

Session 1 – Lecture 2

Working with Short-Read Data Files

Qi Sun

**Computational Biology Service Unit
Cornell University**

OUTLINE

- 1. What are in the data files?**
- 2. Sequencing errors in both platforms.**
- 3. Tools for inspect qualities and pre-processing files.**
- 4. Computational resources available at Cornell.**

Files distributed after Roche 454 run

1. SFF file (Standard flowgram format)

Binary file. Can only be opened by SFF compatible software, eg. gsMapper by Roche. It can be used for depository into NCBI SRA database.

2. FASTA file

Text file, with sequence reads.

```
>FVJHPP401DVJJEI rank=0000010 x=1472.0 y=1128.0 length=245
GTTTCGTGAGGAAATCTAGGCATGTCAGAAAACAATGATTACTTTTCCACCAAAGACTACT
ATTGTATACATTAAGCAAGAAAATACCTGCATATCCTGCCTGTTTAAATGAGAGCTGCAA
CAAAAAGTTAGACAGATATGGGAGATGGCAGTTGGCGCTGCGAAAAGTGTGATGCGAATC
ATCCAAGACCCGAATATCGCTATATTCTATCTCTTAACGTAAATGATCATACTGGTCAAC
TTTGG
>FVJHPP401A56GX rank=0000013 x=363.0 y=1199.0 length=46
GCTGCGACCTGCATCGGCGTCATTCTGCTGGTCATGCTGGTGGAGG
>FVJHPP401C6VFC rank=0000021 x=1191.0 y=1286.0 length=87
TATCCACTAAACTGCGAAGATAGTAGCCTGACGAGTACTTATTAATCCCTTAAGTAGAG
GCCTATTGCGCGTGCCTAACATAAGAG
```

3. QUAL file

Text file, with PHRED based quality.

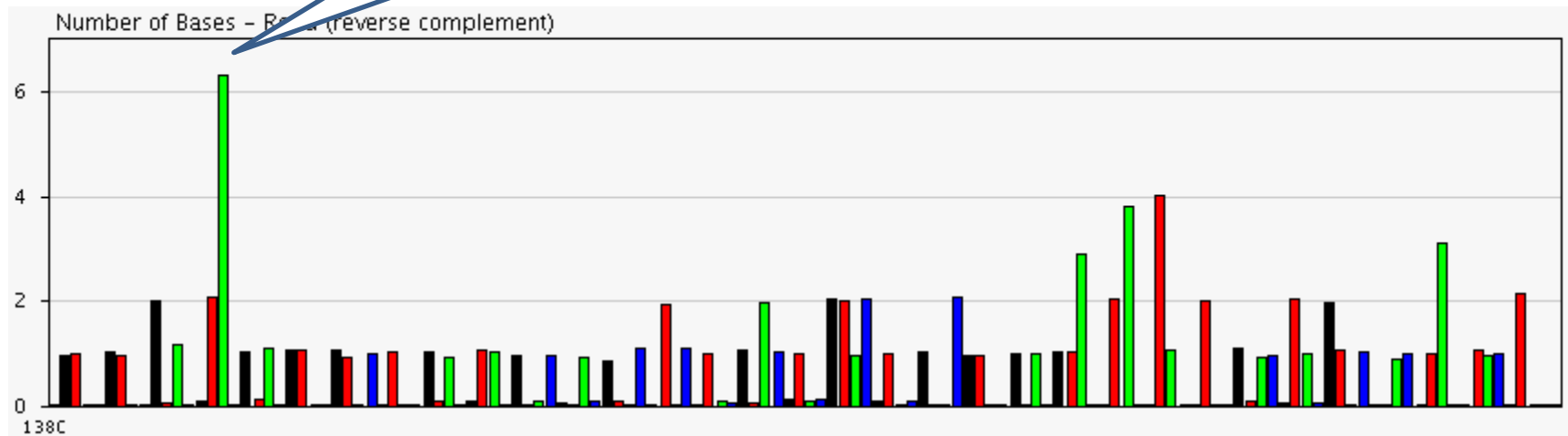
```
>FVJHPP401DVJJEI rank=0000010 x=1472.0 y=1128.0 length=245
32 32 32 32 32 32 32 32 24 24 19 19 19 25 28 32 32 32 32 33 33 33 32 30 32 30 25
 18 18 18 18 18 14 18 27 27 27 27 27 21 23 13 13 13 13 15 15 17 15 20 15 15 15 1
 7 11 23 23 28 28 28
28 28 28 24 23 24 27 28 29 30 30 30 30 27 24 24 24 28 28 30 30 30 30 22 31 28
 23 23 23 30 30 30 30 27 27 27 27 27 27 29 31 31 29 23 17 14 17 17 27 28 2
 5 24 24 24 24 24 19
13 13 13 13 13 13 13 13 26 9 18 18 23 25 28 30 30 30 28 29 26 22 22 19 21 22 29
 29 26 26 26 25 25 28 28 28 22 24 24 24 21 21 15 15 13 13 18 19 23 23 23 23 21 2
 1 25 25 26 26 26 25
27 27 21 21 21 21 21 17 13 12 12 12 12 21 16 23 24 21 21 21 24 26 28 21 21 21
 26 27 27 27 27 27 26 24 24 23 23 23 23 19 19 19 17 17 11 11 11 17 17 20 20 20 1
 1 11 11 11 11 11 11
11 11 11 17 17
>FVJHPP401A56GX rank=0000013 x=363.0 y=1199.0 length=46
```

Homopolymer errors in 454 base calling results

Sequence Read: ... GTGTGGATT AAAAAA GAGT ...

