



CIDC 2.0 BIOINFORMATIC PIPELINES WES VALIDATION

ESSEX MANAGEMENT

NATIONAL CANCER INSTITUTE (NCI)

06/14/2024

VERSION 1.0

DRAFT

SUBMITTED TO:

DAOUD MEERZAMAN
*COMPUTATIONAL GENOMICS AND
BIOINFORMATICS BRANCH (CGBB)*
NATIONAL CANCER INSTITUTE
CENTER FOR BIOMEDICAL INFORMATICS
& INFORMATION TECHNOLOGY
ROCKVILLE, MD 20850

SUBMITTED BY:

NICK RENZETTE, JENNIFER HARVEY
ESSEX MANAGEMENT, LLC
11140 ROCKVILLE PIKE | SUITE 332
ROCKVILLE, MD 20852-3149
DUNS: 829872345
CAGE CODE: 5CYC9

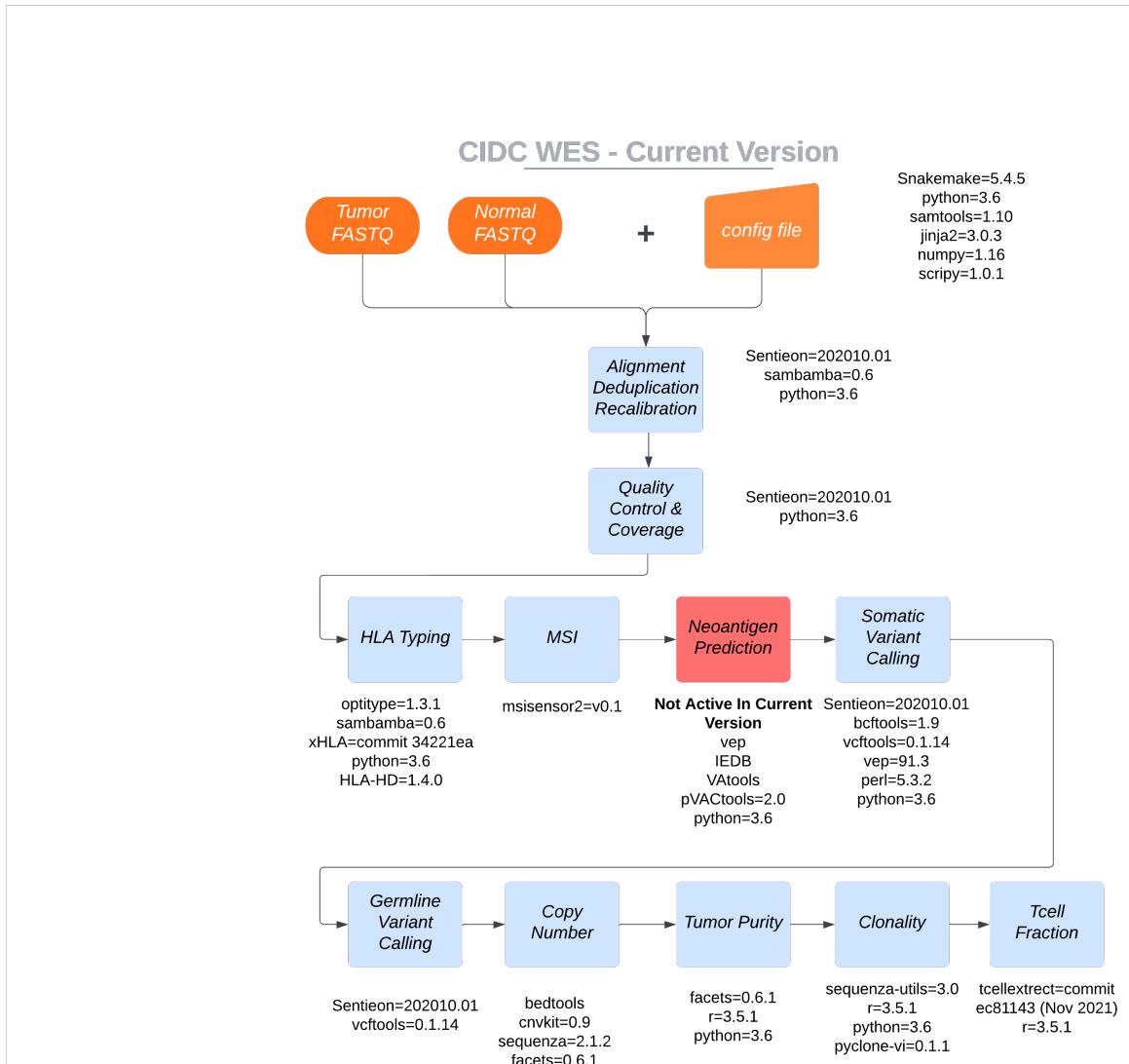
1.	<i>Introduction</i>	3
2.	<i>WES Pipeline – Validation Dataset</i>	6
3.	<i>WES Pipeline – Validation Method</i>	11
4.	<i>WES Pipeline – Validation Results</i>	12
5.	<i>WES Pipeline – Appendix I – Tabular Results</i>	17

1. INTRODUCTION

As part of the planned CIBC enhancements after the migration to the National Cancer Institute (NCI), the bioinformatic pipelines were reviewed by members the National Cancer Institute Computational Genomics and Bioinformatics Branch (NCI-CGBB) and Essex Management (EM). These reviews were carried out to determine which, if any, changes could be made to the pipeline to satisfy the following goals and objectives:

Goals	Objectives
Update and clean up code so that it is easy to read and understand for everyone working in the same code base, thus, making it easier to maintain, debug, and update.	Provide code review to determine areas for clean up and refactoring
Maintain industry standard software to optimize the best combination of biochemistry, mathematics, computer science, data science, and modern data analytics tools.	Review the performance of current packages to determine if current functionality is meeting customer needs (CIMAC Input)
Maintain current software versioning to optimize vendor support and application performance.	Provide gap analysis on current software versions to determine what tools to upgrade
Provide enhancements to pipeline to improve current functionality/performance and better support the stakeholder community the analysis of DNA, while maintaining backward compatibility with previous versions.	Add new functionality and features as desired by

The reviews created schematics of the extant pipelines stood up after the migration from Dana Farber Cancer Institute (DFCI) to NIH STRIDES. The **Whole Exome Sequencing (WES)** pipeline was divided into 11 distinct modules:

