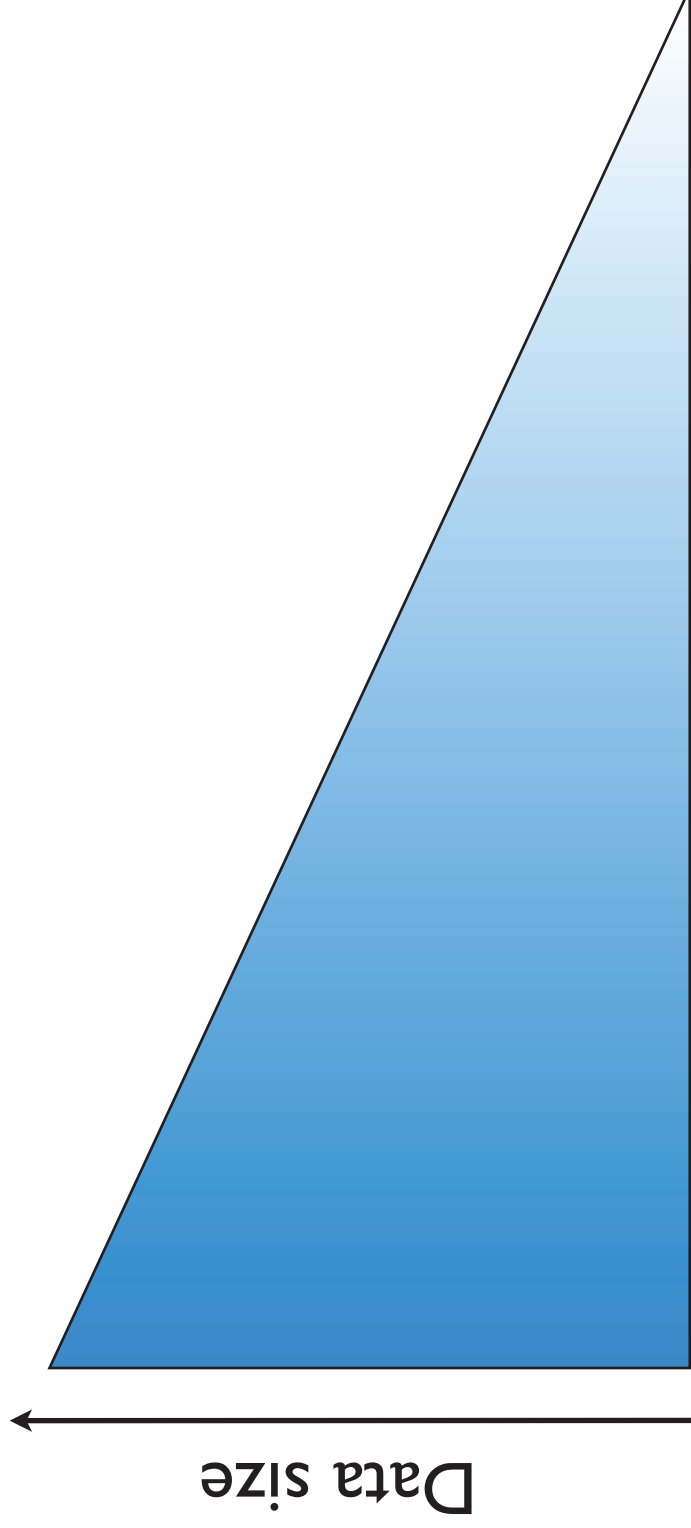


Generalized NGS analysis



Question

Raw
reads

Pre-
processing

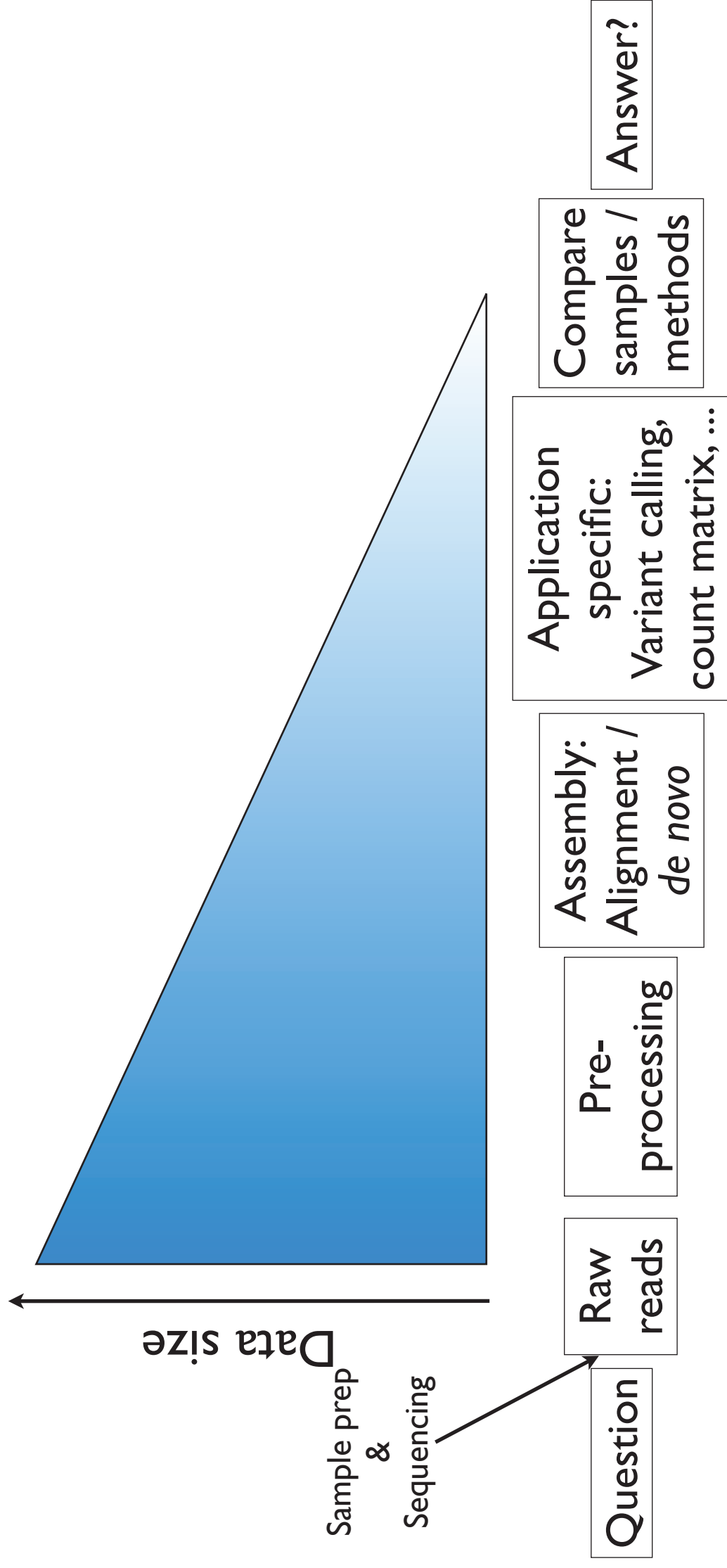
Assembly:
Alignment /
de novo

Application
specific:
Variant calling,
count matrix, ...