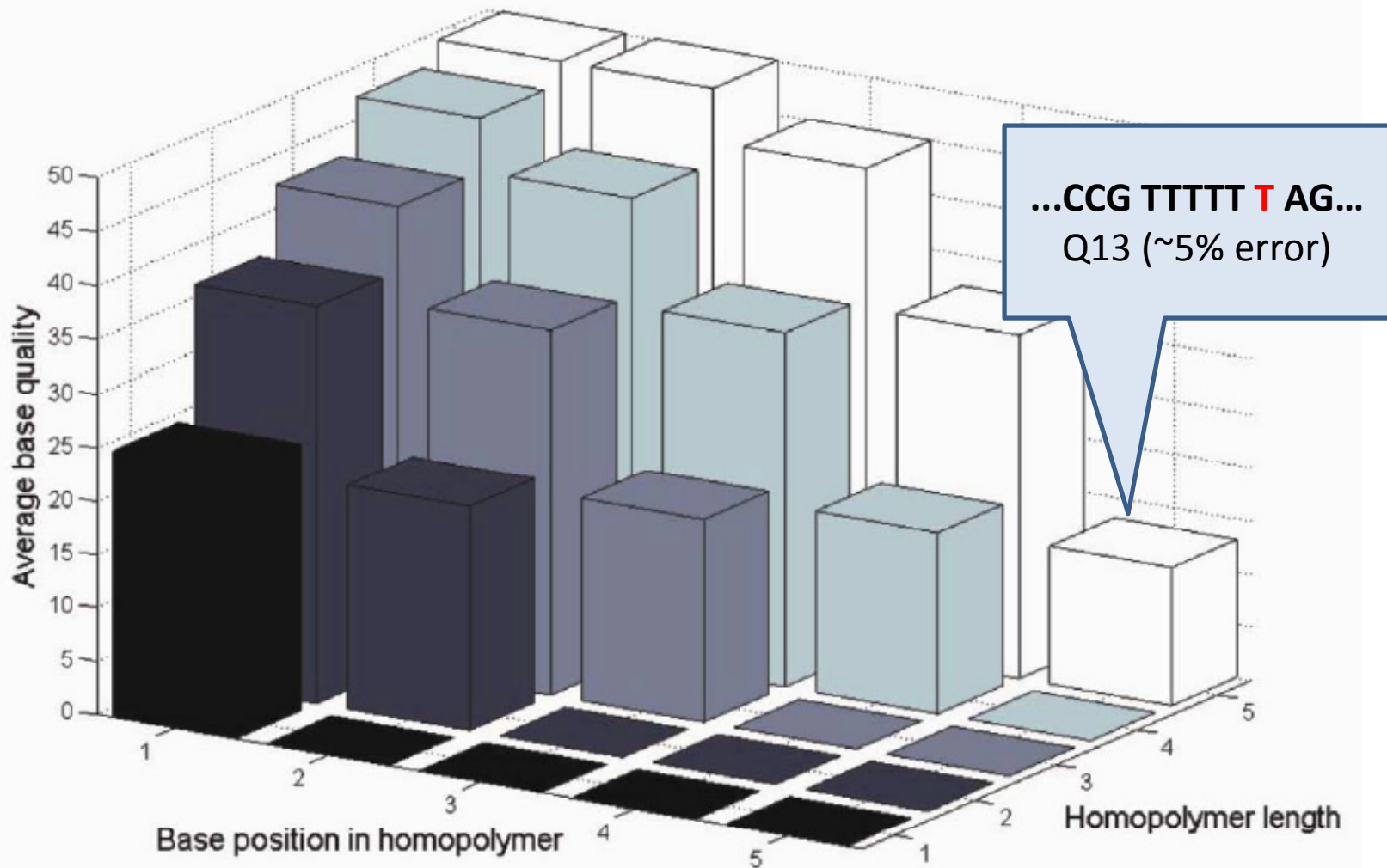
Flowgram:

6.3 A



Distribution of base quality scores within homopolymer runs

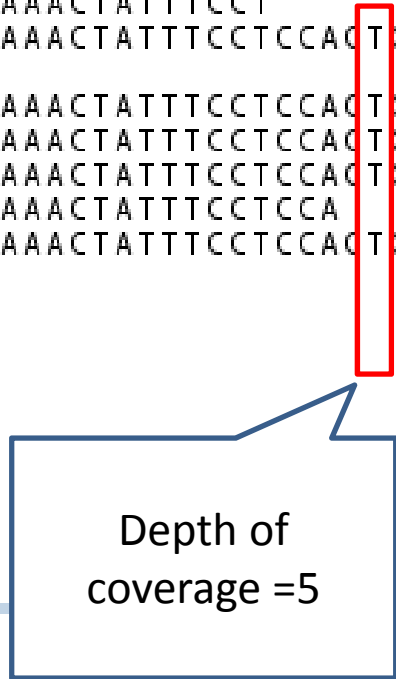
a



Quinlan AR, Marth GT et al. (2008) *Nature Methods* 5:179 - 181

Alignment of 454 reads to the reference

```
gi|170079663|ref|NC_01047... ATGATGAAGGGCAAAA-CAACGTTACCGAAAAC TATTTCTCCACTGATAACTCCTTTC|
contig00057                ATGATGAAGGGCAAAA-CAACGTTACCGAAAAC TATTTCTCCACTGATAACTCCTTTC|
136554                     ATGATGAAGGGCAAAA-CAACGTTACCGAAAAC TATTTCT
189013                     ATGATGAAGGGCAAAA-CAACGT-ACCGAAAAC TATTTCTCCACTGA
129193                     ATGATGAAGGGCAAAA-CAACG
141753                     ATGATGAAGGGCAAAA-CAACGTTACCGAAAAC TATTTCTCCACTGATAACTCCTTTC|
174302                     ATGATGAAGGGCAAAA-CAACGTTACCGAAAAC TATTTCTCCACTGATAACTCCTTTC|
37009                      ATGATGAAGGGCAAAA-CAACGTTACCGAAAAC TATTTCTCCACTGATAACTCCTTTC|
185657                     ATGATGAAGGGCAAAAACAACGTTACCGAAAAC TATTTCTCCA
5163                       AAGGGCAAAA-CAACGTTACCGAAAAC TATTTCTCCACTGATAACTCCTTTC|
169688                     ATA ACTCCTTTC|
129622                     ATA ACTCCTTTC|
```



Depth of
coverage =5

Illumina sequencing data are distributed as FASTQ files

Each sequence read is represented by 4 lines

```
@HWUSI-EAS690_0001:8:3:274:1957#0/1
CAGTAGCAATCCAACTTTGTTACTCGTCAGAAAATCGAAATCATCTTCGGTTAAATCCAAA
+HWUSI-EAS690_0001:8:3:274:1957#0/1
aabaaaaaaaaa_Y_aabaaaaaaaaa`a`[a`a```[``[``^`\\^\\_aZOY^^\\_BB
```

