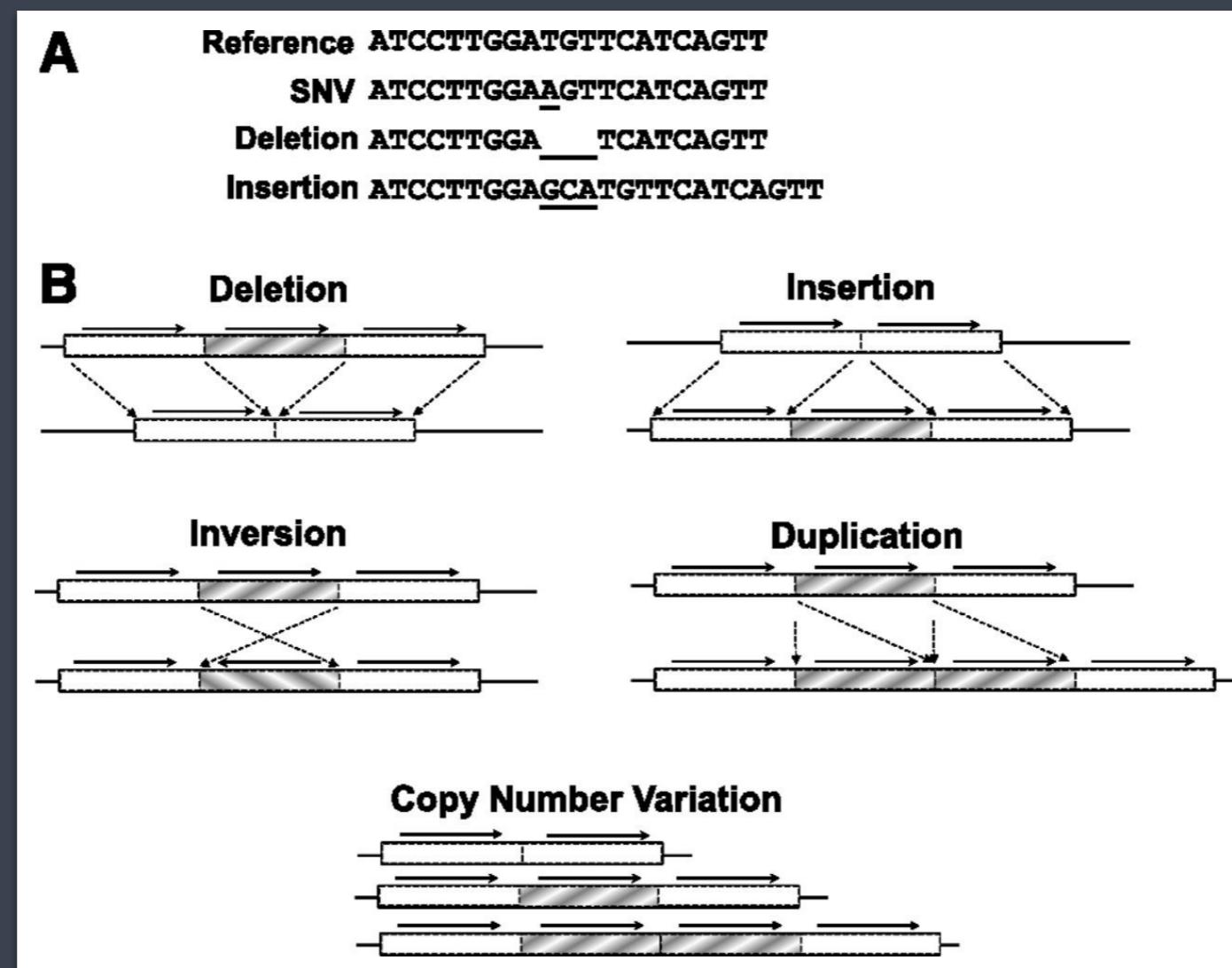
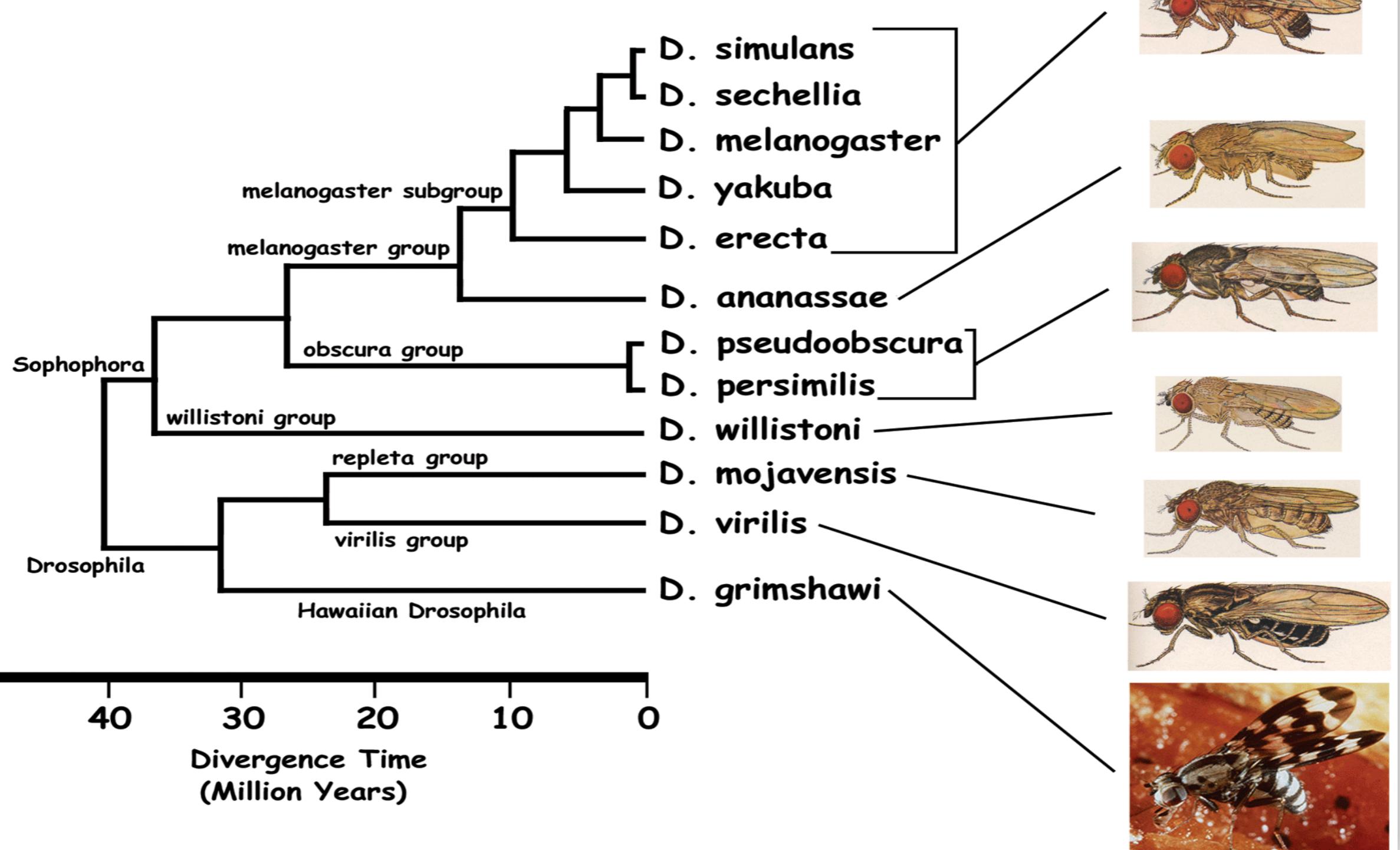
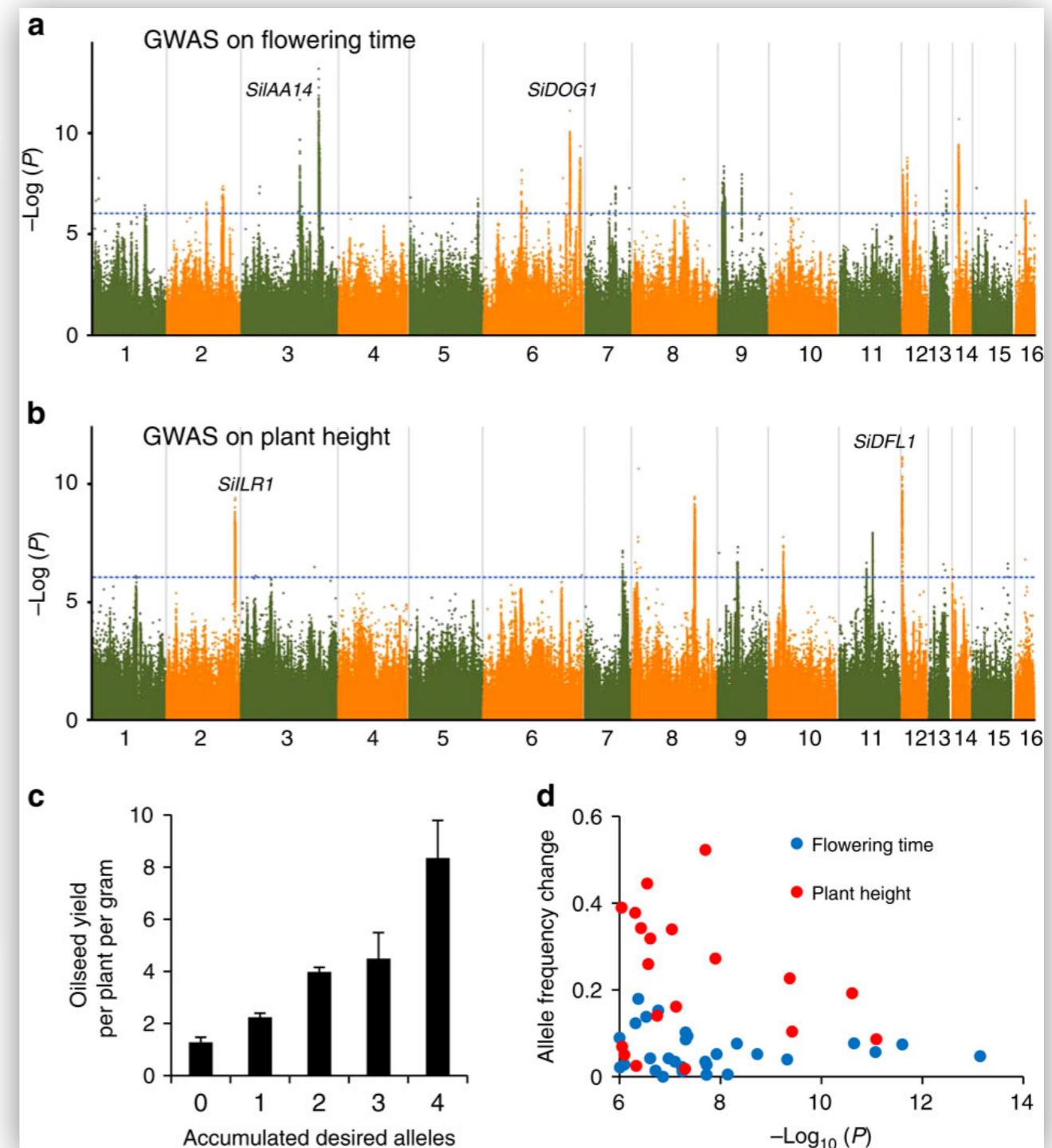
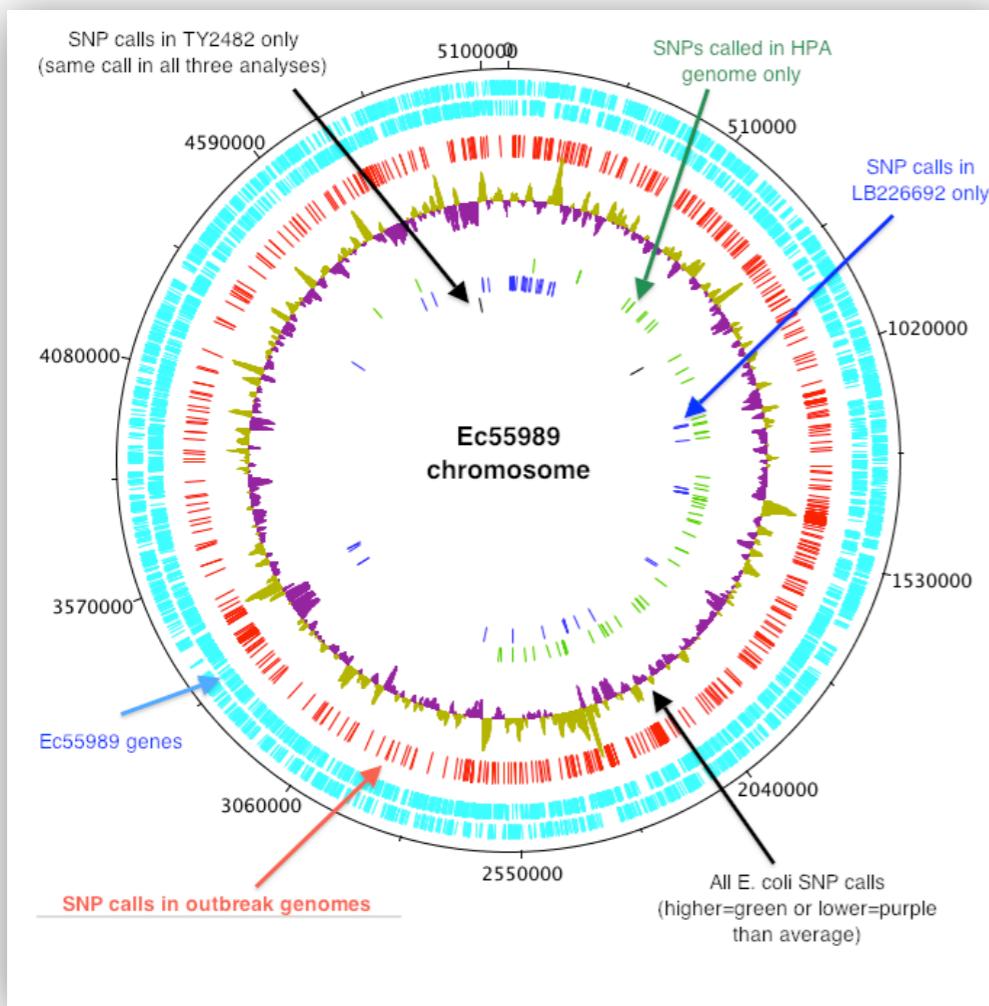


