

Establishing reference samples for detection of somatic mutations and germline variants with NGS technologies

Li Tai Fang^{1*}, Bin Zhu^{2*}, Yongmei Zhao^{3*}, Wanqiu Chen⁴, Zhaowei Yang^{4,31}, Liz Kerrigan⁵, Kurt Langenbach⁵, Maryellen de Mars⁵, Charles Lu⁶, Kenneth Idler⁶, Howard Jacob⁶, Ying Yu⁷, Luyao Ren⁷, Yuanting Zheng⁷, Erich Jaeger⁸, Gary Schroth⁸, Ogan D. Abaan⁸, Justin Lack³, Tsai-Wei Shen³, Keyur Talsania³, Zhong Chen⁴, Seta Stanbouly⁴, Jyoti Shetty⁹, Parimal Kumar⁹, John Bettridge⁹, Bao Tran⁹, Daoud Meerzaman¹⁰, Cu Nguyen¹⁰, Virginie Petitjean¹¹, Marc Sultan¹¹, Margaret Cam¹², Tiffany Hung¹³, Eric Peters¹³, Rasika Kalamegham¹³, Sayed Mohammad Ebrahim Sahraeian¹, Marghoob Mohiyuddin¹, Yunfei Guo¹, Lijing Yao¹, Lei Song², Hugo YK Lam¹, Jiri Drabek^{14,15}, Roberta Maestro^{15,16}, Daniela Gasparotto^{15,16}, Sulev Kõks^{15,17,18}, Ene Reimann^{15,18}, Andreas Scherer^{19,15}, Jessica Nordlund^{20,15}, Ulrika Liljedahl^{20,15}, Roderick V Jensen²¹, Mehdi Pirooznia²², Zhipan Li²³, Chunlin Xiao²⁴, Stephen Sherry²⁴, Rebecca Kusko²⁵, Malcolm Moos²⁶, Eric Donaldson²⁷, Zivana Tezak²⁸, Baitang Ning²⁹, Weida Tong²⁹, Jing Li³⁰, Penelope Duerken-Hughes³¹, Huixiao Hong^{29#}, Leming Shi^{7#}, Charles Wang^{4,30,31#}, Wenming Xiao^{28,29#}, and The Somatic Working Group of SEQC-II Consortium

¹Bioinformatics Research & Early Development, Roche Sequencing Solutions Inc., 1301 Shoreway Road, Suite #300, Belmont, CA 94002; ²Division of Cancer Epidemiology and Genetics, National Cancer Institute, National Institutes of Health, 9609 Medical Center Drive, Bethesda, Maryland 20892, USA.; ³Advanced Biomedical and Computational Sciences, Biomedical Informatics and Data Science Directorate, Frederick National Laboratory for Cancer Research, 8560 Progress Drive, Frederick, MD21701; ⁴Center for Genomics, Loma Linda University School of Medicine, 11021 Campus St., Loma Linda, CA 92350; ⁵ATCC (American Type Culture Collection), 10801 University Blvd, Manassas, VA 20110; ⁶Computational Genomics, Genomics Research Center (GRC), AbbVie, 1 North Waukegan Road, North Chicago, IL 60064; ⁷State Key Laboratory of Genetic Engineering, School of Life Sciences and Shanghai Cancer Center, Fudan University, 2005 SongHu Road, Shanghai, China 200438; ⁸Core applications group, Product development, Illumina Inc, 200 Lincoln Centre Dr. Foster City, CA 94404; ⁹Genomics Laboratory, Cancer Research Technology Program, Frederick National Laboratory for Cancer Research, 8560 Progress Drive, Frederick, MD21701; ¹⁰Computational Genomics Research, Center for Biomedical Informatics and Information Technology (CBIIT), National Cancer Institute, 9609 Medical Center Drive, Rockville, MD 20850; ¹¹Biomarker development, Novartis Institutes for Biomedical Research, Fabrikstrasse 10, CH-4056 Basel, Switzerland; ¹²CCR Collaborative Bioinformatics Resource (CCBR), Office of Science and Technology Resources, Center for Cancer Research, NCI, Bldg 37, Rm 3041C, 37 Convent Drive, Bethesda, MD 20892; ¹³Companion Diagnostics Development, Oncology Biomarker Development, Genentech, 1 DNA Way, South San Francisco, CA 94080; ¹⁴IMTM, Faculty of Medicine and Dentistry, Palacky University Olomouc, Hnevotinska 5, 77900 Olomouc, the Czech Republic; ¹⁵member of EATRIS ERIC-European Infrastructure for Translational Medicine; ¹⁶Centro di Riferimento Oncologico di Aviano (CRO) IRCCS, National Cancer Institute, Unit of Oncogenetics and Functional Oncogenomics, Via Gallini 2, 33081 Aviano (PN), Italy; ¹⁷Perron Institute for Neurological and Translational Science, Verdun St, Nedlands, Western Australia, 6009, Australia; ¹⁸ the Centre for

44 Molecular Medicine and Innovative Therapeutics, Murdoch University, Murdoch, 6150,
45 Western Australia; ¹⁹Institute for Molecular Medicine Finland (FIMM), HiLIFE, P.O.Box 20, FI-
46 00014 University of Helsinki, Finland; ²⁰Department of Medical Sciences, Molecular Medicine
47 and Science for Life Laboratory, Uppsala University, Molekylär Medicin, Box 1432, BMC,
48 Uppsala, 75144 Sweden; ²¹Department of Biological Sciences, Virginia Tech, Life Sciences 1, 970
49 Washington St., Blacksburg, VA 24061; ²²Bioinformatics and Computational Biology Core,
50 National Heart Lung and Blood Institute, National Institutes of Health, 12 SOUTH DR. Bethesda
51 MD 20892; ²³Sentieon Inc., 465 Fairchild Drive, Suite 135, Mountain View CA 94043; ²⁴National
52 Center for Biotechnology Information, National Library of Medicine, National Institutes of
53 Health, 45 Center Drive, Bethesda, Maryland 20894; ²⁵Immuneering Corporation, One Broadway
54 14th Fl Cambridge MA 02142 USA; ²⁶The Center for Biologics Evaluation and Research, U.S.
55 Food and Drug Administration, FDA, Silver Spring, Maryland; ²⁷Division of Antiviral Products,
56 Office of Antimicrobial, Center for Drug Evaluation and Research, FDA, Silver Spring, Maryland;
57 ²⁸The Center for Devices and Radiological Health, U.S. Food and Drug Administration, FDA,
58 Silver Spring, Maryland; ²⁹Bioinformatics branch, Division of Bioinformatics and Biostatistics,
59 National Center for Toxicological Research, Food and Drug Administration, 3900 NCTR Road,
60 Jefferson, AR 72079; ³⁰Department of Allergy and Clinical Immunology, State Key Laboratory of
61 Respiratory Disease, National Clinical Research Center for Respiratory Disease, Guangzhou
62 Institute of Respiratory Health, the First Affiliated Hospital of Guangzhou Medical University,
63 Guangzhou, Guangdong, 510182, P. R. China; ³¹Department of Basic Science, Loma Linda
64 University School of Medicine, Loma Linda, CA 92350

