## Data Cleanup: Quality Filtering

## Learning objectives:

- 1. Become familiar with the Phred Quality scoring system used to assign base call accuracies.
- 2. Identify the need for quality filtering and the basics of the procedure underlying it.

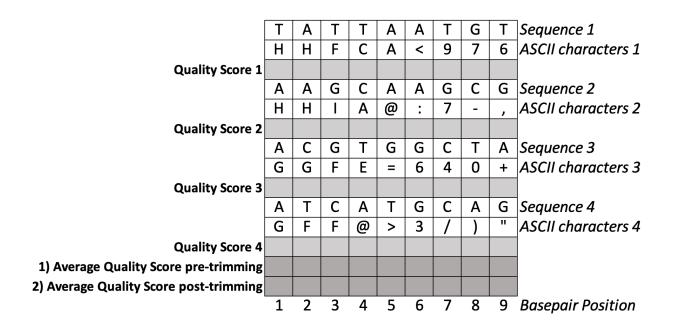
Level of Difficulty: 3/5

After removing contaminants, you'll need to get rid of the data that Illumina wasn't very sure about in a step called "quality trimming". What data is selected for trimming is determined from the data's Phred Quality score (aka Q Score), which is provided in addition to the sequence information in .fastq files (.fasta files only have the sequence information). Low Phred scores correspond to poor base pair reads by Illumina, which can be due to a variety of reasons (for instance Illumina can have trouble distinguishing several of the same base pair in a row). When Illumina is recording Phred scores into a data file, instead of writing numbers it uses ASCII (pronounced "ask-ee") characters, based on the conversion key shown below.

Symbol	Phred Q- Score	Probability of correct base assignment	Illumina Report Card Grade	Symbol	Phred Q- Score	Probability of correct base assignment	Illumina Report Card Grade
"	1	21%	F	6	21	99.2%	A
#	2	37%	F	7	22	99.4%	A
\$	3	50%	F	8	23	99.5%	A
%	4	60%	F	9	24	99.5%	A
&	5	68%	F	:	25	99.7%	Α
'	6	75%	D	;	26	99.75%	Α
(	7	80%	D	<	27	99.80%	Α
)	8	84%	D	=	28	99.84%	Α
*	9	87%	D	>	28	99.84%	Α
+	10	90%	D	?	30	99.90%	Α
,	11	92%	С	@	31	99.92%	Α
-	12	94%	С	Α	32	99.94%	Α
	13	95%	С	В	33	99.95%	Α
1	14	96%	С	С	34	99.96%	Α
0	15	96.8%	В	D	35	99.97%	Α
1	16	97.5%	В	E	36	99.975%	Α
2	17	98.0%	В	F	37	99.980%	Α
3	18	98.4%	В	G	38	99.984%	Α
4	19	98.7%	В	Н	39	99.987%	Α
5	20	99.0%	Α	I	40	99.990%	Α

## To do:

- 1) Using the conversion key provided, convert the ASCII characters for sequences 1-4 below into Phred quality score values in the grey boxes. Then plot the average Phred quality score for each base pair position on the provided plot.
- 2) Now "trim" any quality scores below 20 by crossing out that portion of the sequence to ensure only high quality data remains. Recalculate and replot the average Q-score values for each position.



## **Average Read Quality Along the Sequences**

