

py-gbcms

High-Performance Variant Counting for Liquid Biopsy

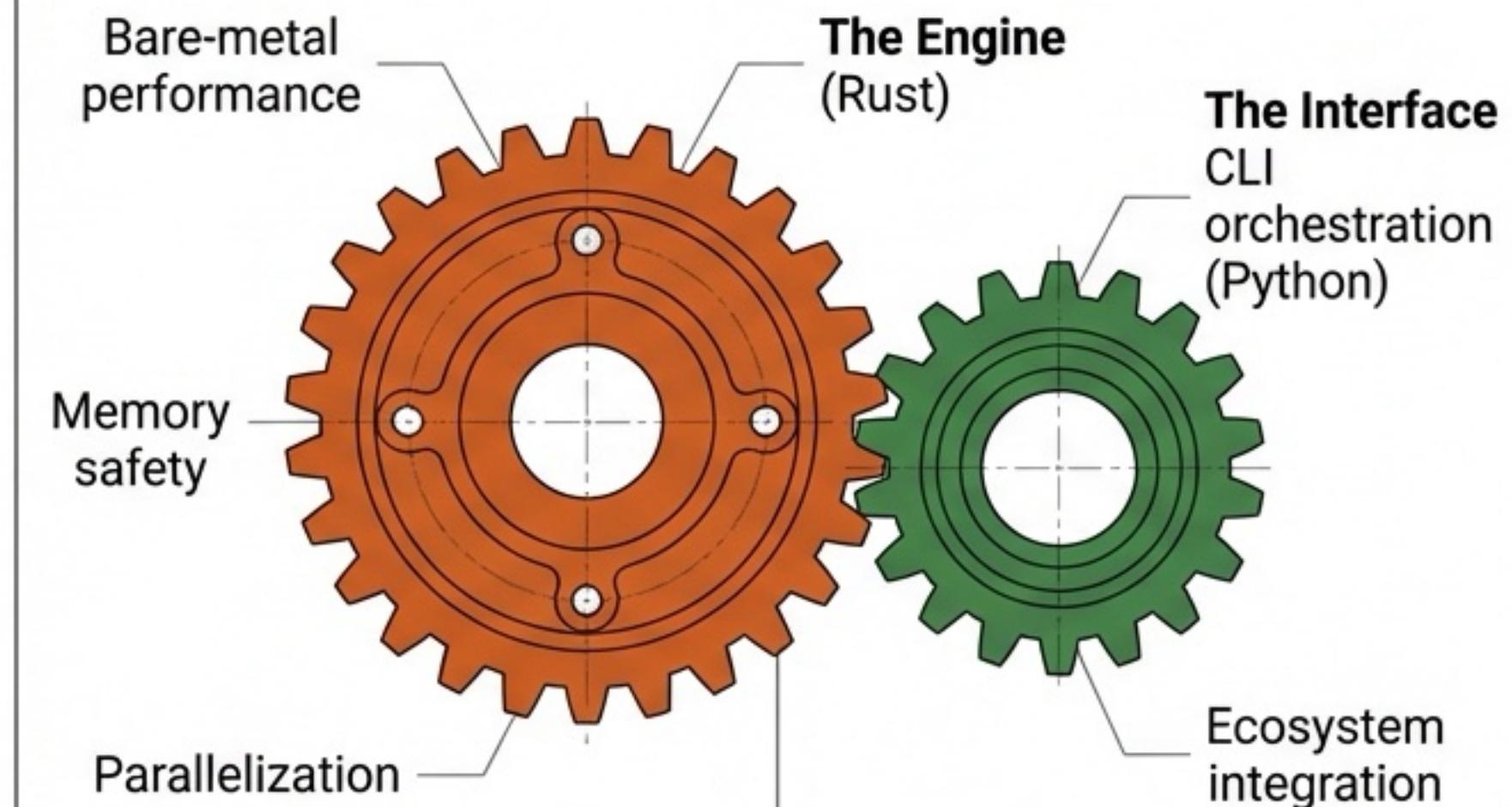
A Rust-Accelerated Architecture for Accurate cfDNA Analysis

Author: MSK-ACCESS | License: AGPL-3.0

Precision Counting at Production Scale

- **Core Function:** Extracts allele counts, VAF, and fragment metrics from BAM files.
- **The Focus:** Designed for Liquid Biopsy (cfDNA).
- **The Output:** Standardized VCF and extended MAF formats.

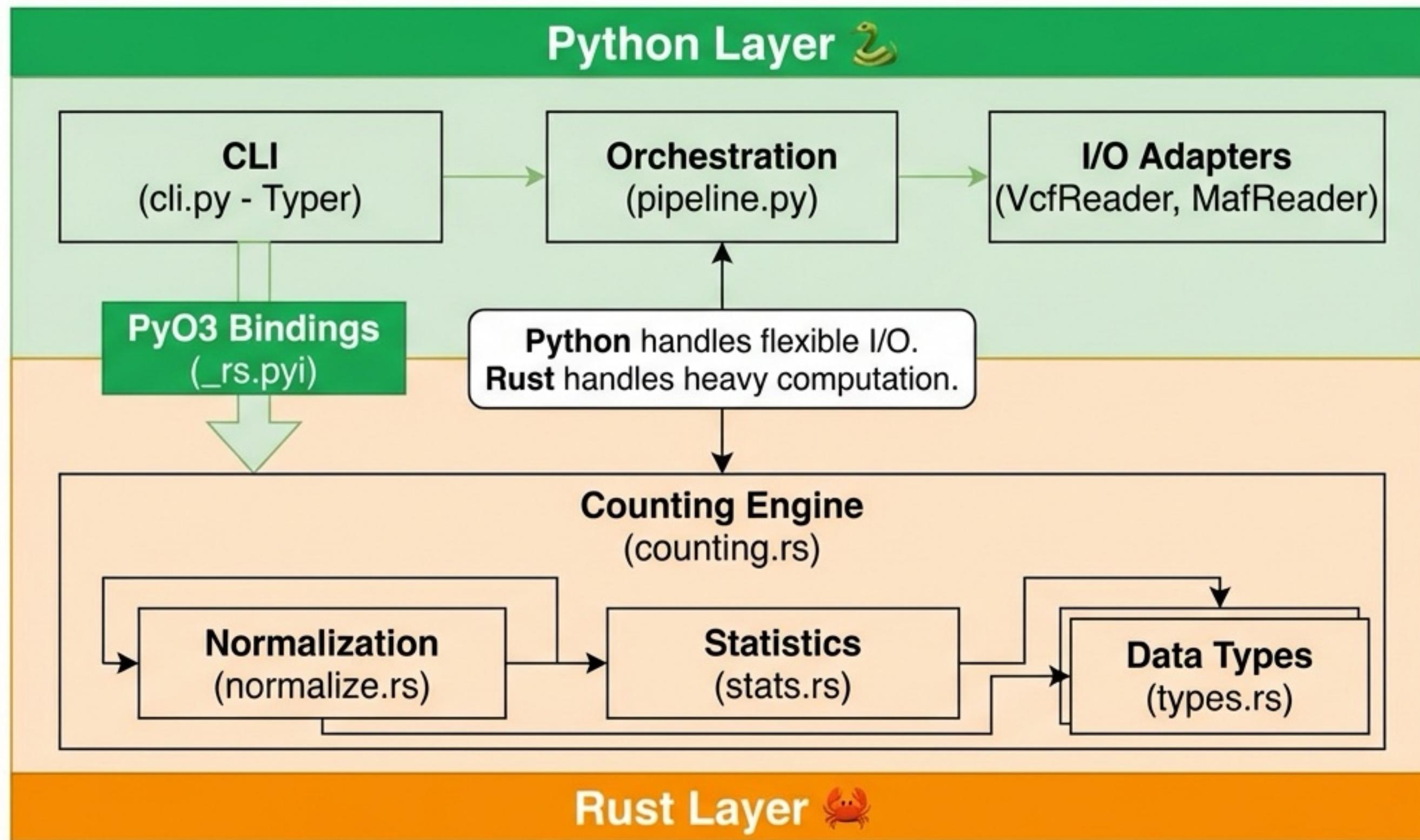
The Hybrid Advantage



Key Features

- PCR-Aware Fragment Counting
- Fisher's Exact Test for Strand Bias
- Smith-Waterman Complex Resolution

