CEITEC

Central European Institute of Technology
BRNO | CZECH REPUBLIC

# Moderní metody analýzy genomu Bioinformatika I

## Mgr. Nikola Tom

Brno,
11.11.2016

# Bioinformatics

Bioinformatics is a quite new field… (first NGS in 2005)
How to analyse data defived from NGS = bottleneck of NGS

**AIM:** clean the data and give them biological sense

Bioinformatics **SOLUTION 1**:

- commercial software and ready to use pipelines
**BUT** they have usually not-transparent settings and/or
not enough of options
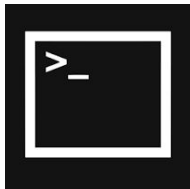(good programs expensive)

# Bioinformatics

Bioinformatics **SOLUTION 2**:

* command-line based tools/software
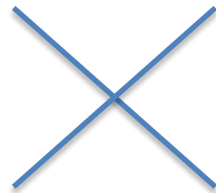Each tools solves only a part of the analysis
* Need for setup the pipeline & tune programs' parameters
(challenging & more precise!!!)

# Bioinformatics

Choice of programs & settings heavily depends on type of experiment, library preparation, biological question

Laptop or PC are usually not enough… need for cluster

# 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, ncRNA, alternative splicing

**Metagenomics** (bacteria, viruses)
- Composition of organisms in the sample, genetic variants
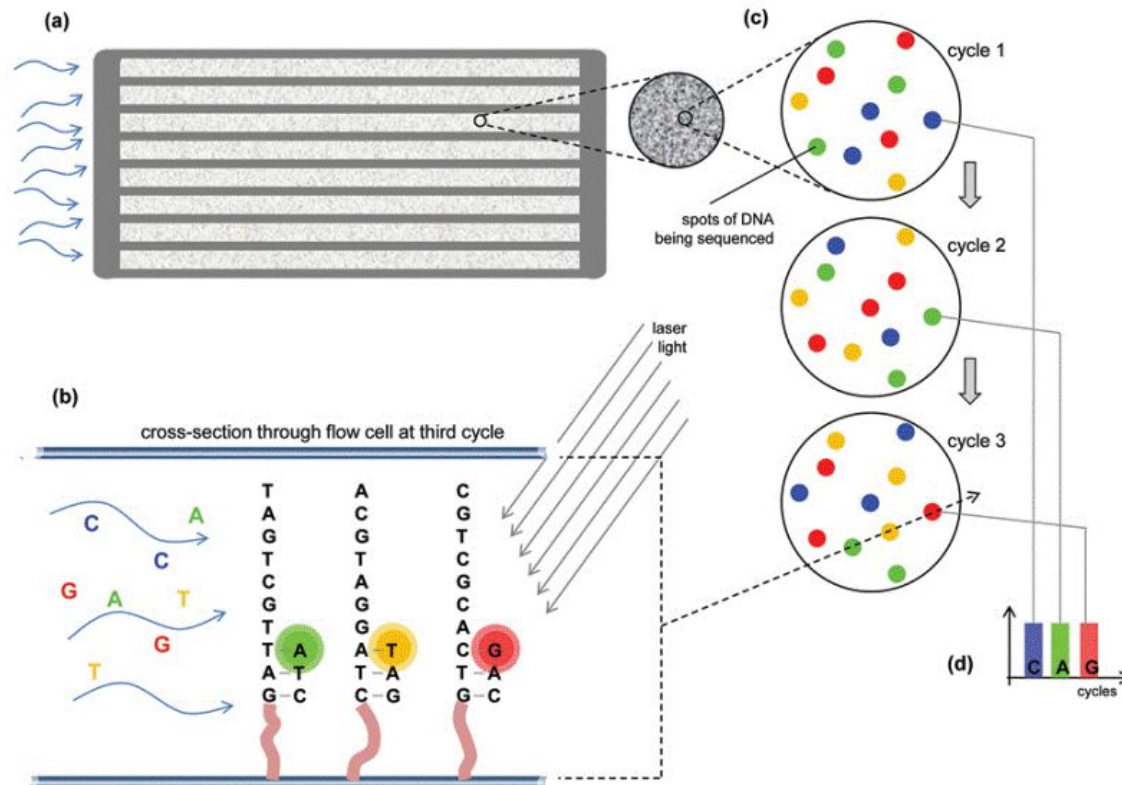
**ChIP sequecing** (DNA-protein interactions)

