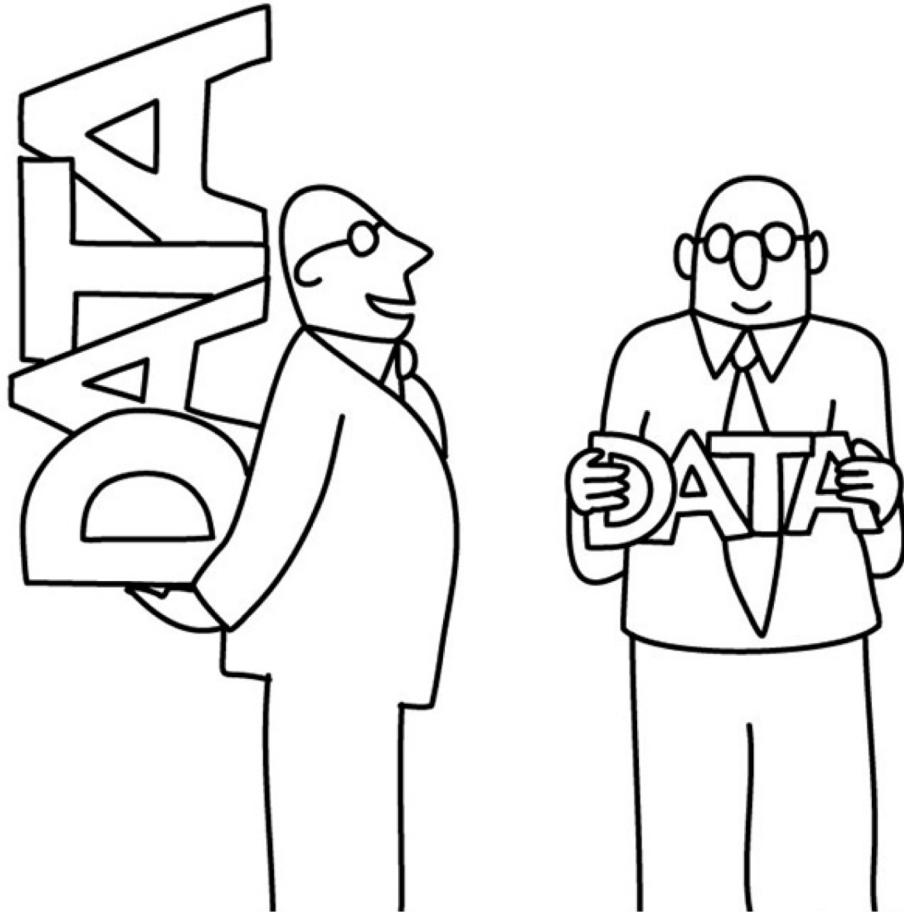


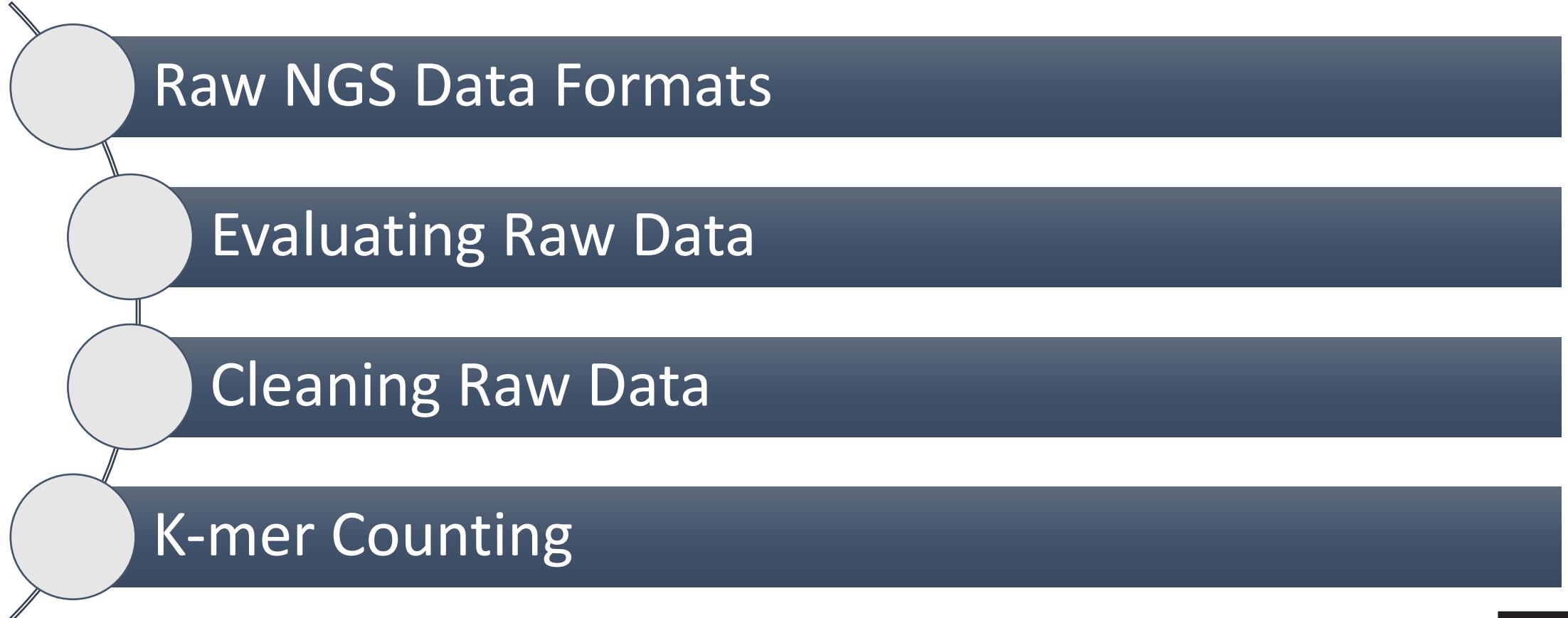
# NGS QC

An introduction

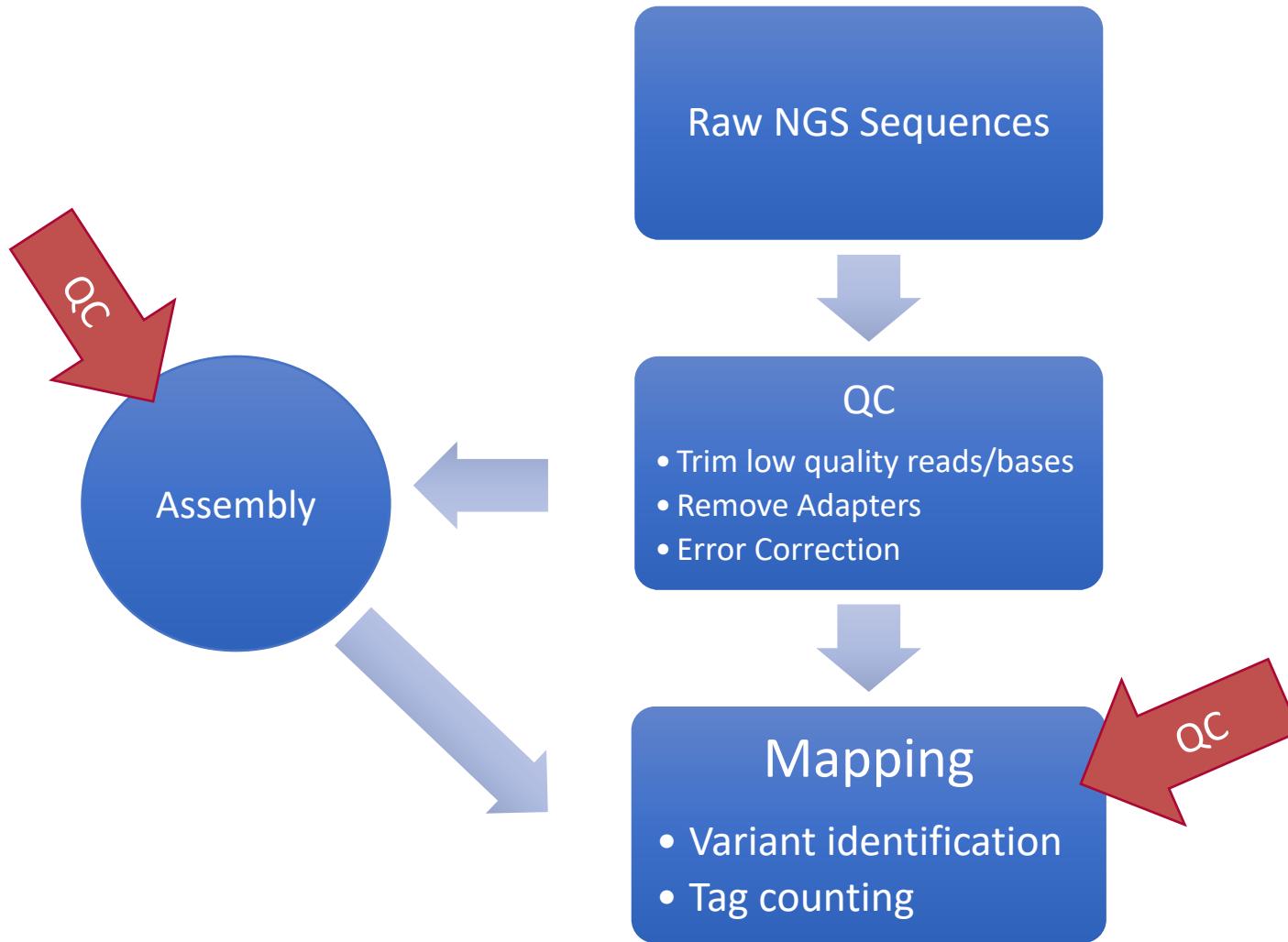


*“I think you'll find that mine is bigger...”*

# Outline



# The Big Picture



# Common Sequence Formats

## FASTA

- Simple
- Nucleotide or amino acid strings
- No quality info
- Compressible (.gz)

## FASTQ

- Mildly complex
- Nucleotide strings (not AA)
- Quality information included
- Compressible (.gz)

## FAST5

- Complex (HDF5)
- Nanopore Data
- Nucleotide strings (not AA)
- Raw *squiggles*
- Natively compressed



# Common Sequence Formats

## FASTA

- Simple
- Nucleotide or amino acid strings
- No quality info
- Compressible (.gz)

## FASTQ

- Mildly complex
- Nucleotide strings (not AA)
- Quality information included
- Compressible (.gz)

## FAST5

- Complex (HDF5)
- Nanopore Data
- Nucleotide strings (not AA)
- Raw *squiggles*
- Natively compressed



# The FASTA format

```
>sequence 1
CATCGATCGCATGCTACTGACTG
CATGCTCGCGCCCCCCCCGATG

>sequence 2
ACTGACTCGCGCGCGCGGGGG
GAGCTGATGTG

>sequence 3
CATCGATCGCATGCTACTGACTG
CATGCTCGCGCCCCCCCCGATG
ACTGACTCGCGCGCGCGGGGG
GAGCTGATGTG
```



# The FASTA format

```
>sequence 1  
CATCGATCGCATGCTACTGACTGCATGCTCGGCCCGATG.....  
>sequence 2  
ACTGACTCGCGCGCGGGGGAGCTGATGTG  
>sequence 3  
CATCGATCGCATGCTACTGACTGCATGCTCGGCCCGATGAC...
```



