# Moderní metody analýzy genomu - analýza

## Mgr. Nikola Tom

Brno, 20.11.2015

# Before we start analysis

We have to know what we are dealing with… and what we want to find out…

**Concept of the project**
DNA/RNA/methylation/…

**DNA**
Targeted sequencing (amplicons, gene panels, exomes)
Whole genome sequencing
- Finding differences to known reference genome = re-sequencing

***De novo* assembly**
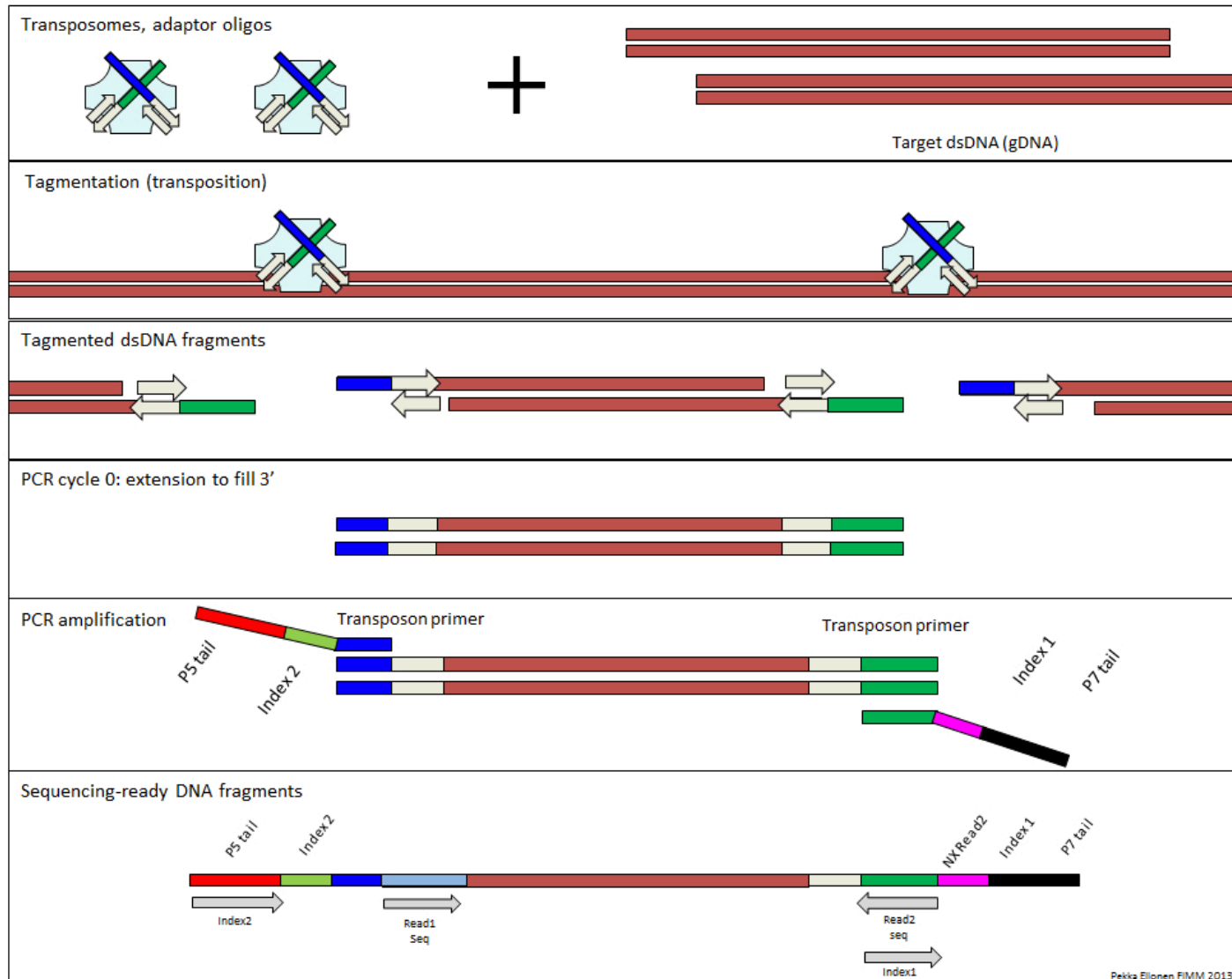- Genome construction

# Before we start analysis…

**RNA**
- Gene expression, alternative splicing
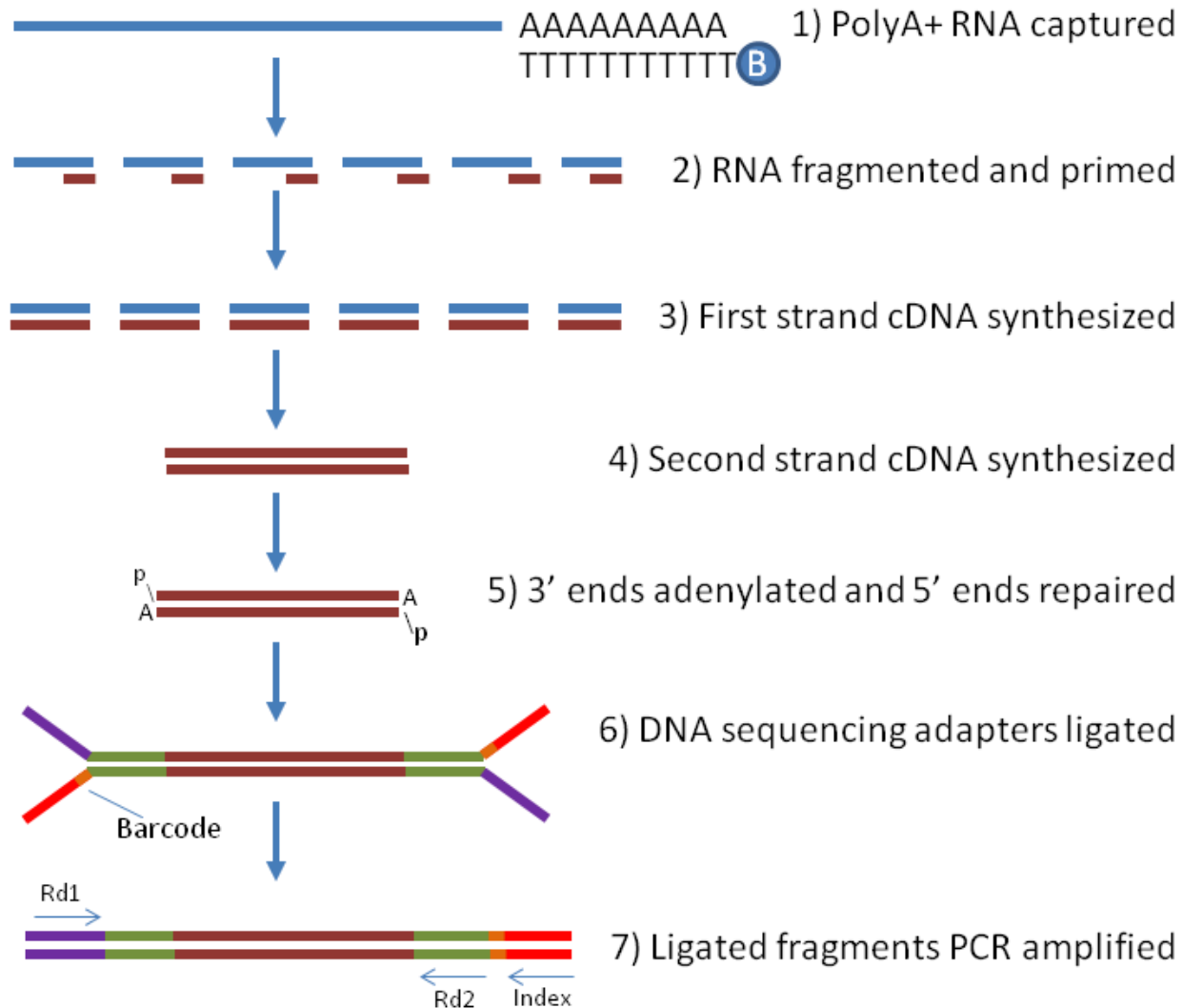
**Metagenomics** (bacteria, viruses)
- Their composition, variants