65
66 * Contributed equally

67 # To whom correspondence should be addressed to: Huixiao.Hong@fda.hhs.gov,
68 lemingshi@fudan.edu.cn, chwang@llu.edu, and Wenming.Xiao@fda.hhs.gov
69

70 Abstract

71

72 We characterized two reference samples for NGS technologies: a human triple-negative
73 breast cancer cell line and a matched normal cell line. Leveraging several whole-genome
74 sequencing (WGS) platforms, multiple sequencing replicates, and orthogonal mutation
75 detection bioinformatics pipelines, we minimized the potential biases from sequencing
76 technologies, assays, and informatics. Thus, our “truth sets” were defined using evidence from
77 21 repeats of WGS runs with coverages ranging from 50X to 100X (a total of 140 billion reads).
78 These “truth sets” present many relevant variants/mutations including 193 COSMIC mutations
79 and 9,016 germline variants from the ClinVar database, nonsense mutations in *BRCA1/2* and
80 missense mutations in *TP53* and *FGFR1*. Independent validation in three orthogonal
81 experiments demonstrated a successful stress test of the truth set. We expect these reference
82 materials and “truth sets” to facilitate assay development, qualification, validation, and
83 proficiency testing. In addition, our methods can be extended to establish new fully
84 characterized reference samples for the community.

85

86

87 Introduction

88

89 In oncology, accurate somatic mutation detection is essential to diagnose cancer, pinpoint
90 targeted therapies, predict survival, and identify resistance mutations. Despite the recent
91 explosion of technological advancements, many studies have reported difficulties in obtaining
92 consistent and concordant somatic mutation calls from individual platforms or pipelines¹⁻³, which
93 hampers clinical validation and advancement of these biomarkers.

94

95 As more sequencing technologies can detect clinically actionable somatic mutations for
96 oncology, the need grows stronger for benchmark samples with known “ground-truth” variants.
97 Such a publicly available sample set would allow platform and pipeline developers to quantify
98 accuracy of somatic mutation calls, study reproducibility across platforms or pipelines, perform
99 validation using orthogonal techniques, and calibrate best practices of protocols and methods.
100 The FDA has released a guidance on the use of NGS technologies for *in vitro* diagnosis of
101 suspected germline diseases⁴, in which well-characterized reference materials are
102 recommended to establish NGS test performance.

103

104 In the absence of well-characterized samples with somatic mutations, normal samples
105 such as the Platinum Genome⁵, HapMap⁶ cell lines, or Genome in a Bottle (GiaB) consortium
106 materials^{7,8} are often used in clinical test development and validation of somatic applications.
107 Also there are some gene-specific reference samples available, such as KRAS in the WHO 1st
108 International Reference Panel⁹, or from synthetic materials¹⁰. Such samples do not adequately
109 address cancer-specific quality metrics such as somatic mutation variant allele frequency (VAF),
110 heterogeneity, tumor mutation burden (TMB), etc. Therefore, cancer reference samples with
111 an abundance of well-defined genetic alterations characterized across the whole genome are
112 highly desirable and urgent needed.

113

114 Previous attempt has characterized a cancer cell line (from metastatic melanoma) that
115 inquired somatic mutations (SNV/indels) in exon regions only. Germline variants and somatic
116 mutations across the rest of the genome were not defined¹¹. In addition, this dataset is
117 distributed under dbGAP-controlled access, limiting its accessibility and utility. In fact, a recent
118 landscape analysis of currently available somatic variant reference samples published by the
119 Medical Devices Innovation Consortium (MDIC) did not identify any reference mutation sets that
120 can be used to evaluate the somatic mutation calling accuracy on a whole-genome basis¹².

121

122 To fulfill this unmet need, we chose a pair of cell lines, HCC1395 (triple-negative breast
123 cancer) and HCC395BL (B lymphocytes) from the same donor, supplied by the American Type
124 Culture Collection (ATCC). These two specific cell lines were chosen because they are rich in
125 testable features (CNVs, SNVs, indels, SVs, and genome rearrangements¹³), and may have a
126 potential to serve as a long-term, publicly available, and renewable reference samples with
127 appropriate consent from donor. Using multiple next generation sequencing (NGS) platforms,
128 sequencing centers, and various bioinformatics analysis pipelines we profiled these tumor-
129 normal matching cell lines. Thus, we minimized biases that were specific to any platform,
130 sequencing center, or bioinformatic algorithm, to create a list of high-confidence mutation calls

131 across the whole genome, here called the “truth set.” A subset of these calls was further
132 confirmed with orthogonal targeted sequencing and Whole Exome Sequencing (WES). We also
133 sequenced a series of titrations between HCC1395 and HCC1395BL genomic DNA (gDNA) to
134 confirm candidate somatic SNV/indels.

135

136 We defined truth sets containing somatic mutations and germline variants in a paired cell
137 lines, HCC1395/HCC1395BL, with methods that minimized potential bias from library preparation,
138 sequencing center, or bioinformatics pipeline. While the “truth set” germline variants in
139 HCC1395BL can be used for benchmarking germline variant detection, the “truth set” somatic
140 mutations in HCC1395 can be used for benchmarking cancer mutation detection with VAF as low
141 as 5%. Many of variants and mutations have clinical implications. In the coding regions, a total of
142 193 somatic mutations are documented in the COSMIC database and 8 germline variants are
143 annotated as pathogenic in the ClinVar database. Interestingly, there is a nonsense somatic
144 mutation in the *BRCA2* gene and a nonsense germline variant in the *BRAC1* gene. Other hotspot
145 somatic mutations are also observed in the *TP53* and *FGFR1* genes. Thus, we believe these paired
146 cell lines may be highly valuable for those looking for reference samples to benchmark products
147 in detection of mutations in these four genes.