# The FASTQ format

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
+ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
hhhhhhhhhhghhhhhhehhhedhhhhfhhhhh
```



# The FASTQ format

Sequence ID

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
+ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
hhhhhhhhhhghhhhhhehhhedhhhhfhhhhh
```



# The FASTQ format

## Sequence

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
+ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
hhhhhhhhhhghhhhhhehhhedhhhhfhhhhh
```



# The FASTQ format

+ description (or empty)

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
+ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
hhhhhhhhhhghhhhhhehhhedhhhhfhhhhh
```



# The FASTQ format

+ description (or empty)

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1  
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
```

+

```
hhhhhhhhhhghhhhhhehhedhhhfhhhhhh
```



# The FASTQ format

Quality score of each base

```
@ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
GTAGAACTGGTACGGACAAGGGGAATCTGACTGTAG
+ILLUMINA-F6C19_0048_FC:5:1:12440:1460#0/1
hhhhhhhhhhghhhhhhehhhedhhhhfhhhhh
```



# Illumina Sequence ID Lines: A Decoder

**@M01137:30:00000000-AA299:1:1101:10929:1966**

M01137	the unique instrument name
30	the run id
00000000-AA29	the flowcell id
1	flowcell lane
1101	tile number within the flowcell lane
10929	'x'-coordinate of the cluster within the tile
1966	'y'-coordinate of the cluster within the tile
1 or 2 (not shown, optional)	the member of a pair, 1 or 2 ( <i>paired-end or mate-pair reads only</i> )
ATCACG (not shown, optional)	index sequence



# Quality Scores

- Phred Score
- $Q = -10^{\star} \log_{10} P$        $P$  = probability the base call is incorrect
- ASCII (character) - 33

Phred Quality Score	Probability of incorrect base call	Base call accuracy
0	1	0 %
10	1 in 10	90 %
20	1 in 100	99 %
30	1 in 1000	99.9 %
40	1 in 10000	99.99 %
50	1 in 100000	99.999 %
93	1 in 2000000000	99.9999995 %

# Why QC NGS Data?

OPEN  ACCESS Freely available online



## An Extensive Evaluation of Read Trimming Effects on Illumina NGS Data Analysis

Cristian Del Fabbro<sup>1</sup>✉, Simone Scalabrin<sup>2</sup>✉, Michele Morgante<sup>1</sup>, Federico M. Giorgi<sup>1,3\*</sup>