**ChIP sequrcing** (DNA-protein interactions)

# Library preparation – example of DNA library

# Library preparation – example of mRNA library

AAAAAAAAA    1) PolyA+ RNA captured
TTTTTTTTTT**B**

2) RNA fragmented and primed

3) First strand cDNA synthesized

4) Second strand cDNA synthesized

p
A    A    5) 3' ends adenylated and 5' ends repaired
p

6) DNA sequencing adapters ligated

Barcode

Rd1

7) Ligated fragments PCR amplified

Rd2    Index

# Bioinformatics

Bioinformatics is a quite new field… (first NGS in 2005)

How to analyse data defived from NGS = bottleneck of NGS

A lot of tools/software for NGS data analysis…

Most of the tools are command-line based

No tool is working perfectly…☹

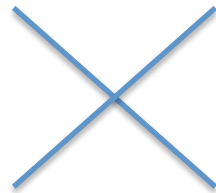Each tools solves only a peace of the cake…
NO tool, that is able to perform analysis from the very beginning
to the end => Need for setup the **pipeline**

# Bioinformatics

Exception: commercial software and ready to use pipelines **BUT** they have usually not-transparent settings and/or not enough of options

Heavily depends on type of experiment, library preparation and project

Laptop or PC are usually not enough… need for cluster

# Pipeline/Workflow

Base calling

Reads pre-processing

Mapping on reference

Quality based variant detection

Post-processing

Local realignment

Beta-binomial- based variant detection

Biological interpretation

results

results

results

Variant annotation

MiSeq
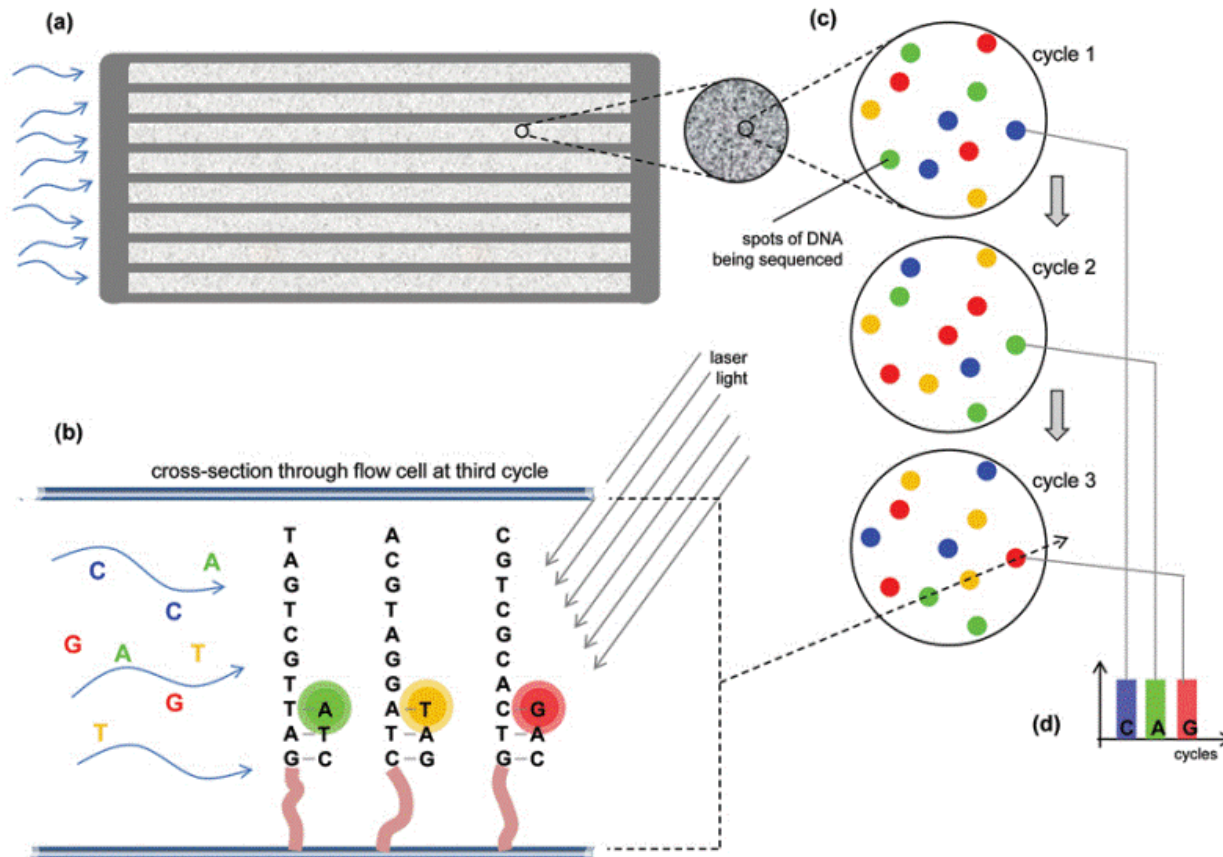
# Base calling

Signal to sequence conversion and assigning base quality scores (fastq file)
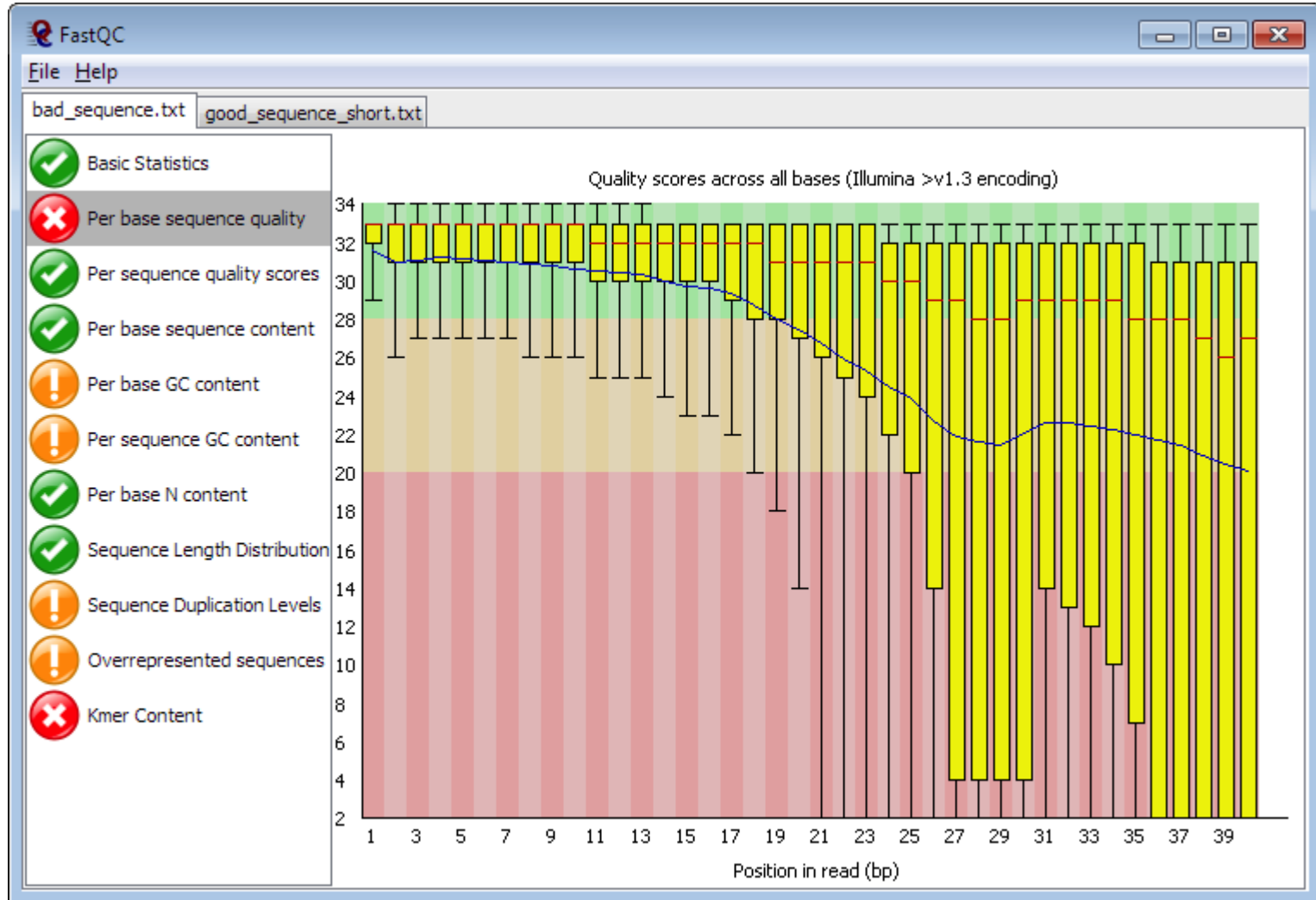**Phred score** – probability of arising an error (log based)

# fastq

- Consists of reads - biological sequences
(each read represents 1 input molecule sequenced on flowcell)
- Corresponding quality score for each base
- ASCII character
- (fasta+ qual, csfasta + csqual, sff)
- Pair-end sequencing – 2 fastq files

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

# Quality control (FastQC)

# Pipeline/Workflow

Base calling

Reads pre-processing

Mapping on reference

Quality based variant detection

Post-processing

Local realignment

Beta-binomial- based variant detection

Biological interpretation

results
results
results

Variant annotation

©2011, Illumina Inc. All rights reserved.
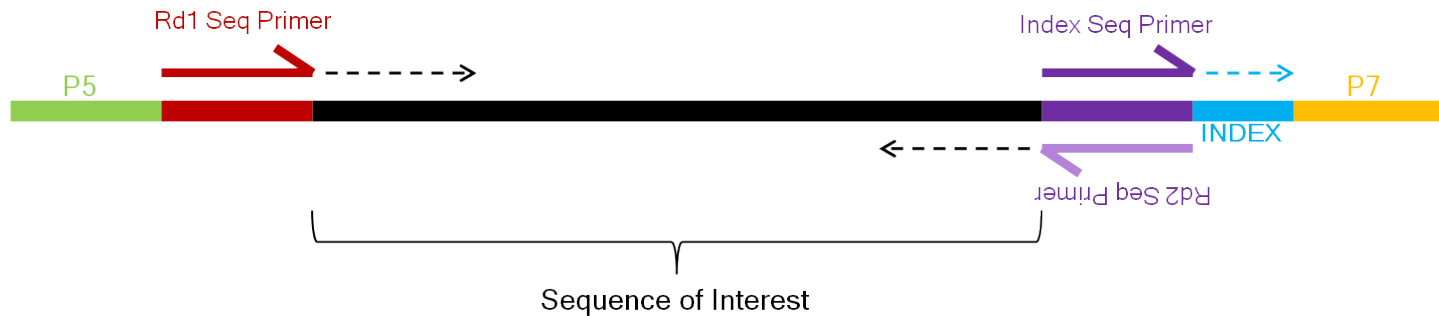
# Cleaning reads (Cutadapt)

