CONTROL 2]	50	32	2	82	52	R	114	72	r
19	13	[DEVICE CONTROL 3]	51	33	3	83	53	S	115	73	S
20	14	[DEVICE CONTROL 4]	52	34	4	84	54	т	116	74	t
21	15	[NEGATIVE ACKNOWLEDGE]	53	35	5	85	55	U	117	75	u
22	16	[SYNCHRONOUS IDLE]	54	36	6	86	56	V	118	76	v
23	17	[ENG OF TRANS. BLOCK]	55	37	7	87	57	W	119	77	w
24	18	[CANCEL]	56	38	8	88	58	X	120	78	x
25	19	[END OF MEDIUM]	57	39	9	89	59	Y	121	79	У
26	1A	[SUBSTITUTE]	58	ЗA		90	5A	Z	122	7A	z
27	1B	[ESCAPE]	59	3B	;	91	5B	[123	7B	{
28	1C	[FILE SEPARATOR]	60	3C	<	92	5C	١	124	7C	
29	1D	[GROUP SEPARATOR]	61	3D	=	93	5D	1	125	7D	}
30	1E	[RECORD SEPARATOR]	62	3E	>	94	5E	^	126	7E	~
31	1F	[UNIT SEPARATOR]	63	3F	?	95	5F	_	127	7F	[DEL]

Illumina Quality Encoding (version > 1.8+)

Decimal	Hex	Char	JDecimal	Hex	Char
32	20	[SPACE]	64	40	@
33	21	1	65	41	Α
34	22		66	42	В
35	23	#	67	43	С
36	24	\$	68	44	D
37	25	%	69	45	E
38	26	&	70	46	F
39	27	1.00	71	47	G
40	28	(72	48	H
41	29)	73	49	1.1
42	2A	*	74	4A	J
43	2B	+	75	4B	K
44	2C	,	76	4C	L
45	2D	-	77	4D	M
46	2E	1.0	78	4E	N
47	2F	1	79	4F	0
48	30	0	80	50	Р
49	31	1	81	51	Q
50	32	2	82	52	R
51	33	3	83	53	S
52	34	4	84	54	т
53	35	5	85	55	U
54	36	6	86	56	V
55	37	7	87	57	W
56	38	8	88	58	X
57	39	9	89	59	Y
58	ЗA	1.0	90	5A	Z
59	3B	;	91	5B	[
60	3C	<	92	5C	λ
61	3D	=	93	5D	1
62	3E	>	94	5E	^
63	3F	?	95	5F	_

Q + 33 = ASCII $20 + 33 = 53 \rightarrow 5$

position	1	2	3	4	•••
nucleotide	A	С	G	Т	•••
quality score Q	20	20	22	21	•••
Ascii	53	53	55	54	
char endcoding	5	5	7	6	•••

G٢

Centre

Diversity

Genetic

Zurich

GDG Genetic Diversity Centre Zurich

SSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSSS	SSSSSSSSSSS	SSSS		
x	*****		*****	
			3333333333	
LLLLLLLLLLLLLLLLLLLLLLLLLLLL	LLLLLLLLL	LLLL		
! #\$%&'()*+/0123456789::	<=>?@ABCDE	FGHIJKLMNOPORSTUVWXYZ(\)^	abcdefghiiklmnopgr	stuvwxvz{ }~
33 50	64	73	104	126
0 26	31	40	101	120
_5	0	9	4.0	
	0	9	40	
	3		40	
0.2		41		
S - Sanger Phred+33,	raw reads	s typically (0, 40)		
X - Solexa Solexa+64,	raw reads	s typically (-5, 40)		
<pre>I = Illumina 1.3+ Phred+64,</pre>	raw reads	typically (0, 40)		
J - Illumina 1.5+ Phred+64,	raw reads	typically (3, 40)		
with 0=unused, 1=unused,	2=Read Se	gment Quality Control Ind	licator (bold)	
(Note: See discussion ab	ove).			
L - Illumina 1.8+ Phred+33,	raw reads	stypically (0, 41)		

