# FASTQ Processing Tools for Data Analysis

Zhongfeng Wang
TAS | 03/26/2020

illumına®

# Outline of the webinar

**Overview of Illumina FASTQ generation**

- Data analysis pipeline overview
- FASTQ format
- Demultiplexing and FASTQ generation

**FASTQ tools**

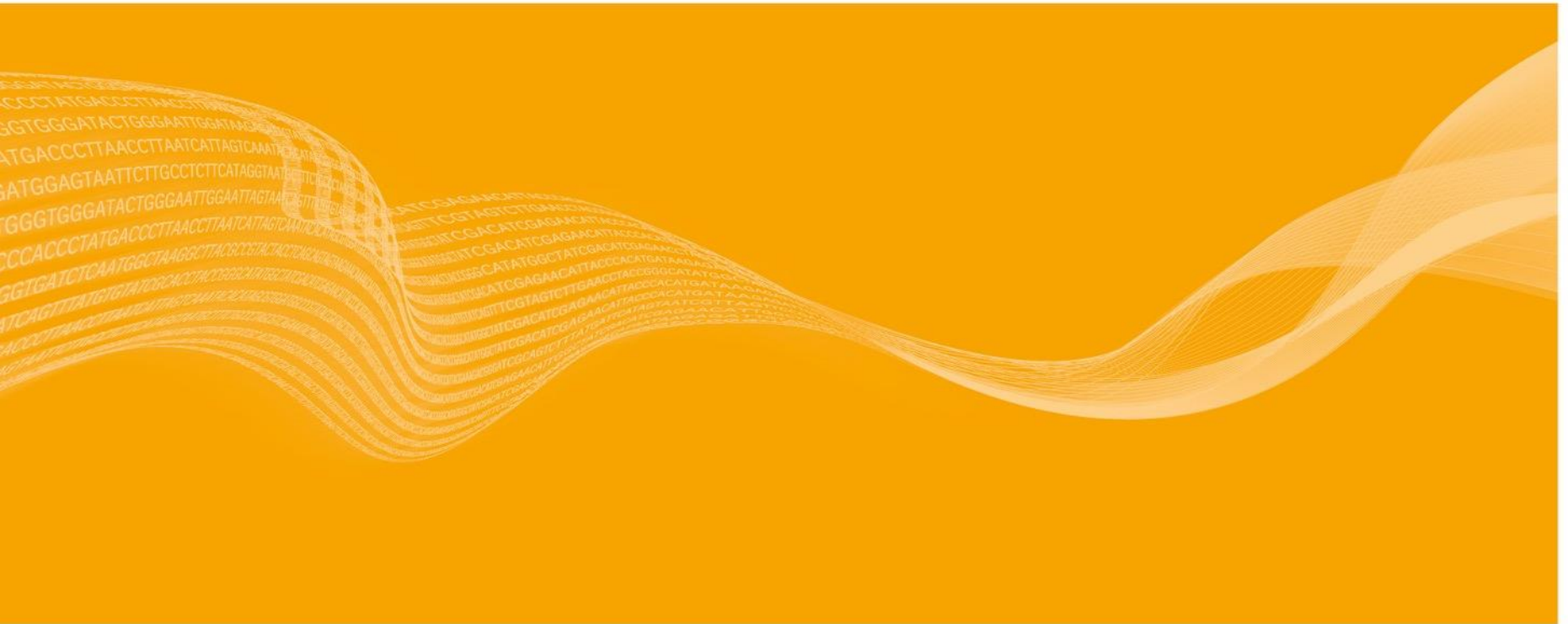- Adapter trimming
- Quality trimming
- Read merging

**Does my run look good vs does my FASTQ look good?**

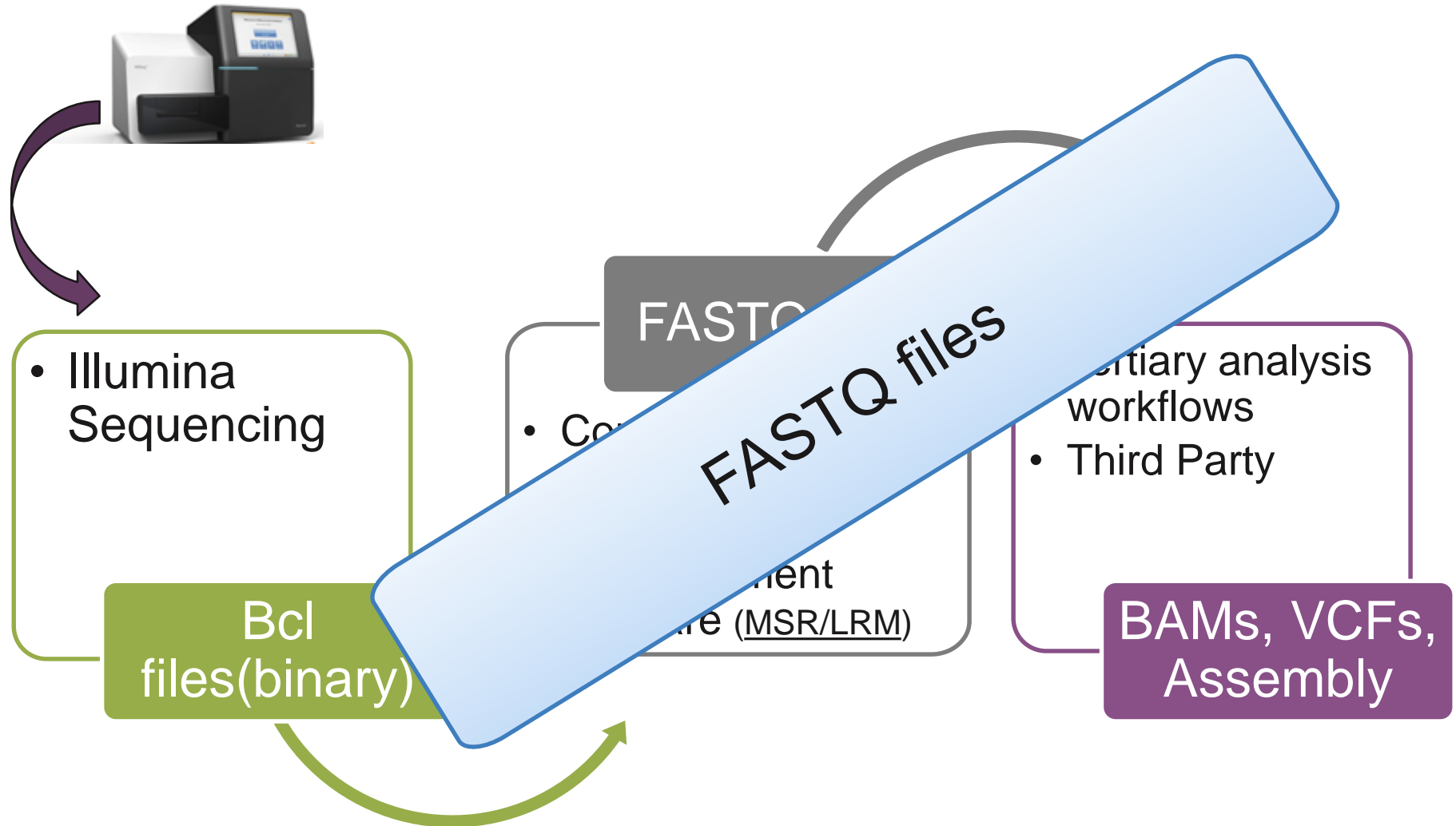**Illumina metrics vs third party QC tools**

