

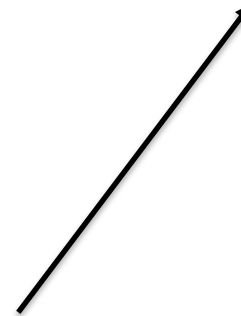
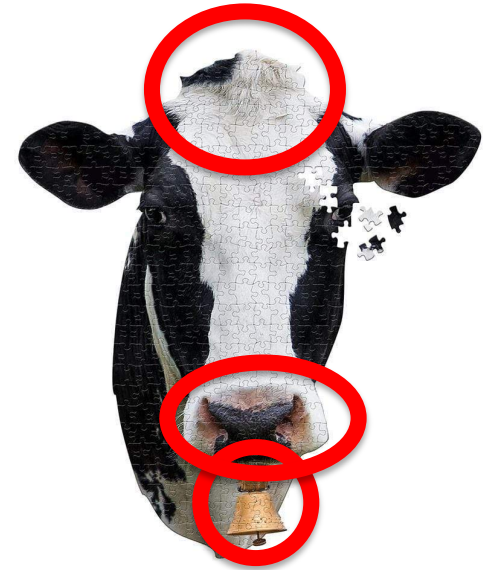
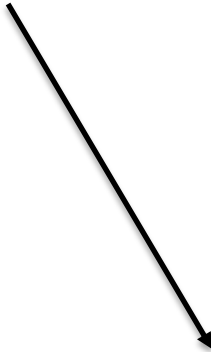
Introduction to Whole Genome Sequencing

From FASTQ to VARIANTS



Rodrigo Pelicioni Savegnago
Department of Animal Science

From Fastq to VCF file



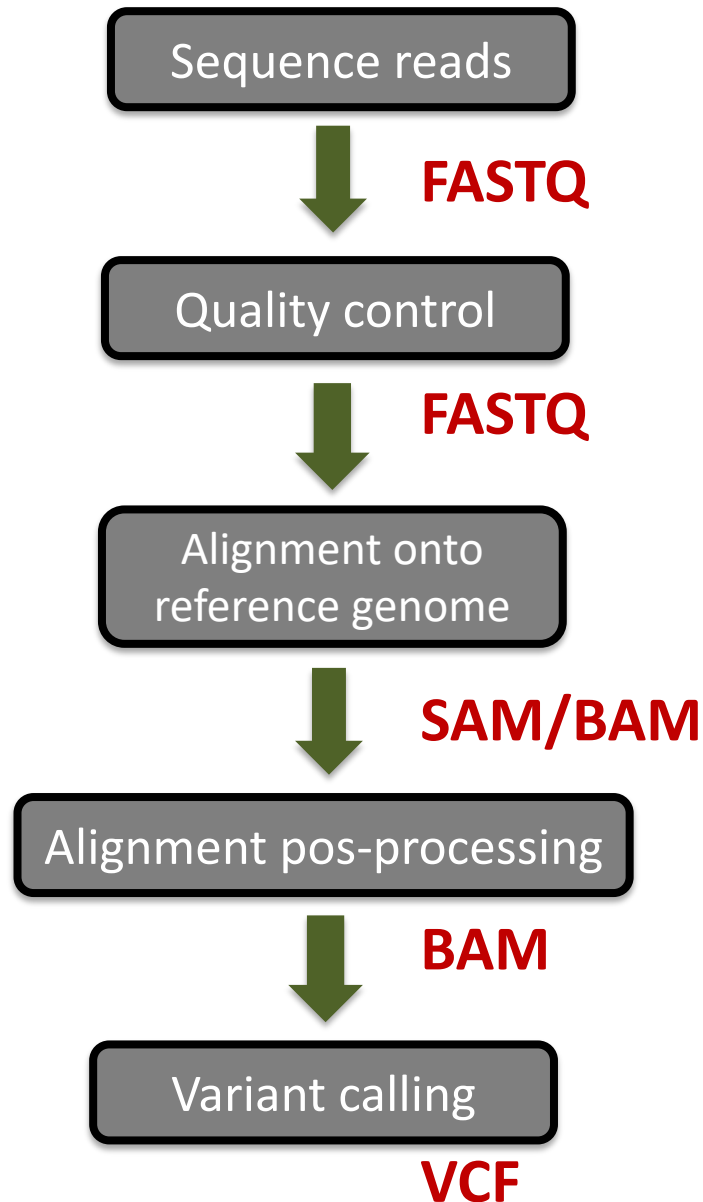
Hum Genet (2012) 131:1541–1554
DOI 10.1007/s00439-012-1213-z

REVIEW PAPER

A beginners guide to SNP calling from high-throughput DNA-sequencing data

**André Altmann · Peter Weber · Daniel Bader ·
Michael Preuß · Elisabeth B. Binder ·
Bertram Müller-Myhsok**

Received: 31 January 2012 / Accepted: 31 July 2012 / Published online: 11 August 2012
© Springer-Verlag 2012



