# Benchmarking workflows to assess performance and suitability of germline variant calling pipelines in clinical diagnostic assays.

Vandhana Krishnan[1,2], [†]Sowmi Utiramerur[3], Zena Ng[3], Somalee Datta[2], Michael P. Snyder[1,2], Euan A. Ashley[1,4,5]

[1]Department of Genetics, School of Medicine, Stanford University, Stanford, CA
[2]Stanford Center for Genomics and Personalized Medicine, Stanford University, Palo Alto, CA
[3]Clinical Genomics Program, Stanford Health Care, Stanford, CA
[4]Department of Cardiovascular Medicine
[5]Department of Biomedical Data Science


Current address of Somalee Datta: School of Medicine, Information Resources and Technology (IRT) – Research IT, Stanford University, Palo Alto, CA

Email address of Vandhana Krishnan: vandhana.krishnan@stanford.edu

Email address of Zena Ng: zng@stanfordhealthcare.org

Email address of Somalee Datta: somalee@stanford.edu

Email address of Michael P. Snyder: mpsnyder@stanford.edu

Email address of Euan A. Ashley: euan@stanford.edu

[†]Corresponding author: Sowmi Utiramerur: Email: sutiramerur@stanfordhealthcare.org

# 1 Abstract

3 Benchmarking the performance of complex analytical pipelines is an essential part of
4 developing Laboratory Developed Assays (LDT). Reference samples and benchmark calls
5 published by Genome in a Bottle (GIAB) Consortium have enabled the evaluation of
6 analytical methods. However, the performance of such methods is not uniform across the
7 different regions of the genome/exome and different variant types and lengths. Here we
8 present a scalable and reproducible, cloud-based benchmarking workflow that can be used by
9 clinical laboratories to rapidly access and validate the performance of LDT assays, across
10 their regions of interest and reportable range, using a broad set of benchmarking samples.

# 12 Keywords

13 Benchmarking, workflow, GIAB reference genomes, precision, recall, truth set, docker,
14 germline variants, laboratory developed assays

# 16 Background

17 Next Generation Sequencing (NGS) and analytical methods developed to detect various
18 forms disease-causing polymorphisms are now routinely being used by clinical laboratories
19 to determine the molecular etiology of complex diseases/disorders and in many cases to make
20 critical treatment course decisions.  In the past two decades, many polymorphisms in the
21 human genome have been identified and validated that serve as predictive, diagnostic, and
22 prognostic markers for complex inherited diseases. These genomic disease markers can be of
23 different forms such as Single Nucleotide Variants (SNVs), small INsertions and DELetions
24 (INDELs), large deletions and duplications (del/dups), and Copy Number Variations (CNVs)
25 and can vary in size from a single base change to several Mega Bases (MB) in length and
26 even whole chromosomal polysomy. Clinically relevant polymorphisms occur in different
27 regions of the genome, including exonic, splice-sites, and deep-intronic regions. These
28 polymorphisms also happen in various forms, including single base changes within high
29 entropic regions, copy number changes to homopolymer repeats and copy number changes to
30 Short Tandem Repeat (STR) regions. NGS platforms used to detect these polymorphisms;
31 owing to their different sequencing chemistry and signal processing methods; have very
32 different error modes and hence very different analytical performance across the different
33 regions of the genome. Consequently, analytical methods specific to various NGS platforms
34 such as Illumina, Ion Torrent, Pacific Biosciences, and Oxford Nanopore have been
35 developed to both account for and correct the errors particular to these sequencing platforms.
36 This has resulted in a dizzying array of combinations of sequencing platforms and analytical
37 methods available to a clinical diagnostic laboratory to develop their LDT assay.

39 Benchmarking methods and pipelines are essential to accurately assess the performance of
40 sequencing platforms and analytical methods before they are incorporated into clinical
41 diagnostic assays. Genome In A Bottle (GIAB) consortium hosted by NIST has characterized
42 the pilot genome (NA12878/HG001) (1) and six samples from the Personal Genome Project
43 (PGP) (2). These benchmark calls for SNVs and small INDELs (1-20bp) from reference
44 samples can be used for optimization, performance estimation, and analytical validation of
45 LDT assays using complex analytical pipelines with multiple methods to detect
46 polymorphisms in the genome. Global Alliance for Genomics and Health (GA4GH)
47 benchmarking team have developed standardized tools (3)  to evaluate the performance
48 metrics of germline variant callers used primarily in research applications.

49

50 Clinical Laboratory Improvement Amendments (CLIA) requires all laboratories using LDT
51 to establish the test's performance specifications such as analytical sensitivity, specificity,
52 reportable range, and reference range (4). College of American Pathologist (CAP) laboratory
53 standards for NGS based clinical diagnostic (5) not only require the laboratories to assess and
54 document the performance characteristics of all variants within the entire reportable range of
55 LDTs but also obtain the performance characteristics for every type and size of variants that
56 are reported by the assay. Laboratories are also required to assess the performance
57 characteristics for clinically relevant variants such as $\Delta F508$ and IVS8-5T (6) mutations in a
58 CFTR assay. CAP guidelines also require laboratories to periodically (determined by the
59 laboratory) assess and document the analytical performance characteristics to ensure that the
60 LDT is continuing to perform as expected over time.

61

62 Benchmarking workflows/pipelines that are highly scalable, reproducible and capable of
63 reporting the performance characteristics using a large number of reference and clinical
64 samples within multiple highly stratified regions of interest are essential for clinical
65 laboratories to optimize and routinely assess the performance of their LDT assays.

66

# 67 Results

68

69 Our goal was to develop a benchmarking workflow that any clinical laboratory could use to
70 quickly evaluate and compare the performance characteristics of all suitable secondary
71 analysis pipelines. Benchmarking workflow should further help optimize the analytical
72 pipeline based on well-defined performance metrics and finally produce a thorough analytical
73 validation report to justify the use of the analytical pipeline in their diagnostic assay to
74 regulatory authorities such as CLIA and CAP.