“Trimming is shown to increase the quality and reliability of the analysis, with concurrent gains in terms of execution time and computational resources needed”



# Types of Trimming

## Quality

- remove low quality bases and reads
  - Q20 (1% error) and Q30 (0.1% error) are standard
- Remove too short reads
- Too many 'N' (uncalled bases)

## Complexity

- simple repeats (e.g. TGTGTGTGTG)
- Homopolymers (e.g., AAAAAAAA)

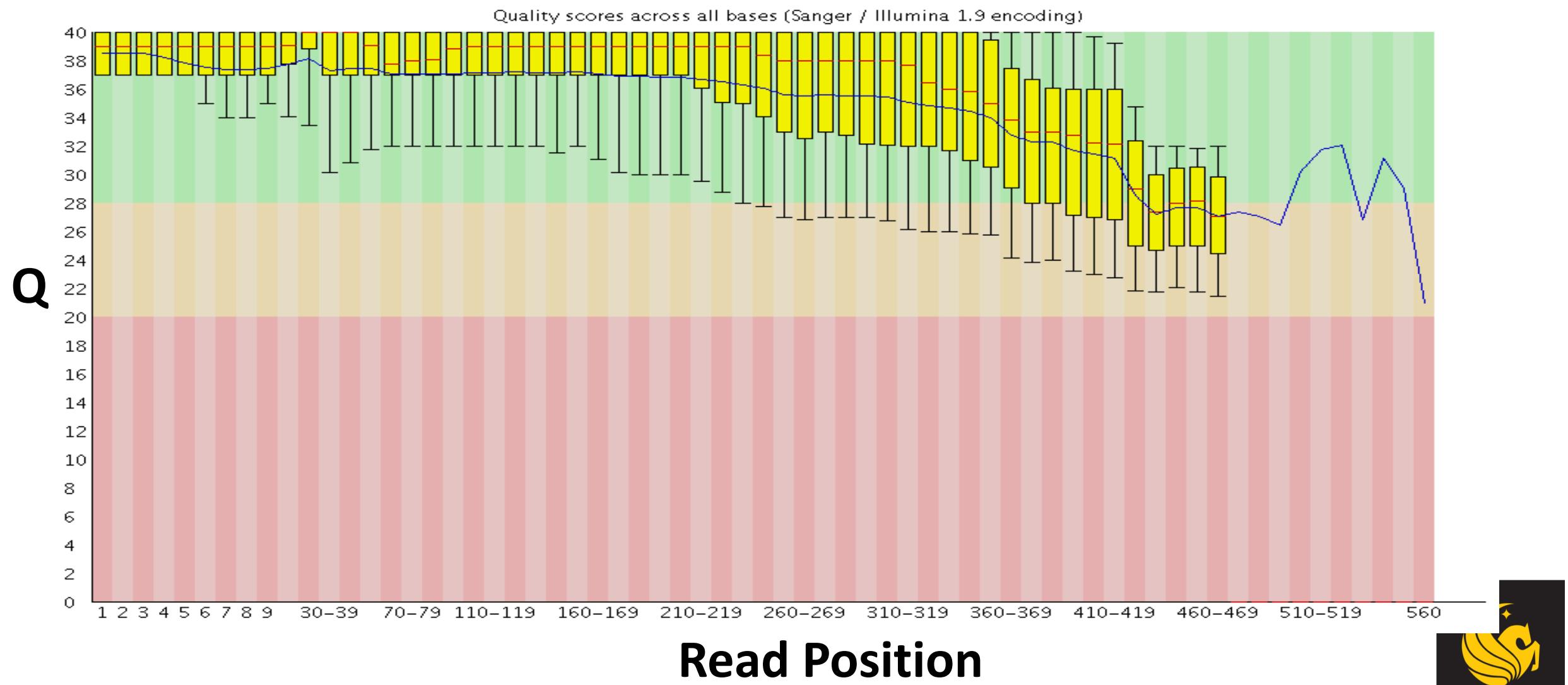
## Contamination

- Sequencing adapters!!!!!!
- lab contamination (human, bacteria)
- environmental

# Low Quality Sequences Before Trimming (Puma 454 sequences)



# Same Sequences After Trimming (Puma 454 sequences)



# Types of Trimming

- 
- High Quality
  - Low quality

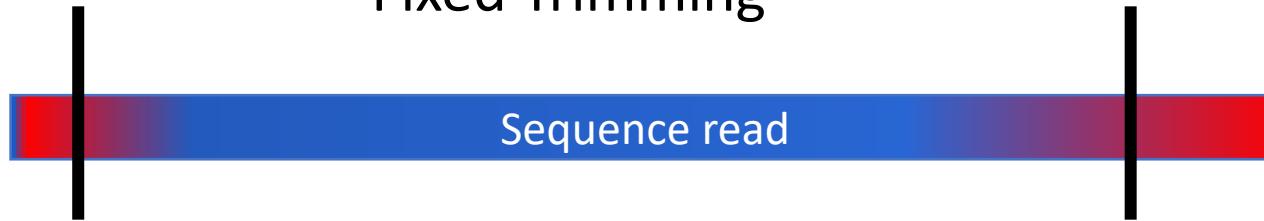


Sequence read

# Types of Trimming

- 
- High Quality
  - Low quality

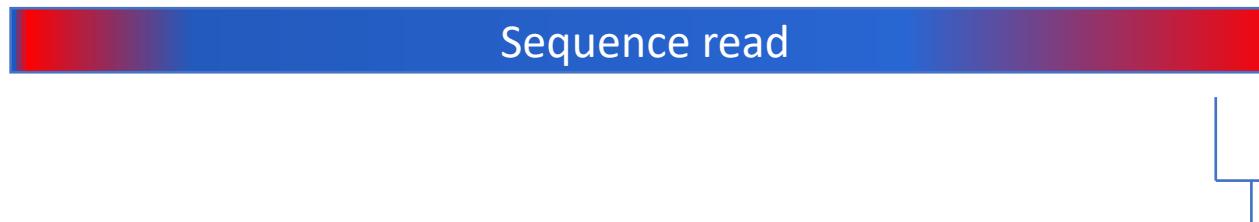
## Fixed Trimming



# Types of Trimming

- 
- High Quality
  - Low quality

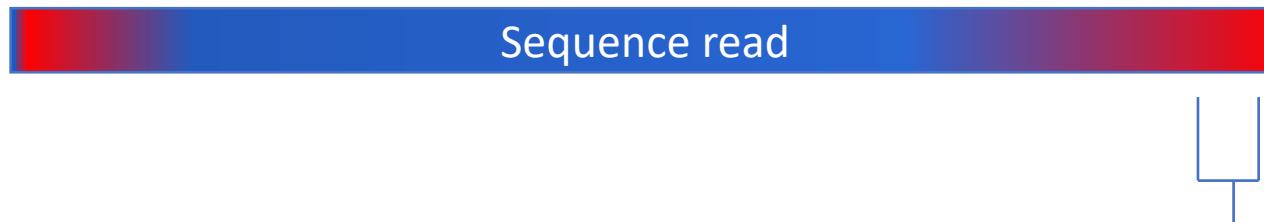
## Sliding Window Trimming



# Types of Trimming

- 
- High Quality
  - Low quality

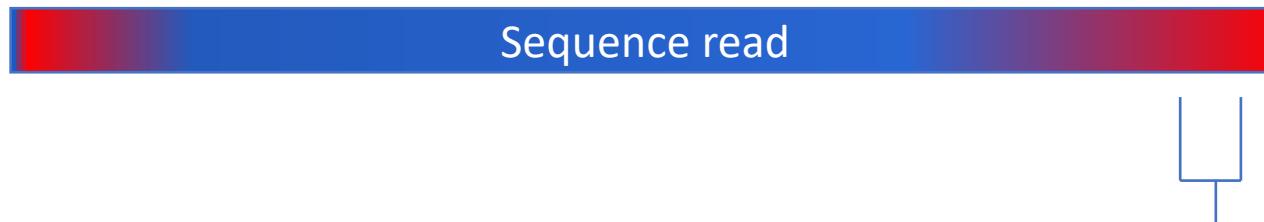
## Sliding Window Trimming



# Types of Trimming

- 
- High Quality
  - Low quality