# Introduction to Variant Detection

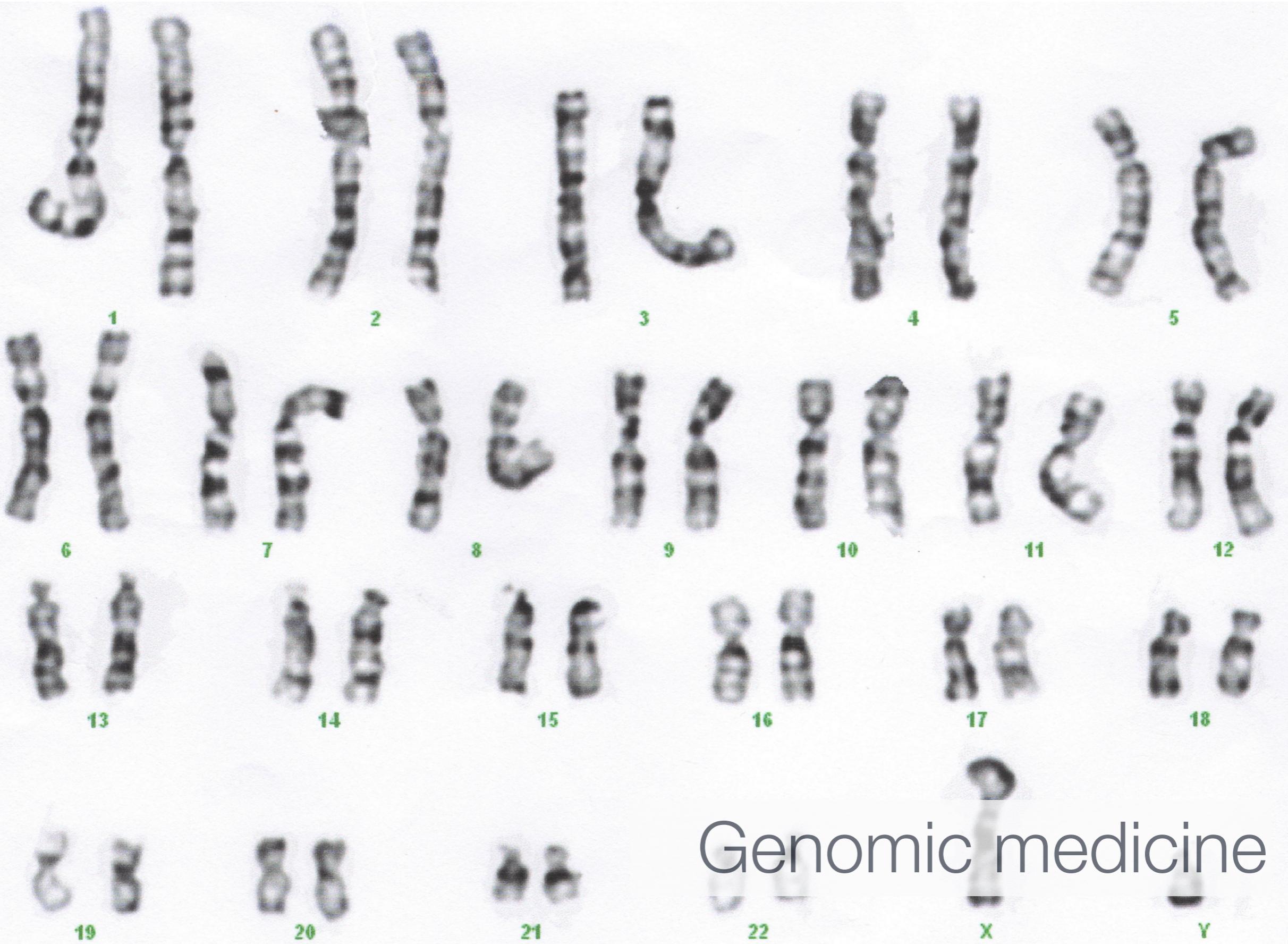




# Evolutionary analysis



# Medicine and Agriculture



Genomic medicine

# Overview

- » Human variations
  - Germline
  - Somatic
- » Types of Variations
- » Sequencing strategies to identify variants
- » Generalized analysis workflow (GATK best practice guidelines)

Any heritable “mutation” is considered a germline variant.

- found in populations, discovered by large-scale population analyses, and contained in databases like dbSNP, HapMap
- most are not deleterious

## Germline vs Somatic mutations

Any heritable “mutation” is considered a germline variant.

- found in populations, discovered by large-scale population analyses, and contained in databases like dbSNP, HapMap
- most are not deleterious

A somatic variant is any mutation that arises in a single cell of an individual and is only present in the descendants of that cell, not all the cells of that individual.

- found in rapidly growing cancer cells
- can be silent or pathogenic

## Germline vs Somatic mutations

# Most human genomic variants have no phenotypic impacts

- Ones that have an impact are either positively selected, i.e. they confer a reproductive advantage
- Or they are neutral. These are often associated with ethnic origin, typically affecting traits like height, facial features, hair or skin color

Phenotypic impacts

## Most human genomic variants have no phenotypic impacts

- Ones that have an impact are either positively selected, i.e. they confer a reproductive advantage
- Or they are neutral. These are often associated with ethnic origin, typically affecting traits like height, facial features, hair or skin color

## Some genomic variants have deleterious effects

- Most of these are recessive: their effect is observed only if both alleles are affected
- Those that are dominant will either be selected against and disappear, or have effects that minimally impact reproductive fitness

Phenotypic impacts

# Types of variations

Single Nucleotide Polymorphisms (SNPs)

Small Insertions/Deletions (Indels)

Copy Number Variations (CNVs)

Structural Variations (SVs)

For SNPs, many different methods have been used:

- Hybridization based, primarily SNP arrays
- Enzyme-based methods, primarily oligonucleotide ligation and RFLP
- Methods measuring physical properties of DNA

How to assess genomic diversity?

For SNPs, many different methods have been used:

- Hybridization based, primarily SNP arrays
- Enzyme-based methods, primarily oligonucleotide ligation and RFLP
- Methods measuring physical properties of DNA

How to assess genomic diversity?

For SNPs, many different methods have been used:

- Hybridization based, primarily SNP arrays
- Enzyme-based methods, primarily oligonucleotide ligation and RFLP
- Methods measuring physical properties of DNA

For CNVs, the main methods are hybridization based

How to assess genomic diversity?

For SNPs, many different methods have been used:

- Hybridization based, primarily SNP arrays
- Enzyme-based methods, primarily oligonucleotide ligation and RFLP
- Methods measuring physical properties of DNA

For CNVs, the main methods are hybridization based

For SVs the most reliable ones used partial sequencing of large clones  
(e.g. fosmids)

How to assess genomic diversity?

For SNPs, many different methods have been used:

- Hybridization based, primarily SNP arrays
- Enzyme-based methods, primarily oligonucleotide ligation and RFLP
- Methods measuring physical properties of DNA

For CNVs, the main methods are hybridization based

For SVs the most reliable ones used partial sequencing of large clones  
(e.g. fosmids)

**NGS can detect all types of variants**  
**(Paired-end data preferred!)**

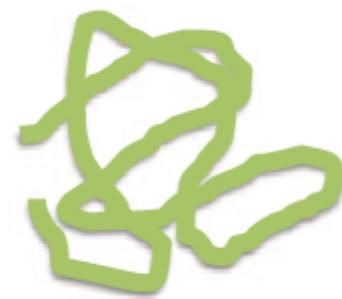
How to assess genomic diversity?

# Sequencing strategies

Whole Genome Sequencing (WGS)  
(for SNPs/Indels, CNVs and SVs)

Exome Sequencing (for SNPs/Indels)

Gene Panels (for SNPs/Indels)



Genomic  
DNA

Next-generation  
DNA sequencing

... CATTCACTAG ...    ... AGCCATTAG ...  
... GGTAGTTAG ...    ... GGTAAACTAG ...  
... TATAATTAG ...    ... CGTACCTAG ...  
...

A diagram showing six short DNA sequence reads, each represented by a box containing a sequence of nucleotides (e.g., CATTCACTAG, AGCCATTAG, GGTAGTTAG, GGTAAACTAG, TATAATTAG, CGTACCTAG) followed by three dots indicating more sequence.