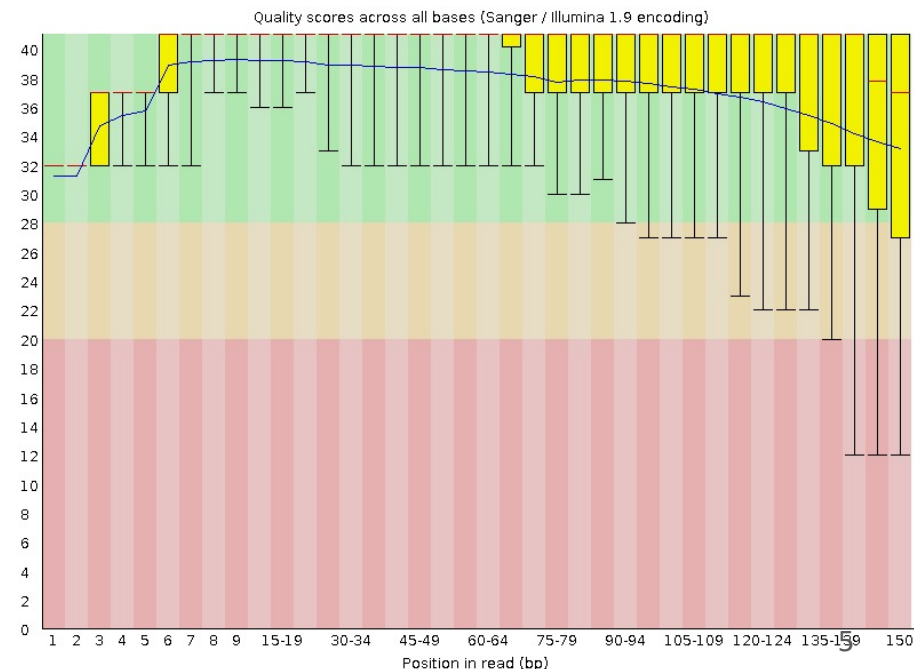
- **Initial visual appraisal**

- Are all files present?
- Are the files of expected size?
- Check the names of FASTQ files

