

Next Generation Sequencing

Short introduction

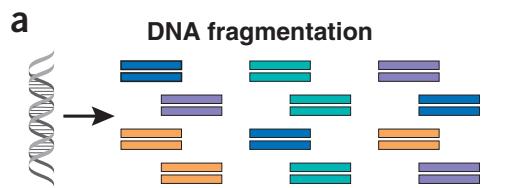
David Gomez-Cabrero

NGS: Short Introduction

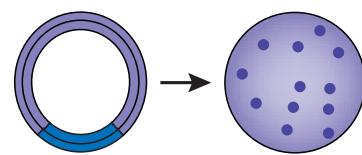
History

NGS

Bioinformatics

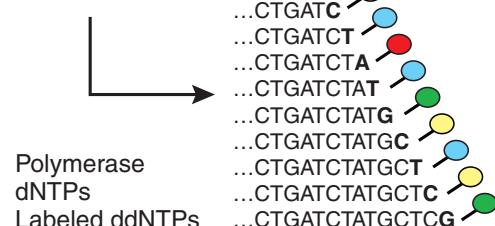


In vivo cloning and amplification



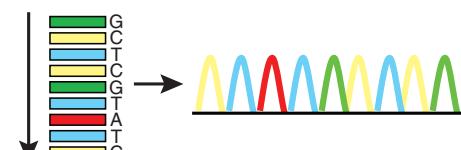
Cycle sequencing

3'-... GACTAGATACGAGCGTGA...-5' (template)
5'-... CTGAT ... (primer)



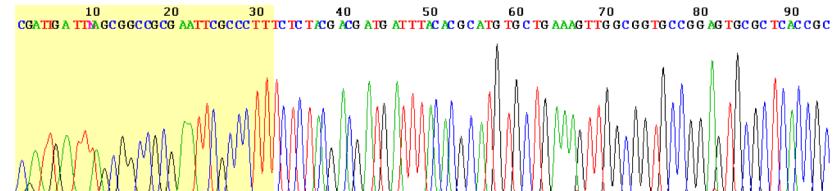
Polymerase
dNTPs
Labeled ddNTPs

Electrophoresis
(1 read/capillary)



Shendure and Ji,
2008, Nat. Biotech.

Sanger Sequencing



start of an example dye-terminator read

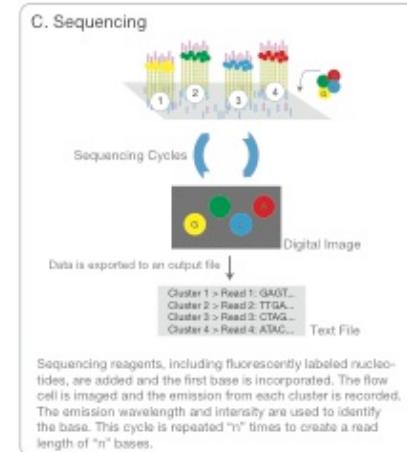
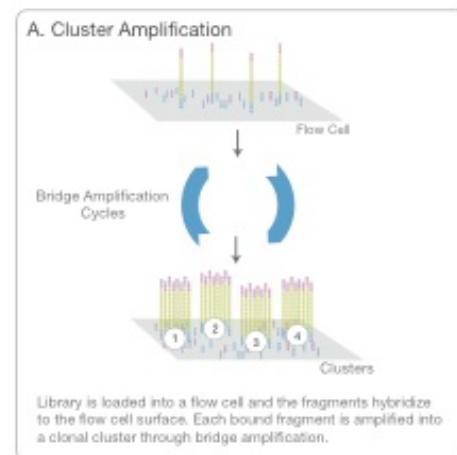
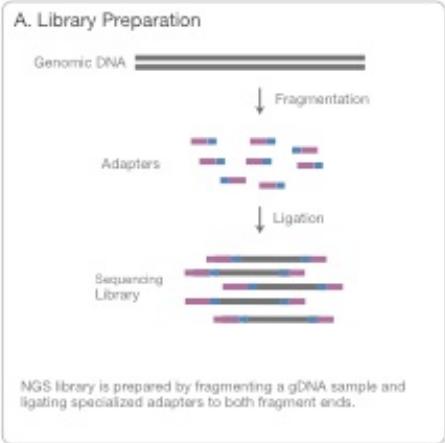
<https://en.wikipedia.org/wiki/User:Loris>

NGS: Short Introduction

History

NGS

Bioinformatics



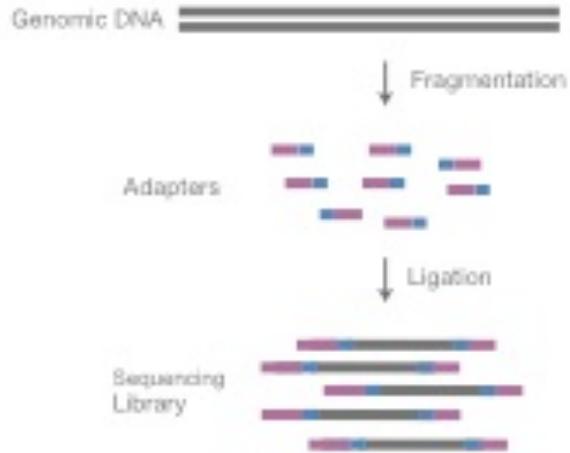
NGS: Short Introduction

History

NGS

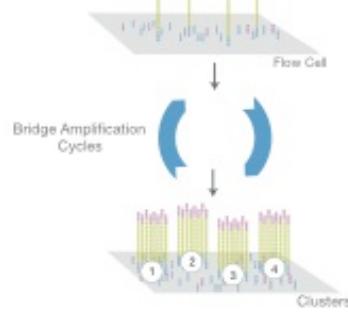
Bioinformatics

A. Library Preparation



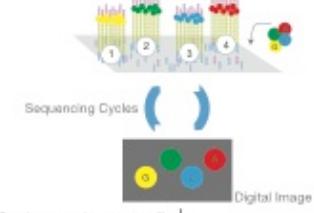
NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

A. Cluster Amplification



Library is loaded into a flow cell and the fragments hybridize to the flow cell surface. Each bound fragment is amplified into a clonal cluster through bridge amplification.

C. Sequencing



Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated "n" times to create a read length of "n" bases.