Line 1. Read Identifier 1

Line 2. Sequence

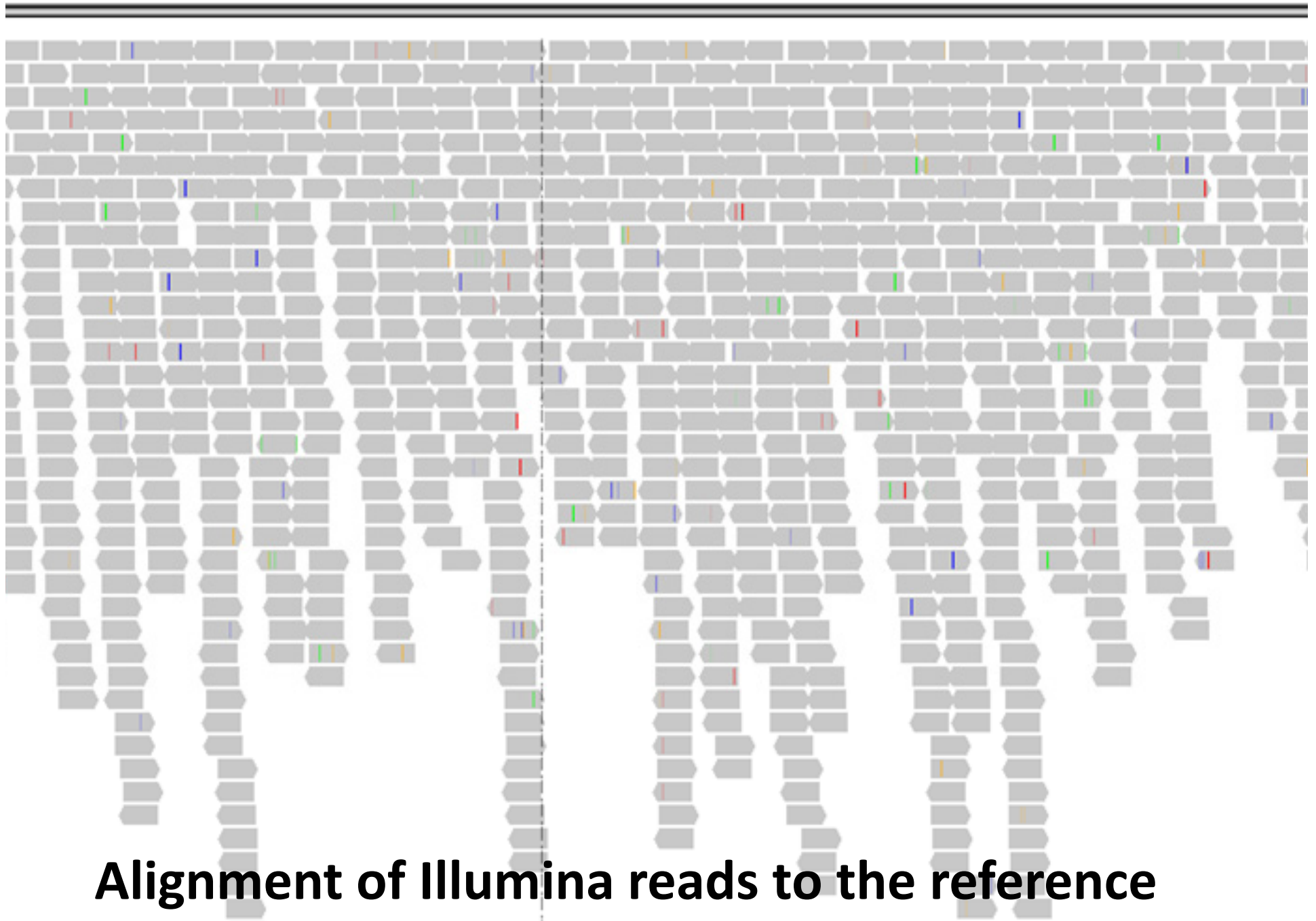
Line 3. Read Identifier 2

Line 4. Quality code

Illumian Quality Code (post v 1.3)

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

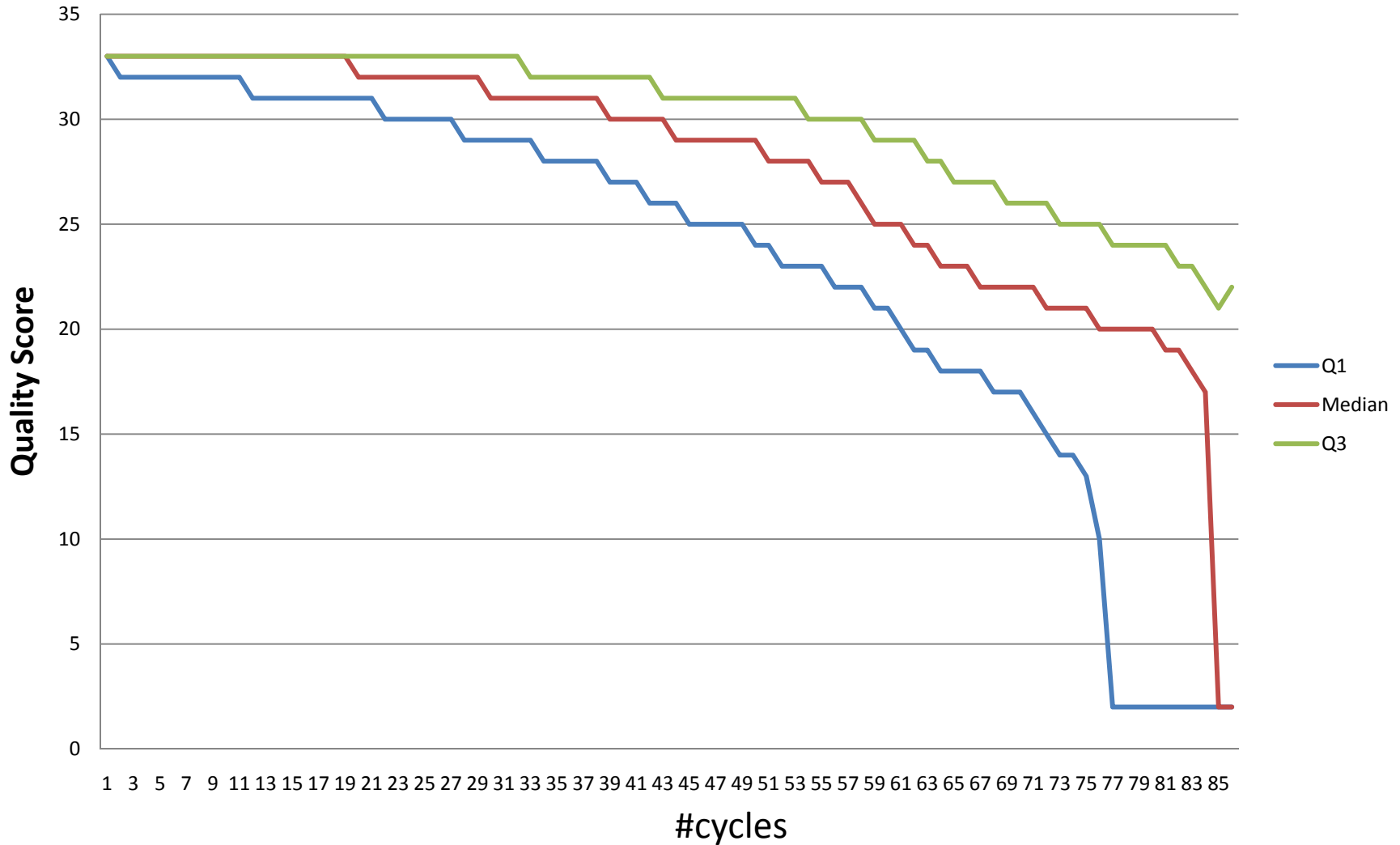
CODE	Quality	P	CODE	Quality	P
@	0	1.0000	T	20	0.0100
A	1	0.7943	U	21	0.0079
B	2	0.6310	V	22	0.0063
C	3	0.5012	W	23	0.0050
D	4	0.3981	X	24	0.0040
E	5	0.3162	Y	25	0.0032
F	6	0.2512	Z	26	0.0025
G	7	0.1995	[27	0.0020
H	8	0.1585	\	28	0.0016
I	9	0.1259]	29	0.0013
J	10	0.1000	^	30	0.0010
K	11	0.0794	_	31	0.0008
L	12	0.0631	`	32	0.0006
M	13	0.0501	a	33	0.0005
N	14	0.0398	b	34	0.0004
O	15	0.0316	c	35	0.0003
P	16	0.0251	d	36	0.0003
Q	17	0.0200	e	37	0.0002
R	18	0.0158	f	38	0.0002
S	19	0.0126	g	39	0.0001



