

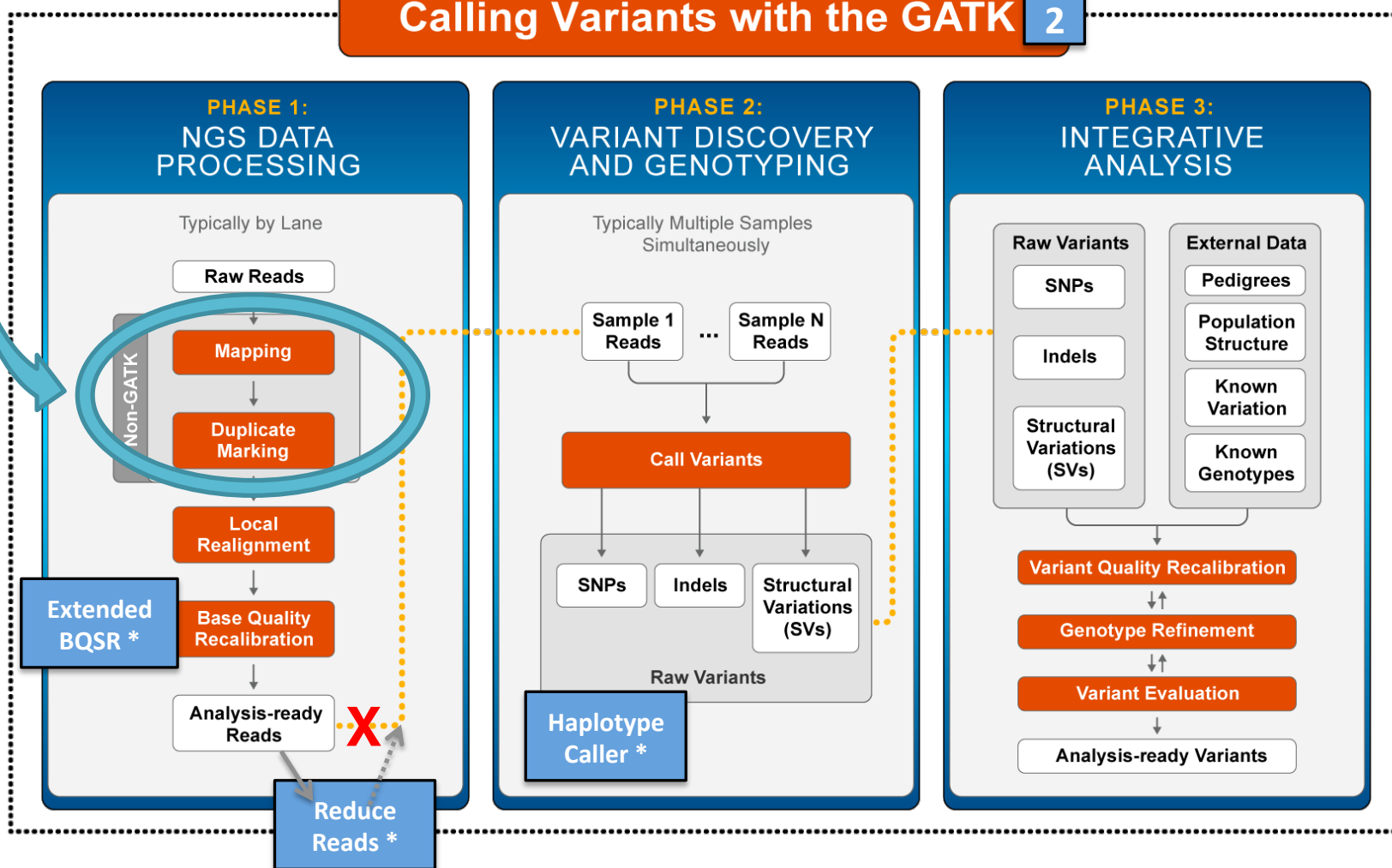
# Mapping and duplicate marking

First step from machine to analysis:  
place the reads in the right place on the  
reference and eliminate duplicates