The Hybrid Architecture



Deployment Strategies

Standard (PyPI)



`pip install py-gbcms`

- Requires: Python 3.10+, glibc 2.34+
- OS: Ubuntu 22.04+, RHEL 9+

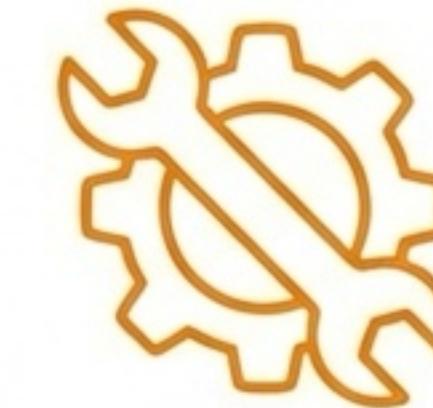
Containerized (Production)



Docker: `ghcr.io/msk-access/py-gbcms`

- Singularity / Apptainer
- Essential for clinical reproducibility

Legacy / HPC Source



Build from source via maturin

- For RHEL 8 / CentOS 8
- Requires: clangdev + rust toolchain

Pre-built wheels available for Linux (x86_64, aarch64), macOS (Intel/Silicon), and Windows.

The CLI Workflow

```
user@host:~$ gbcms run \  
  --variants mutations.vcf \  
  --bam sample:sample.bam \  
  --fasta reference.fa \  
  --output-dir results/ \  
  --min-mapq 30 --filter-duplicates
```

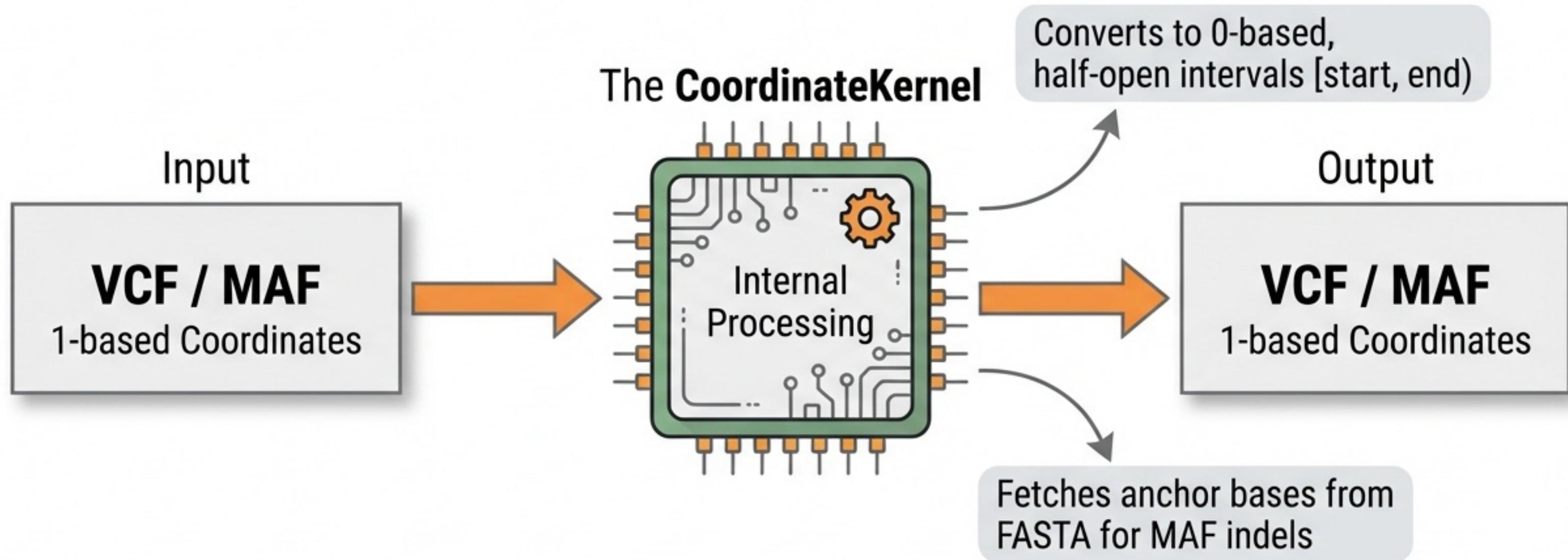
Targets
(VCF or MAF)

Input Alignment
(BAM)

Required for
Normalization

Generates VCF/MAF
with counts

Coordinate Normalization & Integrity



Solves the **off-by-one errors** common in bioinformatics by centralizing logic.

The Metrics That Matter

Primary Metrics

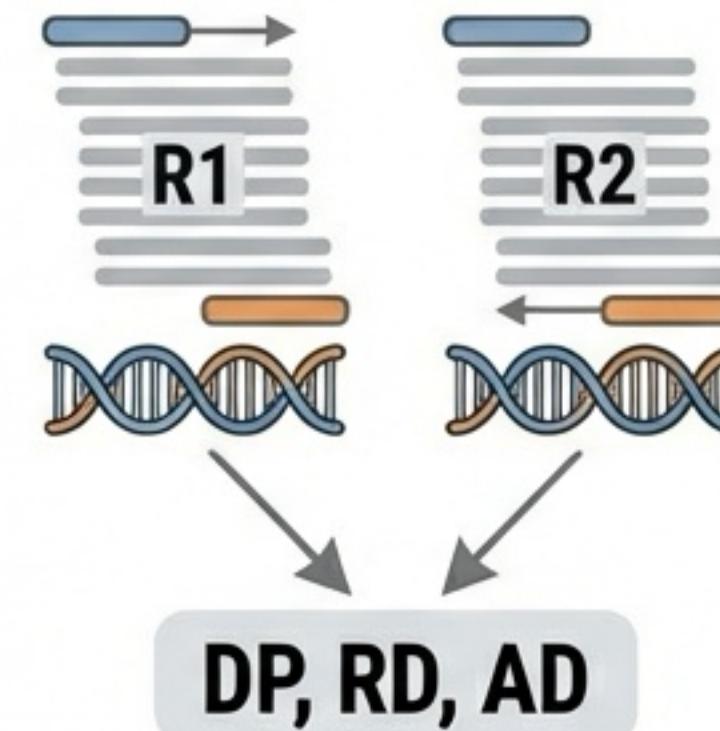
$$\text{VAF} = \frac{AD}{RD + AD}$$

Variant Allele Frequency. The core measure of mutation abundance.