GDC Genetic Diversity Zurich



```
# ascii character > decimal value
asc <- function(x) {
    strtoi(charToRaw(x),16L)
    }
asc("!")</pre>
```

```
# decimal value > ascii character
chr <- function(n) {
        rawToChar(as.raw(n))
        }
chr("33")</pre>
```

GDA | 21.06.2024 | JCW

Encoding	ASCII	Q	Р
J	74	41	0.00008
I	73	40	0.00010
Н	72	39	0.00013
G	71	38	0.00016
F	70	37	0.00020
Е	69	36	0.00025
D	68	35	0.00032
С	67	34	0.00040
В	66	33	0.00050
А	65	32	0.00063
@	64	31	0.00079
?	63	30	0.00100
>	62	29	0.00126
=	61	28	0.00158
<	60	27	0.00200
;	59	26	0.00251
:	58	25	0.00316
9	57	24	0.00398
8	56	23	0.00501
7	55	22	0.00631
6	54	21	0.00794
5	53	20	0.01000
4	52	19	0.01259
3	51	18	0.01585
2	50	17	0.01995
1	49	16	0.02512
0	48	15	0.03162
/	47	14	0.03981
	46	13	0.05012
-	45	12	0.06310
,	44	11	0.07943
+	43	10	0.10000
*	42	9	0.12589
)	41	8	0.15849
(40	7	0.19953
	39	6	0.25119
&	38	5	0.31623
%	37	4	0.39811
\$	36	3	0.50119
#	35	2	0.63096
п	34	1	0.79433
!	33	0	1.00000

Phred Quality Score

 $Q = -10\log_{10}P$

Base-Calling Error Probability

 $P = 10^{\frac{-Q}{10}}$

Q	Р	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%

GDA | 21.06.2024 | JCW

GS

Centre

Zurich

Genetic

Diversity

GD Genetic Diversity Centre Zurich

Phred Scores per Base



GDC Genetic Diversity Centre Zurich



Quality scores across all bases (Sanger / Illumina 1.9 encoding)

GDC Genetic Diversity Centre Zurich





Why do we use a box plot and not a bar plot?



Why do we use a box plot and not a bar plot?

Box plots are generally more informative than bar charts, especially when you need to understand the distribution and variability of your data.
GD Genetic Diversity Centre Zurich



https://pagepiccinini.com/2016/02/23/boxplots-vs-barplots/

GDC Genetic Diversity Centre Zurich

Histogram



mean:100 sd: 20

https://pagepiccinini.com/2016/02/23/boxplots-vs-barplots/

38

GD Genetic Diversity Centre Zurich



https://pagepiccinini.com/2016/02/23/boxplots-vs-barplots/

type

39



GDC

Centre

Zurich

Diversity

Genetic

GSDC Genetic Diversity Centre Zurich

Symmetrical Distribution mean = median = mode



skewed to the left

Quality scores across all bases (Sanger / Illumina 1.9 encoding) <0.1% **Error Probability** Phred Score 0.1% 1% 10% >10% 15-19 25-29 35-39 45-49 55-59 65-69 75-79 85-89 95-99 Position in read (bp)

GDC

Centre

Zurich

Diversity

Genetic

Quality score; across all bases (Sanger / Illumina 1.9 encoding) 5-19 3 4 5 6 7 8 25-29 35-39 45-49 55-59 65-69 75-79 85-89 95-99 Position in read (bp)

GDC

Centre

Zurich

Diversity

Genetic

Phred Score

Phred Score 6 7 8 15-19 25-29 35-39 45-49 55-59 65-69 75-79 85-89 95-99 Position in read (bp)

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Note: Error rate increases along the length of the read.

GDC

iversity

entre

Zurich

Genetic



Q_{mean phred score}: 30

I 50nt

% Q-score >= Q30 (percentage of bases that have a Q-score above or equal to 30; Q30 is a probability of incorrect base calling of 1 in 1000).