millions-billions of *reads*  
~30-1000 nucleotides

**Resequencing**



Align reads to *reference genome* and identify variants

***De novo assembly***

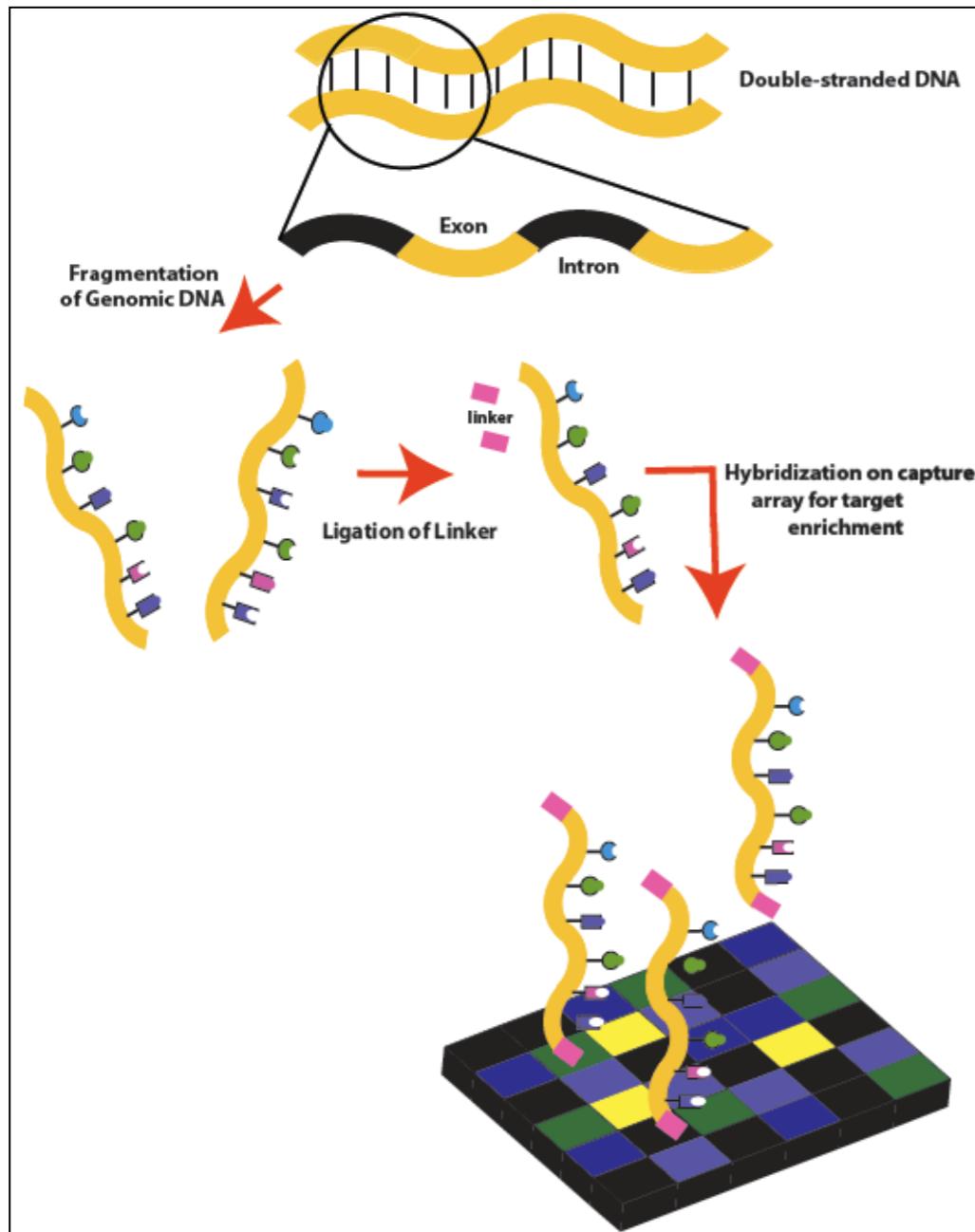


Construct genome sequence from overlaps between reads

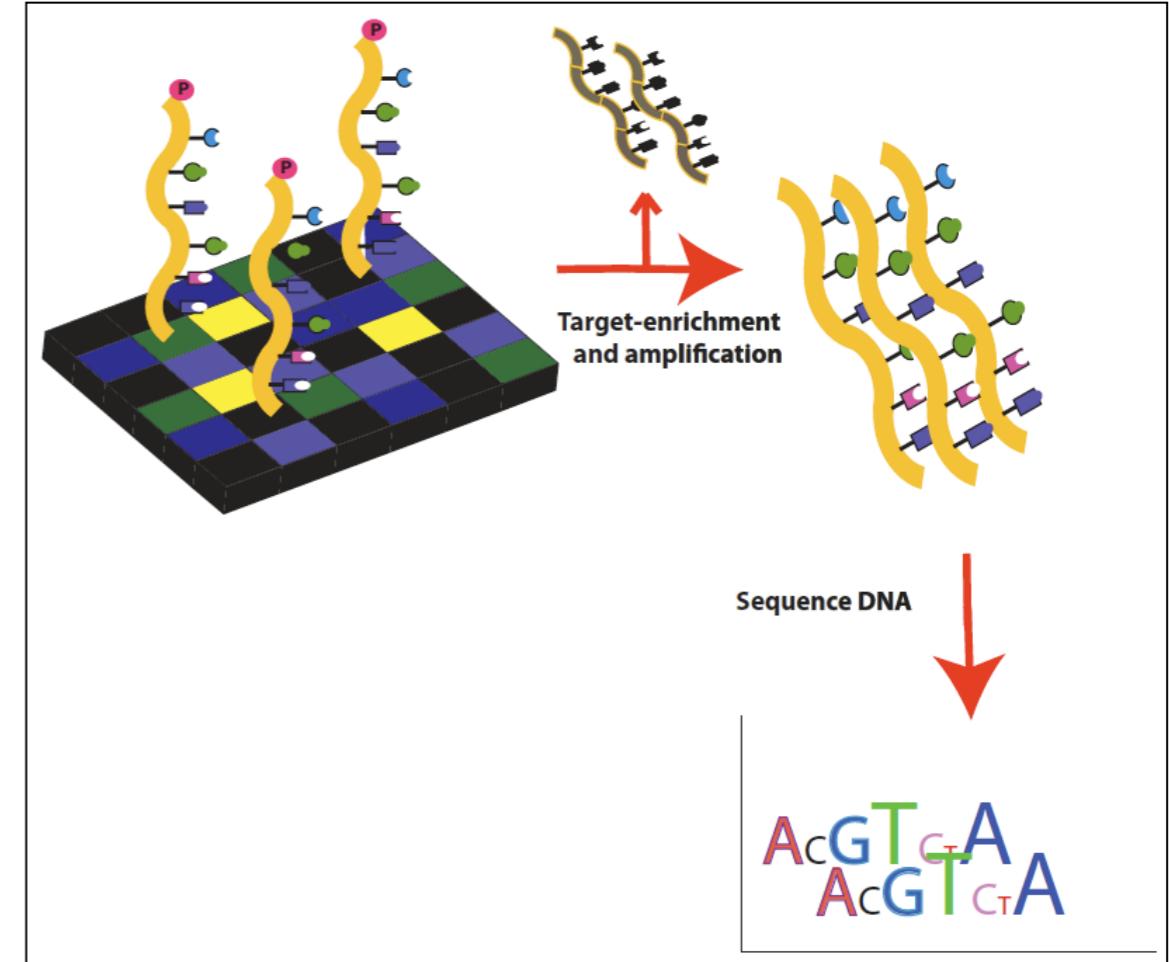
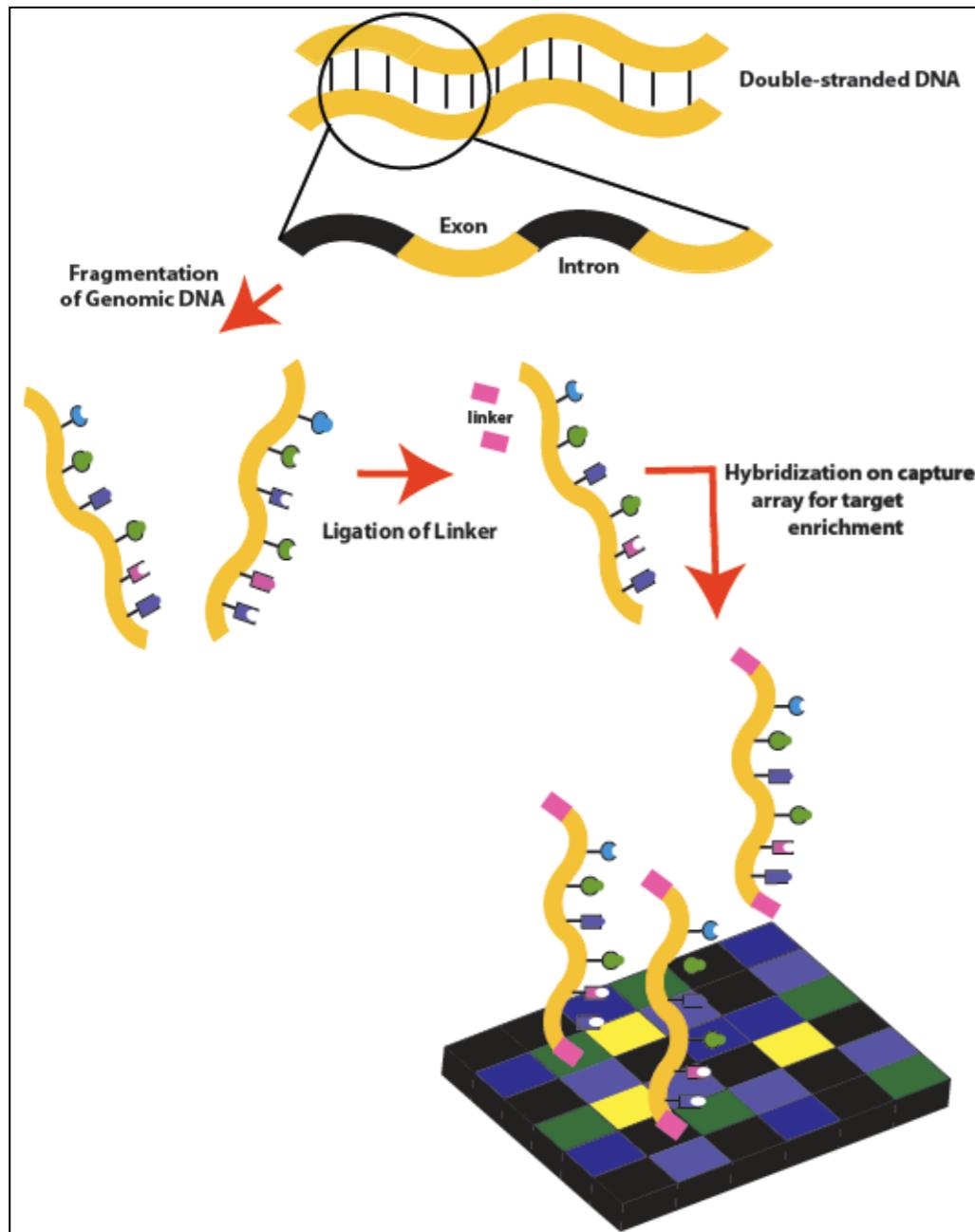
Benjamin J. Raphael\*

Department of Computer Science and Center for Computational Molecular Biology, Brown University, Providence, Rhode Island, United States of America

# Whole genome sequencing

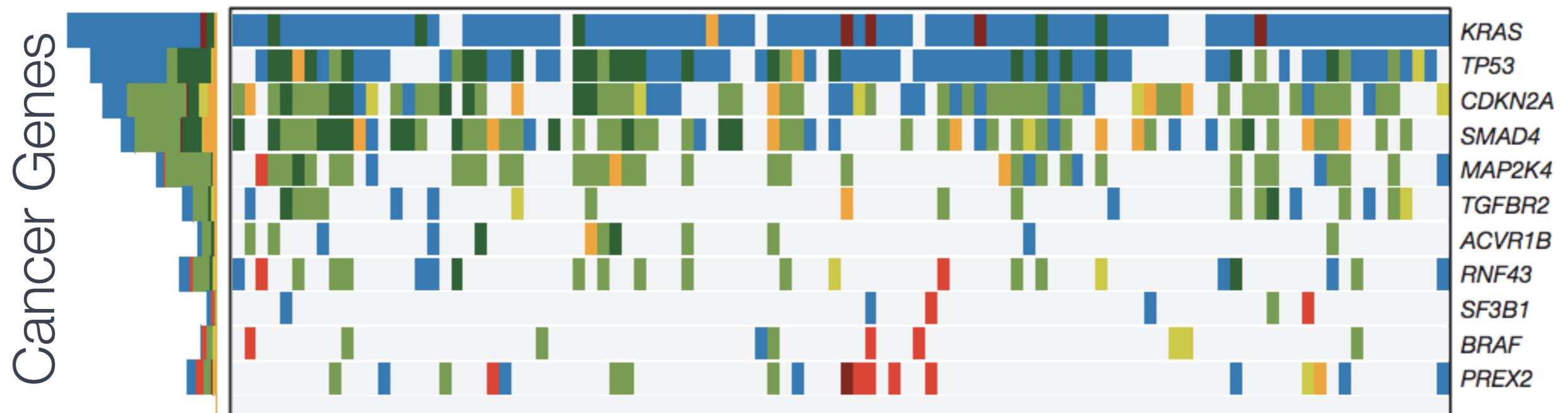


# Exome sequencing



# Exome sequencing

## Patients



A visualization of an analysis using a panel of known cancer genes

Gene panel sequencing for diagnostics

- Targeted gene panels are most commonly used for diagnostics/clinical work
- Coverage: cost considerations for various methods, based on number of samples
- Variants in un-targeted or non-exonic regions will be missed

Gene panels or ES or WGS:  
Which one is “better”?



# Sequencing depth and cost

For WGS

- Haploid genome size => 3.2 Giga base pairs (3.2 billion)

Sequencing depth?

For WGS

- Haploid genome size => 3.2 Giga base pairs (3.2 billion)
- Minimum 30x for WGS

Sequencing depth?

For WGS

- Haploid genome size => 3.2 Giga base pairs (3.2 billion)
- Minimum 30x for WGS

For Exome Sequencing

- Exome size => 33 Mega base pairs (33 million bases)

Sequencing depth?

For WGS

- Haploid genome size => 3.2 Giga base pairs (3.2 billion)
- Minimum 30x for WGS

For Exome Sequencing

- Exome size => 33 Mega base pairs (33 million bases)
- About 100 times smaller than WGS

Sequencing depth?

For WGS

- Haploid genome size => 3.2 Giga base pairs (3.2 billion)
- Minimum 30x for WGS

For Exome Sequencing

- Exome size => 33 Mega base pairs (33 million bases)
- About 100 times smaller than WGS
- 70x-100x for ES, with additional considerations for unevenness of coverage

Sequencing depth?

## For WGS

- Haploid genome size => 3.2 Giga base pairs (3.2 billion)
- Minimum 30x for WGS

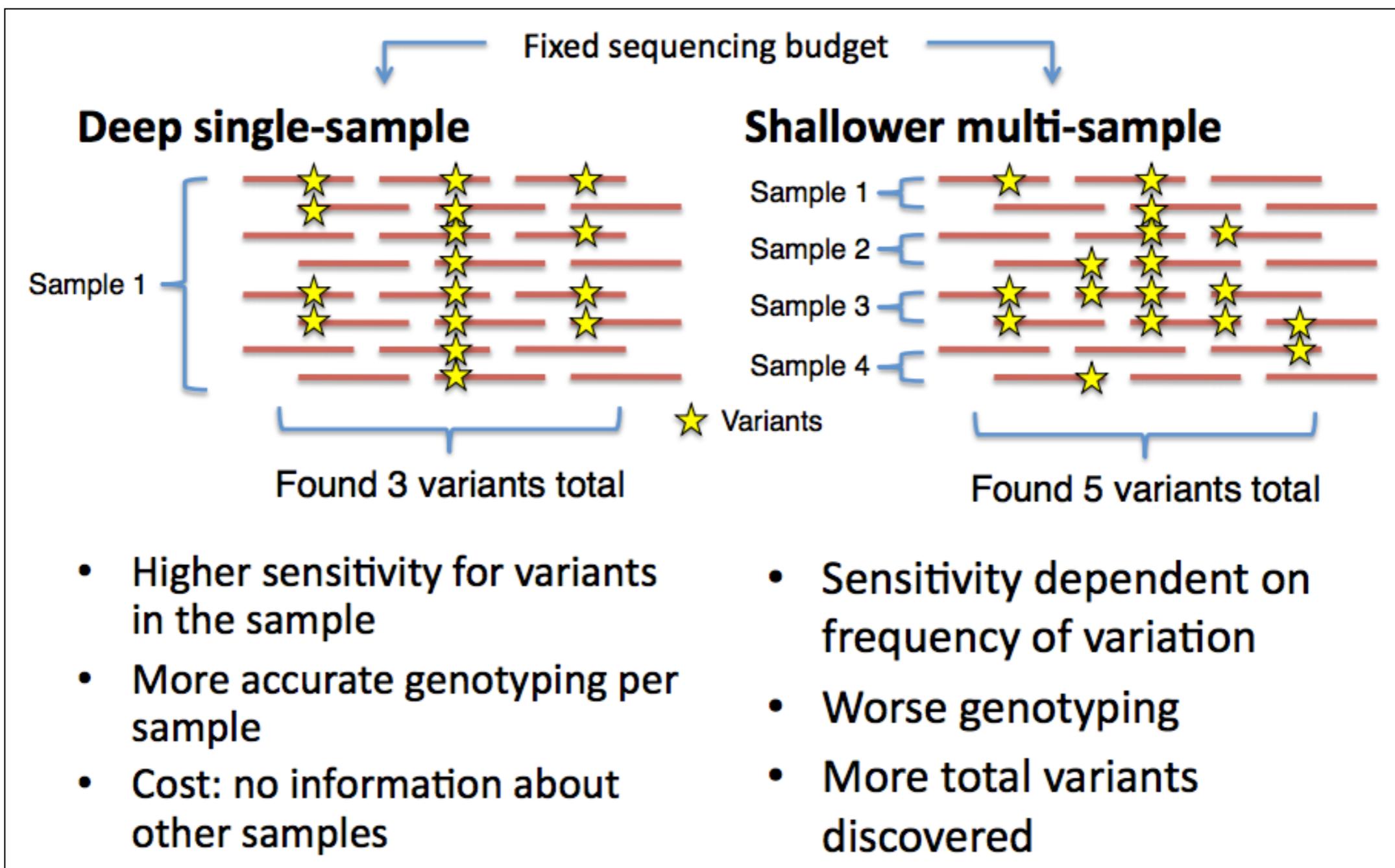
## For Exome Sequencing

- Exome size => 33 Mega base pairs (33 million bases)
- About 100 times smaller than WGS
- 70x-100x for ES, with additional considerations for unevenness of coverage

## For Gene Panels

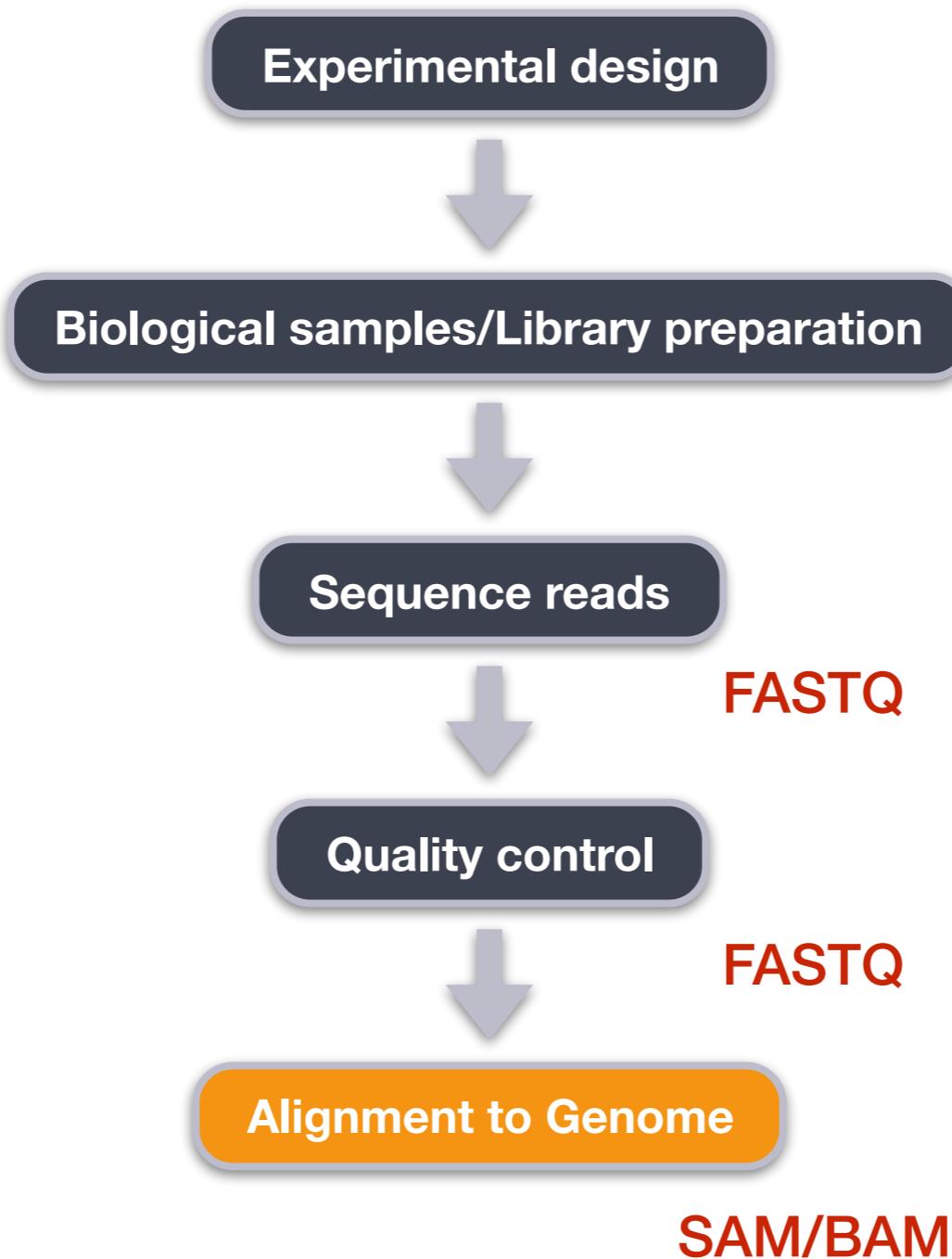
- 10x-20x coverage for gene panels for heterozygous germline variants

Sequencing depth?