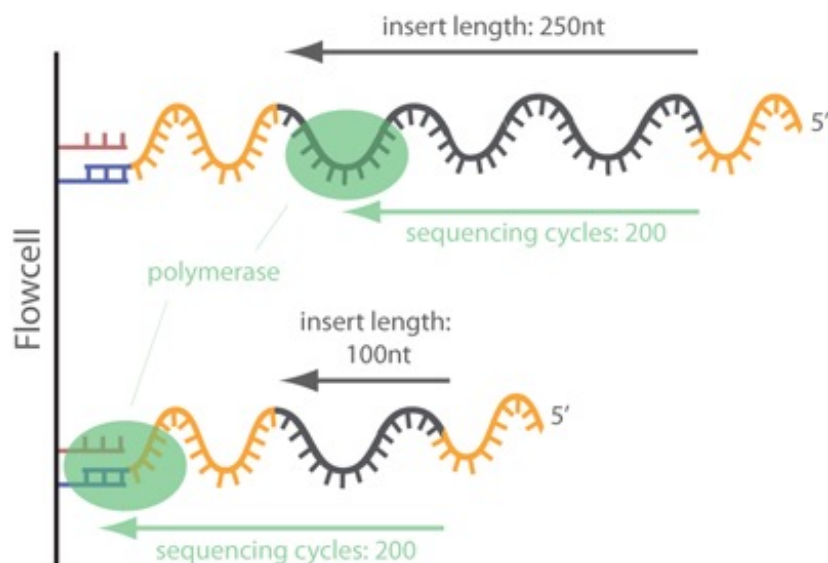
- **Other quality controls**

- Phred score distribution
- GC content

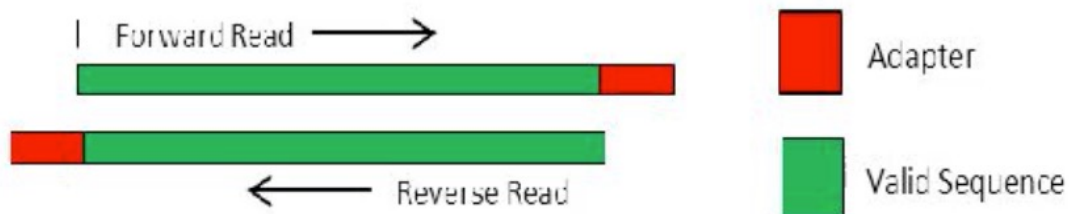
- **FASTQC**

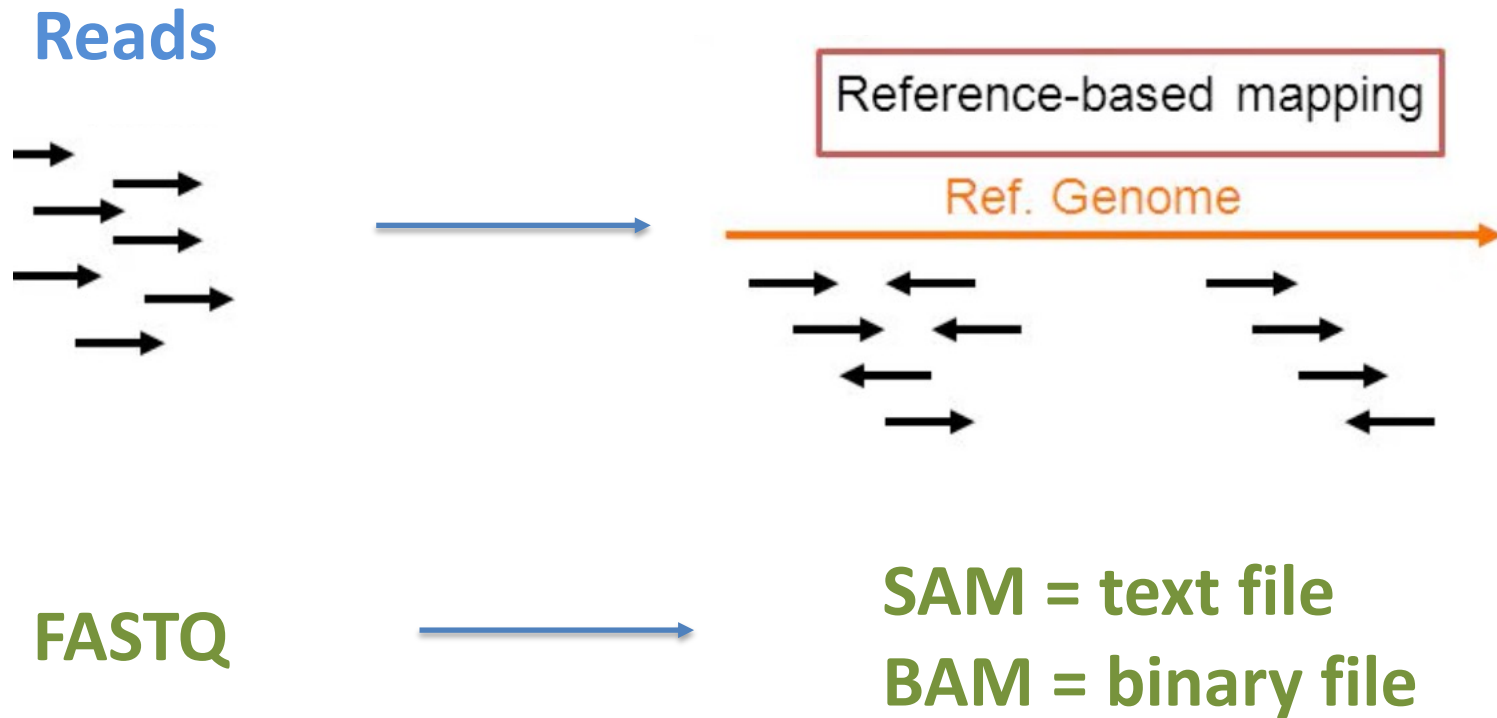


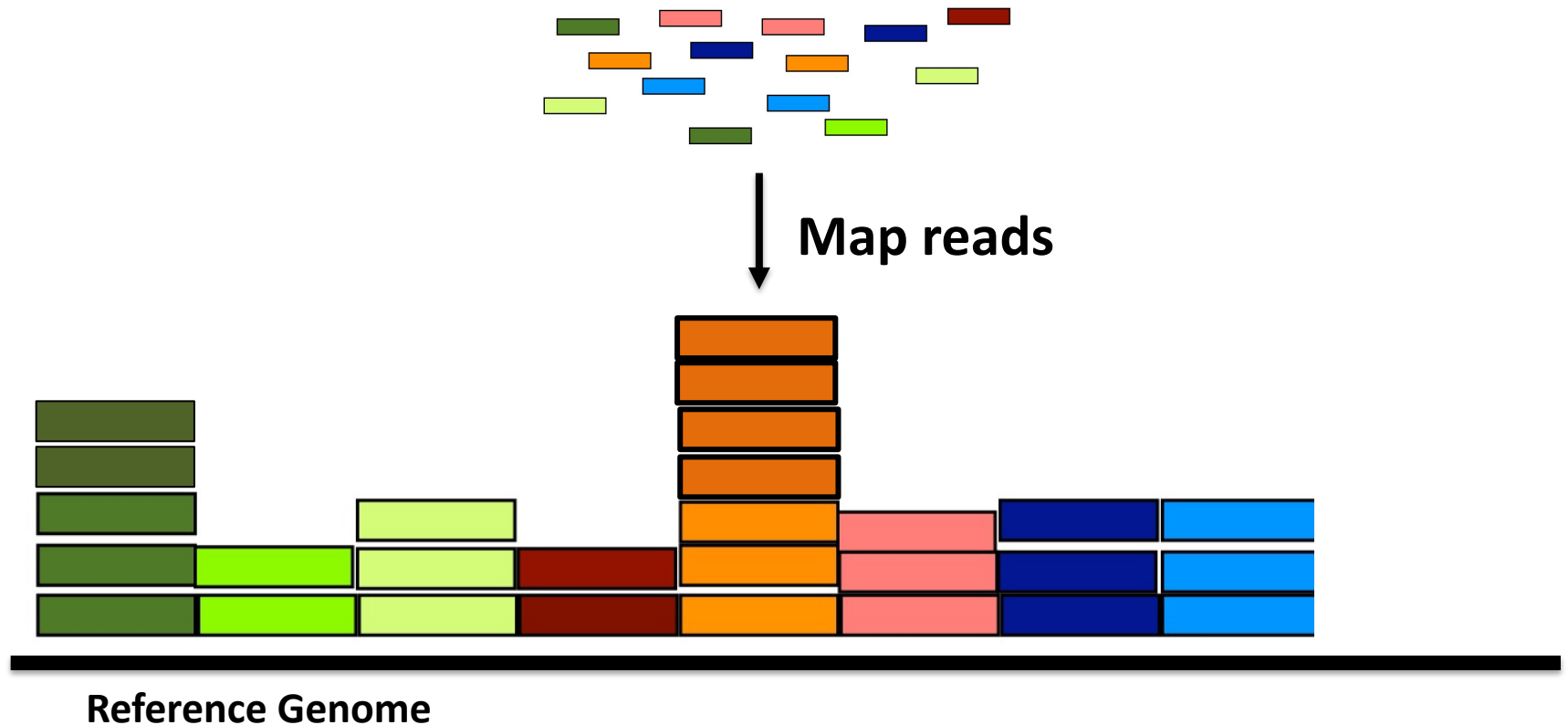
- **Trimmomatic** is one option



- Clip Illumina **adapters**:



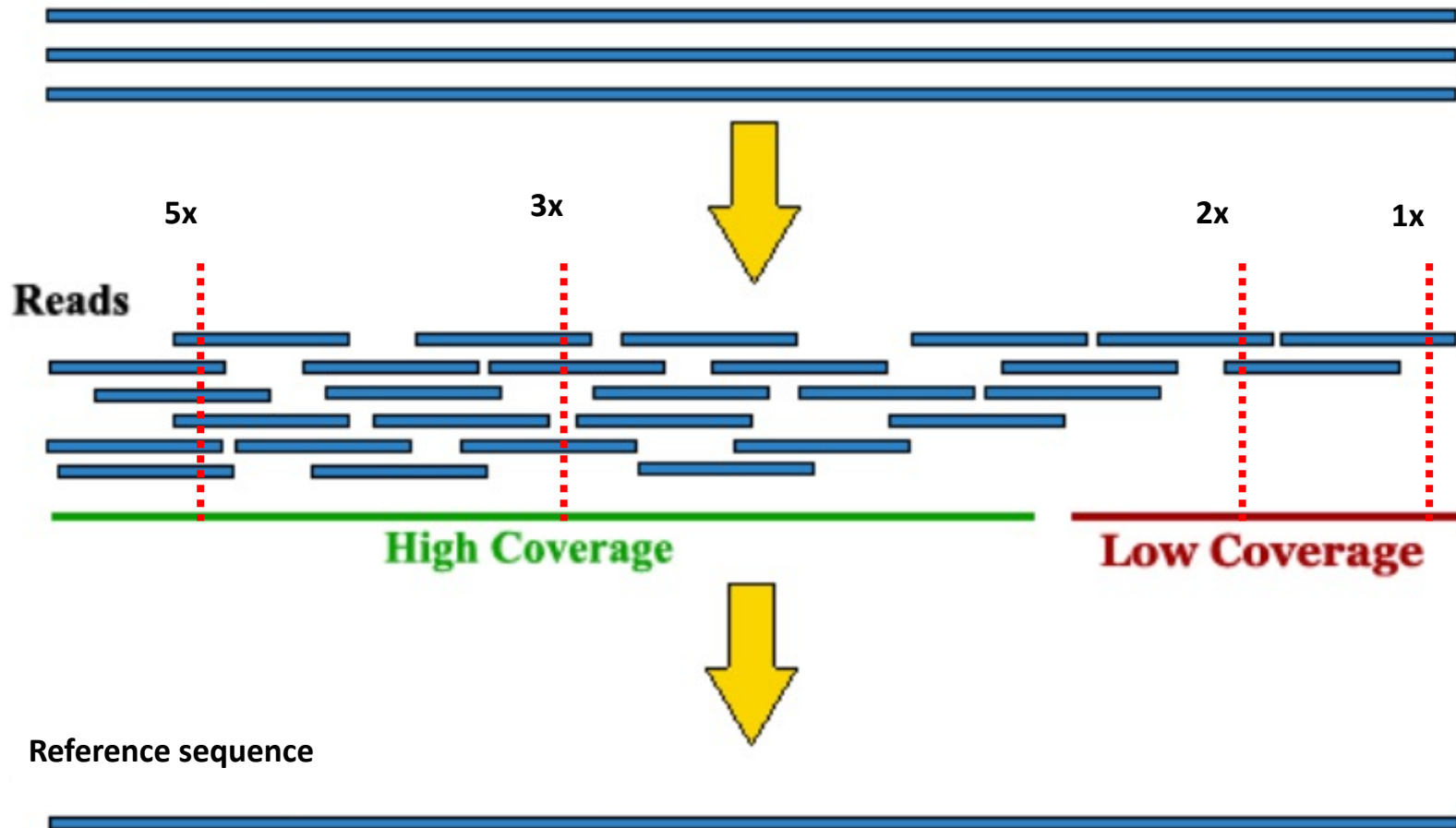




Seesi et al (2016). Genomics-Guided Immunotherapy of Human Epithelial Ovarian Cancer. In Rodriguez-Oquendo (Eds.) **Translational Cardiometabolic Genomic Medicine**. pp. 237-250.