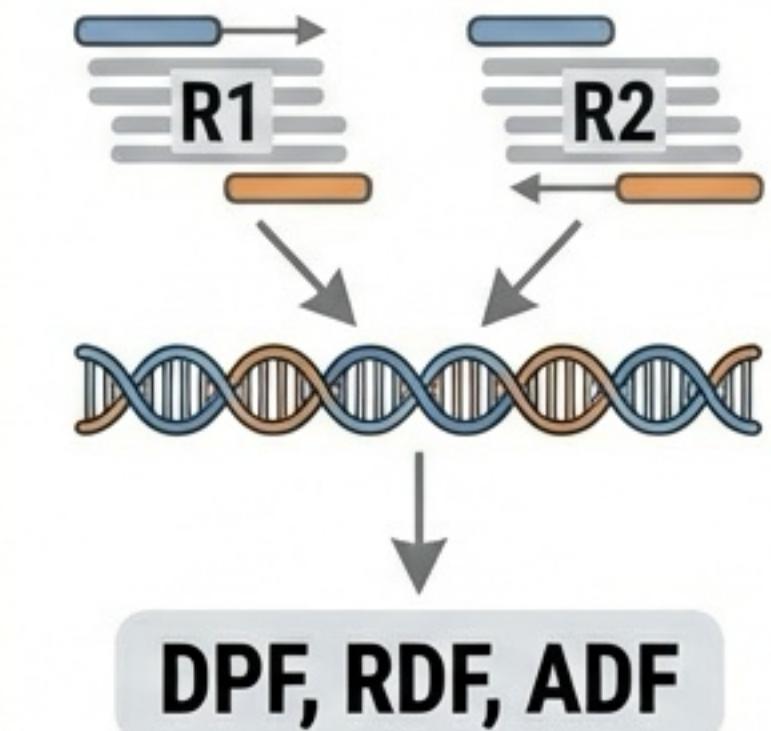
Strand Bias: Fisher's Exact Test (p-value & Odds Ratio). Detects sequencing artifacts.

The cfDNA Distinction

Read Counts

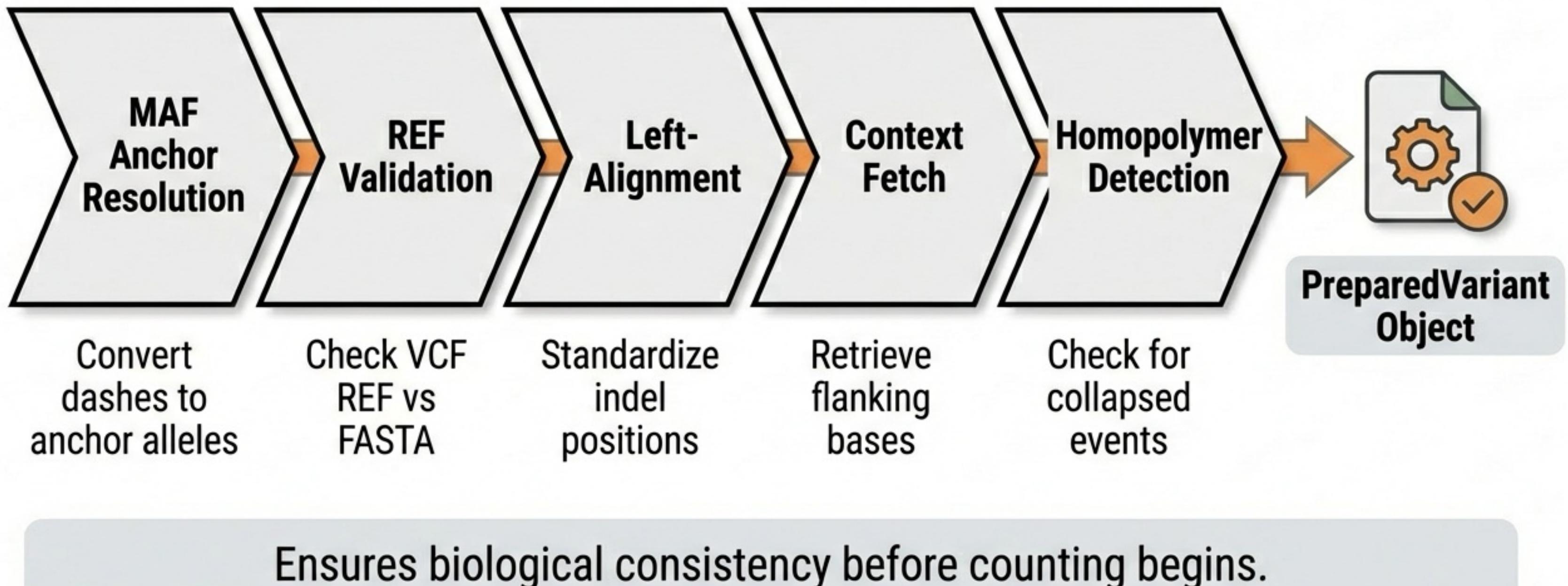


Fragment Counts

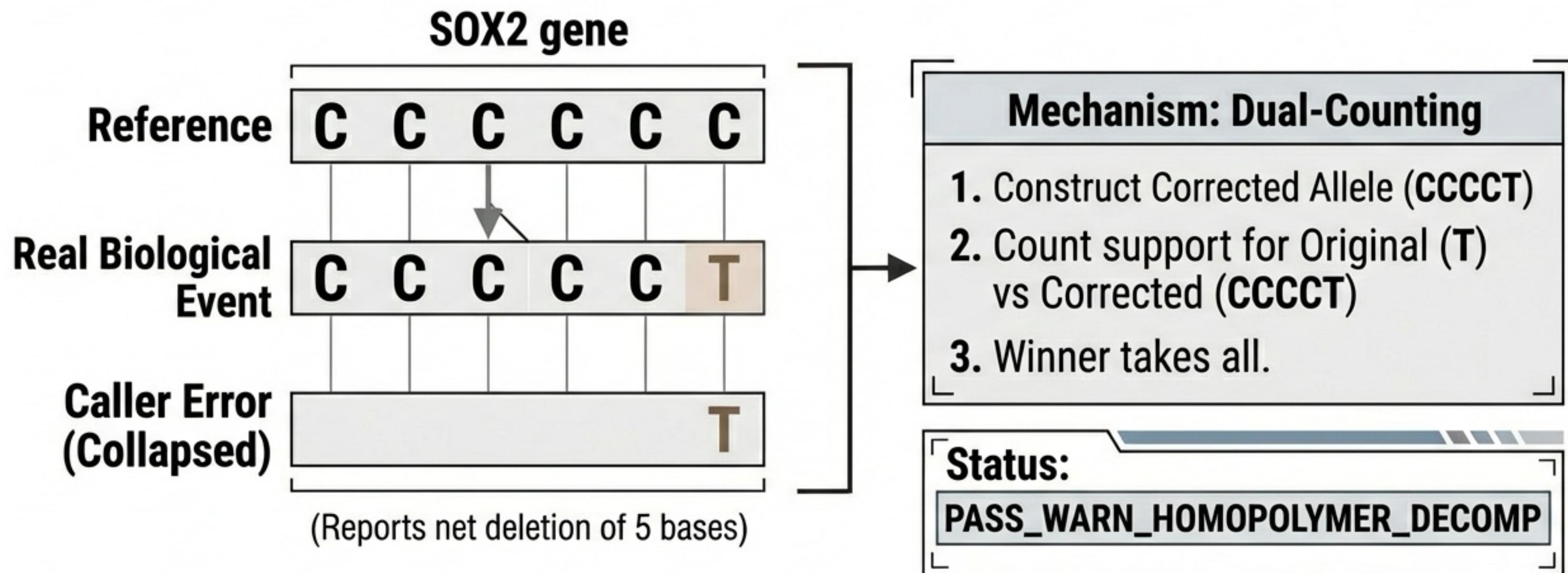


Fragment counting deduplicates read pairs to prevent artificial inflation of evidence.