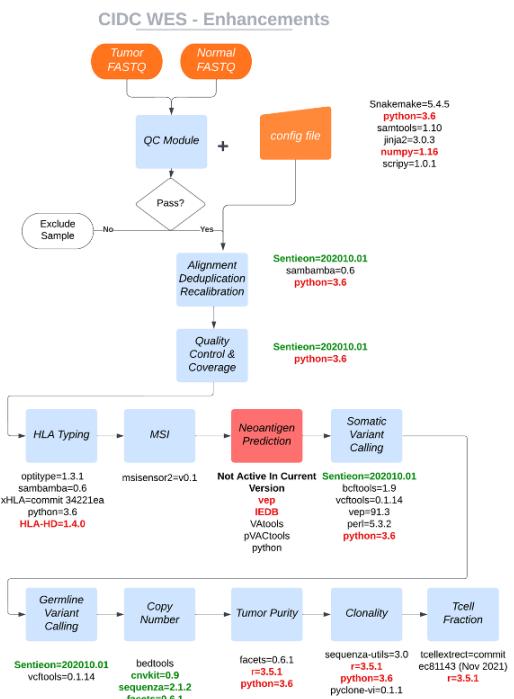


Note that in the original pipeline, the neoantigen prediction module was present but not run.

After thorough review, the following enhancements were planned for the WES Pipeline:

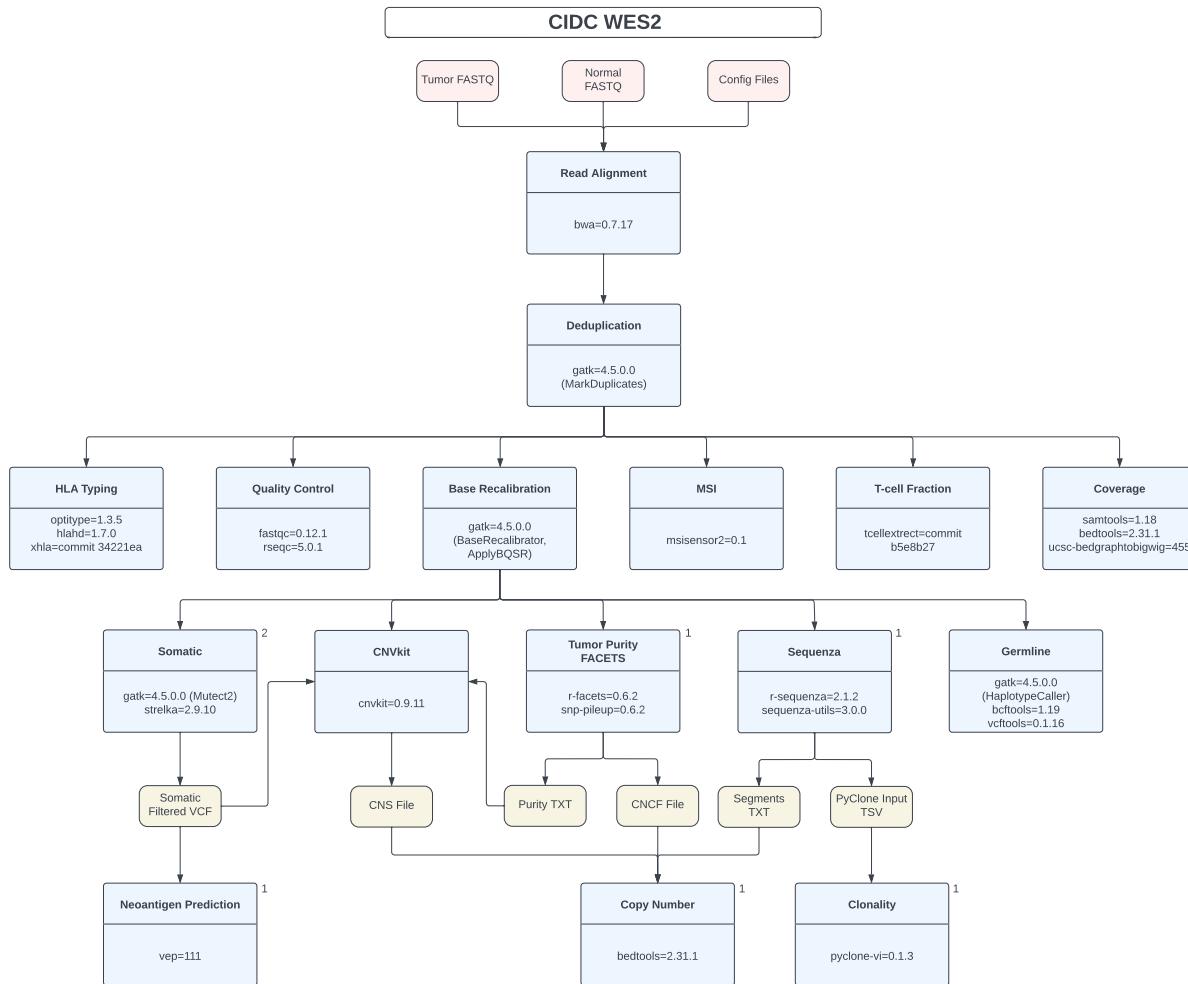
CIDC WES – Recommended Enhancements

WES Pipeline Recommendations	Reason
Add Neoantigen Functionality	Increased utility
Upgrade numpy beyond 1.19	CBIIT identified security vulnerabilities (High - null pointer critical - Numpy Deserialization of Untrusted Data)
Upgrade scipy to 1.10.0	CBIIT identified security vulnerabilities (Moderate - memory leak)
Upgrade Variant Effect Predictor (VEP) to 109.3 (May 2023)	Improved annotation based on the most recent Gencode gene models (Gencode v35).
Upgrade HLA-HD to 1.7.0 released Feb 2023	Reduced Memory Usage: Version 1.7.0 consumes approximately 75% less memory compared to version 1.6.1. We are on 1.4.0
Upgrade r	Current version unsupported; update past v3.6
Upgrade python	Current version unsupported; update past v3.7
Add file validation (QC Module)	To identify corrupt or incorrectly formatted files 8 releases have updated speed and accuracy.
Evaluate use of Senteion	CAVEAT: Difference in SNV/Indel Output
Evaluate CNV Callers	Known source of variability



The specific changes associated with the pipeline are enumerated in the figures above. For each module, code refactoring was also performed to increase readability and ease future maintenance and/or upgrade efforts.