148

149 **Results**

150

151 **Massive data generated to characterize the reference samples**

152

153 To provide reference samples for the community well into the future, a matched pair,
154 HCC1395 and HCC1395BL was selected for profiling¹⁴. Previous studies of this triple negative
155 breast cancer cell line have revealed the existence of many somatic structural and ploidy
156 changes¹³, which are confirmed by our cell karyotype and cytogenetic analysis (Suppl. Fig S1, S2).
157 Several attempts have been made to identify SNVs and small indels¹⁵⁻¹⁷. Given that appropriate
158 consent from the donor has been obtained for tumor HCC1395 and normal HCC1395BL for the
159 purposes of genomic research, we sought to characterize this pair of cell lines as publicly available
160 reference samples for the NGS community. In this manuscript, we focused our efforts on germline
161 and somatic SNVs and indels. By performing numerous sequencing experiments with multiple
162 platforms at different sequencing centers, we obtained high-confidence call sets of both somatic
163 and germline SNVs and indels (Table 1). Larger structural variants and copy number analysis will
164 be included in a separate manuscript that will discuss these findings in greater detail.

165

166 **Initial Determination of Somatic Mutation Call Set**

167

168 High-confidence somatic SNVs and indels were obtained based primarily on 21 pairs of
169 tumor-normal Whole Genome Sequencing (WGS) replicates from six sequencing centers;
170 sequencing depth ranged from 50X to 100X (see manuscript NBT-RA46164). Each of the 21
171 tumor-normal sequencing replicates was aligned with BWA MEM¹⁸, Bowtie2¹⁹, and NovoAlign²⁰
172 to create 63 pairs of tumor-normal Binary Sequence Alignment/Map (BAM) files. Six mutation
173 callers (MuTect2²¹, SomaticSniper²², VarDict²³, MuSE²⁴, Strelka2²⁵, and TNscope²⁶) were applied
174 to discover somatic mutation candidates for each pair of tumor-normal BAM files (Fig. 1).

175 SomaticSeq²⁷ was then utilized to combine the call sets and classify the candidate mutation calls
176 into “PASS”, “REJECT”, or “LowQual”. Four confidence levels (HighConf, MedConf, LowConf, and
177 Unclassified) were determined based on the cross-aligner and cross-sequencing center
178 reproducibility of each mutation call. HighConf and MedConf calls were grouped together as the
179 “truth set” (also known as high-confidence somatic mutations). The call set in its entirety is
180 referred to as the “super set” which includes low-confidence (LowConf) and likely false positive
181 (Unclassified) calls. For low-*VAF* (Variant Allele Frequency) calls, a HiSeq data set with 300×
182 coverage and a NovaSeq data set with 380× coverage were employed to rescue initial LowConf
183 and Unclassified calls into the truth set. The details are described in the Methods.

184
185 A breakdown of the four confidence levels is displayed in Fig. 2a. In the truth set, HighConf
186 calls consist of 94% of the SNVs and 79% of the indels. LowConf calls typically do not have enough
187 “PASS” classifications across the 63 data sets to be included in the truth set. Variants calls labeled
188 as Unclassified are not reproducible and likely false positives, with more “REJECT” classifications
189 than “PASS”. The vast majority of the calls in the super set are either HighConf or Unclassified. In
190 other words, super set calls tend to be either highly reproducible or not at all reproducible.

191
192 In general, HighConf calls were classified as “PASS” in the vast majority of the data sets,
193 with no variant read in the matched normal and high mapping quality scores. MedConf calls
194 tended to be low-*VAF* ($VAF \lesssim 0.10$) variants. Due to stochastic sampling of low frequency variants,
195 MedConf calls were not reproduced as highly across different sequencing replicates as HighConf
196 calls. LowConf calls (not a part of truth set) tended to have *VAF* near or below our detection limits
197 ($VAF \lesssim 0.05$). Distinguishing the LowConf calls with sequencing noise is challenging because they
198 were not reproduced enough to be high-confidence calls (Fig. 2b).

199 200 **Independent AmpliSeq confirmation of Call Set**

201
202 We randomly selected 450 SNV and 21 indel calls of different confidence levels from the
203 super set and performed PCR-based AmpliSeq with approximately 2000× depth for tumor and
204 normal cells on an Illumina MiSeq sequencer. As we treated the AmpliSeq data set as a
205 confirmatory experiment, simple rules were devised to determine whether a variant call was
206 deemed positively confirmed, not confirmed, or uninterpretable based on the presence or
207 absence of somatic mutation evidence in the AmpliSeq data. Overall, positively confirmed calls
208 had at least 100 variant-supporting reads in the tumor but had no variant read in the normal
209 sample, despite sequencing depths of 600× or more in the normal. Not confirmed calls either
210 had no more variant-supporting read than the expected from base call errors, and/or had
211 $VAF \geq 10\%$ in the normal cells. Uninterpretable calls did not satisfy the criteria for either positive
212 or no validation, either because they did not have enough read depth (<50) or had fewer than
213 10 variant-supporting reads. (See Methods for details).

214
215 Both HighConf and MedConf SNV calls had very high coverage in validation and thus had
216 impressive validation rates (99% and 92%) (Table 2). There were only three HighConf SNV calls
217 that were not confirmed by AmpliSeq. Two of them had germline signals below the detection
218 limit of 50× in the WGS, and the third one was likely an actual somatic mutation missed by

219 AmpliSeq. There were only seven “positively confirmed” Unclassified SNV calls. Four of those
220 seven were either a part of di-nucleotide change or had deletions within 1 bp of the call. The
221 other three had low mapping quality scores (MQ), which drove the categorization of
222 “Unclassified”. This result suggests that some of the “positively confirmed” Unclassified calls
223 might be false positives after all, but it also exposes the limitations of our truth set with regard
224 to complex variants and low mappability regions. LowConf and Unclassified calls (not part of the
225 truth set) also had higher fractions of uninterpretable calls, which consist of low-coverage
226 genomic positions or ambiguous variant signals. In addition, there were also 17 HighConf, 2
227 MedConf, 1 LowConf, and 1 Unclassified indel calls re-sequenced by AmpliSeq. The only not
228 confirmed HighConf indel call was caused by a germline signal. The lone Unclassified indel call
229 was not confirmed (we expect Unclassified calls to be not confirmed). For the inquisitive reader,
230 these discrepant calls (i.e., not confirmed HighConf calls and confirmed Unclassified calls) are
231 discussed in greater detail in the Supplementary Material.

232
233 The VAF calculated from the truth set correlated highly with the VAF calculated from
234 AmpliSeq data set, especially for HighConf calls (Fig. 2c). On the other hand, almost all the data
235 points at the bottom of the graph (i.e., VAF = 0 by AmpliSeq) are Unclassified calls (red). It
236 suggests that despite high VAFs (from 21 WGS replicates) for some of the calls, they were
237 categorized correctly as Unclassified (implying likely false positives). In addition, a large number
238 of uninterpretable Unclassified calls (red X’s) lying at the bottom suggest those were correctly
239 labeled as Unclassified in addition to the not confirmed ones (open red circles). Moreover, some
240 of the seven “positively confirmed” Unclassified calls had dubious supporting evidence. Taken
241 together, these results suggest that the actual true positive rate for the Unclassified calls may be
242 even lower than the validation rate (11%) we reported here. The indel equivalent is portrayed in
243 Suppl. Fig S7a.