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

## STRUCTURE DETAILS



Rd1 Seq Primer

Index Seq Primer

P5

P7

INDEX

Rd2 Seq Primer

Sequence of Interest

# Pipeline/Workflow

Base calling

Reads pre-processing

Mapping on reference

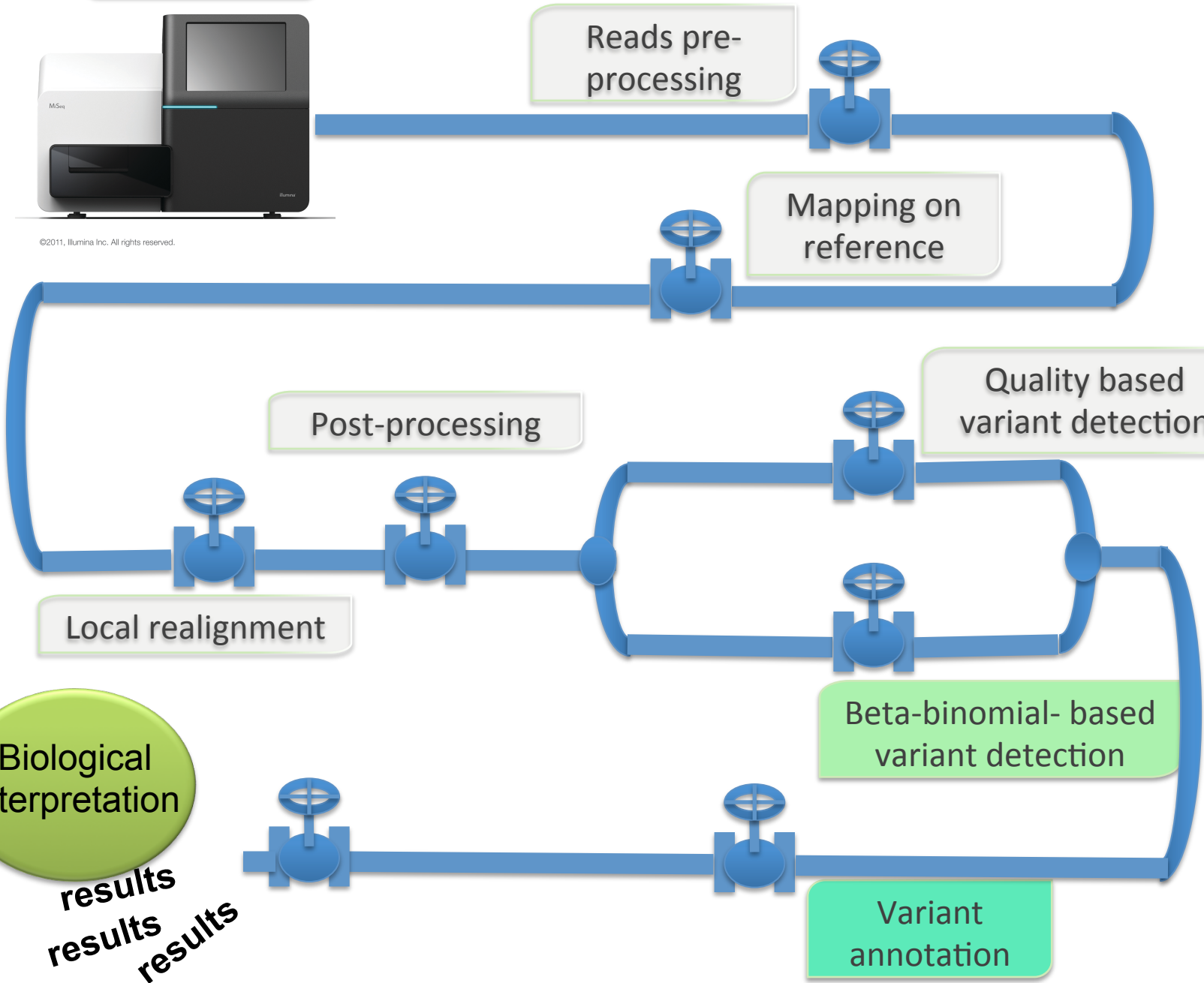Quality based variant detection

Post-processing

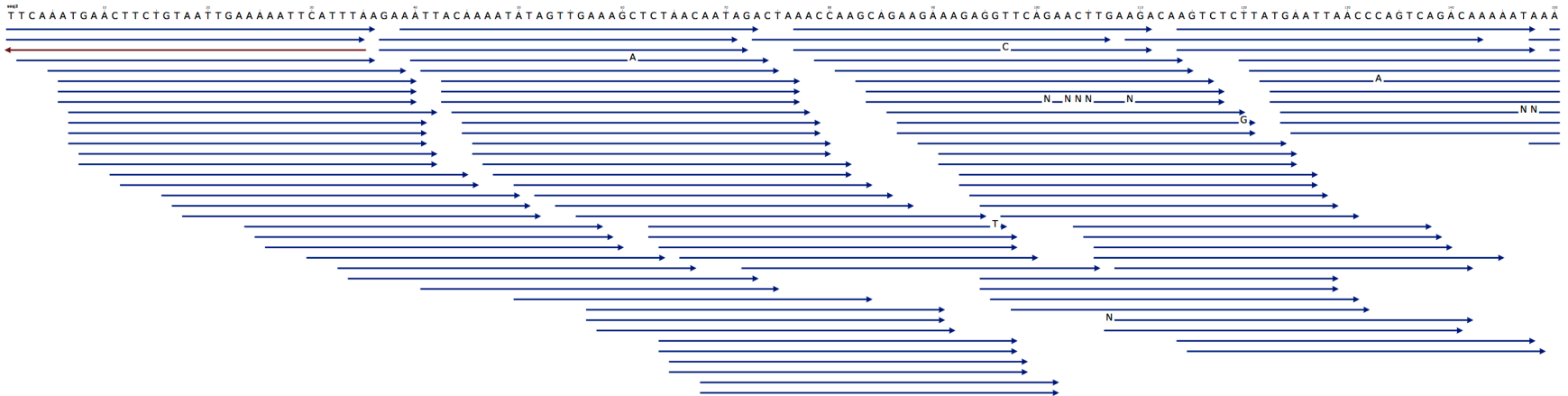Local realignment

Beta-binomial- based variant detection

Biological interpretation

results
results
results

Variant annotation

# Read mapping (alignment)

• Usually mapping reads on reference sequence (DNA/cDNA/16S/other seq) to find corresponding location & differences

• Problem with too many sequences and billions bp long references – need for special algorithms (Burrows-Wheeler transform, hash table indexing)
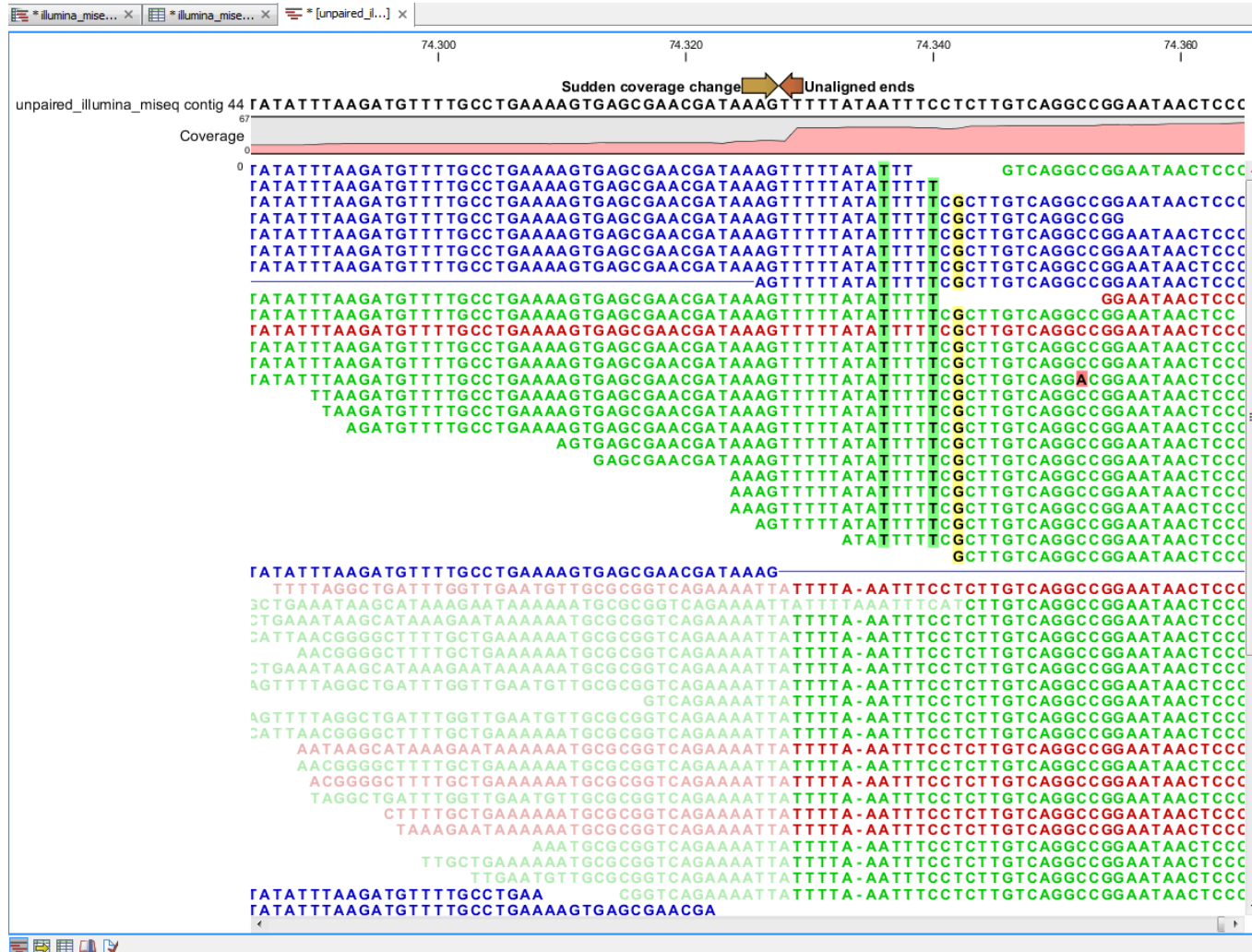
# Mapping of DNA reads