illumına®

# Intended audience

- **Users who are new to Illumina sequencing platforms**
  - End users who use a core facility to sequence their data and need some guidance on data optimization

- **Not intended for bioinformaticians**

illumına®

# Overview of Illumina FASTQ Generation

# Pipeline for data analysis



- Illumina Sequencing

Bcl files(binary)

FASTQ

- Co...

...ent ...re (MSR/LRM)

FASTQ files

...rtiary analysis workflows
- Third Party

BAMs, VCFs, Assembly

illumina®

# BCL Files

- **BCL = Base CalL**

  - The "raw" data
  - Binary file containing base calls and quality scores for each tile for each cycle
  - Produced by on-instrument Real-Time Analysis Software (RTA)

illumına®

# FASTQ Files-format

- **Text file that contains the sequence data from clusters**
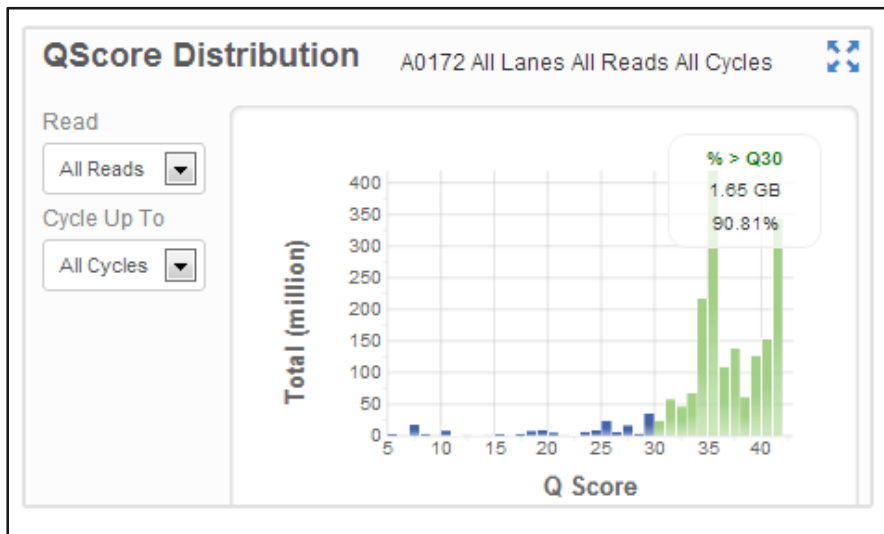
- **Each entry consists of 4 lines:**

1. Sequence identifier with information about the sequencing run and the cluster.
   i. Header format = @<instrument>:<run number>:<flowcell ID>:<lane>:<tile>:<xpos>:<y-pos>:<UMI> <read>:<is filtered>:<control number>:<index>
2. The sequence (the base calls; A, C, T, G and N).
3. A separator, which is simply a plus (+) sign.
4. The base call quality scores.  Phred +33 encoded (ASCII)

```
@NB987655:1:ABCD12345:1:11401:6329:1045 1:N:0:ATCACG
TCGCACTCAACGCCCTGCATATGACAAGACAGAATC
+
<>;##=><9=AAAAAAAAA9#:<#<;<<<????#=
```

# Qscore – Quality Score

- The QScore is based on the Phred scale – Chance that a base call is incorrectly called

  - Q10: 10% chance of incorrect base call
  - Q20: 1% chance of incorrect base call
  - Q30: 0.1% chance of incorrect base call
  - Q40: 0.01% chance of incorrect base call

- **Quality Scores for Next-Generation Sequencing**

**QScore Distribution**  A0172 All Lanes All Reads All Cycles

Read
All Reads

Cycle Up To
All Cycles

% > Q30
1.65 GB
90.81%

Total (million): 400, 350, 300, 250, 200, 150, 100, 50, 0
Q Score: 5, 10, 15, 20, 25, 30, 35, 40

| Symbol | ASCII Code | Q-Score |
|--------|-----------|---------|
| < | 60 | 27 |
| = | 61 | 28 |
| > | 62 | 29 |
| ? | 63 | 30 |
| @ | 64 | 31 |
| A | 65 | 32 |
| B | 66 | 33 |
| C | 67 | 34 |
| D | 68 | 35 |
| E | 69 | 36 |
| F | 70 | 37 |

illumina®

# FASTQ Files- How to view

FASTQ files are large (file size) – best to not open these

**How to view contents without opening?**
- Use a terminal window on Linux or Mac
- zcat *.fastq.gz | less
- zcat *.fastq.gz | head

