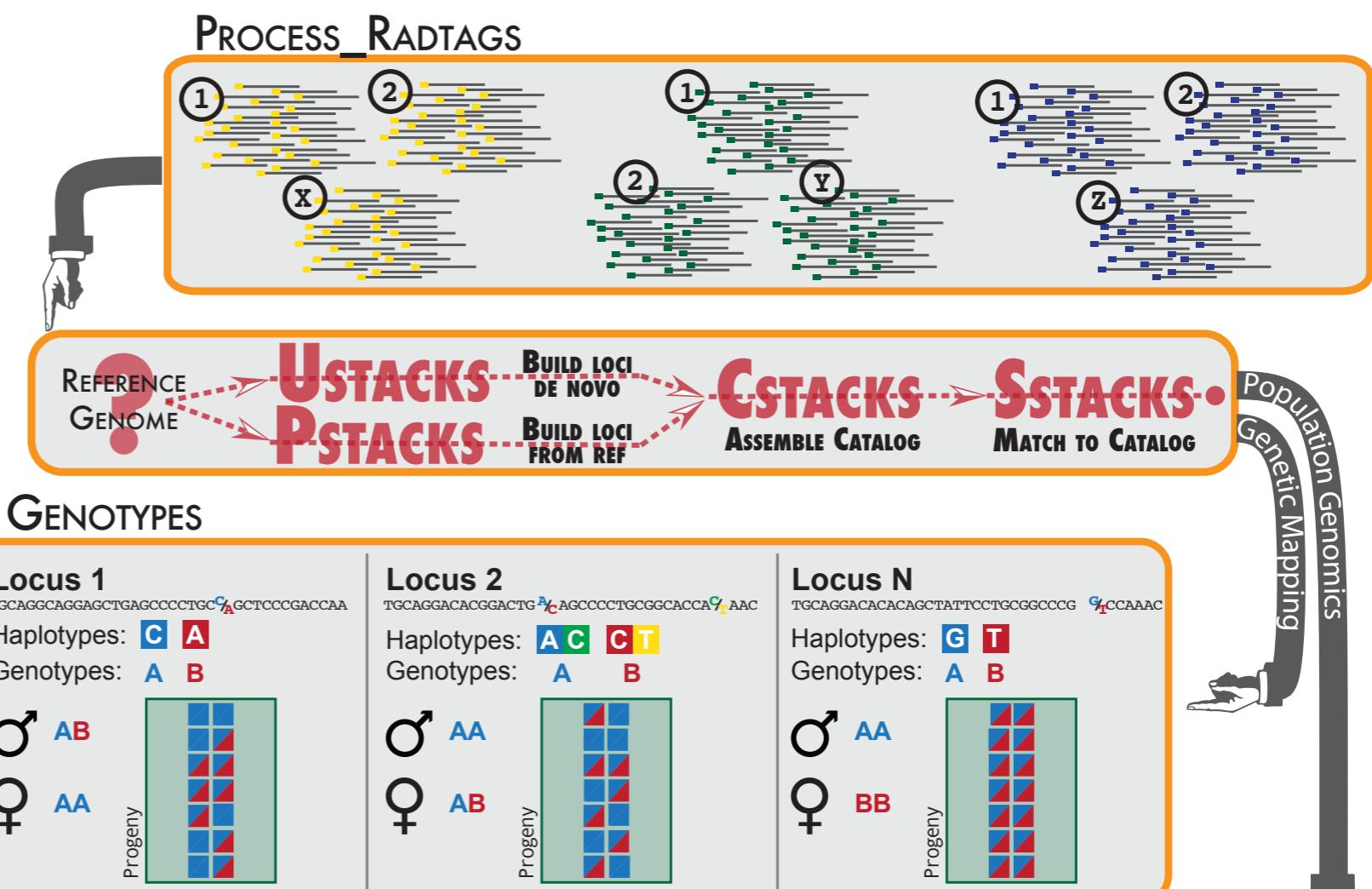
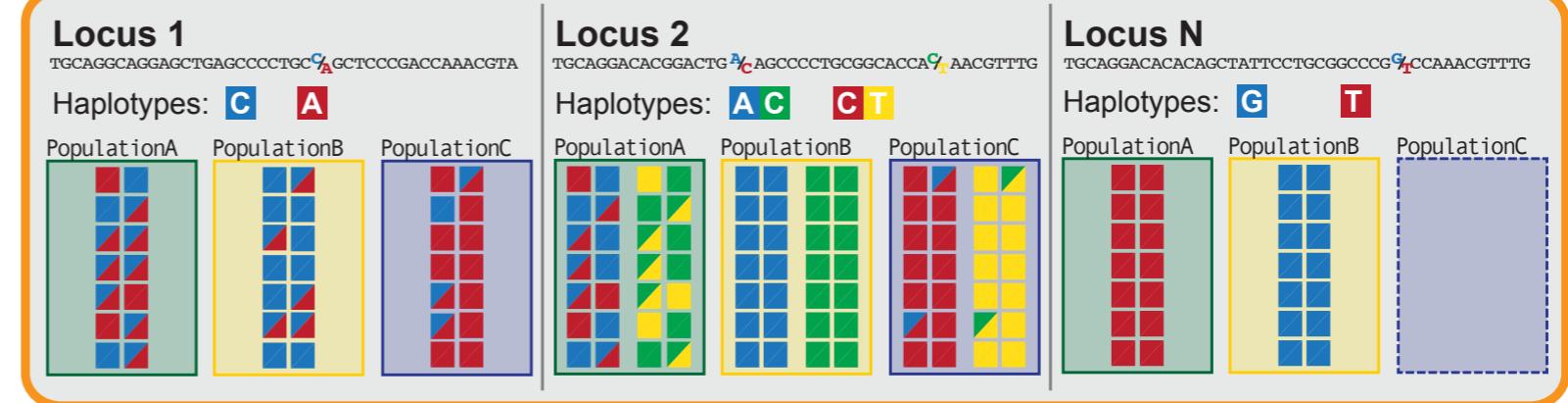


