

Introduction to NGS Variant-Calling

Thibault Latrille, Anouk Necsulea,
Annabelle Haudry, Maud Gautier



Plan

Part 1:
Overview of variant-calling

Part 2:
GATK workflow



Different types of variants

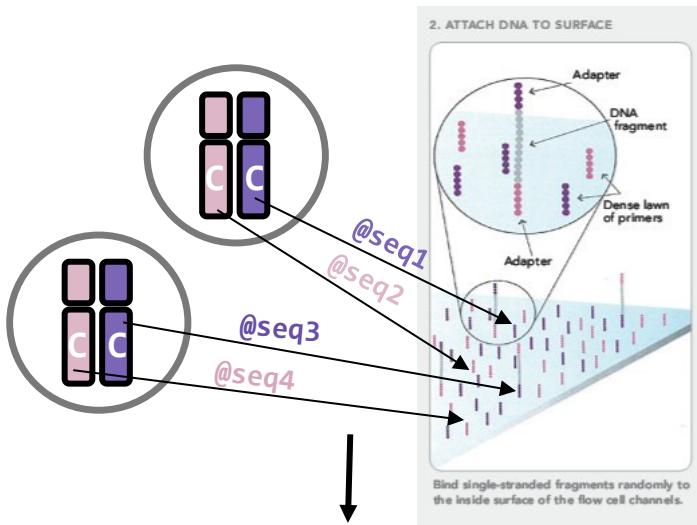
	Substitution	Insertion	Deletion	Indel
Wild-Type:	AACGGCC T GTAAC	AACGGCC T GTAAC	AACGGCC T GTAAC	AACGGCC T GTAAC
Mutant:	AACGGCC A GTAAC	AACGGCC A G C TTAAC	AACGGCC - GTAAC	AACGGCC C TTAAC
	Deletion	Insertion	Translocation	
Wild-Type:				
Mutant:				
	Duplication		Inversion	
Wild-Type:				
Mutant:				
Individual 1:	AACGGCC T GTAAC		Individual 7:	AACGGCC T GTAAC
Individual 2:	AACGGCC T GTAAC		Individual 8:	AACGGCC T GTAAC
Individual 3:	AACGGCC T GTAAC		Individual 9:	AACGGCC T GTAAC
Individual 4:	AACGGCC A GTAAC		Individual 10:	AACGGCC A GTAAC
Individual 5:	AACGGCC T GTAAC		Individual 11:	AACGGCC T GTAAC
Individual 6:	AACGGCC A GTAAC		Individual 12:	AACGGCC A GTAAC