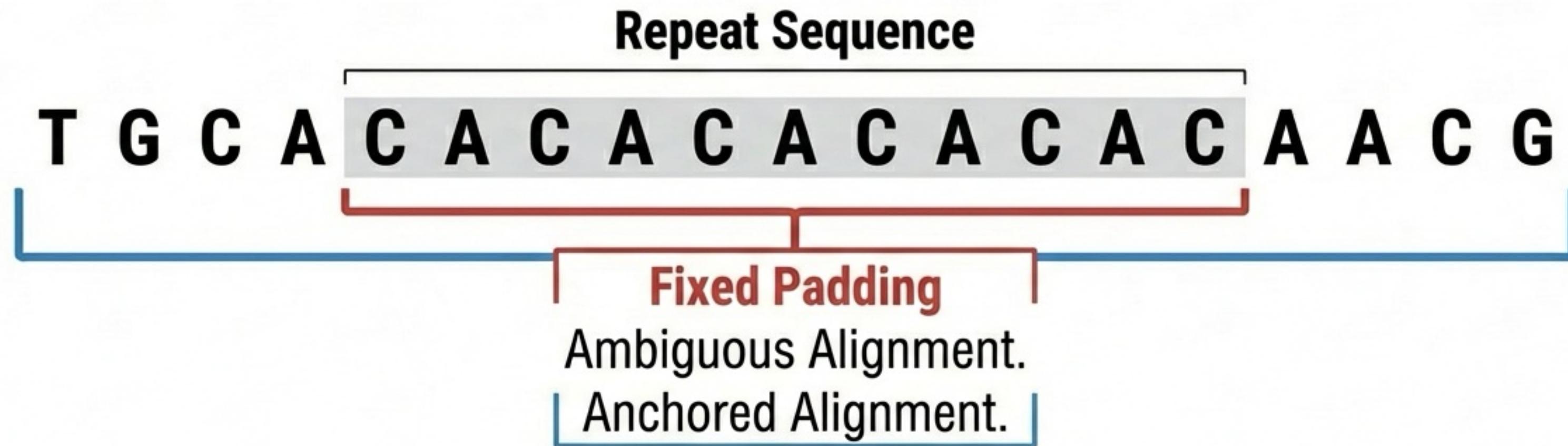
The Variant Preparation Pipeline



Edge Case: Homopolymer Decomposition

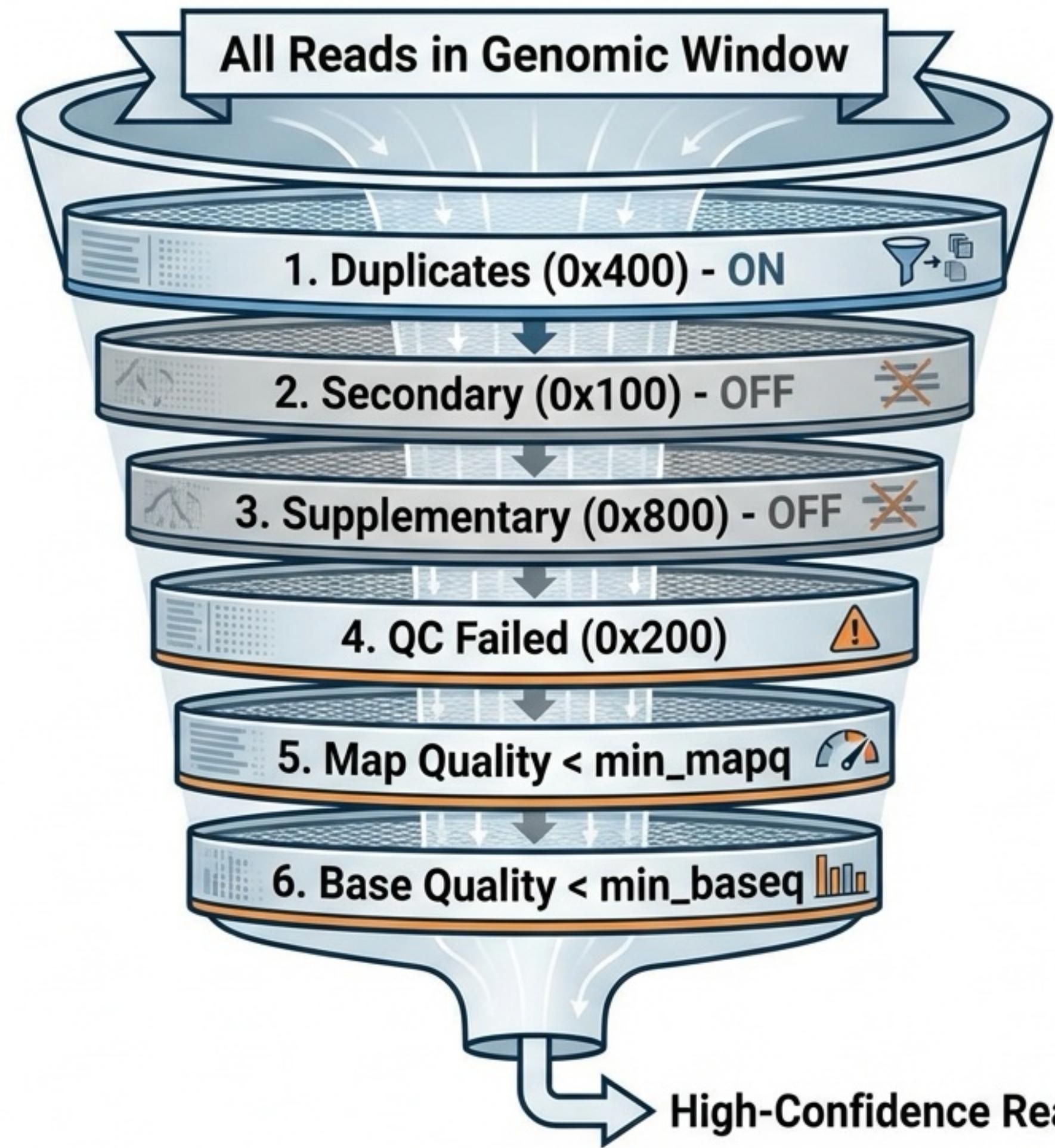


Adaptive Context & Windowed Scanning

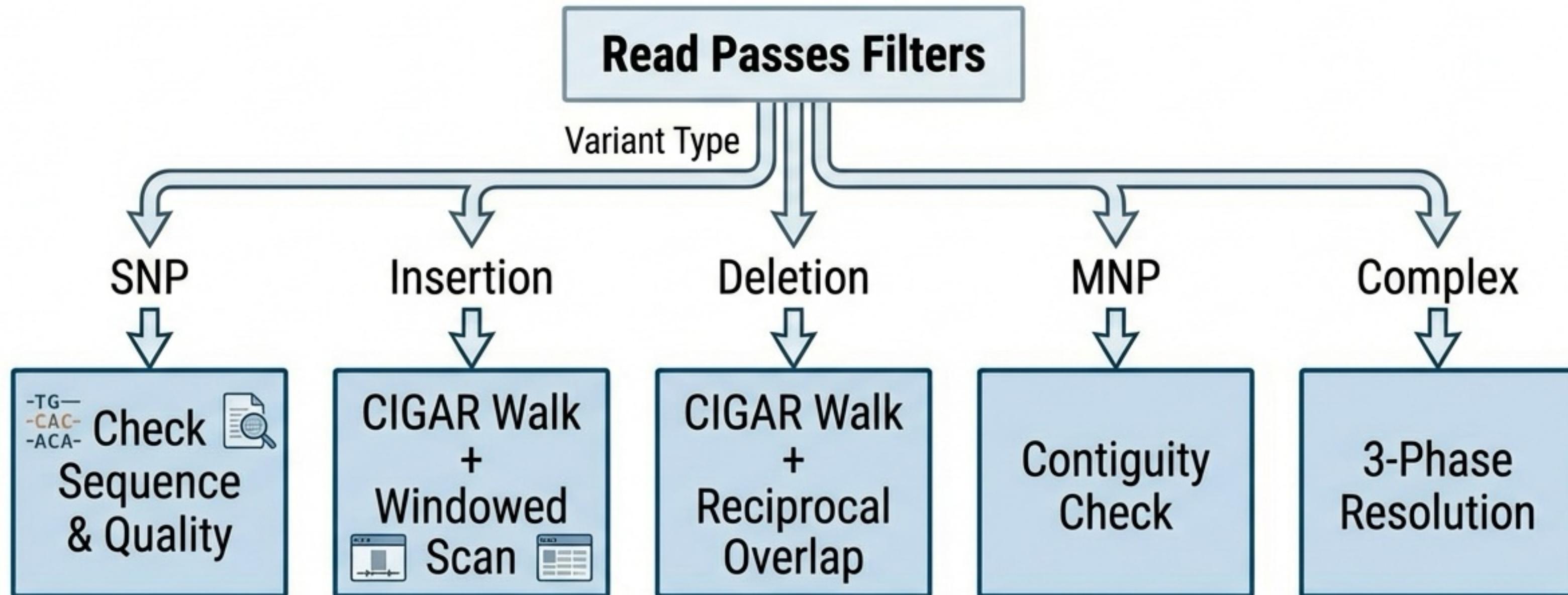


```
padding = max(default, repeat_span / 2 + 3)
```

Result: Expands field of view in repetitive regions.



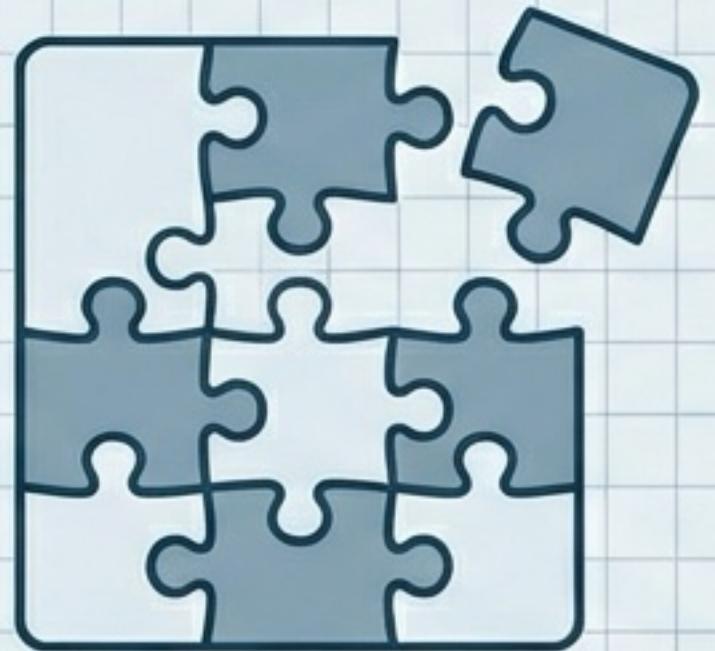
Allele Classification Logic



Auto-detection falls back to Complex check for unknown types.

