

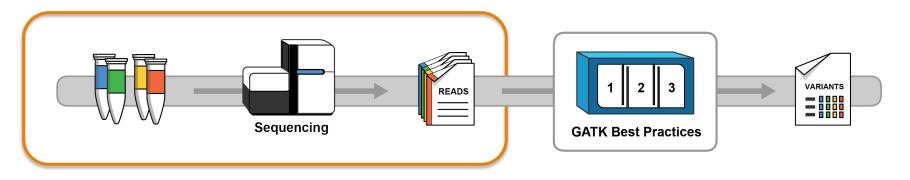
# GATK Best Practices for Variant Discovery

# Introduction to High-Throughput Sequence Data

How it is generated and how we process it



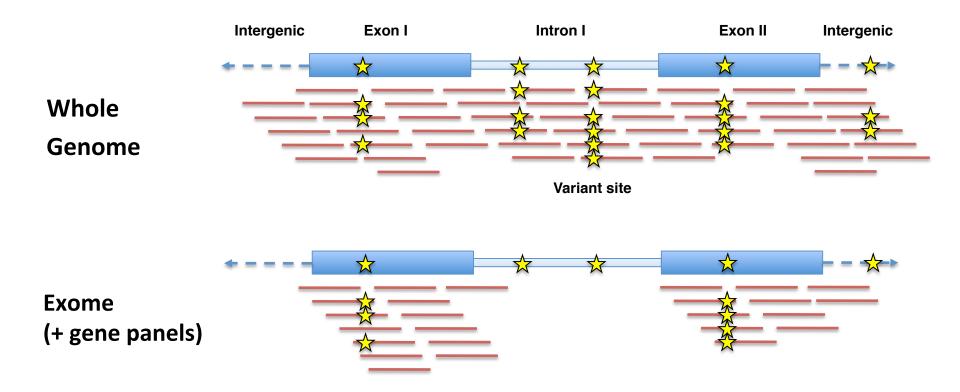




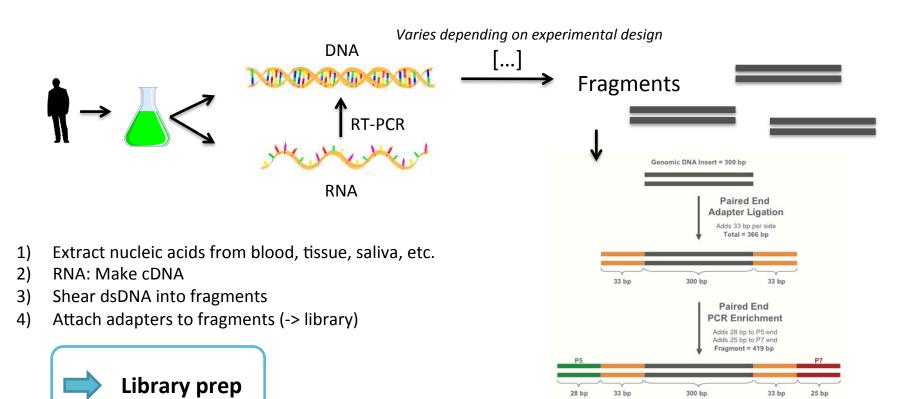
Part 1

# DATA GENERATION

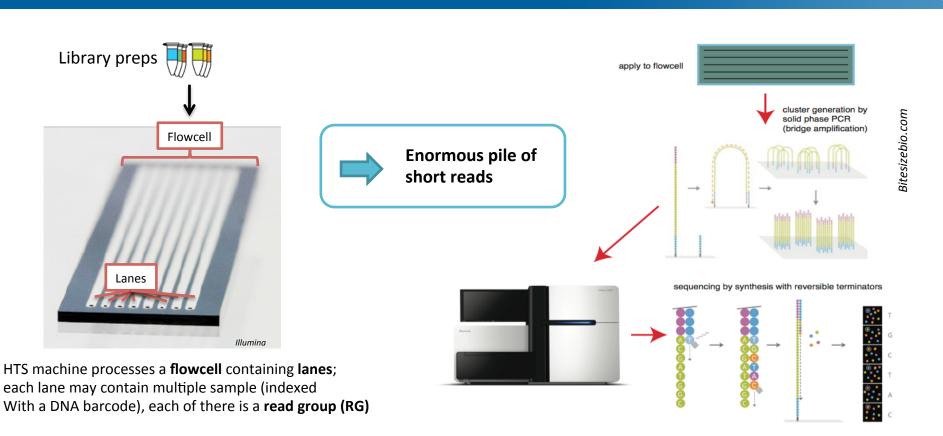
## Different types of experimental design



