# Moderní metody analýzy genomu - analýza

## Mgr. Nikola Tom

Brno, 22.4.2015

# Deep-seq Workflow/Pipeline

### aims
- Detection of rare somatic mutations
- Detailed analysis of clonal evolution

### methods
- Ultra-deep sequencing of amplicons using Illumina Miseq
- Minimum coverage 2500x
- Highly accurate proof-reading polymerase
- 2x150 bp paired-end reads with overlapping ends
- Bayes based statistics

Base calling

Reads pre-processing

Mapping on reference

Local realignment

Reports

Beta-binomial- based variant detection

Quality based variant detection

Variant annotation

# GATK Workflow/Pipeline



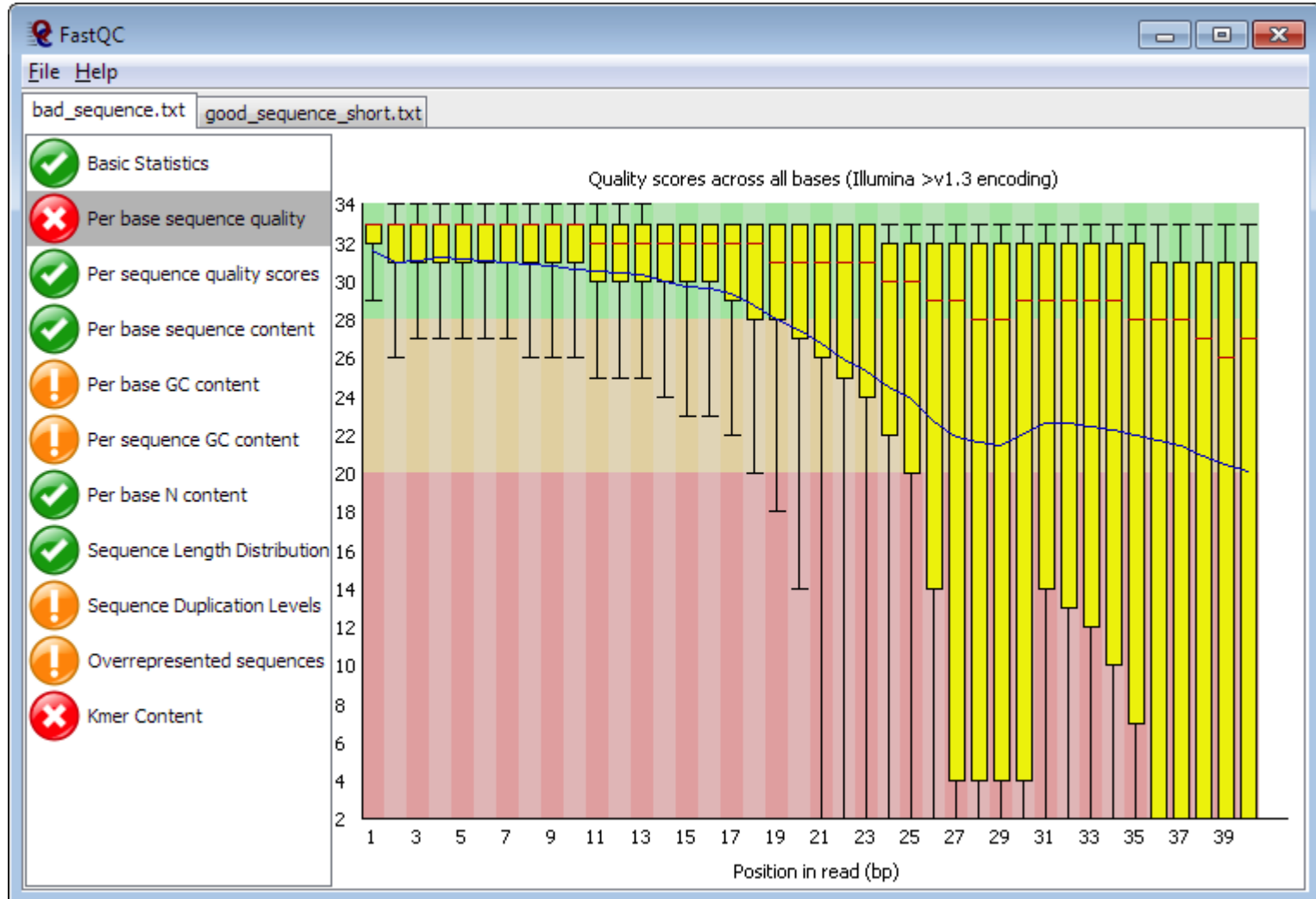https://www.broadinstitute.org/gatk/guide/best-practices