Final Design



2. WES PIPELINE – VALIDATION DATASET

Validation Dataset:

For validation, it would be possible to compare the output of the enhanced pipeline to the original pipeline. While this method would be suitable to determine result continuity associated with the transition, it would not be suitable for determining result accuracy. The aligner and variant caller (Sentieon 202010.01) will be updated to open source components in the enhanced pipeline, which will likely lead to vastly differing results

(<https://support.sentieon.com/versions/202308.01/manual/appendix/releasenotes/> and <https://github.com/lh3/bwa/releases>). Thus, to verify the accuracy of the enhanced pipeline, we will use truth sets obtained from external sources.

The source fastq files for small variant calling will be collected from NIST's [Genome in a Bottle \(GIAB\) project](#), a consortium dedicated to authoritative characterization of the human genome. As such, they distribute high-quality sequencing data and benchmarking datasets for 7 samples (1 HapMap genome [NA12878], and two son/father/mother trios of Ashkenazi Jewish and Han Chinese ancestry).

For CNV benchmarking, data from the triple negative breast cancer line HCC1395 that has been extensively characterized will be used. These data have previously been used in a comprehensive CNV caller benchmarking study¹. The total set is as follows:

Sample	Description	Result Evaluated	Source – Sequencing Reads	Source – Benchmarking Data
NA12878_HG001	HapMap sample	Alignment / Small Variants	https://ftp.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/NA12878/Nebraska_NA12878_HG001_Truseq_Exome/NIST-Hg001-7001-ready.bam	https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/release/NA12878_HG001/latest/GRC38/
HG002_NA24385_son	Ashkenazi Jewish trio – son	Alignment / Small Variants	https://ftp.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/HG002_NA24385_son/OsloUniversityHospital_Exome/151002_7001448_0359_AC7F6GANXX_Sample_HG002-EEogPU_v02-KIT-Av5_AGATGTAC_L008_posiSrt.markDup.bam	https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/release/AshkenazimTrio/HG002_NA24385_son/latest/
HG003_NA24149_father	Ashkenazi Jewish trio – father	Alignment / Small Variants	https://ftp.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/HG003_NA24149_father/OsloUniversityHospital_Exome/151002_7001448_0359_AC7F6GANXX_Sample_HG003-EEogPU_v02-KIT-Av5_TCTTCACA_L008_posiSrt.markDup.bam	https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/release/AshkenazimTrio/HG003_NA24149_father/latest/
HG004_NA24143_mother	Ashkenazi Jewish	Alignment / Small Variants	https://ftp.ncbi.nlm.nih.gov/ReferenceSamples/giab/data/AshkenazimTrio/	https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/

¹ Masood et al. (2024) Under Review

	trio – mother – replicate 1		HG004 NA24143 mother/OsloUniversityHospital_Exome/151002_7001448_0359_AC7F6GANX_X_Sample_HG004-EEogPU_v02-KIT-Av5_CCGAAGTA_L008.posiSrt.markDup.bam	s/giab/release/AshkenaziTrio/HG004_NA24143_mother/latent/
WES_EA_T1	HCC139 5 Replicate 1 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_EA_T_1.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_EA_N_1.bwa.dedup.bam	https://zenodo.org/record/10081574
WES_FD_T1	HCC139 5 Replicate 2 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_FD_T_1.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_FD_N_1.bwa.dedup.bam	https://zenodo.org/record/10081574
WES_FD_T2	HCC139 5 Replicate 3 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_FD_T_2.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_FD_N_2.bwa.dedup.bam	https://zenodo.org/record/10081574

			https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_FD_N_2.bwa.dedup.bam	
WES_FD_T3	HCC139 5 Replicat e 4 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_FD_T_3.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_FD_N_3.bwa.dedup.bam	https://zenodo.org/reCORDS/10081574
WES_IL_T1	HCC139 5 Replicat e 5 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_IL_T_1.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_IL_N_1.bwa.dedup.bam	https://zenodo.org/reCORDS/10081574
WES_IL_T2	HCC139 5 Replicat e 6 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_IL_T_2.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_IL_N_2.bwa.dedup.bam	https://zenodo.org/reCORDS/10081574
WES_IL_T3	HCC139 5	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/dataset/WES/WES_IL_N_3.bwa.dedup.bam	https://zenodo.org/reCORDS/10081574

	Replicate 7 – WES dataset		RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_IL_T_3.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_IL_N_3.bwa.dedup.bam	
WES_LL_T1	HCC139 5 Replicate 8 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_LL_T_1.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_LL_N_1.bwa.dedup.bam	https://zenodo.org/records/10081574
WES_NC_T1	HCC139 5 Replicate 9 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_NC_T_1.bwa.dedup.bam https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_NC_N_1.bwa.dedup.bam	https://zenodo.org/records/10081574
WES_NV_T1	HCC139 5 Replicate 10 – WES dataset	CNV Calling	https://ftp-trace.ncbi.nlm.nih.gov/RerferenceSamples/seqc/Somatic_Mutation_WG/adata/WES/WES_NV_T_1.bwa.dedup.bam	https://zenodo.org/records/10081574

			https://ftp-trace.ncbi.nlm.nih.gov/R_eferenceSamples/seqc/Somatic_Mutation_WG/data/WES/WES_NV_N_1.bwa.dedup.bam	
--	--	--	---	--

Additional details about the GIAB data and analyses:

Source of hg38 reference file used by GIAB:

<https://ftp-trace.ncbi.nlm.nih.gov/ReferenceSamples/giab/release/references/>

GIAB recommended analysis tools for small variant benchmarking:

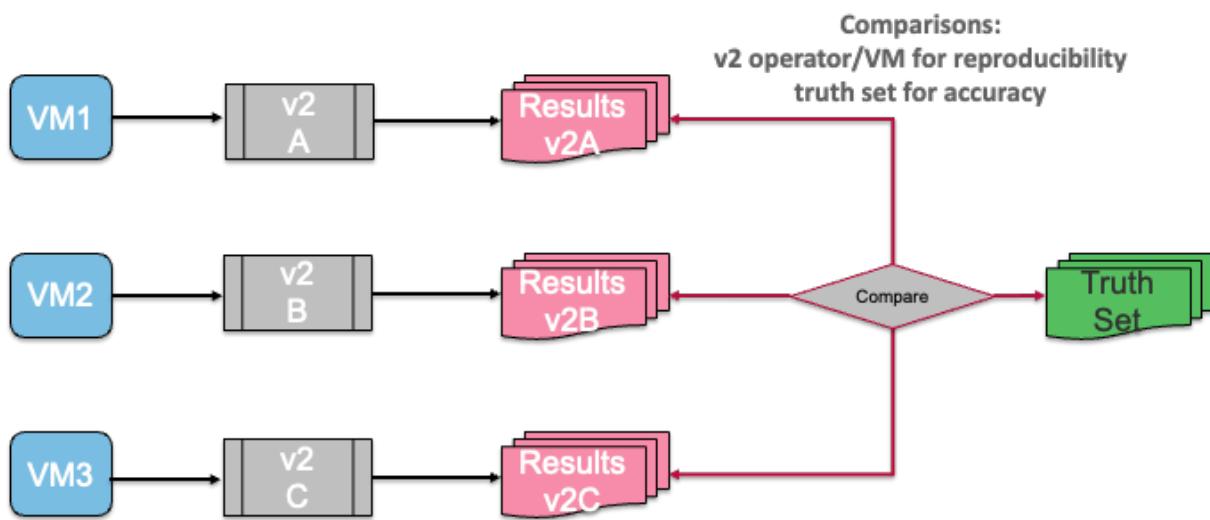
<https://github.com/Illumina/hap.py>

<https://github.com/ga4gh/benchmarking-tools/blob/master/resources/high-confidence-sets/giab.md>

