



Analytical validation of germline small variant detection using long-read HiFi genome sequencing

Nathan Hammond, Linda Liao, Pun Wai Tong, et al.

Genome Res. published online April 11, 2025

Access the most recent version at doi:[10.1101/gr.278836.123](https://doi.org/10.1101/gr.278836.123)

P<P	Published online April 11, 2025 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

An advertisement banner with a teal background. On the left, the text reads "CRISPR and RNAi Genetic Screening. Your new superpower." In the center, there is a white box with the words "LEARN MORE". On the right, there is a photograph of a woman wearing a red superhero mask and cape, and the Cellecta logo, which consists of a cluster of green dots and the word "CELLECTA" below it.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **Analytical validation of germline small variant detection using long-read HiFi**
2 **genome sequencing**

3
4
5 Nathan Hammond^{1*}, Linda Liao¹, Pun Wai Tong¹, Zena Ng¹, Thuy-Mi P. Nguyen², Chandler
6 Ho¹, Yao Yang^{1,2}, Stuart A. Scott^{1,2}

7
8
9 ¹ Clinical Genomics Laboratory, Stanford Medicine, Palo Alto, CA 94304.

10 ² Department of Pathology, Stanford University, Stanford, CA 94305.

11
12 Current affiliation:

13 * Influx Bio, San Francisco, CA 94124.

14
15 **Running Title:** Long-read HiFi genome sequencing validation

16
17 **Manuscript Category:** Regular Article

18 **Abstract:** 234

19 **Tables:** 3

20 **Figures:** 3

21
22 **Key words:** Clinical genome sequencing; long-read sequencing; SMRT sequencing; HiFi
23 sequencing; germline genetic testing; diagnostic testing; medical genomics; LRS Special Issue.

24
25 **Sources of support:** S.A.S was supported in part by NIH/NHGRI grant U01HG011762.

26
27 **Disclosures:** N.H. is currently an employee of Influx Bio; all other authors declare no conflicts
28 of interest.

29
30
31
32
33
34
35
36
37
38
39 **CORRESPONDENCE TO:**

40 Stuart A. Scott, PhD

41 Professor,

42 Department of Pathology

43 Stanford University

44 Palo Alto, CA 94305

45 Telephone: 650-724-0973

46 Email: sascott@stanford.edu

47 **ABSTRACT**

48 Long-read sequencing has the capacity to interrogate difficult genomic regions and phase
49 variants; however, short-read sequencing is more commonly implemented for clinical testing.
50 Given the advances in long-read HiFi sequencing chemistry and variant calling, we analytically
51 validated this technology for small variant detection (single nucleotide variants,
52 insertions/deletions; SNVs/indels; <50bp). HiFi genome sequencing was performed on DNA
53 from reference materials and clinical specimen types, and accuracy results were compared to
54 short-read genome sequencing data. HiFi genome sequencing recall and precision across
55 Genome in a Bottle (GIAB)-defined non-difficult and difficult genomic regions (high
56 confidence) for SNVs were >99.9% and >99.7%, respectively, and for indels were >99.8% and
57 >99.1%, respectively. Moreover, HiFi genome sequencing outperformed short-read genome
58 sequencing on overall SNV/indel F1-score accuracy at all paired sequencing depths, which were
59 further stratified across 100 total GIAB-defined genomic regions for a comprehensive evaluation
60 of performance. Of note, HiFi genome sequencing F1-scores for SNVs and indels surpassed 99%
61 at ~15× and ~25×, respectively. In addition, high confidence small variant concordance across all
62 HiFi genome sequencing reproducibility assessments (two specimens, three independent
63 sequencing datasets) were >99.8% for SNVs and >98.6% for indels, and average high
64 confidence small variant concordance between paired blood, saliva, and swab specimens were all
65 >99.8%. Taken together, these data underscore that long-read HiFi genome sequencing detection
66 of SNVs and indels is very accurate and robust, which supports the implementation of this
67 technology for clinical diagnostic testing.

68

69

70 INTRODUCTION

71 Genome sequencing has evolved from a widely used research platform to a comprehensive
72 clinical test at selected medical centers and laboratories (Belkadi et al. 2015; Costain et al. 2020),
73 with the capacity to sequentially interrogate regions of the genome based on new evidence and/or
74 clinical indication (Rehm 2017; Costain et al. 2018; Bick et al. 2019; Yang et al. 2024). Short-
75 read sequencing is the most commonly implemented platform for genome sequencing; however,
76 long-read sequencing is rapidly emerging as an alternative platform with notable benefits over
77 short-read sequencing (Logsdon et al. 2020; Cohen et al. 2022; Conlin et al. 2022). For example,
78 long-read genome sequencing has improved interrogation of clinically significant regions,
79 including structural variants, repeat expansions, homologous gene families and the *HLA* region,
80 as well as the inherent benefit of variant phasing (Ardui et al. 2018; Ameur et al. 2019). The two
81 primary long-read sequencing chemistries currently available are single molecule real-time
82 (SMRT; HiFi) and nanopore sequencing, which recently have been employed by the Telomere-
83 to-Telomere (T2T) Consortium to more comprehensively characterize the CHM13 human
84 reference assembly (Nurk et al. 2022).

85
86 Long-read SMRT sequencing has been shown to generate highly accurate high-fidelity (HiFi)
87 read lengths of ~10-25 kb using the Sequel II platform (Pacific Biosciences) (Wenger et al. 2019;
88 Hon et al. 2020). More recently, long-read HiFi genome sequencing has been used to expand the
89 small variant benchmarks in the commonly leveraged Genome in a Bottle Consortium (GIAB)
90 reference material samples to include difficult-to-map regions and segmental duplications that
91 are inherently challenging for short reads (Wagner et al. 2022a). In addition, comparisons of
92 sequencing platforms and variant calling strategies have recently been reported by the

93 PrecisionFDA Truth Challenge V2, which found long-read HiFi genome sequencing to
94 outperform both short-read and long-read nanopore sequencing with genome-wide variant
95 calling accuracy (Olson et al. 2022).

96

97 Although long-read HiFi genome sequencing has improved accuracy and haploblock phasing
98 performance compared to short-read sequencing, its adoption into clinical genetic testing
99 laboratories is only now emerging. Of note, resources for validating new clinical sequencing
100 assays are available from the College of American Pathologists (CAP), Association for
101 Molecular Pathology (AMP) (Aziz et al. 2015; Roy et al. 2018), and related professional
102 consortia (Gargis et al. 2012; Matthijs et al. 2016; Santani et al. 2017; Santani et al. 2019); as
103 well as benchmarking reference materials from the GIAB/National Institute of Standards and
104 Technology (NIST) consortium,(Zook et al. 2019) and the Global Alliance for Genomics and
105 Health Benchmarking (GA4GH) Team (Krusche et al. 2019). Therefore, to facilitate the
106 implementation of clinical long-read HiFi genome sequencing, our initial effort was centered on
107 a robust analytical validation of germline small variant (SNV/indel) detection using current best
108 practices and benchmarking resources.

109

110 **RESULTS**

111 **Long-Read HiFi Genome Sequencing Small Variant Accuracy**

112 HiFi genome sequencing accuracy (SNV/indel; <50bp) was evaluated by sequencing seven
113 GIAB/NIST reference material samples (average 31.1×). Benchmarking was performed using
114 tools and practices recommended by the GA4GH (Krusche et al. 2019), and hap.py version
115 v0.3.15 was used to compare observed results with the published truth set version v4.2.1.

116 Accuracy was measured by recall (i.e., sensitivity) and precision (i.e., positive predictive value),
117 which were stratified by high confidence genomic regions as defined by GIAB/NIST (Zook et al.
118 2019). SNV/indel detection across the GIAB reference samples was highly accurate, as the
119 average recall and precision for all seven samples were >99.9% for SNVs and >99.8% for indels
120 across the non-difficult genomic regions (**Table 1**). As expected, recall and precision were
121 slightly lower when interrogating genomic regions with known sequencing challenges (low
122 complexity, low mappability, segmental duplications); however, average recall and precision
123 across all difficult genomic regions were still >98.9% and >98.7% for SNVs and indels,
124 respectively. Of note, in addition to the low complexity, low mappability, and segmental
125 duplication regions highlighted in **Table 1**, the GIAB ‘Difficult regions’ are defined broadly to
126 include low or high GC content (<25% or >65%), bad promoter regions, false duplications, and
127 other difficult genomic regions (Krusche et al. 2019). The average recall and precision across all
128 genomic regions were >99.7% and >99.1% for SNVs for indels, respectively (**Table 1**). Among
129 the discordant HiFi small variants with the GIAB truth set, homopolymers were the most
130 common source of error, with ~75% of indel errors located in or adjacent to a homopolymer run.

