Sequence Formats: FASTA, FASTQ PHRED Scores coded as ASCII Characters

White Space (Space or Tab)

> Sequence_Name_1 Sequence Annotation

MDSKGSSQKGSRLLLLLVVSNLLLCQGVVSTPVCPNGPGNCQVSLRDLFDRAVMVSHYIHDLSS EMFNEFDKRYAQGKGFITMALNSCHTSSLPTPEDKEQAQQTHHEVLMSLILGLLRSWNDPLYHL VTEVRGMKGAPDAILSRAIEIEEENKRLLEGMEMIFGQVIPGAKETEPYPVWSGLPSLQTKDED ARYSAFYNLLHCLRRDSSKIDTYLKLLNCRIIYNNNC

> Sequence Name 2 Sequence Annotation

ADQLTEEQIAEFKEAFSLFDKDGDGTITTKELGTVMRSLGQNPTEAELQDMINEVDADGNGTID FPEFLTMMARKMKDTDSEEEIREAFRVFDKDGNGYISAAELRHVMTNLGEKLTDEEVDEMIREA DIDGDGQVNYEEFVQMMTAK

> Sequence_Name_3 | Sequence Annotation

LCLYTHIGRNIYYGSYLYSETWNTGIMLLLITMATAFMGYVLPWGQMSFWGATVITNLFSAIPY IGTNLVEWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIP FHPYYTIKDFLGLLILLLLLLLLLLSPDMLGDPDNHMPADPLNTPLHIKPEWYFLFAYAILRS VPNKLGGVLALFLSIVILYGLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYP YTIIGQMASILYFSIILAFLPIAGXIENY

> Sequence Name_1 Sequence Annotation

MDSKGSSQKGSRLLLLLVVSNLLLCQGVVSTPVCPNGPGNCQVSLRDLFDRAVMVSHYIHDLSS
EMFNEFDKRYAQGKGFITMALNSCHTSSLPTPEDKEQAQQTHHEVLMSLILGLLRSWNDPLYHL
VTEVRGMKGAPDAILSRAIEIEEENKRLLEGMEMIFGQVIPGAKETEPYPVWSGLPSLQTKDED
ARYSAFYNLLHCLRRDSSKIDTYLKLLNCRIIYNNNC

- > Sequence Name 2 Sequence Annotation
- ADQLTEEQIAEFKEAFSLFDKLGDGTTKELGTVMRSLGQNPTEAELODMINEVDADGNGTIDEFFLTMMARKMKDTDSEEELREAFRVFDKLVGYISAAELRHVMTNLGEKLTDEEVDEMIREADDIDGDGQVNYEEFVQMMTAK
- > Sequence Name 3 Sequence Annotation

LCLYTHIGRNIYYGSYLYSETWNTGIMLILITMATAFMGYVLPWGQM FWGATVITNLFSAIPY IGTNLVEWIWGGFSVDKATLNRFFAFHFILPFTMVALAGVHLTFLHETGSNNPLGLTSDSDKIP FHPYYTIKDFLGLLILILLLLLLSPDMLGDPDNHMPADPLNTFLHIKPEWYFLFAYAILRS VPNKLGGVLALFLSIVILYGLMPFLHTSKHRSMMLRPLSQALFWTLTMDLLTLTWIGSQPVEYP YTIIGQMASILYFSIILAFLPIAGXIENY

> Sequence Name 1 Sequence Annotation

FASTA Format is designed to store **DNA**, **RNA**, **Protein** Sequences with minimal **Annotation**.

Allowing **DNA** & **Protein Ambiguity Codes**, can make it difficult to ascertain **Sequence Type**.

Sequence Type is not specified in the **FASTA** format. Distinction is left to the software.

FASTQ Format is an adaptation designed to store Sequencing Reads with minimal Annotation.

Each sequence in a **FASTQ** file is always the that of a single **Sequencing Read**. Thus it can only ever be **DNA** sequence.

CCGCTAGCTGGGGTATCATCAGCATGCATGGCATGAGCGTTCTTAATTCTCAGGGACTCGGAGCAGGGCATCGAG

The transition from FASTA to FASTQ Format includes a couple of cosmetic "enhancements":

The initial ">" becomes an "@"

An extra **Annotation** line beginning with a "+" is added.

The opportunity to include **Annotation** in this line is rarely used.

Currently, as the length of an individual read is modest, sequences occupy just a single line.

The **Q** in **FASTQ** stands for **Quality**. A definitive purpose of the **FASTQ** Format is to record the **Base Call Quality** of each element of a **Sequencing Read**.