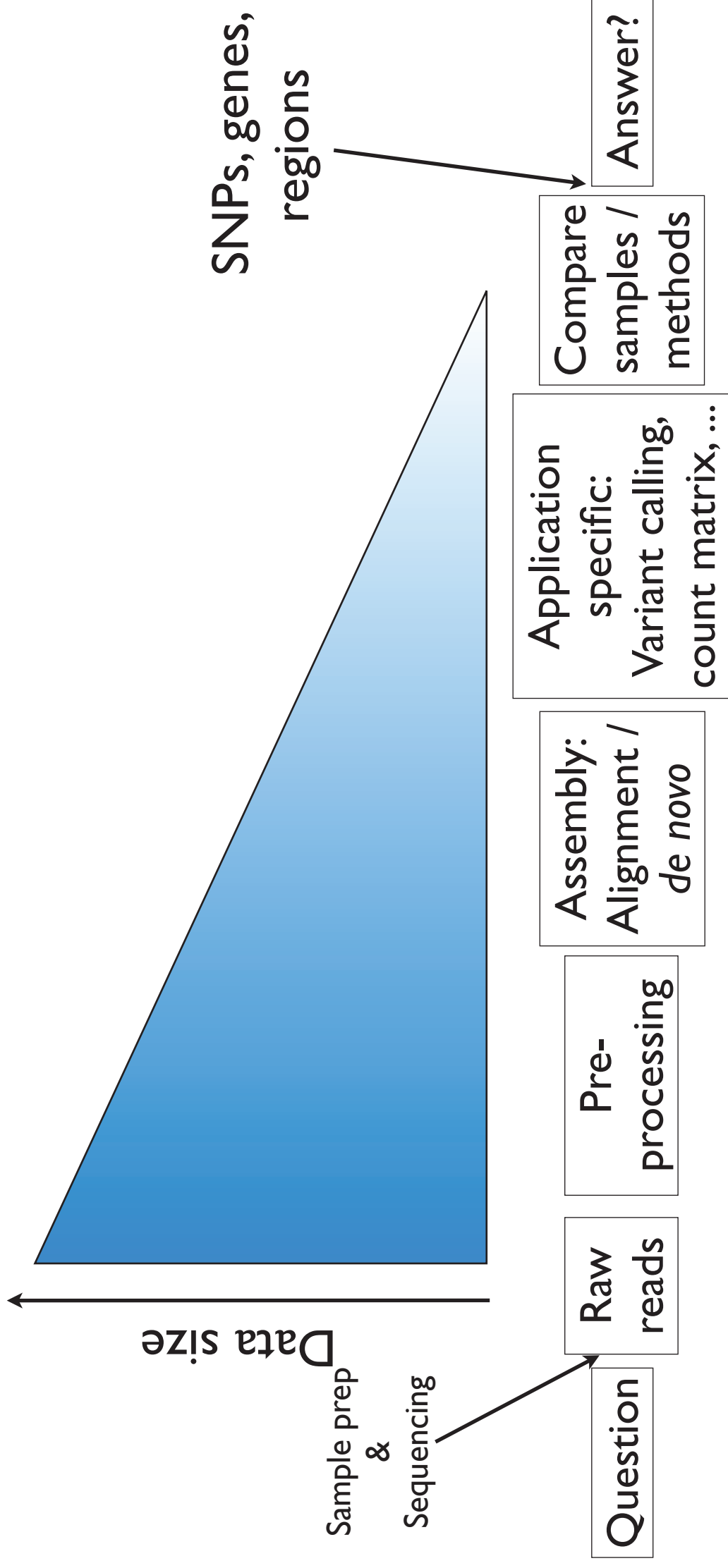
Compare
samples /
methods

Answer?

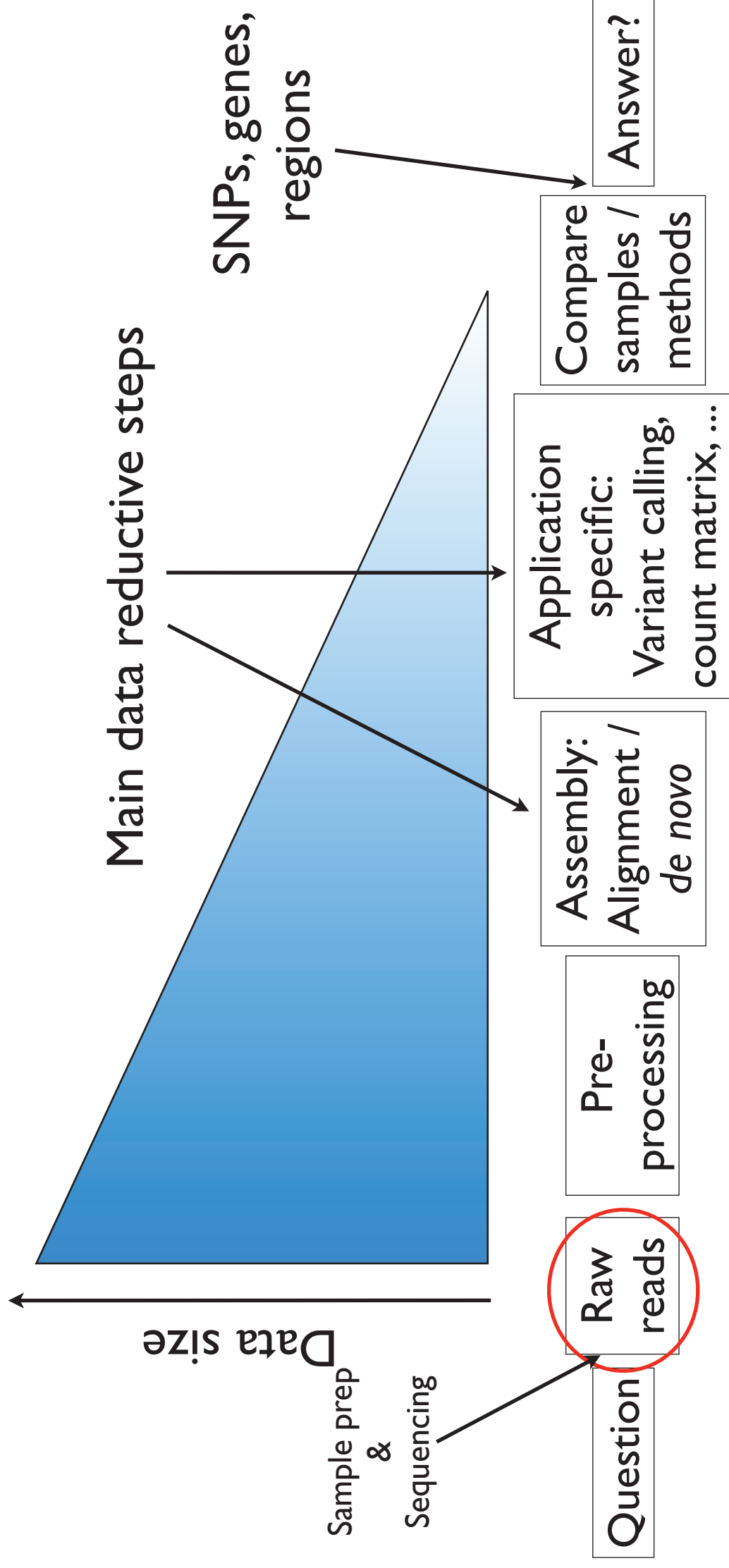
Generalized NGS analysis



Generalized NGS analysis



Generalized NGS analysis



What is sequence data?