• On Existing DNA reference sequence
(ready for many organisms)

• To find substitutions, insertions, deletions, inversions, etc...
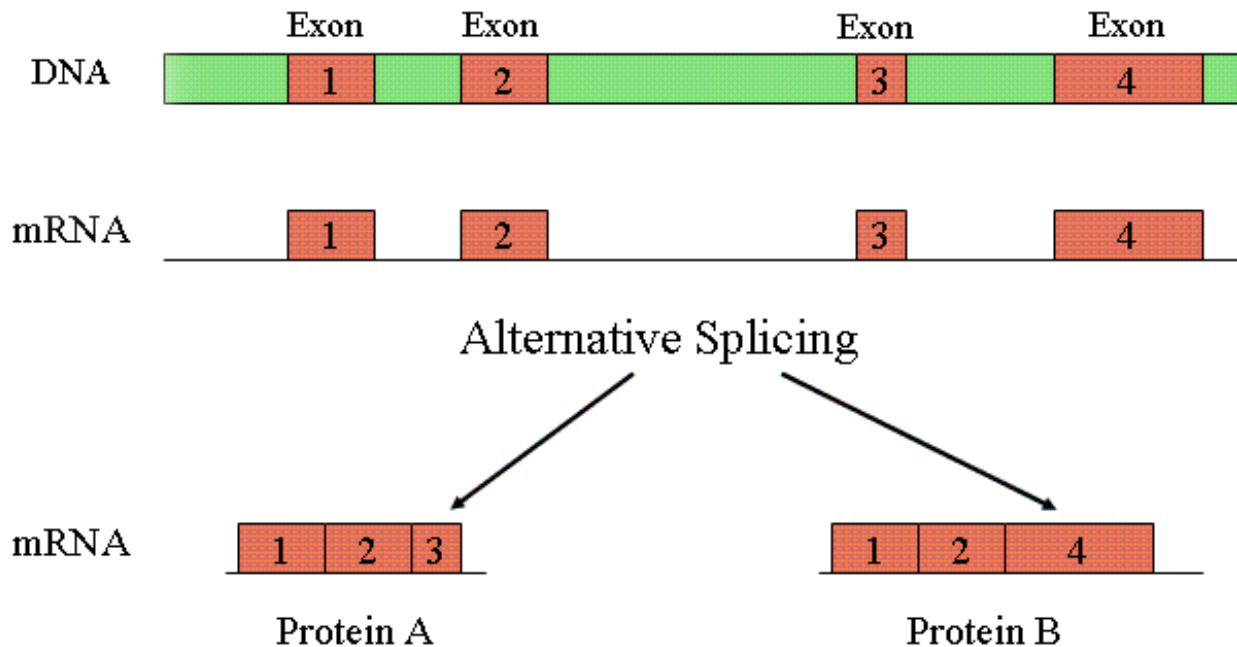**Precisely!**

– BWA, Bowtie, Bfast, SHRiMP

# Example of DNA re-sequencing

# Mapping of RNA reads – alternative splicing

Reads can span exon junctions
- mRNA splicing

# Mapping of RNA reads

- To measure gene expression OR alternative splicing

- On existing **DNA** reference sequence
- To find **alternative splicing**
- More tricky, complicated, slower
- TopHat (*de novo* splice aligner)

On **transcriptome** reference sequences
Reads can map to multiple transcripts (shared exons)
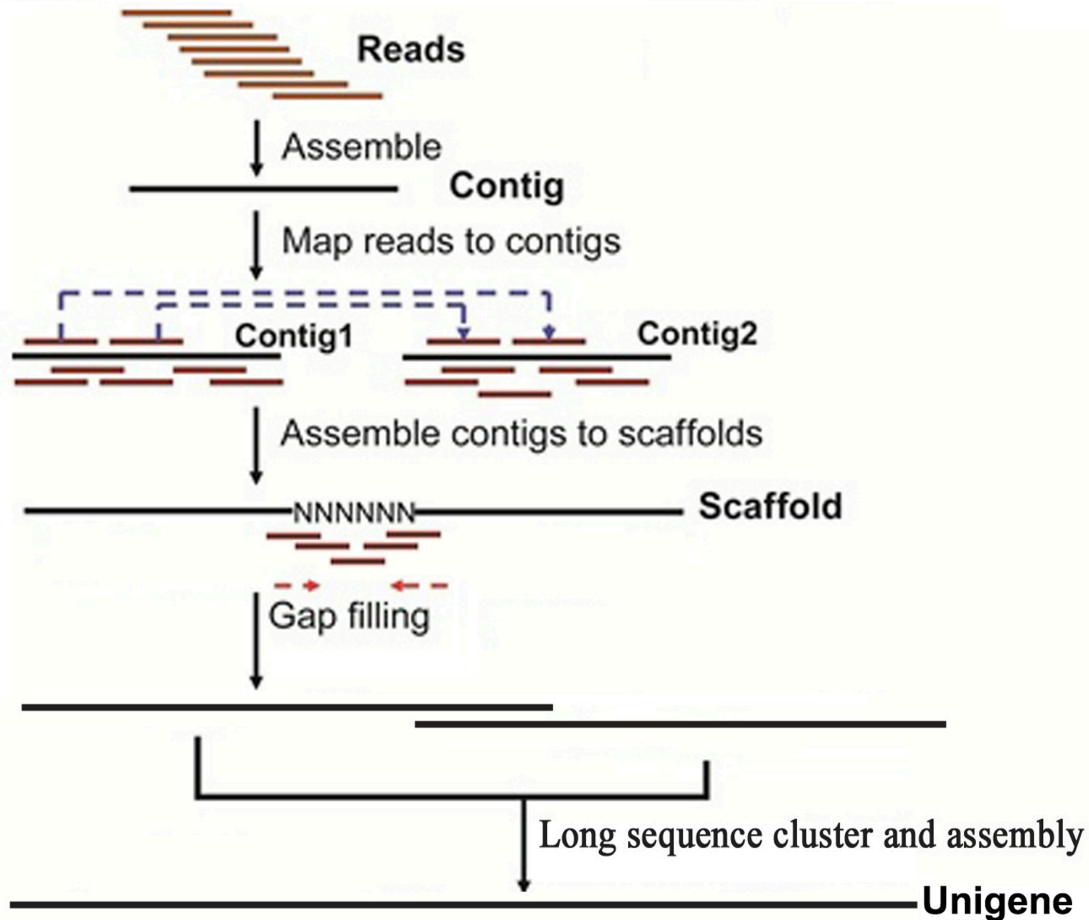Easier, faster, no need for special aligners
- BWA