# Raw sequence = fastq

- Biological sequence
- Corresponding quality scores
- ASCII character
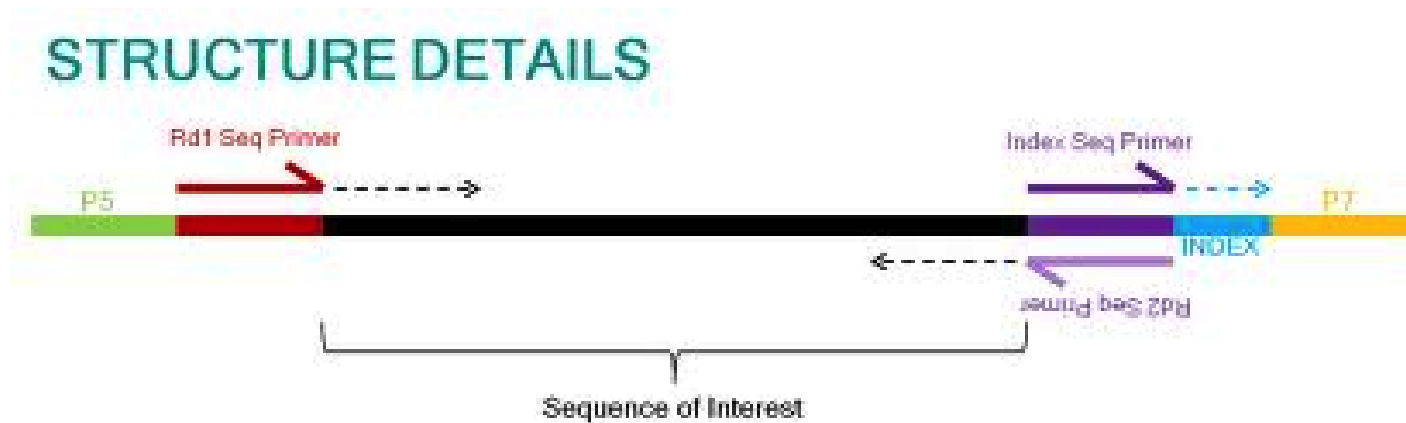- (fasta+ qual, csfasta + csqual, sff)

```
@
SEQ_ID GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
 !''*((((***+))%%%++)(%%%%).1***-+*''))**55CCF>>>>>>CCCCCCC65
```
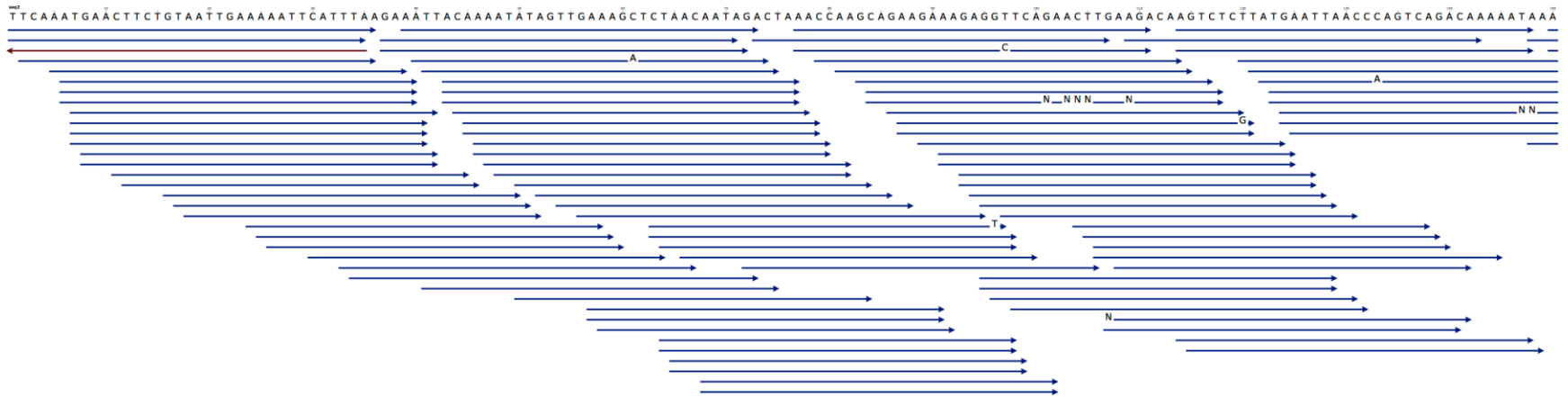
# FastQC