	つ	
	$= \prec$	()
mean		

35 25 Q_{mean} =30 35 15



10 position









error rate:

μ: 0.05 0.1 0.3





error rate:

μ: 0.05 0.1 0.3



5 (25%) error free reads



Q20 - 150nt

99% Accuracy / 1% Error Rate 0.99¹⁵⁰ → **22% Error Free Reads**



For Better or Worse



54







GD Genetic Diversity Centre Zurich

Error Correction



- Read quality
- Number of reads (coverage)

GD Genetic Diversity Centre Zurich

Error Correction



Read quality Number of reads (coverage) Phred score



GDC Genetic Diversity Centre Zurich

Schirmer et al. BMC Bioinformatics (2016) 17:125 DOI 10.1186/s12859-016-0976-y

BMC Bioinformatics

RESEARCH ARTICLE



Illumina error profiles: resolving fine-scale variation in metagenomic sequencing data

Melanie Schirmer^{1,2,4*}, Rosalinda D'Amore³, Umer Z. Ijaz⁴, Neil Hall³ and Christopher Quince⁵

Abstract

Background: Illumina's sequencing platforms are currently the most utilised sequencing systems worldwide. The technology has rapidly evolved over recent years and provides high throughput at low costs with increasing read-lengths and true paired-end reads. However, data from any sequencing technology contains noise and our understanding of the peculiarities and sequencing errors encountered in Illumina data has lagged behind this rapid development.

Results: We conducted a systematic investigation of errors and biases in Illumina data based on the largest collection of in vitro metagenomic data sets to date. We evaluated the Genome Analyzer II, HiSeq and MiSeq and tested state-of-the-art low input library preparation methods. Analysing in vitro metagenomic sequencing data allowed us to determine biases directly associated with the actual sequencing process. The position- and nucleotide-specific analysis revealed a substantial bias related to motifs (3mers preceding errors) ending in "GG". On average the top three motifs were linked to 16 % of all substitution errors. Furthermore, a preferential incorporation of ddGTPs was recorded. We hypothesise that all of these biases are related to the engineered polymerase and ddNTPs which are intrinsic to any sequencing-by-synthesis method. We show that quality-score-based error removal strategies can on average remove 69 % of the substitution errors - however, the motif-bias remains. **Conclusion:** Single-nucleotide polymorphism changes in bacterial genomes can cause significant changes in phenotype, including antibiotic resistance and virulence, detecting them within metagenomes is therefore vital. Current error removal techniques are not designed to target the peculiarities encountered in Illumina sequencing data and other sequencing-by-synthesis methods, causing biases to persist and potentially affect any conclusions drawn from the data. In order to develop effective diagnostic and therapeutic approaches we need to be able to **identify systematic sequencing errors and distinguish these errors from true genetic variation**.



Orig.nuc. ZAZCZGZT

Substitutions

Illumina

Average substitution rates

Platform	R1/R2	A	с	G	т
GAII	R1	0.0015	0.0010	0.0008	0.0018
GAII	R2	0.0035	0.0029	0.0019	0.0026
HiSeq	R1	0.0004	0.0004	0.0004	0.0008
HiSeq	R2	0.0007	0.0007	0.0007	0.0012
MiSeq	R1	0.0012	0.0009	0.0009	0.0012
MiSeq	R2	0.0033	0.0021	0.0015	0.0031

EE(R1) < EE(R2)

 $EE(Sub) > EE(Ins) \approx EE(Del)$

 $EE(HiSeq) < EE(MiSeq) \approx EE(GAII)$

GΣ

Centre

ienetic

iversity

Zurich

59

Orig.nuc.