# Sequencing depth and cost

# Generalized Variant Calling Workflow



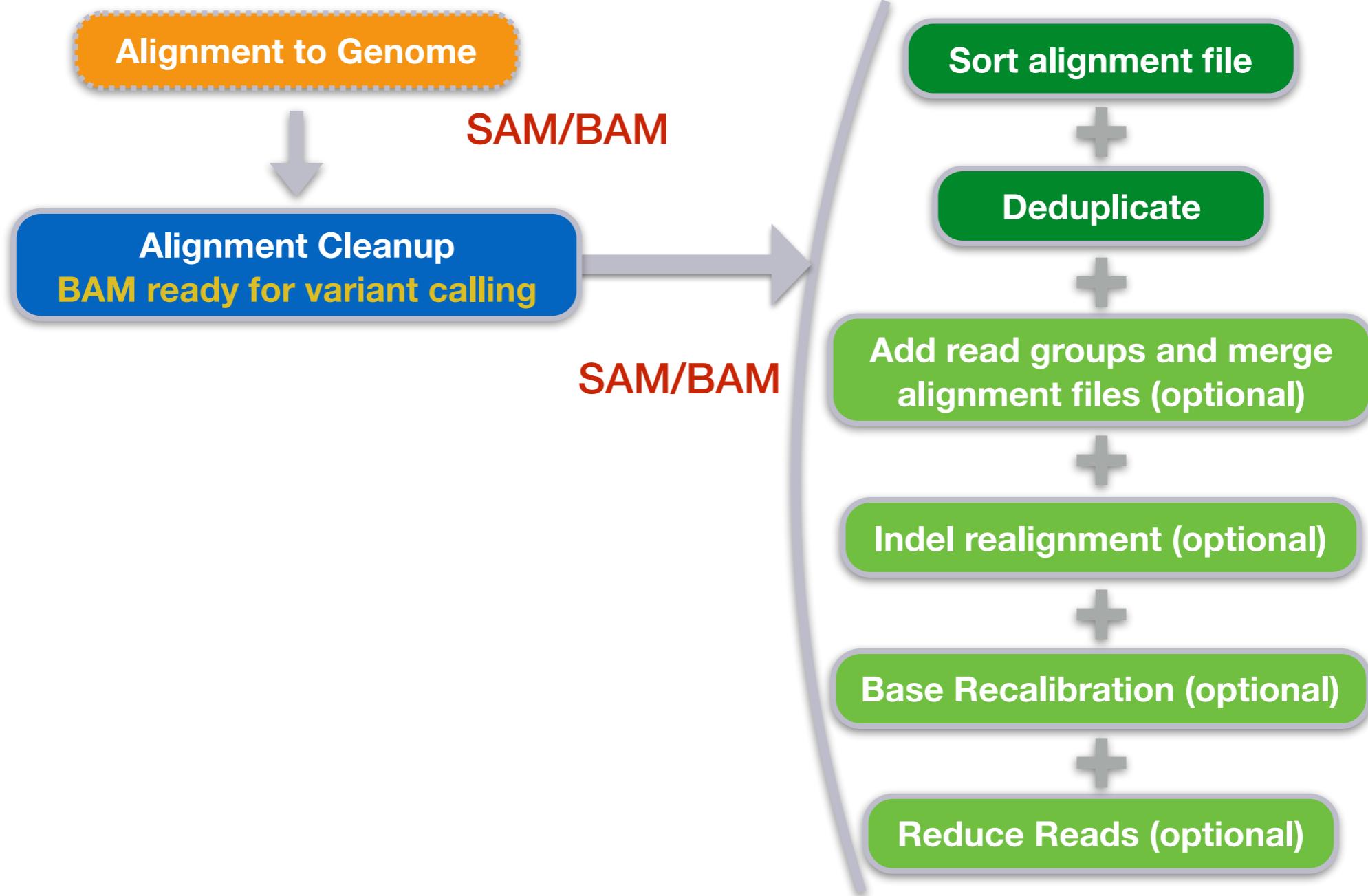
**Alignment to Genome**



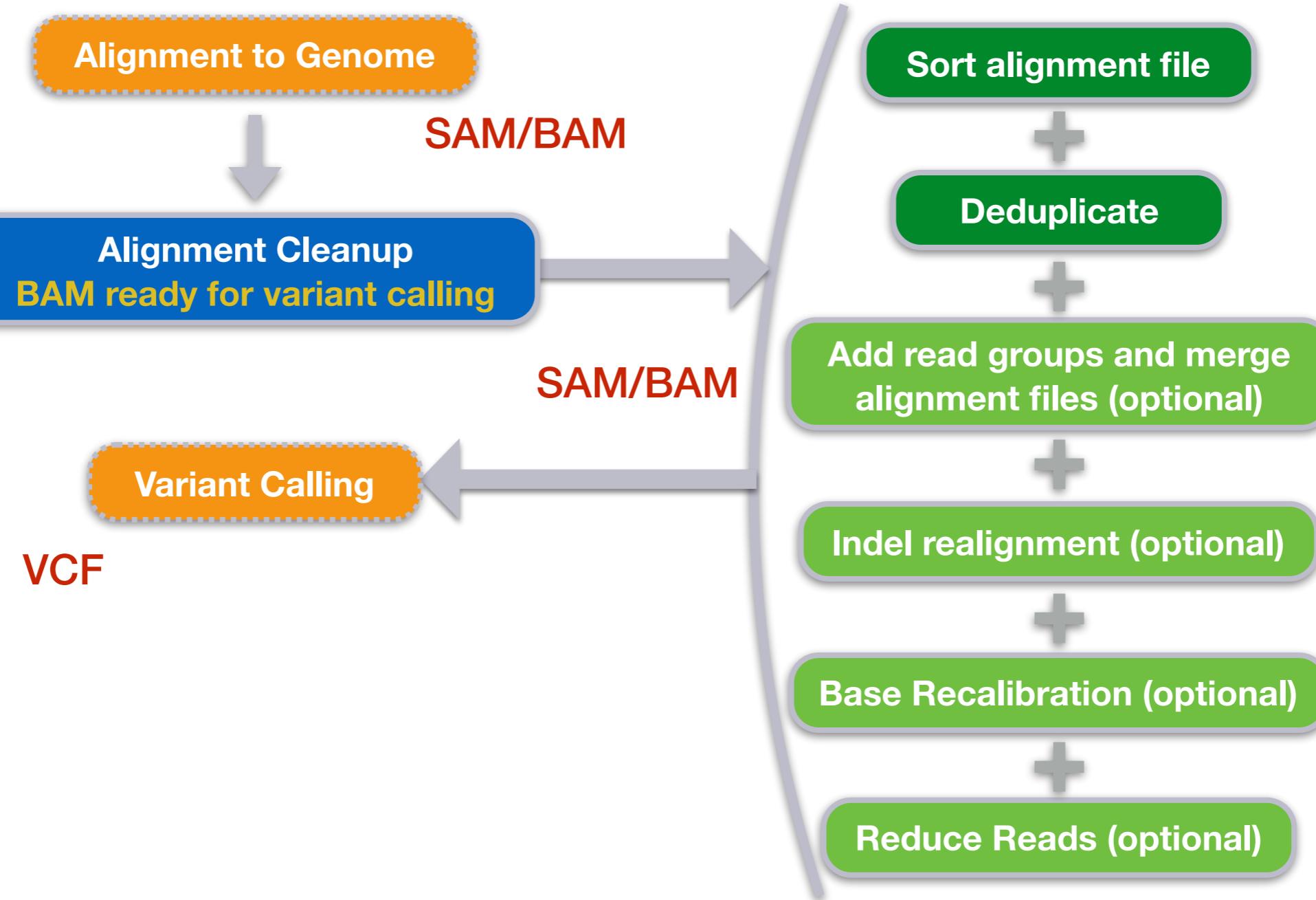
**SAM/BAM**

**Alignment Cleanup**  
**BAM ready for variant calling**

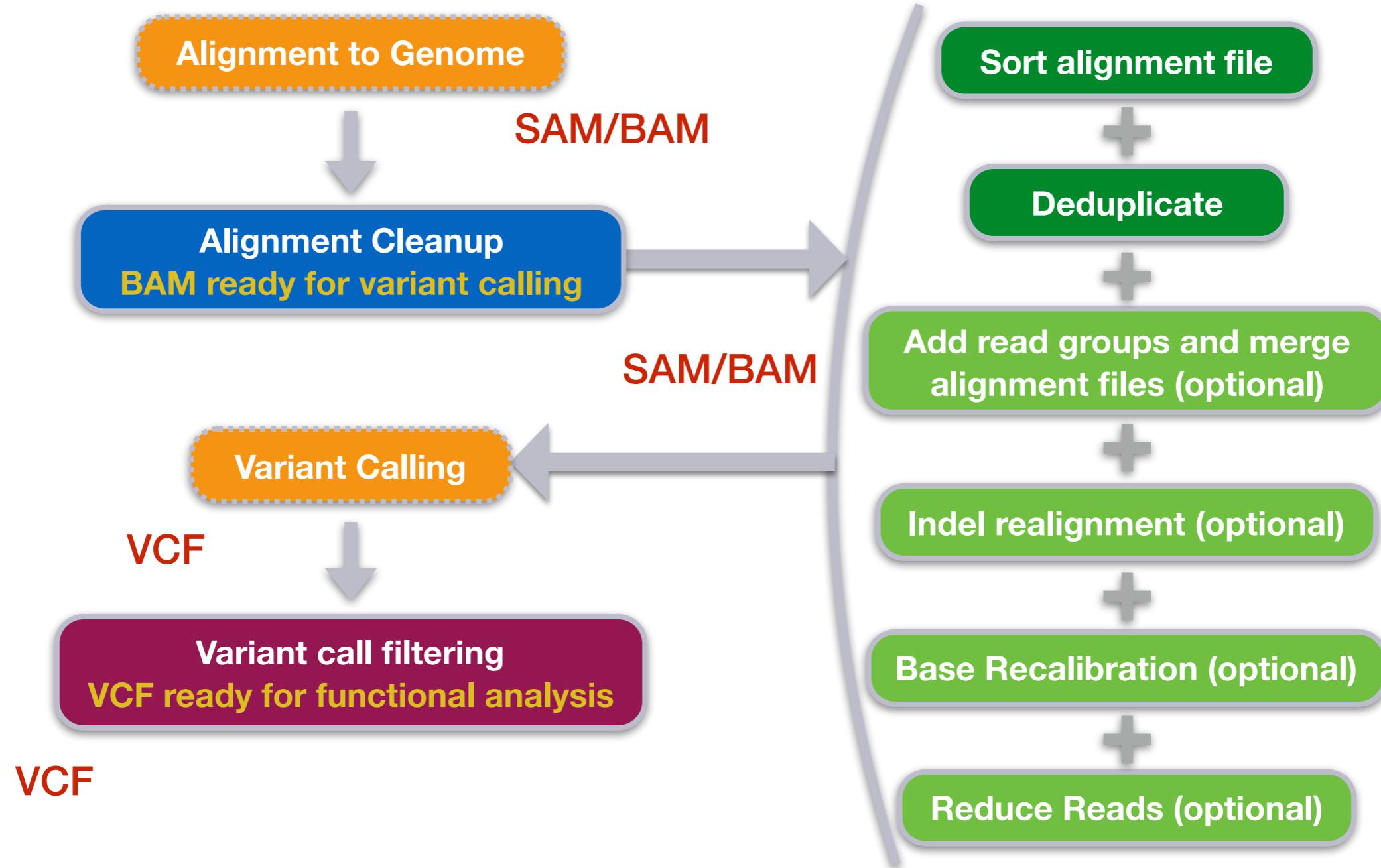
# Generalized Variant Calling Workflow



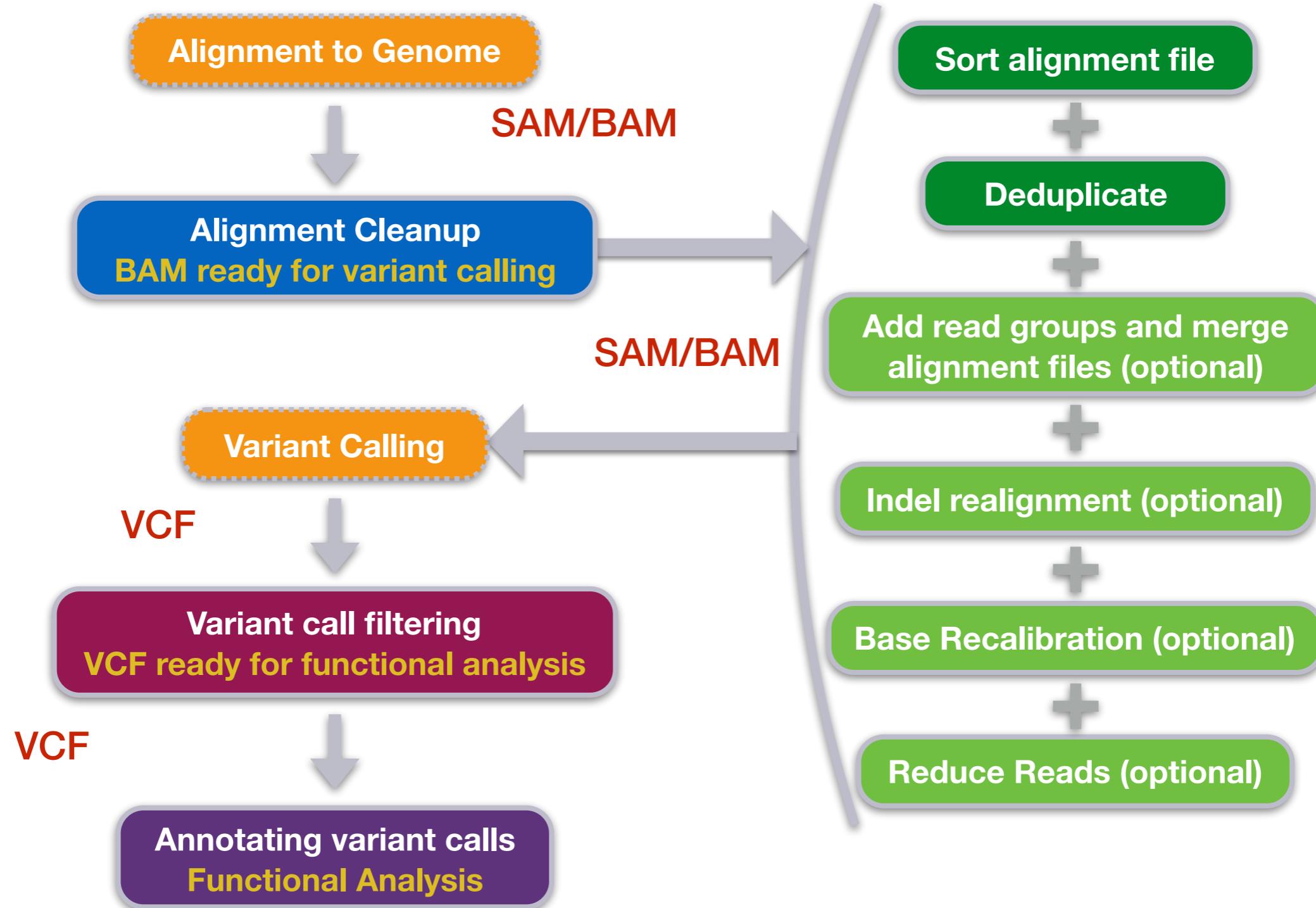