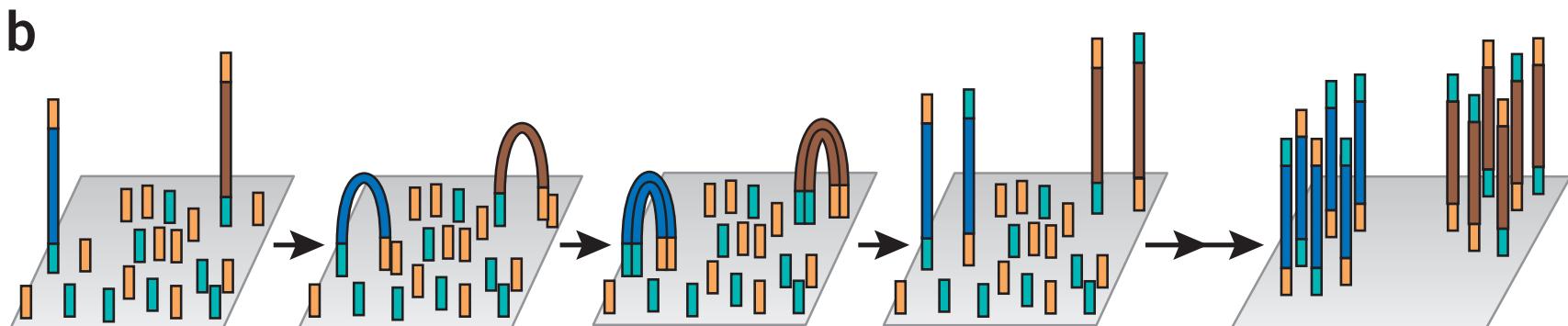
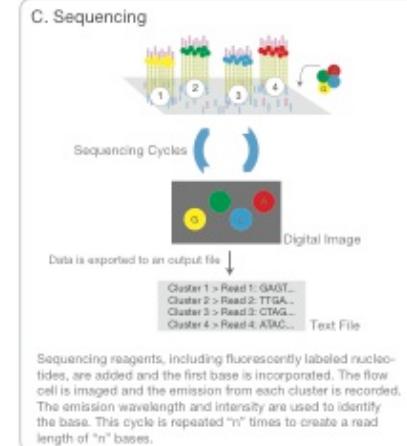
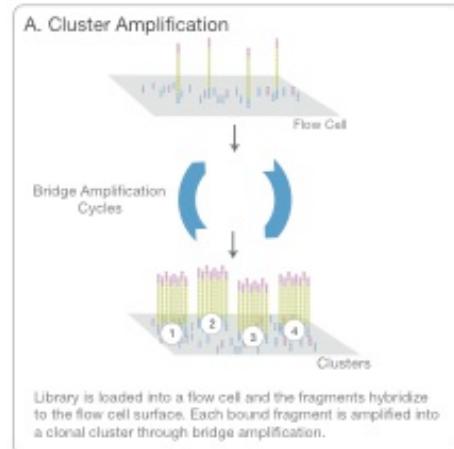
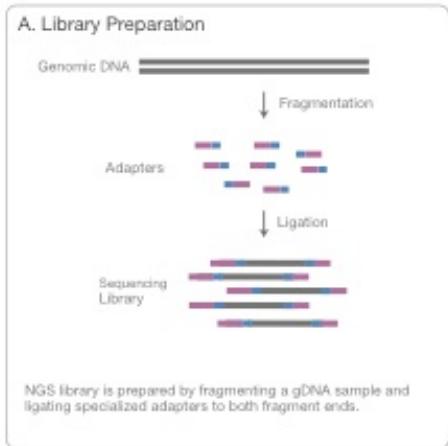


NGS: Short Introduction

History

NGS

Bioinformatics



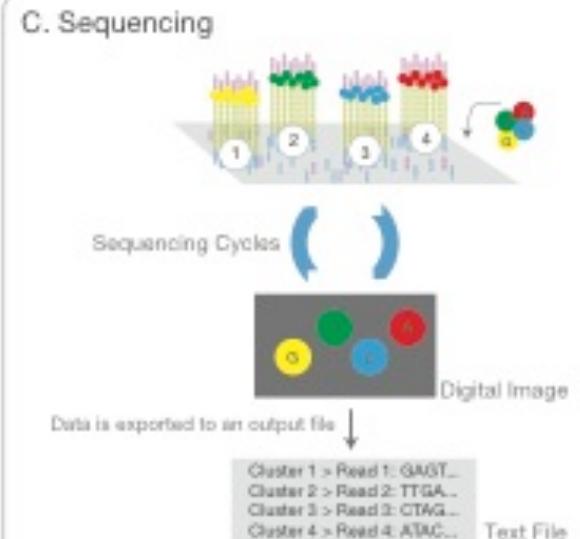
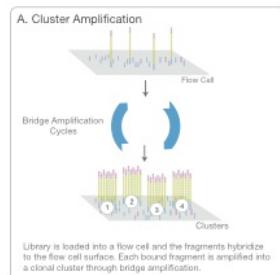
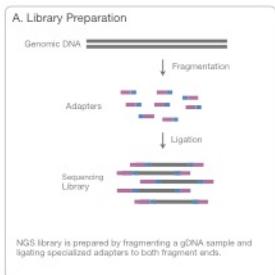
Shendure and Ji,
2008, Nat. Biotech.

NGS: Short Introduction

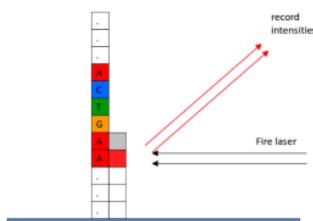
History

NGS

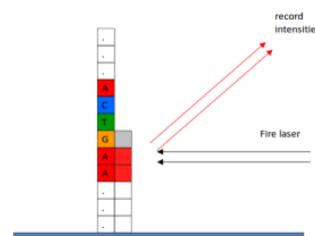
Bioinformatics



Sequencing cycle 1



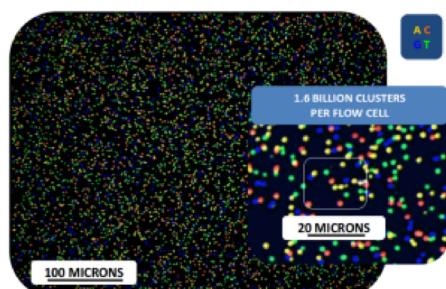
Sequencing cycle 2



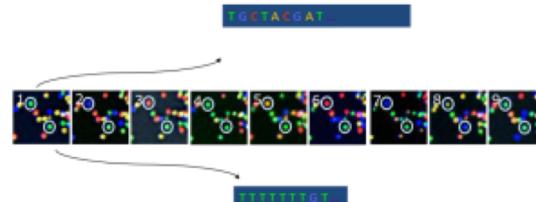
Sequencing reagents, including fluorescently labeled nucleotides, are added and the first base is incorporated. The flow cell is imaged and the emission from each cluster is recorded. The emission wavelength and intensity are used to identify the base. This cycle is repeated "n" times to create a read length of "n" bases.