Base Call Qualities are recorded in a fourth line of the FASTQ file. Each Base Call Quality is recorded as a single printable character. Each Quality character corresponds to one Called Base.

```
Sequence Name 1 Sequence Annotation ... ... ... ... ...
CCGCTAGCTGGGGTATCATCAGCATGCATGGCATGAGCGTTCTTAATTCTCAGGGACTCGGAGCAGGGCA
Spurious Start of new read
!''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>CCCCCCC655676CCGLN6
@5KP455CCF>>%%%%%++)5676CCGL'*(((***+)118477GAAADVV++BBDAGGH7GGC9%%%)
ACGGTA (HG48** (@@HIA90%%++) (%>>>>>CCJJAS!!!!) (^++K''``' ''*(''*(LK*&TT
+F9RSA;;::ACV6^7&&&3£22"(1!DFGH((((<<<<%++)(%%%%).1* HIA90%%++)(%>>>
           Spurious Annotation line
                                    Spurious Read Sequence
```

Most printable characters, including "@" and "+" are allowable Base Call Quality codes.

Multiple line sequences could introduce Base Call Quality lines beginning with "@" or "+".

Such lines could be misinterpreted as: "@" - The start of a new Sequencing Read,

"+" - An additional **Annotation** line.

```
@ Sequence Name 1 Sequence Annotation ... ... ... ... ... ...
CCGCTAGCTGGGGTATCATCAGCATGCATGGCATGAGCGTTCTTAATTCTCAGGGACTCGGAGCAGGGCATCGAG
!''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>CCCCCCC655676CCGLN695KP4 ...
@ Sequence Name 2 Sequence Annotation ... ... ... ... ... ...
CCGCTAGCTGGGGTATCATCAGCATGCATGGCATGAGCGTTCTTAATTCTCAGGGACTCGGAGCAGGGCATCGAG
(**-+((***+))%>>CCCC!CCC655676CCGL!''*(N695KP4%%++)(%%%).1**''))**55CCF>>>> ...
@ Sequence Name 3 Sequence Annotation ... ... ... ... ... ...
CATGAGCGTCAGCATGCATGGTCGGAGCAGGGCATCGAGTCTTAATTCTCAGGGACCCGCTAGCTGGGGTATCAT
***-+*''))**55CCF>>FF>>!''*((((***+))%1CCACC%%++)(%%%).CCC655676CCGLN695KP4 ...
```

To avoid this, even longer **Read Sequences** are confined to one line.

The Base Quality Scores in FASTQ files are PHRED Scores.

The **PHRED Score** for a given **Base Call** is derived from an estimate of the **Probability** of that **Base** being **Called** incorrectly

Where **Q** is the **PHRED Score** for a given **Base Call** and **P** is the estimated **Probability of Error** for that **Base Call**, the following formula applies:

$$Q = -10 * log_{10}(P)$$

P is a **Probability** and so has **Range** $0 \rightarrow 1$

Q is a function of **P** that has **Range** $0 \rightarrow$ **infinity**.

Base Call Error Rate	Probability of Incorrect Call (P)	Log ₁₀ (P)	PHRED Score $Q = -10Log_{10}(P)$
1 in 1	1	0	00
1 in 10	0.1	- 1	10
1 in 100	₩ 0.01	-2	20
1 in 1000	0.001	-3	30
1 in 10,000	0.0001	-4	40
1 in 100,000	0.00001	-5	50
1 in 1,000,000	0.00001	-6	60
1 in 10,000,000	0.000001	-7	70
 O		 -infinity	 infinity

Q is a function of **P** that has **Range** $0 \rightarrow$ **infinity**. and easily represents a useful subset of **P Values**, to an adequate accuracy, as a **Two Digit Integer**.

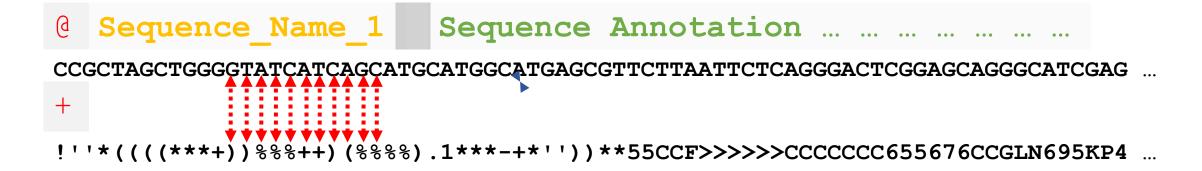
Base Call Error Rate	PHRED Score Q = -10Log ₁₀ (P)
1 in 1	00
1 in 10	10
1 in 100	20
1 in 1000	30
1 in 10,000	40
1 in 100,000	50
1 in 1,000,000	60
1 in 10,000,000	70
 O	 infinity

PHRED Scores can be intuitively thought of as directly representing various **Error Rates**.

An Error Rate worse than 1 in 100 (PHRED = 20) is usually considered a disaster!