# Generalized Variant Calling Workflow



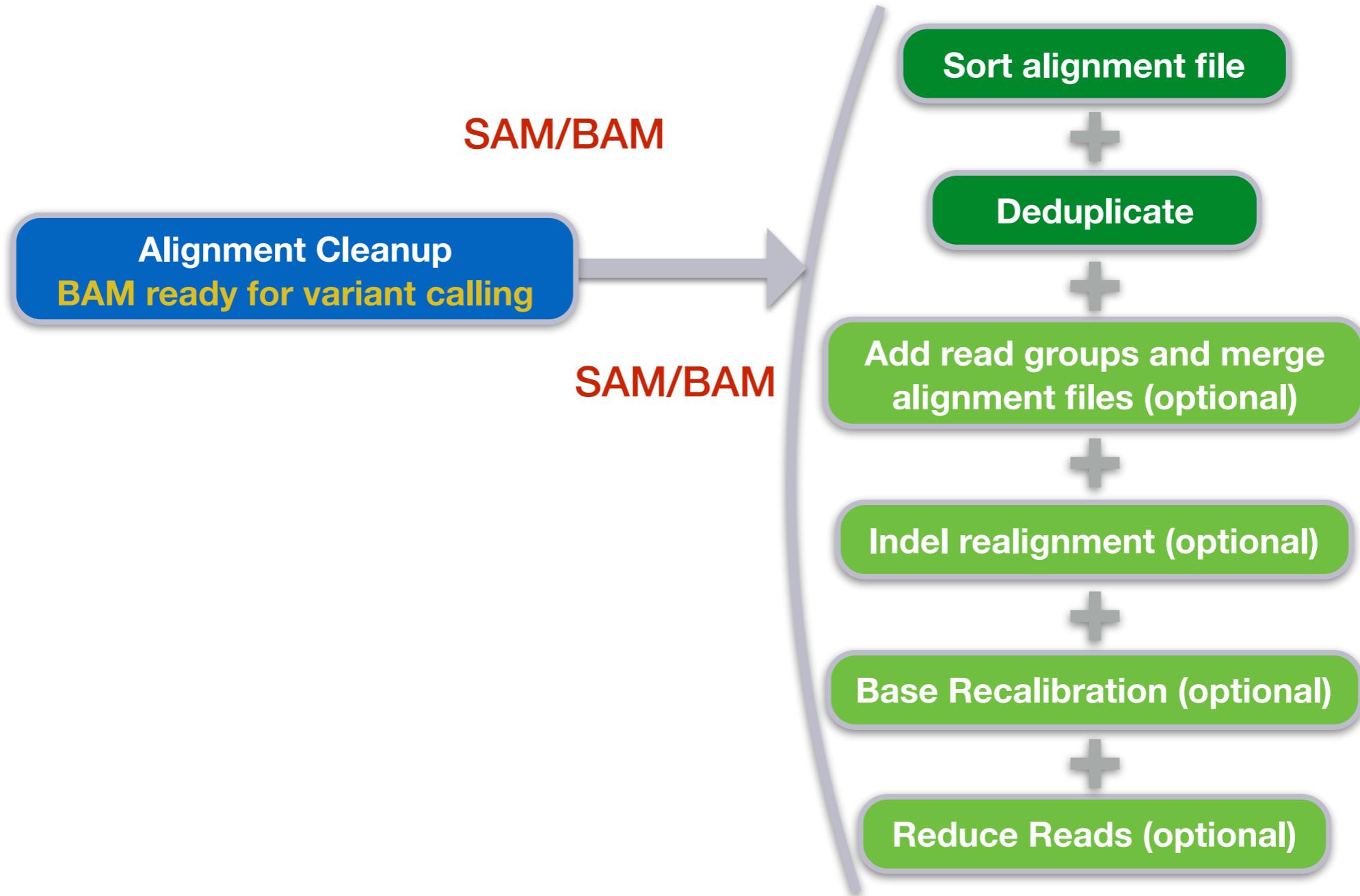
# Generalized Variant Calling Workflow



# Generalized Variant Calling Workflow



# Generalized Variant Calling Workflow



# Generalized Variant Calling Workflow

**Sort alignment file**

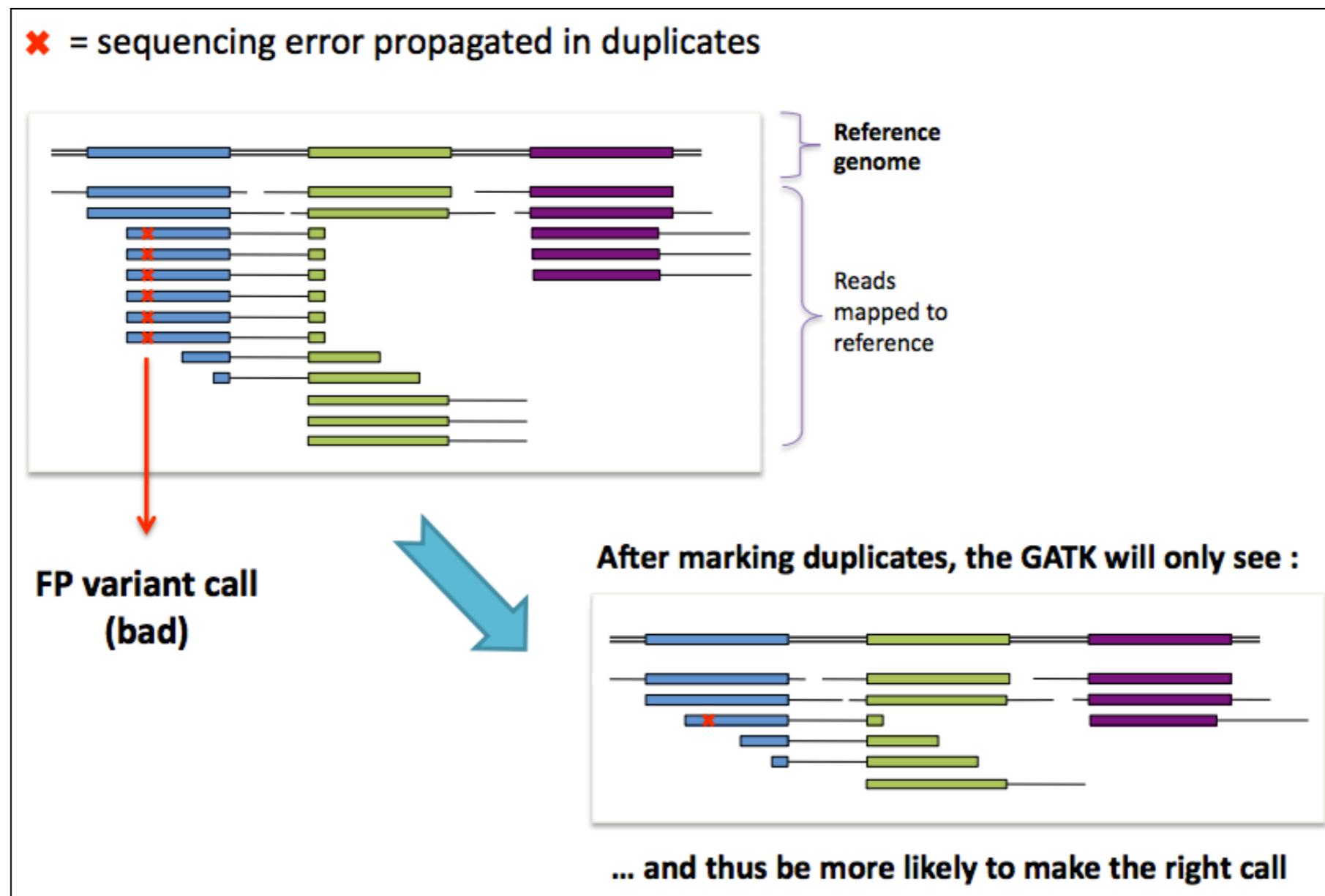


The reads are in no particular order...