3. WES PIPELINE – VALIDATION METHOD

Validation Design: The validation design is shown below. It will be used to compare the results of the enhanced pipeline (v2) to industry-accepted truth sets as well as assess inter-operator and inter-VM variability.

Validation Plan for WES



Acceptance Criteria:

Small Variants

Recall: > 0.90
Precision: > 0.90

CNV

Sensitivity: > 0.8
Specificity: > 0.8

Note: The GIAB BED file defines the high-confidence regions for scoring small variants. Other regions will be excluded from the analysis.

Investigation of Discrepancies:

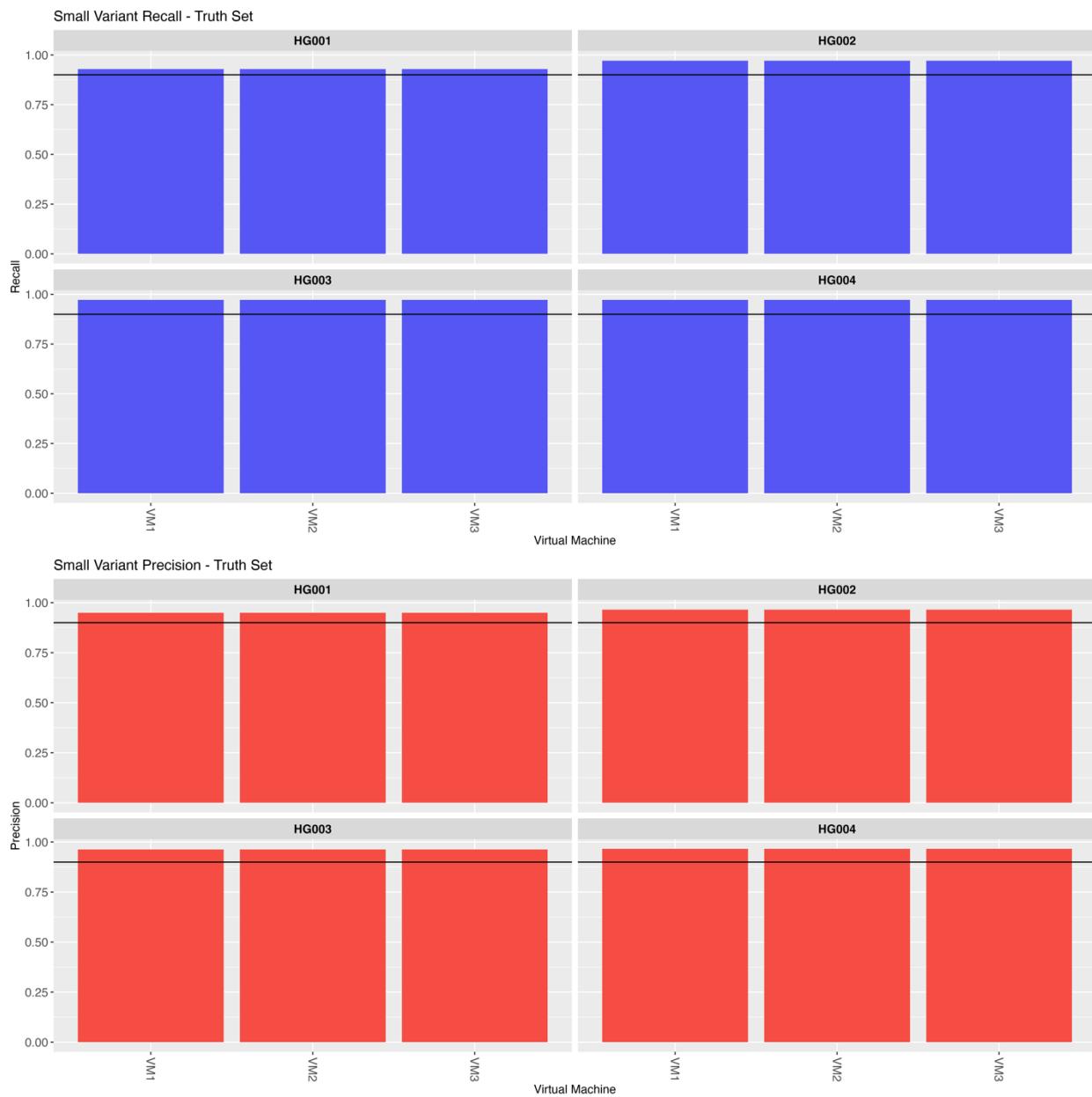
If the comparison metrics do not meet the acceptance criteria, an investigation will be carried out by members of NCI-CGBB and EM. The members to perform the investigation will be designated by Daoud Meerzaman, based on expertise and availability. The investigation should last no longer than 10 business days, at which time a report that outlines the problem and suggests solutions will be presented to Daoud Meerzaman.