Sequences are stored in fasta-files

Header



Sequence



```
>gi|218693476|ref|NC_011748.1| Escherichia coli 55989 chromosome, complete genome
GTAAGTATTTTTCAGCTTTTCATTCTGACTGCAACGGGCAATATGTCTCTGTGTGGATTAAAAAAAGAGT
GTCTGATAGCAGCTTCTGAACCTGTTACCTGCGTGAGTAAATTTAAATTTTATTTGACTTAGGTCACTAA
ATACTTTAAACCAATATAGGCATAGCGCACAGACAGATAAAATTACAGAGTACACAATCCATGAAACG
CATTAGCACCAACCATTAACCAACCATCAACCATTAACAGGTAAACGGTGCGGCTGACGCGTACAGGAA
ACACAGAAAAAGCCCGCACCTGACAGTGGGGCTTTTTCGACCAAAAGTTAACGAGGTAAACAACCAT
GCGAGTGTGAAGTTCGGCGGTACATCAGTGGCAAAATGCAGAAAGTTTCTGCGTGTGCGCGATATTCGTG
GAAAGCAATGCCAGGCAGGGCAGGTGGCCACCGTCTCTGCCCCGCCAAATCACCAACCAACCTGG
TGGCGATGATTGAAAAAACCATTAGCGGCCAGGATGCTTTACCAATATCAGCGATGCCGAAACGTATTTT
TGCCGAACTTTTGACGGGACTCGCCGCCGCCAGCCGGGTTCCTGCTGGCGCAATTGAAAAACTTTCGTC
GATCAGGAAATTGCCCCCAATAAAACATGTCTCATGGCATTAGTTTGTGGGGCAGTGCCCCGGATAGCA
```

E.coli ~ 4.5 - 6 Mbases

Human ~ 3.2 Gbases

Then what is NGS data?

Fastq

Header → @ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1

Sequence → ATCCCGGCCCTTTTCCAGGCCCTGCCCTCGAGC

+ BAAAGECEE<EEDFEDF3DBDBB=A+==>9>>88?

Qualities
(prob. that base call is wrong)

Then what is NGS data?

Fastq

Header → @ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1

Sequence → ATTCCCGGCTTTTCCAGGCCCTGCCCTCGAGC

+ BAAAGECEE<EEDFEDF3DBDBB=A+==>9>>88?

Qualities
(prob. that base call is wrong)

Millions to billions of these

A closer look at the qualities

Header → @ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1

Sequence → ATTCCCGGCTTTTCCAGGCTGCCCTGCTCGAGC

+ BAAAGECFEE<EEDFEDEF3DBDBB=A+==>9>>88?

Qualities
(prob. that base call is wrong)

One character encodes a number
using ascii table (0-255)

This number (Q) can be
converted to P

Phred-scale

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

$$P = 10^{(-Q/10)}$$

Phred scale

@ILLUMINA-C90280_00.
ATTCCCGGCCTTTTCCAG
+
BAAAGECEE<EEDFEF3D]

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	NUL (null)	32	20	040	 	Space	64	40	100	@	@
1	1	001	SOH (start of heading)	33	21	041	!	!	65	41	101	A	A
2	2	002	STX (start of text)	34	22	042	"	"	66	42	102	B	B
3	3	003	ETX (end of text)	35	23	043	#	#	67	43	103	C	C
4	4	004	EOT (end of transmission)	36	24	044	$	\$	68	44	104	D	D
5	5	005	ENQ (enquiry)	37	25	045	%	%	69	45	105	E	E
6	6	006	ACK (acknowledge)	38	26	046	&	&	70	46	106	F	F
7	7	007	BEL (bell)	39	27	047	'	'	71	47	107	G	G
8	8	010	BS (backspace)	40	28	050	((72	48	110	H	H
9	9	011	TAB (horizontal tab)	41	29	051))	73	49	111	I	I
10	A	012	LF (NL line feed, new line)	42	2A	052	*	*	74	4A	112	J	J
11	B	013	VT (vertical tab)	43	2B	053	+	+	75	4B	113	K	K
12	C	014	FF (NP form feed, new page)	44	2C	054	,	,	76	4C	114	L	L
13	D	015	CR (carriage return)	45	2D	055	-	-	77	4D	115	M	M
14	E	016	SO (shift out)	46	2E	056	.	.	78	4E	116	N	N
15	F	017	SI (shift in)	47	2F	057	/	/	79	4F	117	O	O
16	10	020	DLE (data link escape)	48	30	060	0	0	80	50	120	P	P
17	11	021	DC1 (device control 1)	49	31	061	1	1	81	51	121	Q	Q
18	12	022	DC2 (device control 2)	50	32	062	2	2	82	52	122	R	R
19	13	023	DC3 (device control 3)	51	33	063	3	3	83	53	123	S	S
20	14	024	DC4 (device control 4)	52	34	064	4	4	84	54	124	T	T
21	15	025	NAK (negative acknowledge)	53	35	065	5	5	85	55	125	U	U
22	16	026	SYN (synchronous idle)	54	36	066	6	6	86	56	126	V	V
23	17	027	ETB (end of trans. block)	55	37	067	7	7	87	57	127	W	W
24	18	030	CAN (cancel)	56	38	070	8	8	88	58	130	X	X
25	19	031	EM (end of medium)	57	39	071	9	9	89	59	131	Y	Y
26	1A	032	SUB (substitute)	58	3A	072	:	:	90	5A	132	Z	Z
27	1B	033	ESC (escape)	59	3B	073	;	;	91	5B	133	[[
28	1C	034	FS (file separator)	60	3C	074	<	<	92	5C	134	\	\
29	1D	035	GS (group separator)	61	3D	075	=	=	93	5D	135]]
30	1E	036	RS (record separator)	62	3E	076	>	>	94	5E	136	^	^
31	1F	037	US (unit separator)	63	3F	077	?	?	95	5F	137	_	_