- On **miRNA** sequences - miRBase
– Grouping and annotate against mirBase

# *De novo* assembly

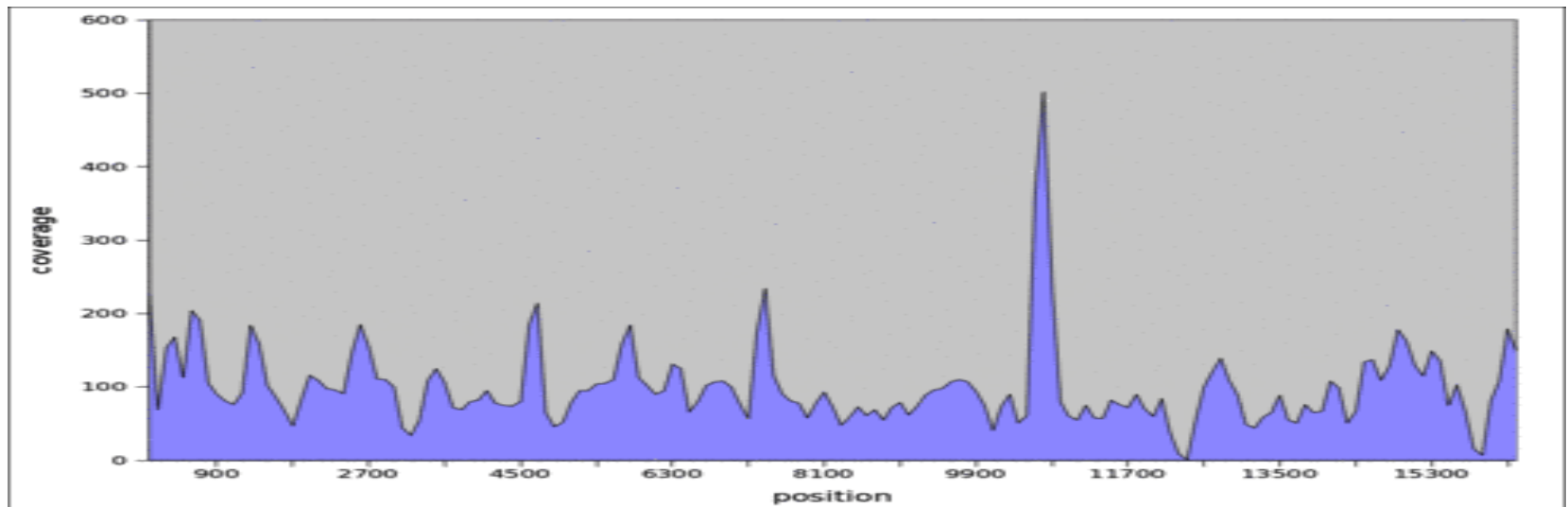– to uncover unknown genomes/transcriptomes
- To detect large structural variants

# SAM/BAM



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

- Repeat alignment/other steps with different criteria?
- Important checkout for lab protocol
- Specificity of PCR
- Settings of variant calling threshold, CNV

# Pipeline/Workflow

Base calling

Reads pre-processing

Mapping on reference

Quality based variant detection

Post-processing

Local realignment

Beta-binomial- based variant detection

Biological interpretation

Variant annotation

results results results

©2011, Illumina Inc. All rights reserved.

MiSeq

# Indel realignment

Usually alignment is not perfect – false positive indels & Substitutions => Need for local indel realignment