244

245 **Orthogonal Confirmation of Call Set with WES on Ion Torrent**

246
247 We have also sequenced the tumor-normal pair with Whole Exome Sequencing (WES) on
248 the Ion Torrent S5 XL sequencer with the Agilent SureSelect All Exon + UTR v6 hybrid capture.
249 The sequencing depths for the HCC1395 and HCC1395BL were 34× and 47×, respectively.
250 Results from this Ion Torrent sequencing were leveraged to evaluate high-VAF SNV calls (Table 1
251 and 2). HighConf and MedConf SNV calls had high positive validation rates (99% and 89%).
252 However, because the Ion Torrent sequencing was performed at much lower depth, nearly 50%
253 of the calls were deemed uninterpretable (compared with 16% for AmpliSeq, despite having
254 AmpliSeq custom target enriched for low-confidence calls vs. WES). The trend of higher
255 uninterpretable fraction with lower confidence level calls was even more pronounced in this data
256 set because the coverage was too low to confirm or invalidate many low-VAF calls. The validation
257 rate for MedConf calls (predominantly low-VAF calls) may have suffered due to low coverage.

258
259 The VAF correlation between truth set and Ion Torrent WES ($R=0.928$) is lower than that
260 between truth set and AmpliSeq ($R=0.958$), although the vast majority of the HighConf SNV calls
261 in Ion Torrent data still stay within the 95% confidence interval area (Fig. 2d).

262

263 There are uninterpretable Unclassified calls (red X's) at the bottom for high-VAF calls,
264 which is again highly suggestive that the true positive rate for Unclassified calls may be lower
265 than the reported validation rate (25%) for Ion Torrent data as well. The indel equivalent is
266 included in Suppl. Fig. S7b.

267

268 **Independent Confirmation of Call Set with WES on HiSeq**

269

270 We used 14 HiSeq WES replicates from six sequencing centers to evaluate the
271 concordance between these data sets and the WGS data sets employed to construct the truth
272 set. While the WES data sets were not sequenced from orthogonal platforms, they provide
273 insights in terms of the reproducibility of our call sets in different library preparations. The scatter
274 plot between the super set derived VAF and medium HiSeq WES-derived VAF is presented in Fig.
275 2e. Almost all truth set (HighConf and MedConf calls) variants had consistent VAFs calculated
276 from both sources.

277

278 Again, simple rules were implemented for validation with the WES data as well (Table 2).
279 The validation rate for HighConf, MedConf, LowConf, and Unclassified SNV calls by WES were
280 100%, 98.4%, 93.1%, and 42.4%. These validation rates are higher than other methods because
281 these WES data were sequenced on the same platform and sequencing centers as those used to
282 build the truth set. Thus, the truth set variant calls are reproducible in WES, though these data
283 sets do not eliminate sequencing center or platform specific artifacts that may exist in both WGS
284 and WES data sets. The indel equivalent is the subject of Suppl. Fig. S7c.

285

286 **Validation with tumor content titration series**

287

288 To evaluate the effects of tumor purity, we pooled HCC1395 DNA with HCC1395BL DNA
289 at different ratios to create a range of admixtures representing tumor purity levels of 100%, 75%,
290 50%, 20%, 10%, 5%, and 0%. For each tumor DNA dilution point, we performed WGS on a HiSeq
291 4000 with 300× total coverage by combining three repeated runs (manuscript NBT-RA46164).
292 We plotted the VAF fitting score between the expected values based on the super set vs. the
293 observed values at each tumor fraction (Fig. 2f). For real somatic mutations, their observed VAF
294 should scale linearly with tumor fraction in the tumor-normal titration series. In contrast, the
295 observed VAF for sequencing artifacts or germline variants will not scale in this fashion. Fig. 2f
296 shows that the fitting scores for HighConf and MedConf calls are much higher than LowConf and
297 Unclassified calls across all VAF brackets, indicating that the HighConf and MedConf calls contain
298 far more real somatic mutations than LowConf and Unclassified calls. The formula [Eq. 2] for the
299 fitting score is described in the Methods.

300

301 **Definition and Confirmation of Germline SNVs/Indels in matched normal**

302

303 For the 21 WGS sequencing replicates of HCC1395BL (aligned with BWA MEM, Bowtie2,
304 and NovoAlign to create 63 BAM files) we employed four germline variant callers, i.e.,
305 FreeBayes²⁸, Real Time Genomics (RTG)²⁹, DeepVariant³⁰, and HaplotypeCaller³¹, to discover
306 germline variants (SNV/indels). To consolidate all the calls, a generalized linear mixed model

307 (GLMM) was fit for each set of SNV calls which are sequenced at different centers on various
308 replicates, aligned by the three aligners, and discovered by the four callers. We estimated the
309 SNVs/indel call probability (SCP) averaged across four factors (sequencing center, sequencing
310 replicate, aligner, and caller), and examined the variance of SCP across these factors. The SNV
311 candidates considered were called at least four times (out of a maximum of $21 \times 3 \times 4 = 252$ times)
312 by various combination of the four factors. The frequency histogram of the averaged SCPs
313 demonstrates a bimodal pattern (Fig. 3a). The vast majority of SNV calls (97%) had SCPs either
314 below 0.1 (57%) or above 0.9 (40%). Only a small minority of calls (3%) lie between 0.1 and 0.9.
315 This indicates when SNVs were repeatedly sequenced and called, only a small proportion of them
316 would be recurrently called as SNVs, and those recurrent calls were in fact highly recurrent.

317

318 Each of our germline SNV or indel calls had annotated SCP. See the Methods and Eq. 2 for
319 details. Suppl. Table S7 demonstrates that, of the highest-confidence calls (SCP=1, i.e. they were
320 called everywhere), the validation rates were approximately 99% for SNV and 98% for indels by
321 Illumina MiSeq, and 98% and 97% for Ion Torrent. Of the 11 SNV with SCP below 0.5, all were not
322 confirmed by MiSeq. Other calls had intermediate validation rates.

323

324 Figs 3b and 3c display that the vast majority of confirmed germline VAF was around 50%
325 and 100%. A considerable number of lower-confidence germline SNV calls clustered around 20%
326 VAF in non-exonic regions (Fig. 3b), with a large proportion of them being uninterpretable during
327 validation. Scatter plots for indels are qualitatively similar (Suppl. Fig. S13).