Phred scale

@ILLUMINA-C90280_00

ATTCCCGGCCTTTTCCAG

+

BAAAGECEE<EEDFE3D



66

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	@ @	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			

Source: www.LookupTables.com

Phred scale

@ILLUMINA-C90280_00

ATTCCCGGCCTTTTCCAG

+

BAAAGECEE<EEDFEF3D

↑↑

66 65

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	@	A	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

Phred scale

@ILLUMINA-C90280_00

ATTCCCGGCCTTTTCCAG

+

BAAAGECEE<EEDFEF3D

↑↑↑

66 65 65

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	NUL (null)	32	20	040	 	Space	64	40	100	@	A
1	1	001	SOH (start of heading)	33	21	041	!	!	65	41	101	A	A
2	2	002	STX (start of text)	34	22	042	"	"	66	42	102	B	B
3	3	003	ETX (end of text)	35	23	043	#	#	67	43	103	C	C
4	4	004	EOT (end of transmission)	36	24	044	$	\$	68	44	104	D	D
5	5	005	ENQ (enquiry)	37	25	045	%	%	69	45	105	E	E
6	6	006	ACK (acknowledge)	38	26	046	&	&	70	46	106	F	F
7	7	007	BEL (bell)	39	27	047	'	'	71	47	107	G	G
8	8	010	BS (backspace)	40	28	050	((72	48	110	H	H
9	9	011	TAB (horizontal tab)	41	29	051))	73	49	111	I	I
10	A	012	LF (NL line feed, new line)	42	2A	052	*	*	74	4A	112	J	J
11	B	013	VT (vertical tab)	43	2B	053	+	+	75	4B	113	K	K
12	C	014	FF (NP form feed, new page)	44	2C	054	,	,	76	4C	114	L	L
13	D	015	CR (carriage return)	45	2D	055	-	-	77	4D	115	M	M
14	E	016	SO (shift out)	46	2E	056	.	.	78	4E	116	N	N
15	F	017	SI (shift in)	47	2F	057	/	/	79	4F	117	O	O
16	10	020	DLE (data link escape)	48	30	060	0	0	80	50	120	P	P
17	11	021	DC1 (device control 1)	49	31	061	1	1	81	51	121	Q	Q
18	12	022	DC2 (device control 2)	50	32	062	2	2	82	52	122	R	R
19	13	023	DC3 (device control 3)	51	33	063	3	3	83	53	123	S	S
20	14	024	DC4 (device control 4)	52	34	064	4	4	84	54	124	T	T
21	15	025	NAK (negative acknowledge)	53	35	065	5	5	85	55	125	U	U
22	16	026	SYN (synchronous idle)	54	36	066	6	6	86	56	126	V	V
23	17	027	ETB (end of trans. block)	55	37	067	7	7	87	57	127	W	W
24	18	030	CAN (cancel)	56	38	070	8	8	88	58	130	X	X
25	19	031	EM (end of medium)	57	39	071	9	9	89	59	131	Y	Y
26	1A	032	SUB (substitute)	58	3A	072	:	:	90	5A	132	Z	Z
27	1B	033	ESC (escape)	59	3B	073	;	;	91	5B	133	[[
28	1C	034	FS (file separator)	60	3C	074	<	<	92	5C	134	\	\
29	1D	035	GS (group separator)	61	3D	075	=	=	93	5D	135]	}
30	1E	036	RS (record separator)	62	3E	076	>	>	94	5E	136	^	^
31	1F	037	US (unit separator)	63	3F	077	?	?	95	5F	137	_	_

Phred scale

@ILLUMINA-C90280_00_

ATTCCCGGCCTTTTCCAG

+

BAAAGECEE<EEDFEF3D]