75

76 To test the abilities of our benchmarking workflow, we used it to compare two analytical
77 pipelines commonly used for germline variant calling 1. Pipeline based on Broad Institute's
78 best practices guidelines using GATK HaplotypeCaller v3.7 and 2. SpeedSeq pipeline (7)
79 based on FreeBayes v0.9.10 (8) as the primary variant calling engine. GATK
80 HaplotypeCaller based pipeline was chosen over the FreeBayes pipeline as it out-performed
81 in the detection of small-INDELs (1 – 20 base pairs).

82

83 The performance characteristics of the analytical pipeline using GATK v3.7 was further
84 optimized using benchmarking metrics generated using the five GIAB reference samples and
85 four GeT-RM samples (see Methods) with known pathogenic variants. Also, it is critical for
86 the clinical laboratories developing NGS based LDT assays to accurately determine the
87 reportable range to avoid misdiagnosis leading to wrong treatment decisions. To this effect,
88 we evaluated the performance metrics using the benchmarking workflow in three distinct
89 genomic regions of interest (see Methods for details).

90

91 Although we have the benchmarking results for the region, including coding exons in all the
92 RefSeq genes, we have omitted those findings in this section and focus on the clinically
93 relevant regions.

94

95 Table 1: Benchmarking metrics for SNPs within coding exons of clinically relevant ~7000
96 genes (as specified in Methods).

97

| GIAB genome / NIST ID | Number of bases | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|---|
| NA12878 | 13728555 | 7803 | 7781 | 4 | 22 | 13720748 | 100 | 99.95 | 99.72 |
| NA24143 | 12549224 | 7470 | 7460 | 14 | 10 | 12541740 | 100 | 99.81 | 99.87 |
| NA24149 | 12538042 | 7495 | 7485 | 19 | 9 | 12530529 | 100 | 99.75 | 99.88 |
| NA24385 | 12626866 | 7452 | 7436 | 0 | 16 | 12619414 | 100 | 100 | 99.79 |
| NA24631 | 12808688 | 7591 | 7581 | 6 | 10 | 12801091 | 100 | 99.92 | 99.87 |

98

99 Table 2: Benchmarking metrics for SNPs in whole exome regions, including non-coding
100 exons, splice sites (+/- 20 bp) and clinically relevant deep intronic regions.

101

| GIAB genome / NIST ID | Number of bases | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|---|
| NA12878 | 71152019 | 57822 | 57024 | 491 | 776 | 71093728 | 100 | 99.15 | 98.66 |
| NA24143 | 65657646 | 55975 | 55340 | 669 | 611 | 65601026 | 100 | 98.81 | 98.91 |
| NA24149 | 65597266 | 55518 | 54827 | 669 | 669 | 65541101 | 100 | 98.79 | 98.79 |
| NA24385 | 65948744 | 56068 | 55329 | 389 | 705 | 65892321 | 100 | 99.30 | 98.74 |
| NA24631 | 66988987 | 56948 | 56303 | 394 | 643 | 66931647 | 100 | 99.31 | 98.87 |

102

103 Tables 1 and 2 show the benchmarking metrics for SNPs in all five GIAB samples within the
104 clinically relevant genes and whole exome regions, respectively. The precision, recall, and
105 NPA metrics for SNPs are uniform across all the reference samples, and there is no sample
106 bias in the results for some of the better-characterized samples such as NA24385 and
107 NA12878. Performance metrics for SNPs within the clinically relevant gene region is
108 significantly better than those within the whole exome region. Recall metrics, in particular,
109 are a percentage point better in the clinically pertinent gene region, across all reference
110 samples. This is attributable to the fact that many genes have isoforms, resulting in higher
111 alignment errors, and some genes have either very high or very low GC content, resulting in
112 higher than average sequencing errors within these regions of the genome. The finding is of
113 great clinical significance, since the reportable region of most inherited disease/disorder,
114 LDT assay is limited to the clinically relevant genes and thereby the overall performance
115 characteristics of the assay is better than that estimated over either the whole genome or
116 whole exome regions.

117

118 Table 3: Benchmarking metrics for indels of different size ranges in NA24385 (truth set
119 NIST v3.3.2, total bases = 12,626,866) for the regions within ~7000 clinically relevant genes
120 (as specified in Methods).

121

| Size of indels in NA24385 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 145 | 136 | 12 | 9 | 12626709 | 100 | 91.89 | 93.79 |
| 11-20 | 9 | 9 | 0 | 0 | 12626857 | 100 | 100 | 100 |
| 21-50 | 3 | 3 | 0 | 0 | 12626863 | 100 | 100 | 100 |
| All Indels | 157 | 148 | 12 | 9 | 12626697 | 100 | 92.50 | 94.27 |

122

123 Table 4: Benchmarking metrics on the number of indels of different size ranges in NA24385
124 (truth set NIST v3.3, total bases = 65,948,744) for the whole exome regions including non-
125 coding exons, splice sites (+/- 20 bp) and clinically relevant deep intronic regions.

| Size of indels in NA24385 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 5169 | 4727 | 872 | 442 | 65942703 | 100 | 84.43 | 91.45 |
| 11-20 | 203 | 188 | 10 | 15 | 65948531 | 100 | 94.95 | 92.61 |
| 21-50 | 67 | 56 | 3 | 11 | 65948674 | 100 | 94.92 | 83.58 |
| All Indels | 5362 | 4920 | 885 | 468 | 65942471 | 100 | 84.75 | 91.27 |

126

127  Tables 3 and 4 provide the indel benchmarking metrics for sample NA24385 in the clinically
128  relevant and whole exome regions, respectively. As expected, the benchmarking workflow
129  reveals that the performance metrics for INDELs are lower than those for SNPs. However,
130  the stratification by INDEL size, helped us determine the reference range for INDELs (1- 20
131  base-pairs). The recall metric for INDELs larger than 20 base-pairs is significantly lower than
132  the recall for INDELs 1 – 20 base-pairs. As in the case of SNPs, performance metrics for
133  INDEL detection within the clinically relevant genes of interest is better than the whole
134  exome region.