AT EACH CLUSTER

GLOBAL IMAGE



Base calling from raw data



NGS: Short Introduction

History

NGS

Bioinformatics

MULTIPLEXING

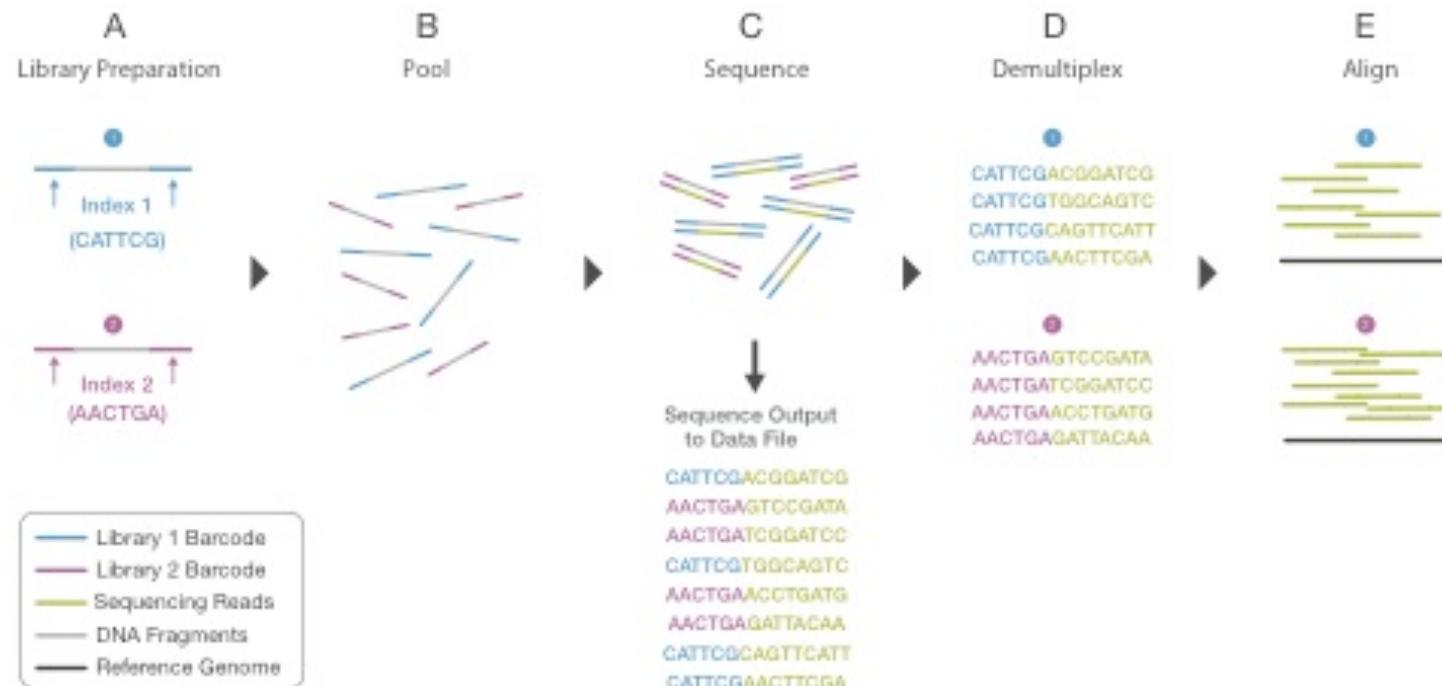


Figure 5: Library Multiplexing Overview.

- Two distinct libraries are attached to unique index sequences. Index sequences are attached during library preparation.
- Libraries are pooled together and loaded into the same flow cell lane.
- Libraries are sequenced together during a single instrument run. All sequences are exported to a single output file.
- A demultiplexing algorithm sorts the reads into different files according to their indexes.
- Each set of reads is aligned to the appropriate reference sequence.