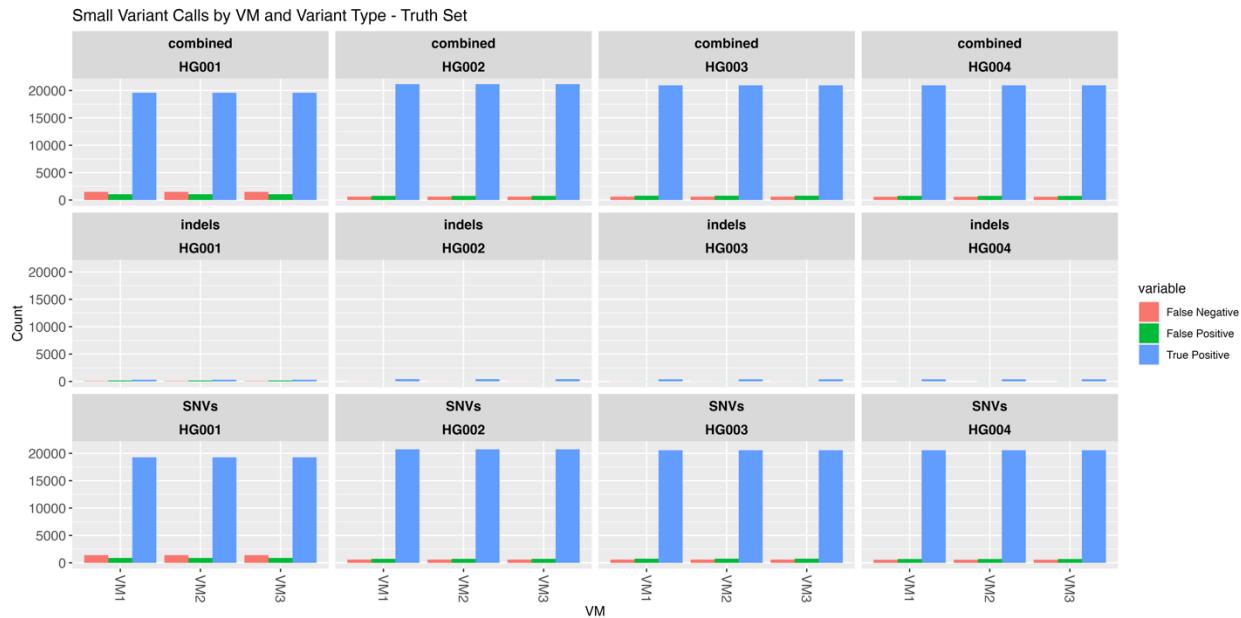
4. WES PIPELINE – VALIDATION RESULTS

Pipeline Performance Specifications

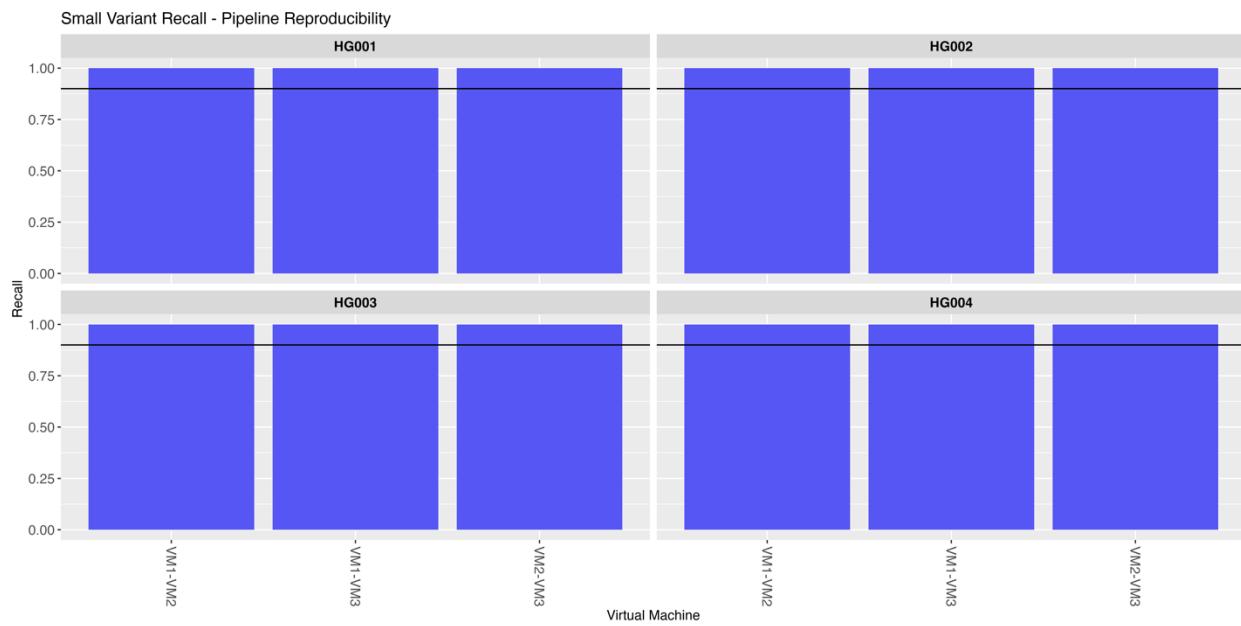
Small Variants		
	Truth Set	Pipeline Reproducibility
Recall	0.9610	1.000
Precision	0.9608	1.000
CNV		
	Truth Set	Pipeline Reproducibility
Sensitivity	0.8201	1.000
Specificity	0.9406	1.000
Accuracy	0.9239	1.000