# Cleaning, Demultiplexing, and Deduplicating with Stacks

# Stacks



## POPULATIONS



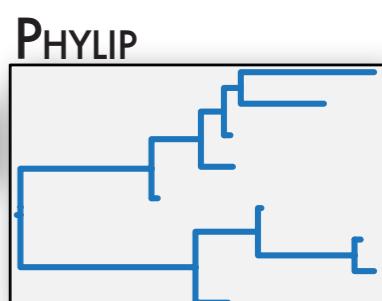
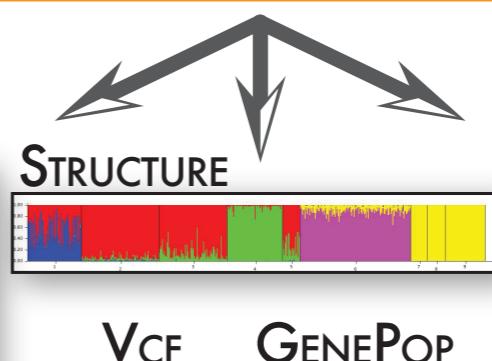
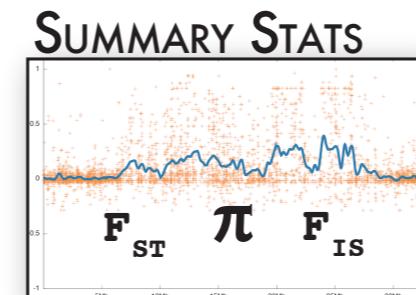
## Stacks: Building and Genotyping Loci De Novo From Short-Read Sequences

Julian M. Catchen,<sup>a</sup> Angel Amores,<sup>b</sup> Paul Hohenlohe,<sup>a</sup> William Cresko,<sup>a</sup> and John H. Postlethwait<sup>a</sup>  
<sup>a</sup>Center for Ecology and Evolutionary Biology and <sup>b</sup>Institute of Neuroscience, University of Oregon, Eugene, Oregon 97403

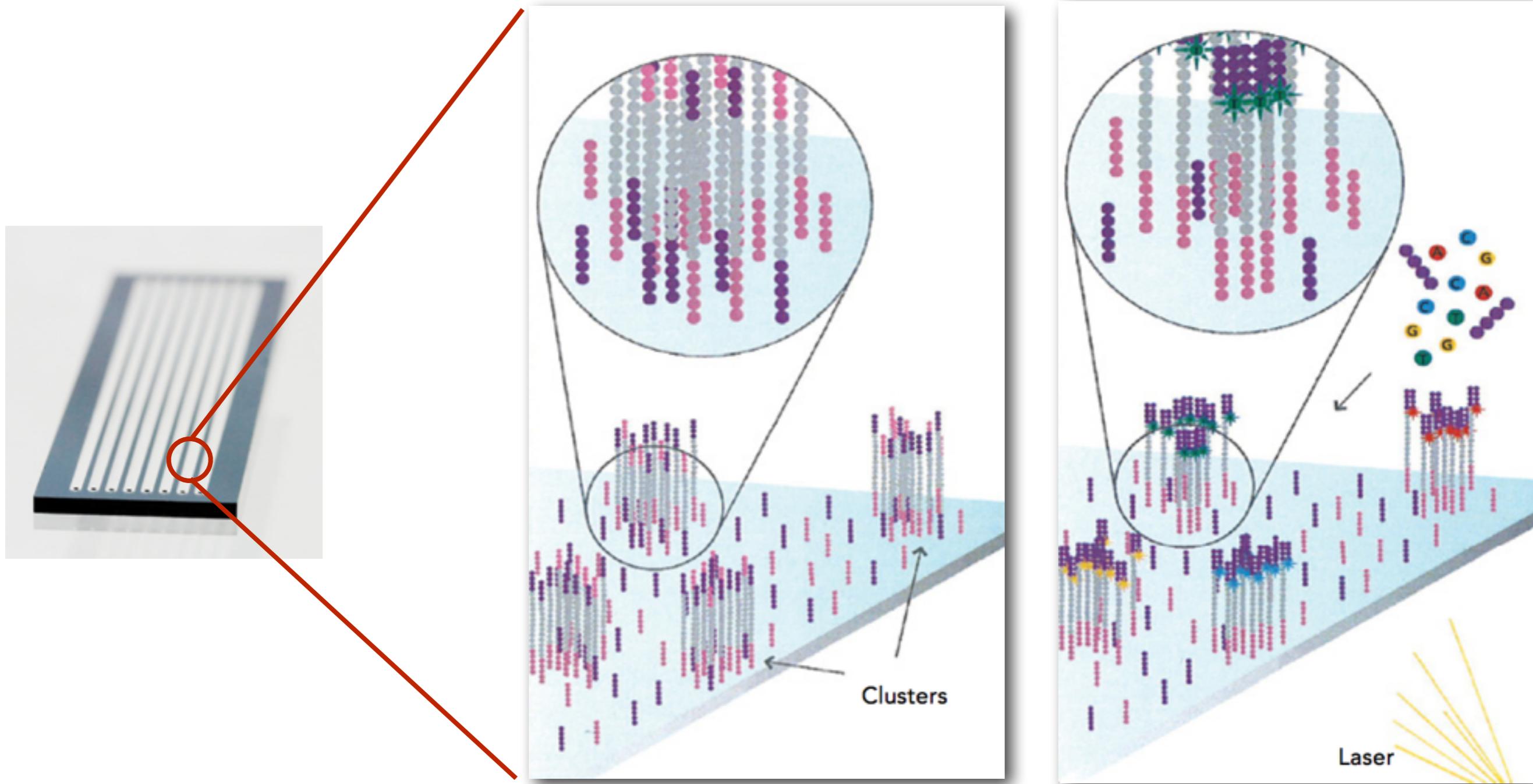
Stacks: an analysis tool set for population genomics