# We are here in the Best Practices workflow

*MAPPING & DEDUPPING*

## Calling Variants with the GATK 2



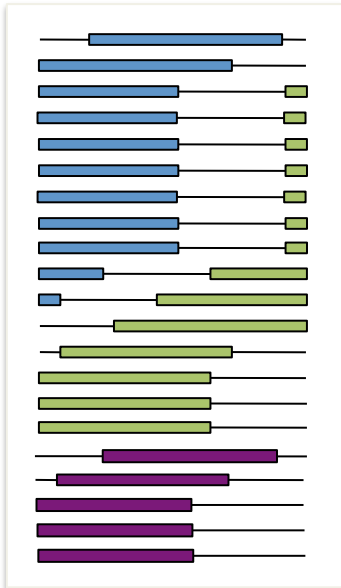
\* New tools or functionalities not available in GATK-Lite

# Overview of this step's goals

Reference genome



Enormous pile of short reads from NGS



→ Map reads to reference with **BWA**

All later steps assume that reads are placed in the right location and represent that region of the genome.

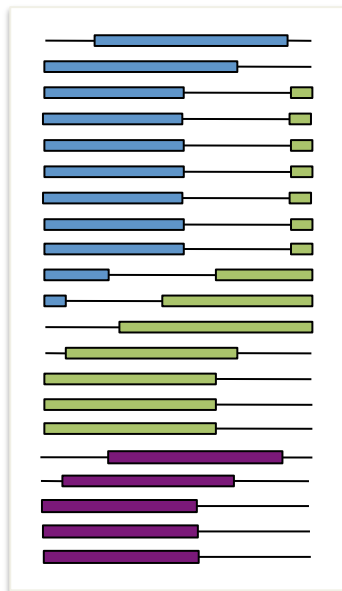
→ Mark duplicates with **Picard tools**

Duplicates originate mostly from DNA prep methods and cause biases that skew variant calling results.