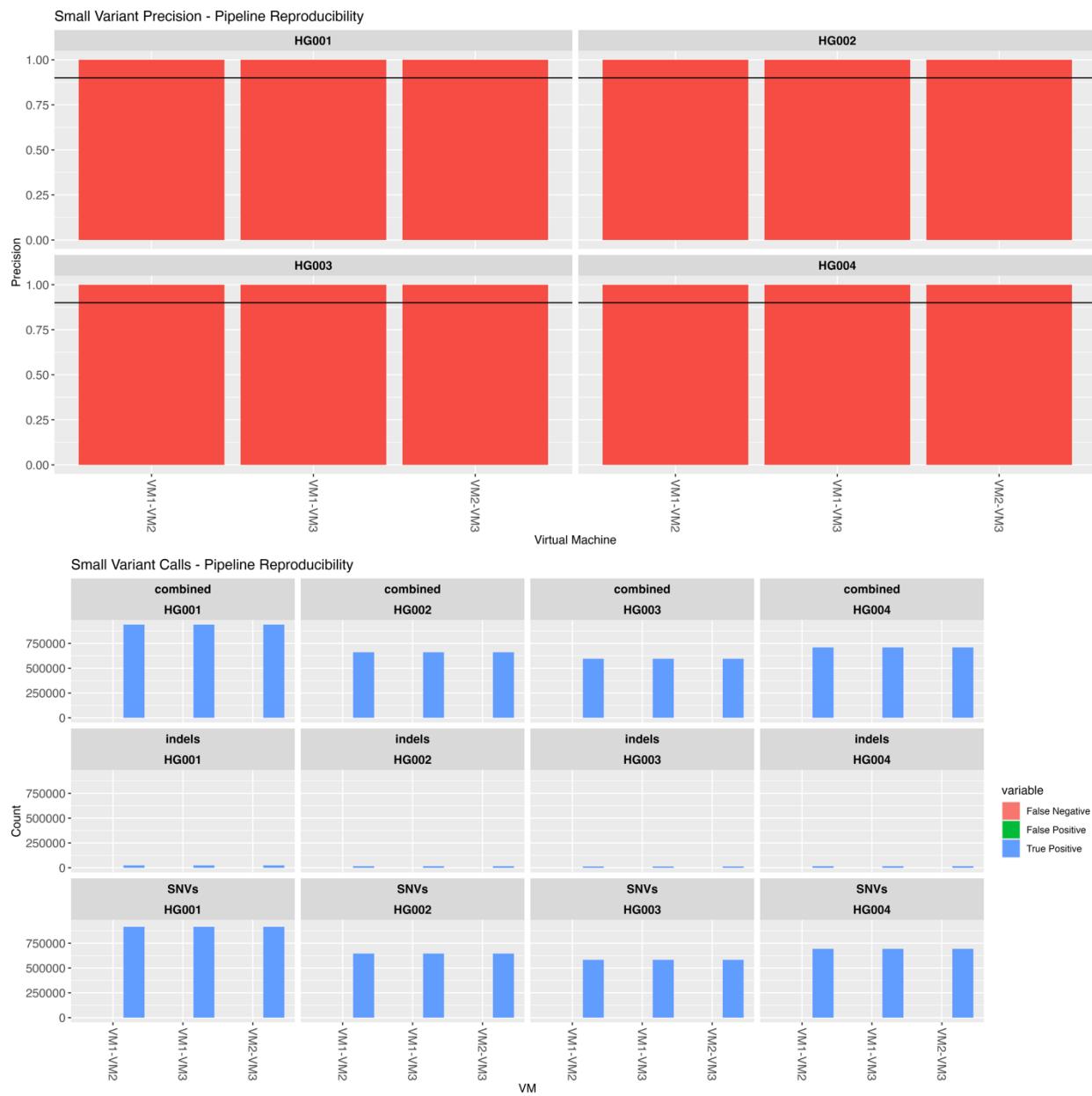
Small Variants - Comparison to Truth Set