Alignment of Illumina reads to the reference

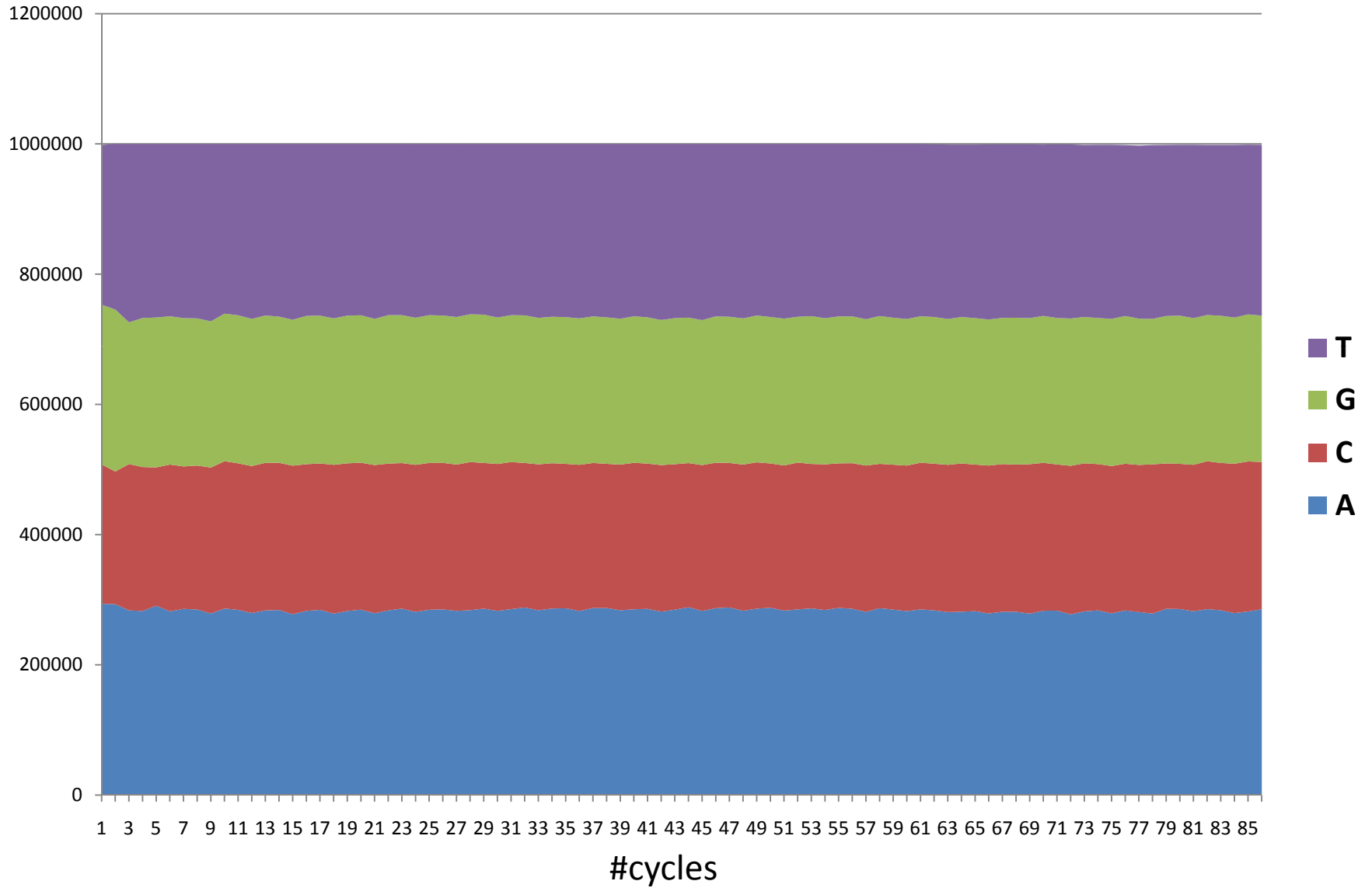
-- the colored lines indicate substitution errors

Quality Score Distribution of Illumina Sequencing Cycles



Q1, Median, Q3 refer to 25%, 50% or 75% of reads with Quality below this level

Nucleotides Distribution of Illumina Sequencing Cycles

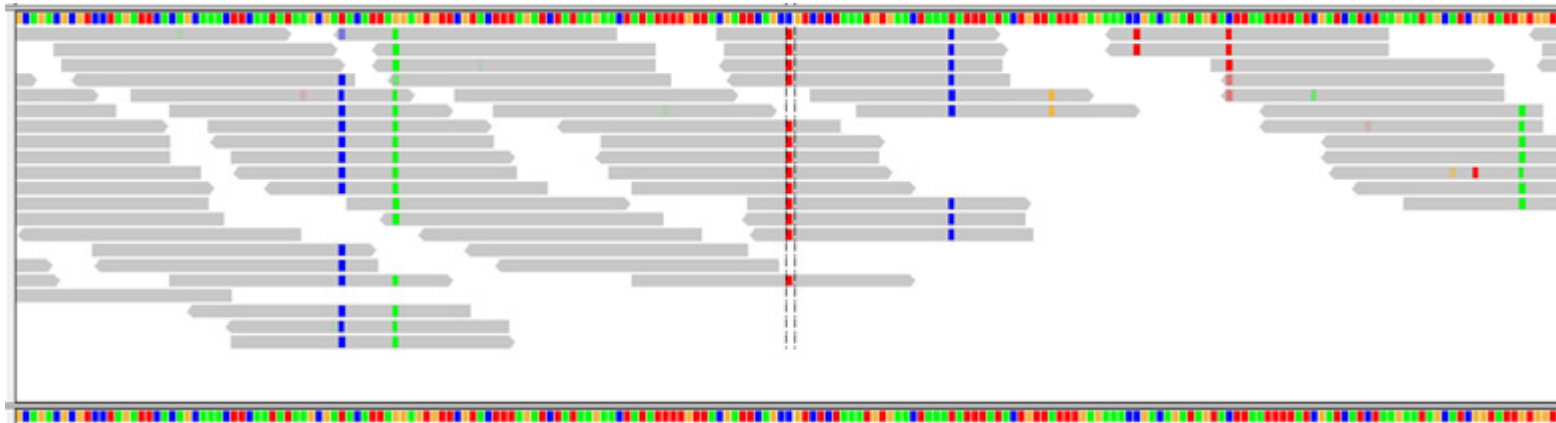


When should sequencing errors be a concern?