↑↑↑

66 65 65

$Q \sim \text{Prob}$

10 ~ 0.1

20 ~ 0.01

30 ~ 0.001

40 ~ 0.0001

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	NUL (null)	32	20	040	 	Space	64	40	100	@	A
1	1	001	SOH (start of heading)	33	21	041	!	!	65	41	101	A	A
2	2	002	STX (start of text)	34	22	042	"	"	66	42	102	B	B
3	3	003	ETX (end of text)	35	23	043	#	#	67	43	103	C	C
4	4	004	EOT (end of transmission)	36	24	044	$	\$	68	44	104	D	D
5	5	005	ENQ (enquiry)	37	25	045	%	%	69	45	105	E	E
6	6	006	ACK (acknowledge)	38	26	046	&	&	70	46	106	F	F
7	7	007	BEL (bell)	39	27	047	'	'	71	47	107	G	G
8	8	010	BS (backspace)	40	28	050	((72	48	110	H	H
9	9	011	TAB (horizontal tab)	41	29	051))	73	49	111	I	I
10	A	012	LF (NL line feed, new line)	42	2A	052	*	*	74	4A	112	J	J
11	B	013	VT (vertical tab)	43	2B	053	+	+	75	4B	113	K	K
12	C	014	FF (NP form feed, new page)	44	2C	054	,	,	76	4C	114	L	L
13	D	015	CR (carriage return)	45	2D	055	-	-	77	4D	115	M	M
14	E	016	SO (shift out)	46	2E	056	.	.	78	4E	116	N	N
15	F	017	SI (shift in)	47	2F	057	/	/	79	4F	117	O	O
16	10	020	DLE (data link escape)	48	30	060	0	0	80	50	120	P	P
17	11	021	DC1 (device control 1)	49	31	061	1	1	81	51	121	Q	Q
18	12	022	DC2 (device control 2)	50	32	062	2	2	82	52	122	R	R
19	13	023	DC3 (device control 3)	51	33	063	3	3	83	53	123	S	S
20	14	024	DC4 (device control 4)	52	34	064	4	4	84	54	124	T	T
21	15	025	NAK (negative acknowledge)	53	35	065	5	5	85	55	125	U	U
22	16	026	SYN (synchronous idle)	54	36	066	6	6	86	56	126	V	V
23	17	027	ETB (end of trans. block)	55	37	067	7	7	87	57	127	W	W
24	18	030	CAN (cancel)	56	38	070	8	8	88	58	130	X	X
25	19	031	EM (end of medium)	57	39	071	9	9	89	59	131	Y	Y
26	1A	032	SUB (substitute)	58	3A	072	:	:	90	5A	132	Z	Z
27	1B	033	ESC (escape)	59	3B	073	;	;	91	5B	133	[[
28	1C	034	FS (file separator)	60	3C	074	<	<	92	5C	134	\	\
29	1D	035	GS (group separator)	61	3D	075	=	=	93	5D	135]	}
30	1E	036	RS (record separator)	62	3E	076	>	>	94	5E	136	^	^
31	1F	037	US (unit separator)	63	3F	077	?	?	95	5F	137	_	_

Phred scale

@ILLUMINA-C90280_00

ATTCCCGGCCTTTTCCAG

+

BAAAGECEE<EEDFEF3D

↑↑↑

66 65 65 ~|e-6

$Q \sim \text{Prob}$

10 ~ 0.1

20 ~ 0.01

30 ~ 0.001

40 ~ 0.0001

Dec	Hx	Oct	Char	Dec	Hx	Oct	Html	Chr	Dec	Hx	Oct	Html	Chr
0	0	000	NUL (null)	32	20	040	 	Space	64	40	100	@	A
1	1	001	SOH (start of heading)	33	21	041	!	!	65	41	101	A	A
2	2	002	STX (start of text)	34	22	042	"	"	66	42	102	B	B
3	3	003	ETX (end of text)	35	23	043	#	#	67	43	103	C	C
4	4	004	EOT (end of transmission)	36	24	044	$	\$	68	44	104	D	D
5	5	005	ENQ (enquiry)	37	25	045	%	%	69	45	105	E	E
6	6	006	ACK (acknowledge)	38	26	046	&	&	70	46	106	F	F
7	7	007	BEL (bell)	39	27	047	'	'	71	47	107	G	G
8	8	010	BS (backspace)	40	28	050	((72	48	110	H	H
9	9	011	TAB (horizontal tab)	41	29	051))	73	49	111	I	I
10	A	012	LF (NL line feed, new line)	42	2A	052	*	*	74	4A	112	J	J
11	B	013	VT (vertical tab)	43	2B	053	+	+	75	4B	113	K	K
12	C	014	FF (NP form feed, new page)	44	2C	054	,	,	76	4C	114	L	L
13	D	015	CR (carriage return)	45	2D	055	-	-	77	4D	115	M	M
14	E	016	SO (shift out)	46	2E	056	.	.	78	4E	116	N	N
15	F	017	SI (shift in)	47	2F	057	/	/	79	4F	117	O	O
16	10	020	DLE (data link escape)	48	30	060	0	0	80	50	120	P	P
17	11	021	DC1 (device control 1)	49	31	061	1	1	81	51	121	Q	Q
18	12	022	DC2 (device control 2)	50	32	062	2	2	82	52	122	R	R
19	13	023	DC3 (device control 3)	51	33	063	3	3	83	53	123	S	S
20	14	024	DC4 (device control 4)	52	34	064	4	4	84	54	124	T	T
21	15	025	NAK (negative acknowledge)	53	35	065	5	5	85	55	125	U	U
22	16	026	SYN (synchronous idle)	54	36	066	6	6	86	56	126	V	V
23	17	027	ETB (end of trans. block)	55	37	067	7	7	87	57	127	W	W
24	18	030	CAN (cancel)	56	38	070	8	8	88	58	130	X	X
25	19	031	EM (end of medium)	57	39	071	9	9	89	59	131	Y	Y
26	1A	032	SUB (substitute)	58	3A	072	:	:	90	5A	132	Z	Z
27	1B	033	ESC (escape)	59	3B	073	;	;	91	5B	133	[[
28	1C	034	FS (file separator)	60	3C	074	<	<	92	5C	134	\	\
29	1D	035	GS (group separator)	61	3D	075	=	=	93	5D	135]	}
30	1E	036	RS (record separator)	62	3E	076	>	>	94	5E	136	^	^
31	1F	037	US (unit separator)	63	3F	077	?	?	95	5F	137	_	_