135

136  The benchmarking results of the other GIAB reference samples in the clinically relevant and
137  whole exome regions can be obtained in the Supplementary Materials Table S1-S4 and Table
138  S5-S8, respectively. The histogram for the indel size distribution in the NA24385 reference
139  sample for the whole exome region is in Supplementary Material as Fig S1. The histograms
140  of indel size distributions for GIAB samples in both the whole exome and clinically relevant
141  regions are available in our github repository - vandhana/stanford-benchmarking-workflows.

142

143  Table 5: Validation of the presence of the truth variants in the GeT-RM samples (as specified
144  in Methods) using our variant calling pipeline.

145

| GeT-RM Sample ID | Chromosome:Position | Truth Variant | Truth Variant Detected |
|---|---|---|---|
| NA04408 | 15:91310152 | TATC -> T | Yes |
| | 15:91310156 | T -> TA | Yes |
| | 15:91310158 | A -> ATTC | Yes |
| NA14090 | 17:41276044 | ACT -> A | Yes |
| NA14170 | 13:32914437 | GT -> G | Yes |
| NA16658 | 10:43609103 | G -> T | Yes |

146

147  Finally, our benchmarking workflow was able to confirm that our variant calling pipeline can
148  detect all the clinical variants in GeT-RM samples listed in Table 5.

149

150  To get all the metrics produced by hap.py and other output files including plots from our
151  benchmarking workflow for each reference sample, please refer to the Supplementary Data
152  files.

153

# Discussion

155

156  GIAB consortium has helped developed standards for genomic data to evaluate the
157  performance of NGS sequencing platforms and analytical methods used for alignment and
158  variant calling. The precisionFDA platform has enabled the genomics community to develop
159  and deploy benchmarking tools that can evaluate the performance of analytical methods
160  against the gold standard datasets. These benchmarking tools, along with accuracy
161  challenges, has led to the development of highly accurate variant calling methods. However,
162  the requirements of a clinical diagnostic laboratory go beyond the simple evaluation of

163  performance characteristics of an analytical pipeline against one or more reference samples.
164  Our purpose was to build a benchmarking workflow to meet the assay optimization and
165  validation needs of a clinical laboratory. The primary benefit of our benchmarking workflow
166  is that it allows for the assay performance to be evaluated using a broad set of both reference
167  samples with a large number of gold-standard variant calls and clinical samples with a small
168  number of clinical variants, that are specific to the diagnostic assay being evaluated. The
169  benchmarking workflows enable the clinical laboratories to establish the reporting range of
170  the diagnostic assay by estimating the performance within multiple regions of interest.
171
172  Unlike web-based benchmarking apps, such as those provided by the precision FDA platform
173  or GA4GH, our benchmarking framework can be seamlessly integrated with any variant
174  calling pipeline in the user's software environment. Thus, our benchmarking workflows
175  enable ease of use and avoid the transfer of sensitive data to different locations, which could
176  be non-Protected Health Information (PHI) compliant.
177
178  Our benchmarking modules if integrated with deployment tools, such as Jenkins (9) and
179  CircleCI (10), that work on the principle of continuous integration and continuous
180  delivery/deployment (CI/CD), it provides a foolproof way of examining consistency in
181  results. In this era where workflows generating reproducible results are gaining attention,
182  easy incorporation of workflows with CI/CD tools is a nice feature to have.
183
184  The benchmarking workflow is distributed using human-readable YAML (11) format, and it
185  might limit direct porting to existing WDL based workflows such as those published by the
186  Broad Institute (12, 13). Similarly, conversion of the benchmarking YAML files to Common
187  Workflow Language (CWL) format is required to run workflows published by GA4GH (14-
188  16). However, since we have used docker images for the software tools used within the
189  benchmarking framework, portability to other runtime environments should not take a
190  significant effort for a bioinformatician.
191

192

193

# Conclusions
194

195

196  Benchmarking variants is a critical part of implementing variant calling pipelines for research
197  or clinical purposes. Here, we have successfully implemented benchmarking workflows that
198  generate metrics such as specificity, precision, sensitivity for germline SNPs, and indels in
199  whole exome sequencing data. Also, indel size distributions even in the form of histograms
200  are provided. Combining these benchmarking results with validation using known variants of
201  clinical significance in publicly available cell lines, we were able to establish our variant
202  calling pipelines in a clinical setting. Our benchmarking workflow can serve as a plug-in to
203  any existing variant calling pipeline to work as an integrated unit or be used as a separate
204  module as well.