**How to open a FASTQ file?**
- gzipped
- 7-zip for windows to unzip
- Open using a notepad

```
@NB00000:84:HN8L6AFXY:4:11401:14537:1100 1:N:0:CGTGTAGG+AGTCCAAC
GTAAACGGCGCGCCCCATGAGANCCCGGNTTGCTTTCCCAAGCCTTCGGGCGTCTGTGTGCGCTCTGTGGATGCCAGGGCCGACCAGAGGAGCCTTTTTAAAACACATGTTTTTTATACAA
+
AAAAAEEEEEEEEEEEEAEEEAA#/EEEE#EEEE/EEEEEAEEEEAEEEEEEEAAEEEEEEEEEAAEEEEEEEEEEEEEAEEEEEE/EAEEEEE/EEE<EEEEEEE/EEEEEEEEEEEEEAEEEEE
```

illumina®

# Anatomy of a FASTQ file

- **Naming**
**SampleName_SampleNumber_Lane_Read_FlowCellIndex.fastq.gz**
**SampleName_S1_L001_R1_001.fastq.gz**
**SampleName_S1_L001_R2_001.fastq.gz**

- **e.g. NextSeq500/550, PE sequencing. 8 fastq files for each sample**

SampleName_S1_L001_R1_001.fastq.gz

SampleName_S1_L002_R1_001.fastq.gz

SampleName_S1_L003_R1_001.fastq.gz

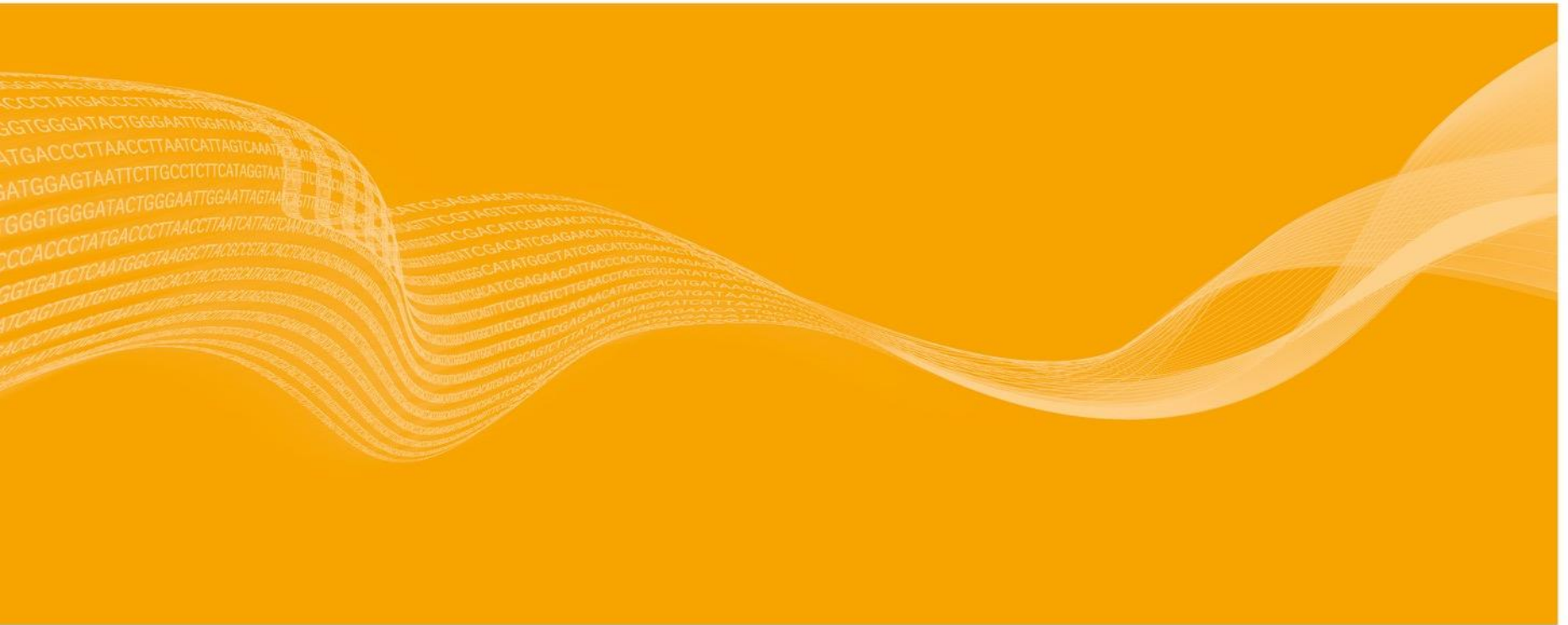SampleName_S1_L004_R1_001.fastq.gz

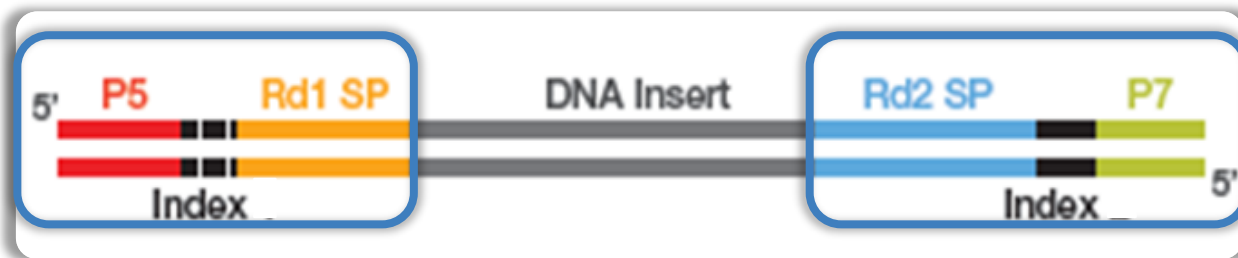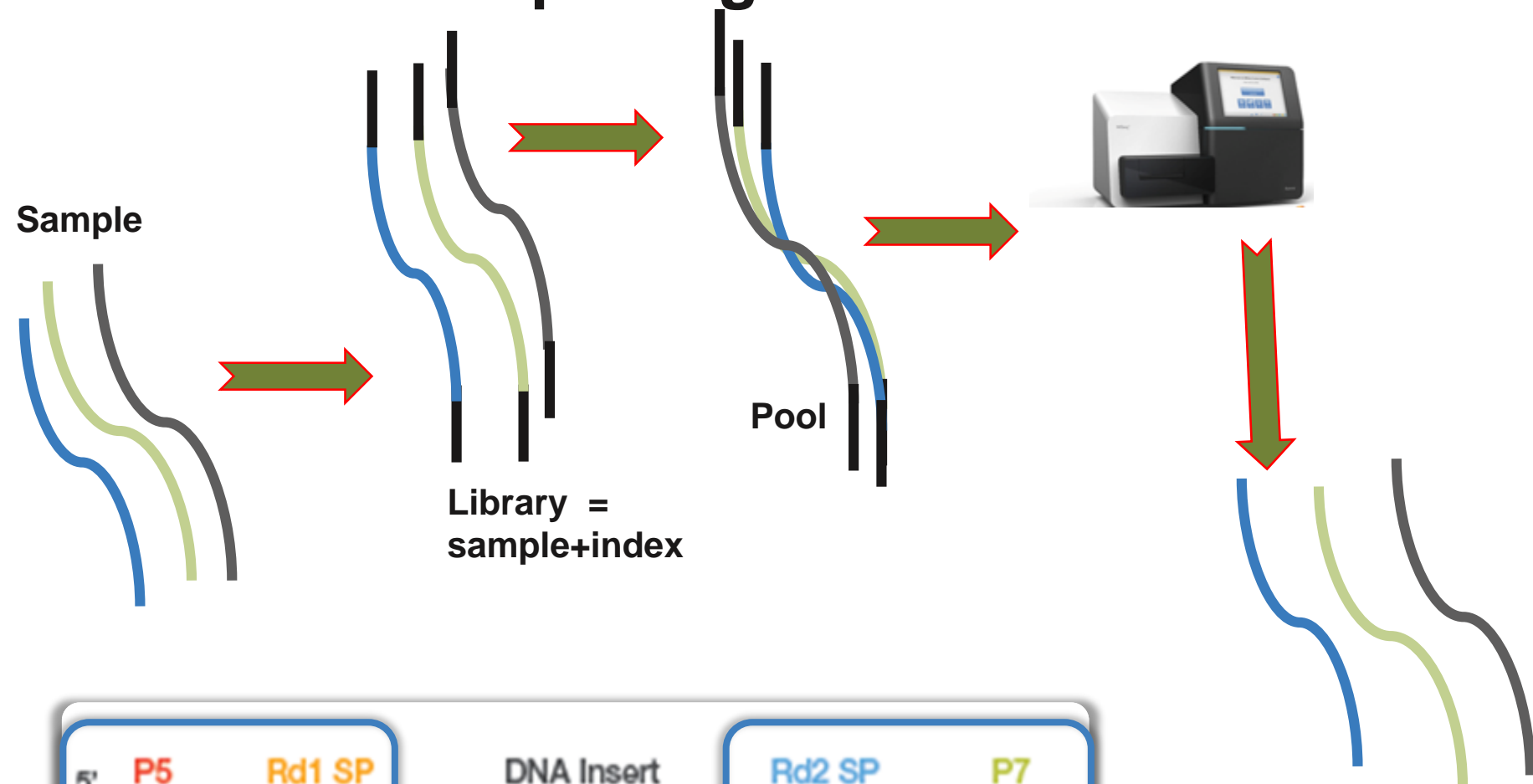SampleName_S1_L001_R2_001.fastq.gz

SampleName_S1_L002_R2_001.fastq.gz

SampleName_S1_L003_R2_001.fastq.gz

SampleName_S1_L004_R2_001.fastq.gz

illumina®

# Demultiplexing and FASTQ generation

illumina®

# What is demultiplexing?



**Sample**

**Library = sample+index**

**Pool**

P5   Rd1 SP   DNA Insert   Rd2 SP   P7
5'
Index   Index
5'