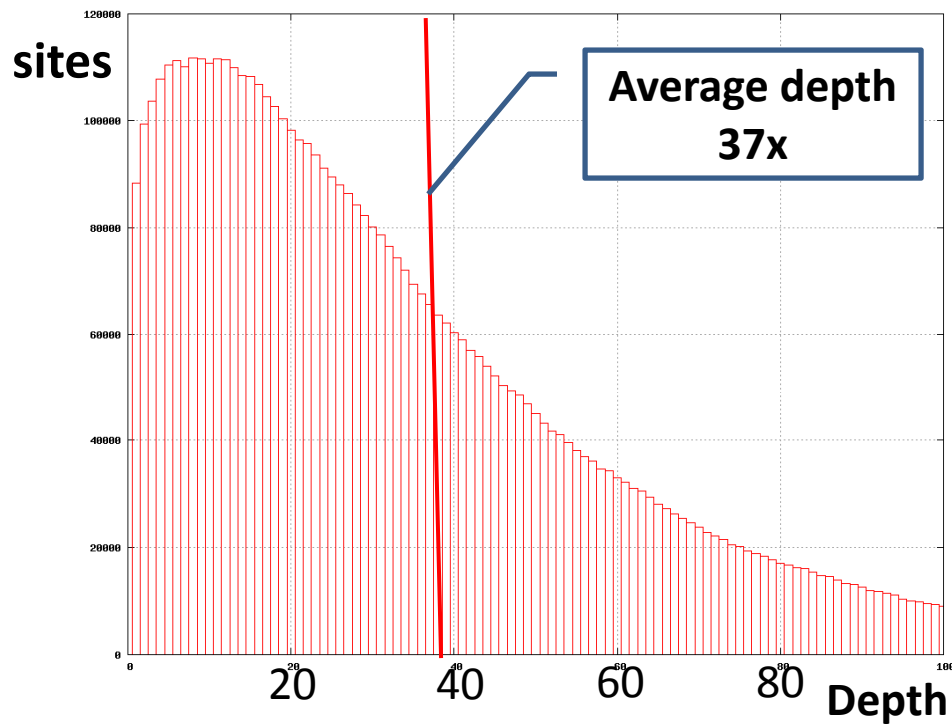
1. For RNAseq/ChIPSeq data, sequencing errors are not a major concern.

A limited number of mismatches can be tolerate. However, for some data files, trimming of low quality data are recommended.

2. For SNP detection, especially when coverage depth is low, the quality score can be used, but not always reliable.



The coverage depth are not evenly distributed across the genome.

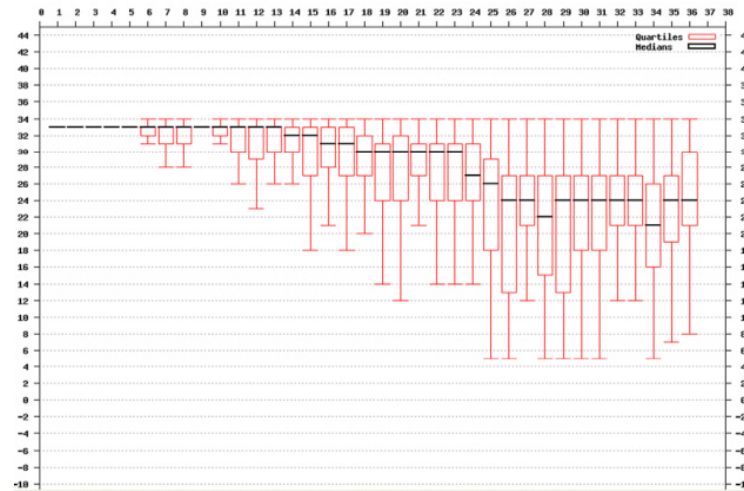


Coverage Depth	% base pairs
≤ 1	1.1 %
≤ 3	4.3 %
≤ 5	10%
≤ 10	20%
≤ 37	63%

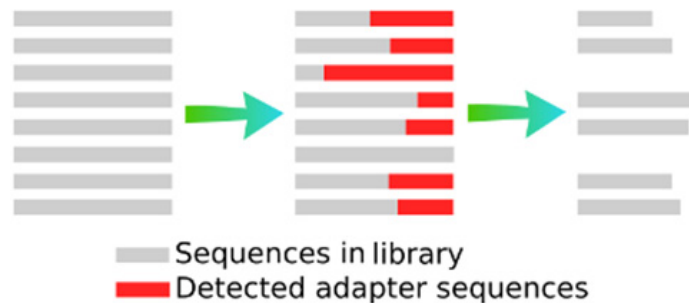
Distribution of Coverage depth in a seq-cap experiment.

FASTX ToolKit – a package for pre-processing Illumina FASTQ files and evaluating data quality (by Hannon Lab of CSHL)

Plot quality scores and trim low-quality data



Clip off the 3' adaptor



FASTX ToolKit – a package for pre-processing Illumina FASTQ files and evaluating data quality (by Hannon Lab of CSHL)

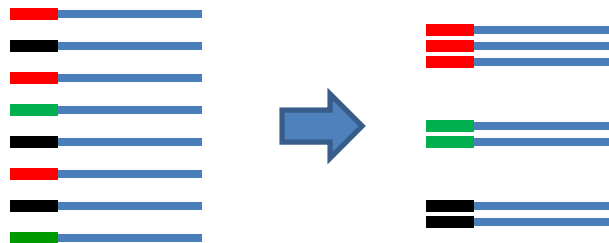
FASTQ -> FASTA

```
@HWUSI-EAS690_0001:8:3:274:1957#0/1  
CAGTAGCAATCCAAACTTTGTTACTCGTCAGAAAATCGAAA  
+HWUSI-EAS690_0001:8:3:274:1957#0/1  
aabaaaaaaaaa_Y_aabaaaaaaaaa`a`[a`a````[`
```