205

# Methods
206

207

208  ***Benchmarking workflow***

209  The benchmarking workflow, as illustrated in Figure 1, is a sequence of steps required to
210  perform a rapid and comprehensive analytical validation of a clinical diagnostic assay based
211  on germline variants. The benchmarking workflow can be easily integrated with any
212  secondary-analysis pipeline used in a diagnostic assay to call germline variants, and our
213  workflow accepts germline variants (SNVs and small INDELs) in Variant Call Format VCF
214  v4.1(17) or higher. The workflow takes one or more stratification files specifying the regions
215  of interest in BED (18) format and generates a comprehensive analytical validation report
216  detailing the performance characteristics of the assay within each of the specified regions of
217  interest. Benchmark variant calls that are considered as ground truths for each of the
218  reference sample used to evaluate the analytical performance can be also be specified in VCF
219  format.
220
221
222  **Figure 1. Schematic diagram of the benchmarking framework used in this study**
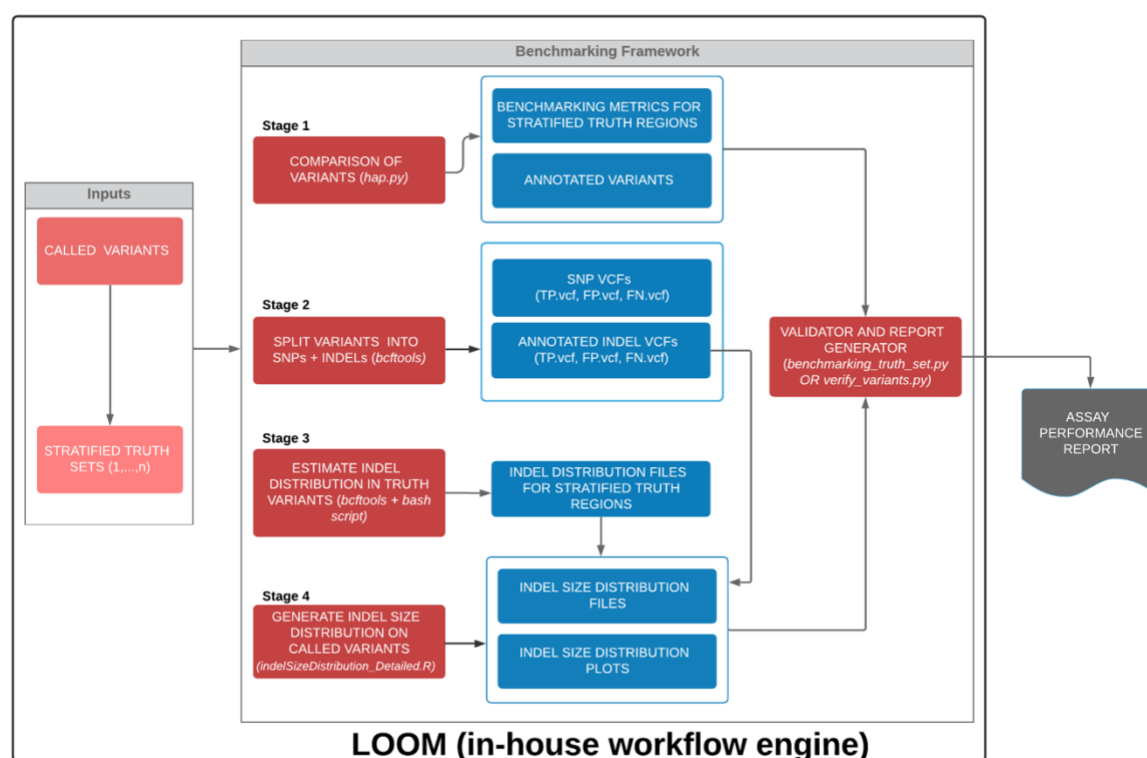223



224
225  Figure 1 legend: All the stages in the benchmarking workflow have been dockerized. The
226  docker images are available in DockerHub as specified in the Methods section.
227
228
229  The first step in the benchmarking process involves the comparison of input variants
230  generated by the analytical pipeline with the benchmark variant calls within each region of
231  interest. The variant calls are compared using hap.py (19, 20), which is capable of haplotype
232  construction from individual genotype calls and is recommended by GIAB consortium and
233  GA4GH. The variant comparison step is performed for each of the stratification or region of
234  interest file specified as input, and hap.py generates a single output VCF file classifying the
235  variant calls defined in the input and truth VCF files as either True Positive (TP), False
236  Positive (FP) or False Negative (FN).
237

238 Step two in the benchmarking workflow splits the variant calls annotated using hap.py by
239 variant type (SNPs and small INDELs) and by variant classification (TP/FP/FN). This step is
240 executed within the workflow for each of the stratification or region of interest file specified.
241 The VCF files are split by variant type using bcftools (21), and a bash script is used to further
242 split the variant calls by the variant classification. This allows the workflow to generate the
243 performance metrics for each of the variant types reported by the diagnostic assay.
244
245 Steps two and three of the benchmarking workflow (see Figure 1.) were used to generate a
246 histogram of small INDELs by size. The bins used for INDEL size histograms were a. 1
247 base-pair, b. 2-5 base-pairs, c. 6-10 base-pairs, d. 11 – 20 base-pairs, e. 21 – 50 base-pairs,
248 and f. Greater than 50 base-pairs. The R script - indelSizeDistribution_Detailed.R (code in
249 Additional File 1) then calculates the performance metrics of the assay for each of the INDEL
250 size bins. The Python script – benchmarking_truth_set.py (Additional File 2) consolidates the
251 benchmarking metrics previously obtained, calculates the NPA related metrics combining
252 some of the bin size ranges (user preferred) for all reference samples provided.
253
254 In addition to benchmarking call sets for well-characterized reference samples published by
255 the GIAB consortium, the benchmarking workflow allows for clinical laboratories to specify
256 addition samples with clinically relevant variants as ground truths to estimate the analytical
257 performance of the assay for specific variant types such as $\Delta F508$ and IVS8-5T in CFTR
258 panels. Python script – verify_variants.py (Additional File 3) accepts the ground-truth variant
259 call sets to confirm the presence/absence of these variants in the VCF files generated by the
260 variant calling pipeline. The details on the usage of the above scripts and associated
261 README file are available in our public repository (also see Supplementary Materials).
262
263 Finally, the benchmarking workflow generates a comprehensive analytical validation report
264 using all the provide benchmarking ground-truth call sets.
265
266
267 *Scalability and Reproducibility of Benchmarking workflow*
268
269 The benchmarking workflow is designed to be repeatable and reproducible by using Docker
270 containers for all software and bioinformatics components used within the workflow (see
271 Table 6.). The workflow is distributed in human-readable data serialization format YAML
272 v1.2, and the workflow can be readily executed using the workflow execution manager –
273 LOOM (22). The workflow definition file – Benchmarking.yaml (see Supplementary
274 Materials) can also be easily ported to Common Workflow Language (CWL) or Workflow
275 Definition Language (WDL) formats and can be executed using workflow execution
276 managers such as Toil (23, 24) and Cromwell (25).
277
278
279 Table 6. Docker containers and DockerHub repository location for each of the individual
280 software components used in the benchmarking workflow.
281

| Software Component | Docker Container |
|---|---|
| hap.py v 0.2.10 | sowmiu/happy:latest |
| bcftools | vandhanak/bcftools:1.3.1 |
| indelSizeDistribution_Detailed.R | vandhanak/rbase:3.3.2 |