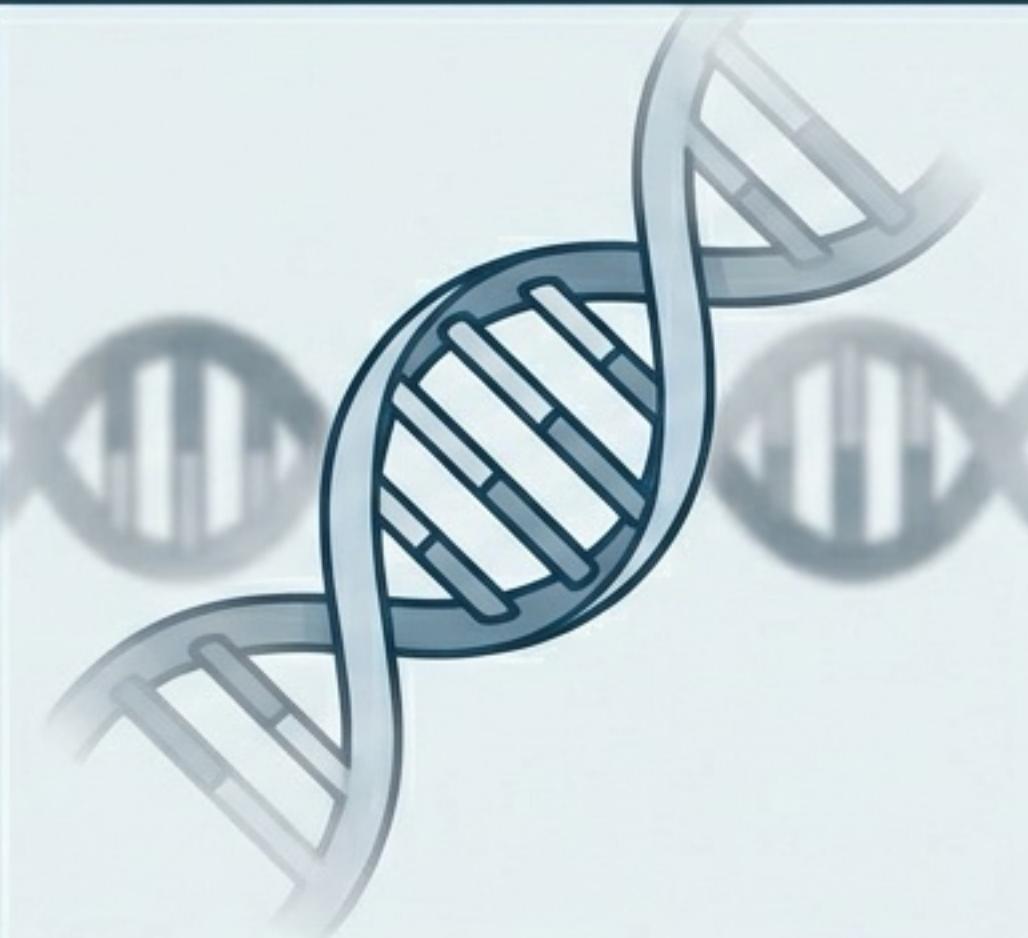
Deep Dive: Complex Variant Resolution

Phase 1: Reconstruction



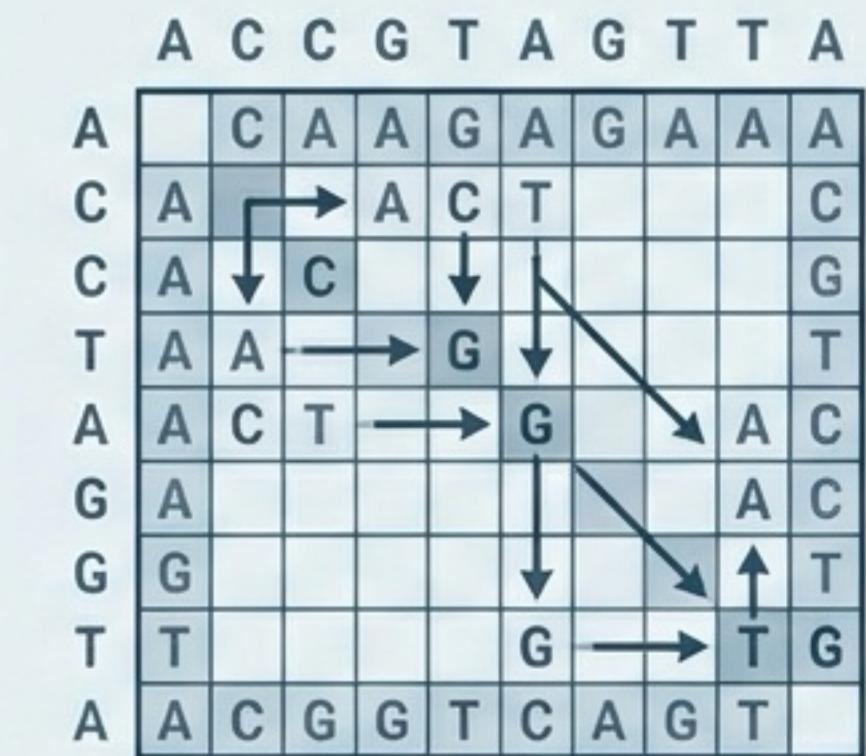
Haplotype Reconstruction
via CIGAR parsing.

Phase 2: Masked Comparison