Error Rate 8-12%

BAM → FASTQ BAM → CCS.FASTX

GDA | 21.06.2024 | JCW

Circular Consensus Sequences (CCS)







Why does it not improve anymore?

p: base-calling error probability

0

GDG Genetic Diversity Centre Zurich



GSDC Genetic Diversity Centre Zurich





	Mappable length (bp)				Error rate (Proportion of overall error) (%)			
Read type	Mean	Median	Standard deviation	Maximum	Overall	Insertion	Deletion	Mismatch
PacBio CCS	1772	1464	1132	8006	1.72	0.087 (5.06)	0.34 (19.48)	1.30 (75.46)
PacBio subread	1570	1299	1076	16040	14.20	5.92 (41.71)	3.01 (21.17)	5.27 (37.12)
ONT 2D	1861	1754	882	9126	13.40	3.12 (23.30)	4.79 (35.70)	5.50 (40.99)
ONT 1D	1695	1602	824	9345	20.19	2.93 (14.51)	7.52 (37.24)	9.74 (48.25)

MITOCHONDRIAL DNA PART B: RESOURCES 2019, VOL. 4, NO. 1, 408–409 https://doi.org/10.1080/23802359.2018.1547133



ARTICLE

OPEN ACCESS

Long-read sequencing of benthophilinae mitochondrial genomes reveals the origins of round goby mitogenome re-arrangements

Silvia Gutnik^a, Jean-Claude Walser^b and Irene Adrian-Kalchhauser^c

^aBiozentrum, Department Growth & Development, University of Basel, Basel, Switzerland; ^bGenetic Diversity Centre Zurich, ETH Zurich, Zurich, Switzerland;^cProgram Man-Society-Environment, Department of Environmental Sciences, University of Basel, Basel, Switzerland



Origin of the re-arranged **tRNA cluster** Gln, Ile, Met. Most Gobiidae carry the arrangement Ile, Gln, Met without spacers. Benthophilinae (subfamily of gobies) however carry the arrangement Gln, Ile, Met, and feature variable length spacers between the genes.

Zurich

Centre

Genetic

iversit





Illumina MiSeq R1 -> EE mean 0.1 Illumina MiSeq R2 -> EE mean 0.2 PacBio Sequel IIe -> EE mean 2.9 PacBio Revio -> EE mean 1.2 ONT MinION -> EE mean 30.7





FastQC

(http://www.bioinformatics.babraham.ac.uk/projects/fastqc/)

FASTX-Toolkit

(http://hannonlab.cshl.edu/fastx_toolkit/)

USEARCH

(https://www.drive5.com/usearch/)

PRINSEQ

(http://edwards.sdsu.edu/cgi-bin/prinseq/prinseq.cgi)

Galaxy

(http://galaxyproject.org)

Rqc

(https://bioconductor.org/packages/release/bioc/vignettes/Rqc/inst/doc/Rqc.html)

CLC Genomic Workbench

(http://www.clcbio.com/products/clc-genomics-workbench/)

Geneious (http://www.geneious.com/)

GD

Centre

Zurich

Diversity

Genetic

68

GDC Genetic Diversity Centre Zurich





Nucleic Acids Research Advance Access published December 16, 2009

Nucleic Acids Research, 2009, 1–5 doi:10.1093/nar/gkp1137

The Sanger FASTQ file format for sequences with quality scores, and the Solexa/Illumina FASTQ variants

Peter J. A. Cock^{1,*}, Christopher J. Fields², Naohisa Goto³, Michael L. Heuer⁴ and Peter M. Rice⁵

ASCII encoded quality scores



GDC Genetic Diversity Centre Zurich

ASCII stands for American Standard Code for Information Interchange. Computers can only understand numbers, so an ASCII code is the numerical representation of a character or an action of some sort.