Separation of reads based on index sequences (or barcodes) using bioinformatics

illumına®

# Illumina Tools for Demultiplexing and Fastq Generation

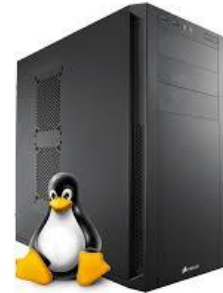**On-Board**

MiSeq™

MiniSeq™

iSeq™

NextSeq™

**MiSeq™ Reporter**

**Local Run Manager**

**Cloud based**

BaseSpace®

**BaseSpace™ Sequence Hub (BSSH)**

**Linux based**

**Bcl2fastq**

**DRAGEN ™ bclconvert**

For Research Use Only.  Not for use in diagnostic procedures.

illumina®

# Supports All Illumina Sequencing Systems



iSeq™  MiniSeq™  MiSeq™  NextSeq™

HiSeq™  NovaSeq™

**Bcl2fastq v2 is required for platforms using RTA 2 & 3.**

**Note:** For sequencing systems running an earlier software version of Real-Time Analysis (RTA) than v1.18.54, use bcl2fastq v1.8.4.

illumina®

# What is bcl2fastq v2?

- **Linux-based software to convert BCLs → FASTQs**

- **Optionally performs demultiplexing per sample sheet (created with IEM) instructions**

- **Command line example**

```
/usr/local/bin/bcl2fastq --runfolder-dir . --output-directory
./Data/Intensities/BaseCalls/  --input-dir
./Data/Intesities/Basecalls  --sample-sheet ./SampleSheet.csv
```

- **Command line configurable options**

  *E.g.* `--no-lane-splitting`
  - SampleName_S1_R1_001.fastq.gz
  - SampleName_S1_R1_001.fastq.gz

For Research Use Only.  Not for use in diagnostic procedures.