328

329 **SNV Functional Relevance and TMB Benchmarks**

330

331 Among the truth set somatic mutations, 186 COSMIC SNVs and 7 COSMIC indels are in
332 the coding region. One hotspot somatic mutation of particular biological significance is a *TP53*
333 *c.128G>A* (COSMIC99023, chr17:7675088 C>T, VAF>99%), which causes an amino acid change
334 *p.Arg43His* that leads to the inactivation of *TP53* tumor suppressive function³². In addition,
335 there is also a stop gain mutation in *BRCA2 c.4777G>T* (COSMIC13843, chr13:32339132 G>A),
336 which causes a nonsense at *p.Glu1593**, though it is only a heterozygous variant with VAF of
337 37.5%. Furthermore, there is a missense mutation in *FGFR1 c.473C>T* (COSM1456963, chr8:
338 38428420 G>A, VAF>99%).

339

340 Of the over 3.5 million high-confidence germline variants discovered in HCC1395BL,
341 9,016 of them are in the ClinVar database. Most of them were annotated as “benign” or “like
342 benign”; however, 8 SNVs were annotated as “pathogenic” (Suppl. Table S9). One germline
343 variant likely to substantially increase the risk of an affected patient to develop breast cancer is
344 a premature stop gain in *BRCA1* (chr17:43057078, *c.5251C>A*, *p.Arg1751**, ClinVar #55480,
345 OMIM Entry #604370). The lifetime risk of breast cancer for carriers of this variant is 80 to
346 90%³³. The premature stop codon deactivates *BRCA1*'s function to repair DNA double-strand
347 breaks. It is one of the most common germline variants among breast cancer patients. HCC1395
348 has both *BRCA1* and *TP53* completely inactivated, one from germline and one acquired
349 somatically. The loss of two critical tumor suppressor genes likely contributed to tumorigenesis.
350 A full list of COSMIC somatic mutations and ClinVar germline variants in the coding region is

351 provided in Supplemental File 2.

352

353 Tumor mutational burden (TMB) is defined as the number of non-synonymous somatic
354 variants per unit area of the genome, i.e., typically the number of non-synonymous mutations
355 per Mb³⁴. Recent literature increasingly has reported correlations between TMB and
356 response to anti-PD(L)-1 immunotherapy treatment³⁵. The “gold standard” to measure TMB is
357 to perform tumor-normal WES and find the total number of non-synonymous mutations (all in
358 the coding regions). Due to the high cost and time required for WES, researchers are trying to
359 infer TMB with much smaller and less expensive targeted oncology panels. One way to increase
360 the statistical power of a much smaller panel is to measure all somatic mutations, including
361 synonymous mutations, which is expected to correlate highly with frequency of non-
362 synonymous mutations if we believe most somatic mutations, especially in high-TMB patients,
363 occur more-or-less randomly. We inferred TMB with various commercially available target
364 panels. The uncertainties of mutation rate (calculated as the 95% binomial confidence interval)
365 inferred by smaller oncology panels are quite large, so we advise caution when attempting to
366 infer TMB from targeted oncology panels (Suppl. Table S10).

367

368 **Defining Genome Callable Regions**

369

370 Accurate variant calling requires an abundance of high-quality reads aligned accurately to
371 the genomic coordinates in question. False positives are overwhelmingly enriched in genomic
372 regions where the alignments are challenging, base call qualities are low, and/or reported
373 coverage is far from the mean or median³⁶. There are parts of the human genome that cannot be
374 covered by current technologies (Fig. 4a). To obtain the callable regions, we ran GATK CallableLoci
375 on each of the 63 HCC1395 and HCC1395BL BAM files to identify regions of low coverage (<10),
376 ultra-high coverage (8× the mean coverage of the sample), difficult to map (MQ<20), poor
377 reads (Base Quality Score BQ<20), or with N in the reference genome. We then created
378 consensus callable regions that we deemed callable for our truth set. A limitation of our callable
379 regions and our truth set is that they were defined and relied on short-read sequencing
380 technologies (i.e., Illumina sequencers), because currently only high-accuracy short-read
381 technologies are fit for somatic variant detection due to their low VAF. Variant calls outside the
382 consensus callable regions were labeled NonCallable in the super set and truth set to warn users
383 of these potential problems (details in Methods). NonCallable regions consisted of approximately
384 8% of the genome but contained over 34% of all Unclassified calls and 23% of all LowConf calls in
385 the super set (Suppl. Table S6).

386

387 The consensus callable regions consist of a total of 2.73 billion bps (Fig. 4b). In comparison
388 with GiaB NA12878 genome’s more strictly defined high-confidence (HC) regions⁷, 88% of our
389 consensus callable regions are in common with GiaB’s HC regions. On the other hand, 98% of
390 GiaB’s HC regions are a part of our consensus callable regions. Unlike GiaB’s HC which exclude
391 regions with structural variations as well as regions where variant calls are inconsistent with
392 pedigree or regions with unexplained pipeline inconsistencies, when there were disagreements
393 in a variant call from various sequencing data, we did not exclude the region. Instead, we
394 attempted to resolve these discrepancies. When there were nearby structural changes, we relied

395 on machine learning algorithms to resolve these challenging events. As a result, our consensus
396 callable regions included some difficult genomic regions, such as human leukocyte antigen (HLA)
397 and olfactory receptor genes which contain high homologous sequences. The confidence (or the
398 lack thereof) we hold for each variant call is annotated on a per call basis. We have demonstrated
399 some benchmarking results with different regions in the Supplementary, section 1.10.

400
401

402 Discussion

403

404 Through a community effort, we generated a high confidence somatic mutation call set
405 with limit of detection (LOD) at 5% VAF (Fig. 2b). To employ as an accuracy benchmark, we
406 recommend considering the variant calls labeled with both HighConf and MedConf as true
407 positives. These true positive variants can be used to assess sensitivity, i.e., the fraction of those
408 variants detected by a pipeline. On the other hand, variant calls labeled as Unclassified plus any
409 unspecified genomic coordinates are likely false positives. LowConf calls could not be confidently
410 determined here and should be blacklisted for current accuracy evaluation. LowConf calls had
411 validation rates around 50%, and often had VAF below our 50× depth detection limit. They
412 represent opportunities for future work to ascertain their actual somatic status.