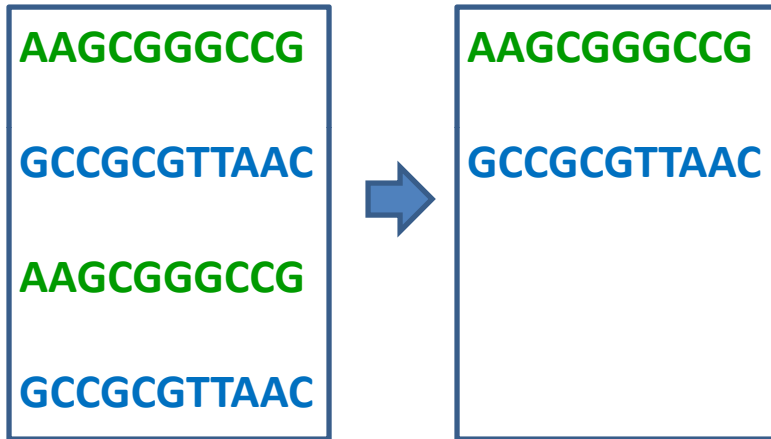
```
>1  
CAGTAGCAATCCAAACTTTGTTACTCGTCAGAAAATCGAAA  
>2  
GATTTAACCGGCCGGAAAACCCTCGACATACGGATACCGAT
```

De-multiplexing



FASTX ToolKit – a package for pre-processing Illumina FASTQ files and evaluating data quality (by Hannon Lab of CSHL)

Remove redundancy



Web Interface

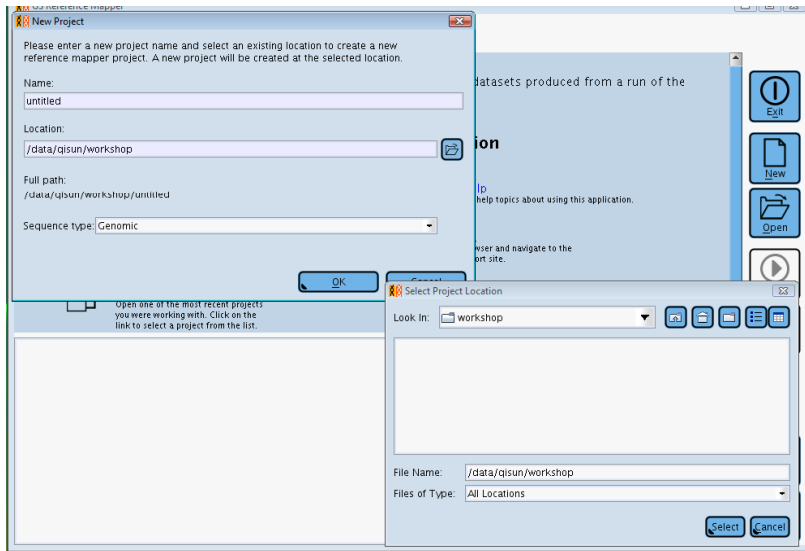
Trim

Library to clip:

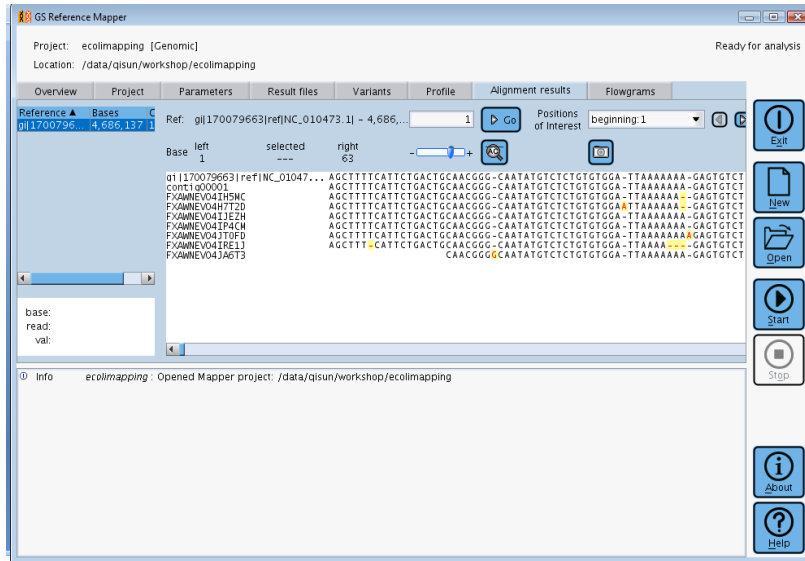
First base to keep:

Last base to keep:

Roche 454 Software

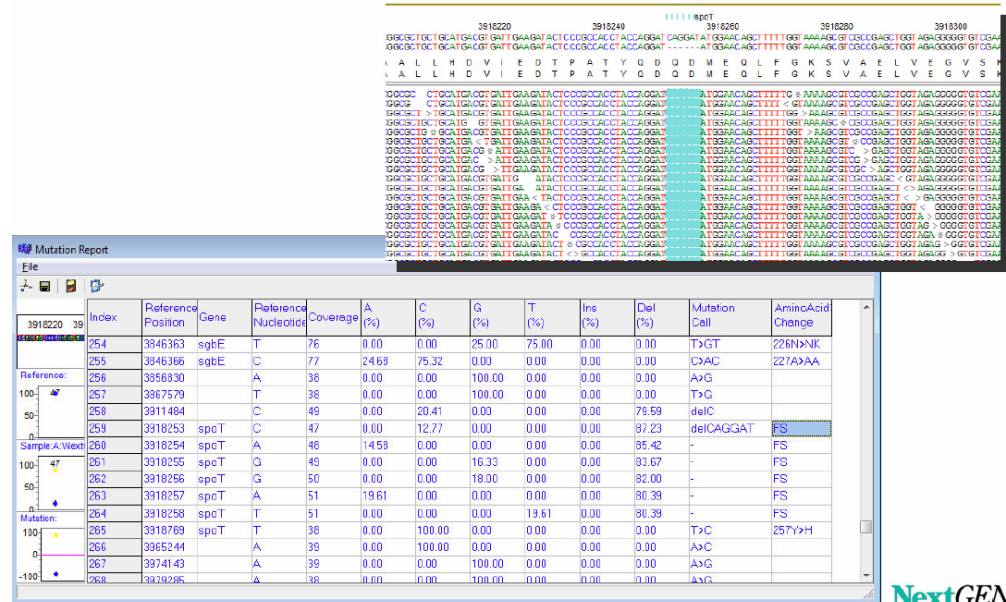
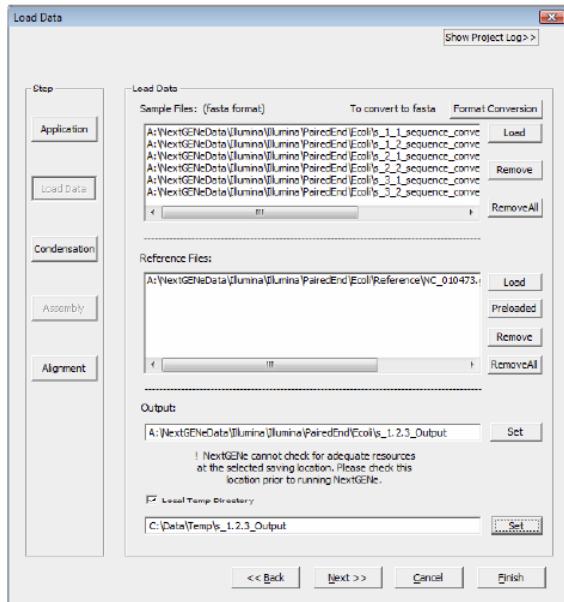