131

132 **Long-Read HiFi and Short-Read Genome Sequencing Small Variant Accuracy**

133 In addition to accuracy benchmarking across the high confidence GIAB/NIST regions (i.e., low
134 complexity, low mappability, segmental duplications, all difficult regions, not in any difficult
135 region, all), HiFi genome sequencing performance was further evaluated across 100 GIAB-
136 defined subregions of the human genome (Zook et al. 2019). F1-scores were generated and
137 stratified by variant type (i.e., SNVs/indels), and results compared to paired analyses with
138 publicly available short-read genome sequencing data (average 40.3×). As illustrated in **Figure**

139 **1A-F**, HiFi genome sequencing small variant F1-scores were superior to short-read genome
140 sequencing small variant F1-scores across 22 informative genomic subregions within the
141 categories of low mappability, homopolymers, tandem repeats and GC content, which were most
142 notable for indels of increasing size. In addition, HiFi genome sequencing variant calling
143 accuracy across the recently reported Challenging Medically Relevant Gene (CMRG) truth set
144 (Wagner et al. 2022b) was also interrogated and compared with short-read genome sequencing
145 accuracy, which identified HiFi's advantage over short-read in detection of >15bp deletions
146 (87.76% vs. 67.81% recall), >15bp insertions (84.44% vs. 74.81% recall) but otherwise showed
147 roughly similar performance in the CMRG regions (**Supplemental Table S1**).

148

149 **Genome Sequencing Depth Stratification and Accuracy**

150 To evaluate HiFi genome sequencing small variant accuracy at different sequencing depths,
151 small variant recall and precision were assessed across a series of sequencing depths
152 downsampled from NA24385 (HG002). As expected, HiFi genome-wide small variant accuracy
153 was reduced at low sequencing depths (**Supplemental Table S2**); however, HiFi genome
154 sequencing F1-scores for SNVs and indels surpassed 99% at ~15× and ~25×, respectfully. In
155 comparison, short-read genome sequencing F1-scores for SNVs and indels surpassed 99% at
156 ~20× and ~35×, respectfully (**Figure 2**). At sequencing depths of ~30×, HiFi genome sequencing
157 F1-scores for SNVs and indels were 99.8% and 99.1%, respectfully; and short-read sequencing
158 F1-scores for SNVs and indels were 99.8% and 98.7%, respectfully (**Figure 2A-C**).

159

160 **Long-Read HiFi Genome Sequencing Concordance and Reproducibility**

161 HiFi genome sequencing library preparation included both manual and automated workflows,
162 and minor updates were introduced to the laboratory procedure to optimize the automated
163 workflow. Workflow and procedure updates were validated by measuring small variant accuracy
164 and/or concordance with reference material and specimen samples as appropriate. As detailed in
165 **Supplemental Tables S3 and S4**, the accuracy and concordance of the workflow updates (e.g.,
166 fragment depletion using the SRE XS and SRE kits) were consistent with paired manual library
167 preparation results, which supported the implementation of these workflow improvements. In
168 addition, a novel validation strategy of two Miro Canvas instruments was accomplished by low-
169 depth NA12878 benchmarking comparisons using a single SMRTcell of data ($\sim 9\times$) (**Figure**
170 **3A-B**), which was supported by consistent quality metrics between manual and automated
171 workflows, and concordant benchmarking results from a subsequent full depth ($27.9\times$) Miro
172 Canvas library preparation (**Supplemental Tables S3 and S4**).

173
174 HiFi genome sequencing reproducibility was evaluated by comparing reference sample results to
175 two independent publicly available datasets (NA12878/HG001, NA24385/HG002; see
176 **Materials and Methods**) and measuring F1-score concordance across all three datasets (average
177 $32.9\times$). Genome-wide SNV/indel non-reference genotype concordance was stratified by genomic
178 context, which ranged from ~ 98 - 99.9% (**Table 2**); however, HiFi genome sequencing
179 reproducibility was reduced when assessed across regions not considered high confidence by
180 GIAB (**Supplemental Table S5**). Non-reference genotype concordance for the high confidence
181 RefSeq CDS regions across all reproducibility and repeatability assessments were $>99.8\%$ and
182 $>99.3\%$ for SNVs and indels, respectively (**Table 2**), indicating that HiFi genome sequencing
183 small variant detection is robust and precise.

184

185 **Long-Read HiFi Genome Sequencing Specimen Validation**

186 Germline specimens were validated by subjecting paired blood, saliva, and swab samples to HiFi
187 genome sequencing and evaluating SNV/indel concordance. As expected, concordance between
188 specimens was reduced for SNVs/indels when evaluating difficult genomic regions as saliva-
189 based specimens are known to harbor bacterial DNA that interferes with sequencing (Troost et al.
190 2019; Yao et al. 2020). However, the average paired SNV/indel concordance between all three
191 specimen types were >99% across all high confidence genomic regions (**Table 3**). Although
192 specimen concordance was reduced when assessed across regions not considered high
193 confidence by GIAB (**Supplemental Table S6**), these validation results indicate that HiFi
194 genome sequencing of saliva and swab specimens are consistent with blood for germline
195 SNV/indel detection.

196

197 **DISCUSSION**

198 To facilitate the implementation of diagnostic long-read HiFi sequencing, we executed an
199 analytical validation plan that was centered on comprehensively evaluating HiFi genome
200 sequencing for germline SNV/indel detection and specimen types that are used for clinical
201 testing in medical genetics. Results were stratified by variant type and GIAB-defined genomic
202 regions to better inform overall performance, which ultimately determined that HiFi genome
203 sequencing is accurate and robust. The accuracy of germline small variant detection in non-
204 difficult genomic regions across reference materials was >99.9% for both SNVs and indels, and
205 small variant detection accuracy in GIAB-defined difficult regions was >99.5% and >98.8% for
206 SNVs and indels, respectively. These analytical validation analyses underscore the accuracy of

207 long-read HiFi genome sequencing for detecting germline SNV/indels (<50 bp), which supports
208 the implementation of this technology for clinical genetic testing. In addition, quality control
209 (QC) thresholds for clinical long-read HiFi genome sequencing based on CAP requirement
210 MOL.36151 are suggested in **Supplemental Table S7**; however, these should be considered
211 preliminary recommendations, as clinical laboratories should leverage their own experience and
212 data to define internal QC metrics.

213

214 Analytical validation is a critical assessment of any new clinical laboratory test, which is defined
215 by CAP Checklists and other state, federal, and/or professional requirements/recommendations.
216 Test performance specifications include reportable range, accuracy, reproducibility/repeatability,
217 sensitivity/specificity, and other relevant performance characteristics. For long-read HiFi
218 genome sequencing we adopted the definitions for ‘Reportable Range’ and ‘Reference Range
219 (Reference Interval)’ based on clinical high-throughput sequencing guidelines (Gargis et al.
220 2012; Santani et al. 2017). However, for a more comprehensive assessment of sequencing
221 performance, reportable range was measured genome-wide but strategically stratified by distinct
222 genomic regions as defined by GIAB/GA4GH. The genomic regions implemented in this
223 validation included high level strata (low complexity, low mappability, segmental duplications,
224 all difficult regions, not in any difficult regions, all high confidence regions, RefSeq CDS
225 regions), as well as the more specific genomic subregions defined by GIAB/NIST and GA4GH
226 (Krusche et al. 2019). These regions were intersected with our SNV/indel performance results, as
227 well as genome sequencing specimen validation data (CAP requirement MOL.31015), as
228 deemed appropriate based on analysis context and intended use.

229

230 Sequencing accuracy is a rapidly evolving area that is driven by continual improvements in
231 available chemistries and informatic algorithms developed for calling germline variants. As an
232 integral component of validating clinical sequencing-based platforms (Roy et al. 2018),
233 benchmarking small variant accuracy (i.e., recall, precision, F1) is supported through the
234 GIAB/NIST/GA4GH resources (Majidian et al. 2023; Olson et al. 2023), which recently has
235 been catalyzed by PrecisionFDA challenges that more comprehensive evaluations of sequencing-
236 based variant calling (Zook et al. 2019; Olson et al. 2022). Our analytical validation of HiFi
237 genome sequencing is consistent with the most recent PrecisionFDA V2 challenge, which
238 concluded that long-read HiFi sequencing coupled with machine learning-based variant calling
239 tools (Pei et al. 2021; Olson et al. 2022) was superior to short-read genome sequencing using
240 graph-based variant calling.