Cardoso et al (2015)
DOI: 10.3389/fbioe.2015.00013



Concrete view: Variants on a flowcell

- Scenario 1: Homozygote



Ref ATCGGG**T**ACCATCCAATCATTACC

GGCACC**AT**CCAAT

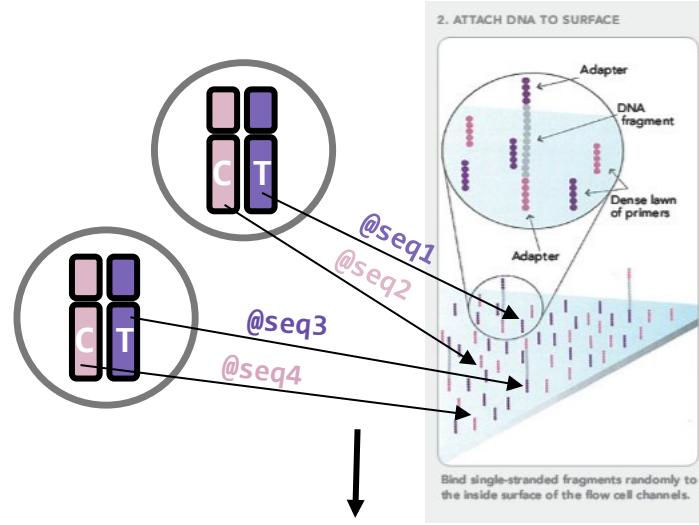
@seq2 ATCGGG**C**ACC**A**T

@seq3 GGGCACC**A**ATCCAA
TCGGG**C**ACC**A**TC

@seq4 CGGG**C**ACC**A**TC

@seq1 CGGGCACC**A**TC

- Scenario 2: Heterozygote



Ref ATCGGG**T**ACCATCCAATCATTACC

GG**T**ACC**AT**CCAAT

@seq2 ATCGGG**C**ACC**A**T

@seq3 GGG**T**ACC**A**ATCCAA
TCGGG**C**ACC**A**TC

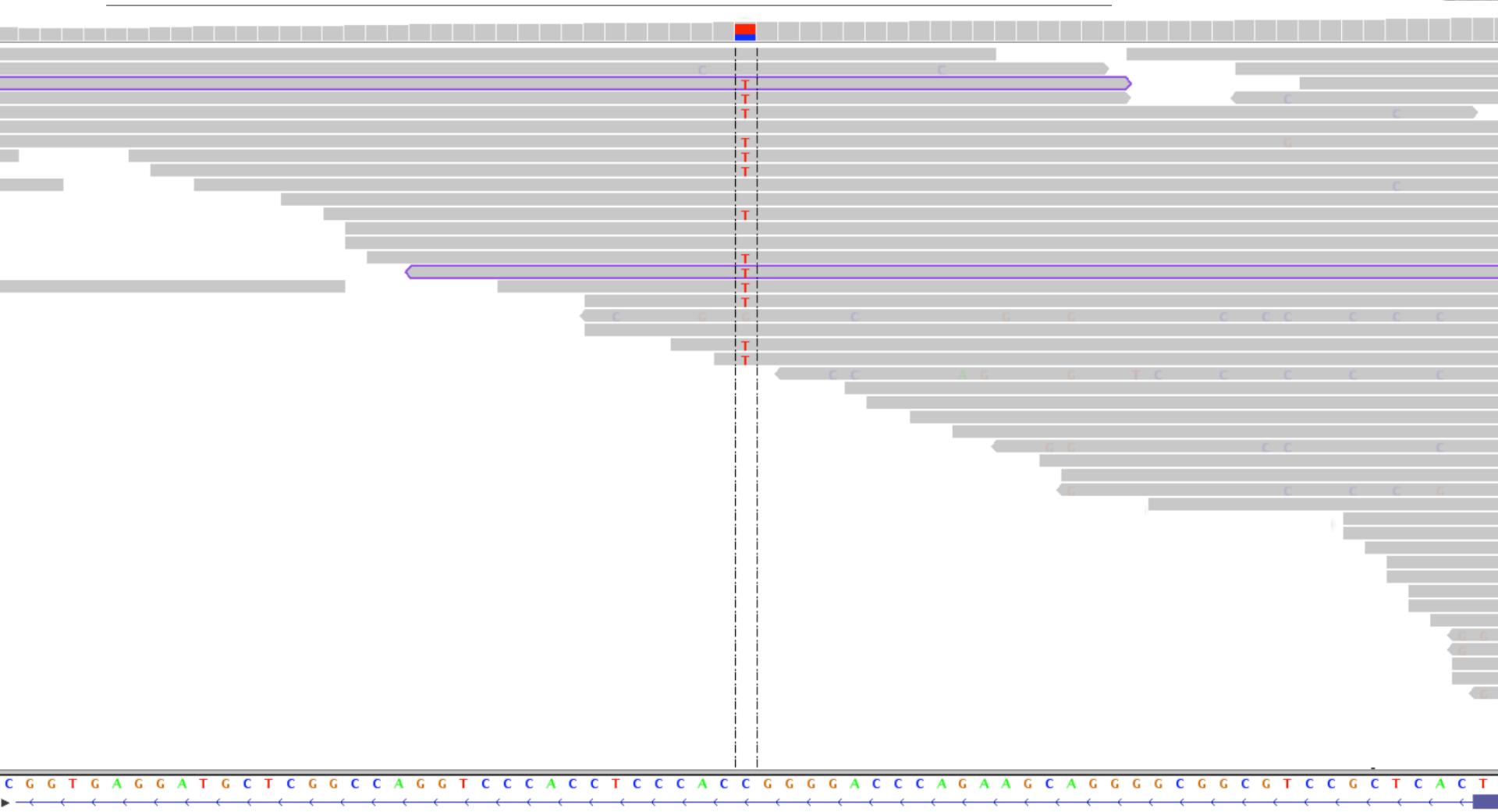
@seq4 CGGG**C**ACC**A**TC

@seq1 CGGG**T**ACC**A**TC

From Aaron Quinlab



Visualising variants in IGV





Complexity of variant-calling

- Identify variants **relative to a reference** (hg38 for humans)
- Complexity of variant-calling: false positives => need to distinguish between true genetic variation and false positives
- False positives may come from:
 - PCR artifacts (→ MarkDuplicates)
 - Sequencing artifacts (→ Remove if low quality)
 - Alignment (→ Locally realign)
 - ...