illumına®

# Demultiplexing reports



Indexing QC

- **IndexQC on SAV**

- **DemultiplexSummaryF1L#.txt**
  - Helps troubleshoot sample dropouts

- **Depending on what software is used for demultiplexing, DemultiplexSummaryF1L1.txt location will vary**

| | A | B | C | D |
|---|---|---|---|---|
| 1 | SampleNu | 0 | 1 | 2 |
| 2 | SampleNa | None | HG200-10n | HG200-10 |
| 3 | L1T1101 | 3.439879 | 1.156535 | 2.958111 |
| 4 | L1T1102 | 3.498839 | 1.135605 | 2.999131 |
| 5 | L1T1103 | 3.558643 | 1.126721 | 2.997476 |
| 6 | L1T1104 | 3.554299 | 1.141691 | 2.957287 |
| 7 | L1T1105 | 3.577636 | 1.125214 | 2.991321 |
| 8 | L1T1106 | 3.640834 | 1.111034 | 2.965273 |
| 9 | L1T1107 | 3.696693 | 1.126197 | 2.936993 |
| 10 | L1T1108 | 3.626516 | 1.119271 | 2.935823 |
| 11 | L1T1109 | 3.610867 | 1.134466 | 2.919896 |
| 12 | L1T1110 | 3.524556 | 1.129786 | 2.87885 |
| 13 | L1T1111 | 3.600484 | 1.106575 | 2.877 |

| ### Most Popular Unknown Index Sequences | |
|---|---|
| ### Columns: Index_Sequence Hit_Count | |
| ACTAAGAT+TGGTGAGT | 11440 |
| ACAGGCGC+TGGTGAGT | 10660 |
| GTGAATAT+CAACAGAT | 8220 |
| TCTCTACT+CAACAGAT | 6580 |
| ACTAGATA+TTGGTGAG | 5540 |
| TGCGAGAC+ATAGCGTC | 5420 |
| GCAGATTA+GGTTATAA | 5260 |
| ACAGGCGC+ATAGCGTC | 5180 |

illumina®

# Tools: Adapter Trimming

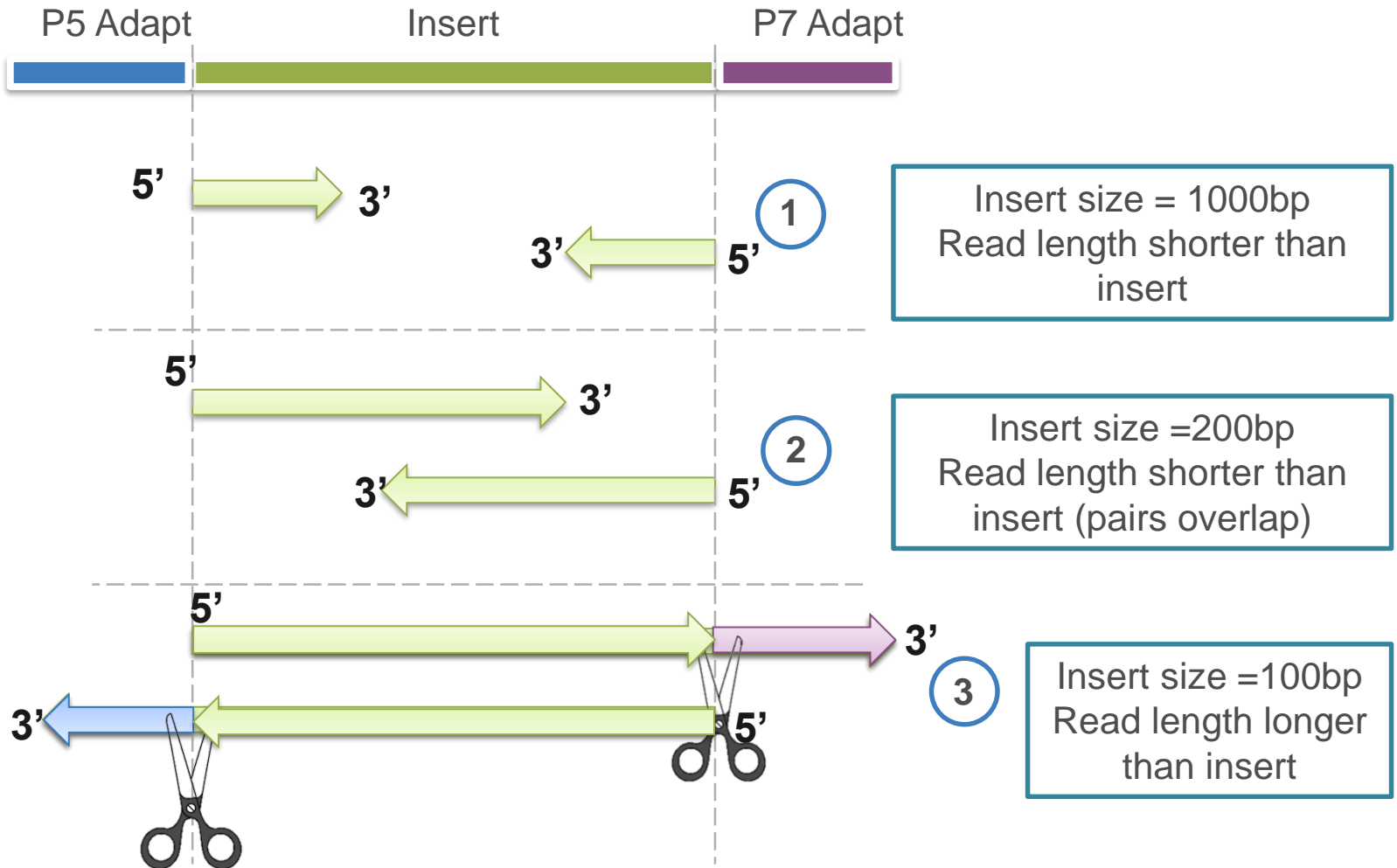# Illumina library design

○ **Library – DNA insert plus full adapter**

○ **Read 1 Sequencing Primer (Rd1 SP)**

○ **Read 2 Sequencing Primer (Rd2 SP)**

# Adapter trimming overview PE 151

# Adapter trimming options
## *MiSeq™ Reporter, Local Run Manager, BaseSpace™ Fastq Generation and Bcl2fastq*

| [Header] | |
|---|---|
| IEMFileVe | 4 |
| Date | 4/11/2017 |
| Workflow | GenerateFASTQ |
| Applicatic | FASTQ Only |
| Assay | Nextera XT |
| Description | |
| Chemistry | Amplicon |

[Reads]
151
151

[Settings]
Adapter    CTGTCTCTTATACACATCT

[Data]

| Sample_II | Sample_Nam | Sample_P | Sample_W | I7_Index_ | index | I5_Index_ | index2 | Sample_P | Description |
|---|---|---|---|---|---|---|---|---|---|
| Test | | | | N701 | TAAGGCG | S502 | CTCTCTAT | | |

**Hard trimming**
(Remove adapter completely from sequence)

If adapter sequence for Read1 and Read2 are different

[settings]
Adapter,.......
AdapterRead2,......