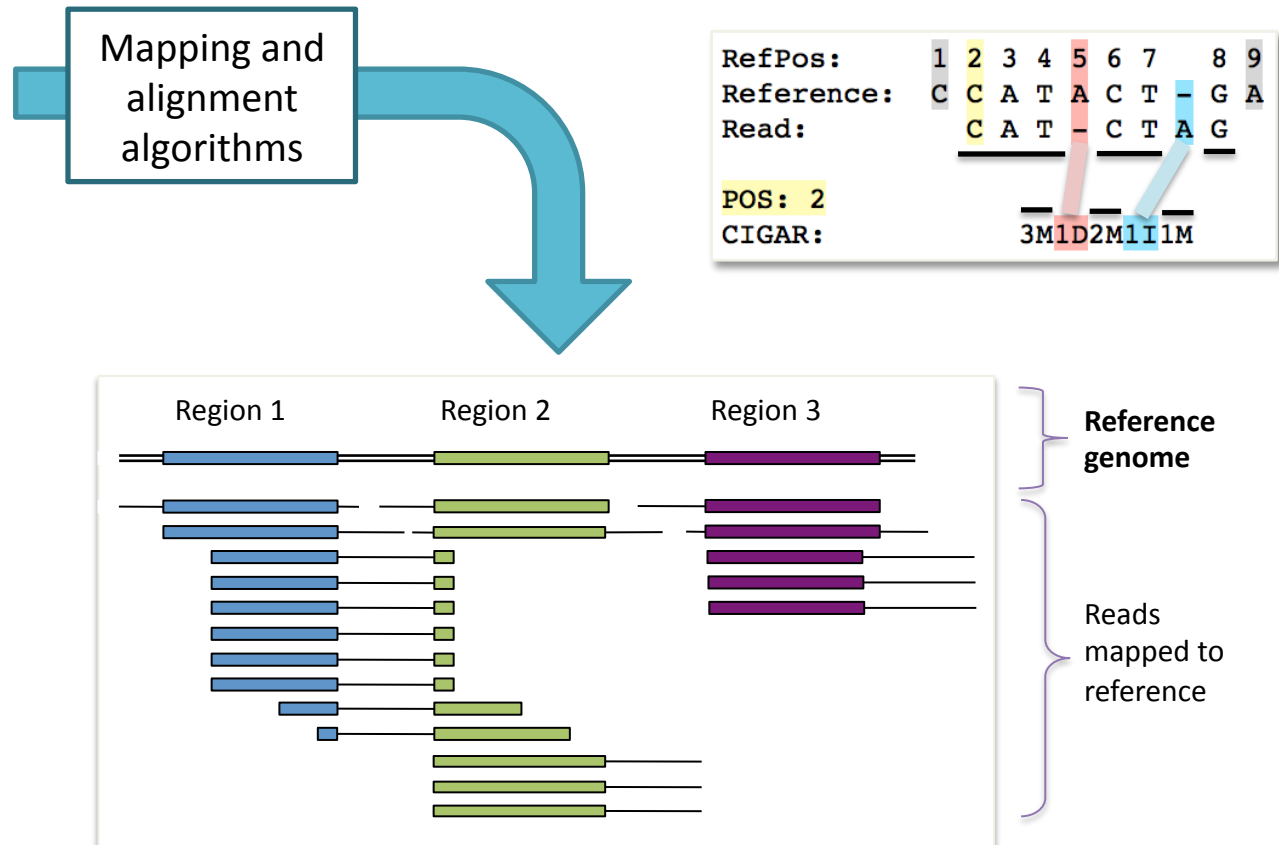
**MAPPING**

# Mapping short reads to a reference is simple in principle

Enormous pile of short reads from NGS

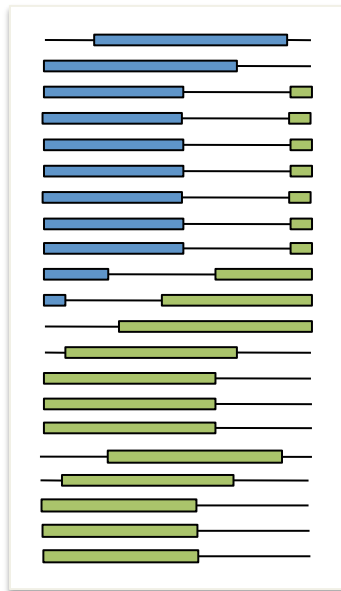


Identify where the read matches the reference sequence and record match details as CIGAR string



But mapping is actually very hard because of mismatches (true mutations or sequencing errors), duplicated regions etc.

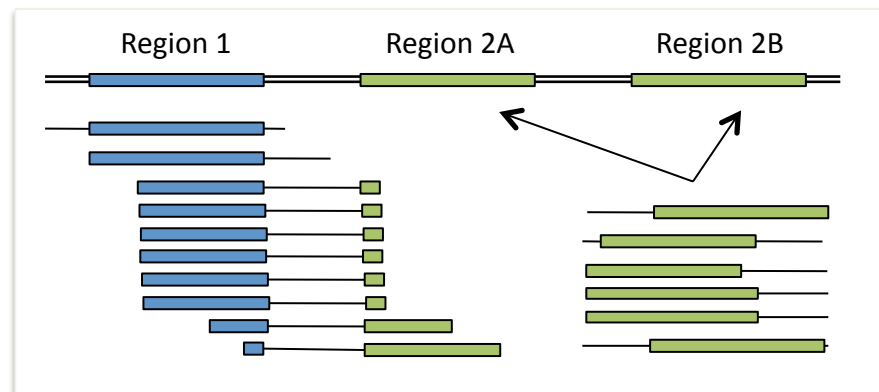
**Enormous pile of short reads from NGS**