**Pipeline/Workflow**

Base calling

Reads pre-processing

Mapping on reference

Quality based variant detection

Post-processing
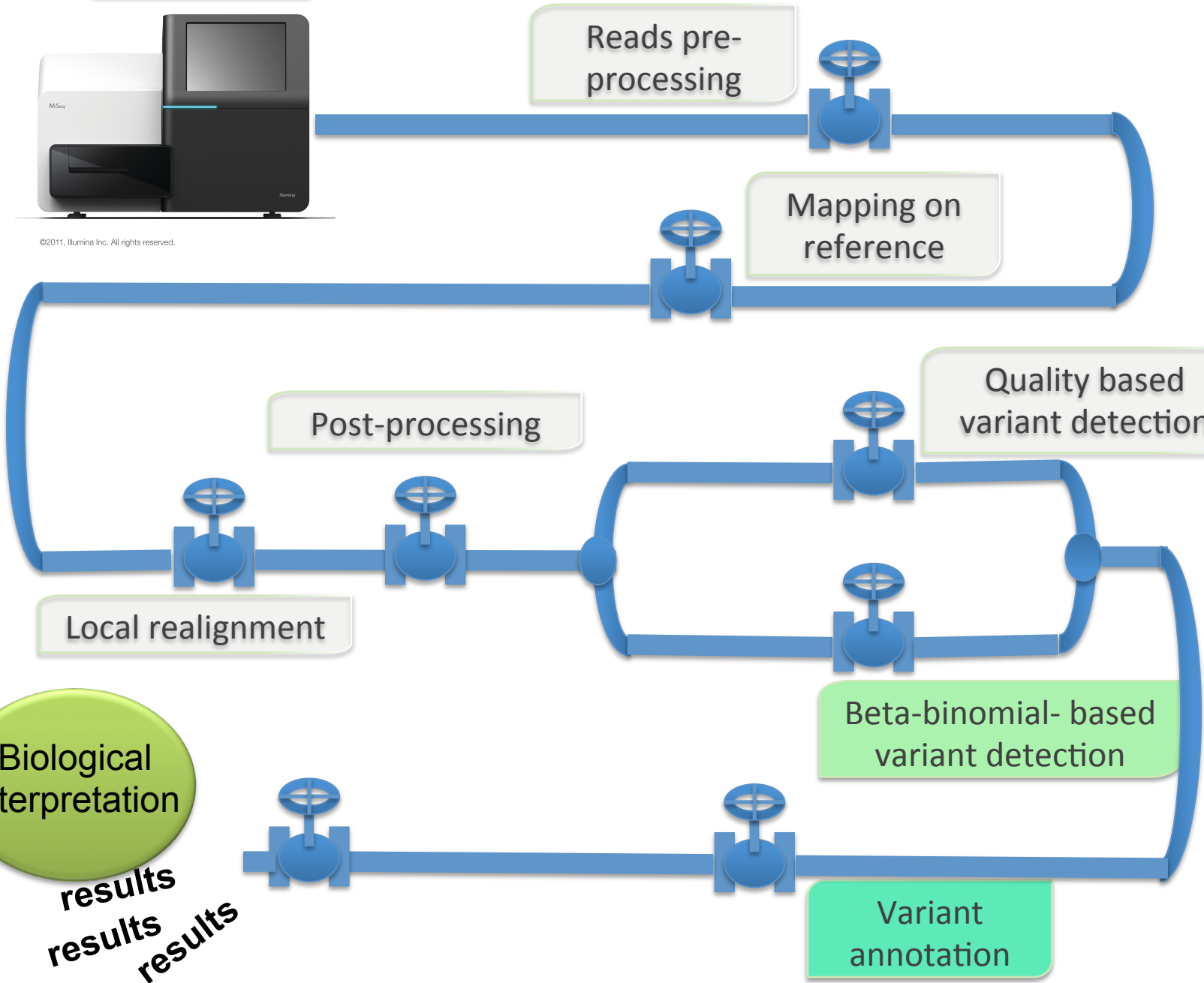
Local realignment

Beta-binomial- based variant detection

Biological interpretation

Variant annotation

results
results
results

# Remove PCR duplicates

Each read represents 1 input molecule

THEORY:
E.g. in case of DNA re-sequencing, 1 diploid cell is represented by 2 reads because of 2 chromosomes
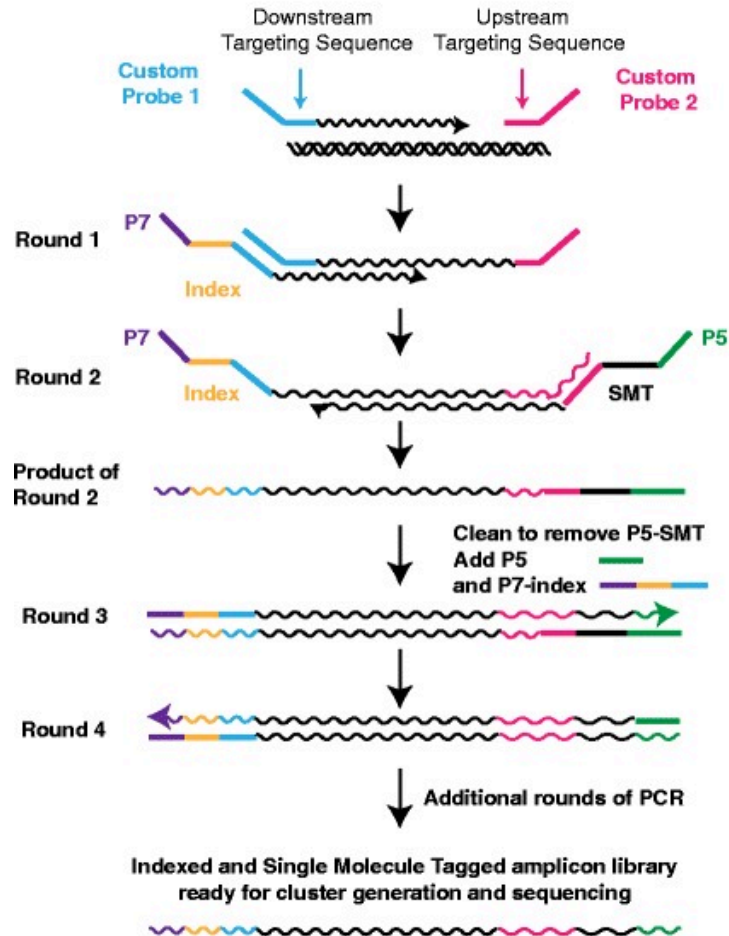
BUT

there is a PCR to amplify genetic material to be analyzable =>
1 input molecule from 1 cell could be after PCR represented by more reads => Biased variant allele frequency
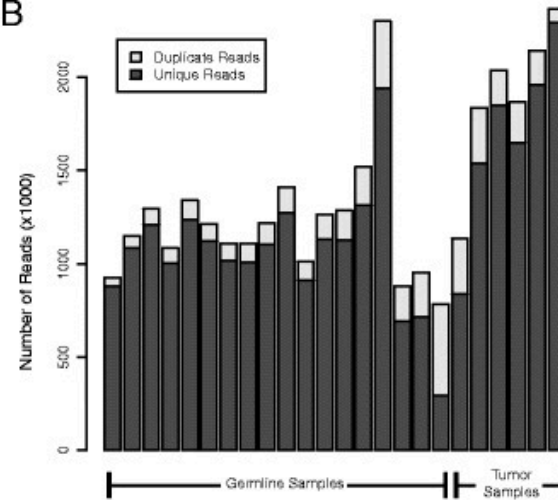