NGS: Short Introduction

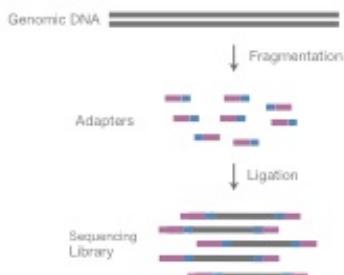
History

NGS

Bioinformatics

PAIR-END vs SINGLE-END

A. Library Preparation



NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

Paired-End Reads

Read 1



Read 2

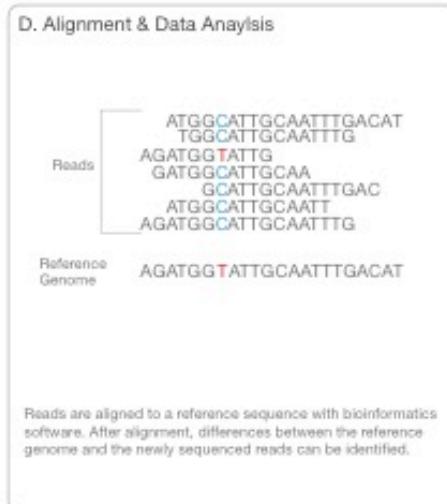
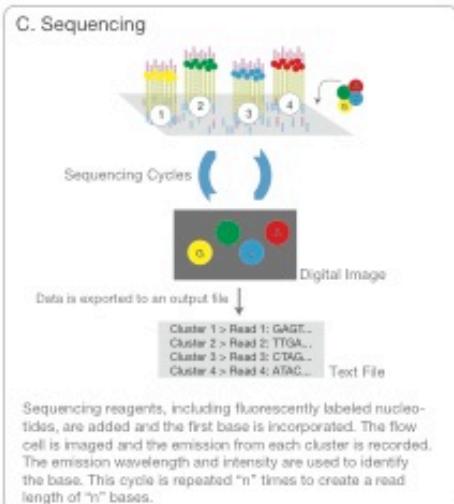
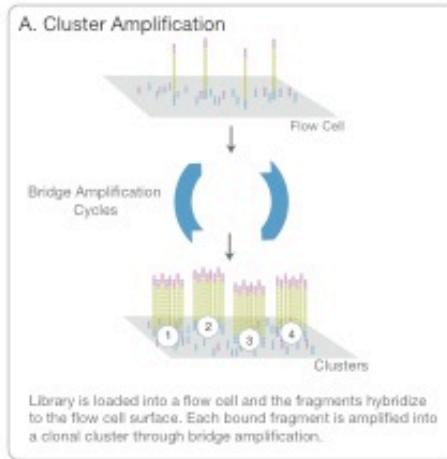
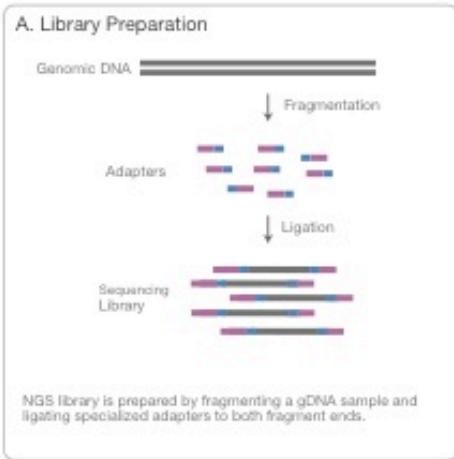


NGS: Short Introduction

History

NGS

Bioinformatics



*Where does
bioinformatics starts?*



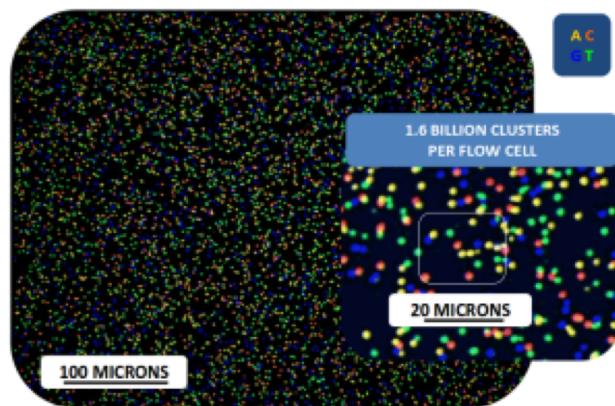
ANALYSIS

NGS: Short Introduction

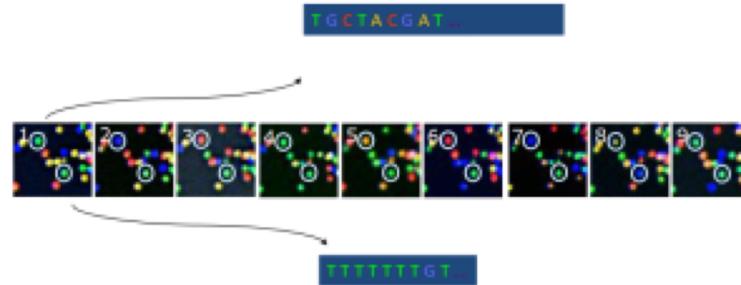
History

NGS

Bioinformatics



Base calling from raw data



ILLUMINA BASE CALL FILE: BCL

This part is hidden for the general user because the output provided by NGS facilities is directly sequences

bcl2fastq Conversion User Guide

Version 1.8.4

FOR RESEARCH USE ONLY

CASAVA

At Illumina, we are dedicated to empowering your research through the continuous development of the most advanced genetic analysis tools. This rapid innovation occasionally results in product discontinuation, as we develop newer, more cutting-edge solutions. We are discontinuing distribution of CASAVA software, and will continue its support until December 31, 2015. We remain committed to providing you with the highest-quality support and service during this transition.

As a replacement for CASAVA, we have added functionality to the BaseSpace Isaac Enrichment and Isaac Whole Genome Sequencing apps available on [BaseSpace](#) and [BaseSpace Onsite](#). In addition, the [bcl2fastq Conversion Software](#) is available to demultiplex and convert BCL files to FASTQ files on your local computer hardware. Illumina is committed to providing easy-to-operate, seamless solutions for the analysis of genomic data. BaseSpace and BaseSpace Onsite offer fully supported software solutions and are continuously optimized through ongoing development efforts.

NGS: Short Introduction

History

NGS

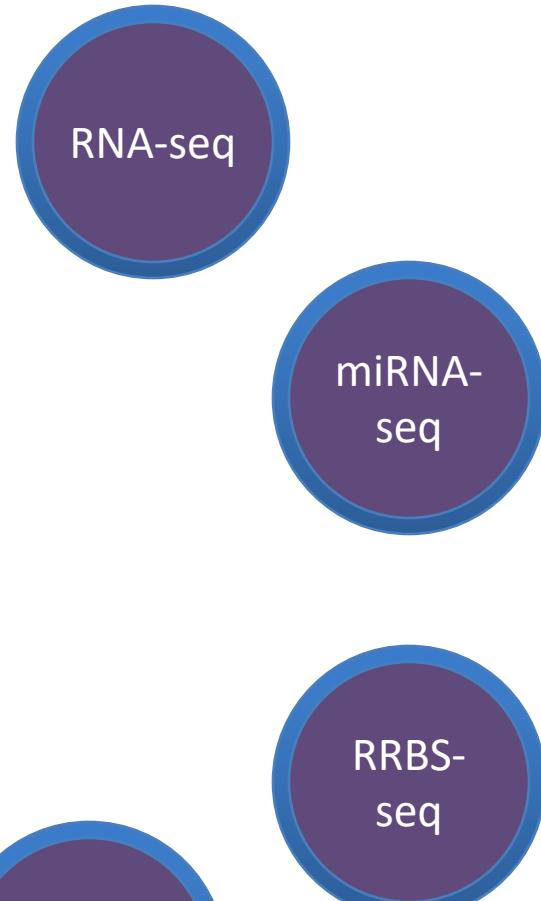
Bioinformatics

GENERAL PIPELINE?