Quality-aware matching.
Low-quality bases are masked.

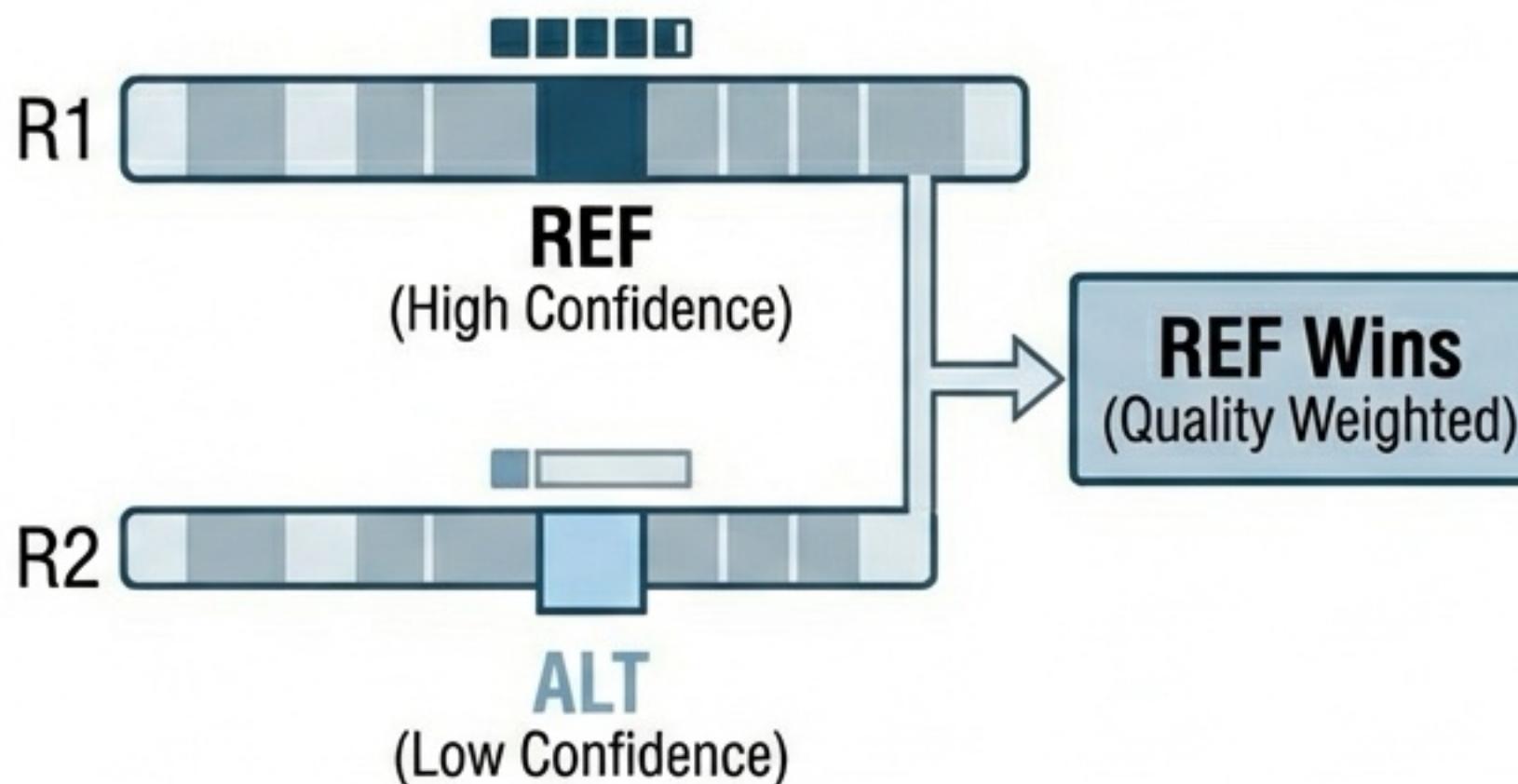
Phase 3: Smith-Waterman



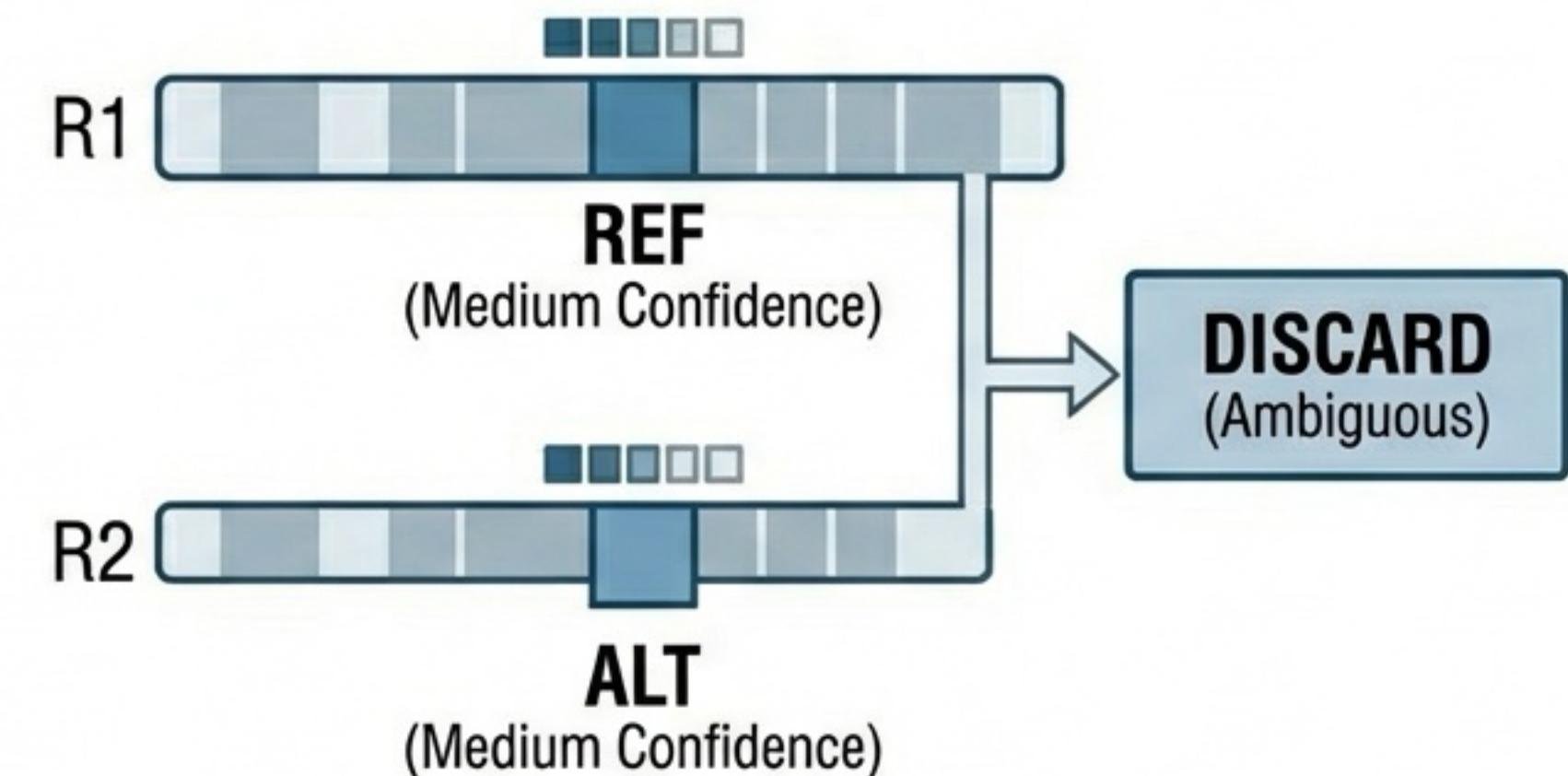
Fallback for difficult alignments.
Semiglobal alignment.