How to solve it?

1) Molecular barcodes (very new method)
2) Identity of start-end positions of read pair

# Molecular barcodes



Smith et al. 2014

# Pipeline/Workflow
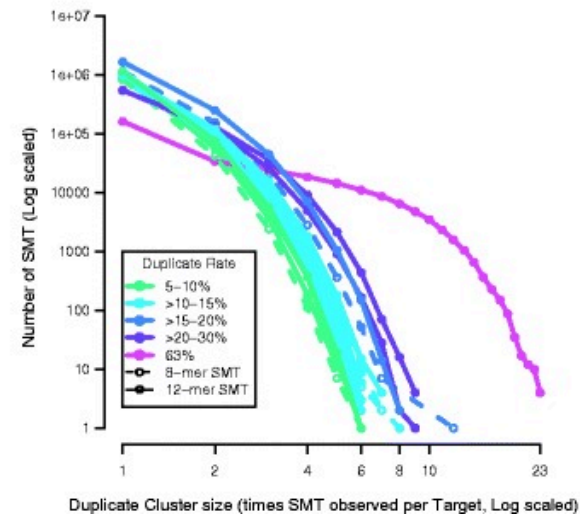
Base calling

Reads pre-processing

Mapping on reference

Quality based variant detection

Post-processing
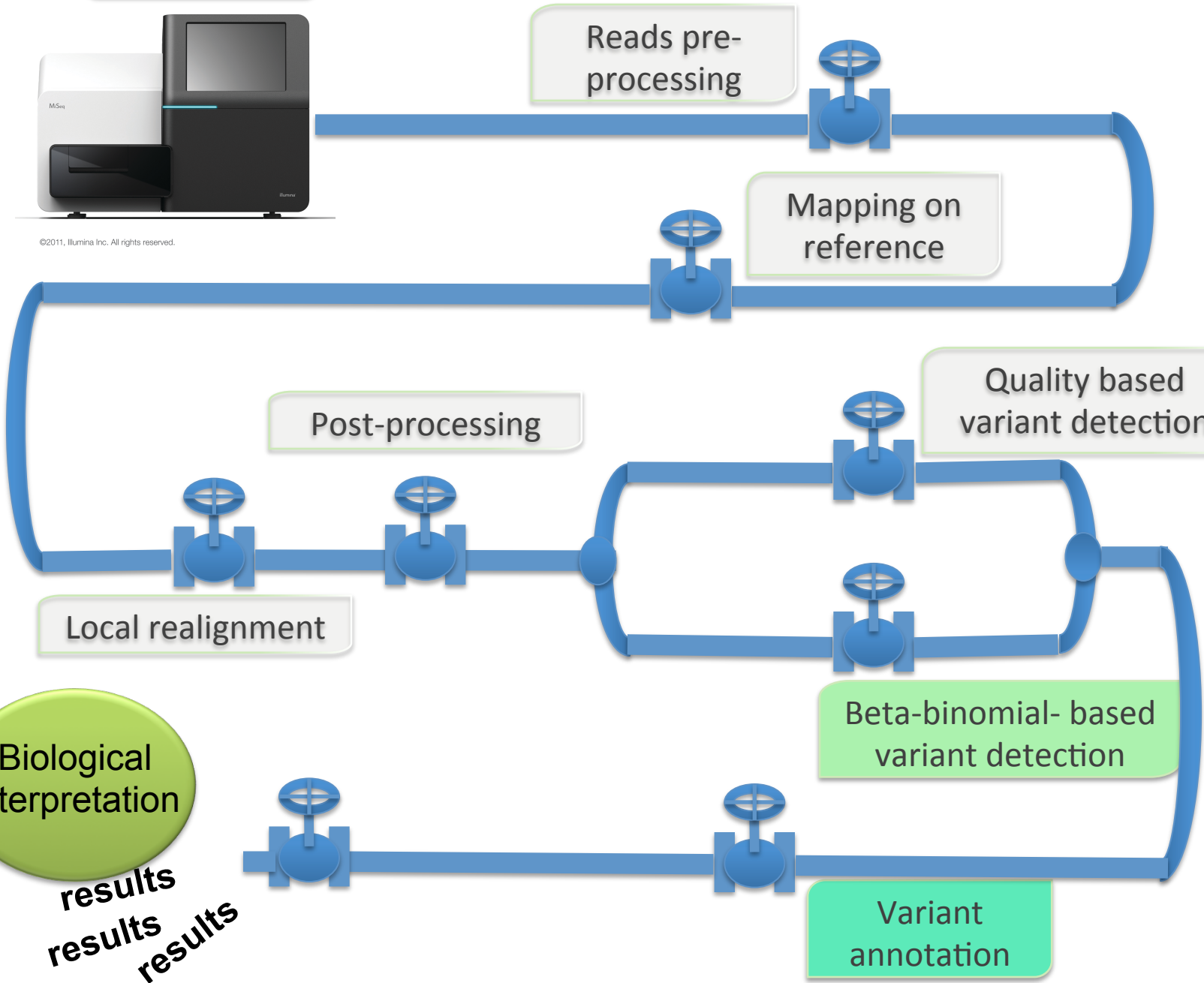
Local realignment

Beta-binomial- based variant detection

Biological interpretation

results
results
results

Variant annotation

# DNA Seq - variant calling

- To detect differences from reference sequence

- Single/multi-nucleotide
- Substitutions
- Insertions
- Deletions

- Inversions
- Large structural variations (translocations, indels)
- Copy number variations