282

283

## *Golden/ground-truth callsets*

The golden/ground-truth sets for five reference and PGP genomes are currently available - NA12878 (CEPH family's daughter), NA24143 (AJ mother), NA24149 (AJ father), NA24385 (AJ son), and NA24631(Chinese son) and these reference call sets were used in this benchmarking study. GIAB provides a high confidence regions file and a high confidence VCF file, and as recommended by GIAB, only the high confidence calls were used in the evaluation of the assay's performance characteristics. The NIST versions and their corresponding FTP site locations used for the above samples in this study can be found in the Supplementary Material.

In addition to the GIAB reference samples, samples with known pathogenic germline variants (see Table 2.) for various inherited diseases/disorders were chosen from Genetic Testing Reference Materials Coordination Program (GeT-RM) (26-30)

Table 7. GeT-RM sample ids and location of ground-truth variants in GRCh37 coordinates.

| GeT-RM Sample ID | Chromosome:Position | Truth Variant |
|---|---|---|
| NA04408 | 15:91310152 | TATC -> T |
| | 15:91310156 | T -> TA |
| | 15:91310158 | A -> ATTC |
| NA14090 | 17:41276044 | ACT -> A |
| NA14170 | 13:32914437 | GT -> G |
| NA16658 | 10:43609103 | G -> T |

## *Stratification or Regions of Interest (ROI) BED files.*

Three stratification files were used to evaluate the performance characteristics of an inherited Whole Exome Sequencing (WES) assay.

1. Coding Exons for all known transcripts in RefSeq genes: RefSeq gene names, transcripts, and coordinates of all coding exons were obtained from the UCSC genome browser(31, 32).
2. Clinically relevant regions of the human genome: Clinically relevant regions were determined by intersecting coordinates of all known pathogenic variants reported in OMIM (33), ClinVar (34) and DECIPHER v9.28 (35) with the all exon regions (Coding and Non-Coding) file for RefSeq genes obtained from UCSC genome browser. The exonic coordinates were later extended by 20 base-pairs on either end to include canonical and non-canonical splice sites. Deep-intronic regions with pathogenic variants were added to the exonic regions to generate the final clinically relevant regions (BED) file.
3. Whole Exome regions file for RefSeq genes was obtained from UCSC genome browser. The exon regions were extended by 20 base-pairs on either end to include splice sites.

## *Benchmarking metrics*

Precision and recall are benchmarking metrics provided as output by hap.py. The true positives (TP), false positives (FP), and false negatives (FN) are counted as described by the

327 developers of hap.py (20). Again, as explained by the authors of hap.py, precision and recall
328 are calculated using the below formulae:

330 *Precision = True Positives/(True Positives + False Positives)*

332 *Recall = True Positives/(True Positives + False Negatives)*

334 Other metrics reported by hap.py such as variants outside the high confidence truth set
335 regions and transition or transversion SNP type can be found in the extended.csv files
336 included in the Supplementary Materials.

338 The total number of bases per sample in a particular region of interest as specified by the
339 corresponding bed file was computed using a bash command provided in the Supplementary
340 Materials.

342 True negatives (TN) and Total Negatives are computed using the following:

344 *TN = Total number of bases in the region of interest – (True Positives + False Positives +*
345 *False Negatives)*

347 *Total Negatives = True Negatives + False Positives*

349 The Negative Percentage Agreement (NPA) or specificity as recommended by the FDA (36)
350 is calculated using

352 NPA = *True Negatives/Total Negatives*

# List of abbreviations

358 NIST – National Institute of Standards and Technology
359 GIAB – Genome in a bottle consortium
360 SNPs – Single nucleotide polymorphisms
361 Indels – insertions/deletions
362 WES – Whole Exome Sequencing
363 NPA – Negative Percent Agreement
364 TN – True Negative
365 TP – True Positive
366 FN – False Negative
367 FP – False Positive
368 OMIM – public database containing the human genes, their genetic phenotypes and
369 associations with genetic disorders (Online Mendelian Inheritance in Man)
370 DECIPHER – public database with genotypic and phenotypic data from ~30,000 individuals
371 ClinVar – public database with information on the relationship between medically important
372 variants and phenotypes.

# Declarations

375
## Ethics approval and consent to participate
Not applicable

378
## Consent for publication
The authors declare that they have no competing interests.

381
## Availability of data and material
The datasets generated and/or analyzed during the current study are available in the GitHub repository - vandhanak/stanford-benchmarking-workflows.

385
## Competing interests
Not applicable

388
## Funding
This work was funded by Stanford HealthCare, Stanford Children's Health and Stanford School of Medicine.

392
## Authors' contributions
VK designed and implemented the benchmarking workflow. SU and VK wrote the manuscript. ZN implemented the scripts to generate the performance assay report including the clinical variant validation. SU, SD, MP and EA conceived, designed and supervised the overall study. All authors read and approved the final manuscript.

398
## Acknowledgements
We thank Amin Zia for providing useful information during the initial phase of the benchmarking work. We thank Nathan Hammond and Issac Liao for the development of the in-house workflow engine "Loom" which was used to run the variant calling pipelines and the subsequent benchmarking workflows.
We are thankful to Chittaranjan Muthumalai for leading the automated benchmarking pipeline testing efforts and Jason Merker for useful discussions in terms of clinical relevance during the benchmarking process.

407
This study makes use of data generated by the DECIPHER community. A full list of centers who contributed to the generation of the data is available from http://decipher.sanger.ac.uk and via email from decipher@sanger.ac.uk. Funding for the project was provided by the Wellcome Trust.