DIFFERENT PIPELINES FOR EACH PROTOCOL

SHARED PARTS



NGS: Short Introduction

History

NGS

Bioinformatics

QUALITY
CONTROL

All-seq

BCL to FASTQ

How many clusters?

How many / Percentage fragments that have been identified?

...

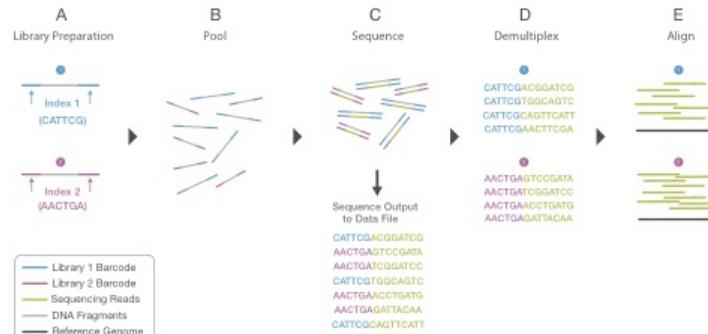


Figure 5: Library Multiplexing Overview.

- Two distinct libraries are attached to unique index sequences. Index sequences are attached during library preparation.
- Libraries are pooled together and loaded into the same flow cell lane.
- Libraries are sequenced together during a single instrument run. All sequences are exported to a single output file.
- A demultiplexing algorithm sorts the reads into different files according to their indexes.
- Each set of reads is aligned to the appropriate reference sequence.

NGS: Short Introduction

History

NGS

Bioinformatics

All-seq

BCL to FASTQ

QUALITY
CONTROL

FASTA FORMAT

> sequence: and information of it
ATCGATCGATCGATCGATCGATCG

Name, Details

Sequence

FASTQ FORMAT

@sequence: and information of it
GATTGGGGTTCAAAG
+
!**(((***+))%%

QUALITY VALUES
FOR NUCLEOTIDES

NGS: Short Introduction

History

NGS

Bioinformatics



FASTQ FORMAT

@sequence: and information of it

GATTGGGGTTCAAAG

+

!'''*((((***)+))%%

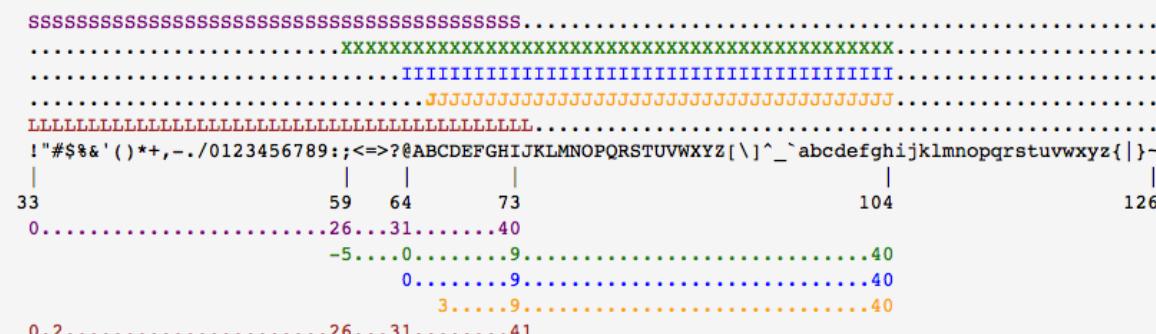
QUALITY VALUES FOR NUCLEOTIDES

p = probability that the corresponding base-call is incorrect

From this is computed a Phred quality score:

$$Q_{sanger} = -10 \log_{10}(p)$$

*This goes from 0 to 40-41,
the larger the better*



S - Sanger Phred+33, raw reads typically (0, 40)

S - Sanger Solexa+35, raw reads typically (-5, 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)

I = Illumina 1.8+ Phred+33 raw reads typically (0-41).