# DNA Seq variant calling

based on many criteria like:
- Coverage
- Variant alelle frequency
- Base quality

- Depends also on:
- Genomic context (homopolymers)
- Nucleotide type
- Position in read (errors at the read end)
- Alignment errors (importance of realignment)
- Presence in both forward and reverse reads

Necessary to take into account type of library preparation (single end; pair end; mate pair)

# DNA Seq variant calling

- Mate-pair library
- Detection of large indels
& translocations

# DNA Seq variant calling

# vcf file



**Example**

```
##fileformat=VCFv4.0                                                        Mandatory header lines
##fileDate=20100707
##source=VCFtools
##reference=NCBI36                                                          Optional header lines (meta-data
##INFO=<ID=AA,Number=1,Type=String,Description="Ancestral Allele">          about the annotations in the VCF body)
##INFO=<ID=H2,Number=0,Type=Flag,Description="HapMap2 membership">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality (phred score)">
##FORMAT=<ID=GL,Number=3,Type=Float,Description="Likelihoods for RR,RA,AA genotypes (R=ref,A=alt)">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##ALT=<ID=DEL,Description="Deletion">
##INFO=<ID=SVTYPE,Number=1,Type=String,Description="Type of structural variant">
##INFO=<ID=END,Number=1,Type=Integer,Description="End position of the variant">
#CHROM POS ID      REF   ALT      QUAL FILTER  INFO                FORMAT    SAMPLE1    SAMPLE2
1      1    .       ACG   A,AT      .   PASS    .                   GT:DP     1/2:13     0/0:29
1      2    rs1     C     T,CT      .   PASS    H2;AA=T             GT:GQ     0|1:100    2/2:70
1      5    .       A     G         .   PASS    .                   GT:GQ     1|0:77     1/1:95
1      100  .       T     <DEL>     .   PASS    SVTYPE=DEL;END=300  GT:GQ:DP  1/1:12:3   0/0:20
```

VCF header

Body

Reference alleles (GT=0)

Alternate alleles (GT>0 is an index to the ALT column)

Deletion

SNP

Large SV

Insertion

Other event

Phased data (G and C above are on the same chromosome)

# DNA Seq variant calling

- Tumor only (amplicon sequencing & diagnostics)
- Tumor & normal (exome sequencing)
-to do variant calling and genotyping more precisely
(somatic, germinal mutations)

- Option is also to analyze tumor vs. group of tumors

Application of many statistical tests:
- negative beta-binomial test
- Bayesian statistics
- Fisher exact test

As higher coverage as higher sensitivity and specificity (but limited)

More about statistics and RNA sequencing in the next courses

# Pipeline/Workflow

Base calling

Reads pre-processing

Mapping on reference

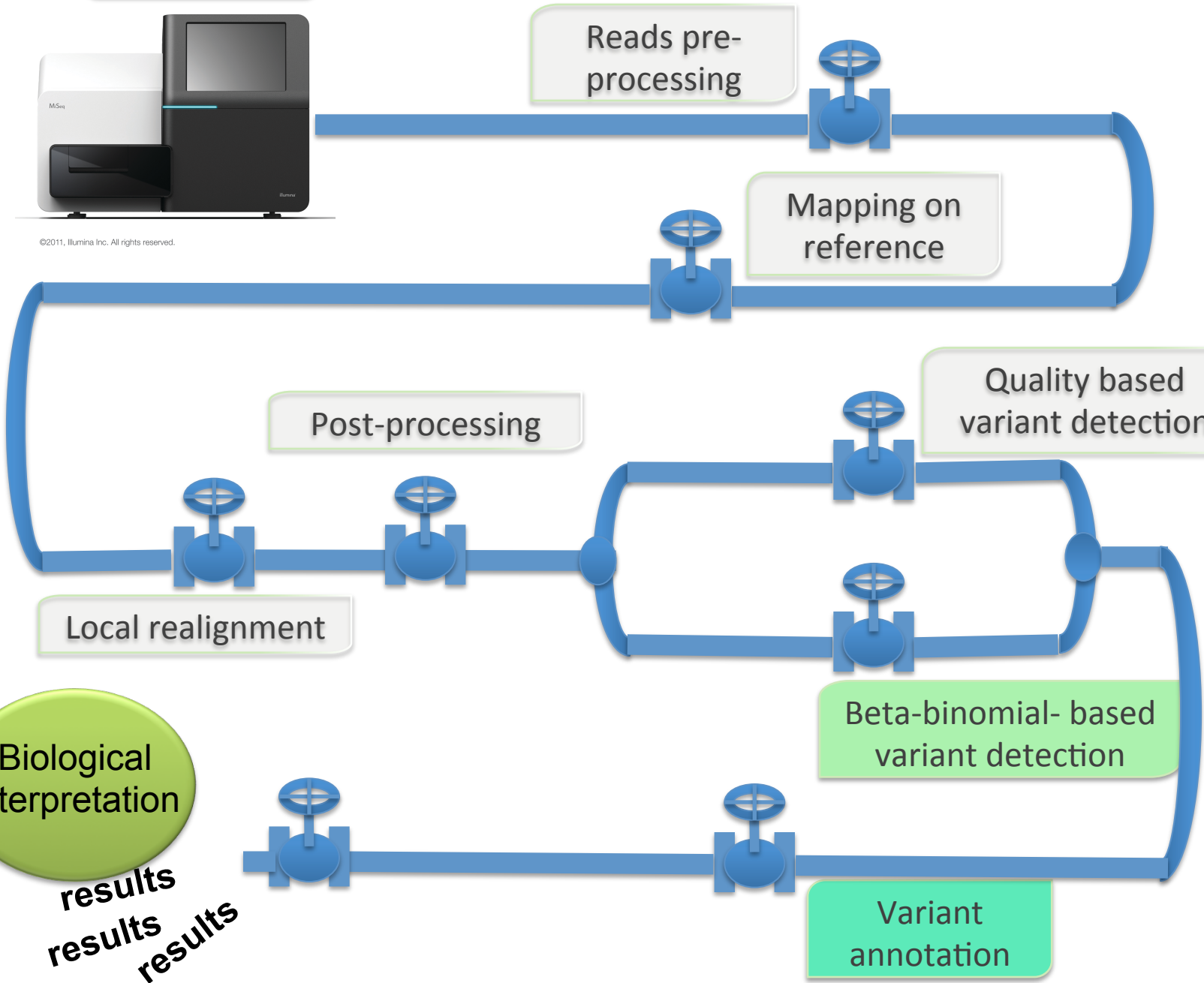Quality based variant detection

Post-processing

Local realignment

Beta-binomial- based variant detection

Biological interpretation

results
results
results

Variant annotation

# Annotating and filtering of detected variants

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