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

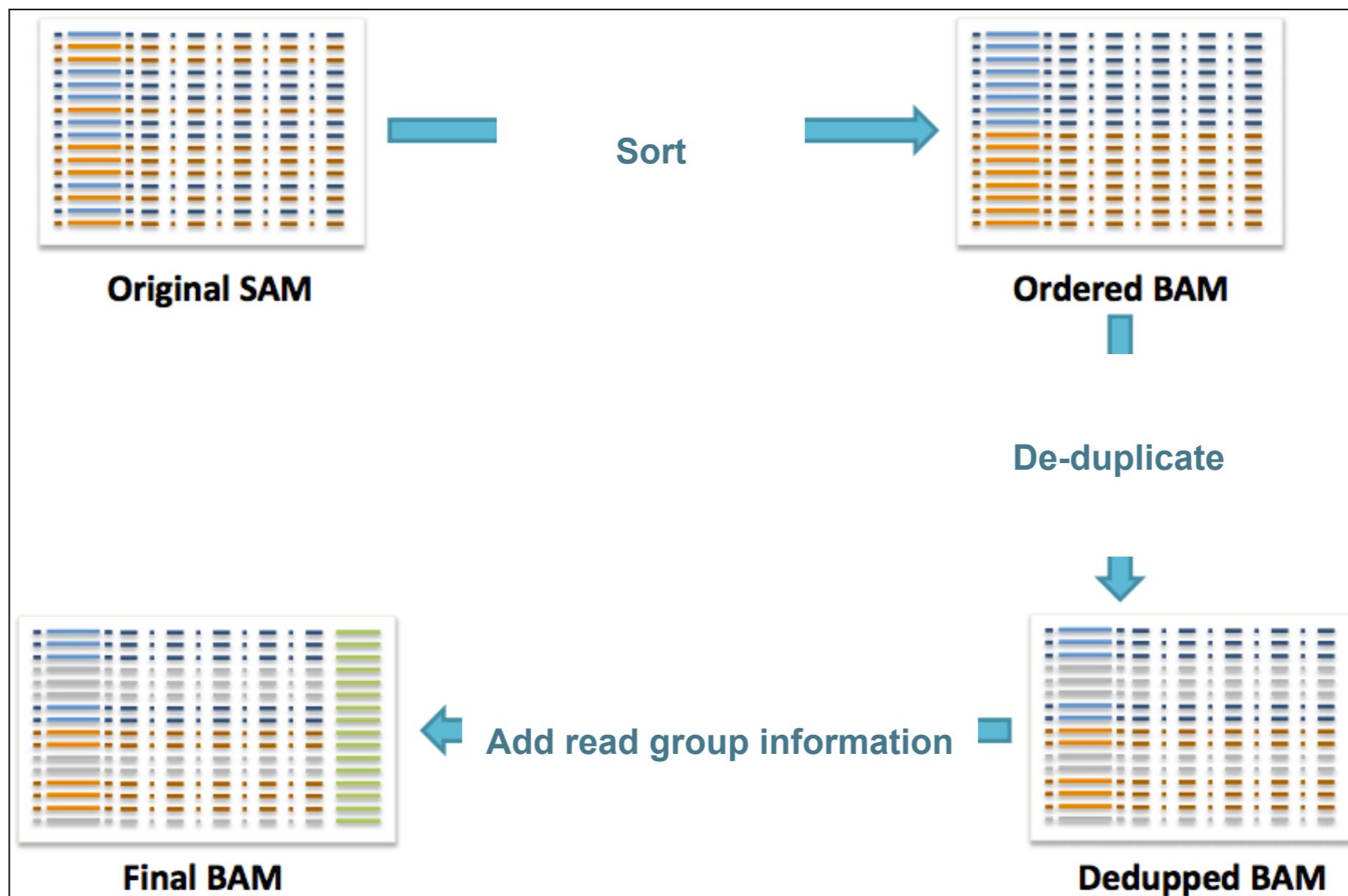
# Generalized Variant Calling Workflow

## Deduplicate



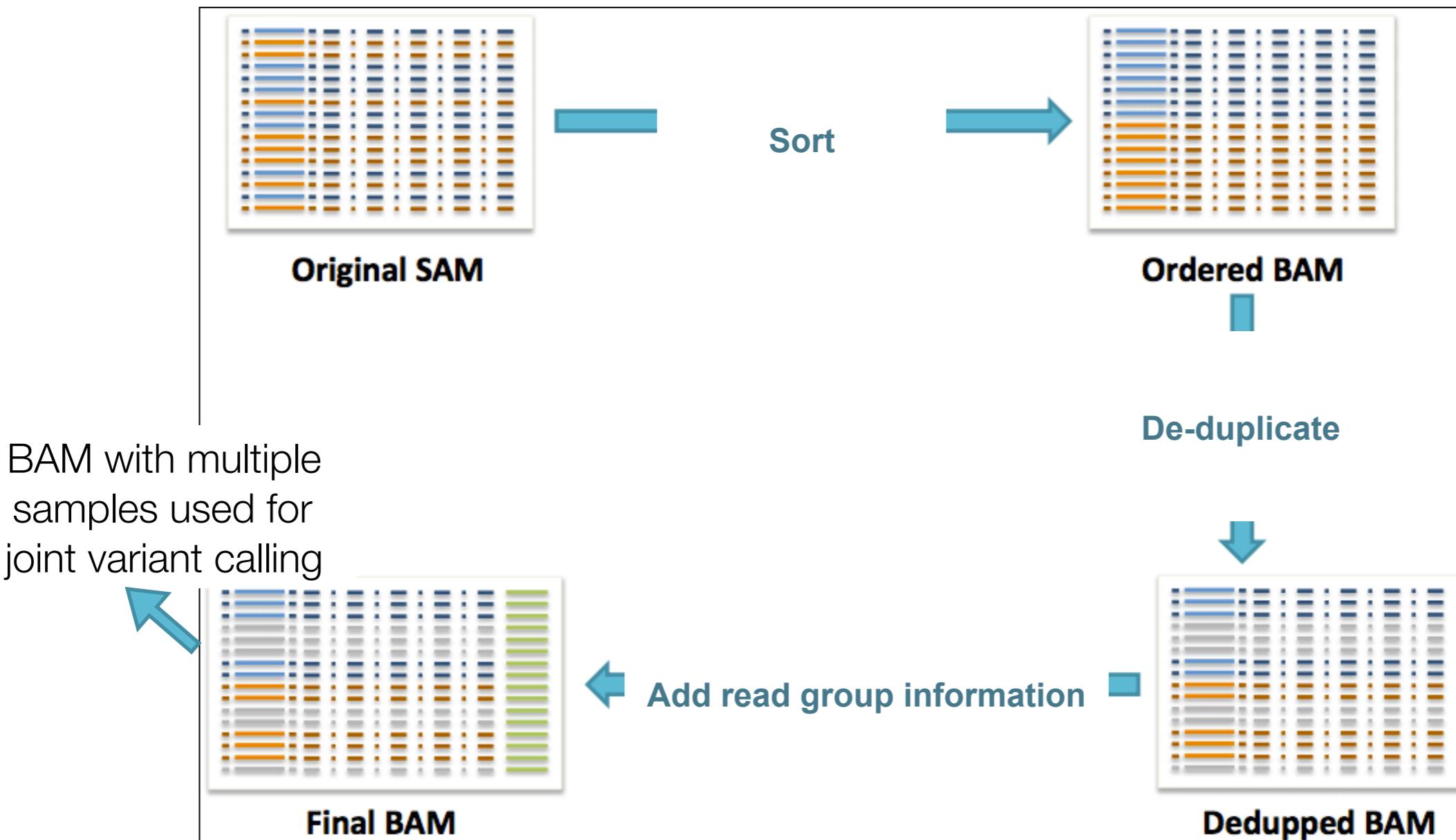
# Generalized Variant Calling Workflow

Add read groups and merge  
alignment files (optional)



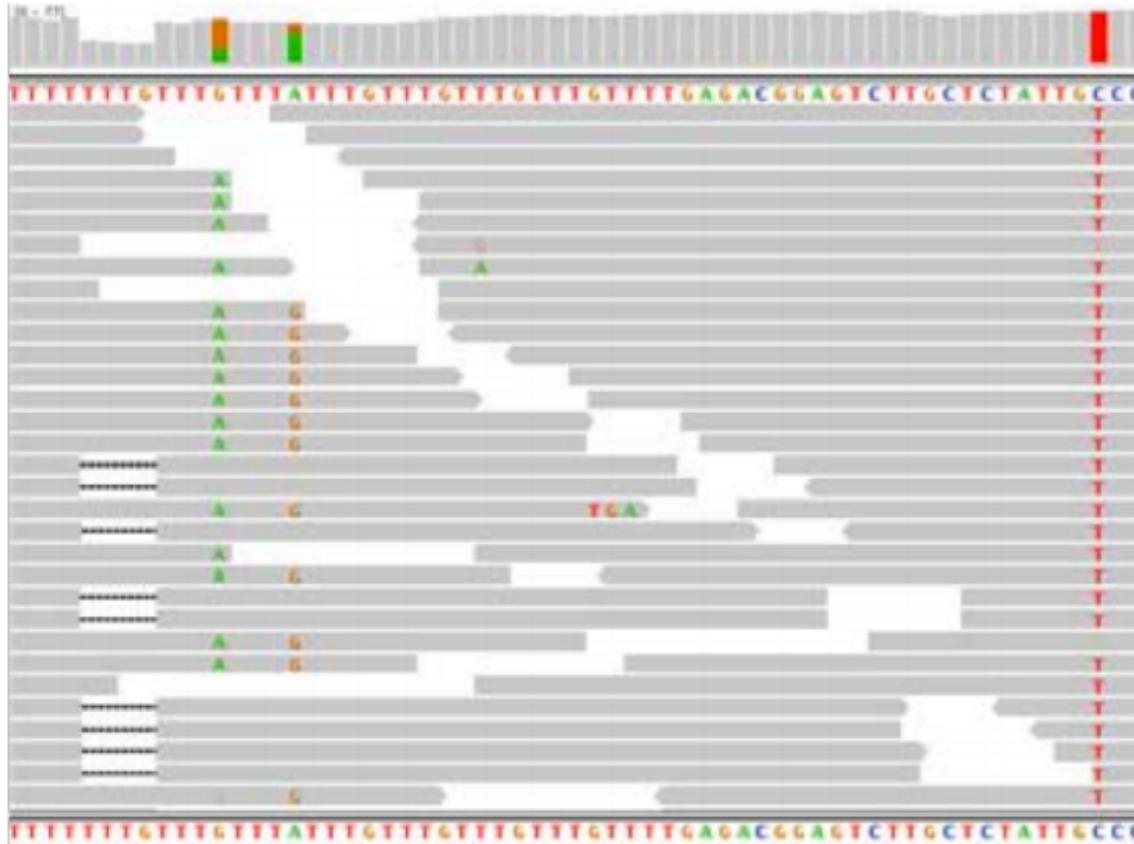
# Generalized Variant Calling Workflow

Add read groups and merge  
alignment files (optional)



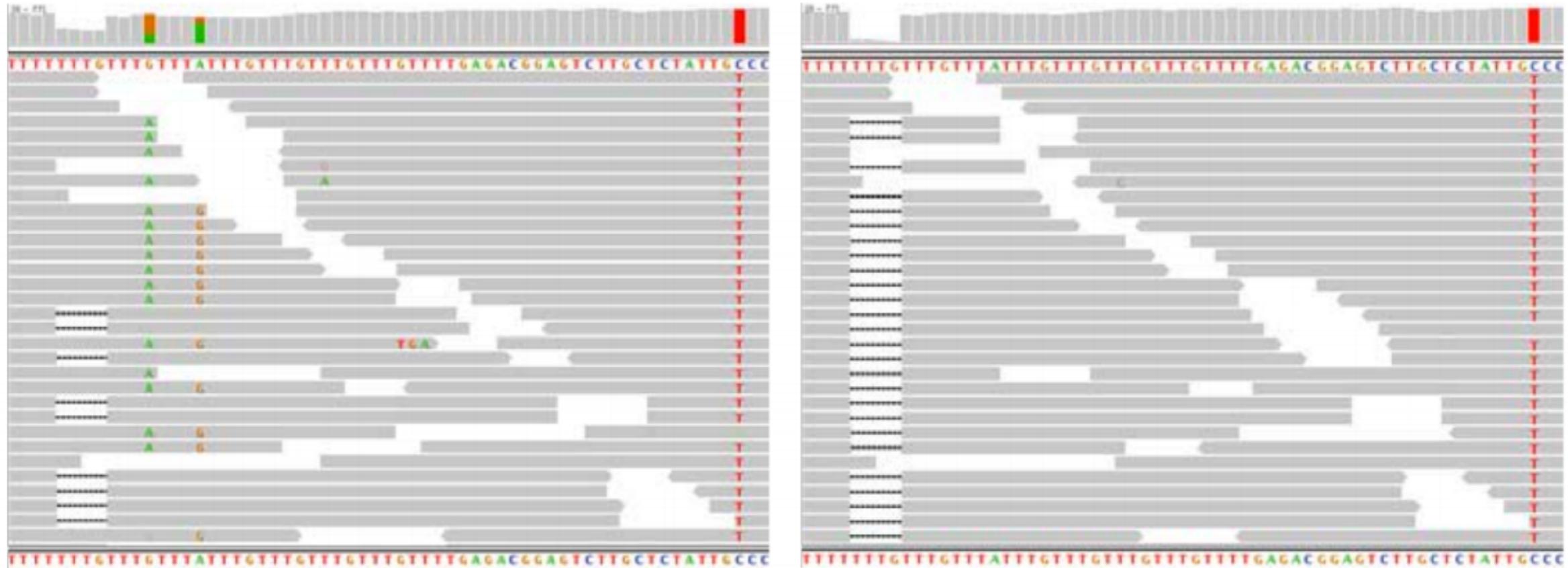
# Generalized Variant Calling Workflow

### Indel realignment (optional)



# Generalized Variant Calling Workflow

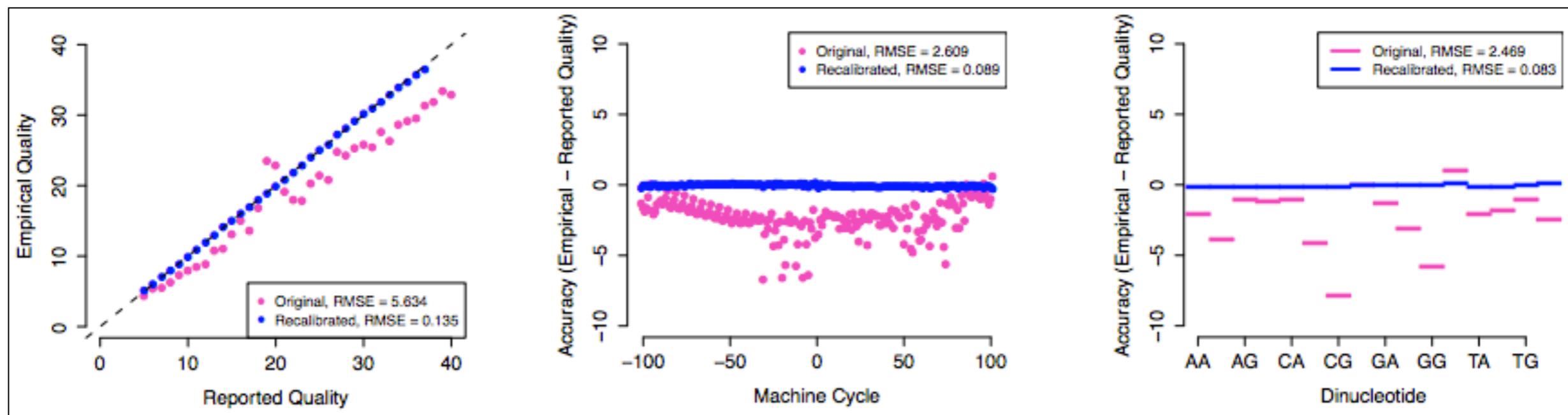
### Indel realignment (optional)



# Generalized Variant Calling Workflow

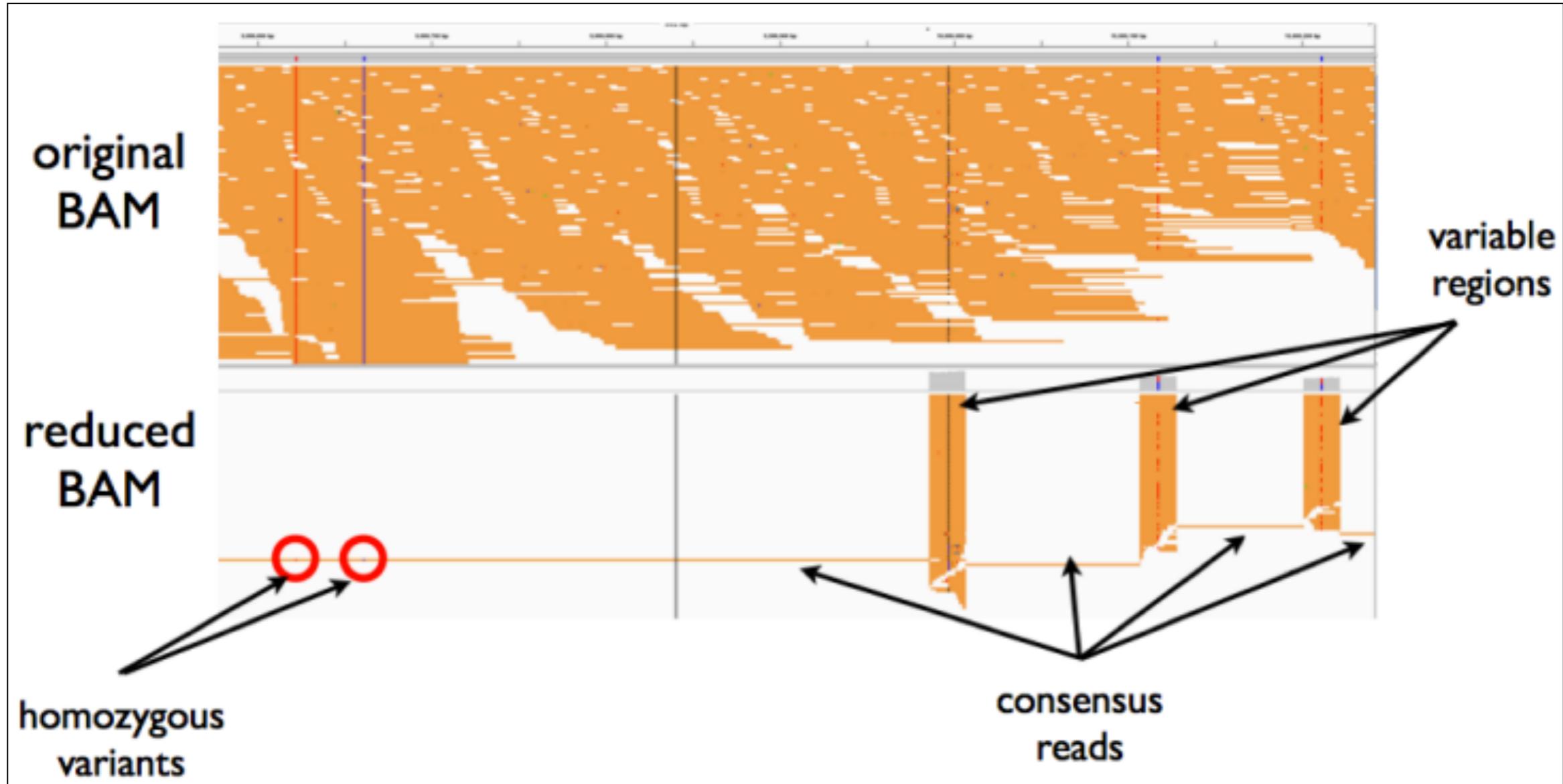
## Base Recalibration (optional)