# Cutadapt

- Adaptor trimming (miRNA)
- Quality filtering
- Length filtering

## STRUCTURE DETAILS

Rd1 Seq Primer

Index Seq Primer

P5    P7

INDEX

Rd2 Seq Primer

Sequence of Interest

# Read mapping => SAM, BAM

- Usually mapping reads on reference
- miRNA - special case
  - Grouping and annotate against mirBase
- **DNA**
  - BWA, Bowtie, Bfast, SHRiMP, CLC
- **RNA**
  - TopHat (*de novo* splice aligner)
- **Commercial**
  - CLC Genomics Workbench
- ***De novo* assembly** – unknown genomes

# Alignment

# SAM



Each row describes a single alignment of a raw read against the reference genome.
Each alignment has 11 mandatory fields, followed by any number of optional fields.

# Mapping, Coverage reports

- Important checkout for lab protocol
- Specificity of PCR
- Normalization
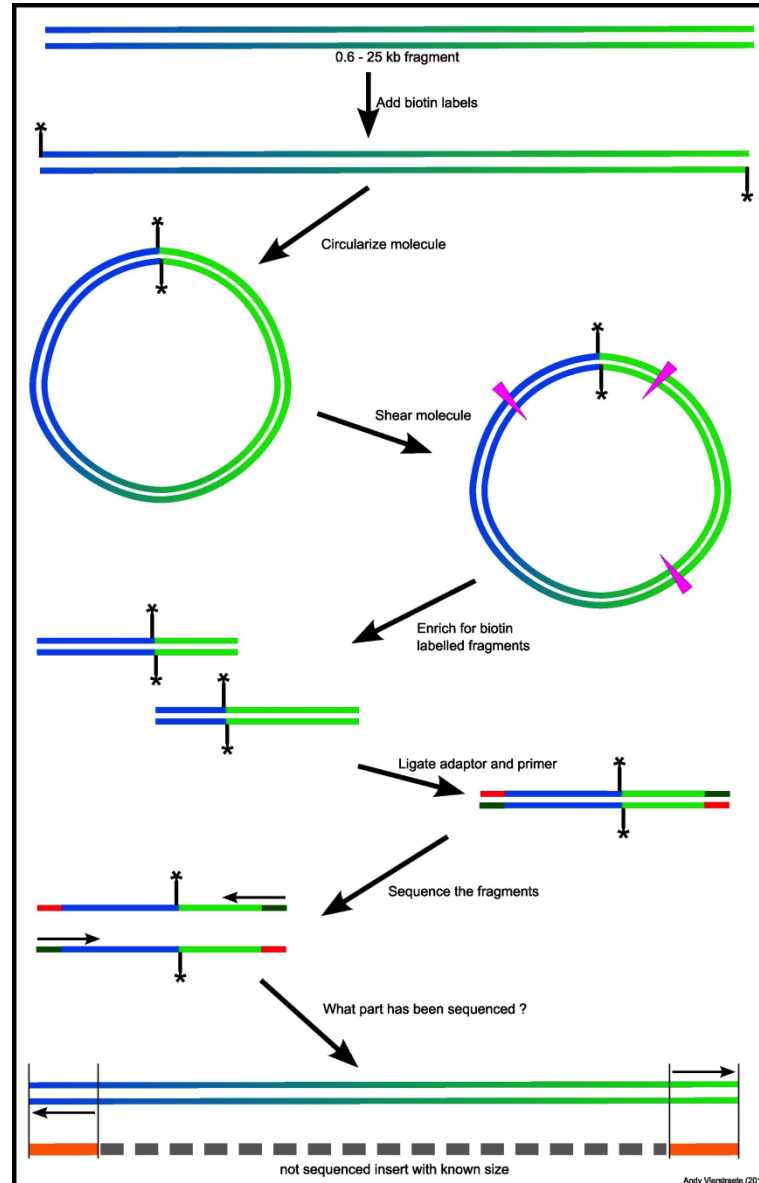- Settings of variant calling threshold, CNV

# SNV and small InDel Calling

- Coverage
- Frequency
- Base quality

- !!!
- Genomic context (homopolymers)
- Nucleotide type
- Position in read (errors at the read end)
- Alignment errors

# Structural variations

- Mate-pair library
- Long InDel
- Translocation



0.6 - 25 kb fragment

Add biotin labels

Circularize molecule

Shear molecule

Enrich for biotin labelled fragments

Ligate adaptor and primer

Sequence the fragments

What part has been sequenced ?

not sequenced insert with known size

Andy Vierstraete (2012)

# Annotating and filtering

- Gene
- Transcript
- dbSNP
- Regulation
- Comparative genomics
- Repeats
- Functional
- Gene ontology
- miRNA targets
- Etc.