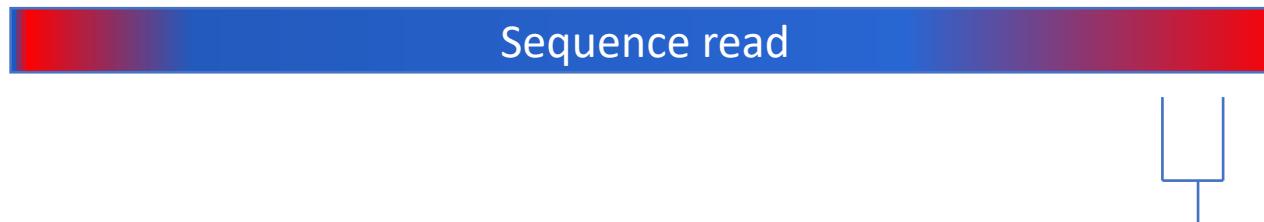
## Sliding Window Trimming



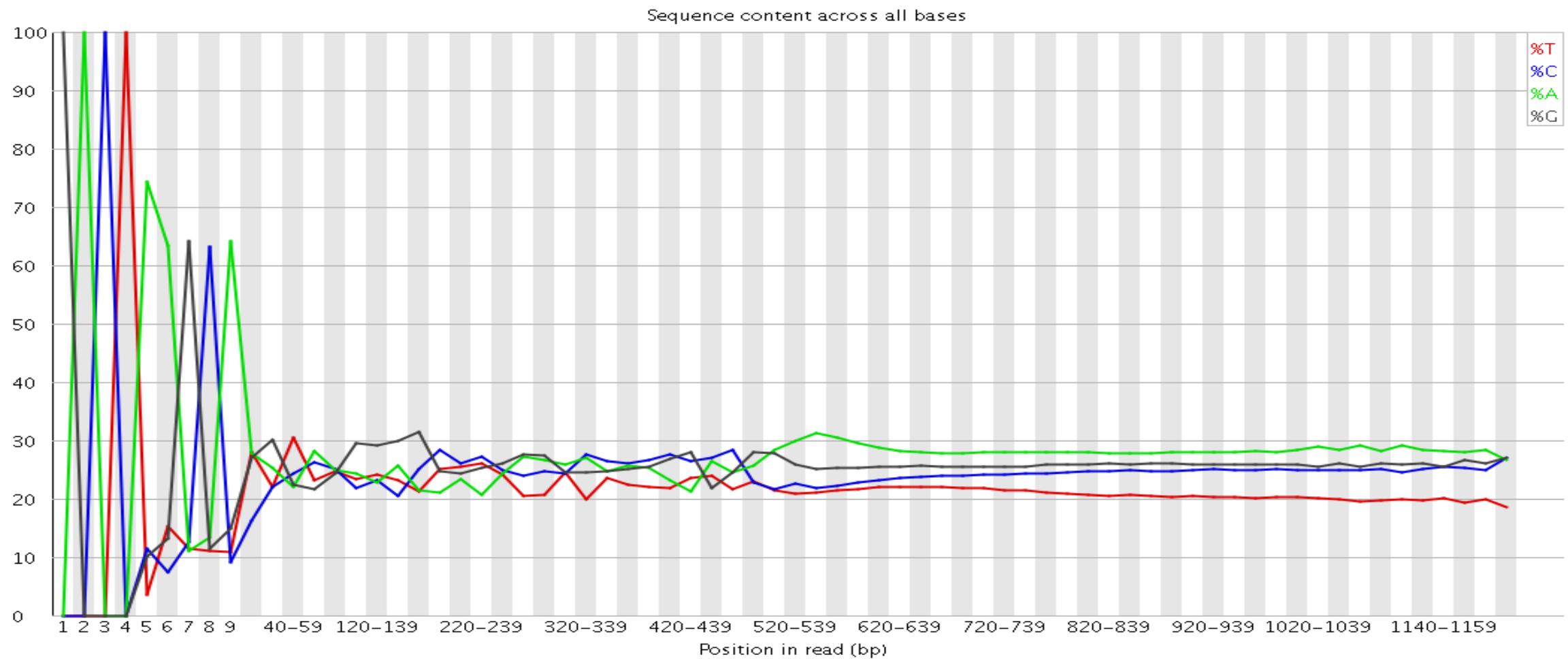
# Types of Trimming

- 
- High Quality
  - Low quality

## Sliding Window Trimming



# Adapter Contamination



# Adapter Contamination

 Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GACTAAGCAGTGGTATCAACGCAGACTACATGGGACACTTGTTCCTGAC	19391	5.415739186535921	No Hit
GACTAAGCAGTGGTATCAACGCAGACTACATGGGACACTTGCTTCCTGAC	11325	3.162974900083508	No Hit
GACTAAGCAGTGGTATCAACGCAGACTACATGGGACACTTGTTCCTGACA	9229	2.5775801636088915	No Hit
.....	.....	.....	.....

 Download [Graphics](#)

gnl|uv|NGB00593.1:1-30 Evrogen Mint PlugOligo-1 adapter  
Sequence ID: Length: 30 Number of Matches: 1

Range 1: 1 to 25 [Graphics](#) ▼ Next Match ▲ Previous Match

Score	Expect	Identities	Gaps	Strand
50.6 bits(25)	0.001	25/25(100%)	0/25(0%)	Plus/Plus

Query 5 AAGCAGTGGTATCAACGCAGAGTAC 29  
Sbjct 1 AAGCAGTGGTATCAACGCAGAGTAC 25

# Error Correction (Illumina data)

## GAGE: A critical evaluation of genome assemblies and assembly algorithms

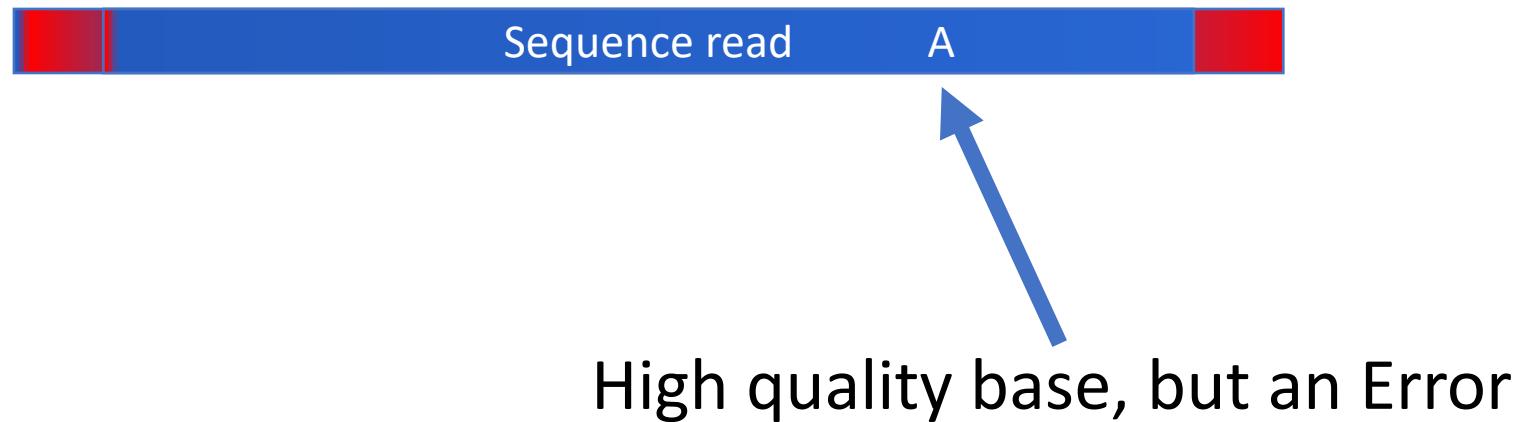