Mapping and alignment algorithms

Mapping algorithms account for this by choosing the most likely placement

→ **mapping quality (MQ)**



Reference genome

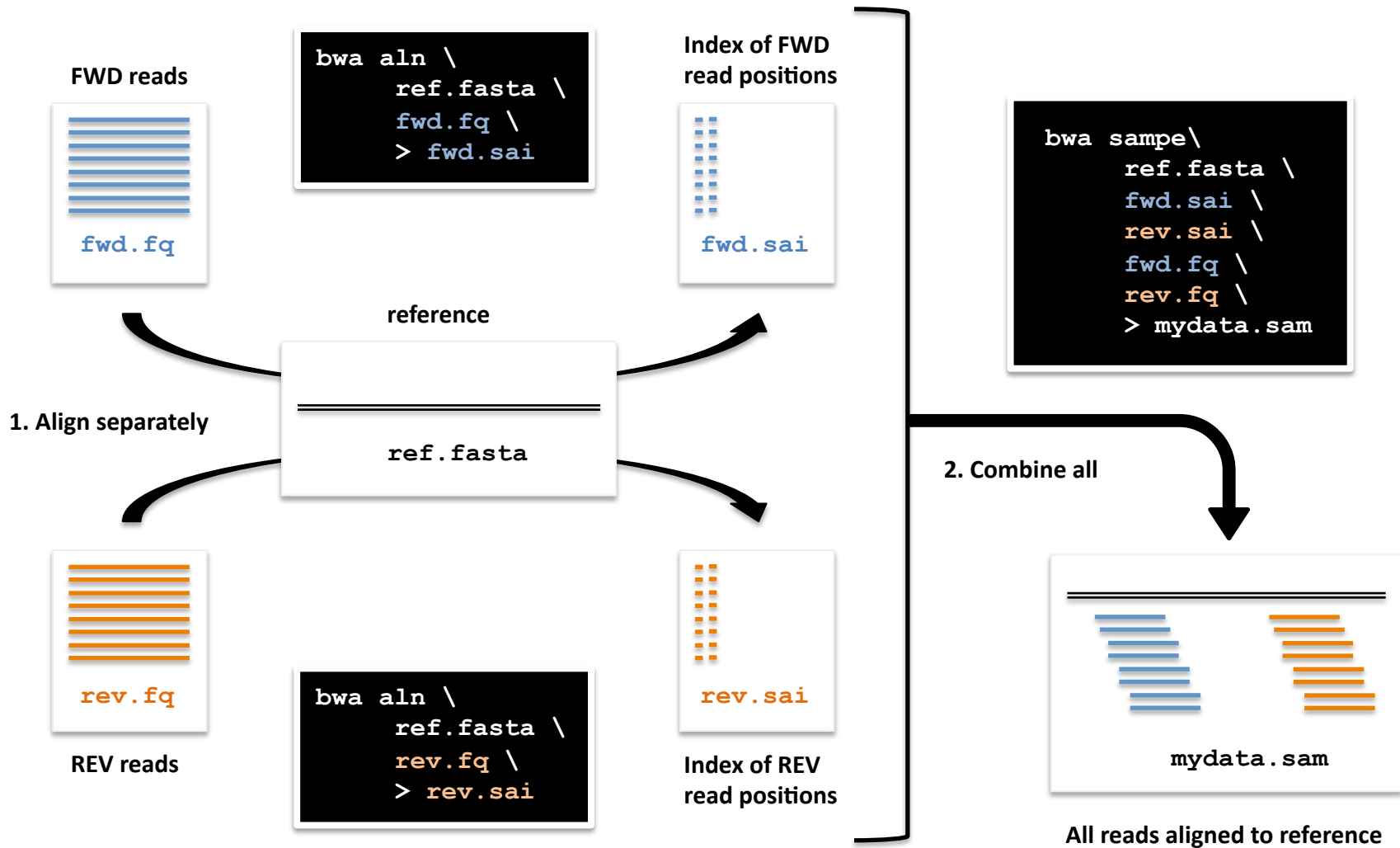
High MQ

Low MQ

For more information see:

Li and Homer (2010). A survey of sequence alignment algorithms for next-generation sequencing. *Briefings in Bioinformatics*.