03. Read mapping: coverage depth

Multiple Copies of a Genome



Reference sequence



$$G = 8''$$



$$\text{length} = 1''$$

$$N = 28 \text{ reads}$$

$$\text{Cov Depth} = \frac{\text{length} * N}{G} = \frac{1 * 28}{8} = 3.5x$$

ALIGNERS

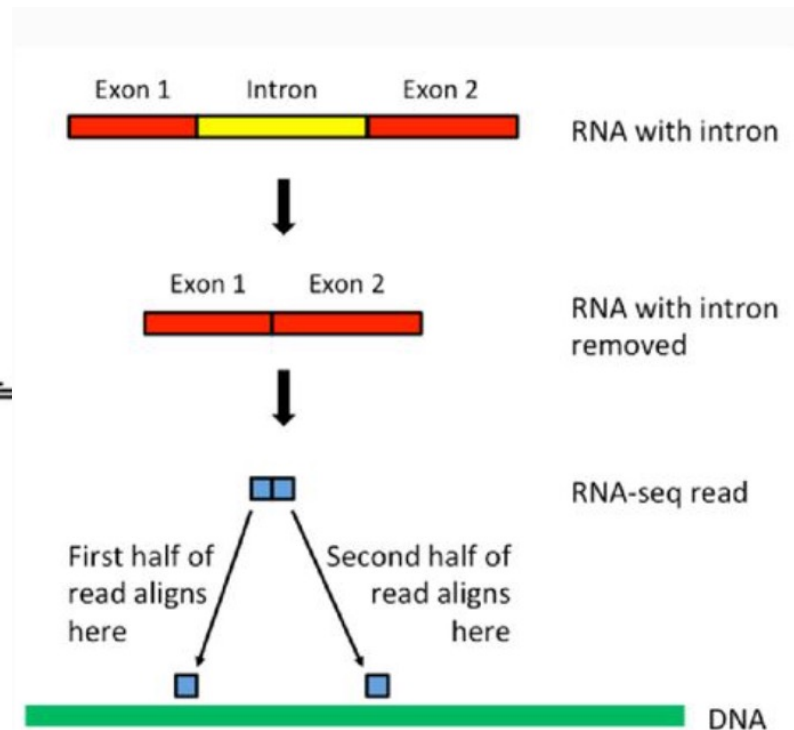
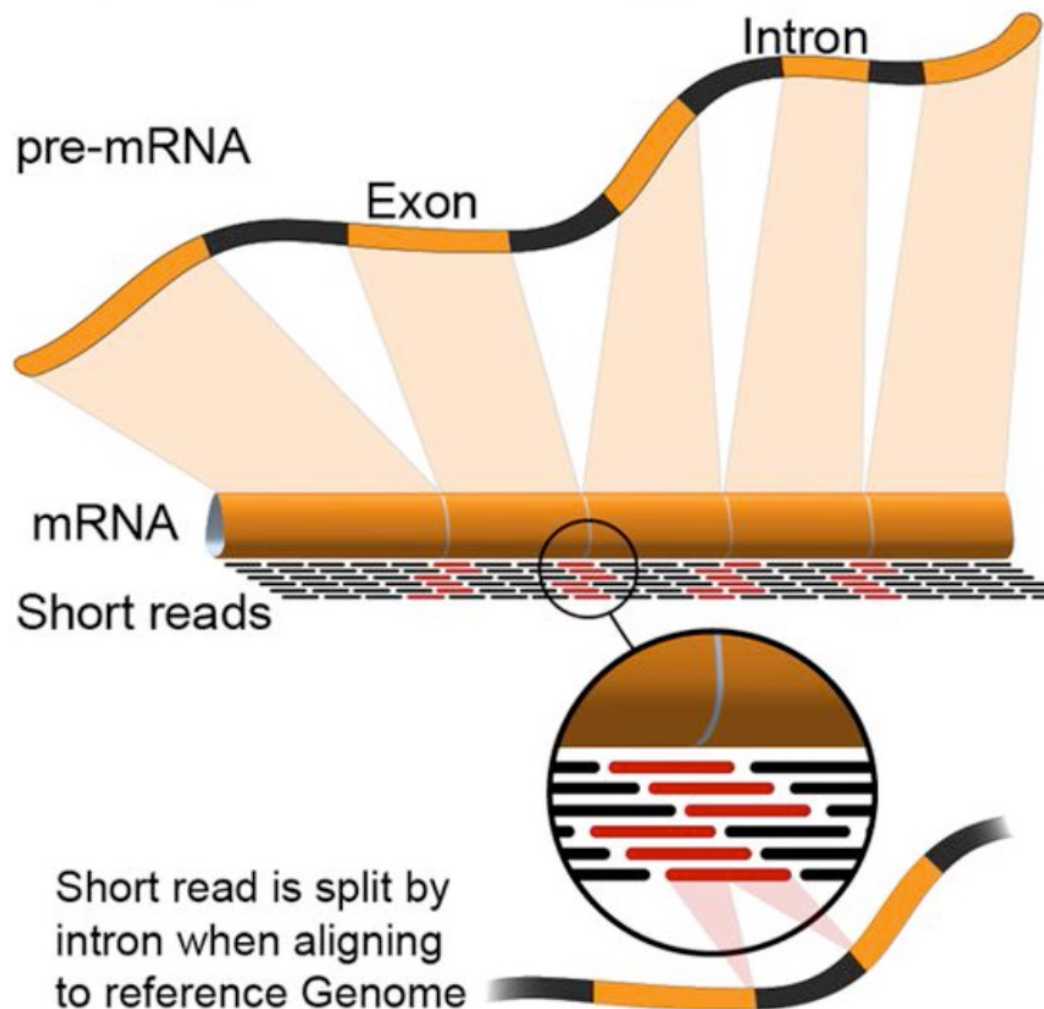
WGS