Plan

Part 1:
Overview of variant-calling

Part 2:
GATK workflow



GATK workflow

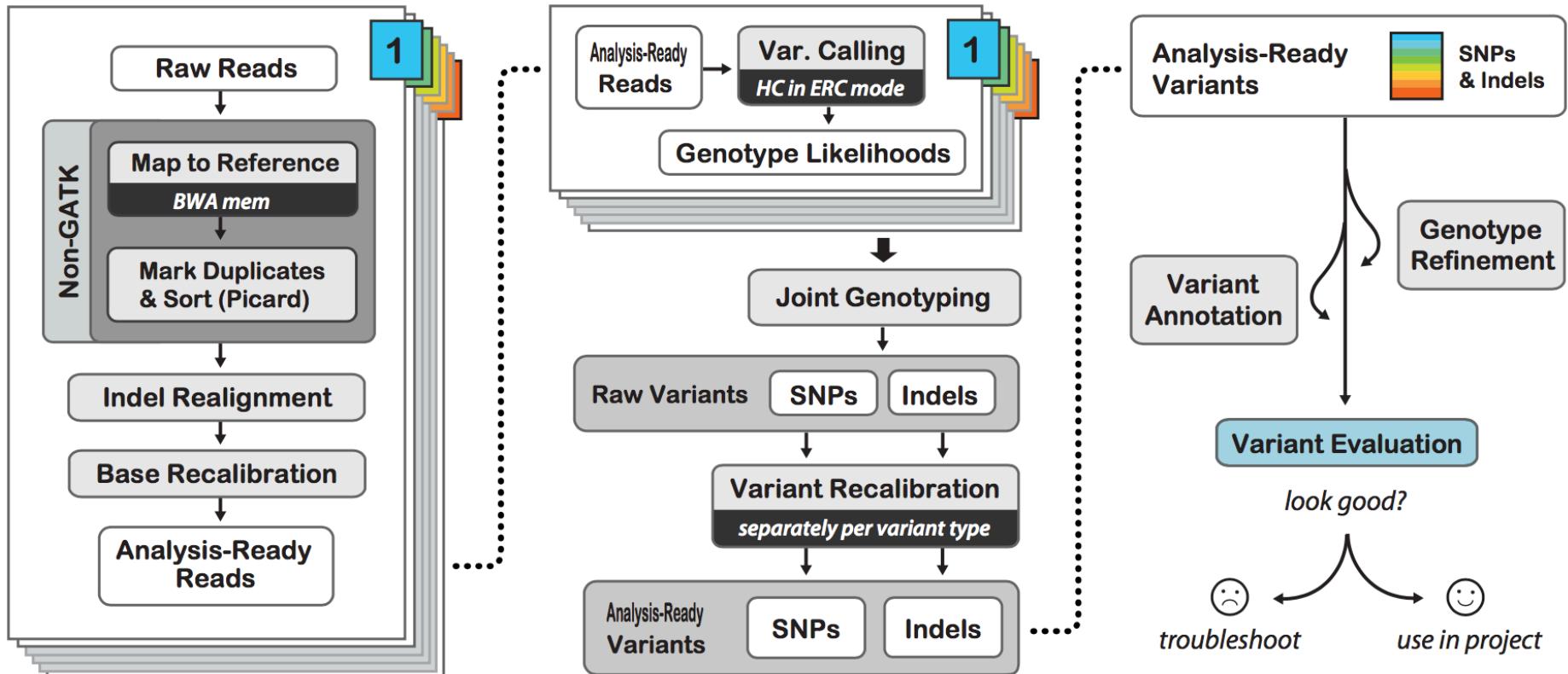
Data Pre-processing

>>

Variant Discovery

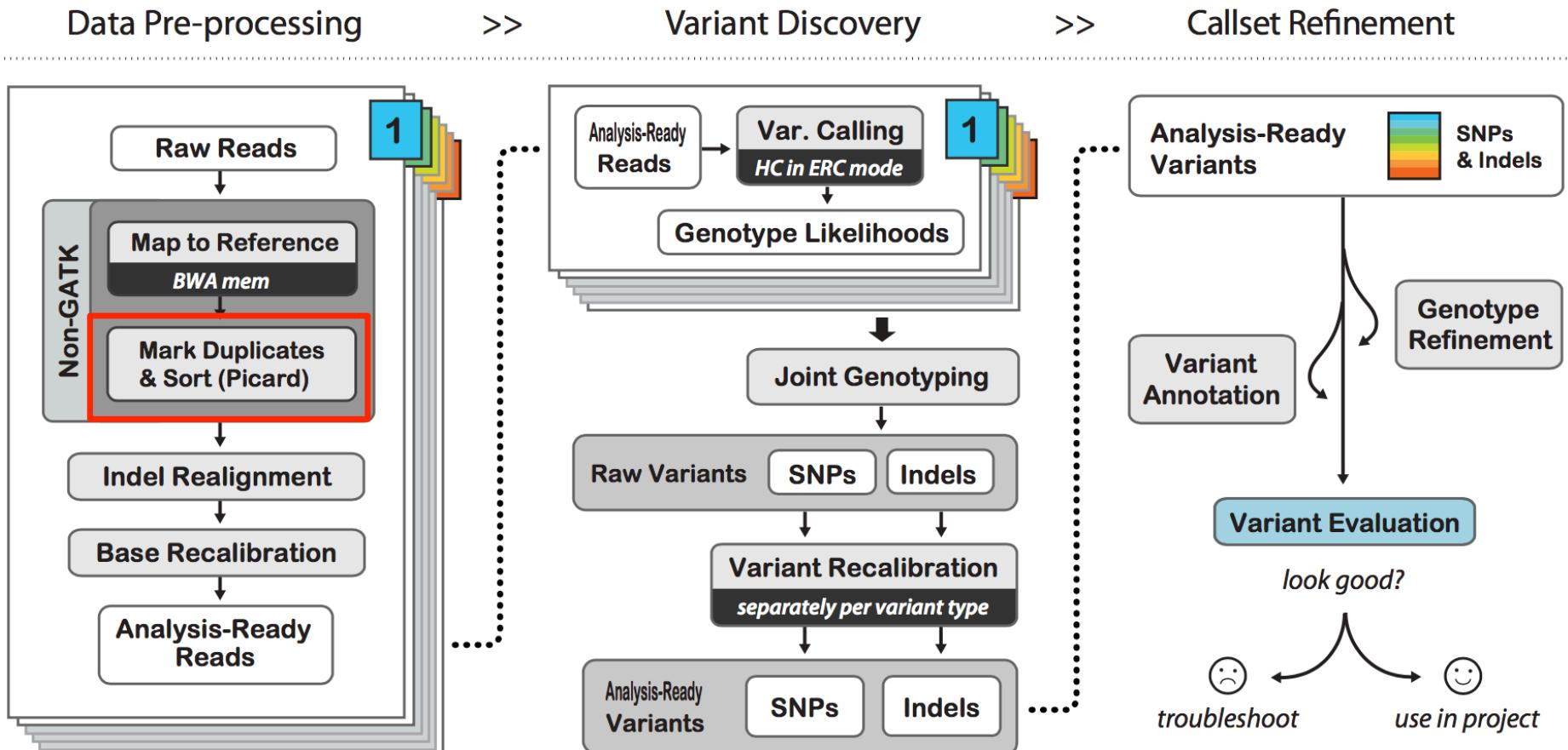
>>

Callset Refinement



GATK Best Practices for Variant Discovery
<https://software.broadinstitute.org/gatk/download/workshops>

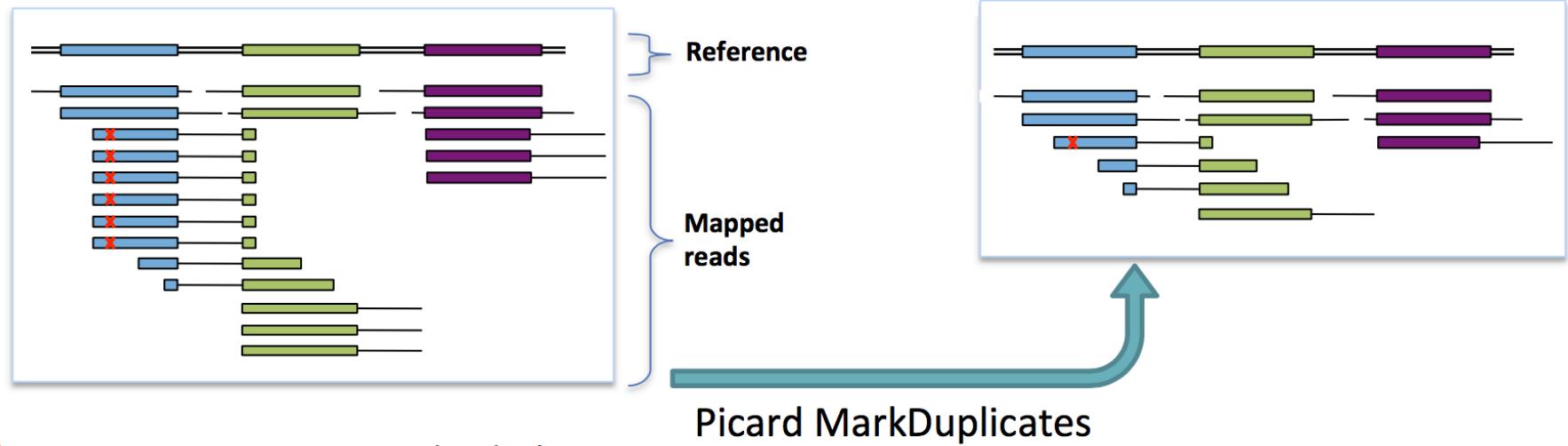
GATK workflow





Data pre-processing - MarkDuplicates

- Duplicates = non-independent measurements of a sequence
→ Must be removed



GATK Best Practices for Variant Discovery
(<https://software.broadinstitute.org/gatk/download/workshops>)