# Typical workflow using BWA to map paired-end data

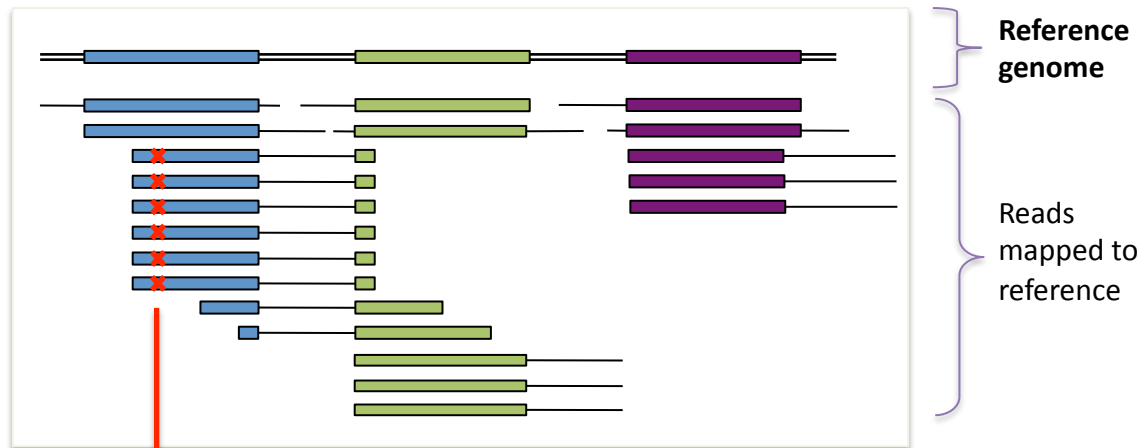


# **MARKING DUPLICATES**

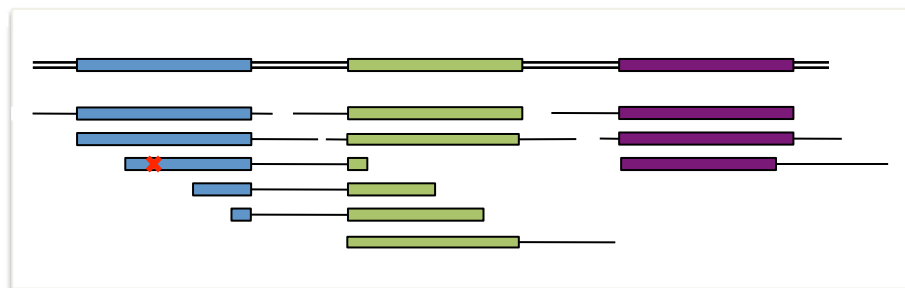


# The reason why duplicates are bad

✗ = sequencing error propagated in duplicates

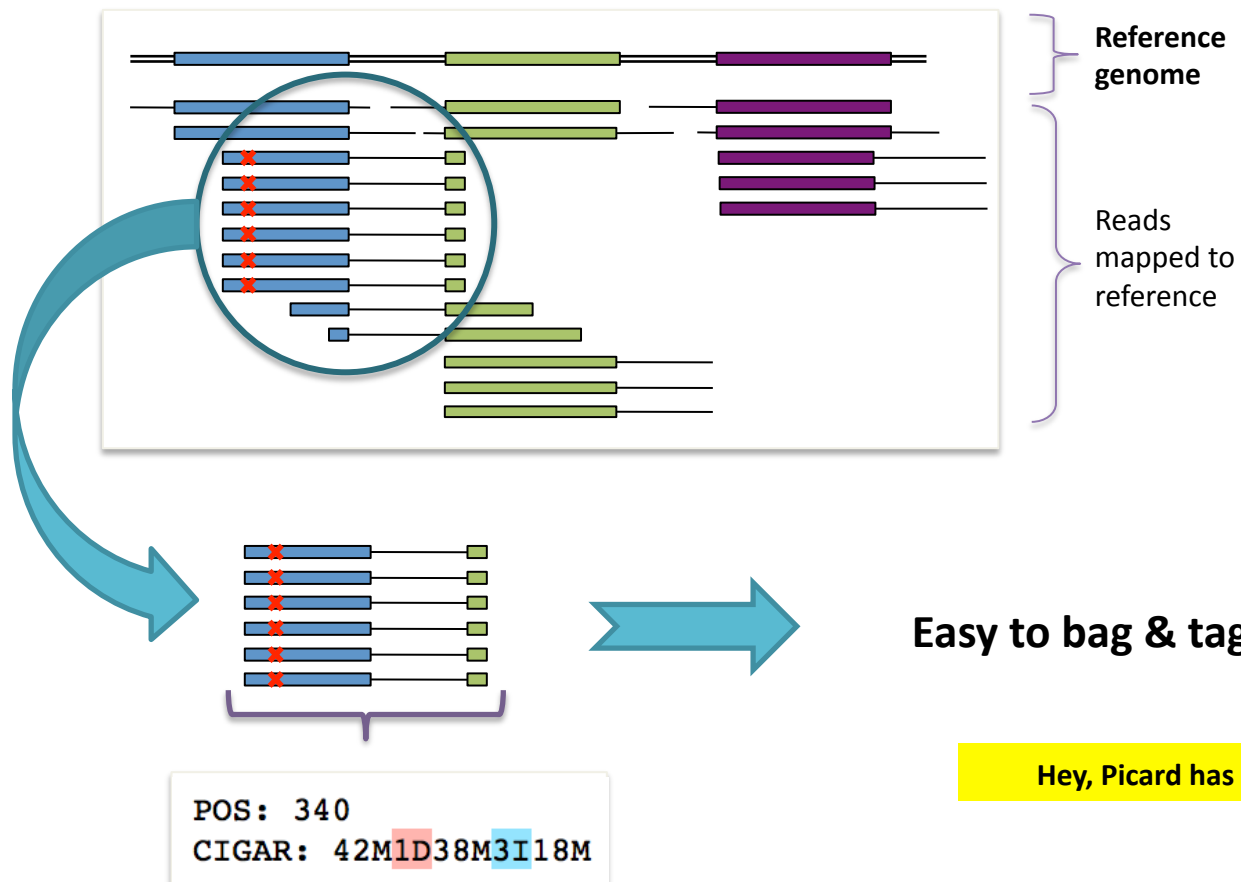


After marking duplicates, the GATK will only see :



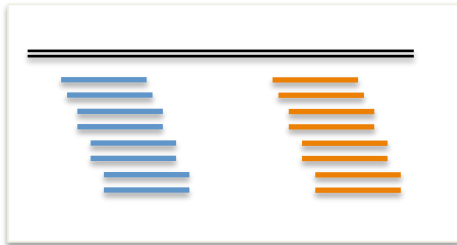
... and thus be more likely to make the right call

Duplicates have the same starting position  
and the same CIGAR string

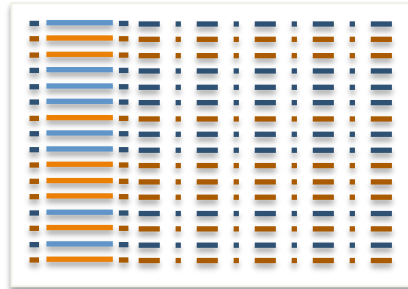


# A quick diversion about sorting and read groups

The information for this:

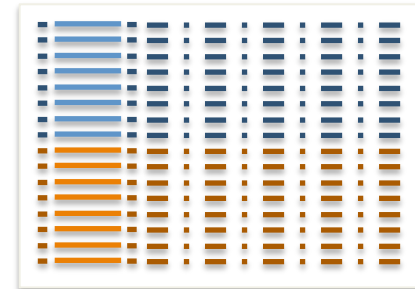


... is actually stored as a text file with one line per read which from far away looks like this:



The reads are in no particular order...

... but the GATK wants reads to be sorted by starting position like this:



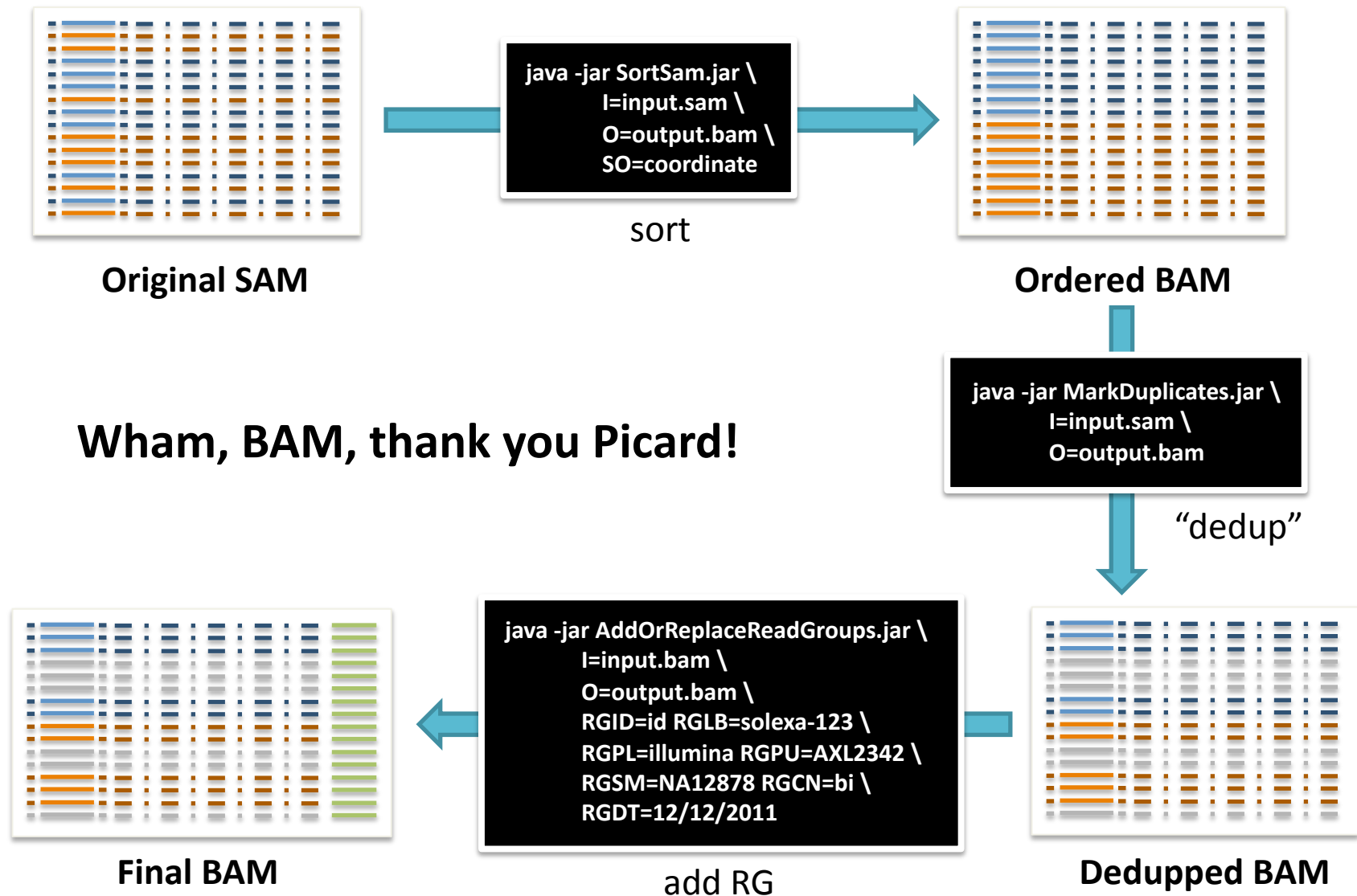
So we need to explicitly sort the SAM file...

... and Picard has an app for that!

And while we're at it, let's add **read group** information if it isn't already there, so **the GATK will know what read belongs to what sample** (that's kind of important).

Hey, Picard has an app for that too!

Typical workflow using Picard tools to mark duplicates *et al.*



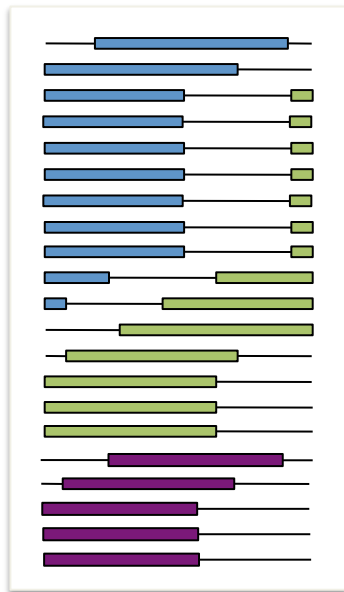
**TO CONCLUDE**

# Recap of what was achieved

Reference genome

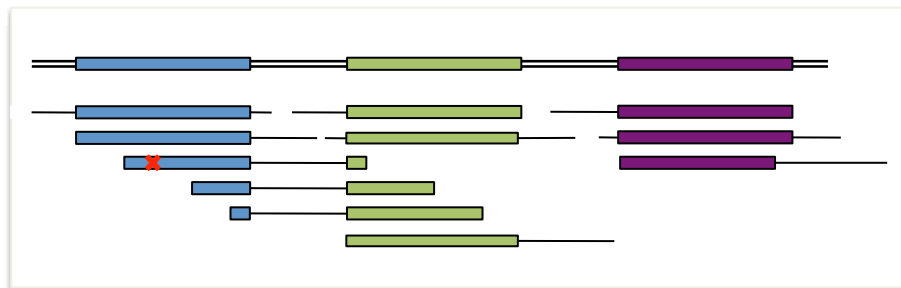


Enormous pile of short reads from NGS



→ Mapped reads to reference with **BWA**

→ Marked duplicates with **Picard tools**  
(and did some additional prep work)