Steven L. Salzberg,<sup>1,7</sup> Adam M. Phillippy,<sup>2</sup> Aleksey Zimin,<sup>3</sup> Daniela Puiu,<sup>1</sup> Tanja Magoc,<sup>1</sup> Sergey Koren,<sup>2,4</sup> Todd J. Treangen,<sup>1</sup> Michael C. Schatz,<sup>5</sup> Arthur L. Delcher,<sup>6</sup> Michael Roberts,<sup>3</sup> Guillaume Marçais,<sup>3</sup> Mihai Pop,<sup>4</sup> and James A. Yorke<sup>3</sup>

“For all four genomes and for all eight assemblers used in GAGE, the best assemblies were created from reads that had been processed through extensive error correction routines”

Illumina Sequencing Errors: ~0.1 - 1%, Substitution errors



# Error Correction



# Error Correction: *K*-mer Counting

$k = 4$

AGCTGTGG



# Error Correction: *K*-mer Counting

$k = 4$

AGCTGTGG  


AGCT

# Error Correction: *K*-mer Counting

$k = 4$

AGCTGTGG



AGCT  
GCTG

# Error Correction: *K*-mer Counting

$k = 4$

AGCTGTGG



AGCT

GCTG

CTGT

# Error Correction: *K*-mer Counting

$k = 4$

AGCTGTGG  


AGCT  
GCTG  
CTGT  
TGTG



# Error Correction: *K*-mer Counting

$k = 4$

AGCTGTGG



AGCT  
GCTG  
CTGT  
TGTG  
GTGG

# Error Correction: *K*-mer Counting

$k = 6$

AGCTGTGG



# Error Correction: *K*-mer Counting

$k = 6$

AGCTGTGG



AGCTGT

# Error Correction: *K*-mer Counting

$k = 6$