If runs are setup using PrepTab for NextSeq / MiniSeq, adapter trimming is on by default.

https://support.Illumina.com/bulletins/2016/12/what-sequences-do-i-use-for-adapter-trimming.html

21

For Research Use Only.  Not for use in diagnostic procedures.

# Adapter trimming example (hard trimming)

**Adapter match > 90% (default)**

```
@M00000:71:000000000-D00LW:1:1101:16265:1658 1:N:0:1
ACTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCCCTGTCTCTTATACACATCTCCGAGCCCA
+
BCCCCFFCCBCCGGGGGGGGGGGGGGGGHHHHHHHHHHHHGHHHHHHHHHHHHHHHHHHHGGGGGGH
```

```
@M00000:71:000000000-D00LW:1:1101:16265:1658 1:N:0:1
ACTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCC
+
BCCCCFFCCBCCGGGGGGGGGGGGGGGGHHHHHHHHHHHHG
```

Bases trimmed (removed) from the start of the Adapter

22

illumına®

# Adapter masking options
*Bcl2fastq*



| [Header] | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| IEMFileVersio | 4 | | | | | | | | |
| Date | 4/11/2017 | | | | | | | | |
| Workflow | GenerateFASTQ | | | | | | | | |
| Application | FASTQ Only | | | | | | | | |
| Assay | Nextera XT | | | | | | | | |
| Description | | | | | | | | | |
| Chemistry | Amplicon | | | | | | | | |
| | | | | | | | | | |
| [Reads] | | | | | | | | | |
| 151 | | | | | | | | | |
| 151 | | | | | | | | | |
| | | | | | | | | | |
| [Settings] | | | | | | | | | |
| MaskAdapter | CTGTCTCTTATACACATCT | | | | | | | | |
| | | | | | | | | | |
| [Data] | | | | | | | | | |
| Sample_ID | Sample_Nam | Sample_P | Sample_W | I7_Index_ | index | I5_Index_ | index2 | Sample_P | Description |
| Test | | | | N701 | TAAGGCG | S502 | CTCTCTAT | | |

**N-masking**
(Replace adapter with "N" and quality "#")
If adapter sequence for Read1 and Read2 are different

[settings]
MaskAdapter,.....
MaskAdapterRead2,.....

illumına®

# Adapter masking example

```
@M00000:71:000000000-D00LW:1:1101:16265:1658 1:N:0:1
ACTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCCCTGTCTCTTATACACATCTCCGAGCCCA
+
BCCCCFFCCBCCGGGGGGGGGGGGGGGGHHHHHHHHHHHHHGGHHHHHHHHHHHHHHHHHHHHGGGGGGH
```

```
@M00000:72:000000000-D00LW:1:1101:16265:1658 1:N:0:1
ACTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCCNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+
BCCCCFFCCBCCGGGGGGGGGGGGGGGGHHHHHHHHHHHHHG#########################
```

Bases masked to "N" and quality score to 2 "#" from the start of the adapter
(Useful if analysis program requires reads to be same read length)

# Masking short adapters – bcl2fastq

```
@M00000:71:000000000-D00LW:1:1101:16265:1658 1:N:0:1
ACCTGTCTCTTATACACATCTCCGAGTCTGCGTTGCGCTTCTGCTCGGCCTCCAGCTCACCCTCCCCCA
+
BCCCCFFCCBCCGGGGGGGGGGGGGGGGHHHHHHHHHHHHHGGHHHHHHHHHHHHHHHHHHGGGGGGH
```



```
@M00000:72:000000000-D00LW:1:1101:16265:1658 1:N:0:1
NNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNNN
+
###################################
```

--minimum-trimmed-read-length 35
--mask-short-adapter-reads 22

Entire read is converted to 35bp Ns and quality score to "#" from the start of the sequence
- Useful if adapter dimers in the sample
- Prevents generation of empty reads

# Why adapter trim?

**(1)** Higher alignment %

**BWA**
(backtrace)

| Sample | Sample Name | Total Aligned Reads | Percent Aligned Reads |
|---|---|---|---|
| 1 | NA12892 | 354,882 | 77.4% |
| 2 | NA12892-trim | 450,007 | 98.2% |

illumına®

# Why adapter trim?

**2** Improved assemblies

Data: 2 x 250bp, *E.coli* (Nextera™ XT)

Velvet de novo Assembly
BASESPACE LABS

| Assembly metrics | Before adapter trimming | After adapter trimming |
|---|---|---|
| N50 | 21 | 29,791 |
| Maximum contig | 553 | 174,326 |
| Assembly length | 18,497,207 | 4,876,437 |
| Number of contigs | 1,387,508 | 1,115 |

illumına®

# Tools: Quality trimming

illumina®

# Quality trimming

- **Filter the end of reads based on read end average quality**

- **When to trim?**
  - Where algorithms sensitive to quality – De Novo assembly, merging reads, Metagenomics (using 16S rRNA gene for classification)

- **When not to trim?**
  - Resequencing. Most aligners take quality scores into account (i.e. BWA, Isaac) will soft clip the ends of reads if low quality