NGS: Short Introduction

History

NGS

Bioinformatics

FASTQ FORMAT

All-seq

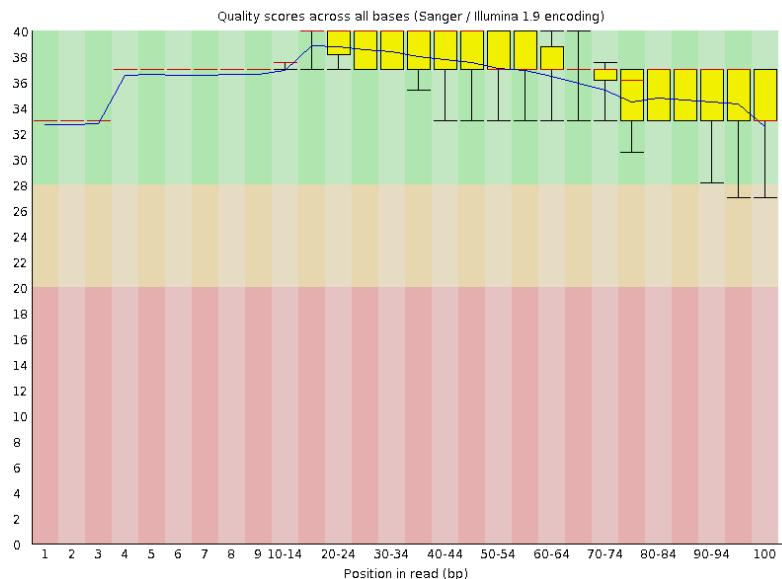
@sequence: and information of it

GATTGGGGTTCAAAG

+

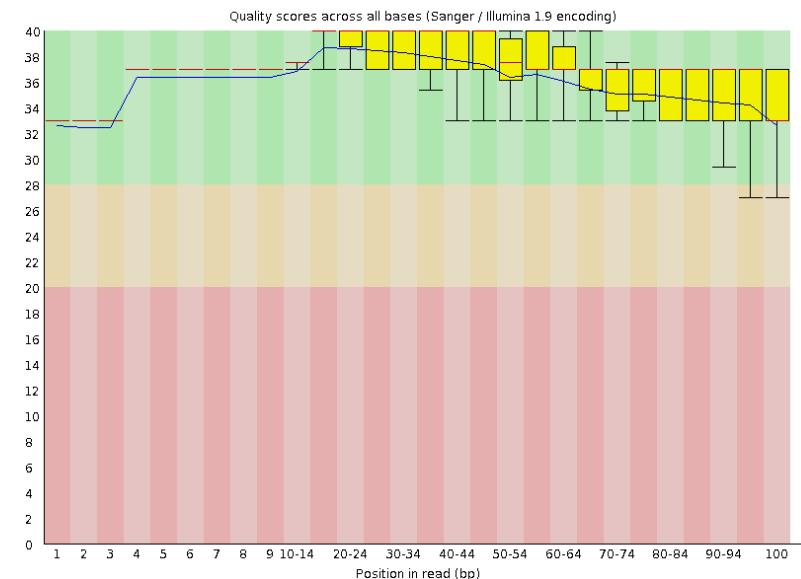
!''*(((****+))%%

$$Q_{sanger} = -10 \log_{10}(p)$$



FRAGMENT

READ 1



FRAGMENT

READ 2

NGS: Short Introduction

History

NGS

Bioinformatics

FASTQ FORMAT

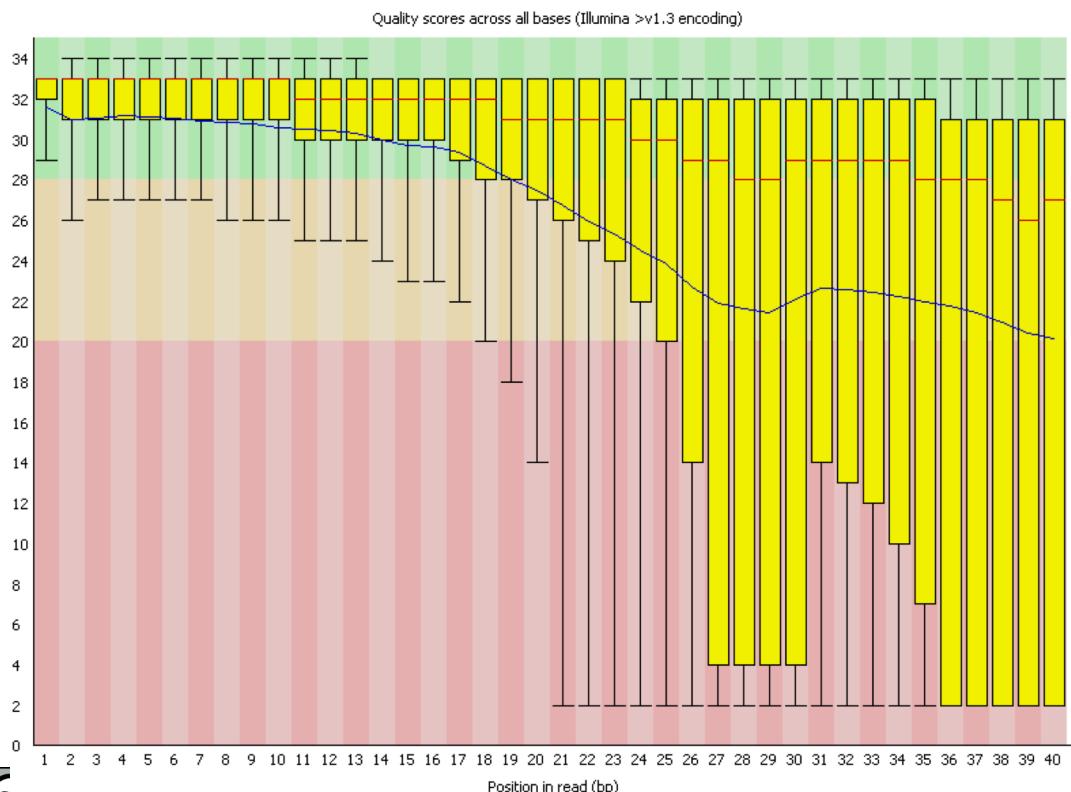
All-seq

@sequence: and information of it

GATTGGGGTTCAAAG

+

!***((***+))%%

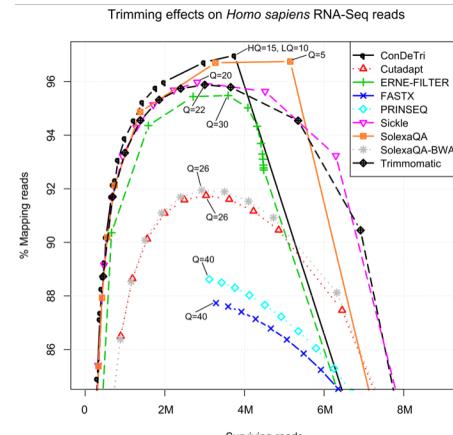


$$Q_{sanger} = -10 \log_{10}(p)$$

TRIMMING

Trim Galore!