# Bioinformatics' starting point

**Raw** sequencing data - READ

Produced during **base calling**
- signal to sequence conversion and assigning base quality scores
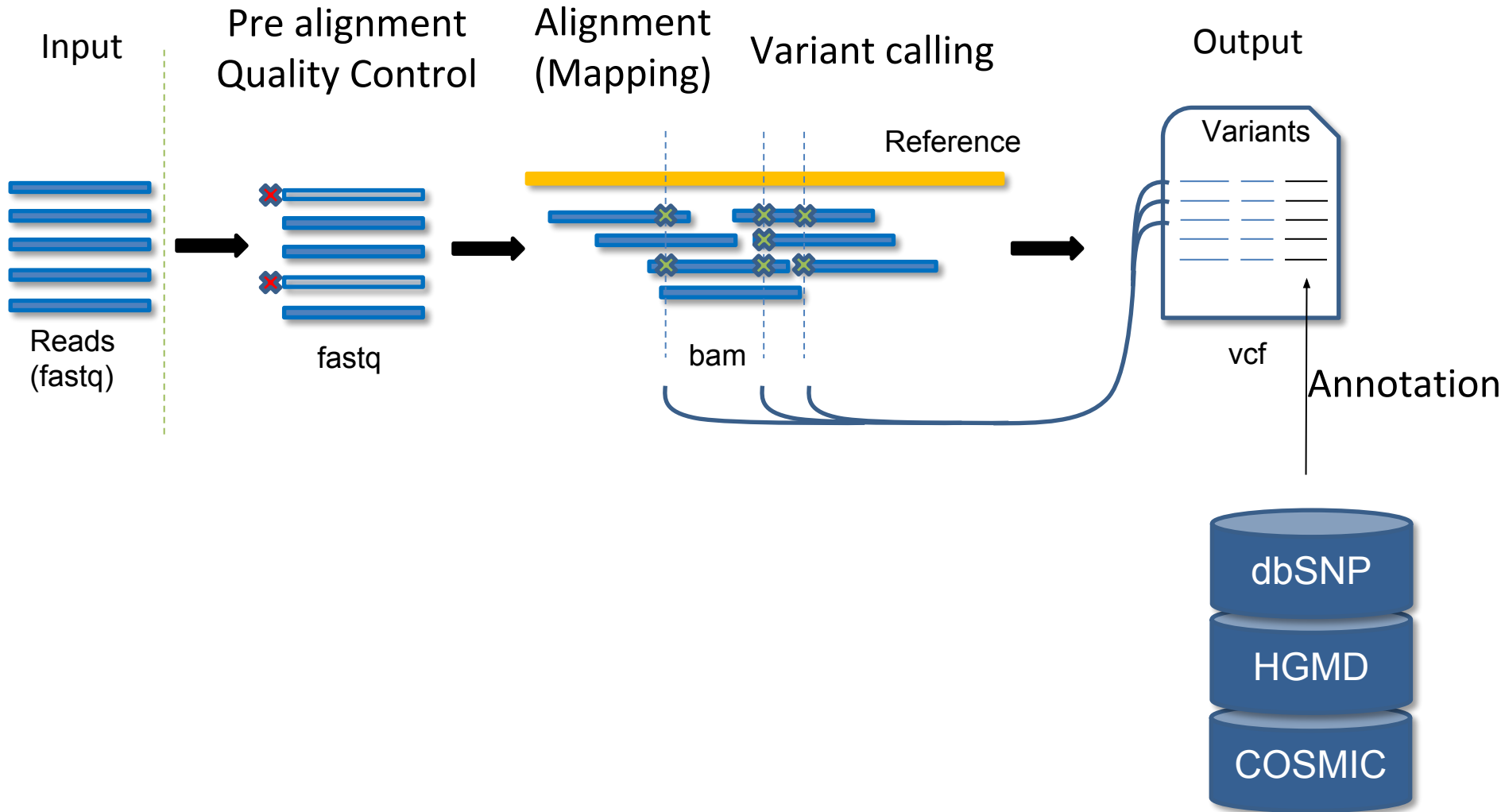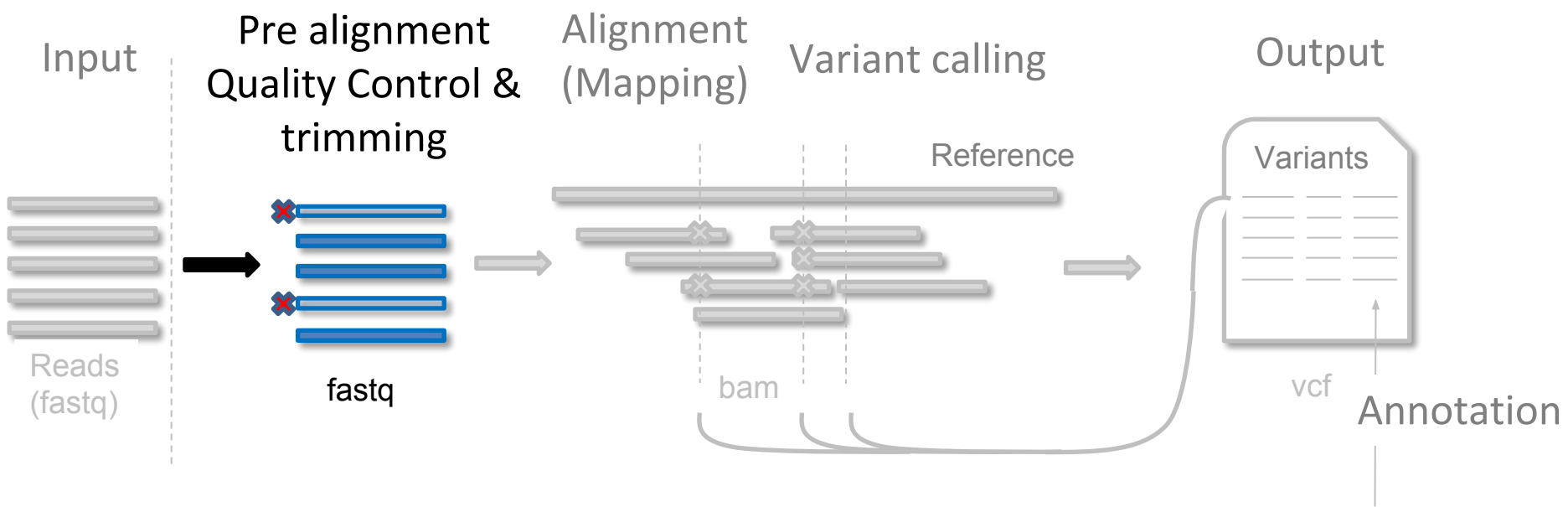(**fastq** file)