GATK workflow

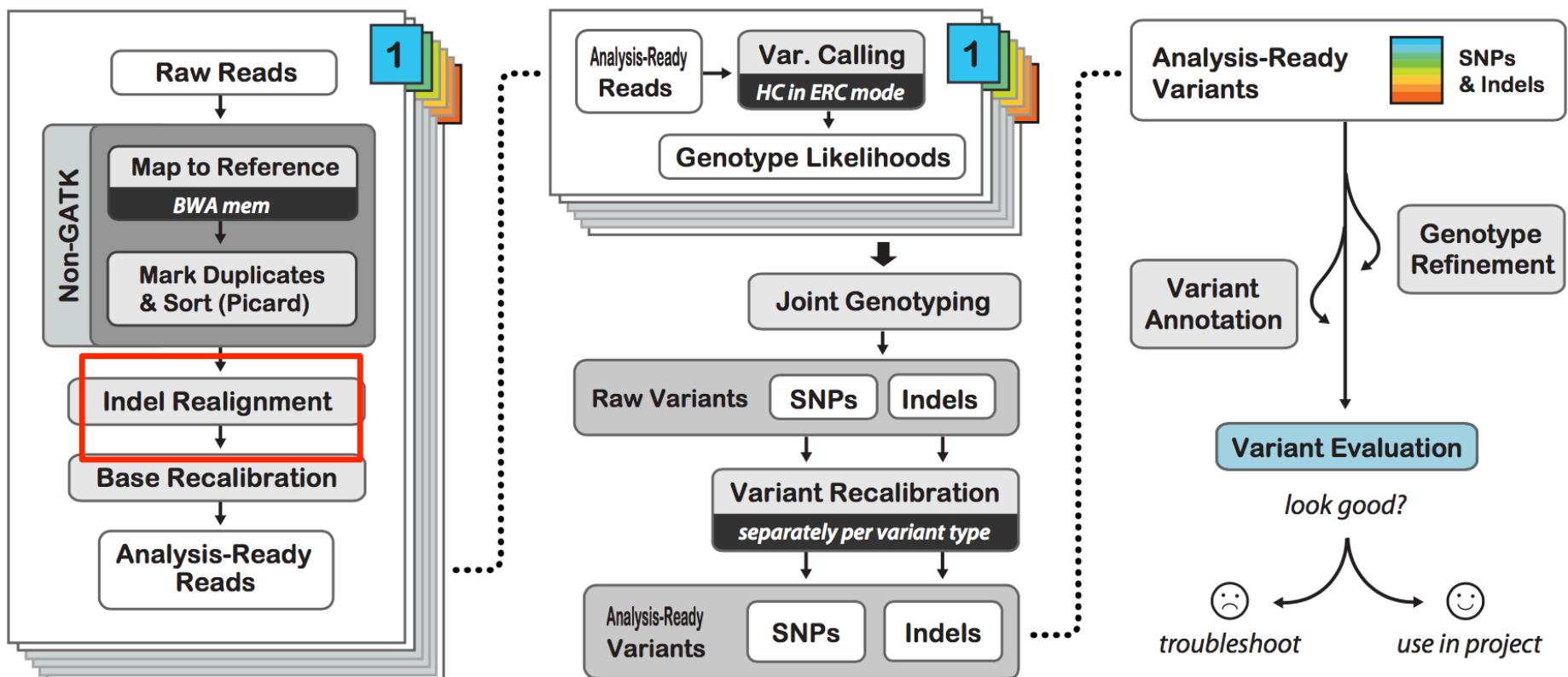
Data Pre-processing

>>

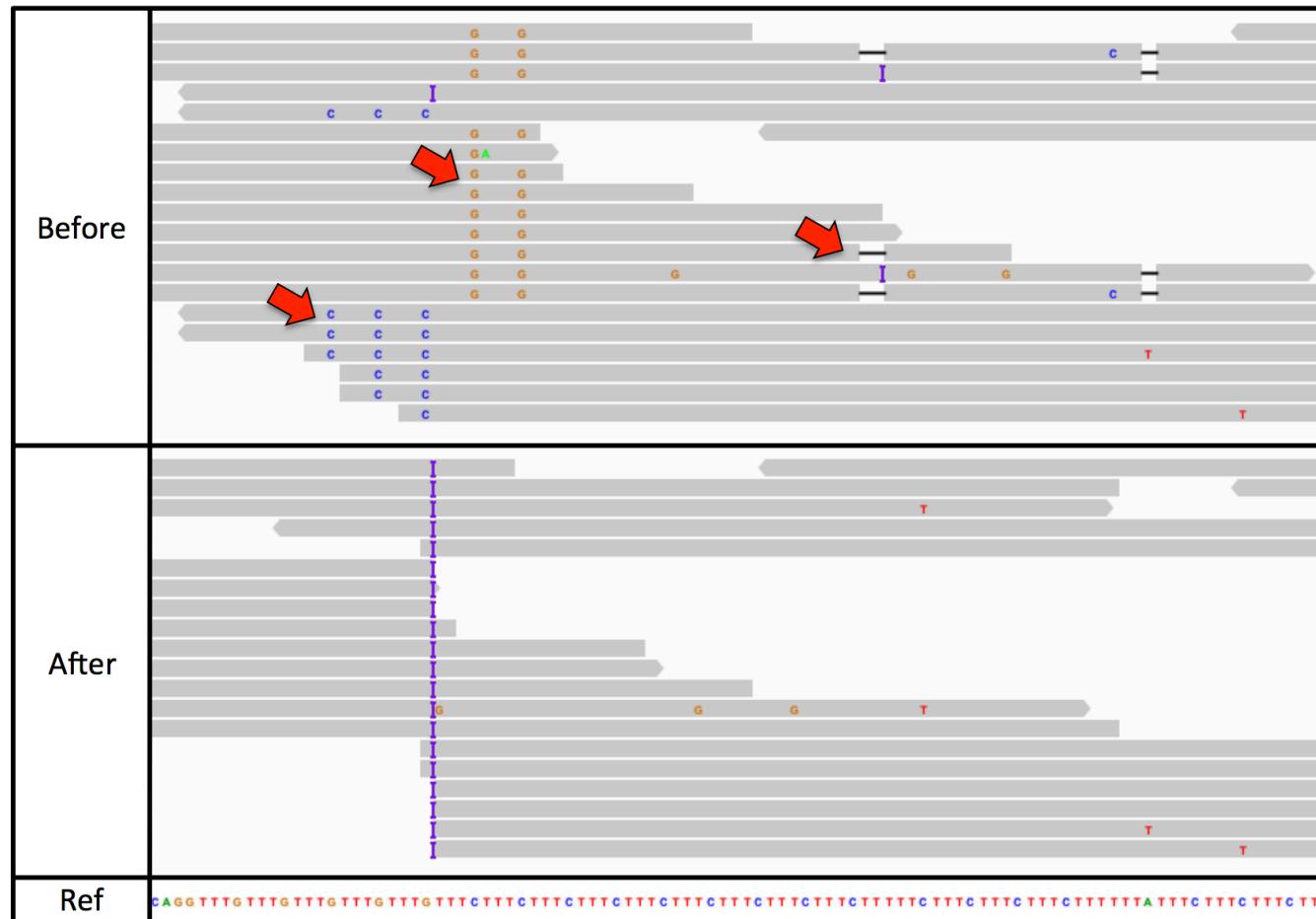
Variant Discovery

>>

Callset Refinement

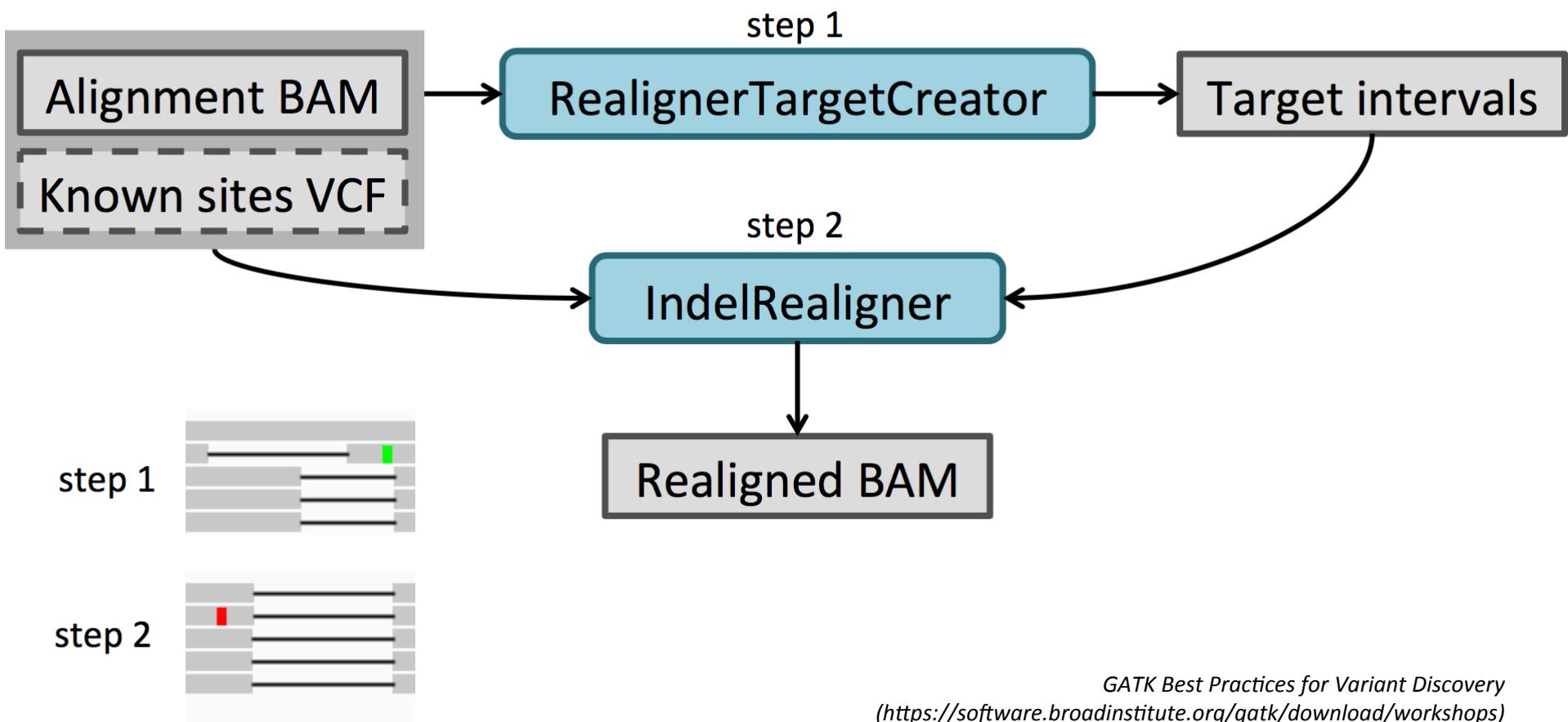


Indel realignment

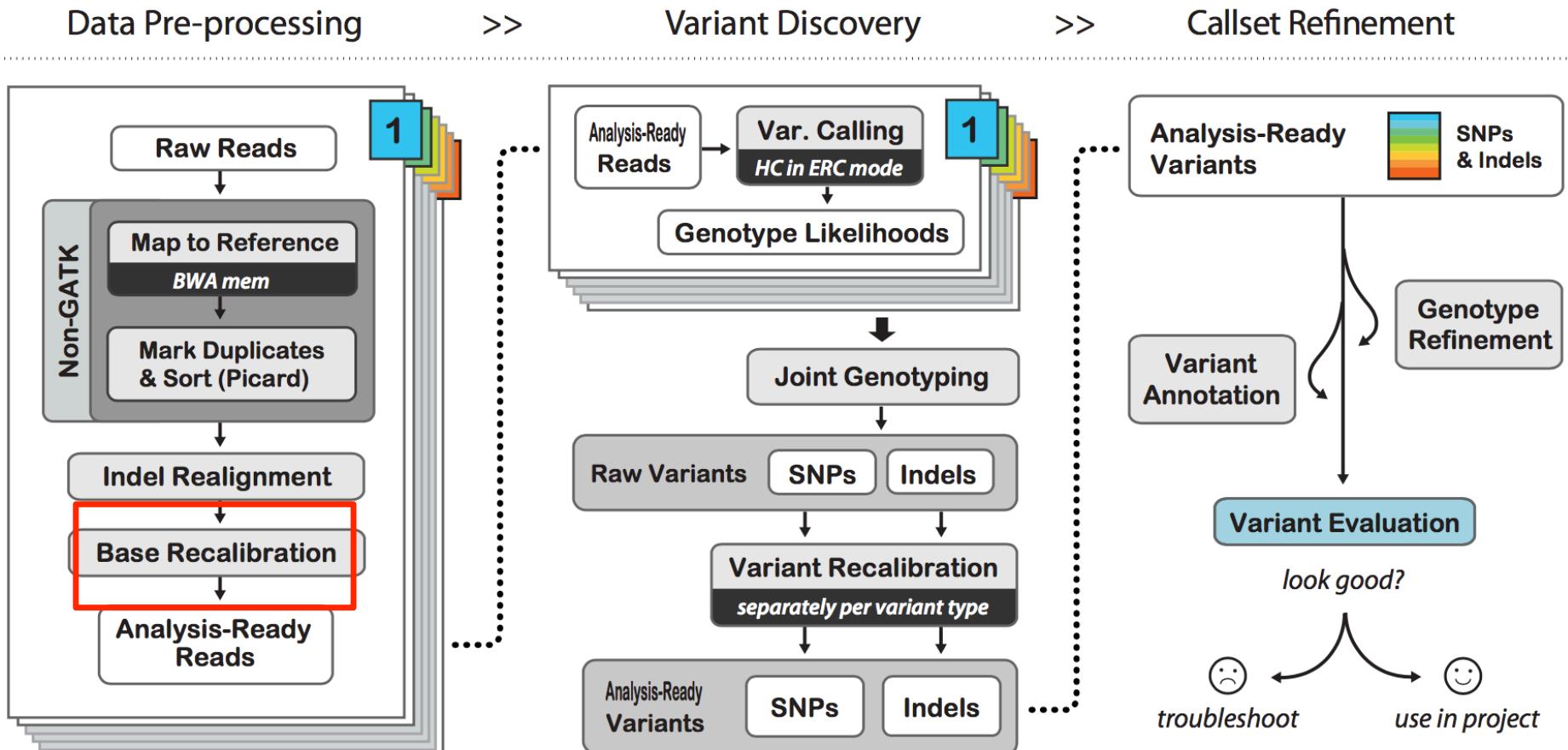


GATK Best Practices for Variant Discovery