This step removes any systematic biases the creep in during sequencing

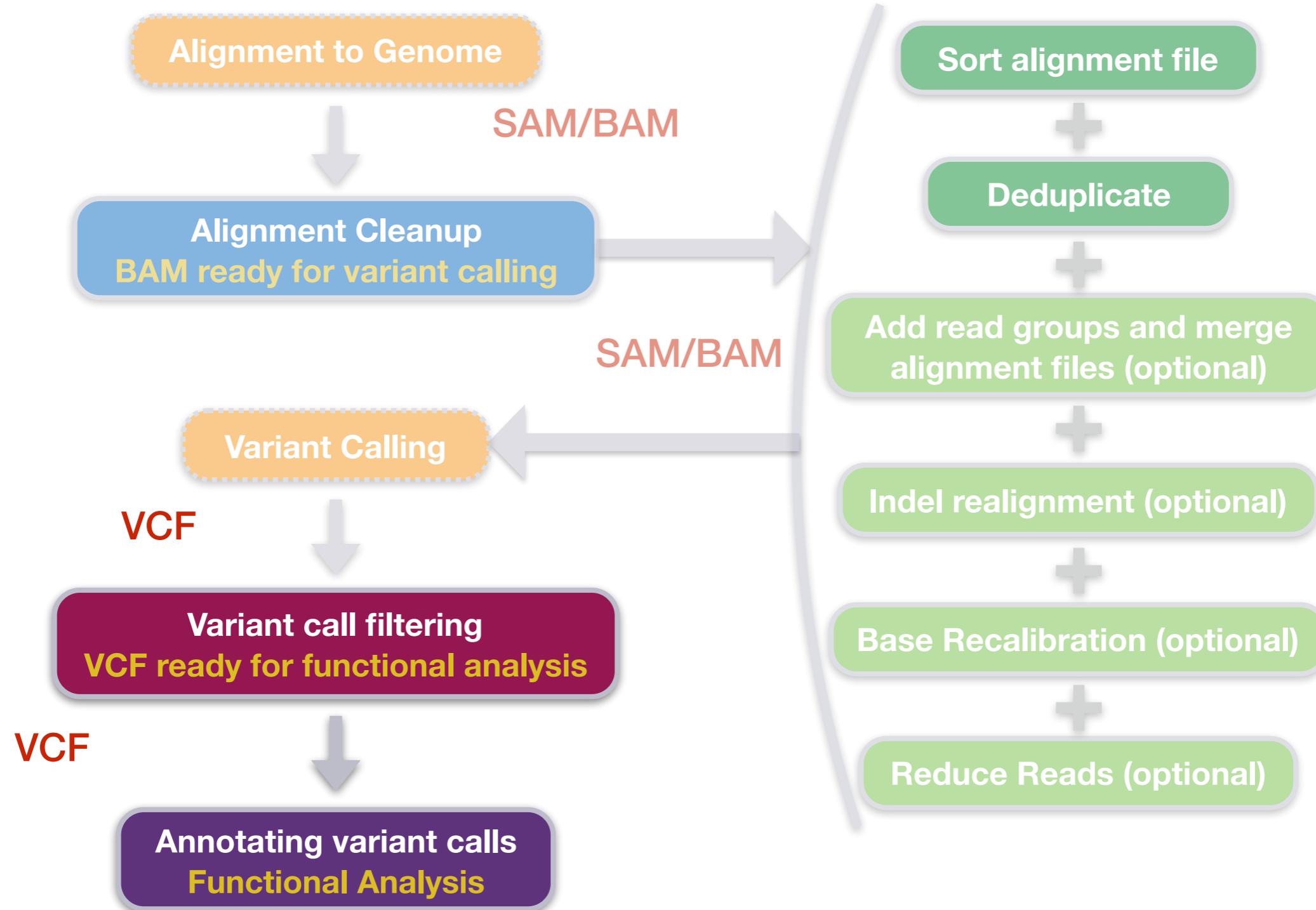


# Generalized Variant Calling Workflow

### Reduce Reads (optional)

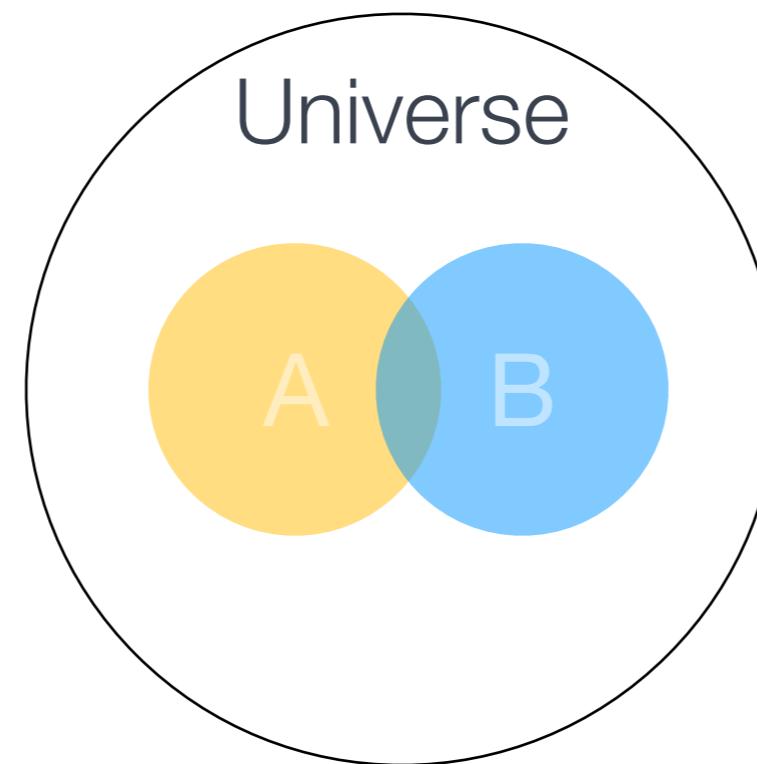


# Generalized Variant Calling Workflow



# Generalized Variant Calling Workflow

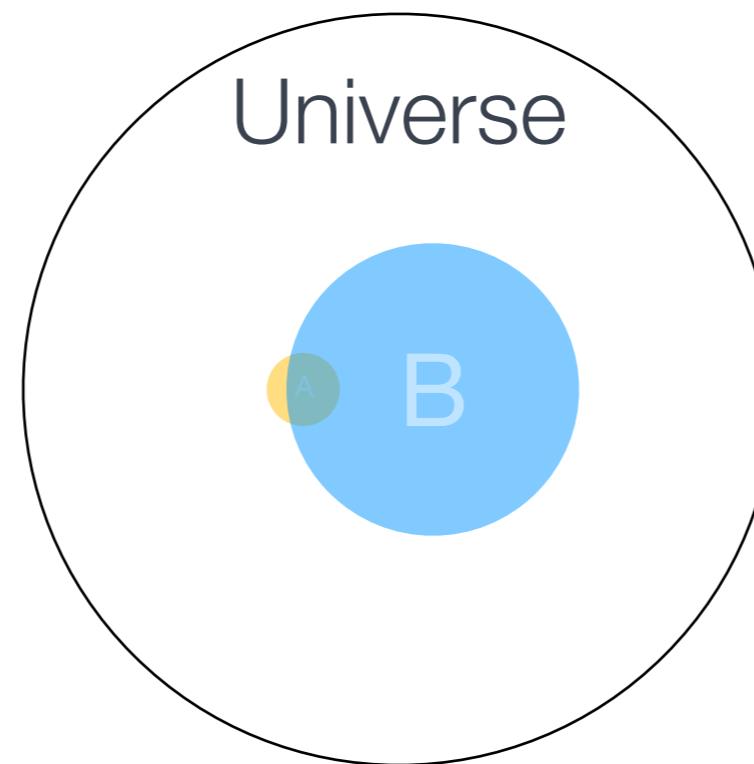
Probability of real variant “A” given observed variants “B” in alignments