# Fastq file

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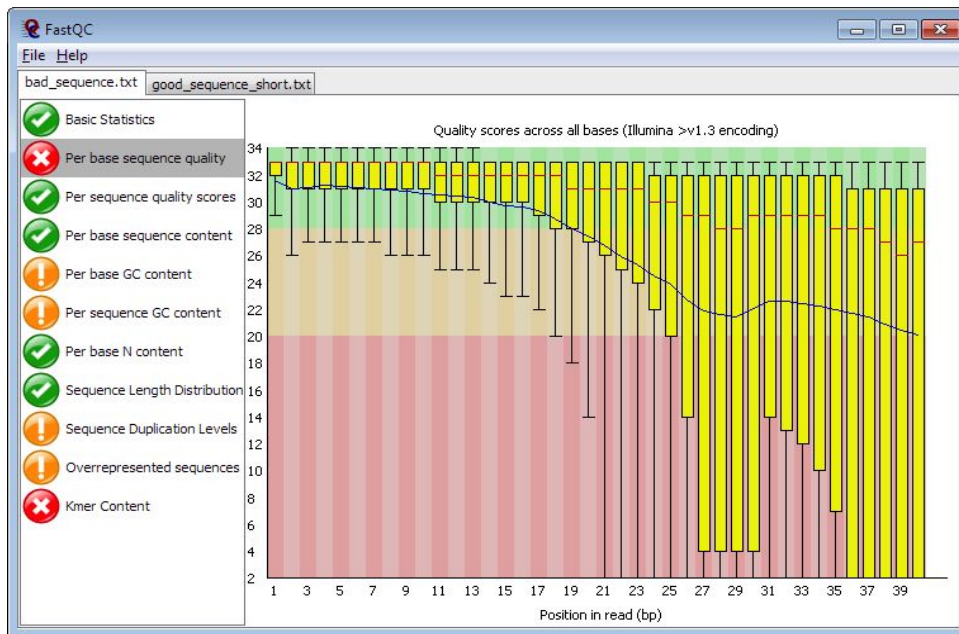- Pair-end sequencing – 2 fastq files