JULIAN CATCHEN,<sup>a</sup> PAUL A. HOHENLOHE,<sup>a,†</sup> SUSAN BASSHAM,<sup>a</sup> ANGEL AMORES,<sup>‡</sup>  
and WILLIAM A. CRESKO<sup>a</sup>

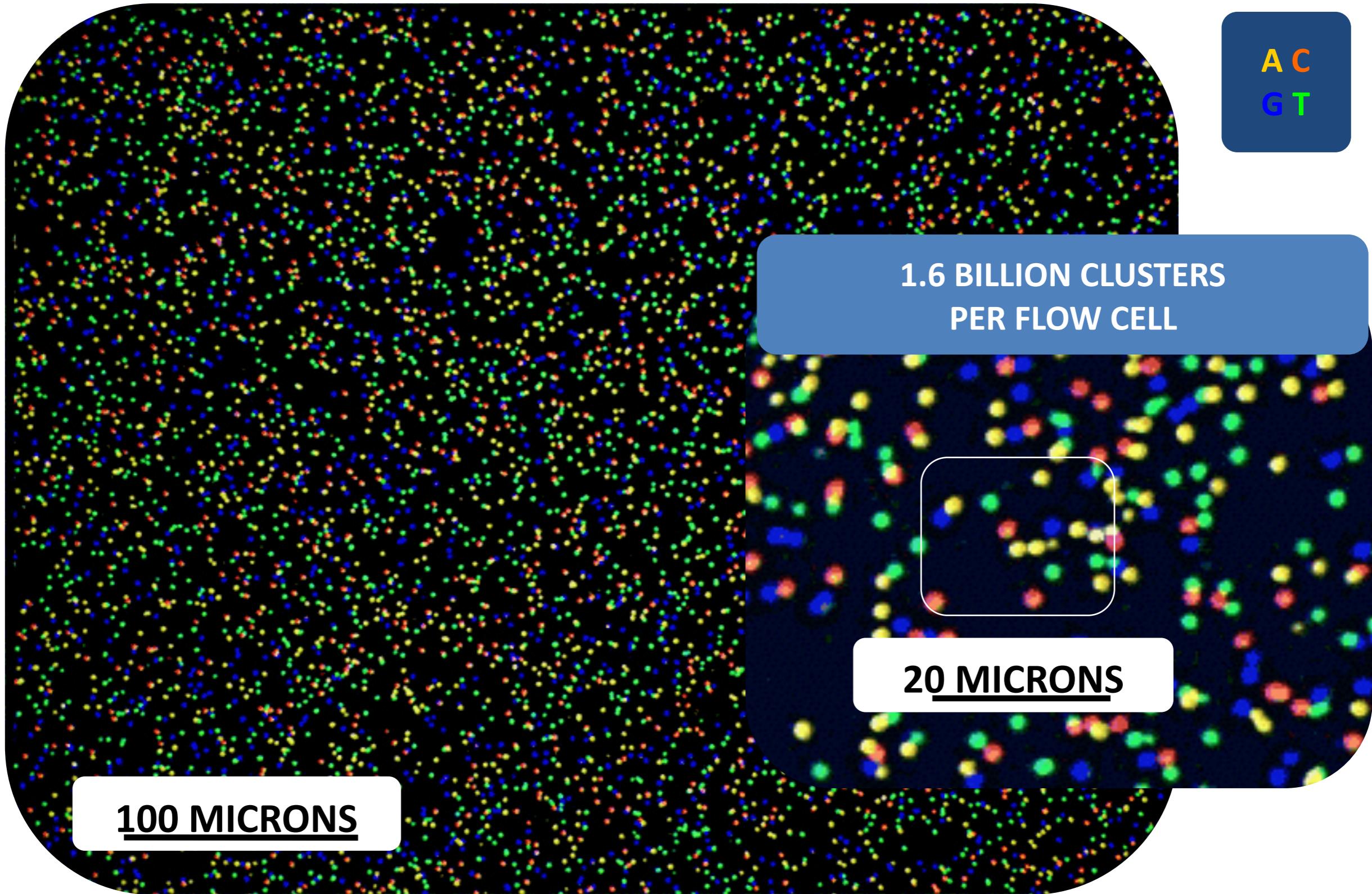
# JOINMAP R/QTL ONE MAP HAPLOTYPES



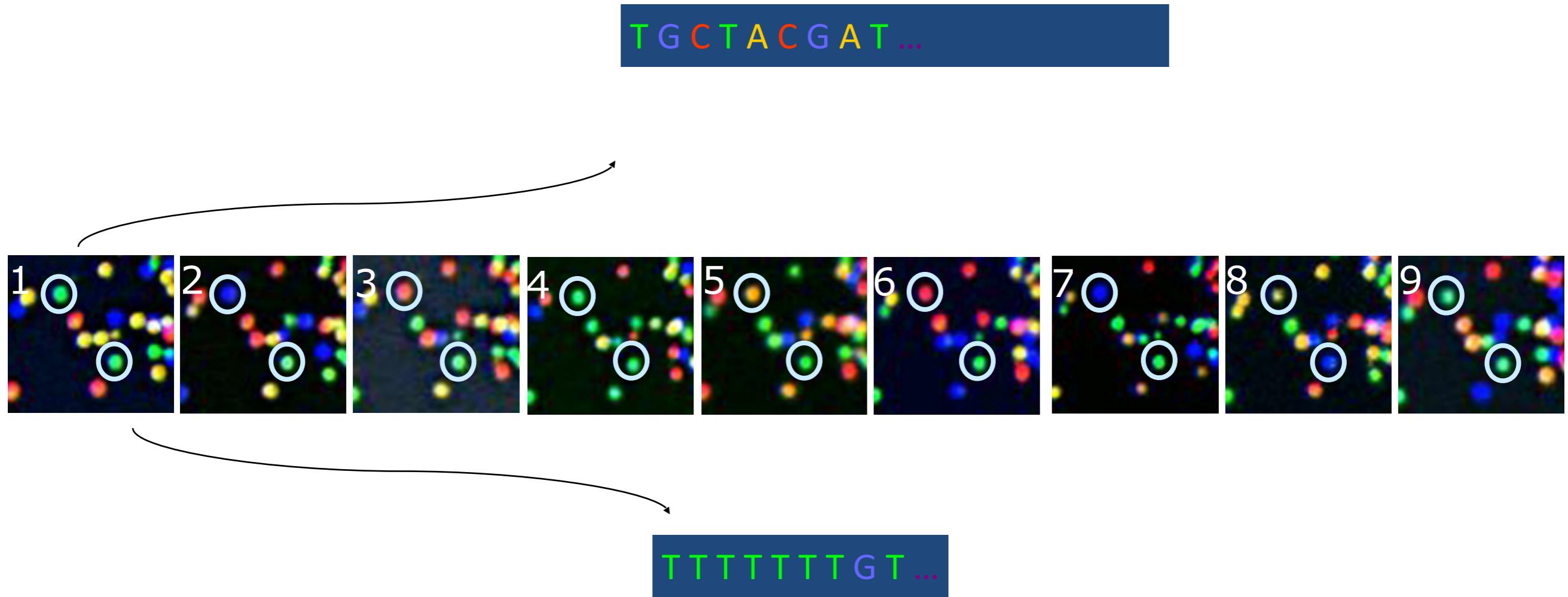