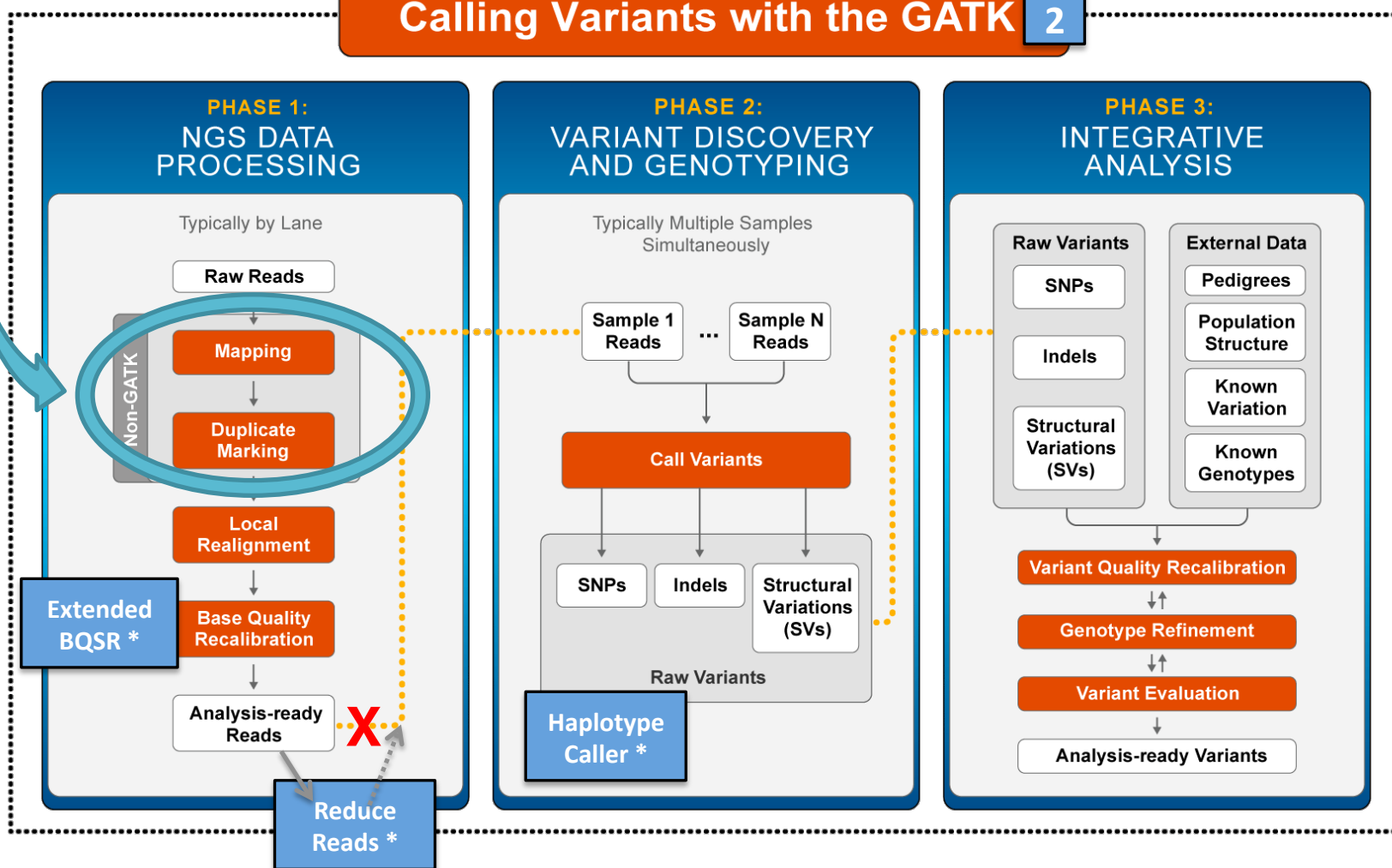


Reads mapped, sorted, deduplicated, +RGs

# We were here in the Best Practices workflow

*NEXT STEP: REALIGNMENT*

## Calling Variants with the GATK 2



\* New tools or functionalities not available in GATK-Lite

## Further reading

<http://www.broadinstitute.org/gatk/guide/topic?name=intro>

<http://www.broadinstitute.org/gatk/guide/topic?name=best-practices>