http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/



Fabbro et al, 2013, PloS One

NGS: Short Introduction

History

NGS

Bioinformatics

FASTQ FORMAT

All-seq

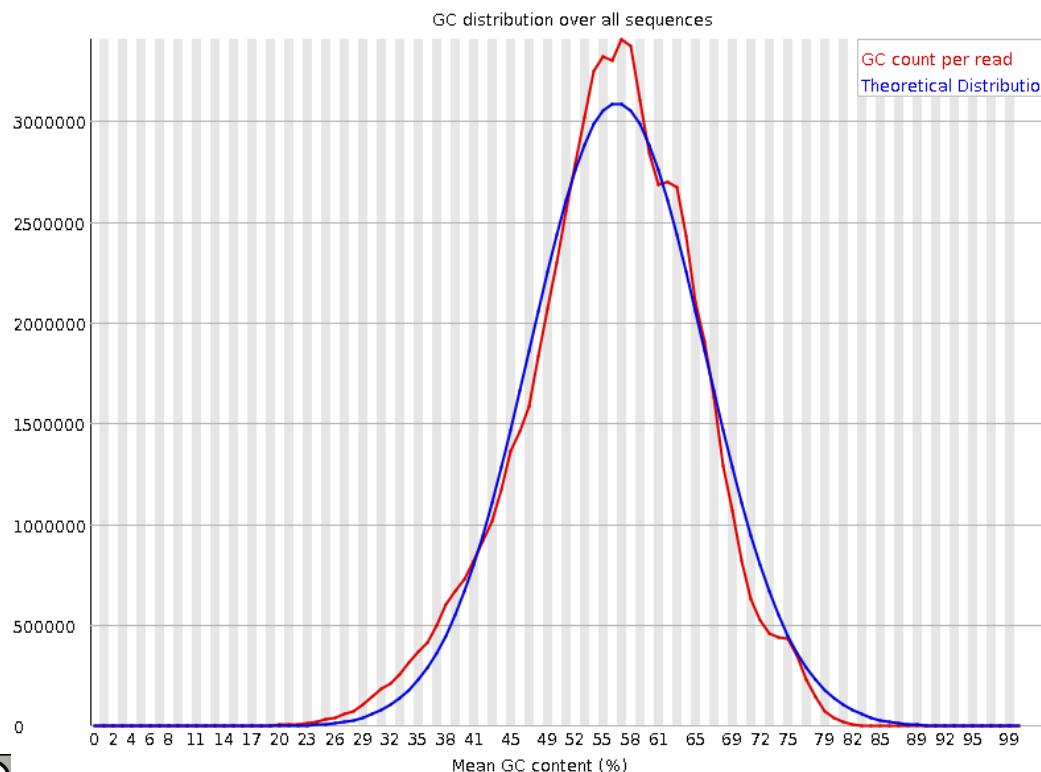
@sequence: and information of it

GATTGGGGTTCAAAG

+

!''*(((****+))%%

Amount of CG content



Theoretical Distribution
is different for RNA and
DNA applications

Important to have
similar proportions for
all samples in a data-set
(if expected)

NGS: Short Introduction

History

NGS

Bioinformatics

FASTQ FORMAT

All-seq

@sequence: and information of it

GATTGGGGTTCAAAG

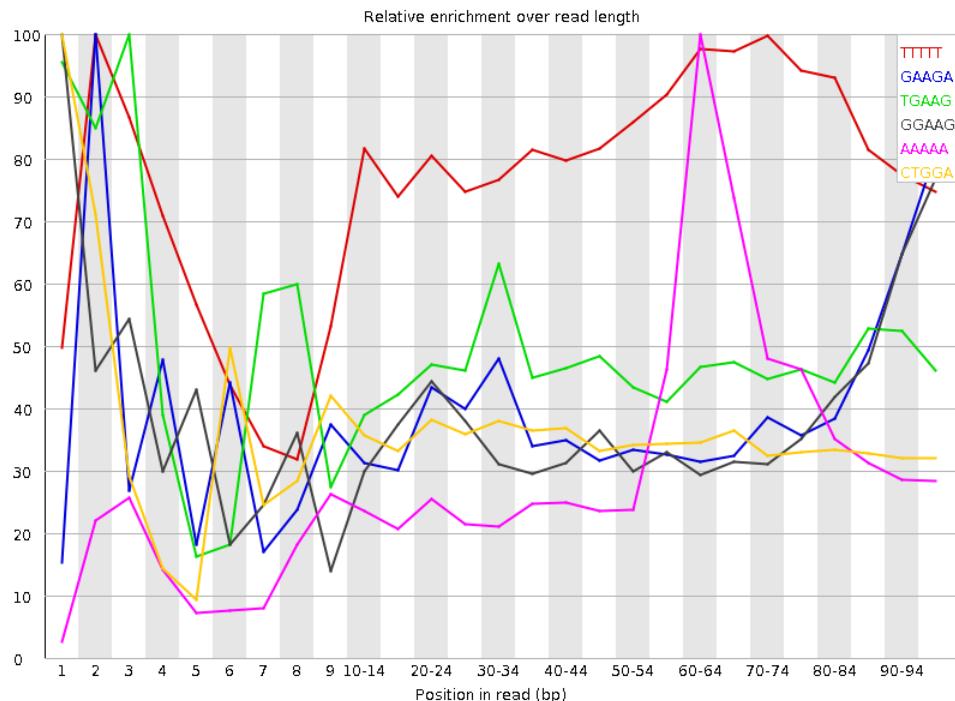
+

!"*(((****+))%%

K-mer over representation

- Over specific locations
- Over entire sequence

Adapters from Illumina?