Phred-scaled probabilities

- Base qualities, read mapping qualities, variant qualities, ...
- Straight-forward, except for when they are used in reads!
- Offset: Sanger = 33, Illumina = 64

```
@ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1
ATTCCGGCCTTTTCCAGGCTGCCTGCTCGAGC
+
BAAAGECEE<EEDFEDF3DBDBB=A+=>9>>88?
```

Phred: 6665 65 ~1e-6

<u>Q ~ Prob</u>
10 ~ 0.1
20 ~ 0.01
30 ~ 0.001
40 ~ 0.0001

Phred-scaled probabilities

- Base qualities, read mapping qualities, variant qualities, ...
- Straight-forward, except for when they are used in reads!
- Offset: Sanger = 33, Illumina = 64

```
@ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1
ATTCCGGCCTTTTCCAGGCTGCCTGCTCGAGC
+
BAAAGECEE<EEDFEDF3DBDBB=A+=>9>>88?
```

Phred: 66 65 65 ~1e-6

Sanger: 33 32 32 ~0.001

<u>Q ~ Prob</u>
10 ~ 0.1
20 ~ 0.01
30 ~ 0.001
40 ~ 0.0001

Phred-scaled probabilities

- Base qualities, read mapping qualities, variant qualities, ...
- Straight-forward, except for when they are used in reads!
- Offset: Sanger = 33, Illumina = 64

```
@ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1
ATTCCGGCCTTTTCCAGGCTGCCTGCTCGAGC
+
BAAAGECEE<EEDFEDF3DBDBB=A+=>9>>88?
```

Phred: 66 65 65 ~|e-6

Sanger: 33 32 32 ~0.001

Illumina: 2 | | ~|

<u>Q ~ Prob</u>
10 ~ 0.1
20 ~ 0.01
30 ~ 0.001
40 ~ 0.0001

Phred-scaled probabilities

- Base qualities, read mapping qualities, variant qualities, ...
- Straight-forward, except for when they are used in reads!
- Offset: Sanger = 33, Illumina = 64

```
@ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1
ATTCCGGCCTTTTCCAGGCTGCCTGCTCGAGC
+
BAAAGECEE<EEDFEDF3DBDBB=A+=>9>>88?
```

Phred: 66 65 65 ~|e-6

Sanger: 33 32 32 ~0.001

Illumina: 2 | | ~|

HUGE difference!

<u>Q ~ Prob</u>
10 ~ 0.1
20 ~ 0.01
30 ~ 0.001
40 ~ 0.0001

Phred-scaled probabilities

- Base qualities, read mapping qualities, variant qualities, ...
- Straight-forward, except for when they are used in reads!
- Offset: Sanger = 33, Illumina = 64

```
@ILLUMINA-C90280_0030_FC:5:1:2675:1090#NNNNNN/1
ATTCCGGCCTTTTCCAGGCTGCCGTCGAGC
+
BAAAGECEE<EEDFEDF3DBDBB=A+=>9>>88?
```

Phred: 66 65 65 ~|e-6

Sanger: 33 32 32 ~0.001

Illumina: 2 | | ~|

HUGE difference!

Exercise today

<u>Q ~ Prob</u>
10 ~ 0.1
20 ~ 0.01
30 ~ 0.001
40 ~ 0.0001

Sanger vs. Illumina vs. Solexa

- 454, Ion Torrent, Pac Bio, Sanger: “Sanger” encoding
- Illumina reads: “Illumina” or “Sanger” encoding. New reads are all “Sanger”
- Solexa data: Solexa encoding (bought by Illumina)
- All data from SRA: “Sanger”

Read types



Single end

Paired end
Ins: 200-800 bp



Mate pair
Ins: 2kb - 40kb (~5kb)