241
242 It is important to note that our reported accuracy results reflect not only the sequencing platforms
243 evaluated but also the variant calling methods used. Given that it was beyond the scope of our
244 study to perform a full comparison of bioinformatics techniques, we selected best-practice tools
245 with high accuracy in challenging genomic regions for each sequencing platform, as
246 demonstrated by PrecisionFDA V2. Our validation also included a detailed evaluation of
247 performance across GIAB-defined genomic stratifications, which highlighted long-read HiFi
248 sequencing accuracy across challenging regions and particularly among indel variants. Of note,
249 our targeted sequencing depth was $\geq 30\times$, consistent with recommendations from the Medical
250 Genome Initiative (Marshall et al. 2020), and as expected small variant accuracy was reduced at
251 lower depths. However, it is notable that 99% accuracy was surpassed at $\sim 15\text{-}25\times$ for long-read

252 HiFi genome sequencing compared to ~20-35× for short-read genome sequencing.

253

254 In addition to accuracy, HiFi genome sequencing small variant reproducibility was also
255 interrogated by measuring non-reference genotype concordance between datasets. Concordance
256 across all replicates in high confidence GIAB-defined regions ranged from 99.84-99.91% for
257 SNVs and 97.66-99.30% for indels, indicating that small variant calling is very robust. Given
258 that exome reproducibility/repeatability is typically higher than that observed with genome
259 sequencing due to the more narrow region interrogated (Linderman et al. 2014), we also
260 stratified our genome results by RefSeq CDS regions, which resulted in highly concordant small
261 variant calling across all replicates in the high confidence regions (all: 99.85%/99.37%; non-
262 difficult: 99.91%/99.73%, for SNVs/indels) and non-high confidence regions (all:
263 98.22%/93.13%; non-difficult: 99.66%/99.19%, for SNVs/indels).

264

265 Of note, the GIAB high confidence regions encompass 81.6% of the autosomal GRCh38 human
266 genome, which translates to ~2.52 Gb across the seven GIAB reference samples. The remaining
267 18.4% of non-high confidence autosomal bases (~567.2 Mb) represent subregions of the genome
268 (and Chromosomes X and Y) that are difficult to benchmark given the uncertainty in underlying
269 truth set (Zook et al. 2016; Zook et al. 2019; Zook et al. 2020). The concordance results across
270 reference materials and specimen types in our validation study were reduced in the genome-wide
271 analyses (i.e., including non-high confidence regions) compared to the concordance results
272 limited to the high confidence genomic regions, most notably in the GIAB-defined difficult
273 regions (low complexity, low mappability, segmental duplications, etc.). These metrics were
274 considered acceptable, as 40% of the variants in these regions had genotype quality scores of

275 <Q20 compared to less than 1% of variants in the non-difficult regions before filtering, and the
276 increase in variant numbers was much greater in the difficult regions than the non-difficult
277 regions in the non-high confidence regions ($2.51\times$ vs $1.08\times$). As such, these thorough
278 reproducibility analyses together indicate that long-read HiFi genome sequencing is highly
279 robust across the high confidence regions of the human genome; however, variants identified in
280 the GIAB non-high confidence difficult regions in a clinical setting would likely require
281 independent confirmation if reportable.

282

283 Another critical CAP requirement for test implementation is validating the specific specimen
284 types used for clinical processing (MOL.31015), which for germline genetic testing typically
285 includes peripheral blood and/or saliva specimens. To satisfy this requirement, paired specimens
286 were subjected to HiFi genome sequencing and concordance was measured across the genome.
287 Despite known challenges with using oral saliva samples for sequencing due to the presence of
288 competing bacterial DNA (Krusche et al. 2019; Trost et al. 2019), concordance between paired
289 blood, saliva, and assisted saliva (swab) specimens ranged from 99.82-99.92% for all
290 SNVs/indels across high confidence genomic regions. However, it is notable that ~99% of reads
291 aligned to the reference genome for blood and cell line specimens, compared to ~93% for
292 assisted saliva specimens and 86% for saliva, resulting in lower average depth (blood, cell lines:
293 $33\times$; assisted saliva: $28\times$; saliva: $25\times$). As such, additional sequencing of oral samples to
294 compensate for unmapped reads (as defined by QC thresholds) may be warranted in clinical
295 production. These specimen validation study results indicate that our HiFi genome sequencing
296 procedure and pipeline generates highly comparable results between peripheral blood and DNA

297 isolated from saliva and assisted saliva, which supports their use as acceptable clinical specimens
298 for this test.

299

300 Of note, copy number variant (CNV) detection by sequencing is routinely implemented among
301 clinical laboratories (Kadalayil et al. 2015; Rajagopalan et al. 2020), and GIAB/NIST has
302 developed consensus germline structural variant (SV) calls from HG002 (NA24385) (Zook et al.
303 2020). However, this dataset is an integration of 68 callsets from multiple algorithms and four
304 different sequencing technologies, each with their own strengths and weaknesses, and as a result
305 it does not currently include robust duplication calls or SVs >100 kb (Whitford et al. 2019).

306 Although long-read HiFi genome sequencing has been shown to be highly effective at CNV/SV
307 detection (Chaisson et al. 2019; Mahmoud et al. 2019; Aganezov et al. 2020), these variants were
308 considered out of scope for this initial analytical validation; however, they are currently being
309 evaluated for a subsequent analytical work product.

310

311 In conclusion, long-read HiFi genome sequencing ($\geq 30\times$) was analytically validated for germline
312 SNV/indel detection, which supports the implementation of this platform as a robust technology
313 for clinical genetic testing. Of note, practical factors for sequencing platform selection were
314 intentionally excluded from this analytical validation, including cost, labor and sequencing time,
315 as these variables were not applicable to analytical performance testing. This validation also did
316 not explicitly include ‘clinical performance characteristics’ as defined by CAP (MOL.31590), as
317 these analyses were reserved for subsequent validation of clinically significant germline variants.
318 As such, these analytical validation data provide the infrastructure for long-read HiFi genome

319 sequencing-based detection of germline variation, which supports the use of this innovative
320 technology for clinical diagnostic testing.

321

322 **MATERIALS and METHODS**

323 **Analytical Validation Specimens**

324 High molecular weight (HMW) reference material DNA samples were acquired from the Coriell
325 Institute for Medical Research (Camden, NJ), which included seven benchmarking samples from
326 the GIAB/NIST consortium. Peripheral blood was collected in EDTA vacutainer tubes using
327 standard practices and DNA isolated using the Maxwell RSC Buffy Coat DNA Kit (Promega
328 Corporation) according to manufacturer instructions. Saliva samples were collected using the
329 Oragene Dx OGD-500 kit (DNA Genotek, Ottawa, ON, Canada) or the assisted saliva (swab)
330 Oragene Dx OGD-575 kit (DNA Genotek). DNA was isolated from saliva specimens using
331 Maxwell RSC Stabilized Saliva DNA Kit (Promega Corporation) according to manufacturer
332 instructions. All validation samples and sequencing metrics are summarized in **Supplemental**
333 **Tables S7-S10.**

334

335 **Long-Read HiFi Genome Sequencing**

336 *Library Preparation and Long-Read HiFi Sequencing*

337 Genomic DNA was analyzed with the Femtopulse Genomic DNA 165 kb kit (Agilent, Santa
338 Clara, CA) to confirm adequate quantity of HMW DNA. Approximately 3-10 ug of DNA was
339 mechanically sheared to 10-20 kb using the Megaruptor Shearing kit (Diagenode, Denville, NJ),
340 with the DNAFluid+ kit (Diagenode) employed for viscous samples. Library preparation was
341 performed through either a manual or an automated workflow with the SMRTbell Prep Kit

342 according to manufacturer instructions, including end repair, A-tailing, adapter ligation,
343 purification with SMRTbell cleanup beads (Pacific Biosciences, Menlo Park, CA), and nuclease
344 treatment.

345
346 Manual library preparation included purification of sheared gDNA with SMRTbell cleanup
347 beads, SMRTbell library generation, followed by size selection on Blue Pippin (Sage Science,
348 Beverly, MA) to remove fragments <10 kb and purification with Ampure PB beads (Pacific
349 Biosciences). Automated library preparation included small fragment depletion using the Short
350 Read Eliminator (SRE) XS kit (<10 kb) or the SRE kit (<25 kb) (Pacific Biosciences) as needed
351 prior to shearing, followed by purification with SMRTbell cleanup beads and library preparation
352 using the Miro Canvas (Miroculus, San Francisco, CA). Manual and automated workflow
353 SMRTbell libraries were both quantified by the Qubit dsDNA assay kit (Invitrogen, Waltham,
354 MA) and bound to sequencing polymerase using the Sequel II Binding Kit (Pacific Biosciences).
355 Long-read HiFi genome sequencing was performed on the Sequel IIe system (Pacific
356 Biosciences) with a 30-hour movie collection time, and each sample was sequenced on three
357 SMRTcells except where otherwise noted.

358

359 ***Publicly Available Data***

360 To evaluate internal long-read HiFi genome sequencing reproducibility, selected publicly
361 available long-read HiFi genome sequencing data for the NA12878/HG001 and
362 NA24385/HG002 reference materials were acquired from the National Center for Biotechnology
363 Information (NCBI) FTP server (<ftp://ftp.ncbi.nlm.nih.gov/giab>): PacBio_SequelII_CCS_11kb

364 (NA12878), HudsonAlpha_PacBio_CCS (NA12878) (Zook et al. 2016),
365 PacBio_CCS_15kb_20kb_chemistry2 (NA24385), PacBio_SequelII_CCS_11kb (NA24385).

