# Strand III

## **An Automated Variant Verification** System To Improve Reporting Efficiency

S.Katragadda, S. Ghosh, S. Nayanala, A. Das Mahapatra, S. Aliya Afreen, A. Singh, A. Janakiraman, V. Veeramachaneni; Strand Life Sciences, Bangalore, India

#### Introduction

As NGS laboratories move to whole exome sequencing (WES), a larger number of SNVs/Indels are shortlisted in each case.

Before these variants can be included in a report, a final check on variants with borderline quality indicators is carried out by experienced bioinformaticians to ensure the variants are not artifacts arising from genome assembly differences, homology, sequencing errors, or pipeline parameters.

This variant verification (VV) process can be time consuming since it often involves the use of external tools and databases. It can also be subjective since it depends on experts reviewing the reads in a genome browser.

Our modular and pluggable system aims to significantly reduce VV time, by computationally determining if the variant in question is an artifact or not.

## Achievements

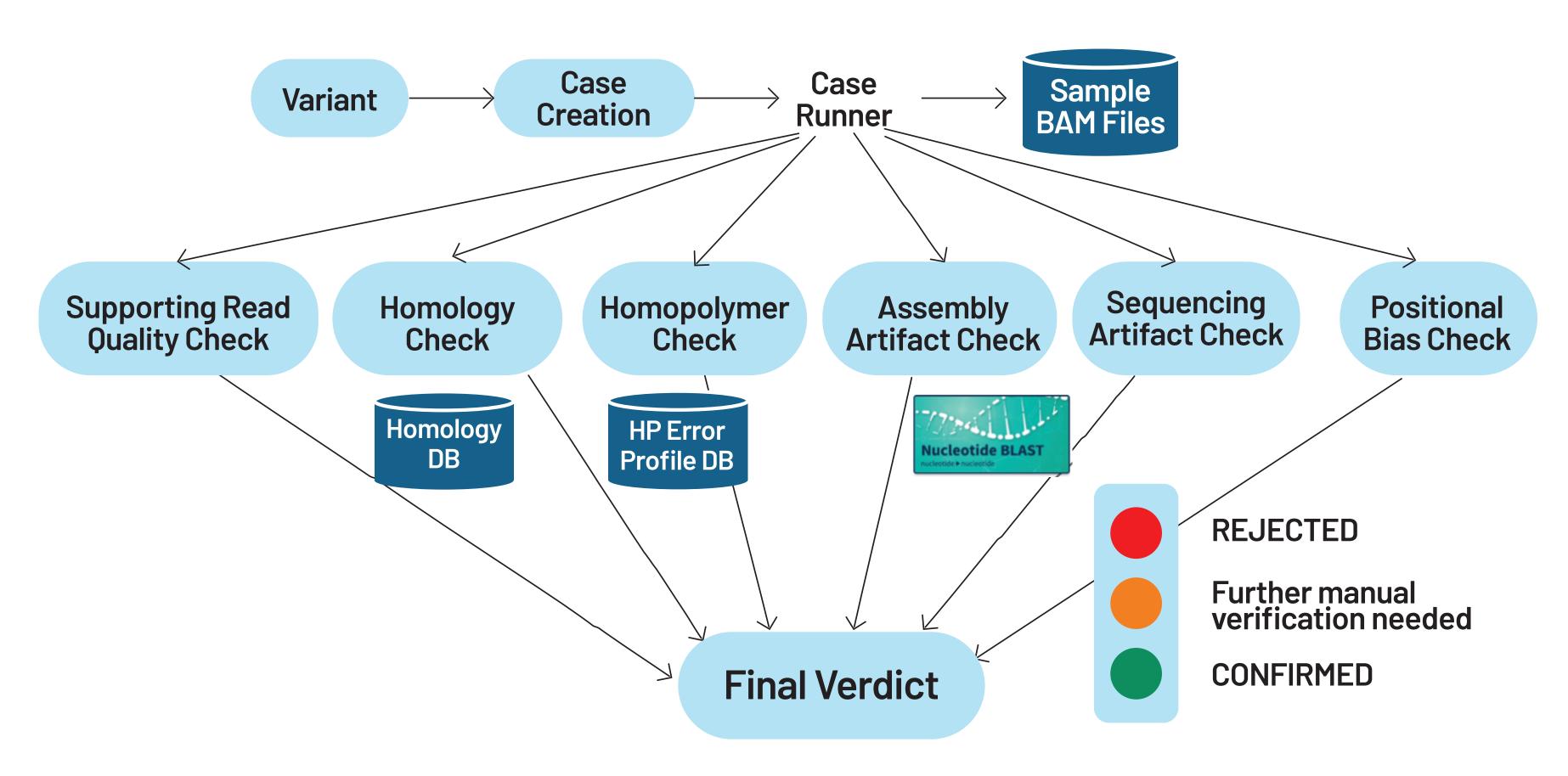
- SNVs/Indels for clinical reporting
- variants are automatically processed.
- Overall reduction of VV effort by ~80%.