- BWA
- BOWTIE
- BOWTIE2

RNAseq

- TOPHAT
- HISAT2
- STAR

03. Read mapping: splice aware aligners

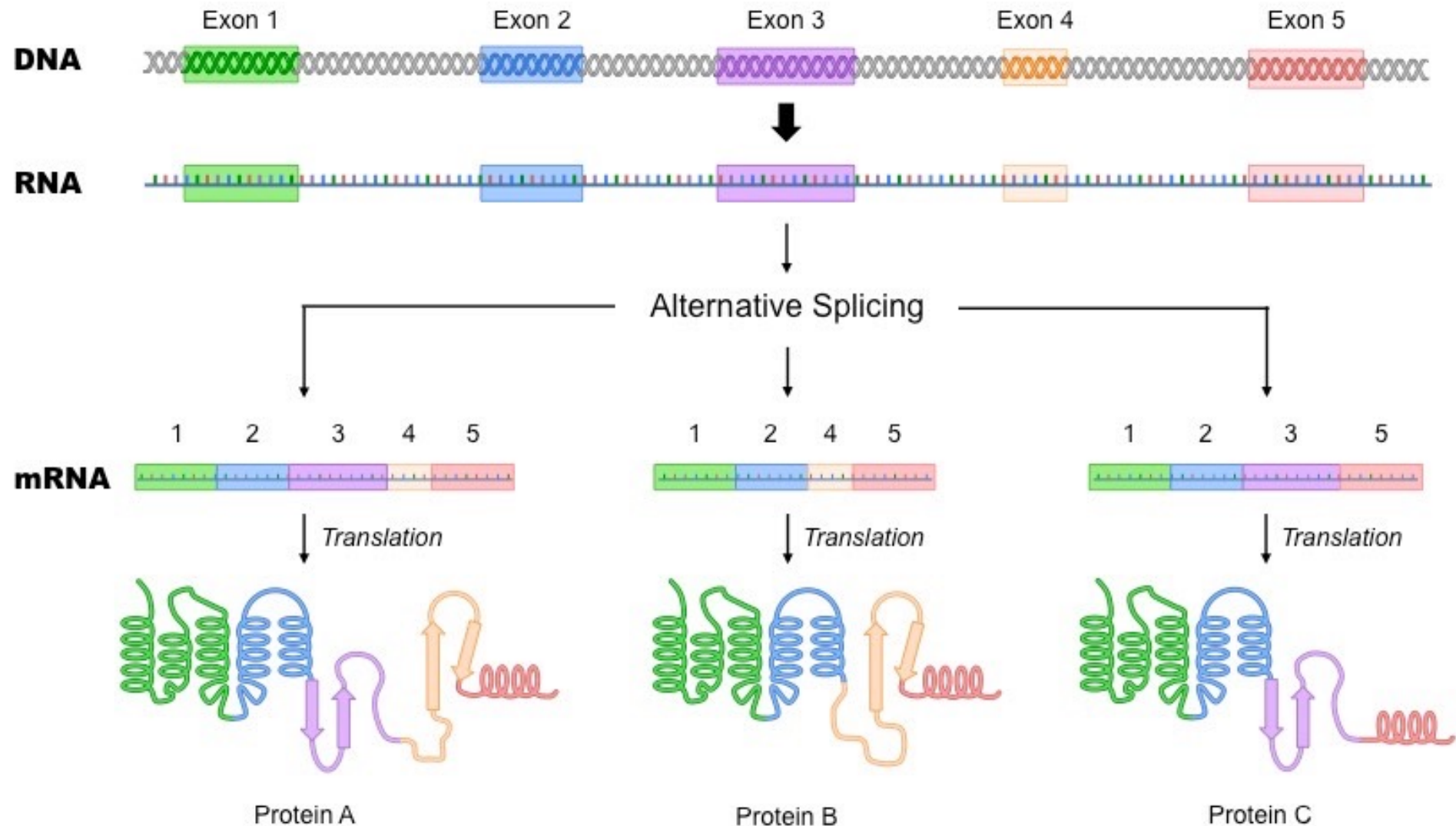


<http://en.wikipedia.org/wiki/File:RNA-Seq-alignment.png>

"Systematic evaluation of spliced alignment programs for RNA-seq data"