412
413
# References

415
1. Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, et al. Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls. Nat Biotechnol. 2014;32(3):246-51.
2. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, et al. Extensive sequencing of seven human genomes to characterize benchmark reference materials. Sci Data. 2016;3:160025.

422   3.      Krusche P, Trigg L, Boutros PC, Mason CE, De La Vega FM, Moore BL, et al. Best
423   practices for benchmarking germline small-variant calls in human genomes. Nat Biotechnol.
424   2019;37(5):555-60.
425   4.      Jennings L, Van Deerlin VM, Gulley ML, Committee CoAPMPR. Recommended
426   principles and practices for validating clinical molecular pathology tests. Arch Pathol Lab
427   Med. 2009;133(5):743-55.
428   5.      N A, Q Z, L B, DK D, B F, JS G, et al. College of American Pathologists laboratory
429   standards for next-generation sequencing clinical tests. Arch Pathol Lab Med.
430   2015;139(4):481-93.
431   6.      MS W, GR C, RJ D, DA D, K K, M M, et al. Cystic fibrosis population carrier
432   screening: 2004 revision of the American College of Medical Genetics
433   mutation panel. Genetics in Medicine. 2004;6(5):387–91.
434   7.      Chiang C, Layer RM, Faust GG, Lindberg MR, Rose DB, Garrison EP, et al.
435   SpeedSeq: ultra-fast personal genome analysis and interpretation. Nat Methods.
436   2015;12(10):966-8.
437   8.      Garrison E, Marth G. Haplotype-based variant detection from short-read sequencing.
438   arXiv:1207.3907v2 [q-bio.GN]; 20 Jul 2012.
439   9.      Jenkins; Available from: https://jenkins.io/doc/.
440   10.     CircleCI; Available from: https://circleci.com/docs/.
441   11.     YAML; Available from: https://yaml.org/Available from:
442   https://www.tutorialspoint.com/yaml/index.htm.
443   12.     OpenWDL: Broad Institute;  [Available from: https://software.broadinstitute.org/wdl/.
444   13.     GATK workflows: Broad Institute;  [Available from: https://github.com/gatk-
445   workflows/.
446   14.     Amstutz P, Crusoe MR, Tijanić  N, Chapman B, Chilton J, Heuer M, et al. Common
447   Workflow Language, v1.0. Specification, *Common Workflow Language working group*. In:
448   Peter Amstutz MRC, Nebojša Tijanić, editor. 2016.
449   15.     Common Workflow Language (CWL). Software Freedom Conservancy, Inc.
450   16.     O'Connor BD, Yuen D, Chung V, Duncan AG, Liu XK, Patricia J, et al. The
451   Dockstore: enabling modular, community-focused sharing of Docker-based genomics tools
452   and workflows. F1000Res. 2017;6:52.
453   17.     Variant Call Format; Available from:
454   http://www.internationalgenome.org/wiki/Analysis/variant-call-format.
455   18.     BED format; Available from: http://genome.ucsc.edu/FAQ/FAQformat#format1.
456   19.     Krusche P, Trigg L, Boutros PC, Mason CE, De La Vega FM, Moore BL, et al.
457   Author Correction: Best practices for benchmarking germline small-variant calls in human
458   genomes. Nat Biotechnol. 2019;37(5):567.
459   20.     Krusche P. Haplotype VCF Comparison Tools; Available from:
460   https://github.com/Illumina/hap.pyAvailable from:
461   https://github.com/Illumina/hap.py/blob/master/doc/happy.md.
462   21.     BCFtools; Available from: http://samtools.github.io/bcftools/.
463   22.     Hammond N. Loom: platform-independent tool to create, execute, track, and share
464   workflows. 2017.
465   23.     Vivian J. Toil; Available from: https://toil.readthedocs.io/en/latest/.
466   24.     Vivian J, Rao AA, Nothaft FA, Ketchum C, Armstrong J, Novak A, et al. Toil enables
467   reproducible, open source, big biomedical data analyses. Nat Biotechnol. 2017;35(4):314-6.
468   25.     Cromwell: Broad Institute; Available from:
469   https://cromwell.readthedocs.io/en/stable/.
470   26.     CDC. GeT-RM Home: Available from:
471   https://wwwn.cdc.gov/clia/Resources/GETRM/default.aspx.

27.    GeT-RM NA04408: Coriell Institute of Medical Research; Available from: https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=NA04408&Product=DNA.

28.    GeT-RM NA14090: Coriell Institute of Medical Research; Available from: https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=NA14090&Product=DNA.

29.    Get-RM NA14170: Coriell Institute of Medical Research; Available from: https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=NA14170&Product=DNA.

30.    GeT-RM NA16658: Coriell Institute of Medical Research; Available from: https://www.coriell.org/0/Sections/Search/Sample_Detail.aspx?Ref=NA16658&Product=DNA.

31.    Fujita PA, Rhead B, Zweig AS, Hinrichs AS, Karolchik D, Cline MS, et al. The UCSC Genome Browser database: update 2011. Nucleic Acids Res. 2011;39(Database issue):D876-82.

32.    Karolchik D, Hinrichs AS, Kent WJ. The UCSC Genome Browser. Curr Protoc Bioinformatics. 2009;Chapter 1:Unit1.4.

33.    Online Mendelian Inheritance in Man, OMIM®. 2017 ed: McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University (Baltimore, MD).

34.    Landrum MJ, Lee JM, Benson M, Brown G, Chao C, Chitipiralla S, et al. ClinVar: public archive of interpretations of clinically relevant variants. Nucleic Acids Res. 2016;44(D1):D862-8.

35.    Firth HV, Richards SM, Bevan AP, Clayton S, Corpas M, Rajan D, et al. DECIPHER: Database of Chromosomal Imbalance and Phenotype in Humans Using Ensembl Resources. Am J Hum Genet. 2009;84(4):524-33.