## Approach

The system consists of multiple tools which perform different types of quality checks.

When a VV request is made for a specific variant in a sample, each tool assesses if the variant could be a specific type of artifact and returns a verdict.

delivers a final verdict.



Developed a modular and extensible system for verifying

Integrated the system into WES workflow so that all shortlisted

The system then aggregates the results from all the individual tools and

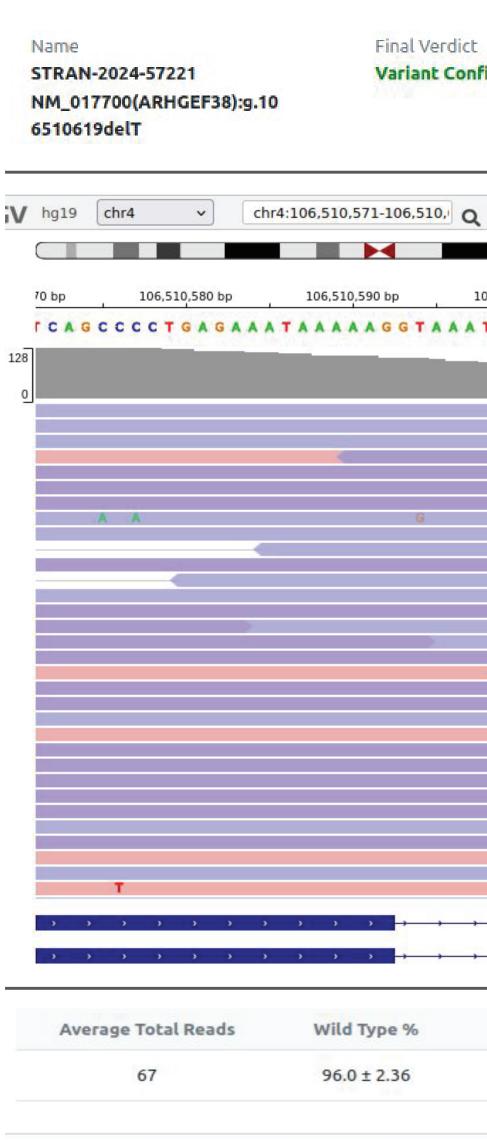
#### Supporting Read Quality Check

If there are at least 30 reads with length > 120 bp and mapping quality = 60 at the variant locus, and the variant is supported by at least 30% of such reads, the check is passed.

#### Homopolymeric Artifact Check

Homopolymer checker assesses if an indel in a homopolymer stretch is likely to be genuine. It uses a pre-computed error profile containing the mean  $(\mu)$  and standard deviation ( $\sigma$ ) of the supporting reads % for indels of different lengths.

If the supporting reads % of the variant is less than  $\mu$  +  $2 * \sigma$ , the variant is labelled a homopolymeric artifact.