AGCTGTGG



AGCTGT  
GCTGTG

# Error Correction: *K*-mer Counting

$k = 6$

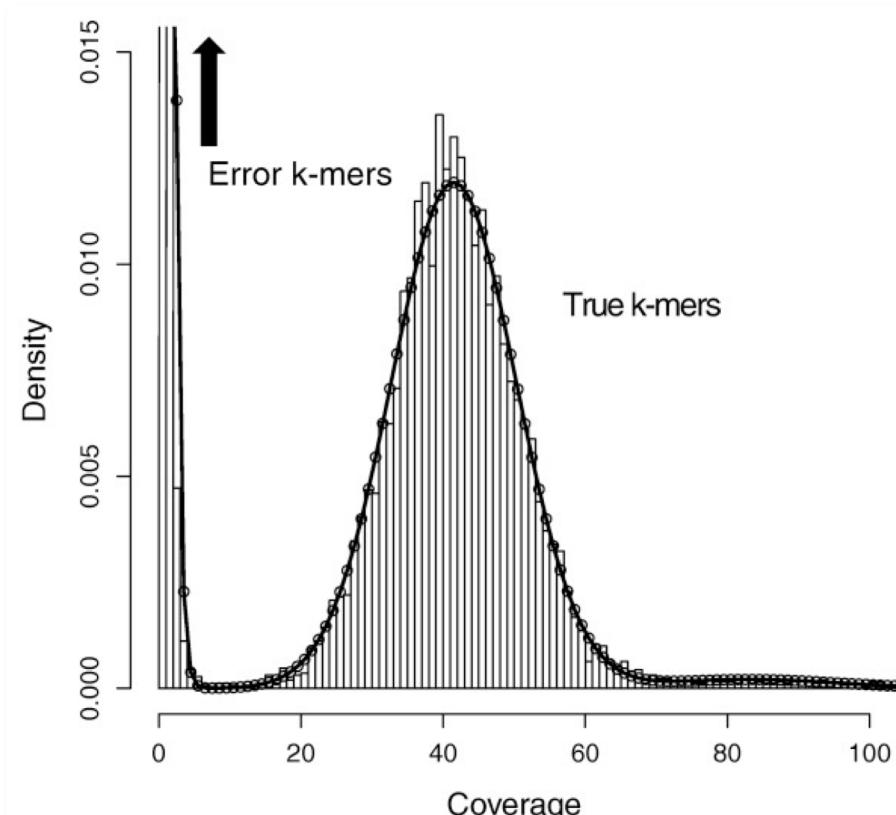
AGCTGTGG



AGCTGT  
GCTGTG  
CTGTGG

# Error Correction: *K*-mer Counting

- Expected Distribution of *k*-mer frequency



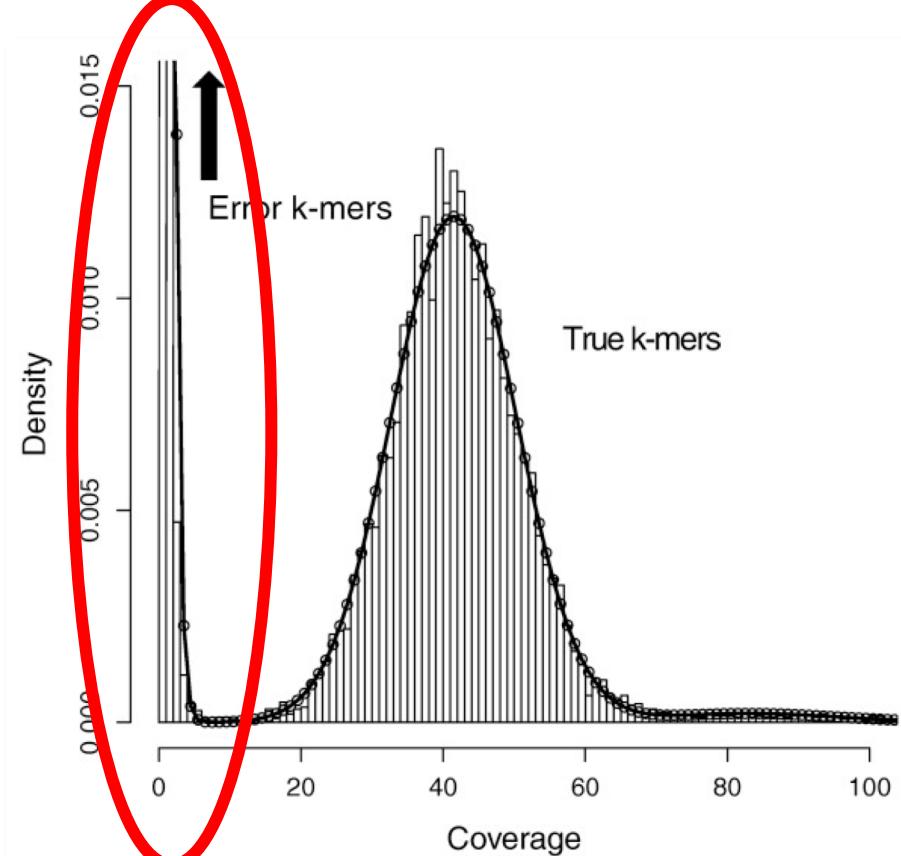
DSK; Rizk et al. 2013



# Error Correction: *K*-mer Counting

- Expected Distribution of *k*-mer frequency

Corrected



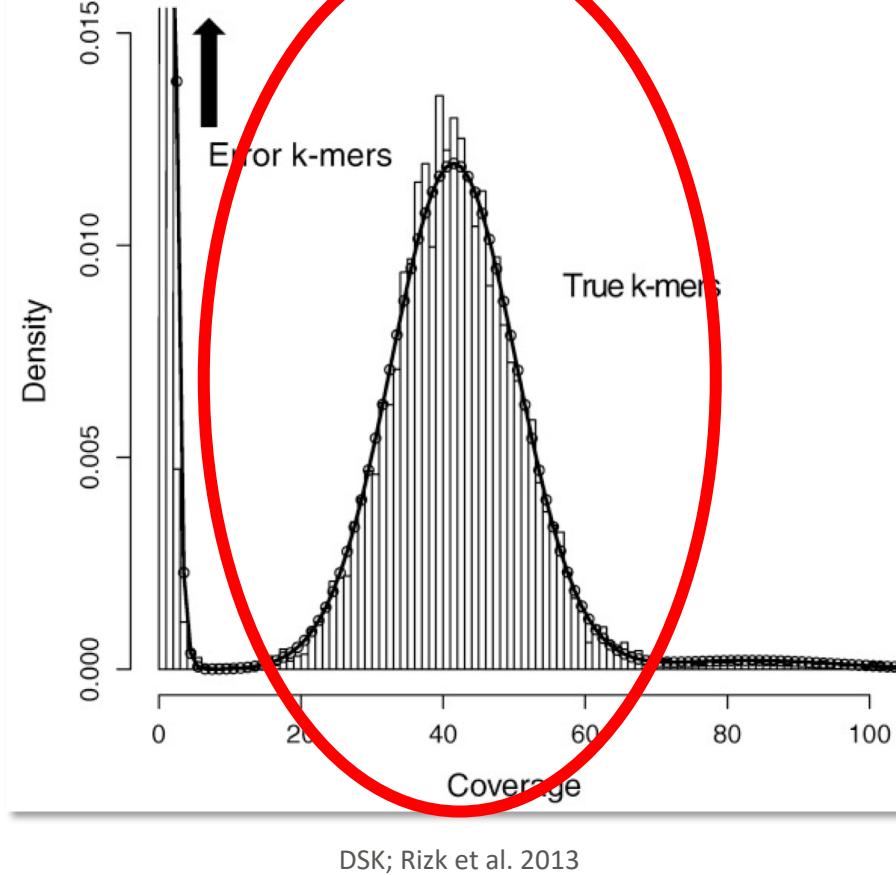
DSK; Rizk et al. 2013



# Error Correction: *K*-mer Counting

- Expected Distribution of *k*-mer frequency

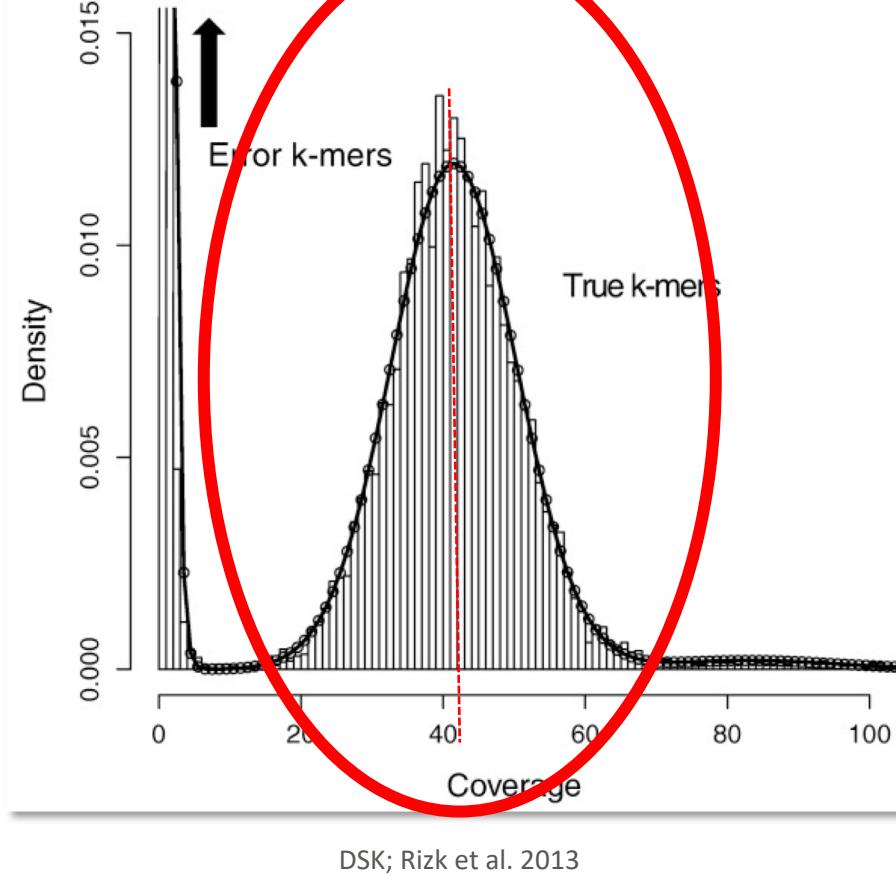
Estimate genome size



# Error Correction: *K*-mer Counting

- Expected Distribution of *k*-mer frequency

Estimate genome size

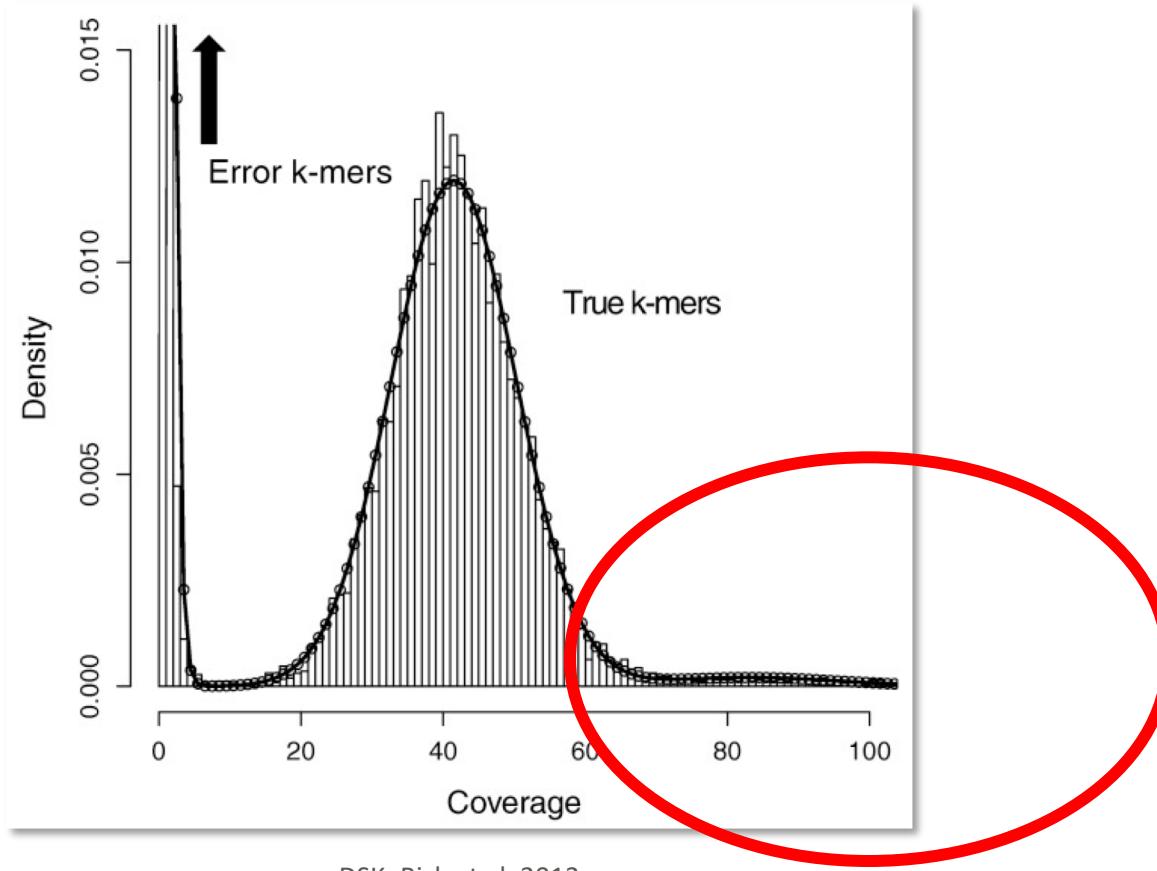


$G = C / P$   
G = genome size  
C = total count of true k-mers  
P = peak coverage

# Error Correction: *K*-mer Counting

- Expected Distribution of *k*-mer frequency

Estimate  
repetitive  
content



# Recap: NGS QC

- Remove low quality bases and reads
- Identify and remove adapter contamination
- Optional: Correct substitution sequencing errors
- Optional: De-duplication

# To Your Terminals!