## Library preparation



## Sequence the library



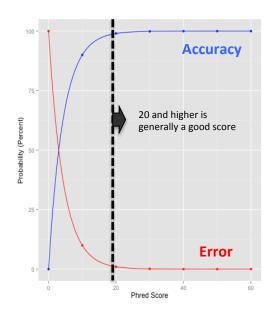
## Raw sequence: typically in FASTQ format

- Sequence Name (read name, group, etc.)
- Sequence
- + (optional: Sequence name again)
- Associated quality score

#### **Example record**

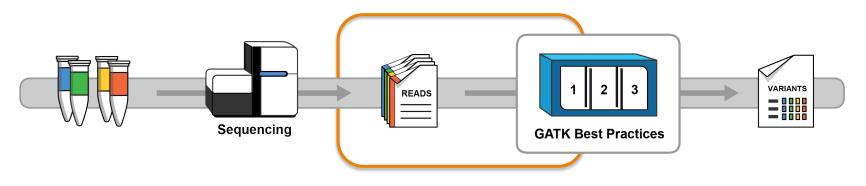
- @EAS54 6 R1 2 1 413 324
- CCCTTCTTGTCTTCAGCGTTTCTCC
- +
- · ;;3;;;;;;;;;;7;;;;88

#### -> ASCII code translates to Phred-scale Q scores



Phred value =  $-10 * log_{10}(\epsilon)$ 

90% confidence (10% error rate) = Q10 99% confidence (1% error rate) = Q20 99.9% confidence (0.1% error rate) = Q30



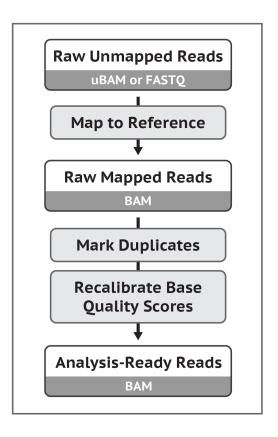
Part 2

## **DATA PRE-PROCESSING**

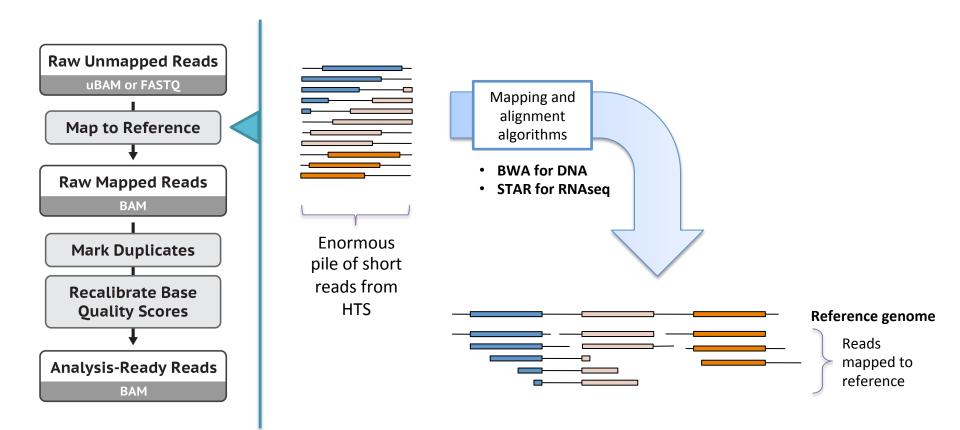
## What is data pre-processing?

#### Data produced by the sequencers:

- A huge pile of paired reads without mapping information
- Afflicted by various technical biases and artifacts
- May include artificial duplicates reads from the same original molecule
- Need to map, sort and clean up