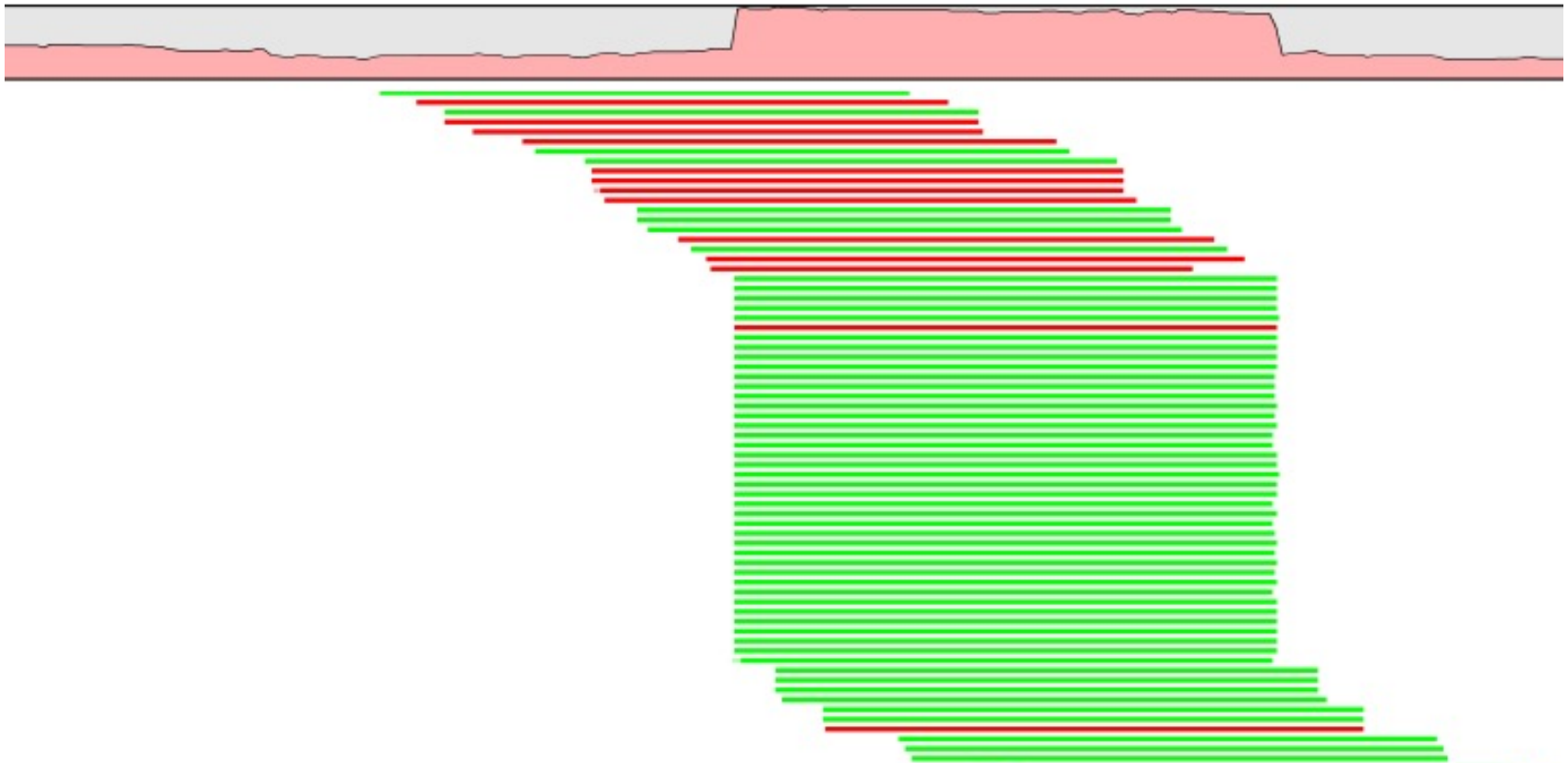
Nature Methods, 2013

03. Read mapping: RNA alternative splicing

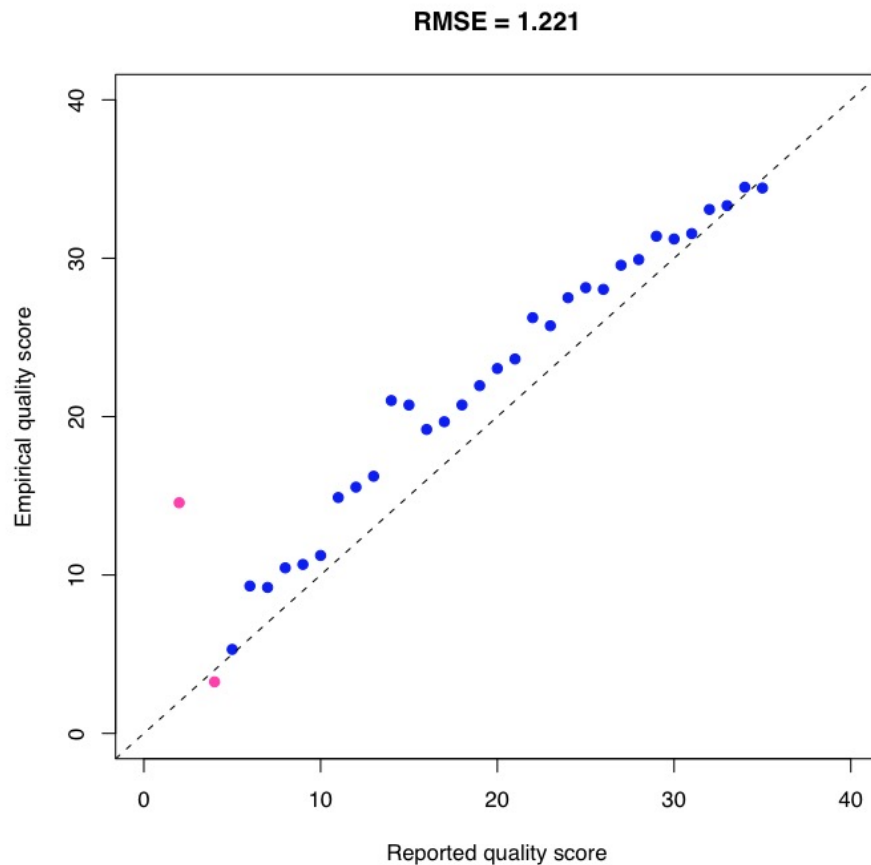


PCR amplification bias

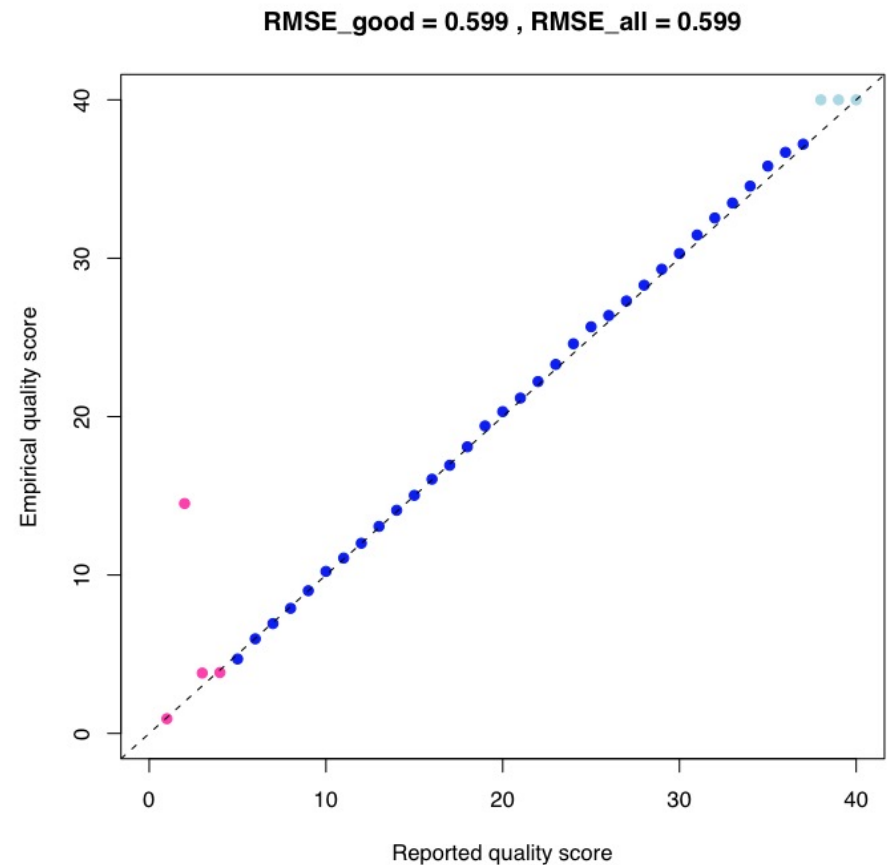
- Some reads are better amplified than others → bias!
- Keep only one (with highest mapping quality)