(<https://software.broadinstitute.org/gatk/download/workshops>)

Indel realignment workflow



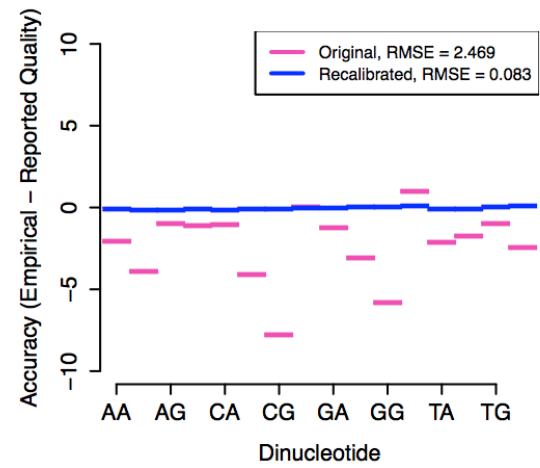
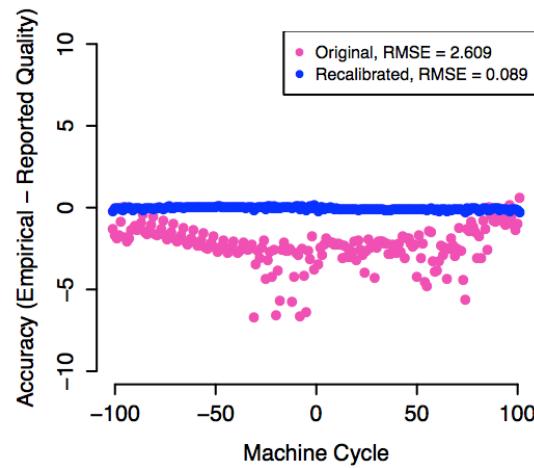
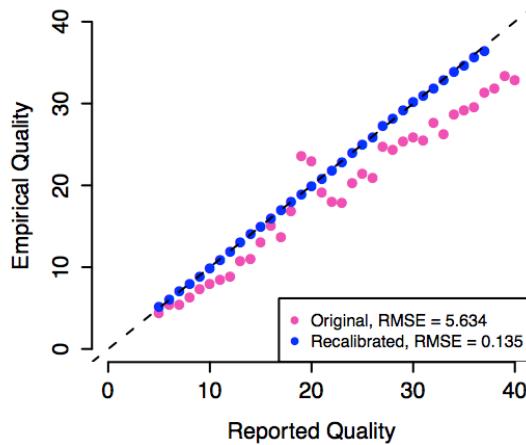
GATK workflow





Base Quality Score Recalibration (BQSR)