366

367 ***Long-Read HiFi Sequencing Bioinformatics Pipeline and Variant Calling***

368 HiFi reads were generated using SMRTLink 10.2 software and the Circular Consensus
369 Sequencing mode. For analyses requiring downsampled data, subsampling was executed with
370 SAMtools v1.18 (Danecek et al. 2021). Alignment and variant calling were performed using a
371 modified version of the PacBio HiFi-human-WGS-WDL pipeline
372 (<https://github.com/PacificBiosciences/HiFi-human-WGS-WDL>), with the following steps:
373 alignment of HiFi sequencing reads to Genome Reference Consortium Human Build 38
374 (GRCh38) using pbmm2 v1.7.0 (Hon et al. 2020); small variant calling using DeepVariant v1.4.0
375 and the PacBio machine learning model included with the software (Poplin et al. 2018); and
376 calculating aligned read depth with mosdepth v0.2.9 (Pedersen and Quinlan 2018). BCFtools
377 v1.20 was used to filter variants, removing all SNV calls with QUAL < 20 (Danecek et al. 2021).
378 Finally, Picard Tools v2.27.4 was used to calculate quality yield metrics
379 (<https://broadinstitute.github.io/picard>), alignment summary metrics, and variant calling metrics.
380 For analysis of CMRG genes, a modified version of the GRCh38 build was used, wherein false
381 duplications were masked and decoy contigs were added for falsely collapsed duplications
382 (Behera et al. 2023).

383

384 **Short-Read Genome Sequencing**

385 ***Publicly Available Data***

386 To compare long-read HiFi and short-read genome sequencing accuracy, publicly available
387 short-read genome sequencing data for the seven GIAB benchmarking reference materials were
388 acquired from the National Center for Biotechnology Information (NCBI) FTP server
389 (<ftp://ftp.ncbi.nlm.nih.gov/giab>): NIST_NA12878_HiSeq_300x,
390 NIST_Illumina_2x250bps (NA24385, NA24143, and NA24149),
391 HG005_NA24631_son_HiSeq_300x, NA24694_Father_HiSeq100x, and
392 NA24695_Mother_HiSeq100x.

393

394 ***Short-Read Sequencing Bioinformatics Pipeline and Variant Calling***

395 For analyses requiring downsampled data, subsampling was executed with SAMtools v1.18
396 (Danecek et al. 2021). Alignment, germline small variant calling, and calculation of quality
397 control metrics were performed with the Illumina DRAGEN Germline Pipeline v4.2.4 in
398 BaseSpace, using the HG38 alt-masked multi-genome graph reference.

399

400 **Analytical Validation Strategy**

401 The HiFi genome sequencing small variant analytical validation plan followed Laboratory
402 Developed Test (LDT) guidelines as defined by the CAP and AMP (Jennings et al. 2009; Aziz et
403 al. 2015; Roy et al. 2018), the American College of Medical Genetics and Genomics (ACMG)
404 (Rehm et al. 2013), high-throughput sequencing recommendations from professional consortia
405 (Gargis et al. 2012; Matthijs et al. 2016; Santani et al. 2017; Santani et al. 2019), and the Clinical
406 Laboratory Evaluation Program at the Wadsworth Center, New York State Department of
407 Health. The plan was centered on determining the analytical performance characteristics of HiFi
408 genome sequencing for use as a diagnostic technology, as well as defining standard operating

409 procedures (SOPs), quality control/quality assurance procedures, and validating small variant
410 detection and specimen types.

411

412 **DATA ACCESS**

413 All GIAB reference material and blood/saliva sample HiFi genome sequencing aligned bam
414 datasets that were generated in this study have been submitted to the NCBI BioProject database
415 (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession number PRJNA1143955.

416

417 **COMPETING INTEREST STATEMENT**

418 N.H. is currently an employee of Influx Bio; all other authors declare no conflicts of interest.

419

420 **ACKNOWLEDGEMENTS**

421 The authors would like to thank Stanford Health Care, Stanford Children’s Health, and Pacific
422 Biosciences (Menlo Park, CA) for their programmatic support. S.A.S was supported in part by
423 NIH/NHGRI grant U01HG011762. Project conceptualization: N.H., Y.Y., S.A.S.; Data
424 acquisition: N.H, L.L., P.W.T., Z.N.; Data analysis, interpretation, and management: N.H,
425 P.W.T., Z.N., C.H., T.P.N., Y.Y., S.A.S.; Drafting and revision: N.H., Y.Y., S.A.S.; Final
426 approval: N.H, L.L., P.W.T., Z.N., C.H., T.P.N., Y.Y., S.A.S.

427

428 **REFERENCES**

429

430 Aganezov S, Goodwin S, Sherman RM, Sedlazeck FJ, Arun G, Bhatia S, Lee I, Kirsche M, Wappel
431 R, Kramer M et al. 2020. Comprehensive analysis of structural variants in breast cancer
432 genomes using single-molecule sequencing. *Genome Res* **30**: 1258-1273.

433 Ameer A, Kloosterman WP, Hestand MS. 2019. Single-Molecule Sequencing: Towards Clinical
434 Applications. *Trends Biotechnol* **37**: 72-85.

- 435 Ardui S, Ameer A, Vermeesch JR, Hestand MS. 2018. Single molecule real-time (SMRT)
436 sequencing comes of age: applications and utilities for medical diagnostics. *Nucleic Acids*
437 *Res* **46**: 2159-2168.
- 438 Aziz N, Zhao Q, Bry L, Driscoll DK, Funke B, Gibson JS, Grody WW, Hegde MR, Hoeltge GA,
439 Leonard DG et al. 2015. College of American Pathologists' laboratory standards for next-
440 generation sequencing clinical tests. *Arch Pathol Lab Med* **139**: 481-493.
- 441 Behera S, LeFaive J, Orchard P, Mahmoud M, Paulin LF, Farek J, Soto DC, Parker SCJ, Smith
442 AV, Dennis MY et al. 2023. FixItFelix: improving genomic analysis by fixing reference
443 errors. *Genome Biol* **24**: 31.
- 444 Belkadi A, Bolze A, Itan Y, Cobat A, Vincent QB, Antipenko A, Shang L, Boisson B, Casanova
445 JL, Abel L. 2015. Whole-genome sequencing is more powerful than whole-exome
446 sequencing for detecting exome variants. *Proc Natl Acad Sci U S A* **112**: 5473-5478.
- 447 Bick D, Jones M, Taylor SL, Taft RJ, Belmont J. 2019. Case for genome sequencing in infants and
448 children with rare, undiagnosed or genetic diseases. *J Med Genet* **56**: 783-791.
- 449 Chaisson MJP, Sanders AD, Zhao X, Malhotra A, Porubsky D, Rausch T, Gardner EJ, Rodriguez
450 OL, Guo L, Collins RL et al. 2019. Multi-platform discovery of haplotype-resolved
451 structural variation in human genomes. *Nat Commun* **10**: 1784.
- 452 Cohen ASA, Farrow EG, Abdelmoity AT, Alaimo JT, Amudhavalli SM, Anderson JT, Bansal L,
453 Bartik L, Baybayan P, Belden B et al. 2022. Genomic answers for children: Dynamic
454 analyses of >1000 pediatric rare disease genomes. *Genet Med*
455 doi:10.1016/j.gim.2022.02.007.
- 456 Conlin LK, Aref-Eshghi E, McEldrew DA, Luo M, Rajagopalan R. 2022. Long-read sequencing
457 for molecular diagnostics in constitutional genetic disorders. *Hum Mutat* **43**: 1531-1544.
- 458 Costain G, Jobling R, Walker S, Reuter MS, Snell M, Bowdin S, Cohn RD, Dupuis L, Hewson S,
459 Mercimek-Andrews S et al. 2018. Periodic reanalysis of whole-genome sequencing data
460 enhances the diagnostic advantage over standard clinical genetic testing. *Eur J Hum Genet*
461 **26**: 740-744.
- 462 Costain G, Walker S, Marano M, Veenma D, Snell M, Curtis M, Luca S, Buera J, Arje D, Reuter
463 MS et al. 2020. Genome Sequencing as a Diagnostic Test in Children With Unexplained
464 Medical Complexity. *JAMA Netw Open* **3**: e2018109.
- 465 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
466 McCarthy SA, Davies RM et al. 2021. Twelve years of SAMtools and BCFtools.
467 *Gigascience* **10**.
- 468 Gargis AS, Kalman L, Berry MW, Bick DP, Dimmock DP, Hambuch T, Lu F, Lyon E,
469 Voelkerding KV, Zehnbauser BA et al. 2012. Assuring the quality of next-generation
470 sequencing in clinical laboratory practice. *Nat Biotechnol* **30**: 1033-1036.