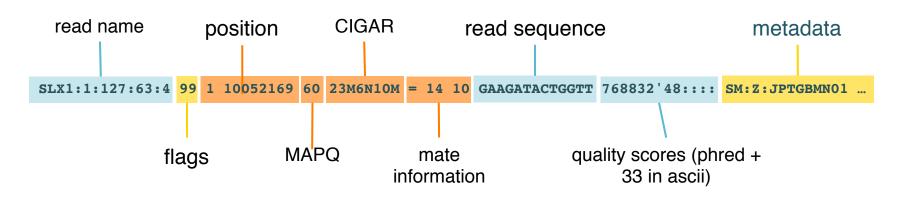


### **Step 1:** Map the reads produced by the sequencer to the reference



## Output format: Sequence/Binary Alignment Map (SAM/BAM)

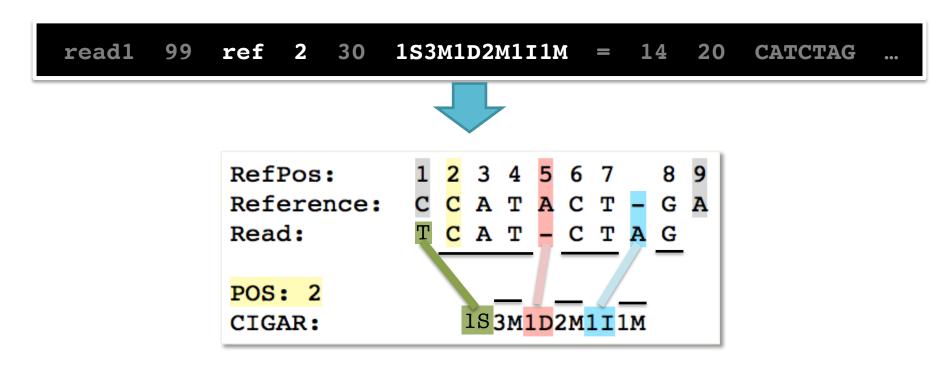
**HEADER** containing metadata (sequence dictionary, read group definitions, etc.) **RECORDS** containing structured read information (1 line per read record)



- Added mapping info summarizes position, quality, and structure for each read
- Mate information points to the read from the other end of the molecule

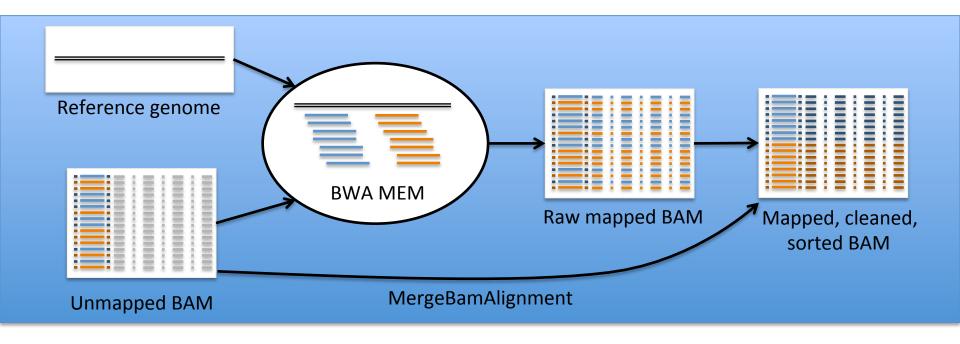
## **CIGAR** summarizes alignment structure

**CIGAR = Concise Idiosyncratic Gapped Alignment Report** 

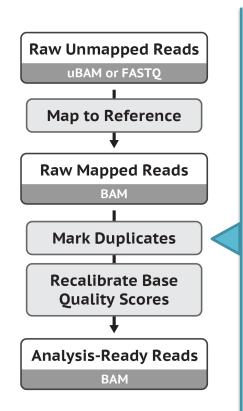


## At Broad: Unmapped BAM instead of FASTQ

Special workflow using Picard tools for improved data management

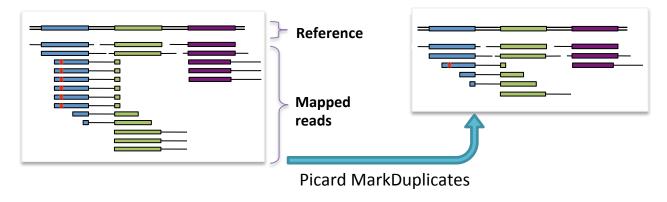


## **Step 2:** Mark duplicates to mitigate duplication artifacts



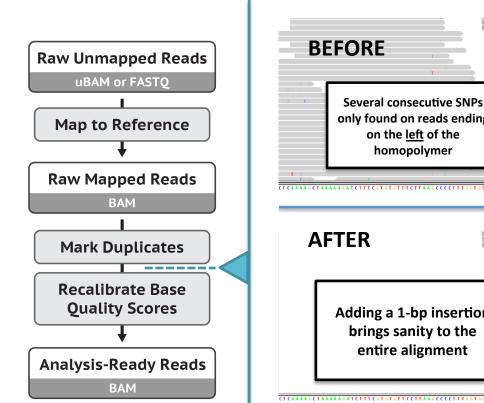
Duplicates = **non-independent measurements**of a sequence fragment