413

414 The confidence level of each variant call was determined by the “PASS” classifications
415 provided by SomaticSeq across different sequencing centers with different aligners (see
416 Methods). If a variant was not detected by any caller in a data set, it was considered “Missing” in
417 that data set, which is common for low-VAF calls due to stochastic sampling. For most calls,
418 however, they either had “PASS” classifications or “REJECT or Missing” classifications, but not
419 both. Few variant candidates had a large number of “PASS *and* REJECT” classifications (Suppl. Fig.
420 S6a). HighConf calls had many “PASS” classifications, very few “REJECT” classifications, and a full
421 range of VAFs. MedConf calls had fewer “PASS” calls (still high), still very few “REJECT”
422 classifications, but were mostly low-VAF, which explains the lower number of “PASS” calls.
423 LowConf calls had even fewer “PASS” calls than MedConf though they overlapped significantly,
424 and also a low number of “REJECT” classifications. LowConf calls tended to have even lower VAF
425 than MedConf, around or below our detection limit (Fig. 2b). Only Unclassified calls suffered a
426 significant number of “REJECT” classifications, and they also displayed a full range of VAF. The
427 performance of Unclassified calls indicated that SomaticSeq labeled them “REJECT” due to poor
428 mapping, poor alignment, germline risk, or causes other than lack of variant reads. HighConf and
429 Unclassified calls are far apart in all of the metrics described above.

430

431 Variant re-sequencing with AmpliSeq (Suppl. Fig. S6c) pointed to a high validation rate for
432 HighConf and MedConf calls. Suppl. Fig. S6c also contains a cluster of Unclassified and LowConf
433 calls in the middle of the XY plane, representing calls with some conflicts (i.e., large number of
434 “PASS” and “REJECT” calls).

435

436 Each time a human cell divides, somatic mutations could be introduced by replication
437 errors. Somatic mutations can occur much more frequently in cancer cells with malfunctioning
438 DNA repair systems. It is not feasible to detect extremely low-VAF somatic mutations because

439 they may appear in few tumor cells. Our "truth set" for somatic mutation was built upon WGS
440 with 50×-100X coverages, and thus it was designed to detect somatic mutations limited to 5% of
441 VAF. Variants with low-VAF ($\leq 12\%$) were cross-referenced with two data sets with depths over
442 300× to ascertain their presence. While we do not expect our truth set to be 100% accurate or
443 100% comprehensive, the AmpliSeq and Ion Torrent data sets demonstrated combined 99% and
444 91% validation rates for HighConf and MedConf SNV calls, respectively. AmpliSeq also showed a
445 94% validation rate for HighConf indel calls. VAF of 5% represents the lower detection limit of the
446 first release of the somatic mutation truth set, even though there are many true mutations with
447 VAF under that threshold. We recommend that if using this truth set as a benchmark, novel
448 variant calls (i.e., variants calls not present in our super set) with VAF<5% should be blacklisted
449 from the accuracy calculations because we cannot confidently determine their status. Due to
450 losses of chr6p, chr16q, and chrX in HCC1395BL (Suppl. Fig S1, S2), somatic mutations in these
451 regions were excluded.

452
453 For the first time, tumor-normal paired "reference samples" with a whole-genome
454 characterized somatic mutation and germline "truth sets" are available to the community. Our
455 samples, data sets, and the list of known somatic mutations can serve as a public resource for
456 evaluating NGS platforms and pipelines. The massive and diverse amount of sequencing data
457 generated from multiple platforms at multiple sequencing centers can help tool developers to
458 create and validate new algorithms and to build more accurate artificial intelligence (AI) models
459 for somatic mutation detection. The reference samples and call set presented here can help in
460 assay development, qualification, validation, and proficiency testing. Such community defined
461 tumor-normal paired reference samples can be helpful in quality assessment by clinical
462 laboratories engaged in NGS, data exchange between laboratories, characterization of gene
463 therapy products, and premarket review of NGS-based products. Furthermore, the
464 methodology used in this study can be extended to establish truth sets for additional cancer
465 reference samples. Other reference sample efforts may be able to build on the data sets we
466 established or consider using these samples as a genomic background for other reference
467 samples.

468
469

470 **Methods**

471 See Online Methods

472

473 **Acknowledgements**

474 We thank Justin Zook of the National Institute of Standard Technology for advice in establishing
475 reference samples and truth set; Sivakumar Gowrisankar of Novartis, Susan Chacko of the Center
476 for Information Technology, the National Institute of Health for their assistance with data
477 transfer; Dr. Jun Ye of Sentieon for providing the Sentieon software package. We also appreciate
478 Dr. Laufey Amundadottir of the Division of Cancer Epidemiology and Genetics, National Cancer
479 Institute (NCI), National Institutes of Health (NIH), for the sponsorship and the usage of the NIH
480 Biowulf cluster; Drs Reena Phillip, Yun-Fu Hu, Sharon Liang, and You Li of the Center for Devices
481 and Radiological Health, U.S. Food and Drug Administration, for their advices on study design and

482 manuscript writing; and Seven Bridges for providing storage and computational support on the
483 Cancer Genomic Cloud (CGC). The CGC has been funded in whole or in part with Federal funds
484 from the National Cancer Institute, National Institutes of Health, Contract No.
485 HHSN261201400008C and ID/IQ Agreement No. 17X146 under Contract No.
486 HHSN261201500003I. Chunlin Xiao and Steve Sherry were supported by the Intramural Research
487 Program of the National Library of Medicine, National Institutes of Health. This work also used
488 the computational resources of the NIH Biowulf cluster (<http://hpc.nih.gov>). Original data was
489 also backed up on the servers provided by Center for Biomedical Informatics and Information
490 Technology (CBIT), NCI. We would also like to thank the partially support from the Ardmore
491 Institute of Health (AIH) and Dr. Charles A. Sims' gift to Loma Linda University (LLU) Center for
492 Genomics. The LLU Center for Genomics is partially supported by AIH grant (2150141) and Charles
493 A. Sims' gift.
494

495 **Disclaimer**

496 This is a research study, not intended to guide clinical applications. The views presented in this
497 article do not necessarily reflect current or future opinion or policy of the US Food and Drug
498 Administration. Any mention of commercial products is for clarification and not intended as
499 endorsement.
500

501 **Data availability**

502 All raw data (FASTQ files) are available on NCBI's SRA database (SRP162370). The truth set for
503 somatic mutations in HCC1395, VCF files derived from individual WES and WGS runs, and
504 source codes are available on NCBI's ftp site ([ftp://ftp-
505 trace.ncbi.nlm.nih.gov/seqc/ftp/Somatic_Mutation_WG/](ftp://ftp-trace.ncbi.nlm.nih.gov/seqc/ftp/Somatic_Mutation_WG/)). Some alignment files (BAM) are also
506 available on Seven Bridges' s Cancer Genomics Cloud (CGC) platform.
507