05. BQSR – Reported vs Empirical Quality

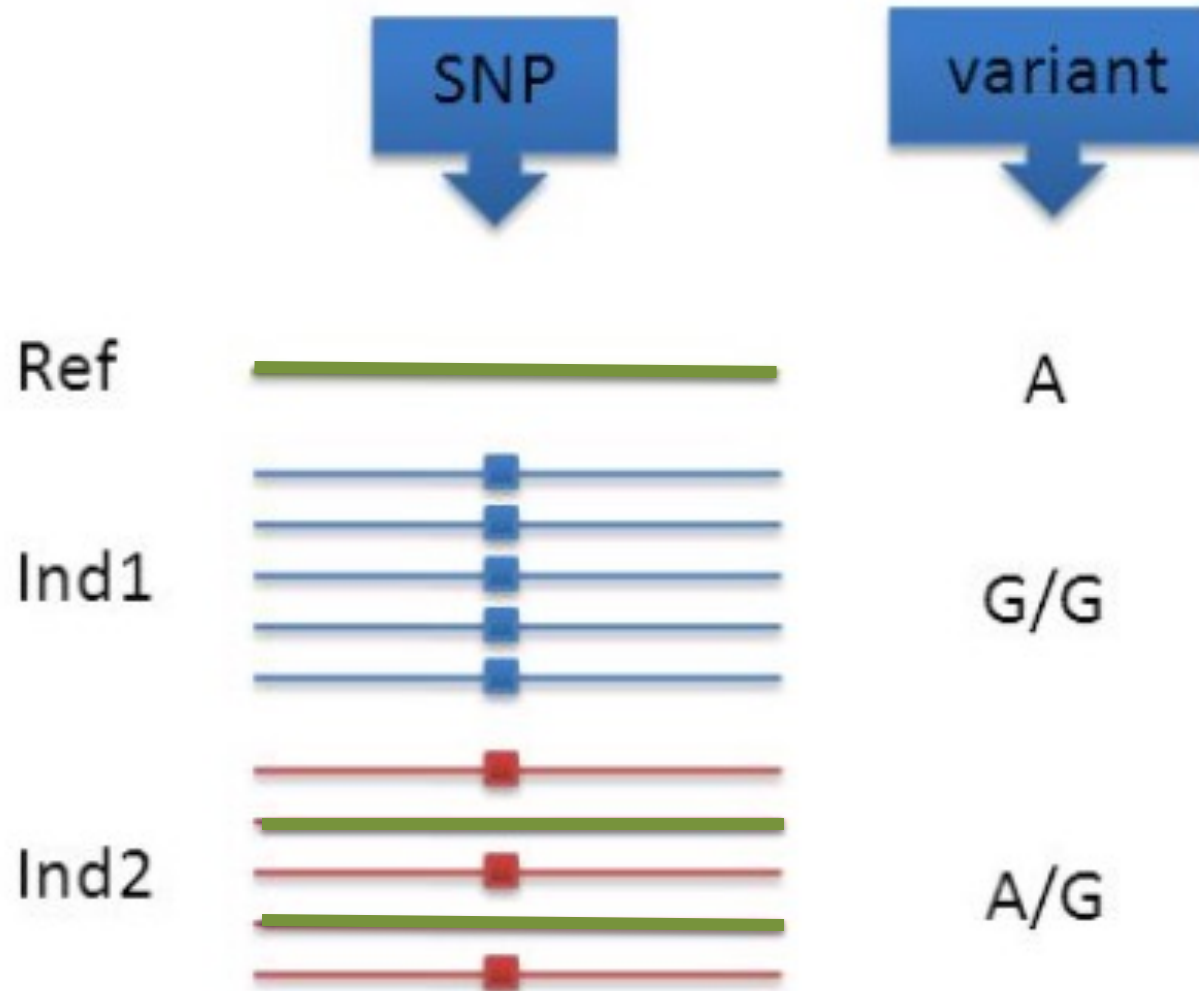


Original Data



After GATK Recalibration

- SAMTOOLS mpileup + BCFTOOLS call (deprecated)
 - BCFTOOLS mpileup + BCFTOOLS call
 - FREEBAYES
 - GATK
 - VARSCAN2
 - DeepVariant
- ... others



Coverage is important!

Assuming a true heterozygous variant C/T (expected 50% C and 50% T)

Reference

... GTGCCAGGACCAGATCG ...

GTGCCAGGACCAGATCG	@1X	50% chance not to sample the T-allele
GTGCCAGGACCAGATCG	@2X	25%
GTGCCAGGACCAGATCG	@3X	12.5%
GTGCCAGGACCAGATCG	@4X	6.25%
GTGCCAGGACCAGATCG	@5X	3.125%
	...	
	@15X	0.003%

Coverage is important!

Assuming a true heterozygous variant C/T (expected 50% C and 50% T)

Increase the coverage will increase detection of all alleles

(Assuming a unbiased approach which is not necessarily the case:
PCR amplification errors, reference errors, mapping errors.

- Step-by-step from FASTQ to VCF
- Specific bioinformatic programs
- Linux language