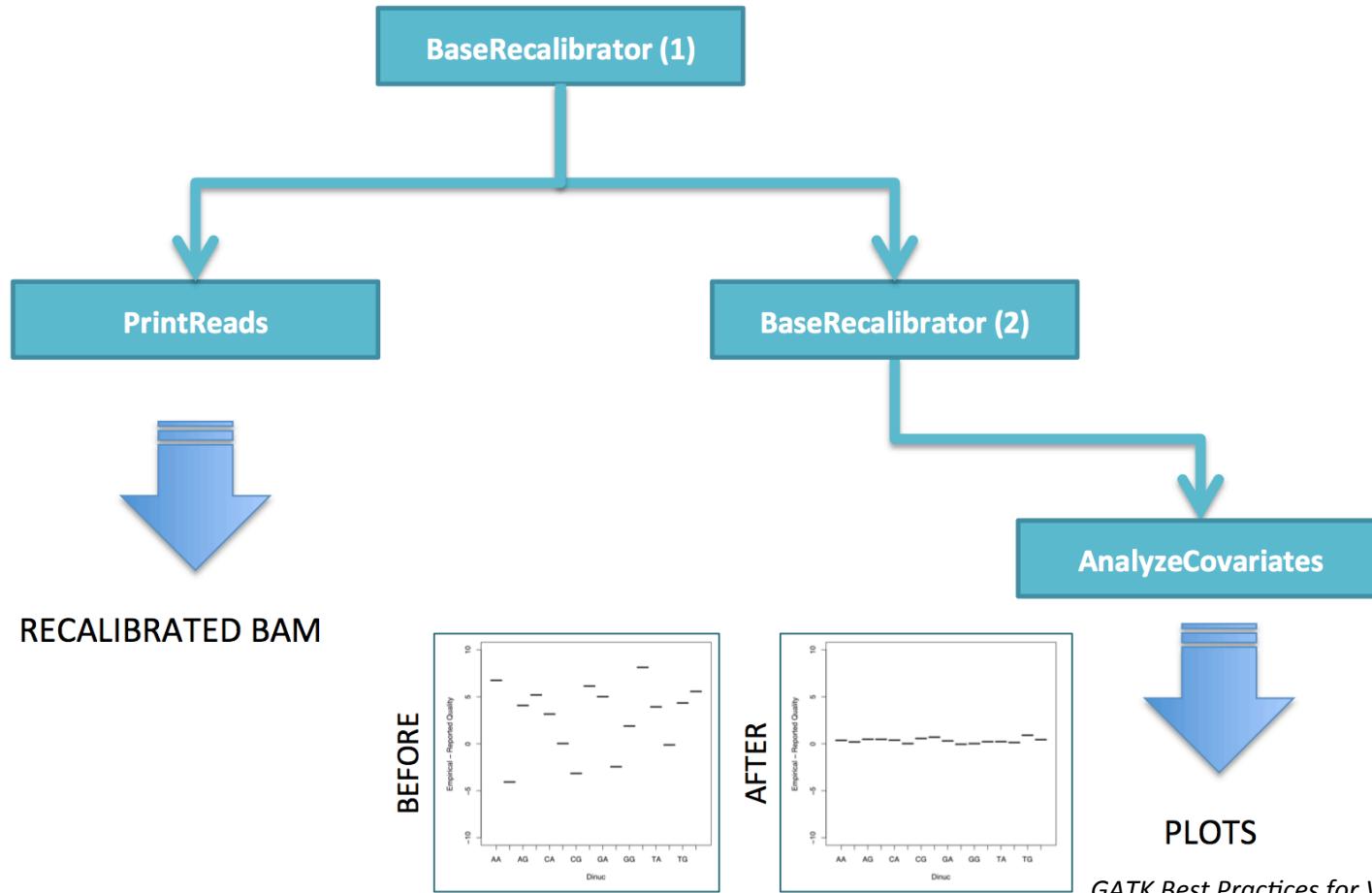
- Systematic biases: quality of base calls depend on nucleotidic context, machine cycle...
- → Machine learning algorithm to find how error varies with basecall features. Then, apply recalibration.



GATK Best Practices for Variant Discovery
<https://software.broadinstitute.org/gatk/download/workshops>



BQSR workflow



GATK Best Practices for Variant Discovery

(<https://software.broadinstitute.org/gatk/download/workshops>)