example.fastq
@
SEQ_ID GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT
+
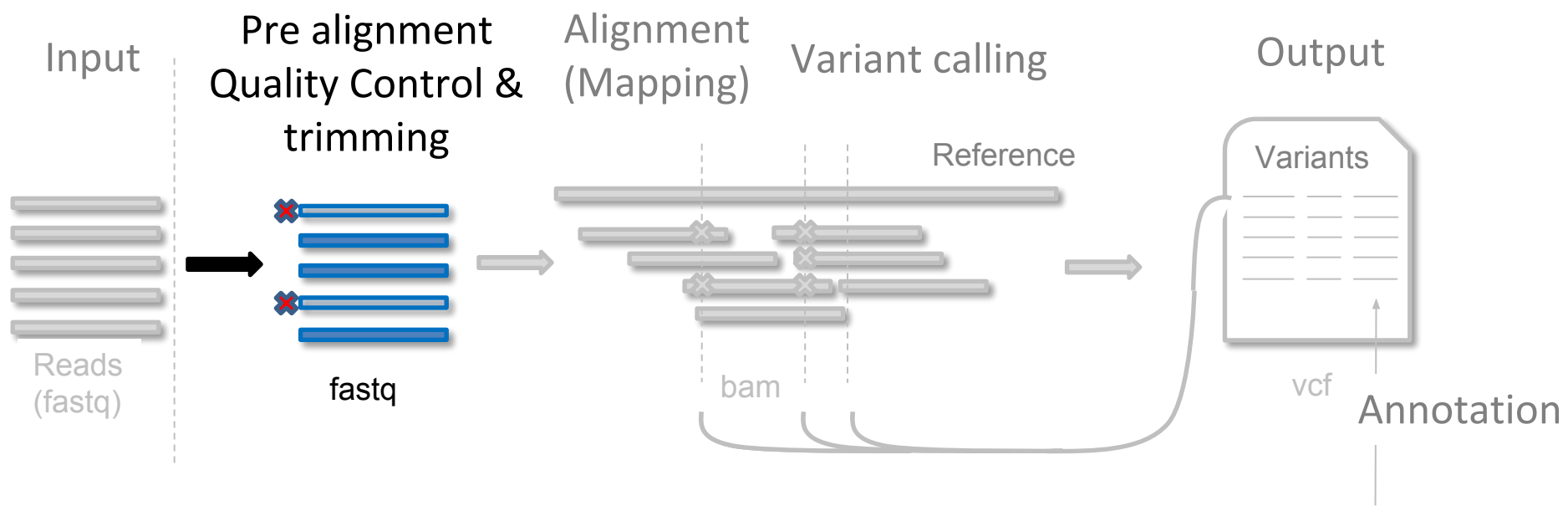 !"*(((((***+))%%%++)(%%%%).1***-+*"))**55CCF>>>>>>CCCCCCC65
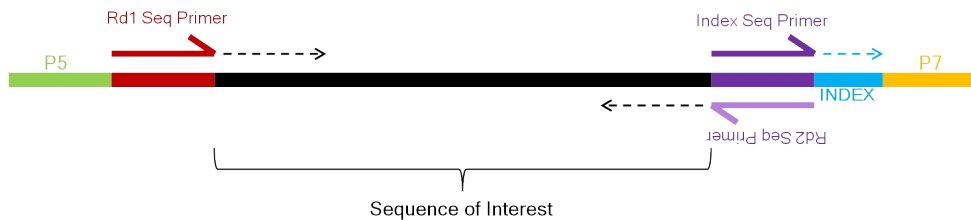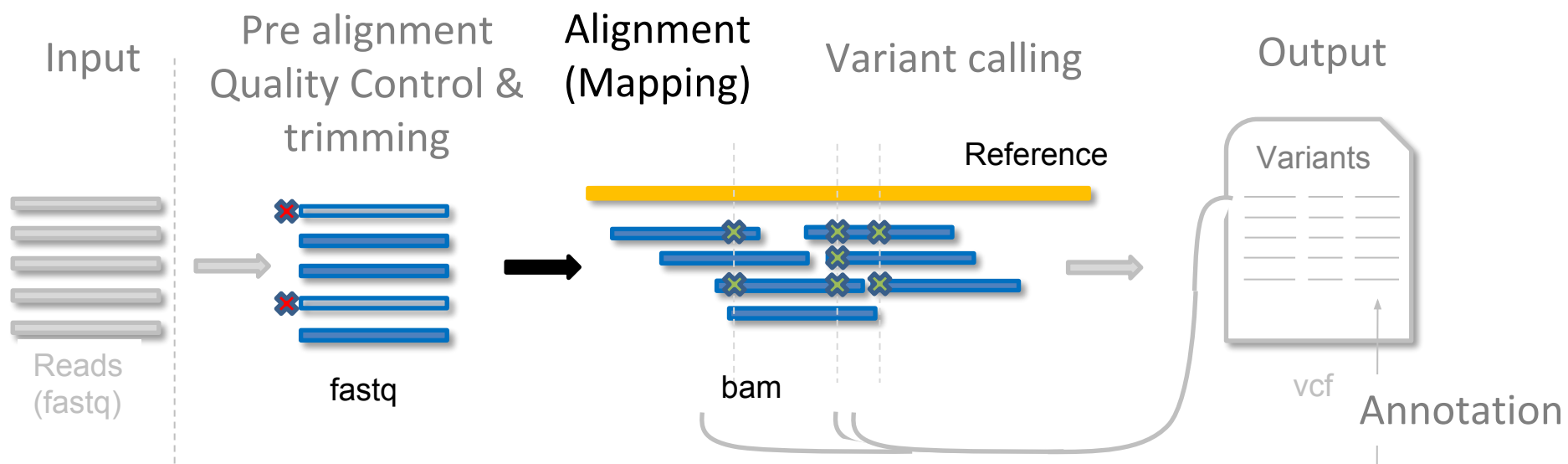
# Input

Pre alignment
Quality Control &
trimming

## Alignment (Mapping)

## Variant calling

# Output

Reads
(fastq)

fastq

Reference

bam

Variants

vcf

Annotation

dbSNP

HGMD

COSMIC

## Quality control (FastQC)

# Input

Reads
(fastq)

# Pre alignment
## Quality Control & trimming

fastq

# Alignment
## (Mapping)

Reference

bam

# Variant calling

# Output

Variants

vcf

Annotation

dbSNP

HGMD

COSMIC

## Cleaning reads (Cutadapt)

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

### STRUCTURE DETAILS

P5

Rd1 Seq Primer

Index Seq Primer

P7

INDEX

Rd2 Seq Primer

Sequence of Interest

Input

Pre alignment
Quality Control &
trimming

Alignment
(Mapping)

Variant calling

Output

Reference

Variants

Reads
(fastq)