# Sequencing on Illumina's Flow cell



# Illumina Sequencing : How it looks



## Base calling from raw data

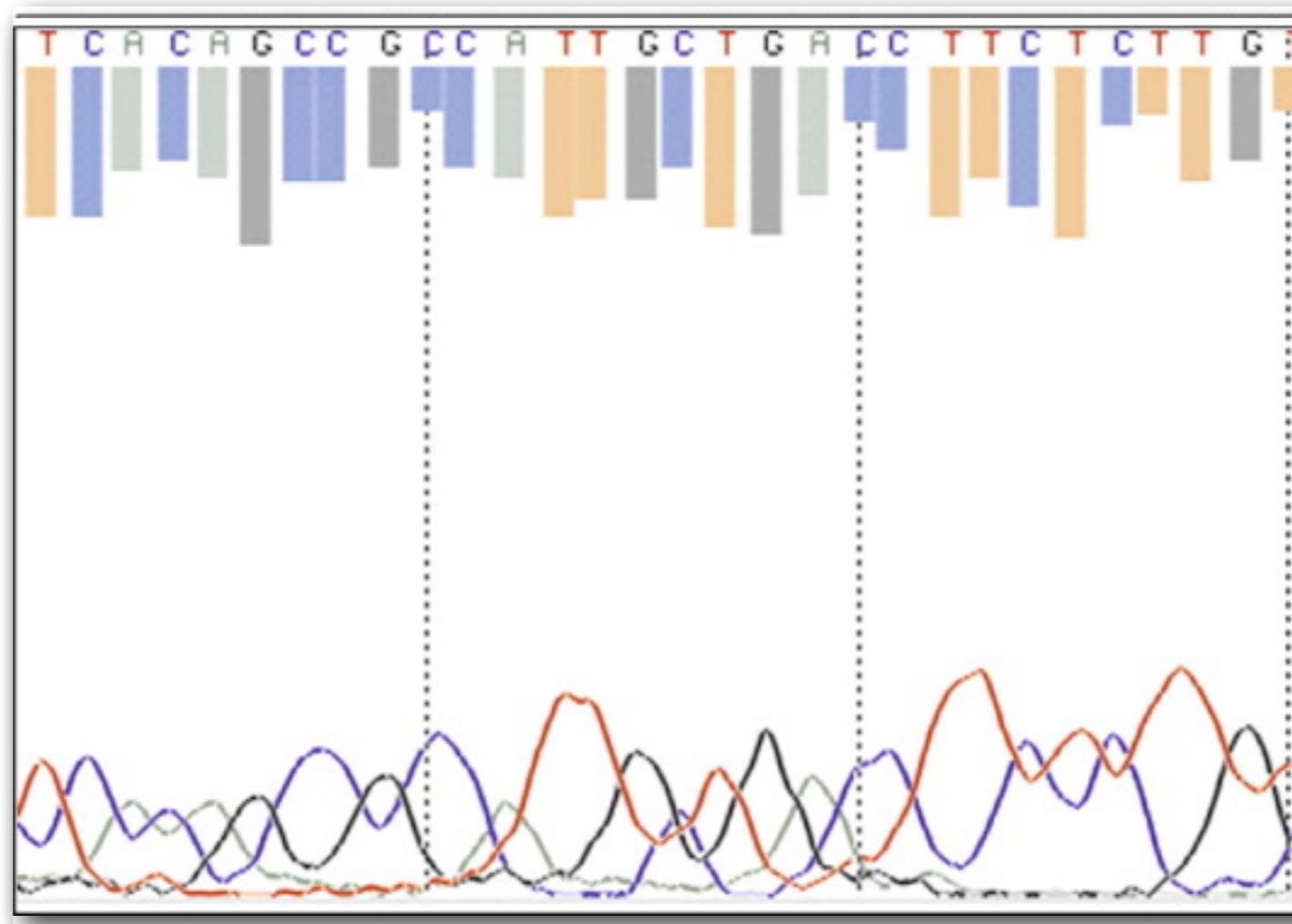


The identity of each base of a cluster is read off from sequential images.