- **FASTQ toolkit app on BSSH (BaseSpace Sequence Hub) can be used to perform quality trimming**

illumına®

# Quality score trimming example
## *QualityScoreTrim,20*

```
@M00000:72:000000000-D00LW:1:1101:22420:18334 1:N:0:1
CACCAAGGGCCTGGGGTGTCAATGGCGGGGCTTGTGACTGCACAAAAGGGGCCTCCCGCAGGGGCTCCCGCC
+
BBBBBBFBBBBBGGGGEEFGGGHHHHGGG00>A0B355@BB3@3BGB?E1///1/11///////////?////
```

```
@M00000:72:000000000-D00LW:1:1101:22420:18334 1:N:0:1
CACCAAGGGCCTGGGGTGTCAATGGCGGGGCTTGTGACTGCACAAAAGG
+
BBBBBBFBBBBBGGGGEEFGGGHHHHGGG00>A0B355@BB3@3BGB?E
```

| Q | ASC |
|---|-----|
| 13 | . |
| 14 | / |
| 15 | 0 |
| 16 | 1 |
| 18 | 3 |
| 20 | 5 |
| 22 | 7 |
| 25 | 9 |
| 30 | ? |
| 31 | @ |
| 32 | A |
| 33 | B |

illumına®

# Tools: Read merging

illumına®

# Read merging

- **Combines paired end read into a single read**

- **Generally requires substantially high overlap**
  - StitchReads,1 option in MiSeq™ Reporter
  - Workflows = Amplicon – DS, GenerateFASTQ workflow, TruSeq™ Amplicon
  - 10 bp overlap
  - Not supported in any other software

- **When to merge?**
  - Where read continuity is important
  - Improve indel detection in some cases (i.e. where Indel is in overlap region)
  - To use tools that only take single end reads (for example some metagenomics software)
  - Where majority of reads overlap

- **When not to merge?**
  - When some proportion of reads do not overlap
  - Where there are simple repeats in overlap region (i.e. sequencing amplicons with expansion repeats)
  - Most downstream aligners will handle unmerged reads

For Research Use Only.  Not for use in diagnostic procedures.

illumina®

# Read merging



| | | |
|---|---|---|
| 5' Adapt | Insert | 3' Adapt |

1 — No merging

2 — Can be merged

3 — Merge and remove adapters

illumına®

# Data optimization with an example: Using FASTQ toolkit app for Small RNA workflow

illumına®

# BaseSpace™

- **BaseSpace™ is the Illumina <span style="color:red">cloud-based</span> genomics computing environment for next-generation sequencing (NGS) data management and analysis.**

- **Is used to analyze or manipulate samples – results stored in the project. There are currently about 80+ apps**

| Core App | BaseSpace Lab App | | Third party App |
|---|---|---|---|
| Small RNA<br>Illumina, Inc. | FASTQ Toolkit<br>BaseSpace Labs | FastQC<br>BaseSpace Labs | Elastic Genome Browser<br>Strand Life Sciences |

**Please Read:**

This is a BaseSpace™ Labs App. BaseSpace Labs Apps are developed using an accelerated development process in order to make them available to BaseSpace users faster than conventional Illumina Apps. Illumina may provide support for BaseSpace Labs Apps at its sole discretion. BaseSpace Labs Apps are provided AS-IS without any warranty of any kind. BaseSpace Labs Apps are used at the user's sole risk. Illumina is not responsible for any loss of data, incorrect results, or any costs, liabilities, or damages that may result from the use of a BaseSpace Labs App.

For Research Use Only.  Not for use in diagnostic procedures.

illumina®

# Small RNA sample prep kit



**Small RNA**
Illumina, Inc.

- **Human small RNA peaks at ~22bp**

- **Typical run setup is Single End = 36 to 50cycles. So reads have to be adapter trimmed**

- **Bcl2fastq**:  --minimum-trimmed-read-length 20 and --mask-short-adapter-reads 20

- **Recommendation: For workflow GenerateFASTQ – do not turn on adapter trimming**

- **Generate FASTQs with 50bp and trim downstream using 3rd party tools or BSSH FASTQ toolkit**

For Research Use Only.  Not for use in diagnostic procedures.

illumına®

# FASTQ toolkit output- Results

# Small RNA app output



Untrimmed FASTQ

Trimmed FASTQ

# Does my run look good?
# vs
# Does my FASTQ look good?

Previous Webinar -Sequence Analysis Viewer (SAV): A Beginner's Guide

illumına®

# Metrics



| Illumina | 3rd party tools *Eg.,* FASTQC |
|---|---|
| • **Based on Sequencing Analysis Viewer – taking a whole run into consideration** | • **Based on individual FASTQ file** |
| • **Yield, Q30, Reads passing filter** | • **Sample dependent and hence library kit** |
| • **Used for assessing run quality – instrument performance** | • **Used for assessing sample performance – primarily for troubleshooting individual sample** |
| • **Illumina.com > Systems > MiSeq™ > Specifications** | |

# Does my FASTQ look good?

- **FASTQC app quick QC step to assess sample quality - subsamples from read pool**

- **FASTQC tool helps understand the nature of your library**

- **Not all graphs are meaningful - Understand documentation**

- **Failure of QC plots in the app does not translate to poor data quality – understand your data**

- **Sometimes it needs more in-depth downstream data analysis**

# What FastQC report looks like

**FastQCFastQC Report**

## Summary

- ✅ Basic Statistics
- ✅ Per base sequence quality
- ✅ Per sequence quality scores
- ❌ Per base sequence content
- ❌ Per base GC content
- ✅ Per sequence GC content
- ✅ Per base N content
- ⚠️ Sequence Length Distribution
- ❌ Sequence Duplication Levels
- ⚠️ Overrepresented sequences
- ⚠️ Kmer Content

**For Research Use Only.  Not for use in diagnostic procedures.**

illumina®

# Per Base Sequence Quality



Quality scores across all bases (Sanger / Illumina 1.9 encoding)