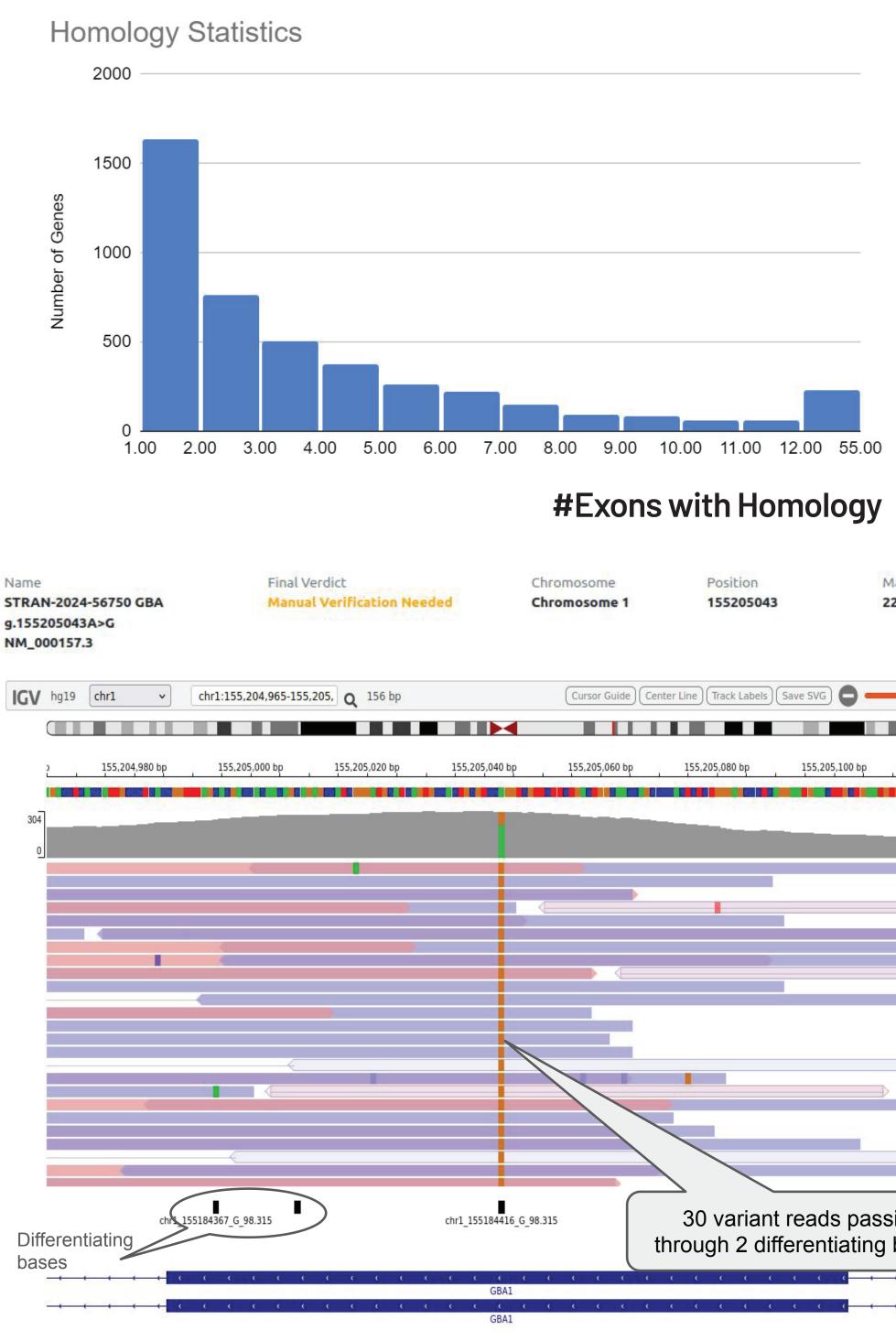


firmed		nosome <b>nosome 4</b>	Position 106510610	Mapping Quality 55.0
78 bp		Cursor Guide Ce	enter Line) (Track Labels) (Sa	ave SVG)
	106,510,610 bp	106,510,620 bp	106,510,630 bp	106,510,640 bp 106,5
	<b>=</b>			
		G .G		G
	===			
	=			
<b>.</b>	E			
				•
	<b>_</b>		C	
			-	
				-
· · · · · ·	ARHGEF38			· · · · · · · · · · · · · · · · · · ·
$\rightarrow \rightarrow \rightarrow \rightarrow$	ARHGEF38	$\rightarrow \rightarrow \rightarrow \rightarrow$	· · · · · · · · ·	· · · · · · · · · · · · · · · · · · ·
One Insertion %	Two lo	sertions %	One Deletion %	Two Deletions %
0.00 ± 0.00		8 ± 1.95	0.657 ± 1.21	0.240 ± 0.678

#### Homology Artifact Check

Homology artifact checker assesses if the variant reads originate from a homologous region. It uses a pre-computed database that contains all differentiating bases between regions in the genome having > 90% similarity with each other. This database covers 16,591 regions from 4,442 genes.

If the variant has at least 5 supporting read pairs that pass through at least 5 differentiating base loci with REF bases at those loci, the check is considered to be passed.



The reads list and the differentiating bases for the reads are as follows

ne summary of the differentiating bases count and corresponding read pair counts are as fo

Differentiating Bases Count	Supporting
1	
2	

Position 155205043		Mapping Q 22.0	Quality
(Track Labels) (Save	e SVG) 🖨 🗕		
205,080 bp	155,205,100 bp	155,20	<u>s.:</u> •
variant re gh 2 differ			<b>*</b>
<ul> <li></li> <li><td>&lt; &lt;<mark></mark></td><td>· · · ·</td><td>*</td></li></ul>	< < <mark></mark>	· · · ·	*

ows 👻	
Reads Pair Count	
5	
30	

#### Assembly Artifact Check

Assembly artifact checker extracts a representative read with the variant allele, aligns it against GRCh38 and T2T-CHM13, and analyzes the alignments to determine if the variant call is a result of build differences.