508 **Author contributions**

509 Study conceived and designed by: W.X., C.W., L.S., H.H., E.D., Z.T., B.N., W.T., R.J.
510

511 Biosample preparation: L.K., K.L., M. M., T.H.
512

513 NGS library preparation and sequencing: W.C., Z.C., S.S., K.I., H.J., E.J., G.S., S.S., J.S., P.K., J.B.,
514 B.T., V.P., M.S., T.H., E.P., R.K., J.D., P.V., R.M., D.G., S.K., E.R., A.S., J.N., U.L., J.W., J.L., P.PH.
515

516 Data analysis: L.T.F., W.X., B.Z., Y.Z., Z.Y., C.L., O.A., L.S., J.L., T.S., K.T., D.M., C.N., M.C., S.M.S.,
517 M.M., Y.G., L.Y., H.L., M.P., Z.L.
518

519 Data management: W.X., C.X., S. S.
520

521 Manuscript writing: L.T.F., W.X., R.K., M.M., C.X., S.S.
522

523 Project management: W.X.

524 **References**

- 525
- 526 1. Hofmann, A. L. *et al.* Detailed simulation of cancer exome sequencing data reveals
- 527 differences and common limitations of variant callers. *BMC Bioinformatics* **18**, 8 (2017).
- 528 2. Krøigård, A. B., Thomassen, M., Lænkholm, A.-V., Kruse, T. A. & Larsen, M. J. Evaluation
- 529 of Nine Somatic Variant Callers for Detection of Somatic Mutations in Exome and Targeted
- 530 Deep Sequencing Data. *PLOS ONE* **11**, e0151664 (2016).
- 531 3. Shi, W. *et al.* Reliability of Whole-Exome Sequencing for Assessing Intratumor Genetic
- 532 Heterogeneity. *Cell Rep.* **25**, 1446–1457 (2018).
- 533 4. Tezak, Z. & Berger, A. Considerations for Design, Development, and Analytical Validation
- 534 of Next Generation Sequencing (NGS) – Based In Vitro Diagnostics (IVDs) Intended to Aid
- 535 in the Diagnosis of Suspected Germline Diseases. (2018).
- 536 5. Eberle, M. A. *et al.* A reference data set of 5.4 million phased human variants validated by
- 537 genetic inheritance from sequencing a three-generation 17-member pedigree. *Genome Res.*
- 538 **27**, 157–164 (2017).
- 539 6. The International HapMap Project. *Nature* **426**, 789 (2003).
- 540 7. Zook, J. M. *et al.* Integrating human sequence data sets provides a resource of benchmark
- 541 SNP and indel genotype calls. *Nat. Biotechnol.* **32**, 246–251 (2014).
- 542 8. Zook, J. M. *et al.* An open resource for accurately benchmarking small variant and reference
- 543 calls. *Nat. Biotechnol.* **1** (2019). doi:10.1038/s41587-019-0074-6
- 544 9. National Institute for Biological Standards and Control. *WHO Reference Panel 1st*
- 545 *International Reference Panel for genomic KRAS codons 12 and 13 mutations NIBSC code:*
- 546 *16/250.* (World Health Organization).
- 547 10. Oncospan Reference Standard HD827. Available at:

- 548 <https://www.horizondiscovery.com/reference-standards/type/oncospan>. (Accessed: 17th
549 April 2019)
- 550 11. Craig, D. W. *et al.* A somatic reference standard for cancer genome sequencing. *Sci. Rep.* **6**,
551 24607 (2016).
- 552 12. Zehnbauer, B. *et al.* MDIC SRS Report: Somatic Variant Reference Samples for NGS.
553 (2019).
- 554 13. Popova, T. *et al.* Ploidy and Large-Scale Genomic Instability Consistently Identify Basal-
555 like Breast Carcinomas with BRCA1/2 Inactivation. *Cancer Res.* **72**, 5454–5462 (2012).
- 556 14. Gazdar, A. F. *et al.* Characterization of paired tumor and non-tumor cell lines established
557 from patients with breast cancer. *Int. J. Cancer* **78**, 766–774 (1998).
- 558 15. Kalyana-Sundaram, S. *et al.* Gene Fusions Associated with Recurrent Amplicons Represent
559 a Class of Passenger Aberrations in Breast Cancer. *Neoplasia* **14(8)**, 702–708 (2012).
- 560 16. Zhang, J. *et al.* INTEGRATE: gene fusion discovery using whole genome and transcriptome
561 data. *Genome Res.* **26**, 108–118 (2016).
- 562 17. Robinson, D. R. *et al.* Functionally recurrent rearrangements of the MAST kinase and Notch
563 gene families in breast cancer. *Nat. Med.* **17**, 1646–1651 (2011).
- 564 18. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
565 *ArXiv13033997 Q-Bio* (2013).
- 566 19. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,
567 357–359 (2012).
- 568 20. NovoCraft Technologies Sdn Bhd. NovoAlign. Available at:
569 <http://www.novocraft.com/products/novoalign/>.
- 570 21. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and

- 571 heterogeneous cancer samples. *Nat. Biotechnol.* **31**, 213–219 (2013).
- 572 22. Larson, D. E. *et al.* SomaticSniper: identification of somatic point mutations in whole
573 genome sequencing data. *Bioinformatics* **28**, 311–317 (2012).
- 574 23. Lai, Z. *et al.* VarDict: a novel and versatile variant caller for next-generation sequencing in
575 cancer research. *Nucleic Acids Res.* **44**, e108 (2016).
- 576 24. Fan, Y. *et al.* MuSE: accounting for tumor heterogeneity using a sample-specific error model
577 improves sensitivity and specificity in mutation calling from sequencing data. *Genome Biol.*
578 **17**, 178 (2016).
- 579 25. Kim, S. *et al.* Strelka2: fast and accurate calling of germline and somatic variants. *Nat.*
580 *Methods* **15**, 591 (2018).
- 581 26. Freed, D., Pan, R. & Aldana, R. TNscope: Accurate Detection of Somatic Mutations with
582 Haplotype-based Variant Candidate Detection and Machine Learning Filtering. *bioRxiv*
583 250647 (2018). doi:10.1101/250647
- 584 27. Fang, L. T. *et al.* An ensemble approach to accurately detect somatic mutations using
585 SomaticSeq. *Genome Biol.* **16**, 197 (2015).
- 586 28. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing.
587 *ArXiv12073907 Q-Bio* (2012).
- 588 29. Real Time Genomics (RTG) Variant Caller. Available at:
589 <https://www.realtimengenomics.com/>.
- 590 30. Poplin, R. *et al.* A universal SNP and small-indel variant caller using deep neural networks.
591 *Nat. Biotechnol.* **36**, 983–987 (2018).
- 592 31. Poplin, R. *et al.* Scaling accurate genetic variant discovery to tens of thousands of samples.
593 *bioRxiv* 201178 (2017). doi:10.1101/201178