In theory

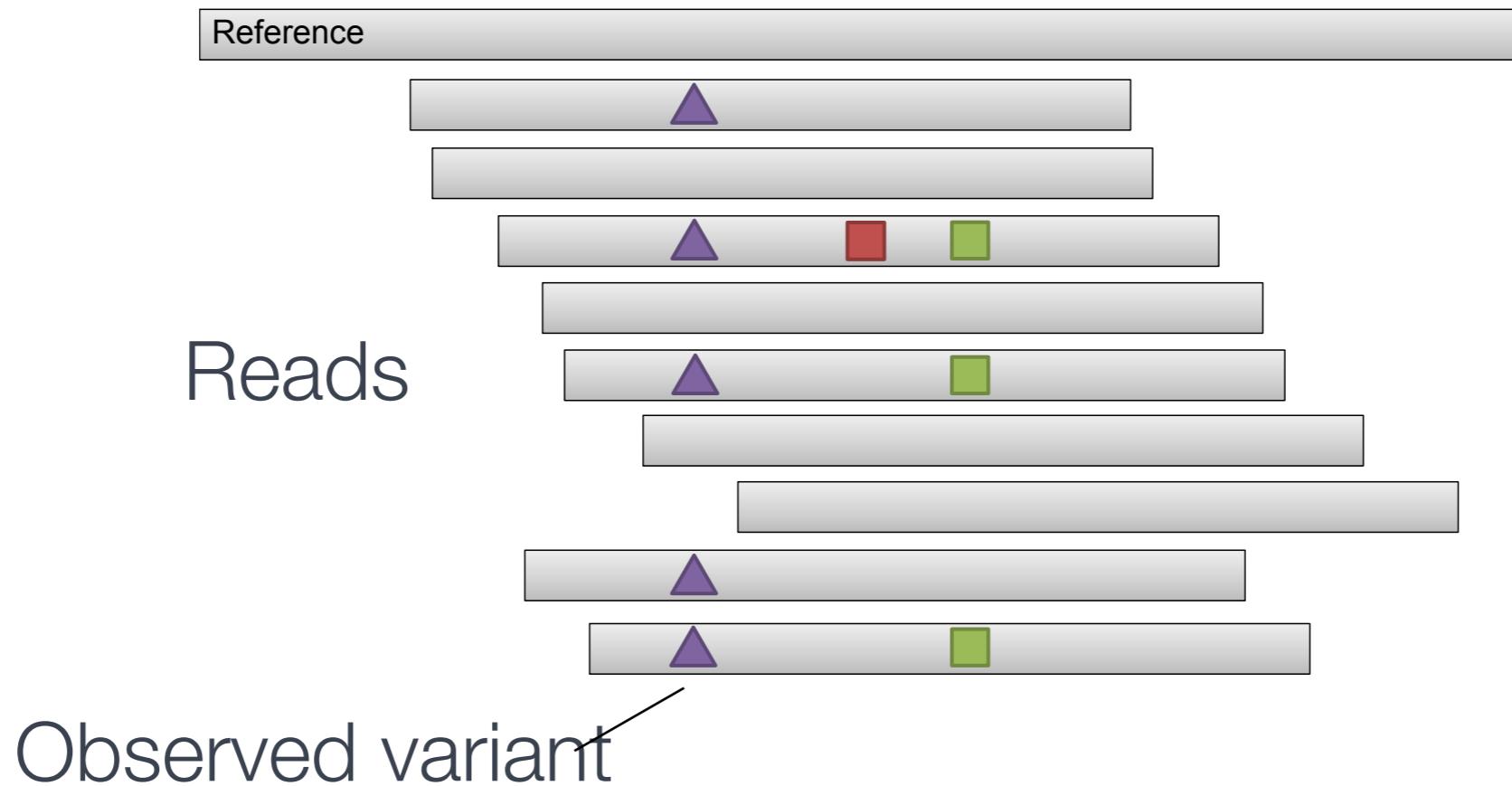
Variant detection

Probability of real variant “A” given observed variants “B” in alignments

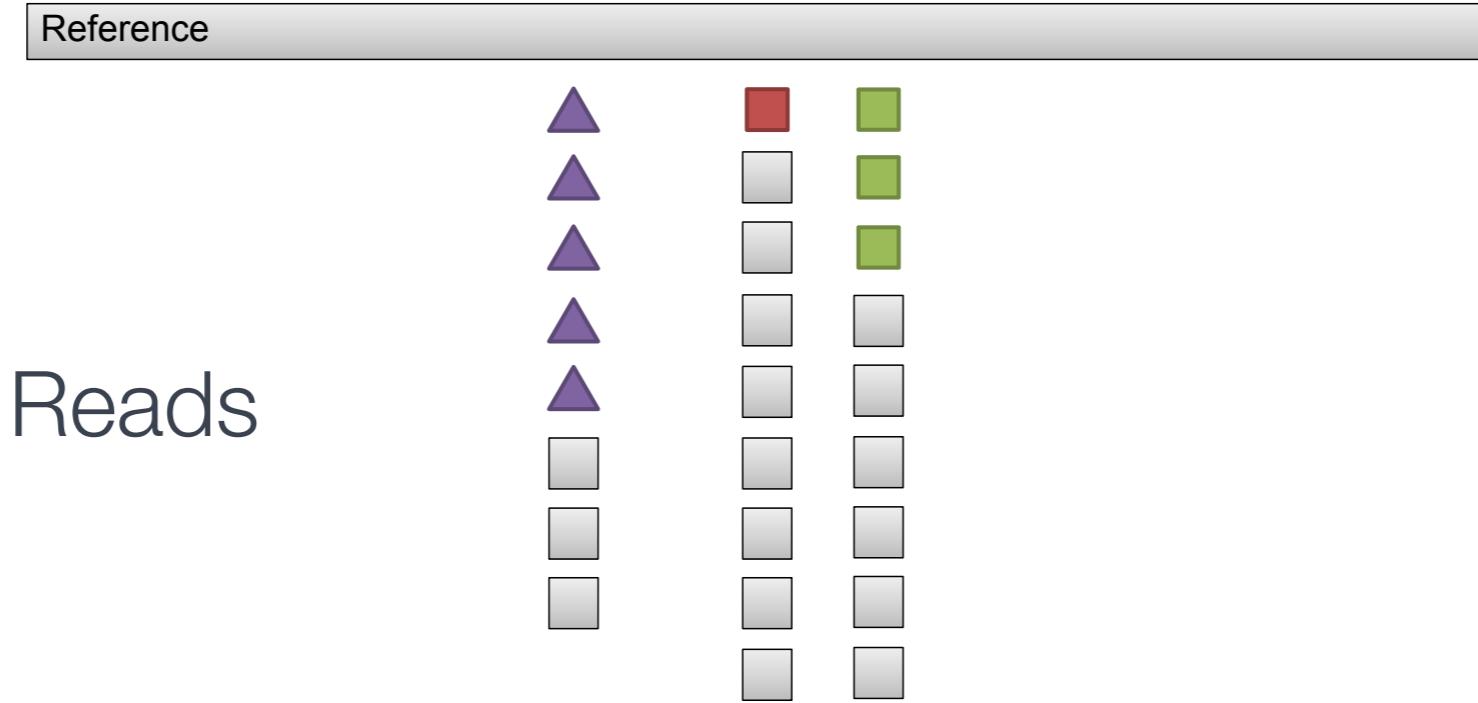


In practice

Variant detection

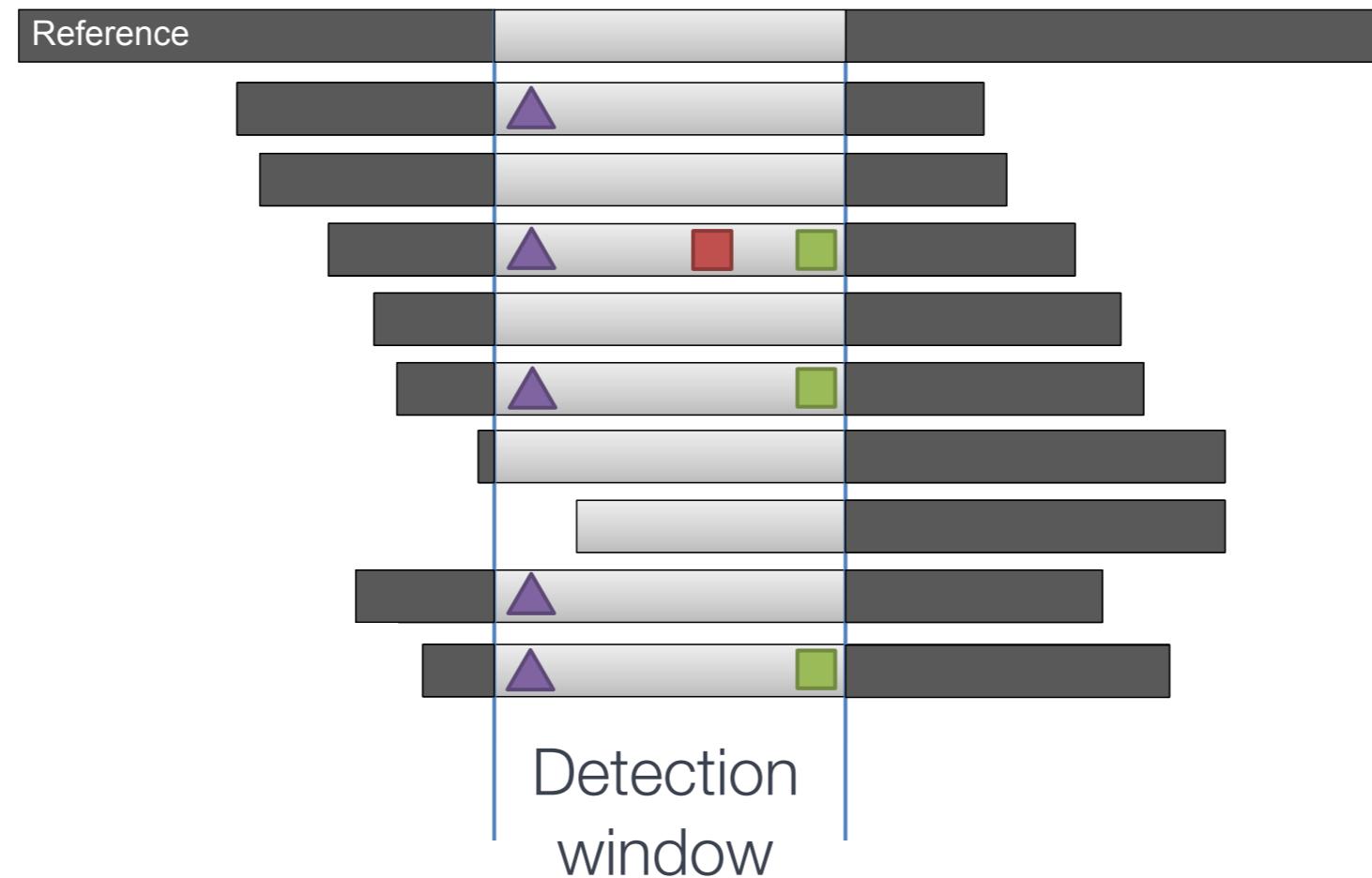


# Finding variants one at a time



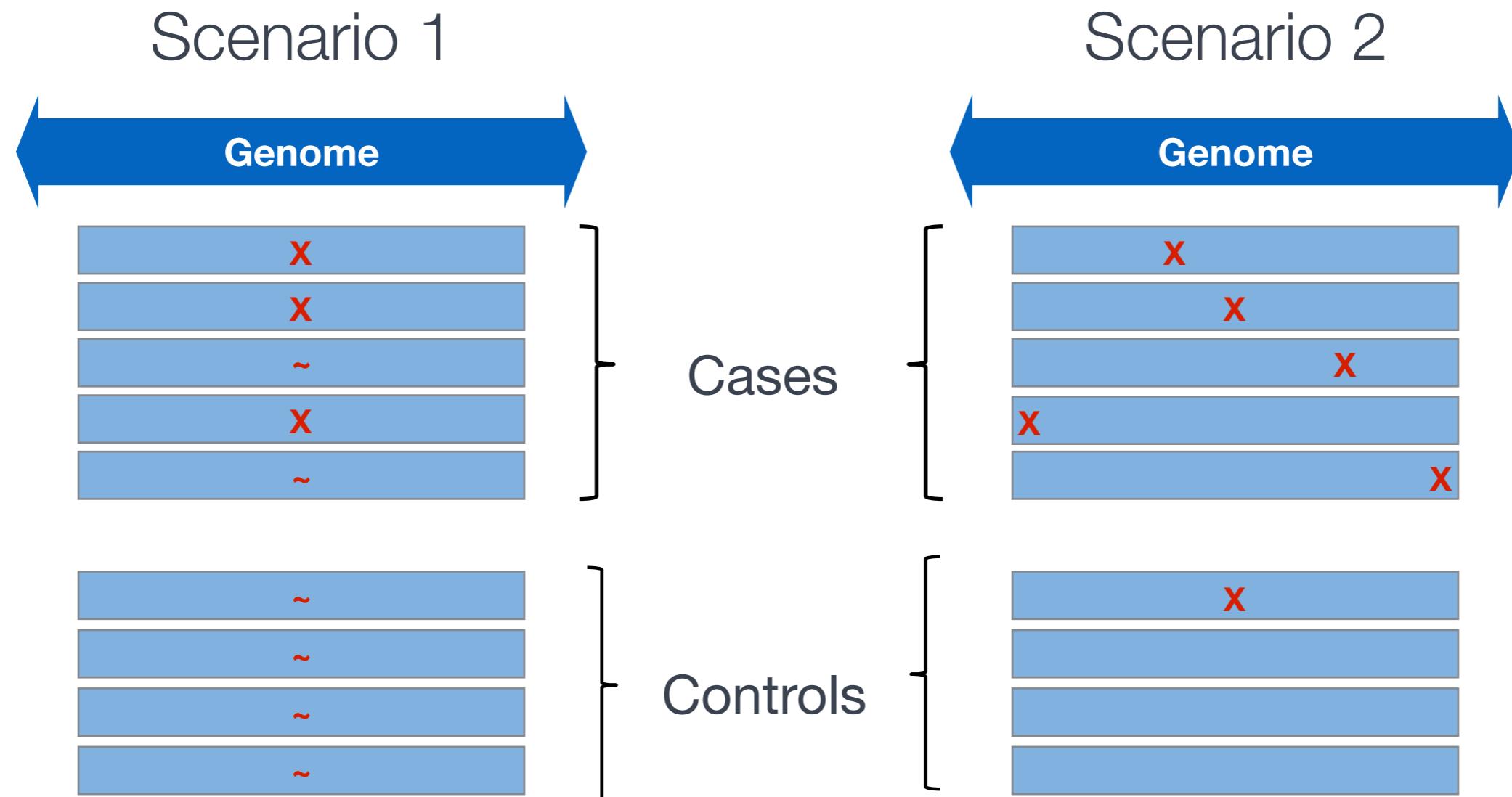
Haplotype information is lost

Finding variants one at a time

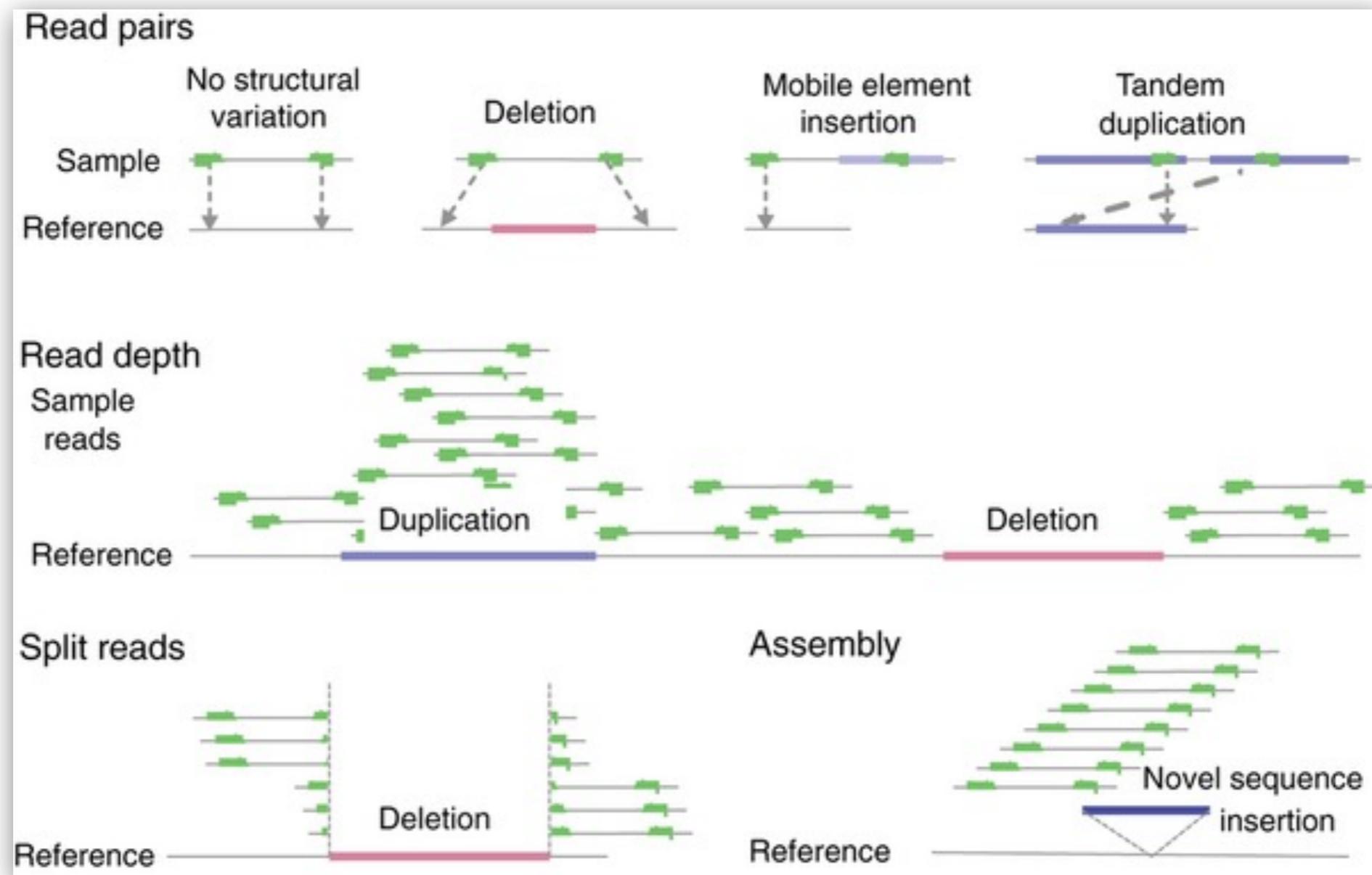


# Finding variants with FreeBayes (or GATK)

# How many samples does one need for a given population study?



# CNVs and SVs



CNV and SV detection is more complex

*These materials have been developed by members of the teaching team at the [Harvard Chan Bioinformatics Core \(HBC\)](#). These are open access materials distributed under the terms of the [Creative Commons Attribution license \(CC BY 4.0\)](#), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.*

