



# Genomics103

BIOL647

Digital Biology

Rodolfo Aramayo

# Genomics103

## FASTQ Format

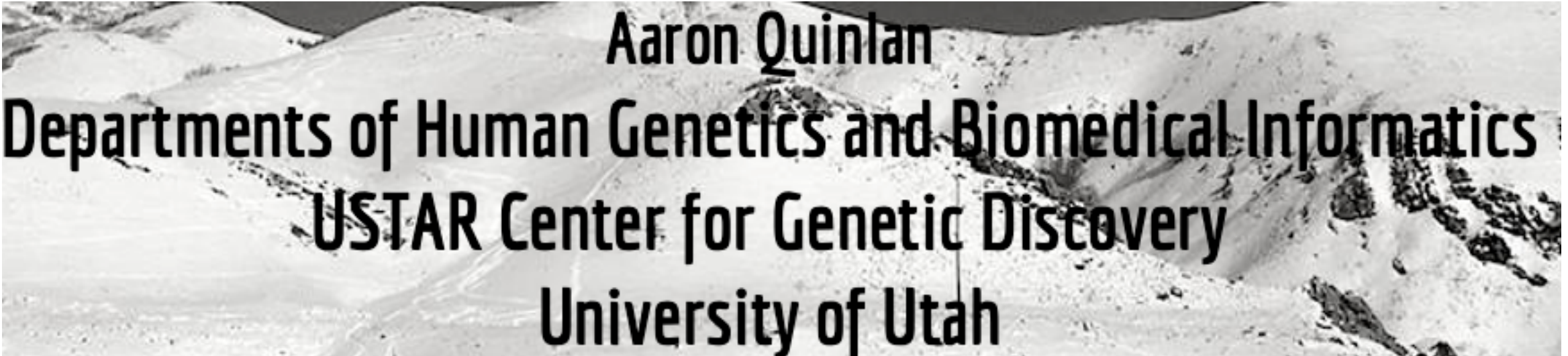
Ben Langmead



JOHNS HOPKINS

WHITING SCHOOL  
*of* ENGINEERING

Department of Computer Science

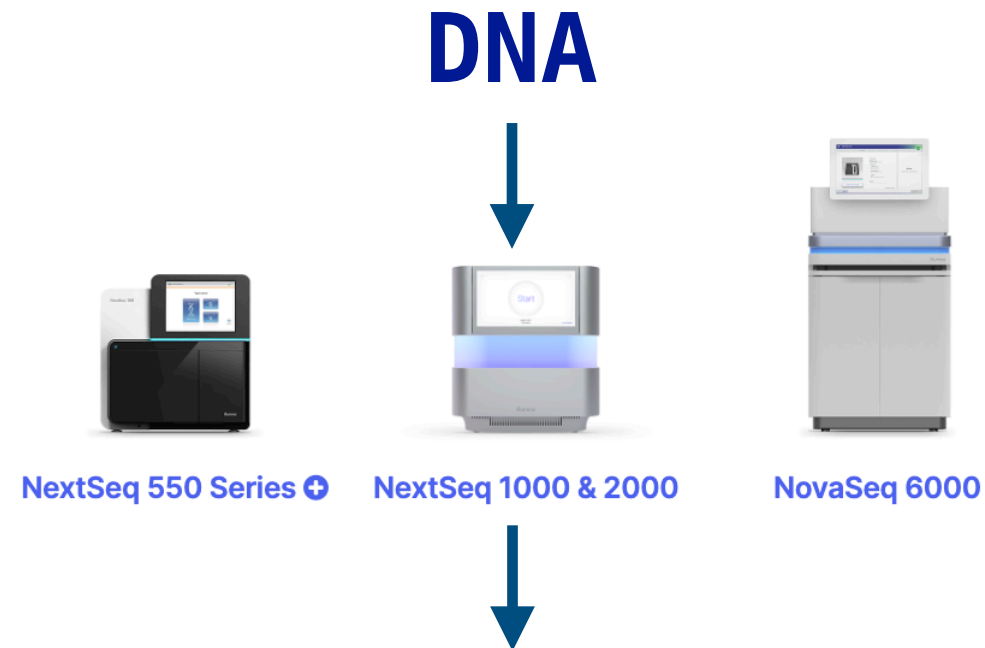
An aerial photograph of a vast desert landscape with rolling sand dunes and some sparse vegetation.