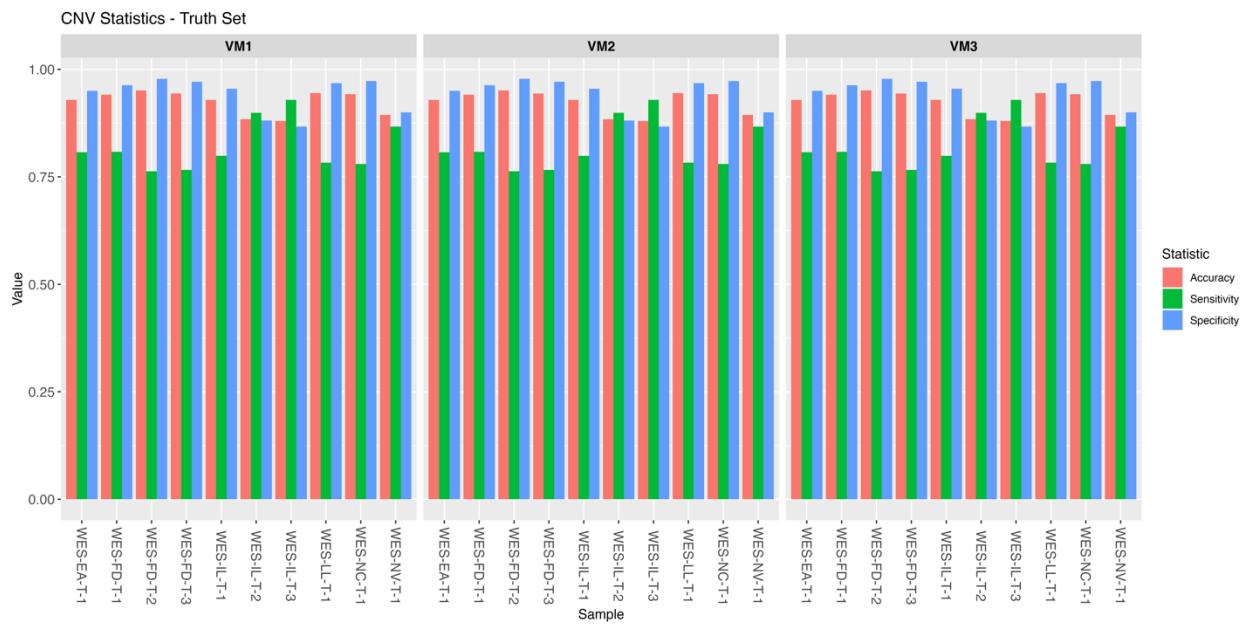


Small Variants – Pipeline Reproducibility

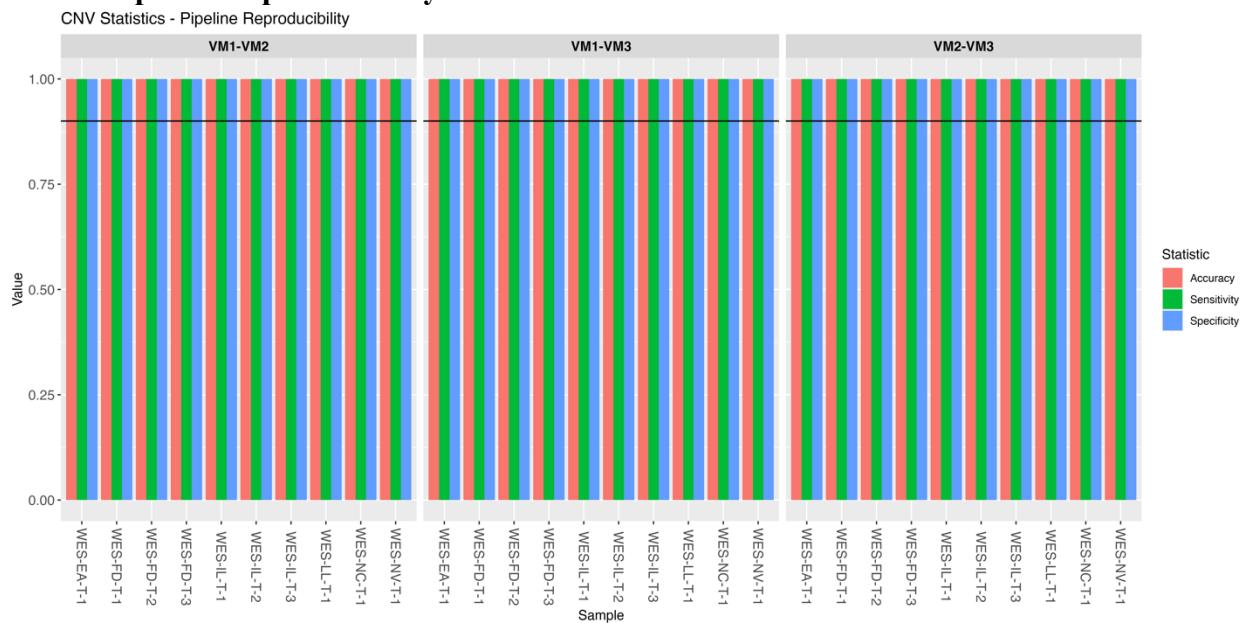




CNV - Comparison to Truth Set



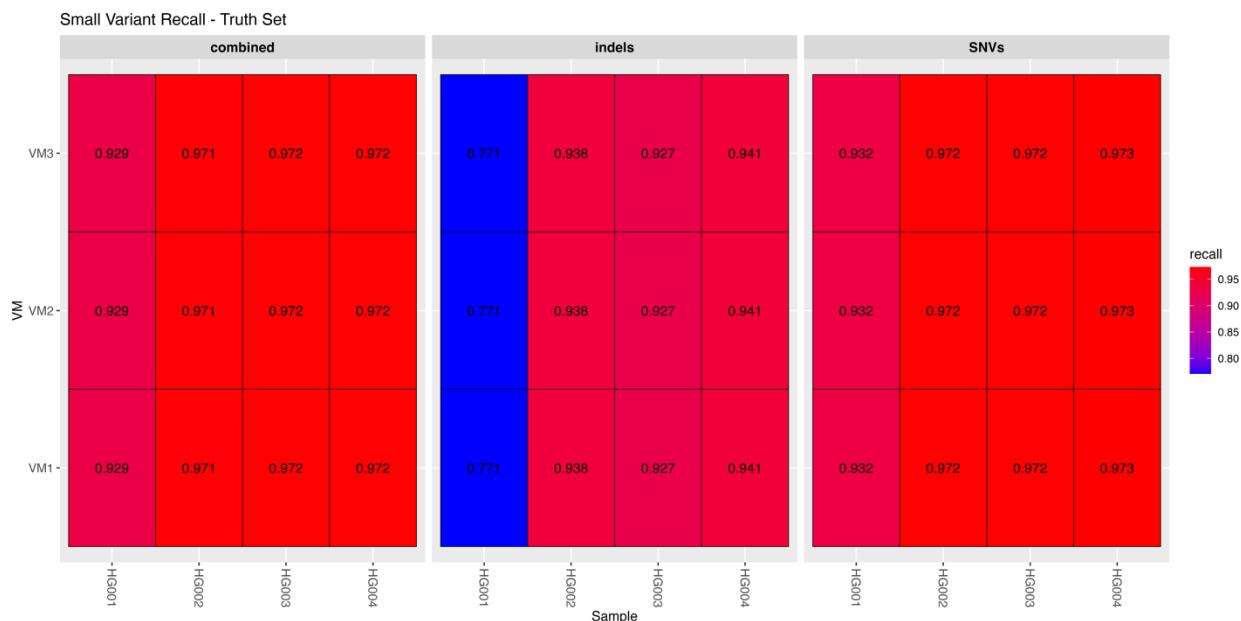
CNV – Pipeline Reproducibility

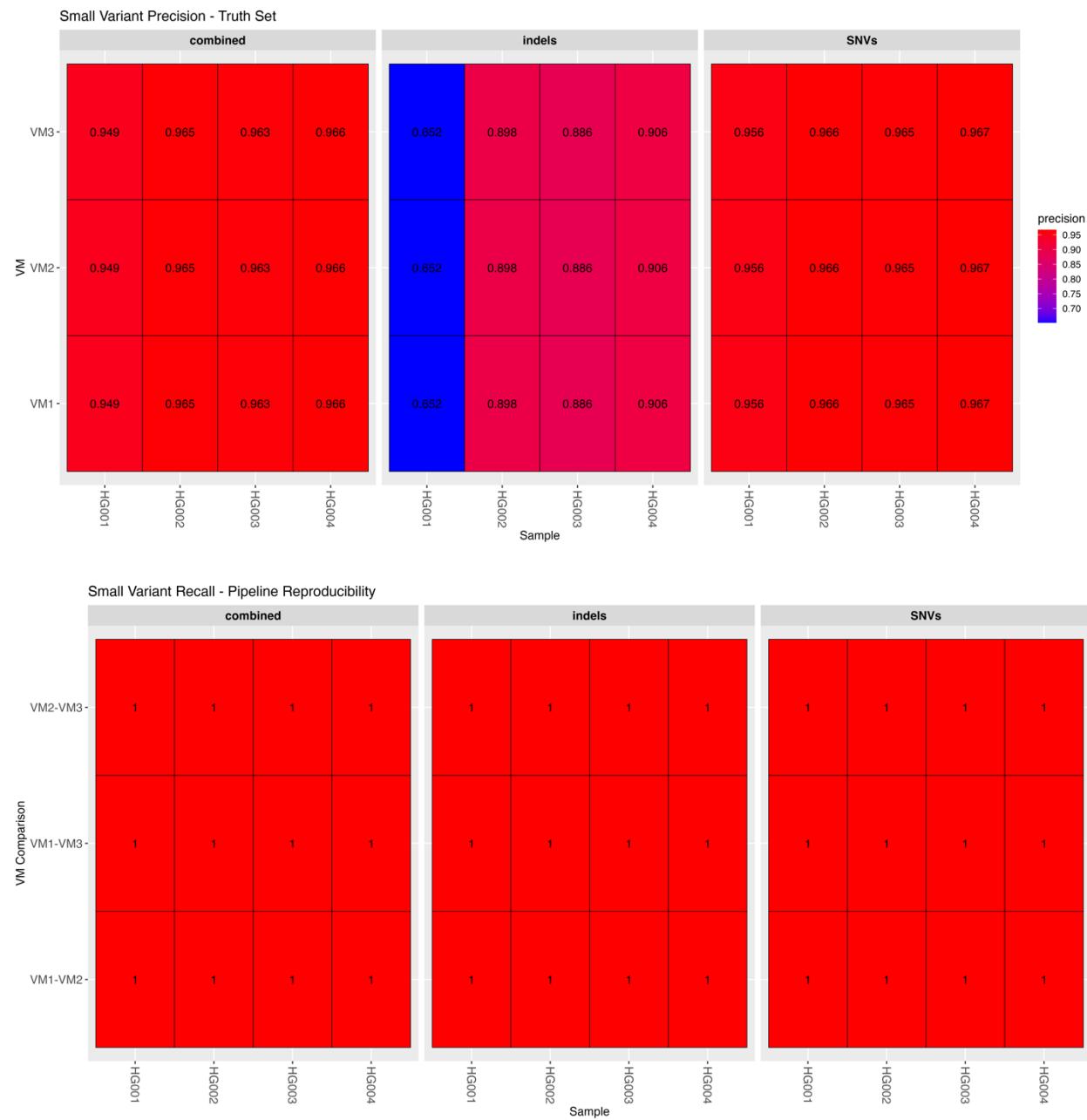


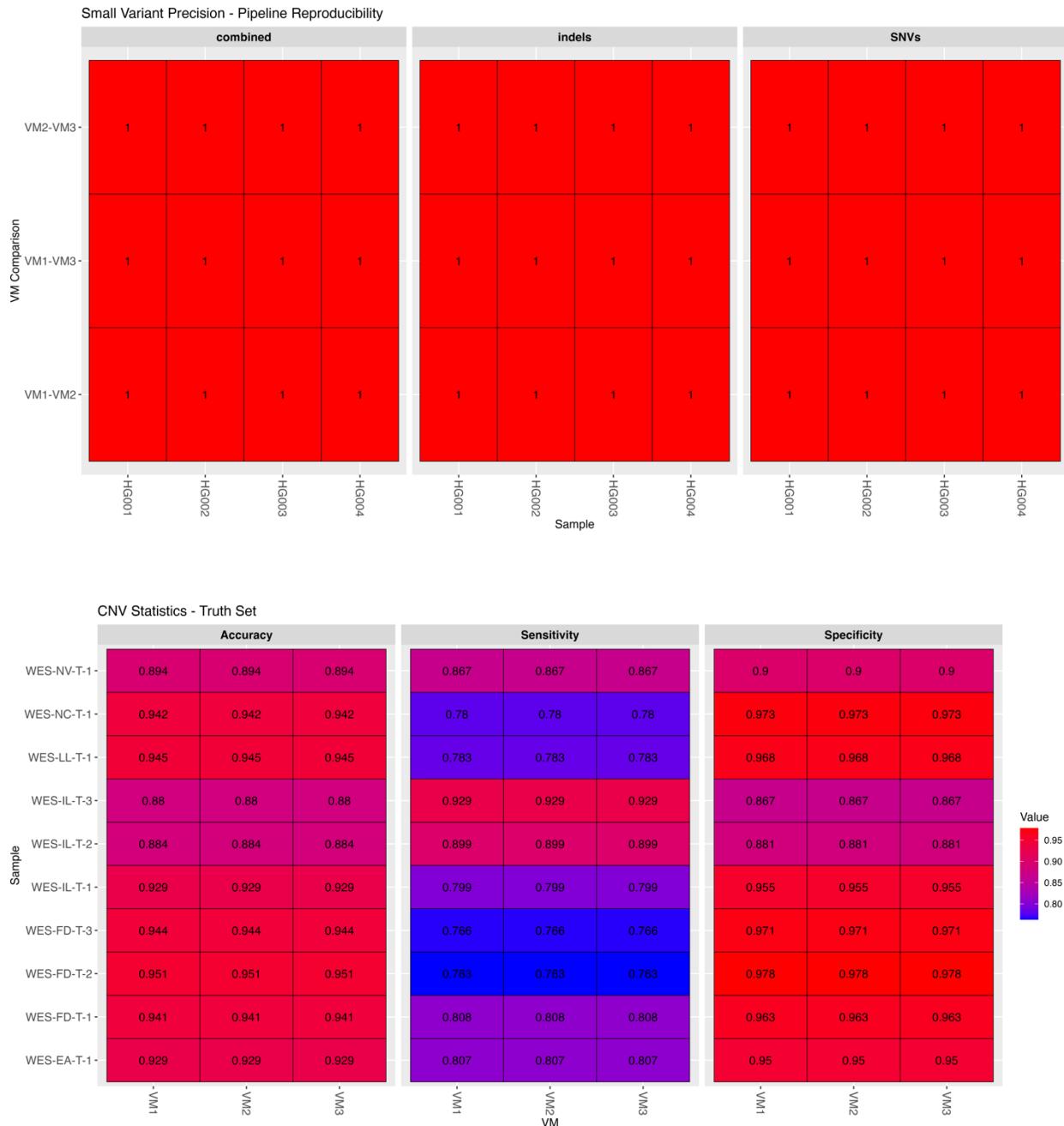
Conclusions: Based on these results, the pipeline meets the acceptance criteria for analytic performance as described prior to commencement of the validation activities. The pipeline generates highly reproducible results, as exemplified by the perfectly matched results created by

three users on three virtual machines. The pipeline also shows excellent performance metrics for calling small variants. The pipeline does show reduced sensitivity (however, above the acceptance criteria) for small indels and CNVs. This characteristic is not unexpected given the known limitations of the short-read whole exome sequencing in general (for example, see *Masood et al. Evaluation of somatic copy number variation detection by NGS technologies and bioinformatics tools on a hyper-diploid cancer genome (2023)*). The specificity of the pipeline is excellent regardless of the variant type, which is a desirable trait for pipelines used for analysis of clinical data. Given these results, it is suggested that this pipeline be released for analysis of CIMAC-CIDC WES data.

5. WES PIPELINE – APPENDIX I – TABULAR RESULTS







CNV Statistics - Pipeline Reproducibility									
Sample	Accuracy			Sensitivity			Specificity		
	VM1-VM2	VM1-VM3	VM2-VM3	VM1-VM2	VM1-VM3	VM2-VM3	VM1-VM2	VM1-VM3	VM2-VM3
WES-NV-T-1-	1	1	1	1	1	1	1	1	1
WES-NC-T-1-	1	1	1	1	1	1	1	1	1
WES-LL-T-1-	1	1	1	1	1	1	1	1	1
WES-IL-T-3-	1	1	1	1	1	1	1	1	1
WES-IL-T-2-	1	1	1	1	1	1	1	1	1
WES-IL-T-1-	1	1	1	1	1	1	1	1	1
WES-FD-T-3-	1	1	1	1	1	1	1	1	1
WES-FD-T-2-	1	1	1	1	1	1	1	1	1
WES-FD-T-1-	1	1	1	1	1	1	1	1	1
WES-EA-T-1-	1	1	1	1	1	1	1	1	1