36.    Administration USFaD. Guidance for industry and FDA staff: statistical guidance on reporting results from studies evaluating diagnostic tests. .

# Supplementary Materials

## *Preparation of truth sets for exome regions*

The NIST version and the ftp site used to download the original data for each of the GIAB samples (before preprocessing) used in this study are listed here.

NA12878
NIST v3.3:
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/NISTv3.3/NA12878_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-Solid-10X_CHROM1-X_v3.3_highconf.bed
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/NISTv3.3/NA12878_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-Solid-10X_CHROM1-X_v3.3_highconf.vcf.gz

NA24143
NIST v3.3:
ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG004_NA24143_mother/NISTv3.

520  3/HG004_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-10X_CHROM1-
521  22_v3.3_highconf.bed
522  ftp://ftp-
523  trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG004_NA24143_mother/NISTv3.
524  3/HG004_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-10X_CHROM1-
525  22_v3.3_highconf.vcf.gz
526
527  NA24149
528  NIST v3.3:
529  ftp://ftp-
530  trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG003_NA24149_father/NISTv3.3/
531  HG003_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-10X_CHROM1-22_v3.3_highconf.bed
532  ftp://ftp-
533  trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG003_NA24149_father/NISTv3.3/
534  HG003_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-10X_CHROM1-
535  22_v3.3_highconf.vcf.gz
536
537  NA24385
538  NIST v3.3:
539  ftp://ftp-
540  trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002_NA24385_son/NISTv3.3/H
541  G002_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-Solid-10X_CHROM1-
542  22_v3.3_highconf.bed
543  ftp://ftp-
544  trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002_NA24385_son/NISTv3.3/H
545  G002_GIAB_highconf_CG-IllFB-IllGATKHC-Ion-Solid-10X_CHROM1-
546  22_v3.3_highconf.vcf.gz
547
548  NA24631
549  NIST v3.3.2:
550  ftp://ftp-
551  trace.ncbi.nlm.nih.gov/giab/ftp/release/ChineseTrio/HG005_NA24631_son/NISTv3.3.2/GRC
552  h37/HG005_GRCh37_highconf_CG-IllFB-IllGATKHC-Ion-SOLID_CHROM1-
553  22_v.3.3.2_highconf_noMetaSV.bed
554  ftp://ftp-
555  trace.ncbi.nlm.nih.gov/giab/ftp/release/ChineseTrio/HG005_NA24631_son/NISTv3.3.2/GRC
556  h37/HG005_GRCh37_highconf_CG-IllFB-IllGATKHC-Ion-SOLID_CHROM1-
557  22_v.3.3.2_highconf.vcf.gz
558
559  ***Bash command to compute total number of bases in a region of interest***
560  awk '{a=$3-$2;print a}' <Consolidated.bed> | paste -sd+ - | bc
561
562  In the above command, <Consolidated.bed> refers to GIAB original high confidence bed file
563  for a sample intersected with the bed file of the region of interest such as coding exons,
564  whole exome or clinically relevant gene regions. The user can use this command to calculate
565  bases with their desired stratified region in the bed format which is required to compute
566  metrics such as true negatives.
567
568  ***Output files generated by Benchmarking workflow***

569 Our benchmarking workflow generates the following output files:
570 `1. <Output file common prefix>_<Sample ID>_CodingExons.vcf.gz`
571 `2. <Output file common prefix>_<Sample ID>_CodingExons.vcf.gz.tbi`
572 `3. <Output file common prefix>_<Sample ID>_CodingExons_counts.csv`
573 `4. <Output file common prefix>_<Sample ID>_CodingExons_counts.json`
574 `5. <Output file common prefix>_<Sample ID>_CodingExons_summary.csv`
575 `6. <Output file common prefix>_<Sample ID>_CodingExons_extended.csv`
576 `7. <Output file common prefix>_<Sample ID>_CodingExons_metrics.json`
577 `8. <Output file common prefix>_<Sample ID>_CodingExons_ConsoleOutput.txt`
578 `9. <Output file common prefix>_<Sample`
579 `ID>_CodingExons_indelSizeDistribution.txt`
580 `10. <Output file common prefix>_<Sample`
581 `ID>_CodingExons_indelSizeDistributionOnPlot.pdf`
582
583 There is a final performance assay report generated in the form of a tab delimited file as
584 below:
585 `Final_benchmarking_metrics_<current_date>.txt`
586
587
588 Another set of 10 files as seen above corresponding to the whole exome regions are
589 generated.
590
591 The benchmarking framework generates the following intermediate files:
592 `1. <Output file common prefix>_<Sample ID>_CodingExons_SNPs_TPonly.vcf.gz`
593 `2. <Output file common prefix>_<Sample ID>_CodingExons_SNPs_FPonly.vcf.gz`
594 `3. <Output file common prefix>_<Sample ID>_CodingExons_SNPs_FNonly.vcf.gz`
595 `4. <Output file common prefix>_<Sample ID>_CodingExons_INDELs_TPonly.vcf.gz`
596 `5. <Output file common prefix>_<Sample ID>_CodingExons_INDELs_FPonly.vcf.gz`
597 `6. <Output file common prefix>_<Sample ID>_CodingExons_INDELs_FNonly.vcf.gz`
598 `7. <Output file common prefix>_<SampleID>_CodingExons_indelDistribution.txt`
599
600 Another set of seven files as seen above corresponding to the whole exome regions are
601 generated.
602
603 **Supplemental Tables**
604
605 Table S1. Benchmarking metrics for indels of different size ranges in NA12878 (truth set
606 NIST v3.3, total bases = 13728555) for the regions within ~7000 clinically relevant genes (as
607 specified in Methods).
608

| Size of indels in NA12878 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 145 | 139 | 10 | 6 | 13728400 | 100 | 93.29 | 95.86 |
| 11-20 | 7 | 7 | 0 | 0 | 13728548 | 100 | 100 | 100 |
| 21-50 | 5 | 5 | 0 | 0 | 13728550 | 100 | 100 | 100 |
| All Indels | 156 | 150 | 10 | 6 | 13728389 | 100 | 93.75 | 96.15 |