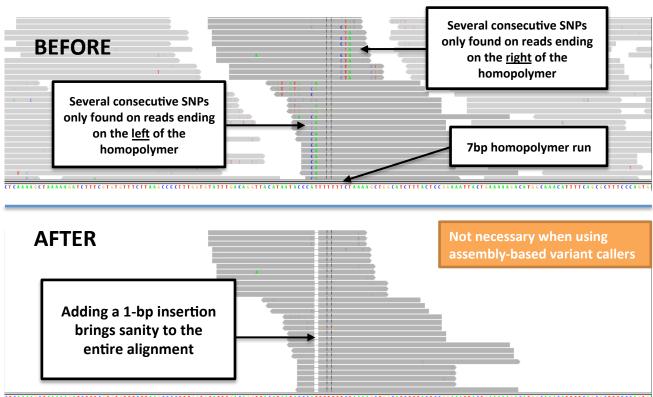
-> Must be removed to assess support for alleles correctly



**x** = sequencing error propagated in duplicates

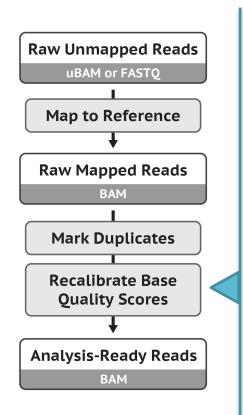
## Step 3: DEPRECATED: Local realignment around indels





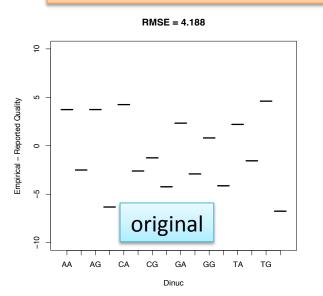
<sup>\*</sup> For implications, see <a href="https://gatkforums.broadinstitute.org/gatk/discussion/7847/changing-workflows-around-calling-snps-and-indels">https://gatkforums.broadinstitute.org/gatk/discussion/7847/changing-workflows-around-calling-snps-and-indels</a>

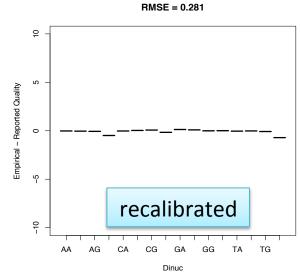
## **Step 4:** Base Recalibration (BQSR) corrects for machine errors



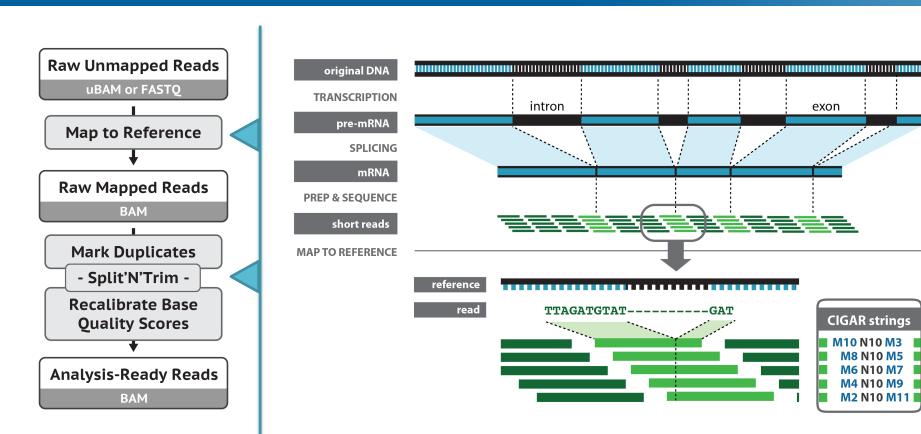
- Sequencers make systematic errors in base quality scores
- Sequencer quality cannot include PCR-based errors
- BQSR corrects the quality scores (not the bases)

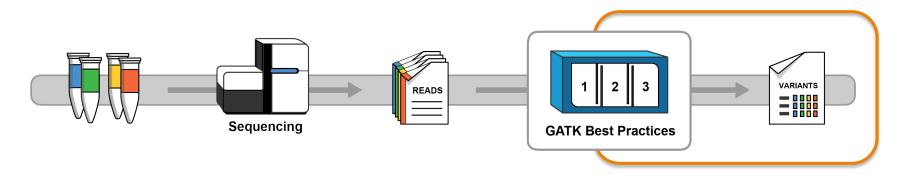
**Example of bias:** qualities reported depending on nucleotide context





## Special handling for RNAseq splice junctions





Next step:

# **VARIANT DISCOVERY**