- Align reads to genome
- De novo assembly
- Amplicon diversity



Computational Resources at Cornell

1. Commercial Software workstation (624 Rhodes Hall)



SoftGenetics NextGENe
DNASTar NGen
Illumina Genome Studio

Computational Resources at Cornell

2.CAC Linux Cluster

- Linux**
- Common open source software preinstalled**

- V4 16GB RAM 19 servers**
- V4-64G 64GB RAM 3 servers**

3. Research collaboration with CBSU

Next-Gen@BioHPC

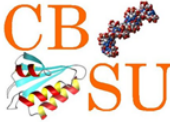
QC

RNA-seq

ChIP-seq

SNP/Indel

Assembly



Next-Gen @ BioHPC

[Log out](#)

[Home](#)

[Run Manager](#)

[Lane Browser](#)

[Reference Mgt](#)

[Change password](#)

[Reset password](#)

[BioHPC @ CBSU](#)

[contact CBSU](#)

BROWSE ALL LANES

This is a list of all lanes configured in the system. To sort results, click on a column header. To filter results, supply templates (all usual wildcards apply) and click **Apply Filters**. Clicking on an entry in **Run ID** column will open the run manager utility, where the lane information can be configured. To download files for a lane, click on the link **(files)** underneath the LaneID. For assistance with download of multiple files in batch mode, use the **Make download script** button.

[Make download script](#)

("check" all lanes you want to include in a download script and click button above)

[Register new lane](#)

(use this option only if you want to manually upload a lane from outside of the CLC sequencing facility)

Filter data by

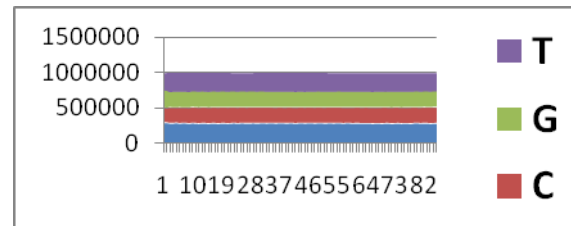
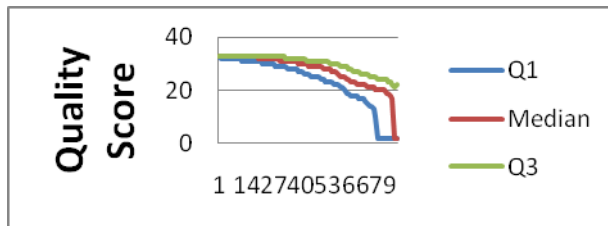
Status: Active Sample name: * Users: *

Page 1 of 8

Lane ID	Run ID	Run Name	Lane#	Type	Sample Name	Status	Annotations	Users	Lab	Order#												
151 <input type="checkbox"/> (files)	28	100301_HWI-EAS339_0004_01GGLAAXX	1	Standard	[img]	ready	<table style="width: 100%; border-collapse: collapse;"> <tr> <td style="text-align: right;">Parameter</td> <td style="text-align: center;">This Lane</td> <td style="text-align: center;">Ctrl Lane</td> </tr> <tr> <td style="text-align: right;">Length</td> <td style="text-align: center;">unknown</td> <td style="text-align: center;">unknown</td> </tr> <tr> <td style="text-align: right;">Clusters_raw</td> <td style="text-align: center;">25.5M</td> <td style="text-align: center;">25.5M</td> </tr> <tr> <td style="text-align: right;">Clusters_PF</td> <td style="text-align: center;">19.8M</td> <td style="text-align: center;">21.3M</td> </tr> </table>	Parameter	This Lane	Ctrl Lane	Length	unknown	unknown	Clusters_raw	25.5M	25.5M	Clusters_PF	19.8M	21.3M	[img]@cornell.edu (owner)	N/A	10214943
Parameter	This Lane	Ctrl Lane																				
Length	unknown	unknown																				
Clusters_raw	25.5M	25.5M																				
Clusters_PF	19.8M	21.3M																				
152 <input type="checkbox"/> (files)	28	100301_HWI-EAS339_0004_01GGLAAXX	2	Standard	[img]	ready	<table style="width: 100%; border-collapse: collapse;"> <tr> <td style="text-align: right;">Parameter</td> <td style="text-align: center;">This Lane</td> <td style="text-align: center;">Ctrl Lane</td> </tr> <tr> <td style="text-align: right;">Length</td> <td style="text-align: center;">unknown</td> <td style="text-align: center;">unknown</td> </tr> <tr> <td style="text-align: right;">Clusters_raw</td> <td style="text-align: center;">24.6M</td> <td style="text-align: center;">25.5M</td> </tr> <tr> <td style="text-align: right;">Clusters_PF</td> <td style="text-align: center;">21.0M</td> <td style="text-align: center;">21.3M</td> </tr> </table>	Parameter	This Lane	Ctrl Lane	Length	unknown	unknown	Clusters_raw	24.6M	25.5M	Clusters_PF	21.0M	21.3M	[img]@cornell.edu (owner)	N/A	10214943
Parameter	This Lane	Ctrl Lane																				
Length	unknown	unknown																				
Clusters_raw	24.6M	25.5M																				
Clusters_PF	21.0M	21.3M																				
							<table style="width: 100%; border-collapse: collapse;"> <tr> <td style="text-align: right;">Parameter</td> <td style="text-align: center;">This Lane</td> <td style="text-align: center;">Ctrl Lane</td> </tr> </table>	Parameter	This Lane	Ctrl Lane												
Parameter	This Lane	Ctrl Lane																				

Exercises:

1. Using FASTX to get statistic report of Illumina sequencing data



2. Using Roch GSMapper to align 454 reads to E.coli genome.

The screenshot shows the Roch GSMapper interface. The reference sequence is `gij170079663|ref|NC_010473.1| - 4,686,...`. The position of interest is `beginning:1`. The alignment shows the reference sequence and several reads (contigs) aligned to it. The reads are:

```
qi|170079663|ref|NC_01047... AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
contig00001 AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
FXAWNEV04IH5MC AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
FXAWNEV04H7T2D AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
FXAWNEV04IJEZH AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
FXAWNEV04IP4CM AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
FXAWNEV04JTOFD AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
FXAWNEV04IRE1J AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
FXAWNEV04JA6T3 AGCTTTTCATTCTGACTGCAACGGG-CAATATGTCTCTGTGTGGA-TTAAAAAAA-GAGTGTCT
```