GATK workflow

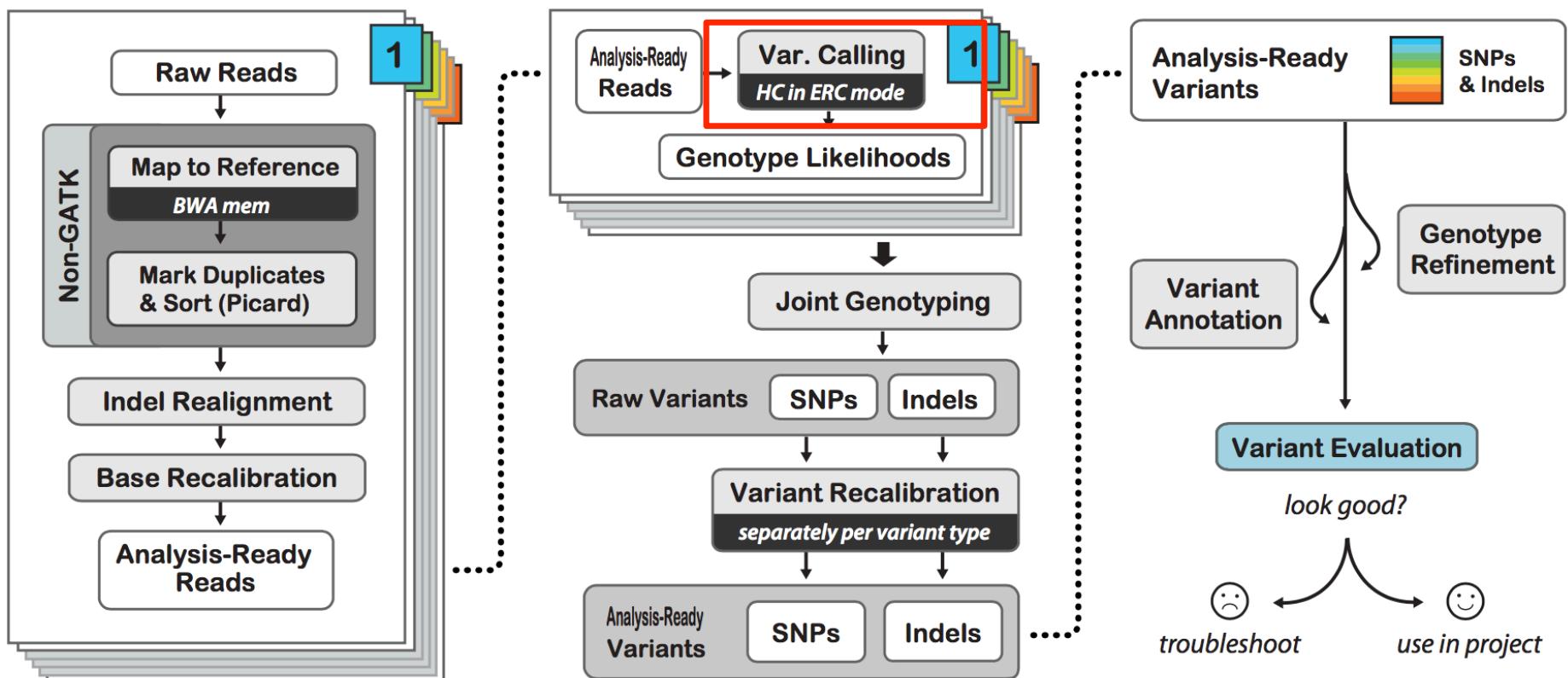
Data Pre-processing

>>

Variant Discovery

>>

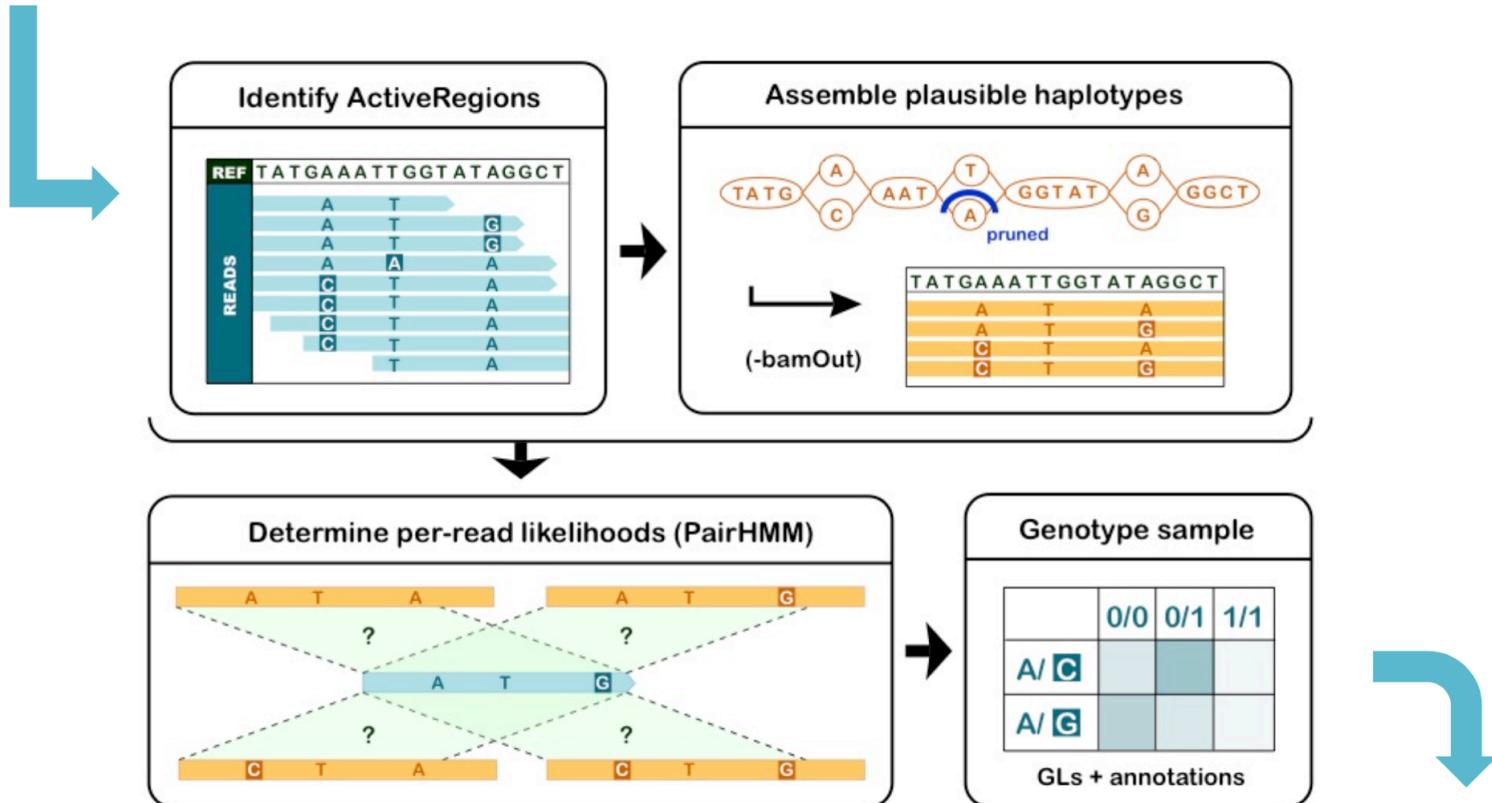
Callset Refinement





Variant-calling: Haplotype-Caller

BAM



VCF & index

GATK Best Practices for Variant Discovery
<https://software.broadinstitute.org/gatk/download/workshops>