Aaron Quinlan  
Departments of Human Genetics and Biomedical Informatics  
USTAR Center for Genetic Discovery  
University of Utah

# Genomics103

## FASTQ Format

*The majority of DNA sequencing technologies produce a FASTQ file*



Name	@ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1
Sequence (ignore)	ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT +
Base qualities	?@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G

# Genomics103

## FASTQ Format

*A “standard” format for storing and defining sequences*

*from next-generation sequencing technologies*

*[http://en.wikipedia.org/wiki/FASTQ\\_format](http://en.wikipedia.org/wiki/FASTQ_format)*



Name	@ERR194146.1 HSQ1008:141:D0CC8ACXX:3:1308:20201:36071/1
Sequence	ACATCTGGTTCCTACTTCAGGGCCATAAAGCCTAAATAGCCCACACGTTCCCCTTAAAT
(ignore)	+
Base qualities	?@@FFBFFDDHHBCEAFGEGIIDHGH@GDHHHGEHID@C?GGDG@FHIGGH@FHBEG:G

# Genomics103

## FASTQ Format

*The FASTQ format's sequence identifier (first line of each record)*

### Old format

```
@HWUSI-EAS100R:6:73:941:1973#0/1
```

<b>HWUSI-EAS100R</b>	the unique instrument name
<b>6</b>	flowcell lane
<b>73</b>	tile number within the flowcell lane
<b>941</b>	'x'-coordinate of the cluster within the tile
<b>1973</b>	'y'-coordinate of the cluster within the tile
<b>#0</b>	index number for a multiplexed sample (0 for no indexing)
<b>/1</b>	the member of a pair, /1 or /2 ( <i>paired-end or mate-pair reads only</i> )

### New format

```
@EAS139:136:FC706VJ:2:2104:15343:197393 1:Y:18:ATCACG
```

<b>EAS139</b>	the unique instrument name
<b>136</b>	the run id
<b>FC706VJ</b>	the flowcell id
<b>2</b>	flowcell lane
<b>2104</b>	tile number within the flowcell lane
<b>15343</b>	'x'-coordinate of the cluster within the tile
<b>197393</b>	'y'-coordinate of the cluster within the tile
<b>1</b>	the member of a pair, 1 or 2 ( <i>paired-end or mate-pair reads only</i> )
<b>Y</b>	Y if the read is filtered, N otherwise
<b>18</b>	0 when none of the control bits are on, otherwise it is an even number
<b>ATCACG</b>	index sequence

# Genomics103

## FASTQ Format

*FASTQ quality scores: estimate of confidence in each base*

*(sequencing technologies make errors!)*

@SEQ\_ID

GATTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT

+

! ' '\*(( (\*\*\*) )%%%++) (%%%) .1\*\*\*-+\* ' ' ) \*\*55CCF>>>>>CCCCCCCC65



Qualities are based on the Phred scale and are encoded:



$$Q = -10 \cdot \log_{10}(P_{\text{err}})$$

**Note:**

The Ph in Phred comes from Phil Green, the inventor of the encoding

# Genomics103

## FASTQ Format

*Phred quality score calculation*

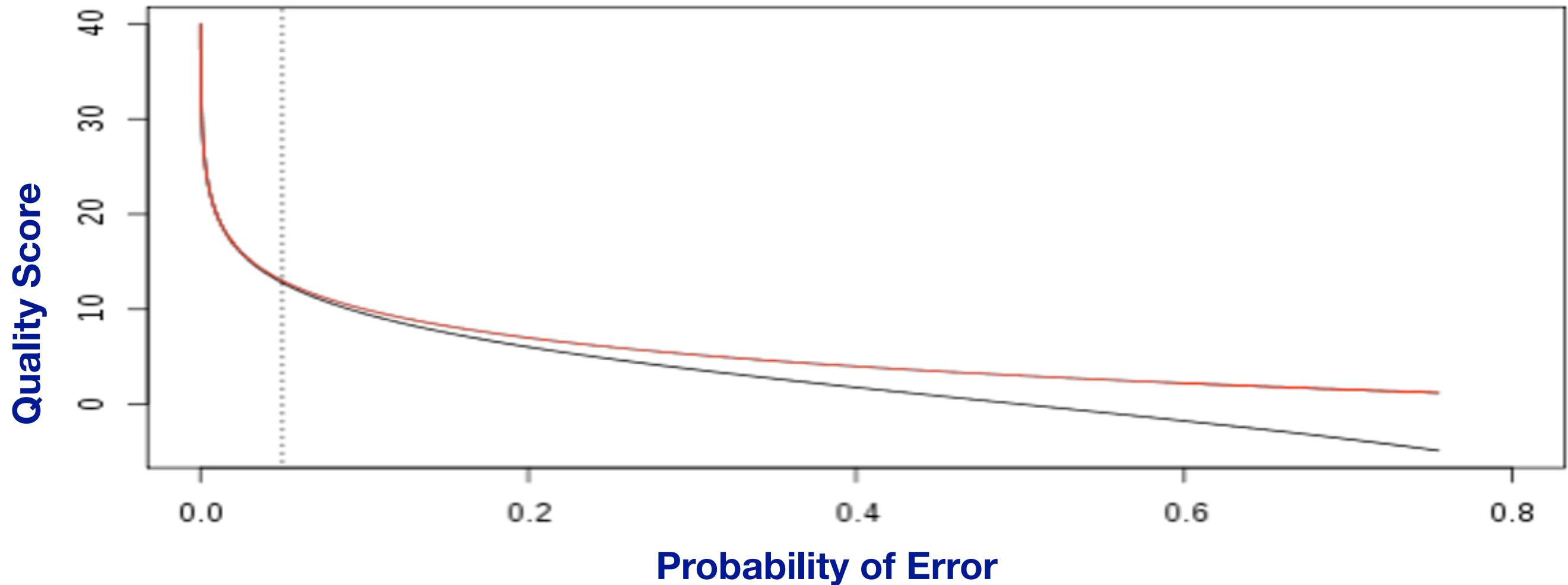
$$Q = -10 \cdot \log_{10}(P_{\text{err}})$$

Error probability	$\log_{10}(P_{\text{err}})$	Phred quality score
1	0	0
0.1	-1	10
0.01	-2	20
0.001	-3	30
0.0001	-4	40

# Genomics103

## FASTQ Format

*A higher quality score is better ( $\geq 20$  is considered "good")*





# Genomics103

# FASTQ Format

***Historically, FASTQ has had different encoding schemes for encoding PHRED quality scores***



S - Sanger Phred+33, raw reads typically (0, 40)  
X - Solexa Solexa+64, raw reads typically (-5, 40)  
I - Illumina 1.3+ Phred+64, raw reads typically (0, 40)  
J - Illumina 1.5+ Phred+64, raw reads typically (3, 40)  
with 0=unused, 1=unused, 2=Read Segment Quality Control Indicator (bold)  
(Note: See discussion above).  
L - Illumina 1.8+ Phred+33, raw reads typically (0, 41)