1. Download workshop slides and exercise instructions from the workshop web site.
<http://cbsu.tc.cornell.edu>

2. Setup your desktop computer to work with Linux server.

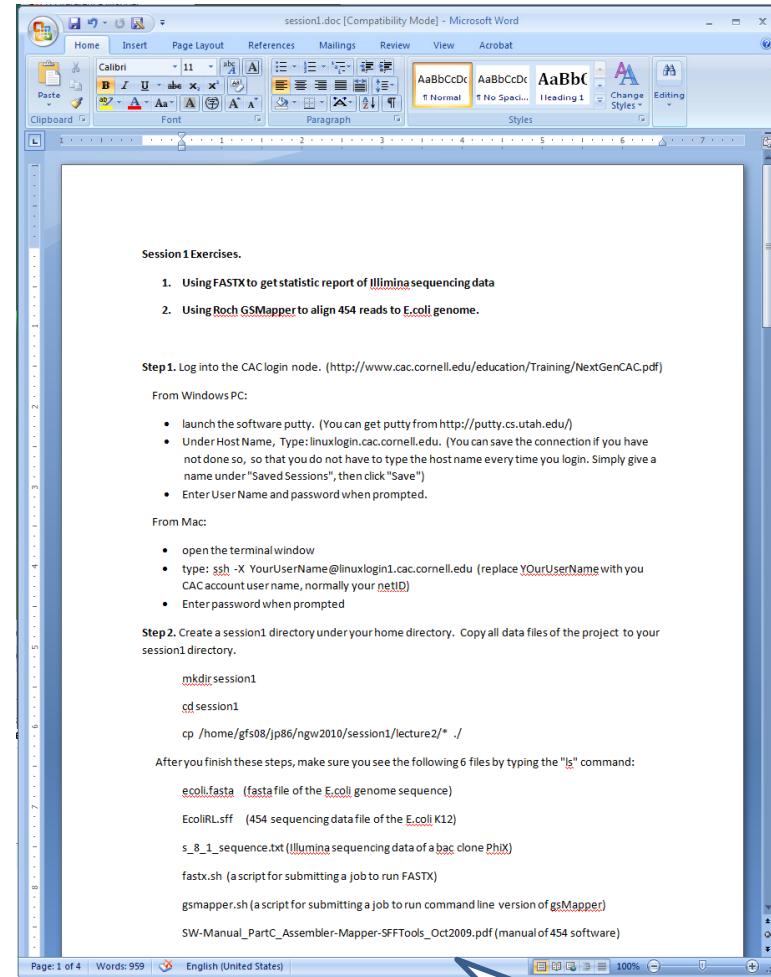
Window Users: Install Putty, WinSCP, Xming, EditPlus

Mac users: Install TextWrangler, FETCH

3. Log into the linuxlogin.cac.cornell.edu

4. Modify the fastx.sh script as described.

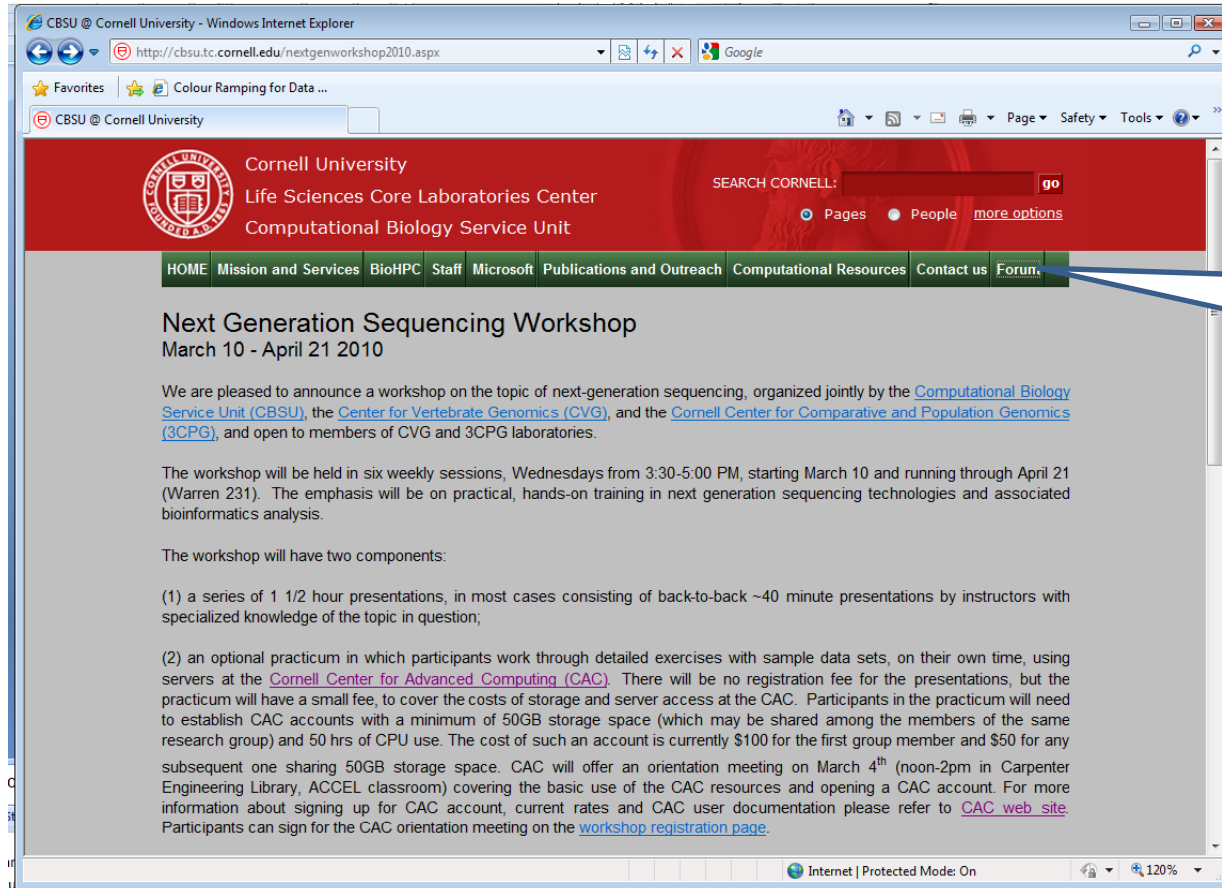
5. Submit the job.



A few tips of using Linux system is available at the end of the document.

Forum

Please post your questions to the online forum. You need register to post questions.



Click to enter
forum

Office hours:

Friday 3:00PM. 102 Weill (small conference room)