- 471 Hon T, Mars K, Young G, Tsai YC, Karalius JW, Landolin JM, Maurer N, Kudrna D, Hardigan
472 MA, Steiner CC et al. 2020. Highly accurate long-read HiFi sequencing data for five
473 complex genomes. *Sci Data* **7**: 399.
- 474 Jennings L, Van Deerlin VM, Gulley ML, College of American Pathologists Molecular Pathology
475 Resource C. 2009. Recommended principles and practices for validating clinical molecular
476 pathology tests. *Arch Pathol Lab Med* **133**: 743-755.
- 477 Kadalayil L, Rafiq S, Rose-Zerilli MJ, Pengelly RJ, Parker H, Oscier D, Strefford JC, Tapper WJ,
478 Gibson J, Ennis S et al. 2015. Exome sequence read depth methods for identifying copy
479 number changes. *Brief Bioinform* **16**: 380-392.
- 480 Krusche P, Trigg L, Boutros PC, Mason CE, De La Vega FM, Moore BL, Gonzalez-Porta M,
481 Eberle MA, Tezak Z, Lababidi S et al. 2019. Best practices for benchmarking germline
482 small-variant calls in human genomes. *Nat Biotechnol* **37**: 555-560.
- 483 Linderman MD, Brandt T, Edelman L, Jabado O, Kasai Y, Kornreich R, Mahajan M, Shah H,
484 Kasarskis A, Schadt EE. 2014. Analytical validation of whole exome and whole genome
485 sequencing for clinical applications. *BMC Med Genomics* **7**: 20.
- 486 Logsdon GA, Vollger MR, Eichler EE. 2020. Long-read human genome sequencing and its
487 applications. *Nat Rev Genet* **21**: 597-614.
- 488 Mahmoud M, Gobet N, Cruz-Davalos DI, Mounier N, Dessimoz C, Sedlazeck FJ. 2019. Structural
489 variant calling: the long and the short of it. *Genome Biol* **20**: 246.
- 490 Majidian S, Agostinho DP, Chin CS, Sedlazeck FJ, Mahmoud M. 2023. Genomic variant
491 benchmark: if you cannot measure it, you cannot improve it. *Genome Biol* **24**: 221.
- 492 Marshall CR, Chowdhury S, Taft RJ, Lebo MS, Buchan JG, Harrison SM, Rowsey R, Klee EW,
493 Liu P, Worthey EA et al. 2020. Best practices for the analytical validation of clinical whole-
494 genome sequencing intended for the diagnosis of germline disease. *NPJ Genom Med* **5**: 47.
- 495 Matthijs G, Souche E, Alders M, Corveleyn A, Eck S, Feenstra I, Race V, Sistermans E, Sturm M,
496 Weiss M et al. 2016. Guidelines for diagnostic next-generation sequencing. *Eur J Hum*
497 *Genet* **24**: 1515.
- 498 Nurk S, Koren S, Rhie A, Rautiainen M, Bizkadze AV, Mikheenko A, Vollger MR, Altemose N,
499 Uralsky L, Gershman A et al. 2022. The complete sequence of a human genome. *Science*
500 **376**: 44-53.
- 501 NYSDOH. Laboratory Standards: Clinical Laboratory Evaluation Program. In
502 <https://www.wadsworth.org/regulatory/clep/clinical-labs/laboratory-standards>,
503 doi:<https://www.wadsworth.org/regulatory/clep/clinical-labs/laboratory-standards>.
- 504 Olson ND, Wagner J, Dwarshuis N, Miga KH, Sedlazeck FJ, Salit M, Zook JM. 2023. Variant
505 calling and benchmarking in an era of complete human genome sequences. *Nat Rev Genet*
506 **24**: 464-483.

- 507 Olson ND, Wagner J, McDaniel J, Stephens SH, Westreich ST, Prasanna AG, Johanson E, Boja
508 E, Maier EJ, Serang O et al. 2022. PrecisionFDA Truth Challenge V2: Calling variants
509 from short and long reads in difficult-to-map regions. *Cell Genom* **2**.
- 510 Pedersen BS, Quinlan AR. 2018. Mosdepth: quick coverage calculation for genomes and exomes.
511 *Bioinformatics* **34**: 867-868.
- 512 Pei S, Liu T, Ren X, Li W, Chen C, Xie Z. 2021. Benchmarking variant callers in next-generation
513 and third-generation sequencing analysis. *Brief Bioinform* **22**.
- 514 Poplin R, Chang PC, Alexander D, Schwartz S, Colthurst T, Ku A, Newburger D, Dijamco J,
515 Nguyen N, Afshar PT et al. 2018. A universal SNP and small-indel variant caller using
516 deep neural networks. *Nat Biotechnol* **36**: 983-987.
- 517 Rajagopalan R, Murrell JR, Luo M, Conlin LK. 2020. A highly sensitive and specific workflow
518 for detecting rare copy-number variants from exome sequencing data. *Genome Med* **12**:
519 14.
- 520 Rehm HL. 2017. Evolving health care through personal genomics. *Nat Rev Genet* **18**: 259-267.
- 521 Rehm HL, Bale SJ, Bayrak-Toydemir P, Berg JS, Brown KK, Deignan JL, Friez MJ, Funke BH,
522 Hegde MR, Lyon E et al. 2013. ACMG clinical laboratory standards for next-generation
523 sequencing. *Genet Med* **15**: 733-747.
- 524 Roy S, Coldren C, Karunamurthy A, Kip NS, Klee EW, Lincoln SE, Leon A, Pullambhatla M,
525 Temple-Smolkin RL, Voelkerding KV et al. 2018. Standards and Guidelines for Validating
526 Next-Generation Sequencing Bioinformatics Pipelines: A Joint Recommendation of the
527 Association for Molecular Pathology and the College of American Pathologists. *J Mol*
528 *Diagn* **20**: 4-27.
- 529 Santani A, Murrell J, Funke B, Yu Z, Hegde M, Mao R, Ferreira-Gonzalez A, Voelkerding KV,
530 Weck KE. 2017. Development and Validation of Targeted Next-Generation Sequencing
531 Panels for Detection of Germline Variants in Inherited Diseases. *Arch Pathol Lab Med*
532 **141**: 787-797.
- 533 Santani A, Simen BB, Briggs M, Lebo M, Merker JD, Nikiforova M, Vasalos P, Voelkerding K,
534 Pfeifer J, Funke B. 2019. Designing and Implementing NGS Tests for Inherited Disorders:
535 A Practical Framework with Step-by-Step Guidance for Clinical Laboratories. *J Mol Diagn*
536 **21**: 369-374.
- 537 Trost B, Walker S, Haider SA, Sung WWL, Pereira S, Phillips CL, Higginbotham EJ, Strug LJ,
538 Nguyen C, Raajkumar A et al. 2019. Impact of DNA source on genetic variant detection
539 from human whole-genome sequencing data. *J Med Genet* **56**: 809-817.
- 540 Wagner J, Olson ND, Harris L, Khan Z, Farek J, Mahmoud M, Stankovic A, Kovacevic V, Yoo
541 B, Miller N et al. 2022a. Benchmarking challenging small variants with linked and long
542 reads. *Cell Genom* **2**.

- 543 Wagner J, Olson ND, Harris L, McDaniel J, Cheng H, Functammasan A, Hwang YC, Gupta R,
544 Wenger AM, Rowell WJ et al. 2022b. Curated variation benchmarks for challenging
545 medically relevant autosomal genes. *Nat Biotechnol* **40**: 672-680.
- 546 Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, Ebler J, Functammasan
547 A, Kolesnikov A, Olson ND et al. 2019. Accurate circular consensus long-read sequencing
548 improves variant detection and assembly of a human genome. *Nat Biotechnol* **37**: 1155-
549 1162.
- 550 Whitford W, Lehnert K, Snell RG, Jacobsen JC. 2019. Evaluation of the performance of copy
551 number variant prediction tools for the detection of deletions from whole genome
552 sequencing data. *J Biomed Inform* **94**: 103174.
- 553 Yang Y, del Gaudio D, Santani A, Scott SA. 2024. Applications of genome sequencing as a single
554 platform for clinical constitutional genetic testing. *GIM Open*
555 doi:10.1016/j.gimo.2024.101840.
- 556 Yao RA, Akinrinade O, Chaix M, Mital S. 2020. Quality of whole genome sequencing from blood
557 versus saliva derived DNA in cardiac patients. *BMC Med Genomics* **13**: 11.
- 558 Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, Weng Z, Liu Y, Mason CE, Alexander
559 N et al. 2016. Extensive sequencing of seven human genomes to characterize benchmark
560 reference materials. *Sci Data* **3**: 160025.
- 561 Zook JM, Hansen NF, Olson ND, Chapman L, Mullikin JC, Xiao C, Sherry S, Koren S, Phillippy
562 AM, Boutros PC et al. 2020. A robust benchmark for detection of germline large deletions
563 and insertions. *Nat Biotechnol* **38**: 1347-1355.
- 564 Zook JM, McDaniel J, Olson ND, Wagner J, Parikh H, Heaton H, Irvine SA, Trigg L, Truty R,
565 McLean CY et al. 2019. An open resource for accurately benchmarking small variant and
566 reference calls. *Nat Biotechnol* **37**: 561-566.
- 567

Table 1. Long-read HiFi genome sequencing SNV/indel (<50 bp) accuracy.