fastq

bam

vcf

Annotation

dbSNP

HGMD

COSMIC

•Usually mapping reads on reference sequence
(DNA/cDNA/16S/other seq) to find corresponding
location & differences
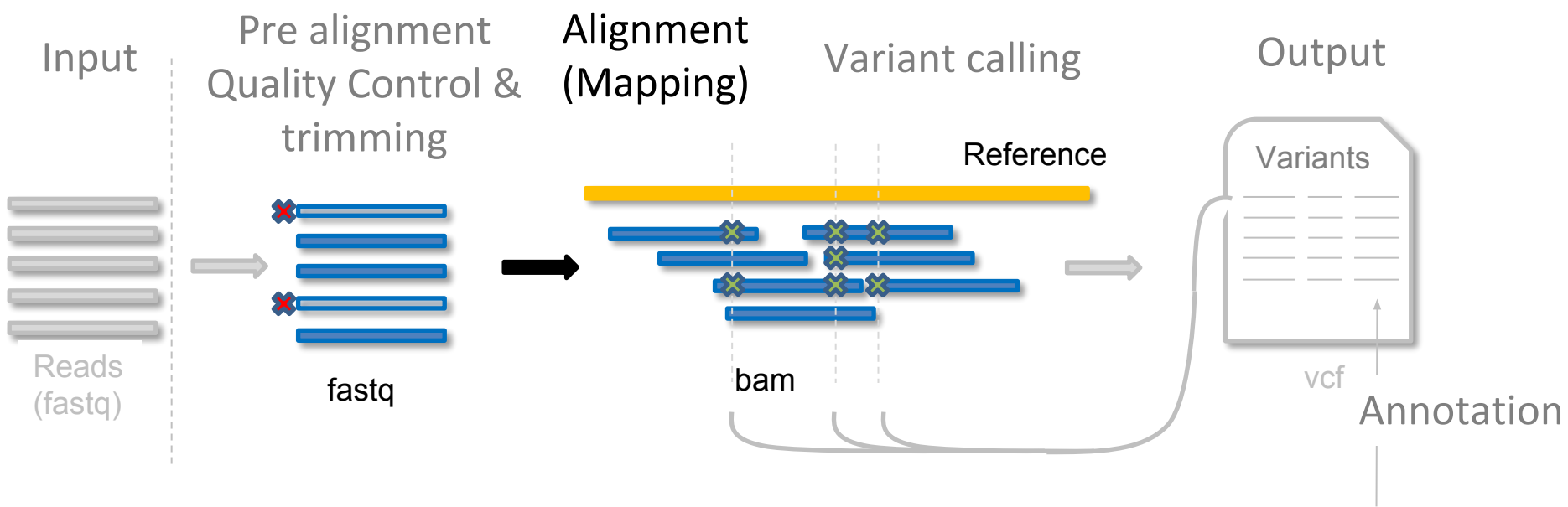(substitutions, insertions, deletions, inversions,
etc… )

•Problem with too many sequences and billions bp
long references – need for special algorithms
(Burrows-Wheeler transform, hash table indexing)
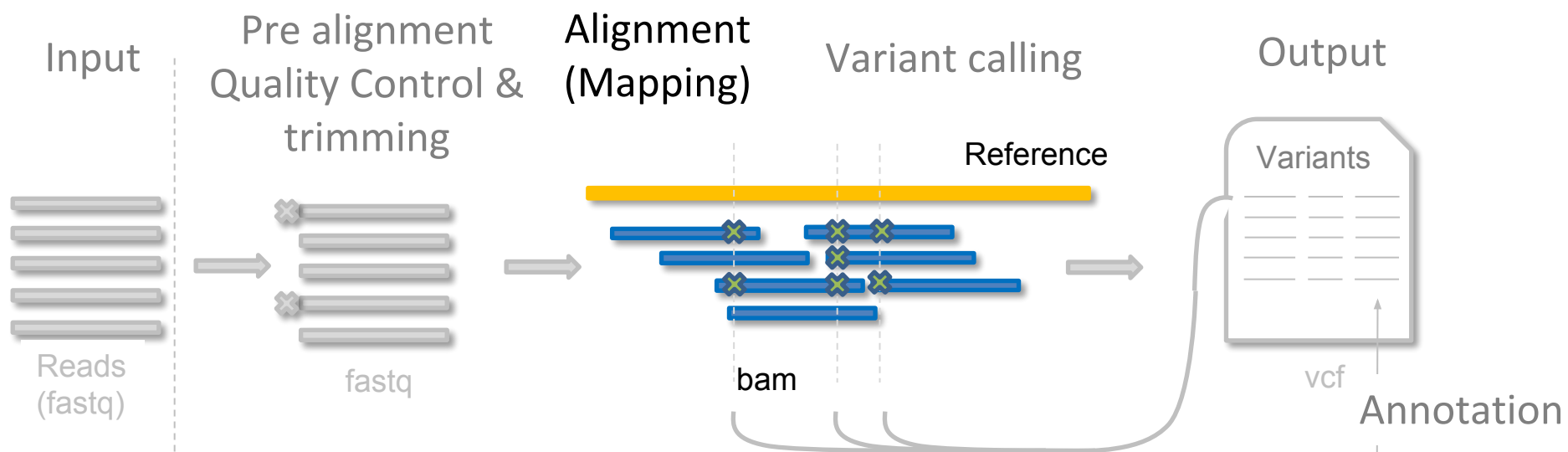
•BWA, Bowtie, Bfast, SHRiMP (BAM format)

# Example of read mapping

Input

Pre alignment
Quality Control &
trimming

Alignment
(Mapping)

Variant calling

Output

Reference

Variants

Reads
(fastq)