Fragment Counting & Consensus

Case 1: Conflict Resolution



Case 2: Ambiguity



Does NOT default to REF. Prevents VAF deflation.

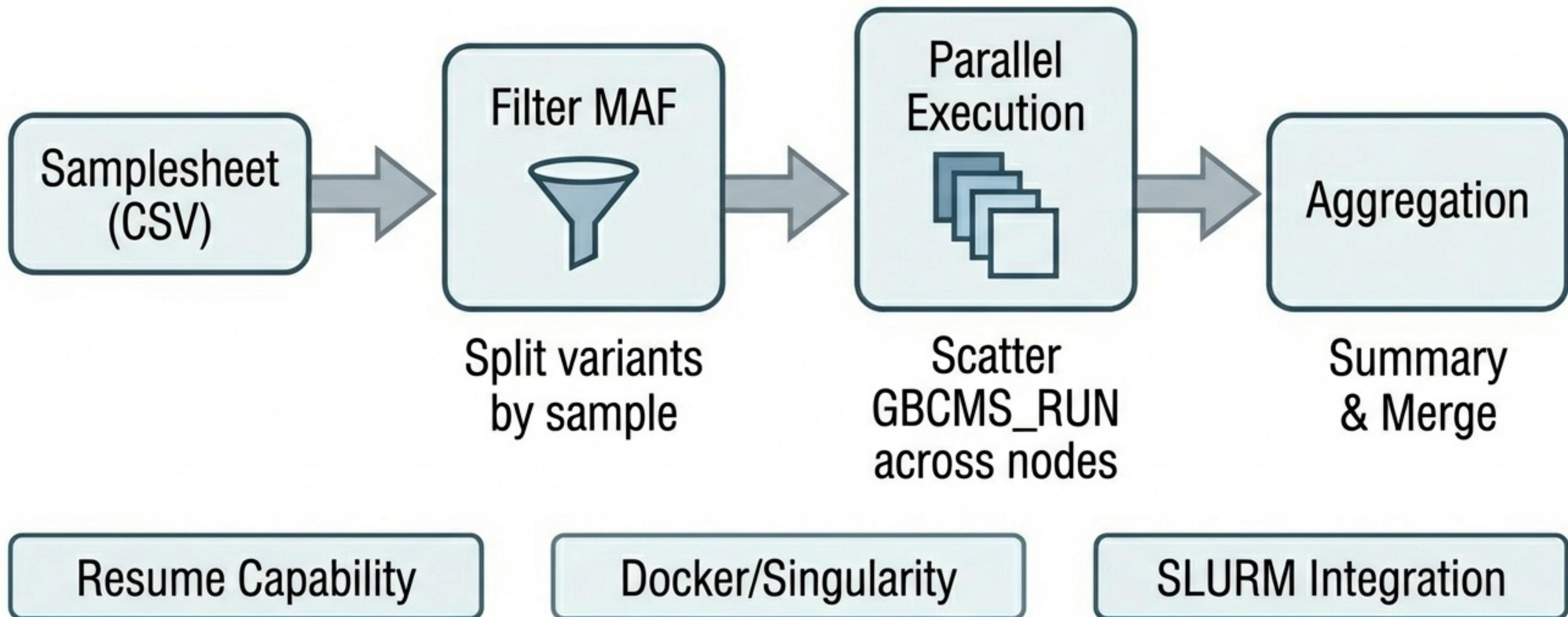
Statistical Rigor: Strand Bias

	Forward Strand	Reverse Strand
Reference Support	a	b
Alternate Support	c	d

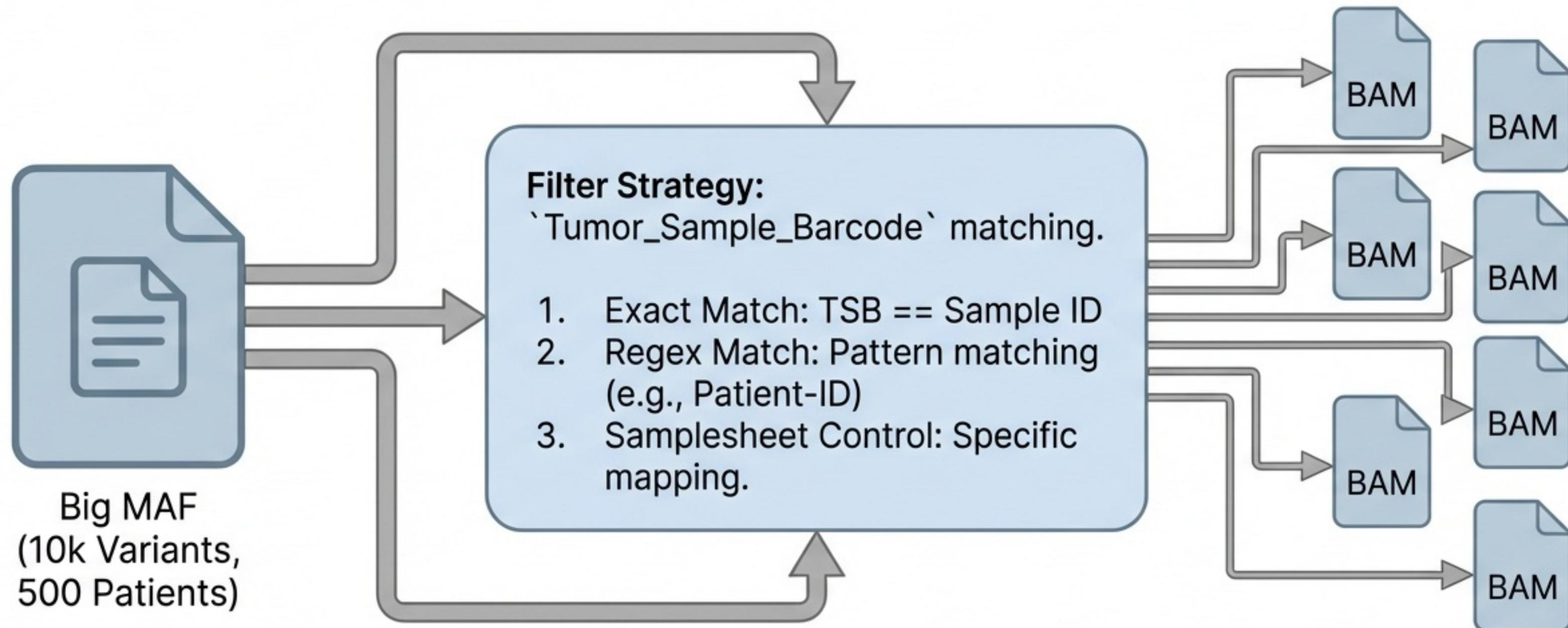
Fisher's Exact Test

- Goal: Detect artifacts where variants appear on only one strand.
- Calculated for: Reads (SB_PVAL) and Fragments (FSB_PVAL).

Scaling Up: The Nextflow Pipeline



Multi-Sample Filtering



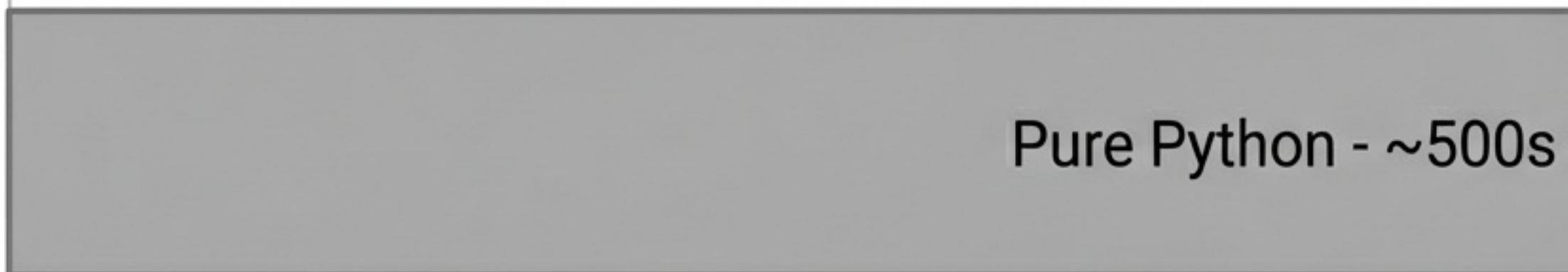
Prevents processing irrelevant variants, saving massive compute time.

Performance & Benchmarks



py-gbcms (Rust) - ~25s