Anything better than **1** in **10,000** (**PHRED** = **40**) is usually considered as near perfect as makes no difference.

Achieving a Consensus PHRED Score of 30 (1 in 1000) is a common target for an Assembly of Reads.



The Problem

Ideally at least, each base in a **FASTQ** file should correspond with one, *two digit* **PHRED Score**, represented as a **Single Character**.

The Solution

Is to map the *two digit* Phred Scores onto a single element of the ASCII Character Set.

The ASCII Character Codes are Stable and Universally Accepted.

They provide a mapping of the integers 000 → 255 to a set of Individual Characters.

Not all of these characters are visibly printable. Specifically, $00 \rightarrow 31$ are not printable.

32 is printable, but not visible (it is a **Space**).

33 → 126 <u>are</u> all printable however, and could be used to represent the PHRED Scores 00 → 93 nicely. More than sufficient for practical purposes.

0	<nul></nul>	32	<spc></spc>	64	@	96	`		128	Ä	160	+	192	خ	224	‡
1	<soh></soh>	33	!	65	Α	97	а		129	Å	161	0	193	i	225	
2	<stx></stx>	34	11	66	В	98	b		130	Ç	162	¢	194	\neg	226	,
3	<etx></etx>	35	#	67	С	99	С		131	É	163	£	195	\checkmark	227	,,
4	<eot></eot>	36	\$	68	D	100	d		132	Ñ	164	§	196	f	228	‰
5	<enq></enq>	37	%	69	Е	101	е		133	Ö	165	•	197	≈	229	Â
6	<ack></ack>	38	&	70	F	102	f		134	Ü	166	\P	198	Δ	230	Ê Á
7	<bel></bel>	39		71	G	103	g		135	á	167	ß	199	«	231	Á
8	<bs></bs>	40	(72	Н	104	h		136	à	168	R	200	>>	232	ËÈ
9	<tab></tab>	41)	73	I	105	i		137	â	169	©	201		233	È
10	<lf></lf>	42	*	74	J	106	j		138	ä	170	TM	202		234	Í
11	<vt></vt>	43	+	75	K	107	k		139	ã	171	,	203	À	235	Î
12	<ff></ff>	44	,	76	L	108	-1		140	å	172		204	Ã	236	Ϊ
13	<cr></cr>	45	-	77	M	109	m		141	Ç	173	≠	205	Õ	237	Ì
14	<s0></s0>	46		78	Ν	110	n		142	é	174	Æ	206	Œ	238	Ó
15	<si></si>	47	/	79	0	111	0		143	è	175	Ø	207	œ	239	Ô
16	<dle></dle>	48	0	80	Р	112	р		144	ê	176	∞	208	-	240	É
17	<dc1></dc1>	49	1	81	Q	113	q		145	ë	177	±	209	_	241	Ò
18	<dc2></dc2>	50	2	82	R	114	r		146	ĺ	178	≤	210	**	242	Ú
19	<dc3></dc3>	51	3	83	S	115	S		147	ì	179	≥	211	"	243	Û
20	<dc4></dc4>	52	4	84	Т	116	t		148	î	180	¥	212	`	244	Ù
21	<nak></nak>	53	5	85	U	117	u		149	Ϊ	181	μ	213	,	245	1
22	<syn< th=""><th>54</th><th>6</th><th>86</th><th>V</th><th>118</th><th>V</th><th></th><th>150</th><th>ñ</th><th>182</th><th>9</th><th>214</th><th>÷</th><th>246</th><th>^</th></syn<>	54	6	86	V	118	V		150	ñ	182	9	214	÷	246	^
23	<etb></etb>	55	7	87	W	119	W		151	ó	183	Σ	215	\Diamond	247	~
24	<can></can>	56	8	88	X	120	X		152	Ò	184	Π	216	ÿ	248	_
25		57	9	89	Υ	121	У		153	ô	185	П	217	Ϋ	249	v
26		58	:	90	Z	122	Z		154	Ö	186	ſ	218	/	250	•
27	<esc></esc>	59	;	91	[123	{		155	õ	187	а	219	€	251	٥
28	<fs></fs>	60	<	92	\	124			156	ú	188	0	220	<	252	,
29	<gs></gs>	61	=	93]	125	}		157	ù	189	Ω	221	>	253	"
30	<rs></rs>	62	>	94	^	126	~		158	û	190	æ	222	fi	254	·
31	<us></us>	63	?	95	_	127	<del< th=""><th>.></th><th>159</th><th>ü</th><th>191</th><th>Ø</th><th>223</th><th>fl</th><th>255</th><th>v</th></del<>	.>	159	ü	191	Ø	223	fl	255	v

Sequence Name 1 Sequence Annotation CCGCTAGCTGGGGTATCATCAGCATGCATGGCATGAGCGTTCTTAATTCTCAGGGACTCGGAGCAGGGCATCGAG !''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>CCCCCCC655676CCGLN695KP4 ... **ASCII Codes** → 62 53 PHRED Scores→ 10 38 30 21 46

Thus it is **ASCII Codes** that represent the **Base Call Qualities** (**PHRED Scores**) in **FASTQ** files.