# Per Base GC Content
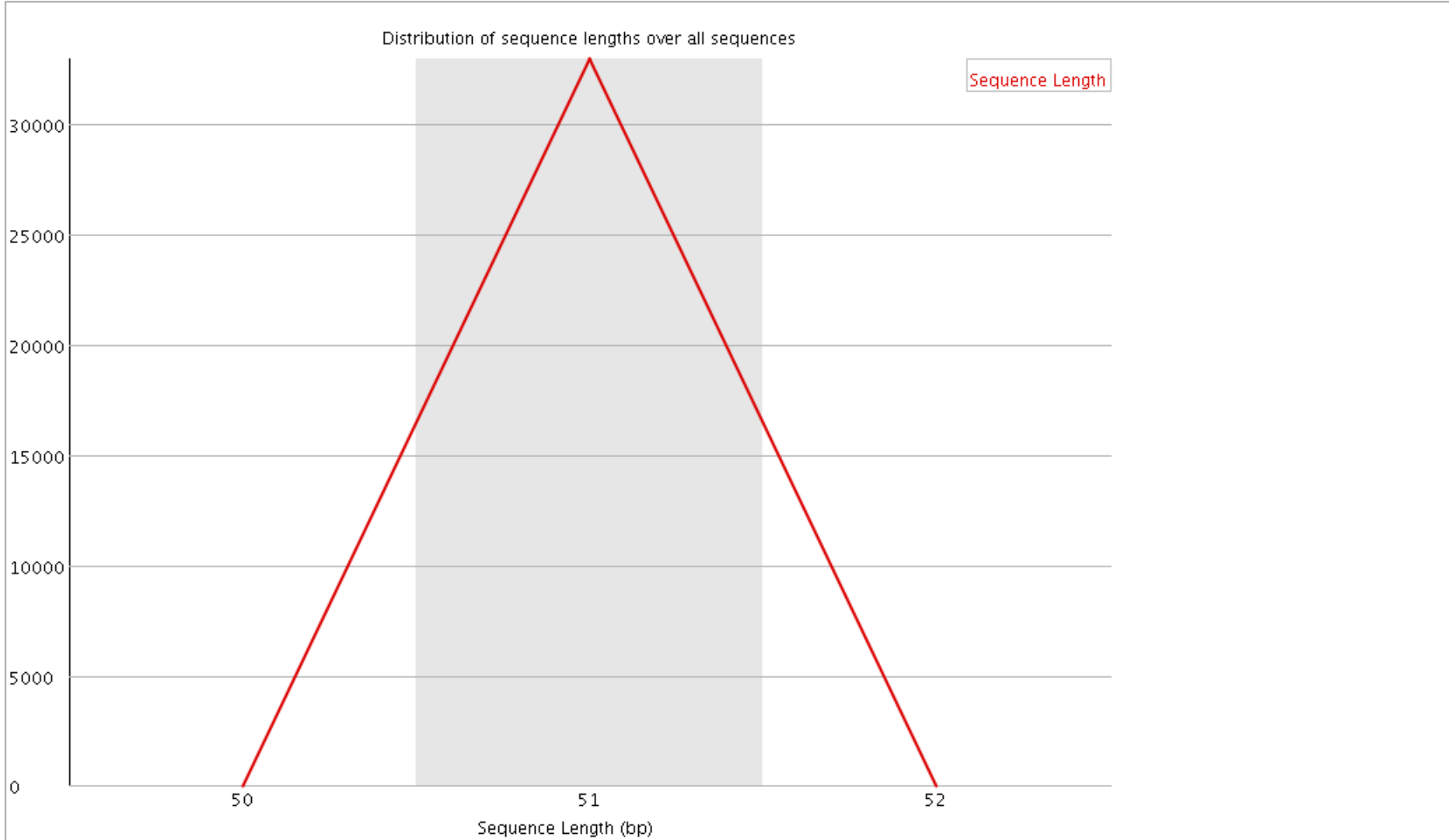
# Per Base Sequence Content

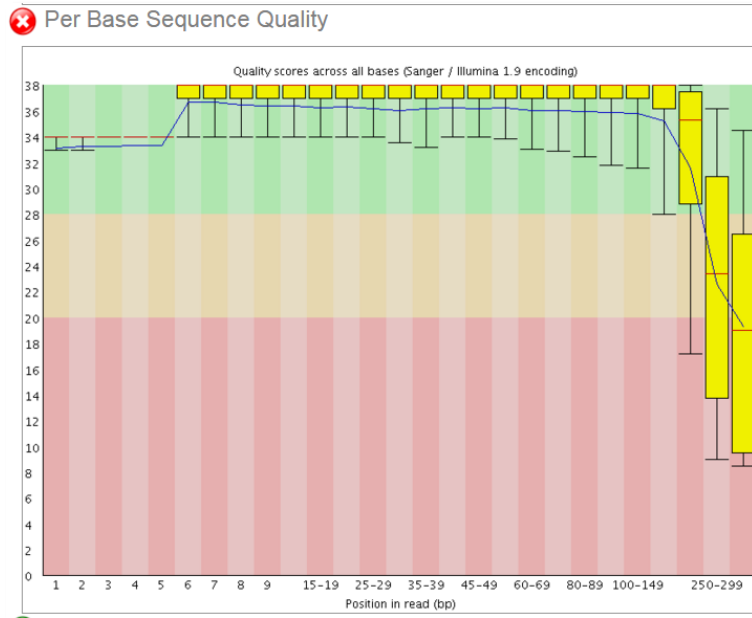# Per Sequence Quality Scores



Quality score distribution over all sequences

Average Quality per read

# Per Sequence GC Content



GC distribution over all sequences

Number of reads (y-axis)
Mean GC content (%) (x-axis)

GC count per read
Theoretical Distribution

illumina®

# ✅ Sequence Length Distribution



Distribution of sequence lengths over all sequences

Sequence Length

Sequence Length (bp)

illumına®

# Example of less than ideal data



▶ Indication of potential run problem, check SAV data

▶ Indication of mixed library types or potential library problem?

illumına®

# Example of less than ideal data



- Indication of poor sequence quality, using non-PF data?
- Using adapter masks, short reads?



- Indication of adapter presence
- adapter-dimer or adapter contamination

illumına®

# Summary

- **FASTQ format, demultiplexing and FASTQ generation using Illumina tools**

- **FASTQ processing tools for adapter trimming, quality trimming, read merging** *et.al.*

- **Illumina metrics are used for the overall run performance whereas third party tools such as FASTQC are used to look at sample performance**

illumına®

# Resources

- **FASTQC detailed documentation**

- What sequences do I use for adapter trimming?

- **Demultiplexing resource bulletin**

- FASTQ files explained

- Adapter trimming: Why are adapter sequences trimmed from only the 3' ends of reads?

- **Illumina在线技术培训研讨会 — Assessing Run Quality With SAV And FastQC**

- **Illumina在线技术培训研讨会 — Introduction To Bcl2fastq V2+**

illumına®

# QUESTIONS?