#### **Positional Bias Check**

Positional bias checker detects if the variant is a result of noisy alignment towards the end of the reads.

#### Sequencing Artifact Check

Sequencing artifact checker flags variants which are supported primarily by mates having disagreeing bases at the variant location.

ame : STRAN-2024-	57783 NM_0010376	75(NBPF9):g.14461523	7T>G Export C	ase Details Export Reads
/iew Case Details				
Name STRAN-2024-57783:chr1	Final Verdict Variant Rejected	Chromosome Chromosome 1	Position 144615237	Mapping Quality <b>7.5</b>
Reference <b>T</b>	Alternate <b>G</b>	Sample STRAN-2024-57783_NXP-LR0142.b	am	
View StrandOmics Metadata				
esults	porting Read Quality? MAYBE (1	This variant has poor Supporting Read Quali	ity.)	~
<b>Results</b> s this variant likely due to poor Supp		This variant has poor Supporting Read Quali ably artifact observed in HG19 build.)	ity.)	~
esults s this variant likely due to poor Supp s this variant due to assembly artifa	act? <b>YES</b> (This variant is an assem			
esults s this variant likely due to poor Supp s this variant due to assembly artifa s this variant due to homology?	act? <b>YES</b> (This variant is an assem MAYBE (This variant may be due to a	nbly artifact observed in HG19 build.)	s in this region.)	~ ~
esults s this variant likely due to poor Supp s this variant due to assembly artifa s this variant due to homology?	act? <b>YES</b> (This variant is an assem MAYBE (This variant may be due to region? <b>NO</b> (This variant is not el	nbly artifact observed in HG19 build.) homology. There are 4 differentiating bases ligible for noise model data only one or two	s in this region.)	~ ~



#### Contact

Shanmukh Katragadda Chief Technology Officer -Bioinformatics +91-9880598107

shanmukh@strandls.com