To compute the **PHRED Score** from **the ASCII Character**, simply look up the **ASCII Code** and **SUBTRACT 32**.

A Warning

Long ago ... stupid people at the Sanger Centre coded **PHRED Scores** starting from **ASCII Code 64** (@).

Should you across any of this older data that has not been updated to reflect the current standards, you will need to instruct the software to subtract 63 instead of 32 in order to compute the PHRED Score from the ASCII Code.

0	<nul></nul>	32	<spc></spc>	64	@	96	`		128	Ä	16	0	†	192	خ	224	#
1	<soh></soh>	33	!	65	Α	97	а		129	Å	16	1	0	193	i	225	
2	<stx></stx>	34	11	66	В	98	b		130	Ç É	16	2	¢	194	\neg	226	,
3	<etx></etx>	35	#	67	С	99	С		131	É	16	3	£	195	\checkmark	227	"
4	<eot></eot>	36	\$	68	D	100	d		132	Ñ	16	4	§	196	f	228	‰
5	<enq></enq>	37	%	69	Е	101	е		133	Ö	16		•	197	≈	229	Â
6	<ack></ack>	38	&	70	F	102	f		134	Ü	16	6	\P	198	Δ	230	Ê
7	<bel></bel>	39	1	71	G	103	g		135	á	16	7	ß	199	«	231	Á
8	<bs></bs>	40	(72	Н	104	h		136	à	16	8	R	200	>>	232	Ë È
9	<tab></tab>	41)	73	I	105	i		137	â	16	9	©	201		233	
10	<lf></lf>	42	*	74	J	106	j		138	ä	17	0	TM	202		234	Í
11	<vt></vt>	43	+	75	Κ	107	k		139	ã	17	1	,	203	À	235	Î
12	<ff></ff>	44	,	76	L	108	-1		140	å	17	2		204	Ã	236	Ϊ
13	<cr></cr>	45	-	77	Μ	109	m		141	ç	17	3	≠	205	Õ	237	Ì
14	<s0></s0>	46		78	Ν	110	n		142	é	17	4	Æ	206	Œ	238	Ó
15	<si></si>	47	/	79	0	111	0		143	è	17	5	Ø	207	œ	239	Ô
16	<dle></dle>	48	0	80	Р	112	р		144	ê	17	6	∞	208	_	240	É
17	<dc1></dc1>	49	1	81	Q	113	q		145	ë	17	7	±	209	_	241	Ò
18	<dc2></dc2>	50	2	82	R	114	r		146	ĺ	17	8	≤	210	**	242	Ú
19	<dc3></dc3>	51	3	83	S	115	S		147	ì	17	9	≥	211	"	243	Û
20	<dc4></dc4>	52	4	84	Т	116	t		148	î	18	0	¥	212	`	244	Ù
21	<nak></nak>	53	5	85	U	117	u		149	Ϊ	18	1	μ	213	,	245	1
22	<syn< td=""><td>54</td><td>6</td><td>86</td><td>V</td><td>118</td><td>V</td><td></td><td>150</td><td>ñ</td><td>18</td><td>2</td><td>9</td><td>214</td><td>÷</td><td>246</td><td>^</td></syn<>	54	6	86	V	118	V		150	ñ	18	2	9	214	÷	246	^
23	<etb></etb>	55	7	87	W	119	W		151	ó	18	3	Σ	215	\Diamond	247	~
24	<can></can>	56	8	88	Χ	120	X		152	ò	18	4	Π	216	ÿ	248	_
25		57	9	89	Υ	121	У		153	ô	18	5	П	217	Ϋ	249	U
26		58	:	90	Z	122	Z		154	Ö	18	6	ſ	218	/	250	•
27	<esc></esc>	59	;	91	[123	{		155	õ	18	7	a	219	€	251	0
28	<fs></fs>	60	<	92	\	124			156	ú	18	8	0	220	<	252	,
29	<gs></gs>	61	=	93]	125	}		157	ù	18	9	Ω	221	>	253	"
30	<rs></rs>	62	>	94	^	126	~		158	û	19	0	æ	222	fi	254	
31	<us></us>	63	?	95	_	127	<di< td=""><td>EL></td><td>159</td><td>ü</td><td>19</td><td>1</td><td>Ø</td><td>223</td><td>fl</td><td>255</td><td>•</td></di<>	EL>	159	ü	19	1	Ø	223	fl	255	•

TO DO:

Add Notes
Make video