## Current encoding:

**! = quality 0**

**J = quality 41**

# Genomics103

## FASTQ Format

*Quality score encoding based on ASCII table*

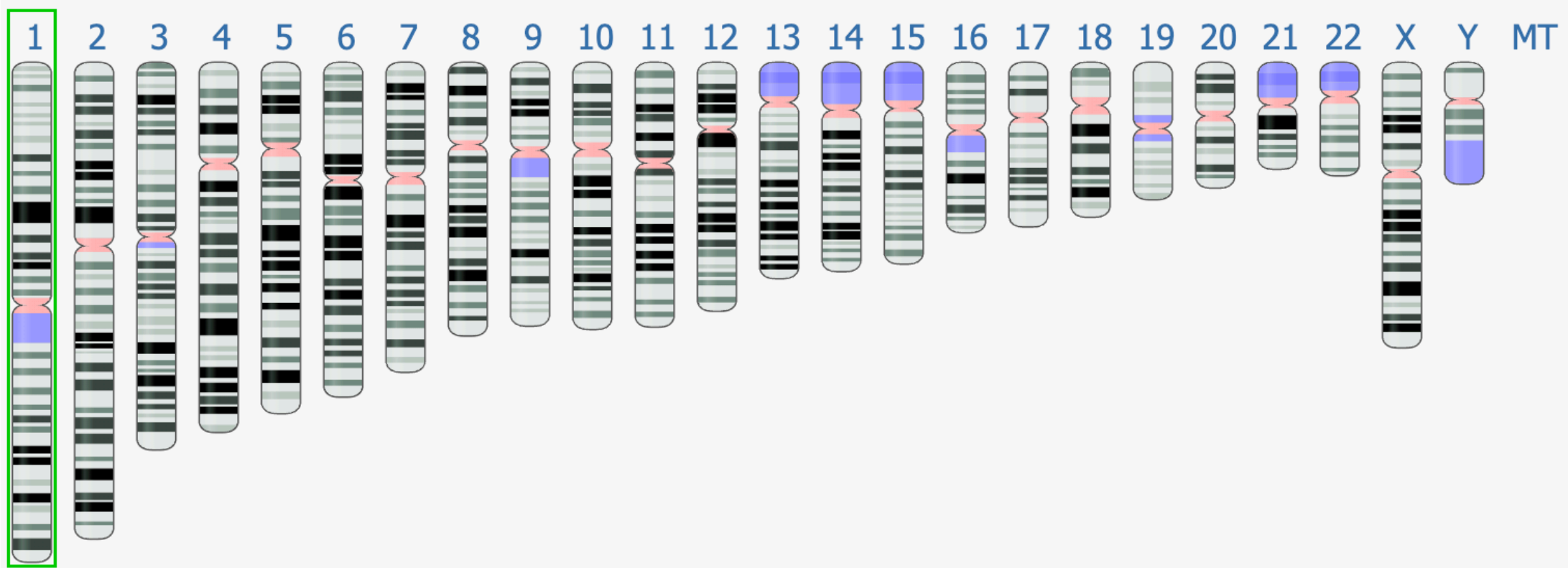
Formula for getting PHRED  
quality from encoded quality:

$$Q = (\text{ascii(char)} - 33)$$

Example:

	!	+	E	J
	↓	↓	↓	↓
ASCII	33	43	69	74
-33	0	10	36	41

Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char	Dec	Hex	Char
0	00	Null	32	20	Space	64	40	@	96	60	`
1	01	Start of heading	33	21	!	65	41	A	97	61	a
2	02	Start of text	34	22	"	66	42	B	98	62	b
3	03	End of text	35	23	#	67	43	C	99	63	c
4	04	End of transmit	36	24	\$	68	44	D	100	64	d
5	05	Enquiry	37	25	%	69	45	E	101	65	e
6	06	Acknowledge	38	26	&	70	46	F	102	66	f
7	07	Audible bell	39	27	'	71	47	G	103	67	g
8	08	Backspace	40	28	(	72	48	H	104	68	h
9	09	Horizontal tab	41	29	)	73	49	I	105	69	i
10	0A	Line feed	42	2A	*	74	4A	J	106	6A	j
11	0B	Vertical tab	43	2B	+	75	4B	K	107	6B	k
12	0C	Form feed	44	2C	,	76	4C	L	108	6C	l
13	0D	Carriage return	45	2D	-	77	4D	M	109	6D	m
14	0E	Shift out	46	2E	.	78	4E	N	110	6E	n
15	0F	Shift in	47	2F	/	79	4F	O	111	6F	o
16	10	Data link escape	48	30	0	80	50	P	112	70	p
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	T	116	74	t
21	15	Neg. 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	End 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	y
26	1A	Substitution	58	3A	:	90	5A	Z	122	7A	z
27	1B	Escape	59	3B	;	91	5B	[	123	7B	{
28	1C	File separator	60	3C	<	92	5C	\	124	7C	
29	1D	Group separator	61	3D	=	93	5D	]	125	7D	}
30	1E	Record separator	62	3E	>	94	5E	^	126	7E	~
31	1F	Unit separator	63	3F	?	95	5F	_	127	7F	□



# Genomics103

BIOL647

Digital Biology

Rodolfo Aramayo