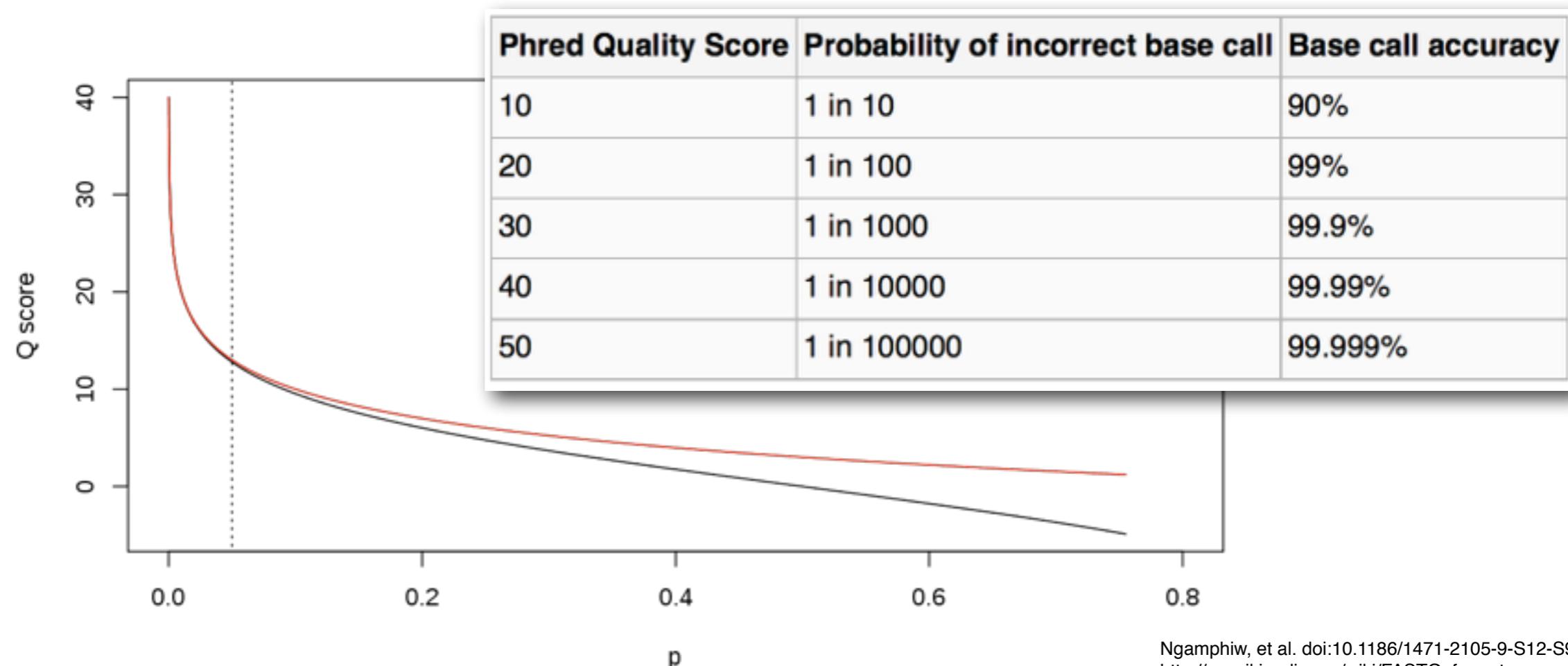
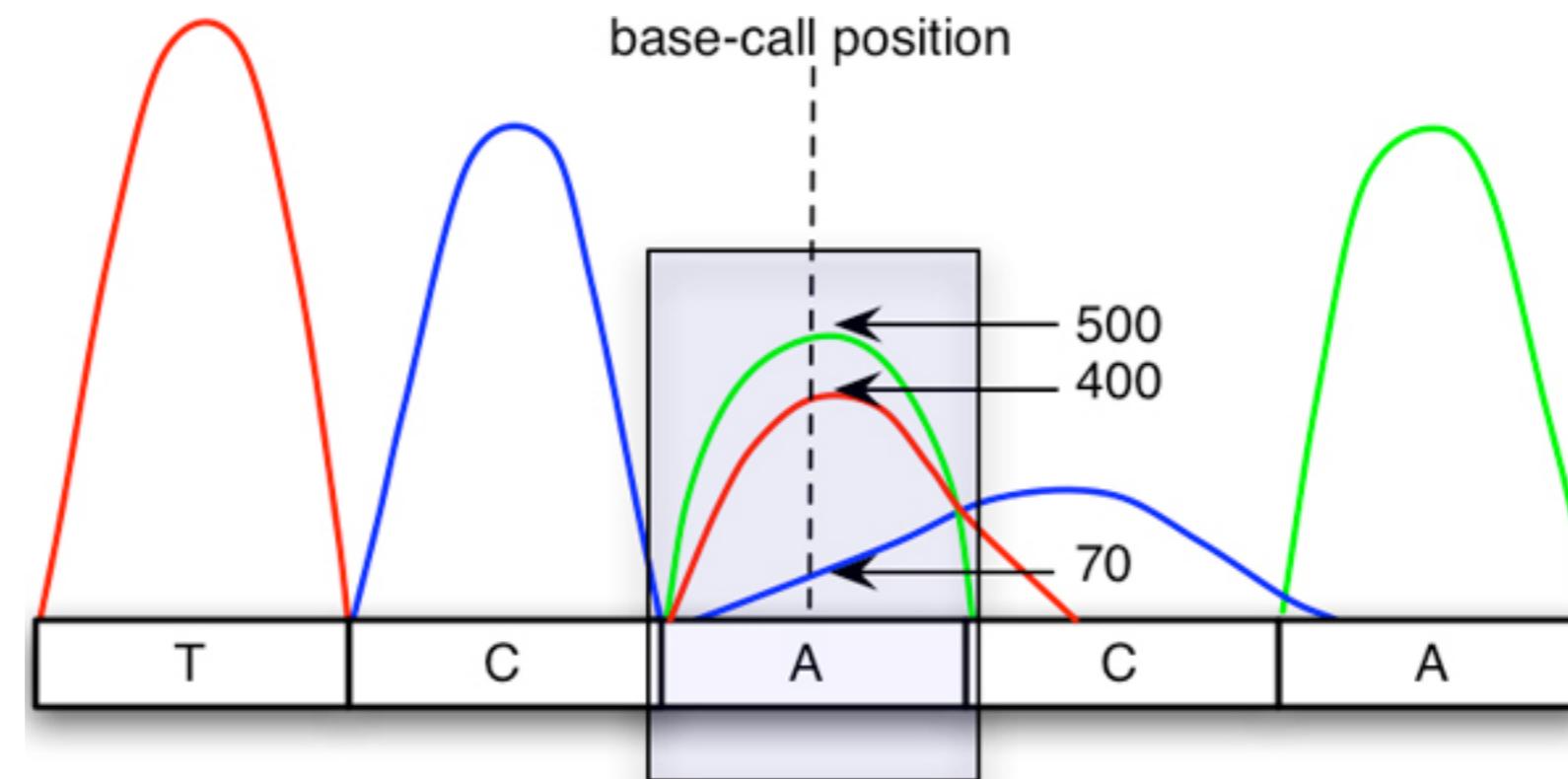
**Current read lengths = 36-150 nt**

**Total sequence data for 1 paired-end lane with 125bp = 500Gb!**

# Sequencing on Illumina's Flow cell, ctd.



# Phred Quality Score



# The FASTQ File Format

## FASTQ

```
@Sequence_137
TTTGTCTGCAGGGGGACACGTCAAAGTCAAACGCAGGCAAGTTGTGTTATGTCCAGTGGATCTTTGATT
+Sequence_137
<?@DDDDDFHHFBB@GGIACFHGGHBGHGCDHBEAHACHI=@CH.=7ACAHHADECDBCC66(6>@C>5@CACCA
```

```
@HWI-ST0747:162:C03AJACXX:3:1108:19763:106771 1:N:0:
TTTGTCTGCAGGGGGACACGTCAAAGTCAAACGCAGGCAAGTTGTGTTATGTCCAGTGGATCTTTGATT
+
<?@DDDDDFHHFBB@GGIACFHGGHBGHGCDHBEAHACHI=@CH.=7ACAHHADECDBCC66(6>@C>5@CACCA
```