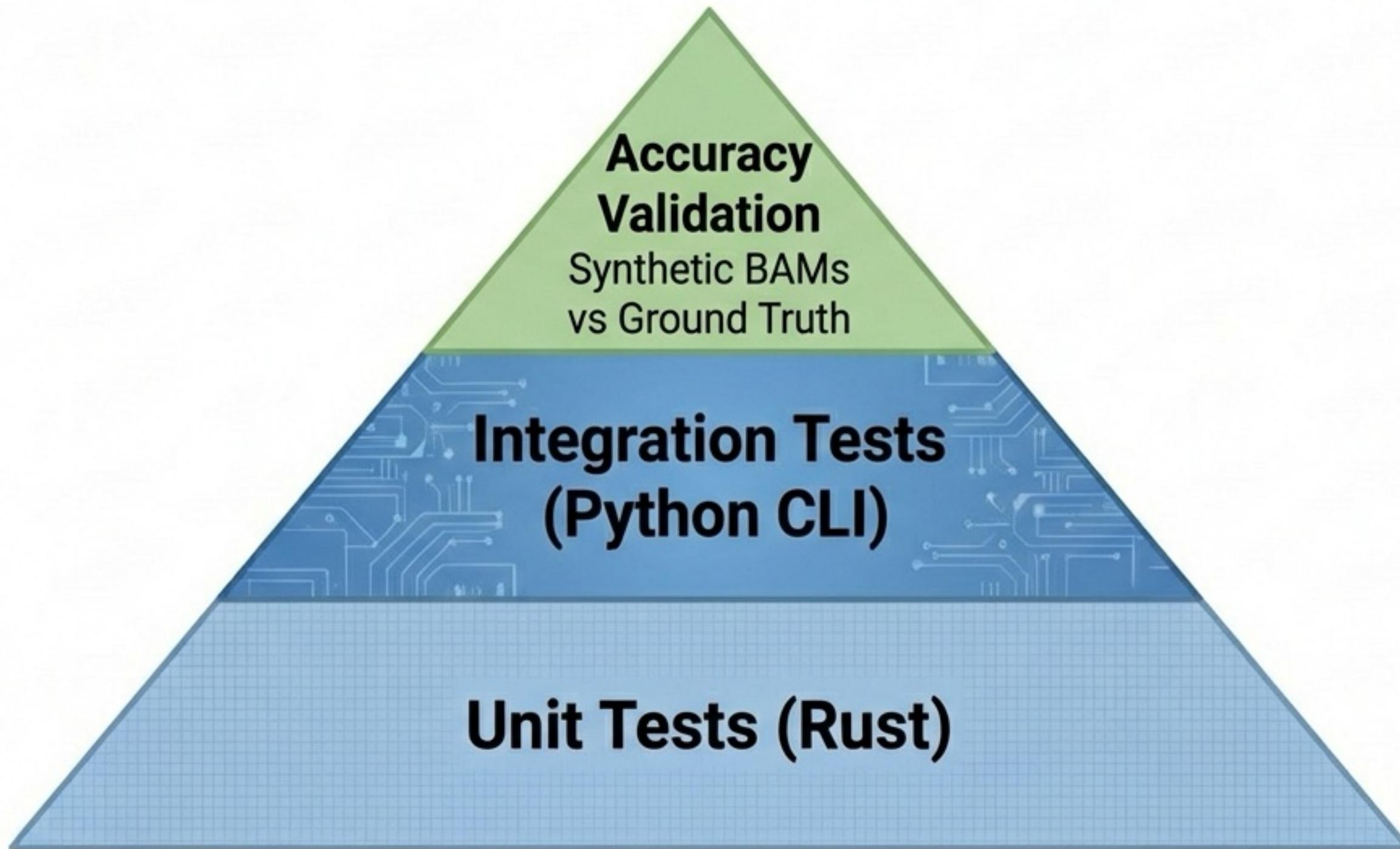
20x Faster Execution



Resource Details Box

- **Data:** 1.3GB cfDNA BAM (608 variants)
- **Compute:** 4 CPUs
- **Efficiency:** Rayon parallelization maximizes CPU saturation.

Validation & Quality Assurance



Specific Test Suites

- `test_shifted_indels`
(Windowed detection)
- `test_fuzzy_complex`
(Masked matching)
- `test_fragment_consensus`
(R1/R2 resolution)

Get Started

-  github.com/msk-access/py-gbcms
-  msk-access.github.io/py-gbcms
-  ghcr.io/msk-access/py-gbcms

```
pip install py-gbcms
```

Shah, Ronak H., et al. “py-gbcms: High-performance variant counting from BAM files.” Memorial Sloan Kettering Cancer Center, 2024.