fastq

bam

vcf

Annotation

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

dbSNP

HGMD

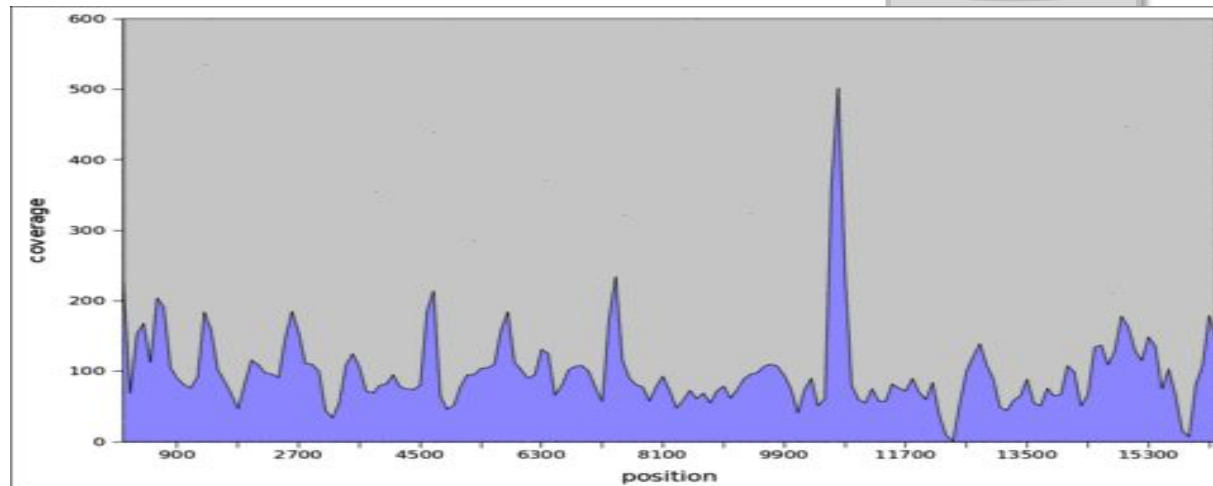COSMIC

Input  Pre alignment Quality Control & trimming  Alignment (Mapping)  Variant calling  Output

Reads (fastq)  fastq  Reference  bam  Variants  vcf  Annotation

dbSNP

HGMD

**Mapping, Coverage reports**

- Repeat alignment/other steps with different criteria?
- Important checkout for lab protocol
- Specificity of PCR
- Settings of variant calling threshold, CNV
- Target bed file (Browser Extensible Data)

chr7  127471196  127472363
chr7  127472363  127473530
chr7  127473530  127474697
chr7  127474697  127475864
chr7  127475864  127477031

(bed format)

# *De novo* assembly – alternative for mapping on reference sequence

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

Input

Pre alignment
Quality Control &
trimming

Alignment
(Mapping)

Variant calling

Output

Reference

Variants

Reads
(fastq)

fastq

bam

vcf

Annotation

dbSNP

HGMD

COSMIC

**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

# Introduction of Molecular barcodes during library preparation



A

B

C

Smith et al. 2014

Input | Pre alignment Quality Control & trimming | Alignment (Mapping) | Variant calling | Output

Reference

Variants

Reads (fastq)

fastq

bam

vcf

Annotation

36,661,660    36,661,680    36,66

**Mutation types:**
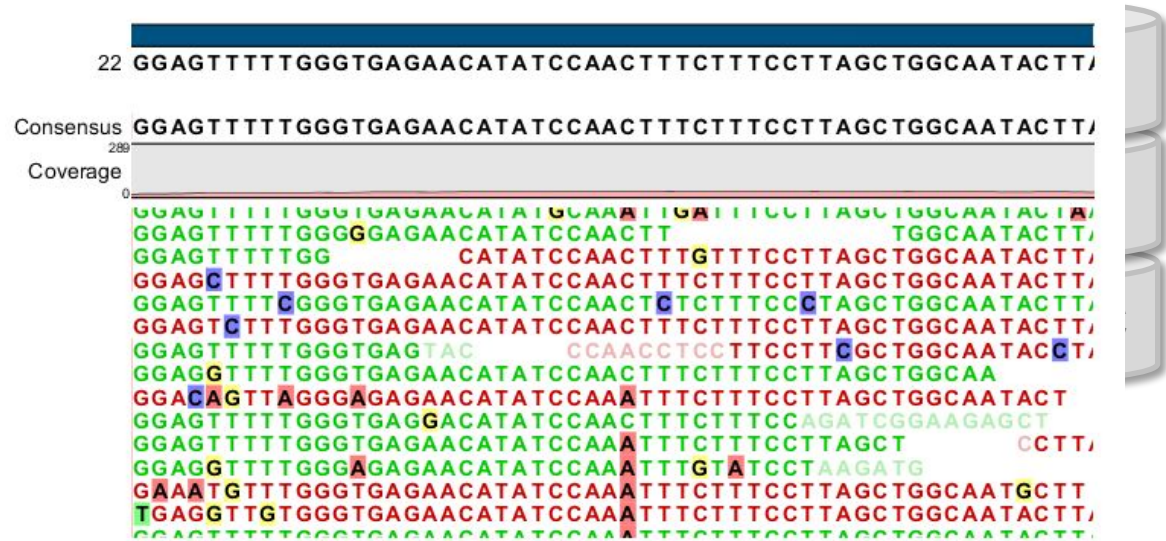Germinal mutations
Somatic mutations

Substitutions
Insertions
Deletions
Complex variants
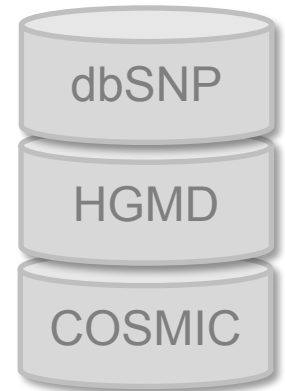
Inversions
Large structural variations (translocations, indels)
Copy number variations

**Experimental designs (also depends on types of samples available):**

Normal only (genotyping)
Tumor only (genotyping, somatic mutations)