		Genome-Wide					RefSeq CDS		
		Low Complexity	Low Mappability	Segmental Duplications	All Difficult Regions	Not in Any Difficult Region	All	Not in Any Difficult Region	All
SNVs	Count	160865	190416	120916	584743	2718604	3303346	14059	20593
	Recall	99.05%	97.84%	96.63%	98.93%	99.91%	99.74%	99.86%	99.56%
	Precision	99.54%	99.81%	99.67%	99.80%	99.99%	99.95%	99.98%	99.96%
Insertions (1-5bp)	Count	128237	4241	4596	140934	69118	210639	46	109
	Recall	98.57%	98.50%	98.62%	98.64%	99.90%	99.05%	99.71%	99.44%
	Precision	98.42%	99.14%	98.98%	98.53%	99.94%	98.99%	100.00%	99.35%
Insertions (6-15bp)	Count	14435	342	486	15460	6496	22004	5	42
	Recall	98.67%	98.39%	98.77%	98.70%	99.73%	99.00%	100.00%	99.58%
	Precision	98.98%	98.52%	98.87%	99.02%	99.94%	99.29%	100.00%	99.25%
Insertions (≥16bp)	Count	2415	122	136	2747	2016	4776	2	13
	Recall	98.76%	96.91%	97.86%	98.80%	99.80%	99.22%	85.71%	97.52%
	Precision	99.16%	97.23%	98.15%	99.17%	99.86%	99.46%	100.00%	96.63%
Deletions (1-5bp)	Count	136024	4600	4507	148953	69677	218129	59	142
	Recall	98.78%	98.45%	98.26%	98.82%	99.87%	99.15%	99.48%	99.00%
	Precision	98.89%	99.18%	98.91%	98.95%	99.93%	99.25%	99.72%	99.49%
Deletions (6-15bp)	Count	16664	466	544	17758	6647	24190	10	50
	Recall	98.68%	98.11%	98.29%	98.68%	99.57%	98.93%	100.00%	99.74%
	Precision	98.88%	98.64%	98.40%	98.90%	99.86%	99.15%	100.00%	99.74%
Deletions (≥16bp)	Count	3424	152	116	3681	1621	5170	2	12
	Recall	99.33%	98.53%	98.15%	99.32%	99.63%	99.42%	85.71%	100.00%
	Precision	99.49%	98.93%	98.15%	99.49%	99.88%	99.61%	100.00%	100.00%
All indels	Count	285386	9840	10200	313715	155551	469067	125	362
	Recall	98.78%	98.43%	98.42%	98.82%	99.86%	99.16%	99.65%	99.26%
	Precision	98.70%	99.09%	98.89%	98.78%	99.93%	99.14%	99.88%	99.38%
SNVs and indels	Count	462064	200340	131301	914275	2874179	3788255	14184	20961
	Recall	98.81%	97.87%	96.78%	98.85%	99.91%	99.65%	99.86%	99.55%
	Precision	98.99%	99.77%	99.61%	99.42%	99.99%	99.85%	99.98%	99.95%

indels: insertions/deletions; RefSeq CDS: NCBI Reference Sequence gene coding sequence; SNVs: single nucleotide variants.

Table 2. Long-read HiFi genome sequencing SNV/indel (<50 bp) reproducibility (GIAB high confidence).

		Genome-Wide						RefSeq CDS	
		Low Complexity	Low Mappability	Segmental Duplications	All Difficult Regions	Not in Any Difficult Region	All	Not in Any Difficult Region	All
SNVs	Count	170332	181786	113753	584487	2721331	3305818	14023	20467
	Concordance	99.10%	99.51%	99.30%	99.54%	99.95%	99.88%	99.91%	99.85%
Insertions (1-5bp)	Count	139069	4085	4495	151456	69083	221146	48	115
	Concordance	97.41%	99.38%	99.26%	97.61%	99.94%	98.35%	99.66%	99.17%
Insertions (6-15bp)	Count	15335	335	490	16341	6487	22882	5	43
	Concordance	98.13%	98.99%	99.31%	98.23%	99.96%	98.74%	100.00%	100.00%
Insertions (≥16bp)	Count	2527	118	119	2805	1847	4664	2	12
	Concordance	98.22%	97.70%	98.01%	98.31%	99.84%	98.93%	100.00%	95.22%
Deletions (1-5bp)	Count	150504	4459	4313	163229	69693	232407	58	141
	Concordance	98.13%	99.35%	99.21%	98.25%	99.97%	98.77%	99.72%	99.41%
Deletions (6-15bp)	Count	17426	457	552	18499	6652	24940	10	47
	Concordance	98.38%	99.17%	98.96%	98.46%	99.96%	98.86%	100.00%	100.00%
Deletions (≥16bp)	Count	3525	160	101	3771	1583	5220	3	16
	Concordance	99.21%	98.92%	99.01%	99.26%	99.96%	99.50%	100.00%	100.00%
All indels	Count	311652	9518	9877	339360	155318	494492	125	370
	Concordance	97.90%	99.32%	99.21%	98.05%	99.95%	98.64%	99.73%	99.37%
SNVs and indels	Count	498720	191398	123824	940587	2876675	3817077	14148	20841
	Concordance	98.27%	99.50%	99.30%	98.96%	99.95%	99.71%	99.91%	99.84%

indels: insertions/deletions; RefSeq CDS: NCBI Reference Sequence gene coding sequence; SNVs: single nucleotide variants.

Table 3. Long-read HiFi genome sequencing SNV/indel (<50 bp) paired specimen concordance (GIAB high confidence).

		Genome-Wide						RefSeq CDS	
		Low Complexity	Low Mappability	Segmental Duplications	All Difficult Regions	Not in Any Difficult Region	All	Not in Any Difficult Region	All
Blood vs Saliva	Count	246228.5	158890	94436.5	639126	2753708.5	3392621.5	13785	19515.5
	Concordance	99.40%	99.54%	99.34%	99.60%	99.94%	99.88%	99.94%	99.84%
Blood vs Swab	Count	246296.5	159100	94574	639544.5	2754294.5	3393623	13781	19527.5
	Concordance	99.52%	99.57%	99.37%	99.66%	99.96%	99.90%	99.94%	99.85%
Swab vs Saliva	Count	246296.5	159100	94574	639544.5	2754294.5	3393623	13781	19527.5
	Concordance	99.30%	99.47%	99.34%	99.54%	99.94%	99.87%	99.93%	99.84%