Read types



Single end

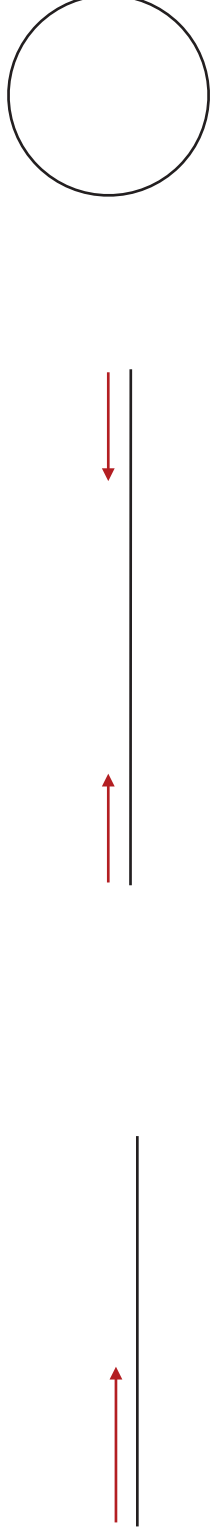


Paired end
Ins: 200-800 bp



Mate pair
Ins: 2kb - 40kb (~5kb)

Read types



Single end

Paired end
Ins: 200-800 bp

Mate pair
Ins: 2kb - 40kb (~5kb)

Read types



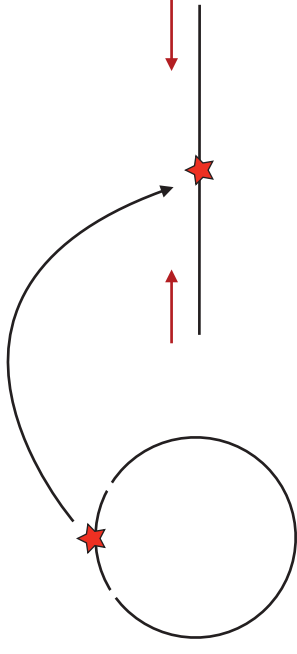
Fragment DNA:



Single end



Paired end
Ins: 200-800 bp

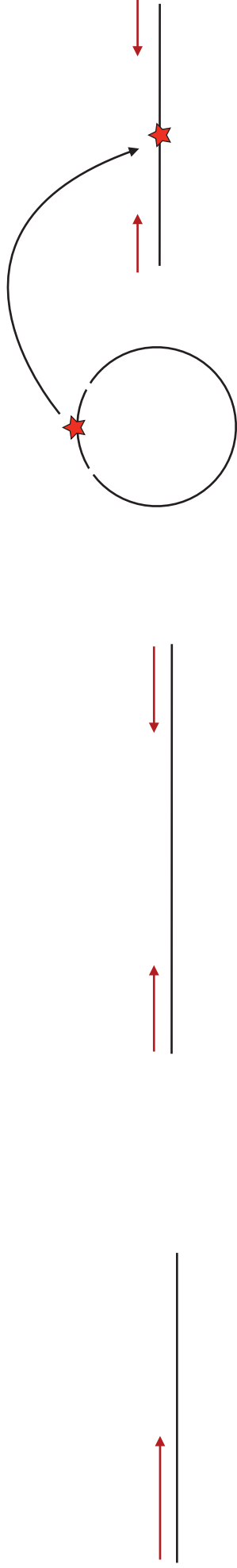


Mate pair
Ins: 2kb - 40kb (~5kb)

Read types



Fragment DNA:



Paired end
Ins: 200-800 bp

Mate pair
Ins: 2kb - 40kb (~5kb)
Protocol/technology dependent

Single end

Read orientation

Single end



Forward

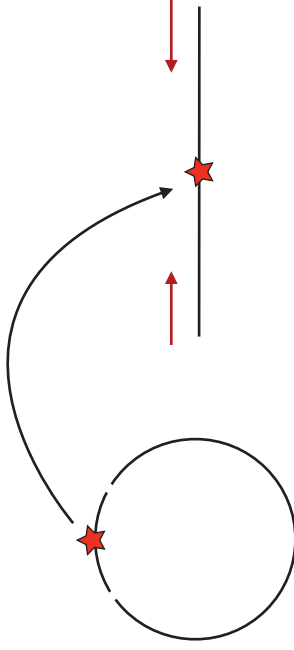
Paired end



Illumina: Forward - Reverse



Mate pair



Illumina: Reverse - Forward




Different for other technologies!

Special applications

- Single end reads:
 - Sometimes the only possibility (small DNA fragments / ancient DNA)
- Paired end reads:
 - More precise mapping/alignment/variation calls
 - Medium/Large indels (insertion/deletion)
 - Structural variations
 - Scaffolding in *de novo* assembly
- Mate pairs:
 - Scaffolding in *de novo* assembly
 - Structural variations

The reads comes in ?

- Illumina: fastq
 - 454: sff, fasta, fasta+qual, fastq
 - Ion torrent: sff, fastq
 - Pac Bio: fastq
 - Solid: csfasta, qual, xseq
- 
- Never keep data in fasta+qual/csfasta+qual - the qual format is horrible (and large)
 - sff is large, convert to fastq unless you use flowgram-capable software
 - Convert csfasta+qual to csfastq (at least if you are doing alignment)

Question

- What does it mean to have paired end reads?
- Discuss with neighbor for 2-3 mins, we discuss