Tumor + related normal control
Tumor + unrelated normal controls
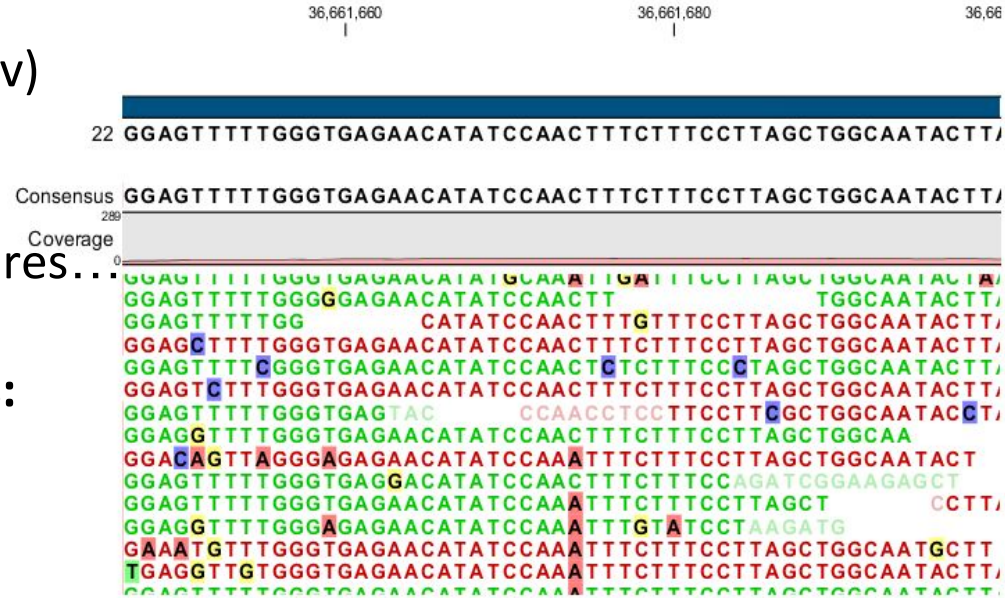Tumor in time

Family (rare diseases, genotyping)

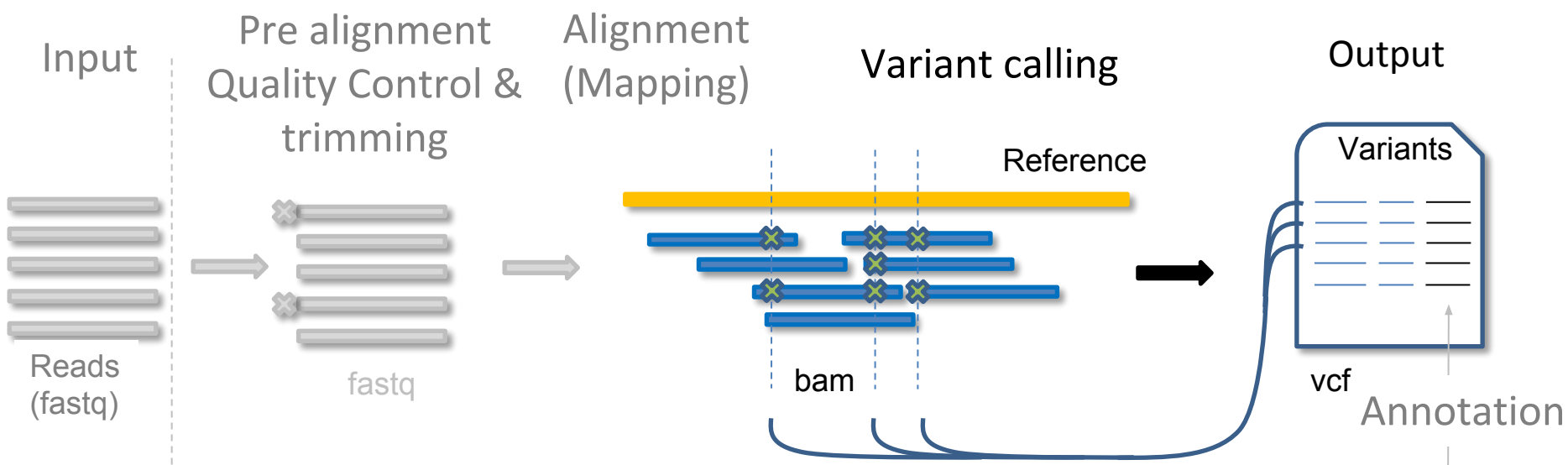Input  Pre alignment Quality Control & trimming  Alignment (Mapping)  Variant calling  Output



Reference

Reads (fastq)  fastq  bam  vcf  Variants  Annotation

**Program algorithms:**

- Bayesian statistics (Mutect, DeepSnv)
- Fisher exact test (Varscan, Vardict)
- …

Giving p-value based on different features...

**Options for many parameters & filters:**

- Minimum coverage
- Variant allele frequency
- Base quality
- Genomic context (homopolymers)
- Position in read (errors at the reads end)
- Mapping quality
- Presence in both forward and reverse reads (strand bias)

**To distinguish real mutation from ERROR**
(library preparation, sequencing, alignment)

Usually 1 approach is not enough =>
to combine more variant callers (aligners) & different settings

Specific pipeline for each type of mutations
(SNV, INDELS, CNV…)

O'Rawe, J. *et al.* Low concordance of multiple variant-calling pipelines: practical implications for exome and genome sequencing. *Genome Medicine* 5, 28 (2013).