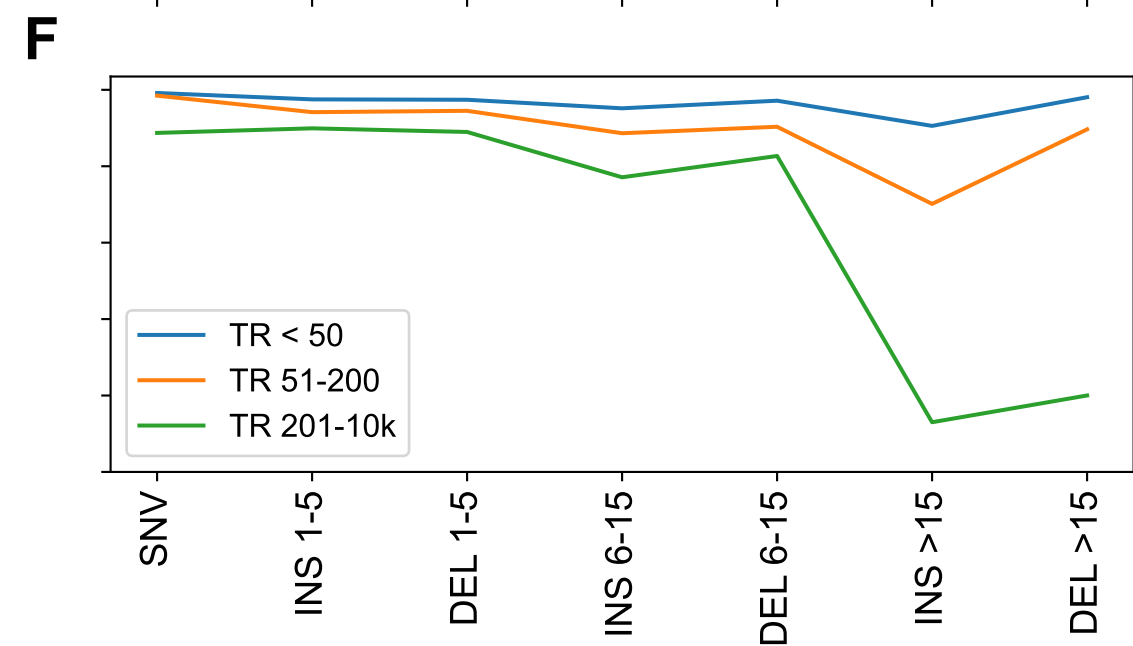
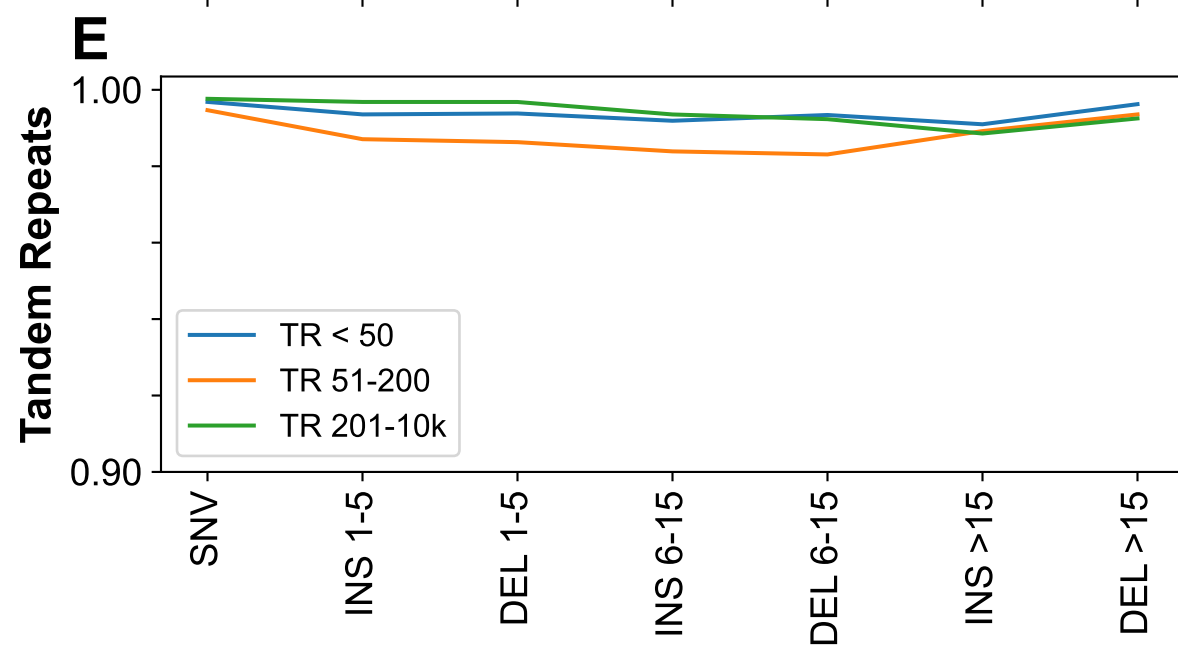
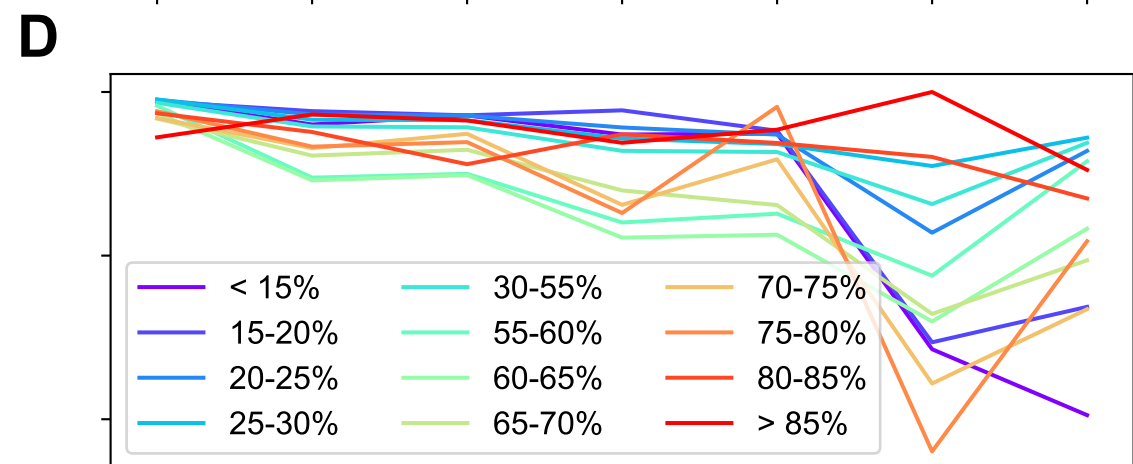
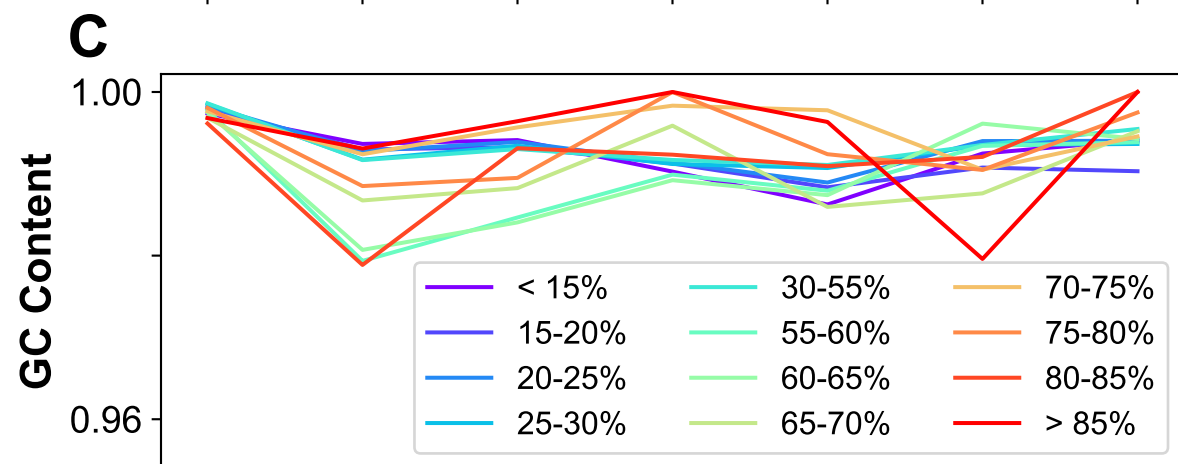
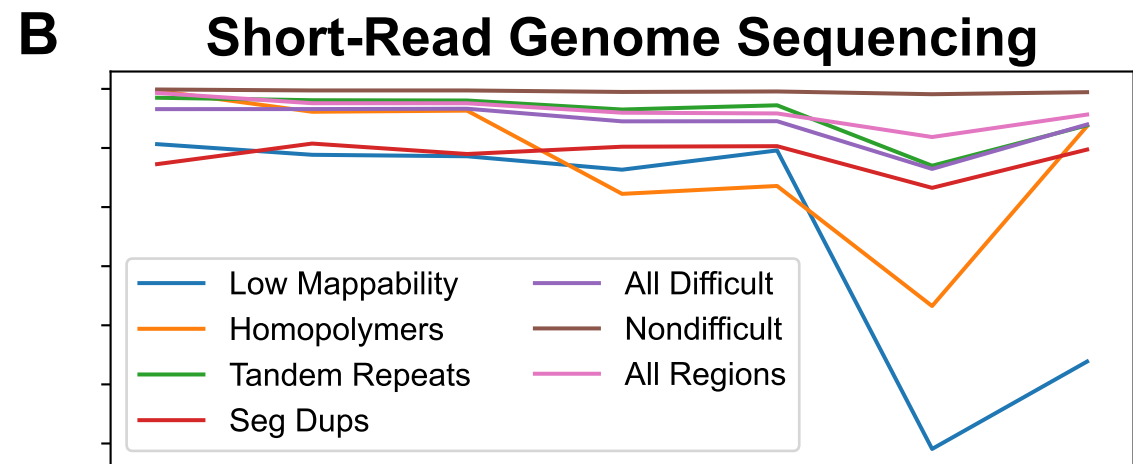
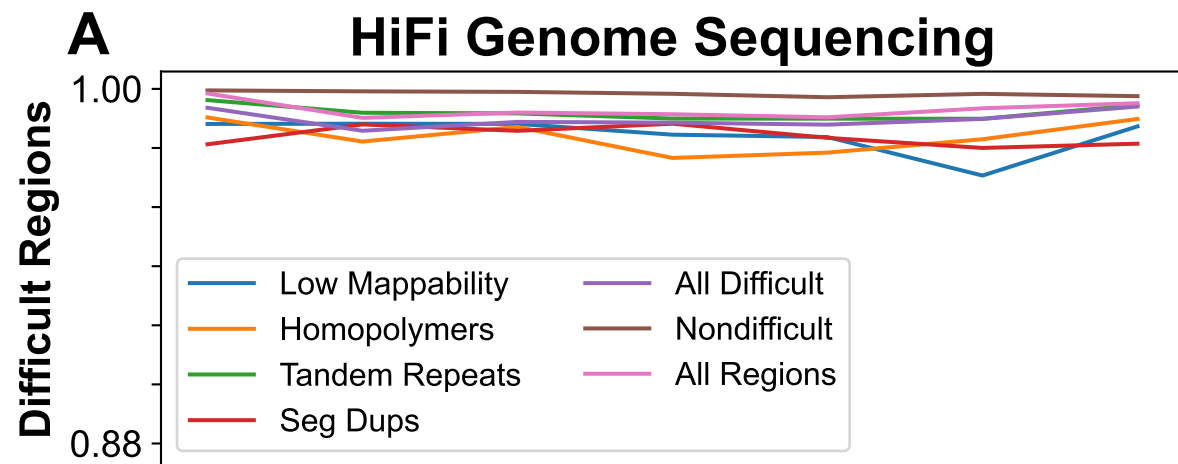
RefSeq CDS: NCBI Reference Sequence gene coding sequence.