GATK workflow

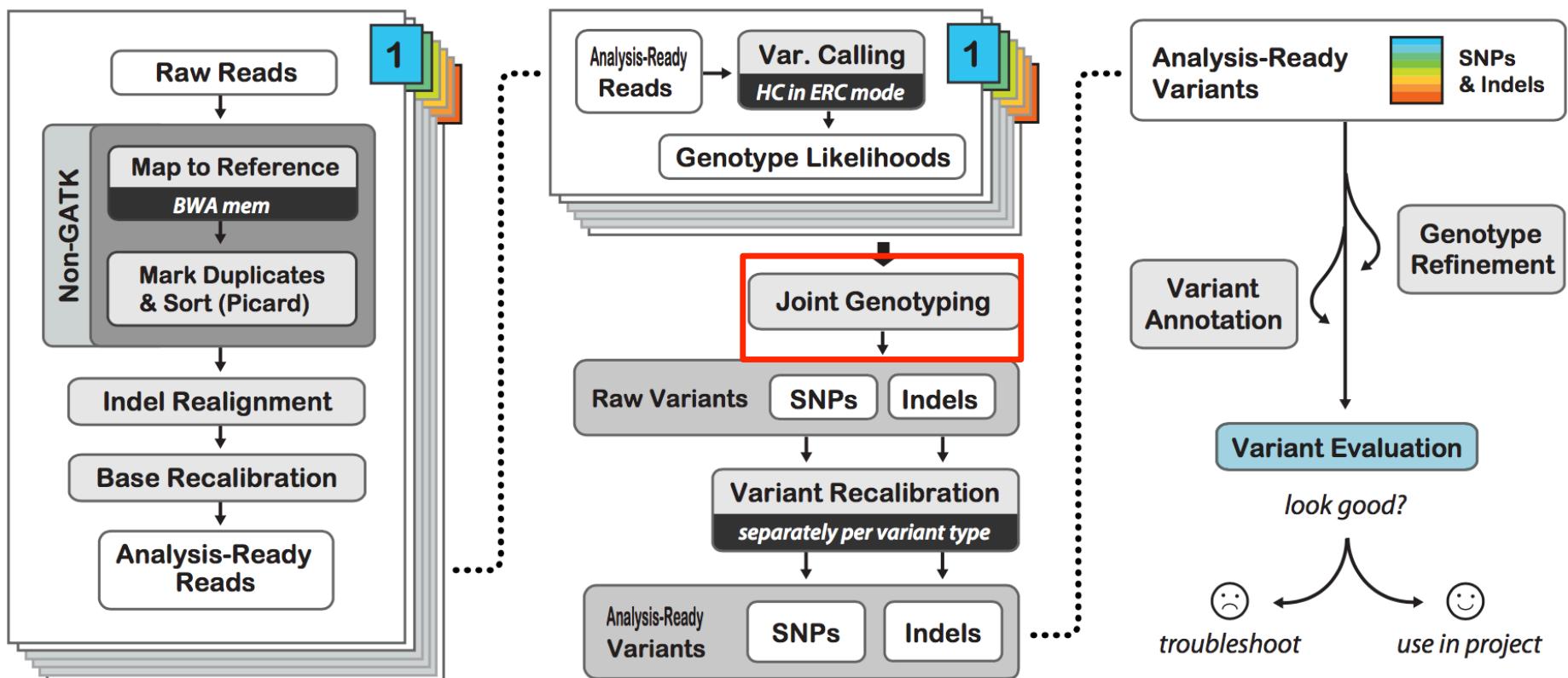
Data Pre-processing

>>

Variant Discovery

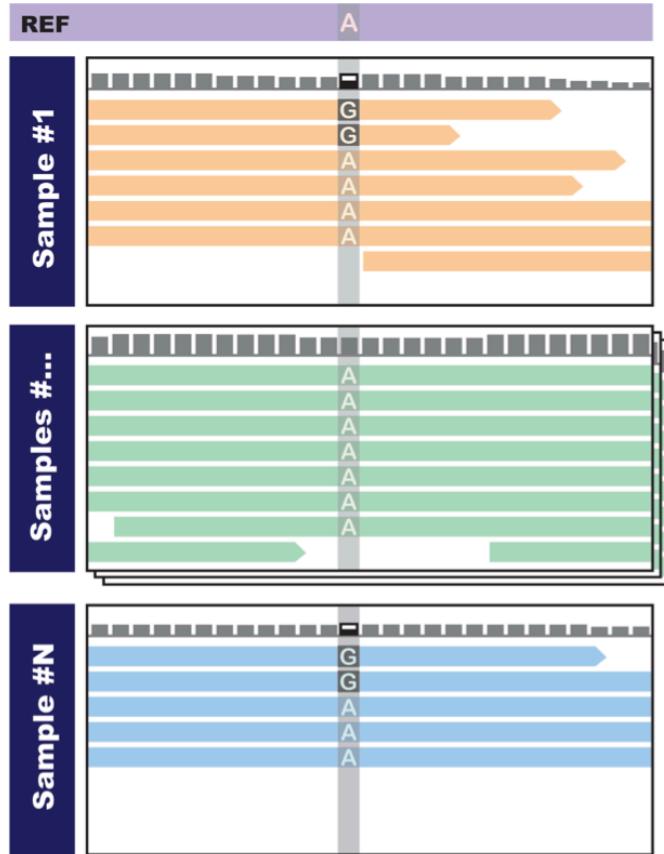
>>

Callset Refinement





Joint-genotyping



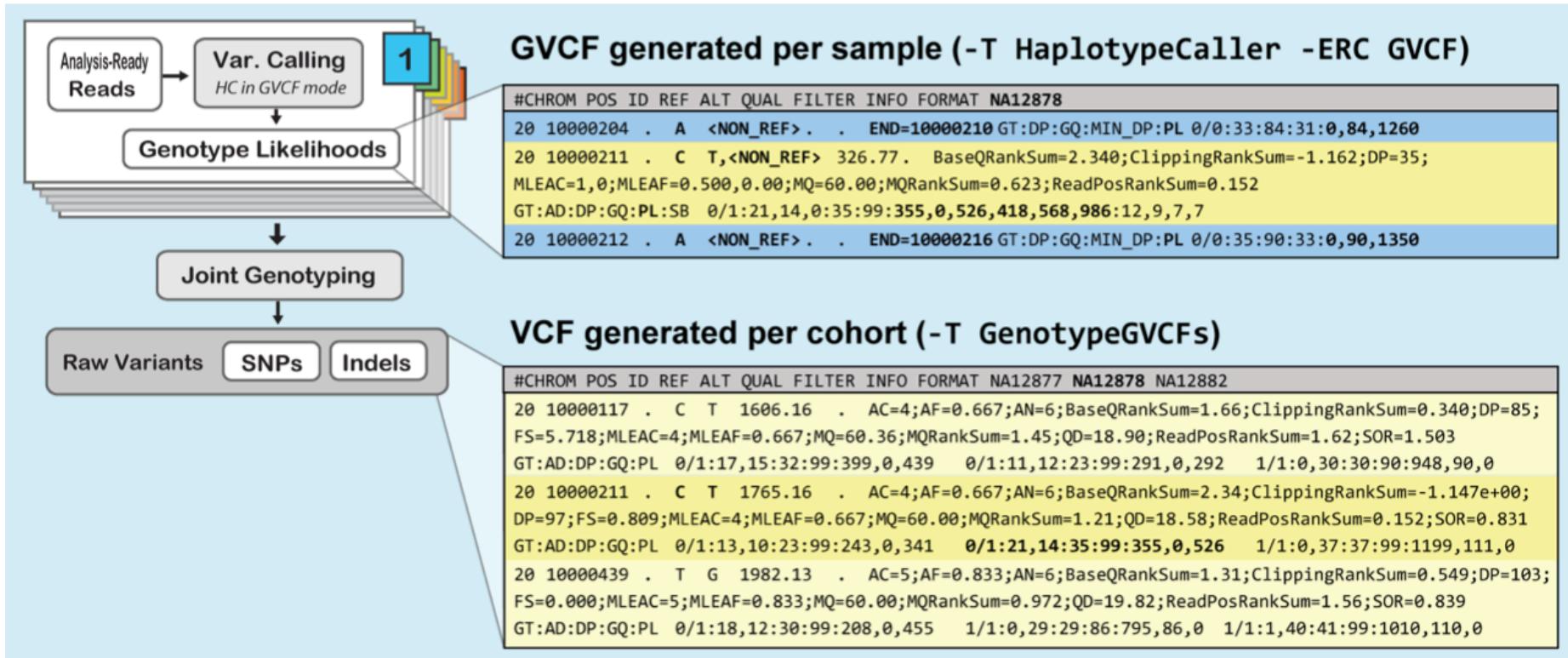
- When only 1 sample => "G" may be considered errors.
- When all samples => More confidence in the calling of A/G variant.

Joint callset → empowered analysis

GATK Best Practices for Variant Discovery
(<https://software.broadinstitute.org/gatk/download/workshops>)



gVCF VS regular VCF





Conclusion on the GATK workflow

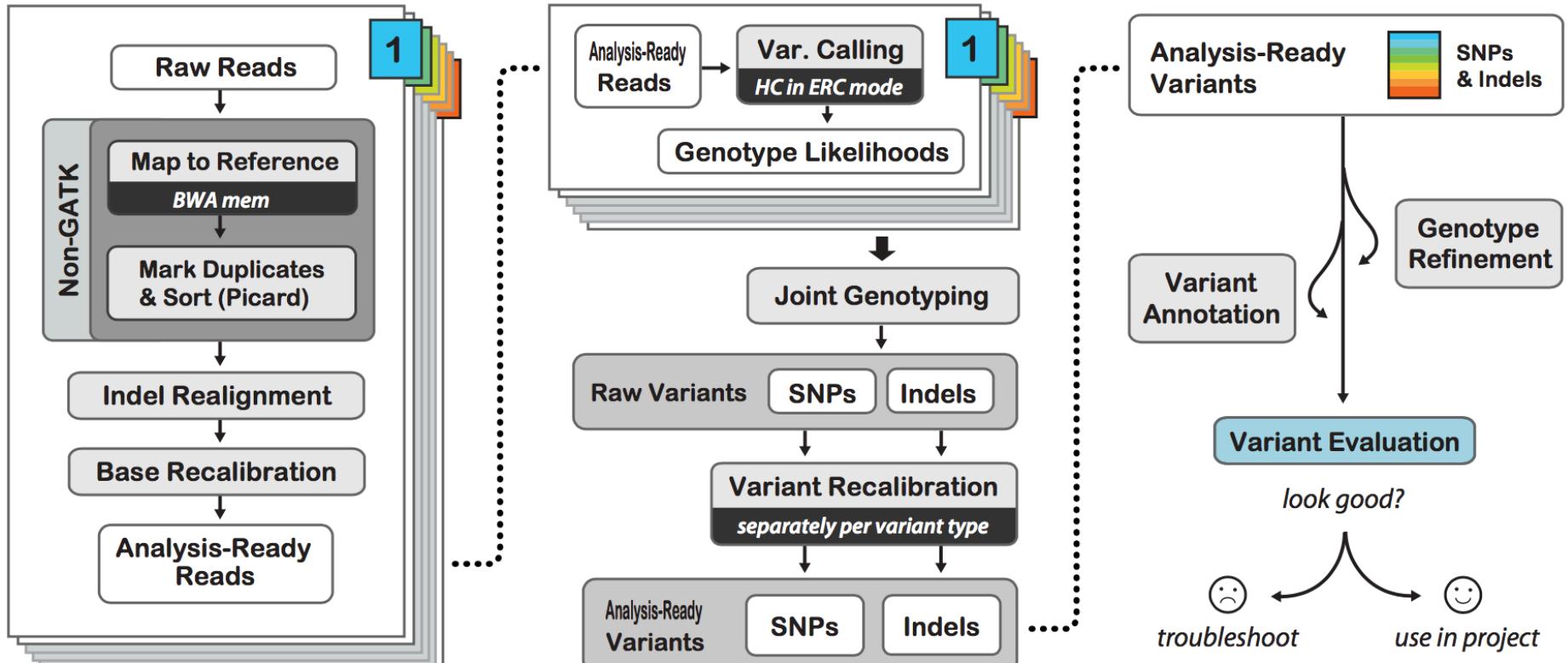
Data Pre-processing

>>

Variant Discovery

>>

Callset Refinement



Thanks for your attention