TRIMMING

NGS: Short Introduction

History

NGS

Bioinformatics

FASTQ FORMAT

All-seq

@sequence: and information of it

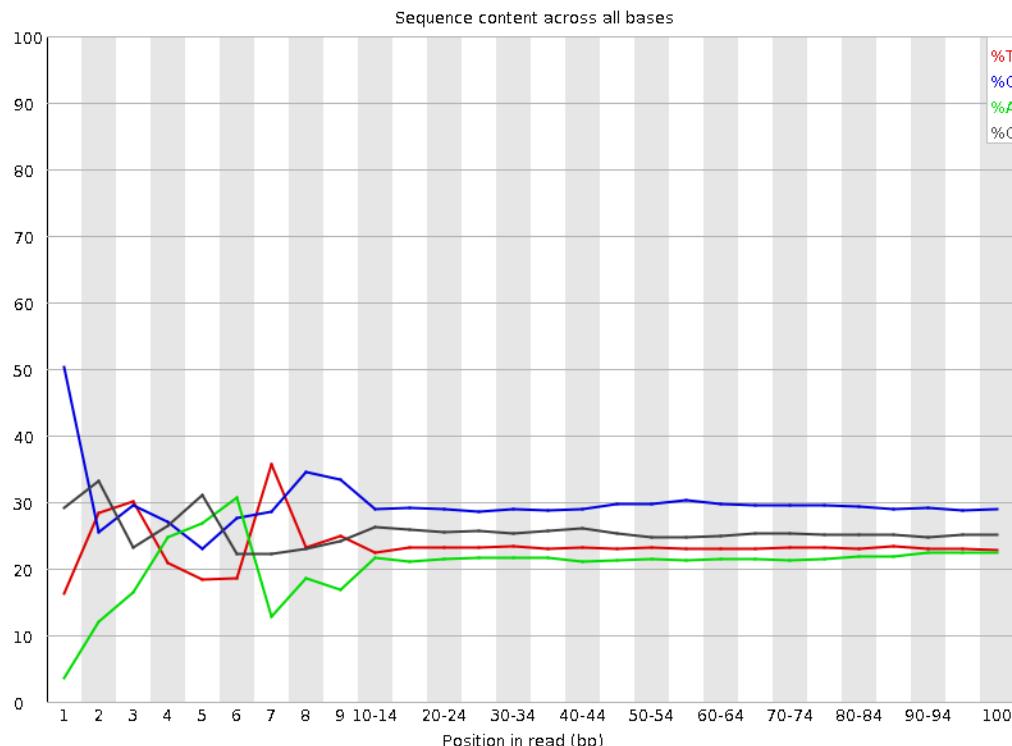
GATTGGGGTTCAAAG

+

!''*(((****+))%%

Nucleotide percentage per-base

From FastQC



It's worth noting that some types of library will always produce biased sequence composition, normally at the start of the read. Libraries produced by priming using random hexamers (including nearly all RNA-Seq libraries) and those which were fragmented using transposases inherit an intrinsic bias in the positions at which reads start. This bias does not concern an absolute sequence, but instead provides enrichment of a number of different K-mers at the 5' end of the reads. Whilst this is a true technical bias, it isn't something which can be corrected by trimming and in most cases doesn't seem to adversely affect the downstream analysis.

NGS: Short Introduction

History

NGS

Bioinformatics

FASTQ FORMAT



@sequence: and information of it

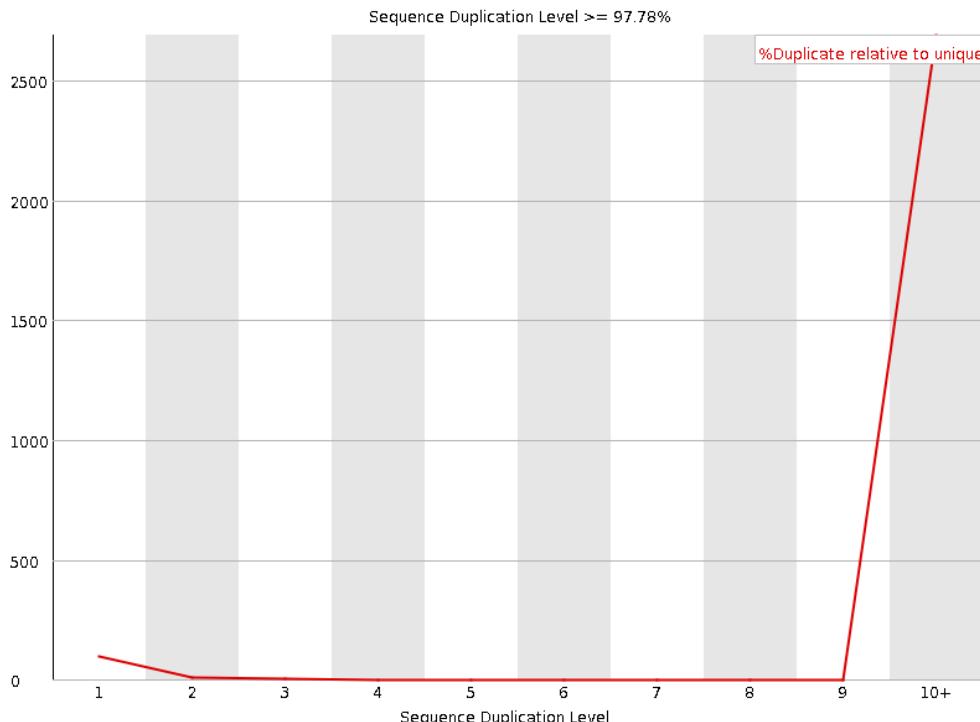
GATTGGGGTTCAAAG

+

!***((***+))%%

Duplication level:

- PCR artifact
- Cluster scanned twice
- Biological duplication



DNA and RNA very different expectations

Important to have similar proportions for all samples in a data-set (if expected)

NGS: Short Introduction

History

NGS

Bioinformatics

FASTQ FORMAT

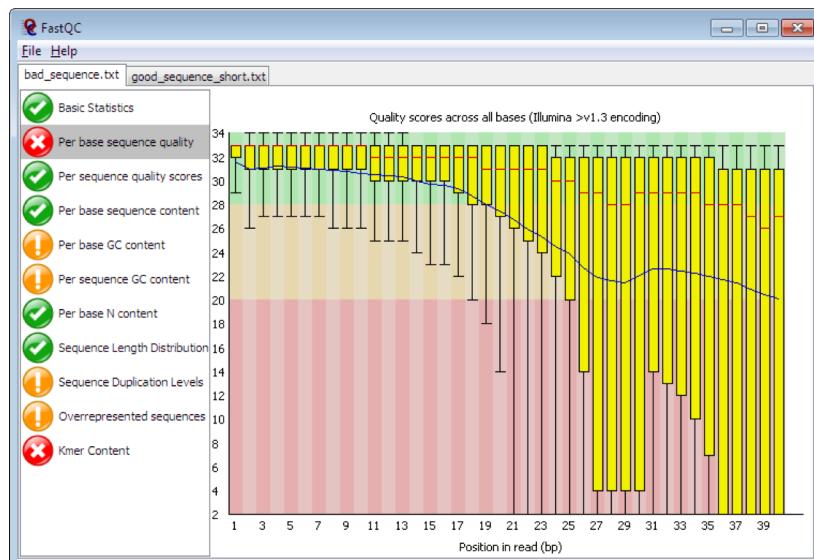


@sequence: and information of it

GATTGGGGTTCAAAG

+

!''*(((****+))%%

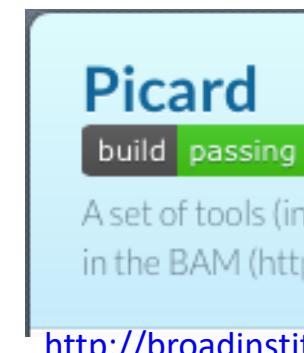


FastQC

Bioinformatics



<http://www.bioconductor.org/packages/devel/bioc/vignettes/sbgr/inst/doc/sbgr.html>



NGS: Short Introduction

History

NGS

Bioinformatics

All-seq

BCL to FASTQ

QUALITY
CONTROL

FASTQ to SAM/BAM

RNA-seq

miRNA-
seq

MAP TO A REFERENCE GENOME

Methodologies to map to the entire genome: BWA, Bowtie,...

Methodologies to map to selected parts of the genome: Tophat,...

More to be shown in the
RNA-seq session

single-cell
RNA-seq

DNase-
seq

RRBS-
seq

NGS: Short Introduction