FIGURE LEGENDS

Figure 1. SNV/indel accuracy across GIAB-defined genomic regions. Small variant F1-scores for HiFi genome sequencing (**A, C, E**) and short-read genome sequencing (**B, D, F**) across 22 of the 100 interrogated GIAB-defined genomic regions. Results summarized by difficulty (mappability, homopolymers, tandem repeats, segmental duplications, all difficult, not in any difficult, all regions), GC content (<15% to >85%) and tandem repeats (<50 bp, 51-200 bp, 201-10,000 bp), and stratified by variant type (SNV, indels 1-5 bp, 6-15 bp, and ≥ 16 bp).

Figure 2. SNV/indel accuracy and depth stratification. Plots of F1-scores across sequencing depths and comparing HiFi genome sequencing and short-read genome sequencing: (**B**) single nucleotide variants (SNVs); (**B**) insertions/deletion variants (indels); and (**C**) SNVs and indels combined.

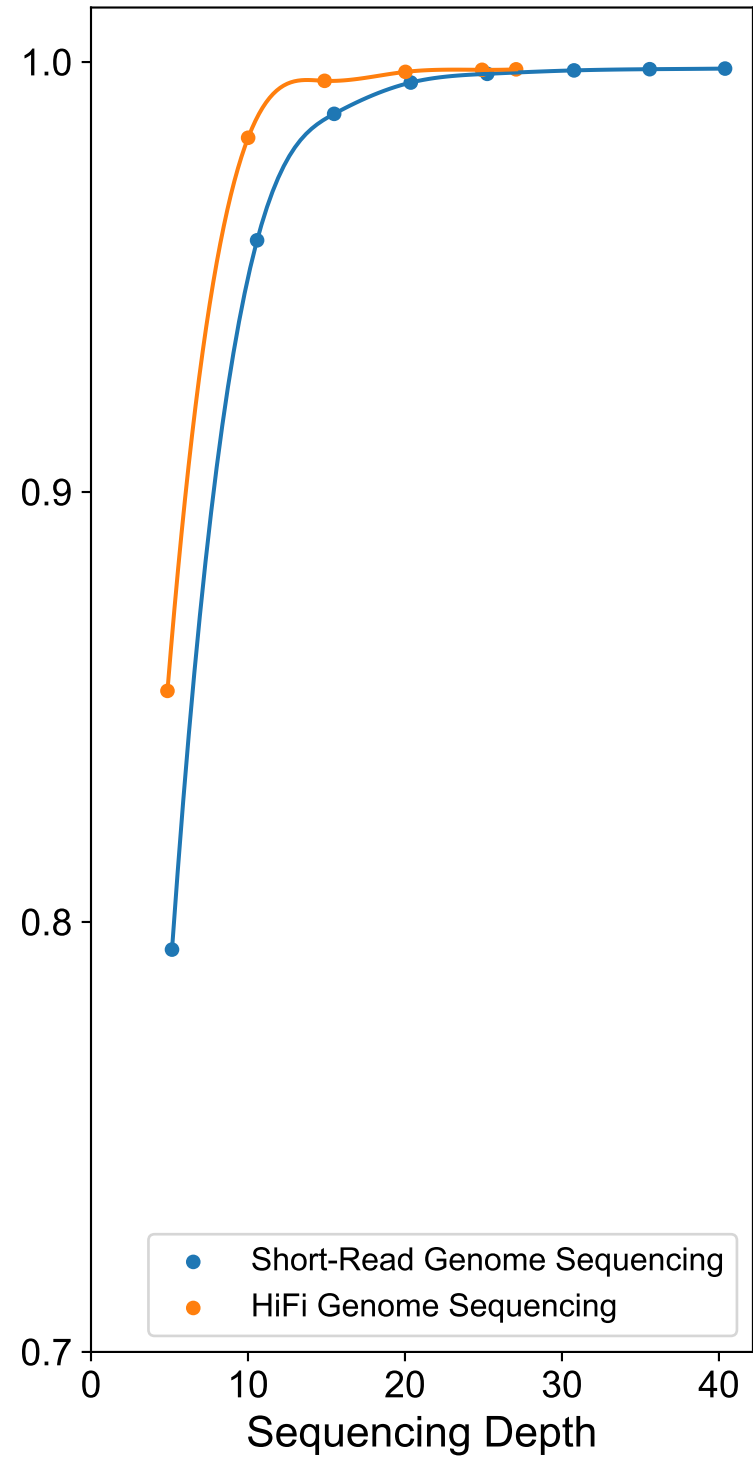
Figure 3. Automated library preparation validation with low depth HiFi genome sequencing. Two Miro Canvas instruments were validated using a single SMRTcell of data ($\sim 9\times$ each) with small variant (SNV/indel) benchmarking of NA12878 and fitting results to reference curves defined by manual library preparation of three GIAB reference samples (NA12878, NA24385, NA24631). Error bars represent standard deviations. Automated library preparation results were considered acceptable if Miro Canvas recall (**A**) and precision (**B**) values were equivalent or greater than the average manual preparation reference material accuracy results at comparable depths.



F1 score vs Depth

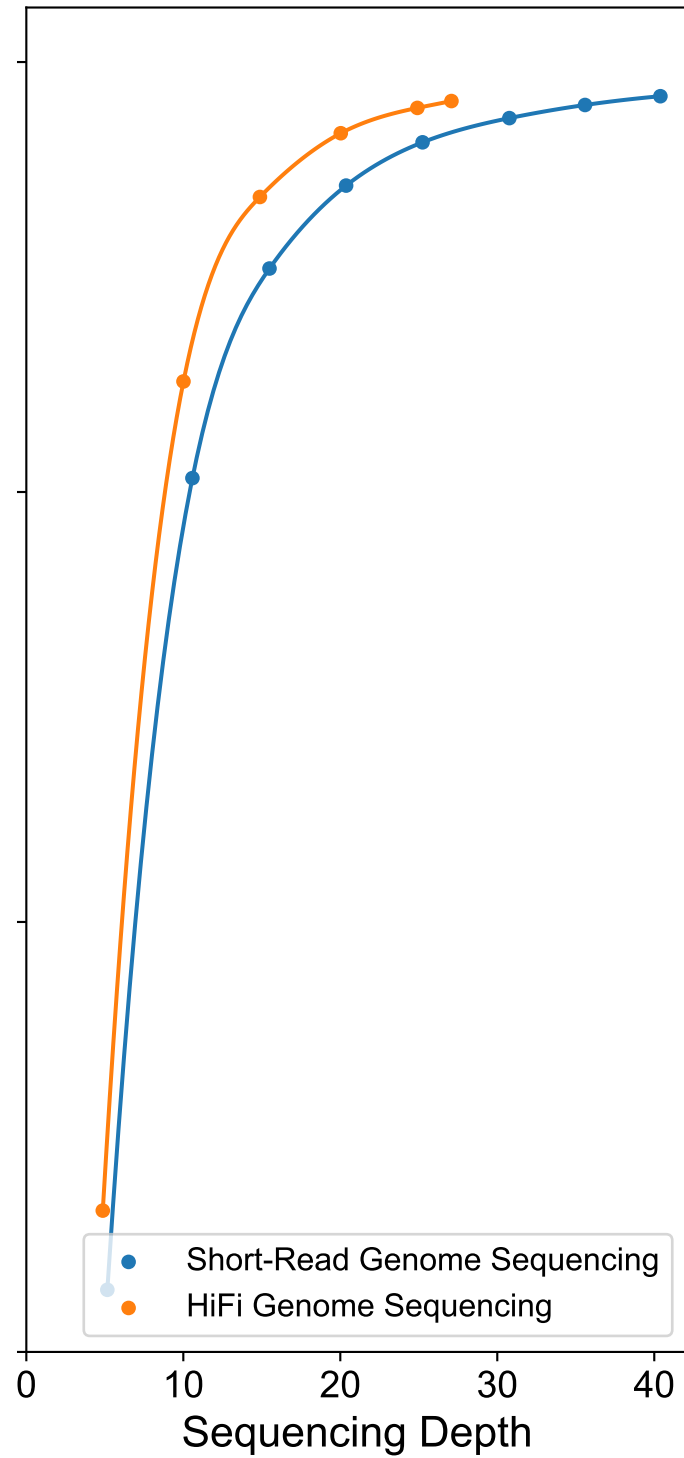
A

SNV



B

indel



C

SNV+indel