609
610
611 Table S2. Benchmarking metrics for indels of different size ranges in NA24143 (truth set
612 NIST v3.3, total bases = 12549224) for the regions within ~7000 clinically relevant genes (as
613 specified in Methods).
614

| Size of indels in NA24143 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 153 | 143 | 16 | 10 | 12549055 | 100 | 89.94 | 93.46 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 11-20 | 8 | 8 | 0 | 0 | 12549216 | 100 | 100 | 100 |
| 21-50 | 3 | 3 | 0 | 0 | 12549221 | 100 | 100 | 100 |
| All Indels | 163 | 153 | 16 | 10 | 12549045 | 100 | 90.53 | 93.87 |

615

616

617 Table S3. Benchmarking metrics for indels of different size ranges in NA24149 (truth set
618 NIST v3.3, total bases = 12538042) for the regions within ~7000 clinically relevant genes (as
619 specified in Methods).

620

| Size of indels in NA24149 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 156 | 153 | 8 | 3 | 12537878 | 100 | 95.03 | 98.08 |
| 11-20 | 8 | 8 | 1 | 0 | 12538033 | 100 | 88.89 | 100 |
| 21-50 | 1 | 1 | 0 | 0 | 12538041 | 100 | 100 | 100 |
| All Indels | 163 | 161 | 9 | 3 | 12537869 | 100 | 94.71 | 98.16 |

621

622

623 Table S4. Benchmarking metrics for indels of different size ranges in NA24631 (truth set
624 NIST v3.3, total bases = 12808688) for the regions within ~7000 clinically relevant genes (as
625 specified in Methods).

626

| Size of indels in NA24631 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 153 | 146 | 16 | 7 | 12808519 | 100 | 90.12 | 95.42 |
| 11-20 | 5 | 5 | 0 | 0 | 12808683 | 100 | 100 | 100 |
| 21-50 | 5 | 4 | 0 | 1 | 12808683 | 100 | 100 | 80 |
| All Indels | 162 | 154 | 16 | 8 | 12808510 | 100 | 90.59 | 95.06 |

627

628

629 Table S5. Benchmarking metrics on the number of indels of different size ranges in NA12878
630 (truth set NIST v3.3, total bases = 71152019) for the whole exome region s including non-
631 coding exons, splice sites (+/- 20 bp) and clinically relevant deep intronic regions.

632

| Size of indels in NA12878 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 5108 | 4704 | 781 | 404 | 71146130 | 100 | 85.76 | 92.09 |
| 11-20 | 209 | 194 | 13 | 15 | 71151797 | 100 | 93.72 | 92.82 |
| 21-50 | 52 | 47 | 5 | 5 | 71151962 | 100 | 90.38 | 90.38 |
| All Indels | 5318 | 4910 | 800 | 424 | 71145885 | 100 | 85.99 | 92.03 |

633

634

635 Table S6. Benchmarking metrics on the number of indels of different size ranges in NA24143
636 (truth set NIST v3.3, total bases = 65657646) for the whole exome regions including non-
637 coding exons, splice sites (+/- 20 bp) and clinically relevant deep intronic regions.

638

| Size of indels in NA24143 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 5168 | 4676 | 681 | 492 | 65651797 | 100 | 87.29 | 90.48 |
| 11-20 | 206 | 184 | 13 | 22 | 65657427 | 100 | 93.40 | 89.32 |
| 21-50 | 84 | 72 | 5 | 12 | 65657557 | 100 | 93.51 | 85.71 |
| All Indels | 5388 | 4878 | 700 | 526 | 65651542 | 100 | 87.45 | 90.24 |

639

640

641 Table S7. Benchmarking metrics on the number of indels of different size ranges in NA24149
642 (truth set NIST v3.3, total bases = 65597266) for the whole exome regions including non-
643 coding exons, splice sites (+/- 20 bp) and clinically relevant deep intronic regions.
644

| Size of indels in NA24149 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 5096 | 4578 | 628 | 518 | 65591542 | 100 | 87.94 | 89.84 |
| 11-20 | 188 | 167 | 17 | 21 | 65597061 | 100 | 90.76 | 88.83 |
| 21-50 | 68 | 62 | 5 | 6 | 65597193 | 100 | 92.54 | 91.18 |
| All Indels | 5290 | 4763 | 651 | 545 | 65591307 | 100 | 87.98 | 89.70 |

645

646 Table S8. Benchmarking metrics on the number of indels of different size ranges in NA24631
647 (truth set NIST v3.3, total bases = 65657646) for the whole exome regions including non-
648 coding exons, splice sites (+/- 20 bp) and clinically relevant deep intronic regions.
649

| Size of indels in NA24631 | Truth total | TP | FP | FN | TN | NPA | Precision | Recall |
|---|---|---|---|---|---|---|---|---|
| 1-10 | 5555 | 5089 | 656 | 466 | 66982776 | 100 | 88.58 | 91.61 |
| 11-20 | 187 | 178 | 6 | 9 | 66988794 | 100 | 96.74 | 95.19 |
| 21-50 | 82 | 68 | 8 | 14 | 66988897 | 100 | 89.47 | 82.93 |
| All Indels | 5805 | 5316 | 671 | 489 | 66982511 | 100 | 88.79 | 91.58 |

650
651
652 **Supplemental Data files for:**
653 **Benchmarking workflows to assess performance and suitability of germline variant**
654 **calling pipelines in clinical diagnostic assays**
655
656 The benchmarking workflow file and relevant scripts (listed below as additional files) and all
657 output files for five GIAB samples per stage are available in our public repository:
658 vandhanak/stanford-benchmarking-workflows
659
660
661 **Supplemental Files**
662
663    1.  indelSizeDistribution_Detailed.R
664    2.  benchmarking_truth_set.py
665    3.  verify_variants.py