	DEC	ост	HEX	BIN	Symbol	HTML Number	HTML Name	Description
ĺ	32	040	20	00100000				Space
	33	041	21	00100001	!	!		Exclamation mark
	34	042	22	00100010	"	"	"	Double quotes (or speech marks)
	35	043	23	00100011	#	#		Number
	36	044	24	00100100	\$	\$		Dollar
	37	045	25	00100101	%	%		Procenttecken
	38	046	26	00100110	8.	&	&	Ampersand
	39	047	27	00100111	1.1	'		Single quote
	40	050	28	00101000	((Open parenthesis (or open bracket)
	41	051	29	00101001))		Close parenthesis (or close bracket)
	42	052	2A	00101010	*	*		Asterisk
	43	053	2B	00101011	+	+		Plus
	44	054	2C	00101100	,	,		Comma
	45	055	2D	00101101	-	-		Hyphen
	46	056	2E	00101110		.		Period, dot or full stop
	47	057	2F	00101111	1	/		Slash or divide
	48	060	30	00110000	0	0		Zero
	49	061	31	00110001	1	1		One
	50	062	32	00110010	2	2		Two
	51	063	33	00110011	3	3		Three
	52	064	34	00110100	4	4		Four
	53	065	35	00110101	5	5		Five
	54	066	36	00110110	6	6		Six
	55	067	37	00110111	7	7		Seven

GDA | 21.06.2024 | JCW



Nucleic Acids Research Advance Access published January 13, 2015

Nucleic Acids Research, 2015 1 doi: 10.1093/nar/gku1341

Insight into biases and sequencing errors for amplicon sequencing with the Illumina MiSeq platform

Melanie Schirmer^{1,*}, Umer Z. Ijaz¹, Rosalinda D'Amore², Neil Hall², William T. Sloan¹ and Christopher Quince¹

¹School of Engineering, University of Glasgow, Glasgow, UK and ²Functional and Comparative Genomics, University of Liverpool, Liverpool, UK

With read lengths of currently up to 2 × 300 bp, high throughput and low sequencing costs Illumina's MiSeq is becoming one of the most utilized sequencing platforms worldwide. The platform is manageable and affordable even for smaller labs. This enables quick turnaround on a broad range of applications such as targeted gene sequencing, metagenomics, small genome sequencing and clinical molecular diagnostics. However, **Illumina error profiles are still poorly understood and programs are therefore not designed for the idiosyncrasies of Illumina data.** A better knowledge of the error patterns is essential for sequence analysis and vital if we are to draw valid conclusions. Studying true genetic variation in a population sample is fundamental for understanding diseases, evolution and origin. We conducted a large study on the error patterns for the MiSeq based on 16S rRNA amplicon sequencing data. We tested state-of-the-art library preparation methods for amplicon sequencing and showed that the library preparation method and the choice of primers are the most significant sources of bias and cause distinct error patterns. Furthermore we tested the efficiency of various error correction strategies and identified quality trimming (Sickle) combined with error correction (BayesHammer) followed by read overlapping (PANDAseq) as the most successful approach, reducing substitution error rates on average by 93%.
MPS > **Quality Control**





Wenger et al. (2019) Highly-accurate long-read sequencing improves variant detection and assembly of a human genome.



A SMRTbell library tightly distributed at 15 kb was chosen for circular consensus sequencing based on estimates of 150 kb polymerase read length and a requirement of 10 passes to achieve Q30 read accuracy. CCS reads with a predicted accuracy of at least Q20 (99%) were retained. The total CCS read yield was 89 Gb, an average of 2.3 Gb per SMRT Cell, with an average read length of 13.5 kb ± 1.2 kb. The predicted accuracy of the CCS reads has a median of Q30 (99.9%) and a mean of Q27 (99.8%).

MPS > **Quality Control**

MITOCHONDRIAL DNA PART B: RESOURCES 2019, VOL. 4, NO. 1, 408–409 https://doi.org/10.1080/23802359.2018.1547133

ARTICLE



Long-read sequencing of benthophilinae mitochondrial genomes reveals the origins of round goby mitogenome re-arrangements

Silvia Gutnik^a, Jean-Claude Walser^b and Irene Adrian-Kalchhauser^c

^aBiozentrum, Department Growth & Development, University of Basel, Basel, Switzerland; ^bGenetic Diversity Centre Zurich, ETH Zurich, Zurich, Switzerland;^cProgram Man-Society-Environment, Department of Environmental Sciences, University of Basel, Basel, Switzerland



Zurich

entre

iversity

Genetic