# ASCII Code

0	<NUL>	32	<SPC>	64	@	96	'	128	Ä	160	†	192	�	224	‡
1	<SOH>	33	!	65	A	97	a	129	Å	161	°	193	i	225	.
2	<STX>	34	"	66	B	98	b	130	Ç	162	¢	194	¬	226	,
3	<ETX>	35	#	67	C	99	c	131	�	163	£	195	�	227	"
4	<EOT>	36	\$	68	D	100	d	132	�	164	§	196	f	228	%o
5	<ENQ>	37	%	69	E	101	e	133	�	165	•	197	�	229	�
6	<ACK>	38	&	70	F	102	f	134	�	166	¶	198	�	230	�
7	<BEL>	39	'	71	G	103	g	135	�	167	�	199	�	231	�
8	<BS>	40	(	72	H	104	h	136	�	168	�	200	»	232	�
9	<TAB>	41	)	73	I	105	i	137	�	169	�	201	...	233	�
10	<LF>	42	*	74	J	106	j	138	�	170	�	202	...	234	�
11	<VT>	43	+	75	K	107	k	139	�	171	�	203	�	235	�
12	<FF>	44	,	76	L	108	l	140	�	172	�	204	�	236	�
13	<CR>	45	-	77	M	109	m	141	�	173	�	205	�	237	�
14	<SO>	46	.	78	N	110	n	142	�	174	�	206	�	238	�
15	<SI>	47	/	79	O	111	o	143	�	175	�	207	�	239	�
16	<DLE>	48	0	80	P	112	p	144	�	176	�	208	-	240	�
17	<DC1>	49	1	81	Q	113	q	145	�	177	�	209	-	241	�
18	<DC2>	50	2	82	R	114	r	146	�	178	�	210	�	242	�
19	<DC3>	51	3	83	S	115	s	147	�	179	�	211	�	243	�
20	<DC4>	52	4	84	T	116	t	148	�	180	�	212	�	244	�
21	<NAK>	53	5	85	U	117	u	149	�	181	�	213	�	245	�
22	<SYN>	54	6	86	V	118	v	150	�	182	�	214	�	246	�
23	<ETB>	55	7	87	W	119	w	151	�	183	�	215	�	247	�
24	<CAN>	56	8	88	X	120	x	152	�	184	�	216	�	248	�
25	<EM>	57	9	89	Y	121	y	153	�	185	�	217	�	249	�
26	<SUB>	58	:	90	Z	122	z	154	�	186	�	218	/	250	�
27	<ESC>	59	;	91	[	123	{	155	�	187	�	219	�	251	�
28	<FS>	60	<	92	\	124		156	�	188	�	220	<	252	�
29	<GS>	61	=	93	]	125	}	157	�	189	�	221	>	253	�
30	<RS>	62	>	94	^	126	~	158	�	190	�	222	fi	254	�
31	<US>	63	?	95	_	127	<DEL>	159	�	191	�	223	fl	255	�

# The FASTQ File Format, ctd

```
@HWI-ST0747:162:C03AJACXX:3:1108:19763:106771 1:N:0:  
TTTGTCTGCAGGGGGACACGTCAAAGTCAAACGCAGGCAAGTTGTGTTATGTCCAGTGGATCTTGATT  
+  
<?@DDDDDHFHHFBB@GGIACFHGGHBGHGCDHBEAHACHI=@CH.=7ACAHHADECDBCC66 (6>@C>5@CACCA
```

# 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)

**ASCII values 33 - 73 = 0 - 40**

**‘F’ = 70**

$$70 - 33 = 37$$

# Single-end, Inline Barcodes

```
@HWI-ST0747:188:C09HWACXX:1:1101:2968:2083 1:N:0:  
TTATGATGCAGGACCAGGATGACGTCAGCACAGTGCGGGTCCATGGATGCTCCTCGGTGGTTGGGGAGGAGGCA  
+  
@@@DDDDDBHHFBF@CCAGEHHBFGIIFGIIGIEDBBGFHCGIIGAEEEDCC;A?;;5,:@A?=B5559999B@BBBBBA  
@HWI-ST0747:188:C09HWACXX:1:1101:2863:2096 1:N:0:  
TTATGATGCAGGCAAATAGAGTTGGATTTGTGTCAGTAGGCGGTAATCCCACAATTTACACTTATTCAAGGTGGA  
+  
CCCFFFFFHHHJJGHIGGAHHIIGGIJDHIGCEGHIFIJIH7DGIIIAHIJGEDHIDEHJJHFEECEFFDECDDD  
@HWI-ST0747:188:C09HWACXX:1:1101:2837:2098 1:N:0:  
GTGCCTTGCAGGCAATTAAGTTAGCCGAGATTAAGCGAAGGTTGAAAATGTCGGATGGAGTCCGGCAGCGAATGTAAA
```

# Paired-end, Inline Barcodes



@9432NS1:54:C1K8JACXX:7:1101:5584:1725 1:N:0:

**ACTGG**CATGATGATCATAGTATAACGTGGGATACATATGCCTAACGGCTAAAGATGCCTGAAGCTTGGCTATGTT

+

DDDBHHFBF@CCAGEHHHBFGIIFGIIGIEDBBGFHCGIIGAEEDCC;A?;;5,:@A?=B5559999B@BBBBBA

@9432NS1:54:C1K8JACXX:7:1101:5708:1737 1:N:0:

**TTCGA**CATGTGTTACAACCGCGAACGGACAAAGCATTGAAAATCCTTGGTTCGTTACTCTCTCCTAGCAT

+

CCCCFFFFHHHHJJGHIGGAHHIIGGIJDHIGCEGHIFIJIH7DGIIIAHIJGEDHIDEHJJHFEEECEFEFFD



@9432NS1:54:C1K8JACXX:7:1101:5584:1725 2:N:0:

**AATTT**ACTTTGATAGAAGAACAAACATAAGCCAAGCTCAAGGCATCTTAGCCTAGGCATATGTATCCCACGTTA

+

@@@DDDDDBHHFBF@CCAGEHHHBFGIIFGIIGIEDBBGFHCGIIGAEEDCC;A?;;5,:@A?=B5559999B@B

@9432NS1:54:C1K8JACXX:7:1101:5708:1737 2:N:0:

**AGTCT**TGTAAAAACGAAATCTTCAAAATGCTAGGAGAGAGTAACGAAACCAAGGATTTCATGCTTG

+

CCCCFFFFHHHHJJGHIGGAHHIIGGIJDHIGCEGHIFIJIH7DGIIIAHIJGEDHIDEHJJHFEEECEFEFFD

# Index Barcodes

## Genomic DNA

XXXXXXXXXXXXXX  
XXXXXXXXXXXXXX

## Enzyme digest

CTAGXXXXXXXXXXX  
XXXXXXXXXXXXTTAA

## Ligation

ACGACGCTTTCCGATC TCCGAATGCTAGX XXXXXXXXXX AATT ACGTTAGA GATCGGAAGAGCACACGT  
TGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGAGGCTTACGATC XXXXXXXXXX XTTAA TGCAATCT CTAGCCTCTCGTGCAGACTTGAGGTCAGTG  
P X

## Limited cycle PCR

### iTru5 primer

AATGATA CGGC GACC ACCGAG ATCTACAC ACCGACAA ACAC TTTCCCTAC



AATGATA CGGC GACC ACCGAG ATCTACAC ACCGACAA ACAC TTTCCCTAC AGCAGC TCTCCGATC TCCGAATGCTAGX XXXXXXXXXX AATT ACGTTAGA GATCGGAAGAGCACACGTCTGAAC TCCAGTCAC AGGTCACT ATCTCGTATGCCGTCTCTGCTTG  
TTACTATGCCGCTGGCTAGATGTGTTGTGAGAAAGGGATGTGCTGCGAGAAGGCTAGAGGCTTACGATC XXXXXXXXXX XTTAA TGCAATCT CTAGCCTCTCGTGCAGACTTGAGGTCAGTG TCCAGTGATAGAGCATA CGGCAGAACGAAAC

### iTru7 primer

GACTTGAGGT CAGTGTCCAGTGATAGAGCATA CGGCAGAACGAAAC



# Index Barcodes

# Paired-end, Index Barcodes

@NS500216:197:HF7C5BGXX:1:11101:1693:1047 1:N:0: **CTCCATGT+GGGGGGGG**  
TTAGGNAATCGGGCCTTCAGGAATGCTCACGGCAAGGGGCCACTAGGAACCTCCAGTCTTACCATGTGGTGAC  
+  
AAAAA#E/EEEEEEEAEEEEEEEEEE/EEEEEEEEEEAEEEEEEEEEEAEEEEEEEEEE  
@NS500216:197:HF7C5BGXX:1:11101:4119:1051 1:N:0: **TCGCTGTT+CTTCGTTTC**  
GATACNGCTAGAAGTGTATAAAATGCTAACTTGTAATGAATGTAAGAAAGGAGGATGAAGATGTGGTAGTGGT  
+  
AAAAA#EEEEEEEEEAEEEEEEEEEEAEEEEEEEAEEEEEEEAEEEEEEEAEEEEEEE

# PCR Duplicates

## Round 1

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

## Round 2

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

## Round 3

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

## Round 4

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTATATTGACCATGCCGCTTTCTTGGC

# PCR Duplicates

## Round 1

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**

## Round 2

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**



## Round 3

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**

## Round 4

AAGGATGGTGT →  
AAGGATGGTGTAAATGCCACTCCTCTCCGACTGTTAACACTACTGGTTA**GATTGACCATGCCGCTTTCTTGGC**

# process\_radtags

```
@Sequence_137
TTTGTCTGCAGGGGGACACGTCAAAGTCAAACGCAGGCAAGTTGTGTTATGTCCAGTGGATCTTGATT
+Sequence_137
<?@DDDDDFHHFBB@GGIACFHGGHBGHGCDHBEAHACHI=@CH.=7ACAHHADECDBCC66(6>@C>5@CACCA
```



```
<?@DDDDDFHHFBB@GGIACFHGGHBGHGCDHBEAHACHI=@CH.=7ACAHHADECDBCC66(6>@C>5@CACCA
```