REALTIME GENOMICS

Stanford University

# Input · Pre alignment Quality Control & trimming · Alignment (Mapping) · Variant calling · Output

**Input** — Reads (fastq)

**Pre alignment Quality Control & trimming** — fastq

**Alignment (Mapping)** — Reference — bam

**Variant calling**

**Output** — Variants — vcf — Annotation

dbSNP

HGMD

COSMIC

## Visualization of genotypes by IGV

Input vcf



EUR.UMICH.201...otypes.vcf.gz

NA06984
NA06986
NA06989
NA06994
NA07000
NA07037
NA07048
NA07051
NA07056
NA07346
NA07347
NA07357
NA10847
NA10851
NA11829
NA11830
NA11831
NA11832
NA11840
NA11843
NA11881
NA11892
NA11893
NA11894
NA11918
NA11919
NA11920
NA11930

1) Each bar across the top of the plot shows the allele fraction for a single locus.
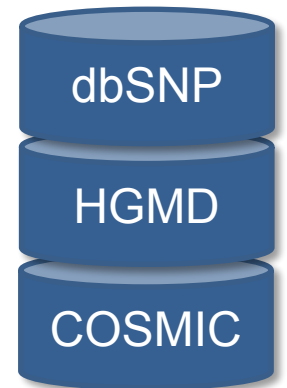2) The genotypes for each locus in each sample. Dark blue = heterozygous, Cyan = homozygous variant, Grey = reference.
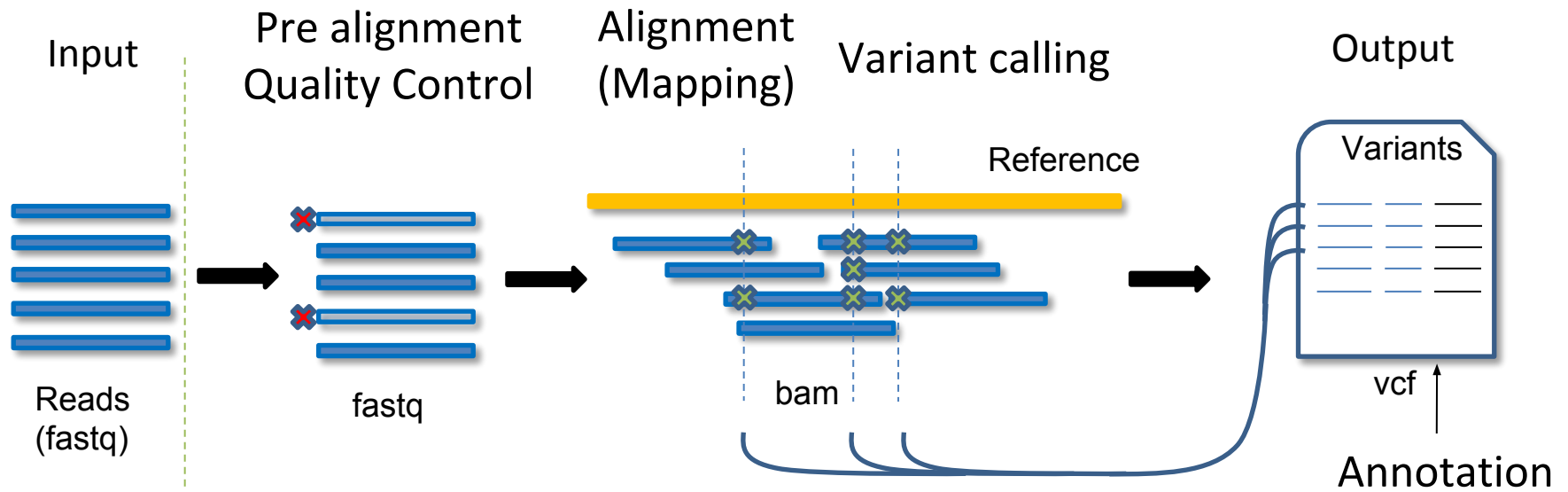
**Annotation**
**From genomic coordinate to biological meaning**
Provide links to various databases (RefSeq, dbSNP, etc.)
To distinguish significant variant from non-significant
(synonymous vs. non-synonymous, gene, exon, intron, cDNA, codon,
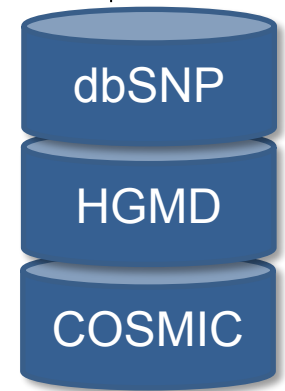transcript, freq in population, presence in other diseases…)
- RefSeq
- dbSNP
- Regulation
- Comparative genomics
- Repeats
- Functional
- Gene ontology
- Etc.

**Sensitivity & Specificity as a matter of:**

- Experiment design

(library preparation + NGS technology + number of samples
+ amount of data)

- Data processing

(pre-processing + alignment + variant calling
+ annotations + filtering)

# Courses

http://meetings.embo.org/event/17-genome

http://www.embo.org/events/practical-courses