- 594 32. Soussi, T., Leroy, B. & Taschner, P. E. M. Recommendations for Analyzing and Reporting
595 TP53 Gene Variants in the High-Throughput Sequencing Era. *Hum. Mutat.* **35**, 766–778
596 (2014).
- 597 33. OMIM Clinical Synopsis - #604370 - BREAST-OVARIAN CANCER, FAMILIAL,
598 SUSCEPTIBILITY TO, 1; BROVCA1. Available at:
599 <https://www.omim.org/clinicalSynopsis/604370>. (Accessed: 22nd March 2019)
- 600 34. Meléndez, B. *et al.* Methods of measurement for tumor mutational burden in tumor tissue.
601 *Transl. Lung Cancer Res.* **7**, 661–667 (2018).
- 602 35. Goodman, A. M. *et al.* Tumor Mutational Burden as an Independent Predictor of Response
603 to Immunotherapy in Diverse Cancers. *Mol. Cancer Ther.* **16**, 2598–2608 (2017).
- 604 36. Li, H. Toward better understanding of artifacts in variant calling from high-coverage
605 samples. *Bioinformatics* **30**, 2843–2851 (2014).

606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627

NGS technologies		platforms	# of reads (coverage)	
			HCC1395	HCC1395BL
Initial	WGS	HiSeq	57 billion (2,800X)	57 billion (2,800X)
		NovaSeq	13 billion (650X)	13 billion (650X)
		10X Genomics	20 billion (1,000X)	20 billion (1,000X)
		PacBio	20 million (50X)	20 million (50X)
Validation	WGS-tumor content	HiSeq	7.6 billion (380X)	7.6 billion (380X)
	WES	HiSeq	5 billion (12,500X)	5 billion (12,500X)
		Ion Torrent	67 million (34X)	82 million (47X)
	AmpliSeq	MiSeq	3.3 million (2000X)	3.3 million (2000X)

628 **Table 1.** Massive data from multiple NGS platforms was obtained to derive and confirm
 629 germline and somatic variants in HCC1395 and HCC1395BL
 630

Validation Platform	Variant Type	Category	Total Number	Fraction Interpretable	Validation Rate (Interpretable)	Validation Rate (Total)
AmpliSeq Deep Sequencing	SNV	HighConf	247	(237/247) 96.0%	(234/237) 98.7%	94.7%
		MedConf	40	(37/40) 92.5%	(34/37) 91.9%	85.0%
		LowConf	58	(41/58) 70.7%	(22/41) 53.7%	37.9%
		Unclassified	105	(62/105) 59.0%	(7/62) 11.3%	6.7%
	INDEL	HighConf	17	(17/17) 100.0%	(16/17) 94.1%	94.1%
		MedConf	2	(2/2) 100.0%	(2/2) 100.0%	100.0%
		LowConf	1	(0/1) 0.0%	(0/0) NA	nan%
		Unclassified	1	(1/1) 100.0%	(0/1) 0.0%	0.0%
Ion Torrent WES	SNV	HighConf	703	(629/703) 89.5%	(623/629) 99.0%	88.6%
		MedConf	43	(27/43) 62.8%	(24/27) 88.9%	55.8%
		LowConf	134	(39/134) 29.1%	(28/39) 71.8%	20.9%
		Unclassified	802	(155/802) 19.3%	(39/155) 25.2%	4.9%
	INDEL	HighConf	31	(25/31) 80.6%	(22/25) 88.0%	71.0%
		MedConf	15	(7/15) 46.7%	(6/7) 85.7%	40.0%
		LowConf	8	(0/8) 0.0%	(0/0) NA	nan%
		Unclassified	36	(8/36) 22.2%	(6/8) 75.0%	16.7%
WES	SNV	HighConf	1074	(1068/1074) 99.4%	(1068/1068) 100%	99.4%
		MedConf	64	(63/64) 98.4%	(62/63) 98.4%	96.9%
		LowConf	197	(144/197) 73.1%	(134/144) 93.1%	68.0%
		Unclassified	1218	(436/1218) 35.8%	(184/436) 42.4%	15.1%
	INDEL	HighConf	45	(43/45) 95.6%	(43/43) 100%	95.6%
		MedConf	17	(17/17) 100.0%	(17/17) 100%	100.0%
		LowConf	13	(10/13) 76.9%	(9/10) 90%	69.2%
		Unclassified	54	(19/54) 35.2%	(14/19) 73.7%	25.9%

631 **Table 2:** Validation of SNVs of different confidence levels by three different methods
 632

633 **Figure legends**

634

635 **Figure 1:** Schematic of the bioinformatics pipeline used to define the confidence levels of the
636 super set and truth set (see Online Methods for detail)

637

638 **Figure 2:** Initial definition of somatic mutation truth set and subsequent validation. (a) A
639 breakdown of the four confidence levels in the super set. (b) Histograms of VAF for SNVs (top)
640 and Indels (bottom) calls. (c) Validation of initial definition of somatic mutation truth set with
641 AmpliSeq. Solid circles are variant calls that were positively confirmed. Open circles are variants
642 that were not confirmed. X's are when validation data were deemed uninterpretable due to
643 low depth or unclear signal. The dashed lines at the diagonal represent the 95% binomial
644 confidence-interval of observed VAF given the actual VAF, calculated based on 2000× depth
645 for AmpliSeq. The figure shows very high correlation between VAF estimated from super set
646 data and validation data for HighConf calls ($R=0.958$). Many Unclassified data points lie at the
647 bottom, implying that those calls were not real mutations despite the large number of apparent
648 variant-supporting reads in the super set data. X-axis: VAF calculated from the super set. Y-axis:
649 VAF calculated from AmpliSeq data. (d) Validation of the initial definition of the somatic
650 mutation truth set with Ion Torrent WES. The 95% binomial confidence-interval dash lines were
651 calculated based on 34× depth for Ion Torrent. $R=0.928$ for HighConf calls. (e) Validation of
652 initial definition of somatic mutation truth set with 12 repeats of WES on the HiSeq platform. Y-
653 axis: median VAF calculated based on 12 HiSeq WES replicates. The 95% binomial confidence-
654 interval dashed lines were calculated based on 150× depth for HiSeq WES. $R=0.992$ for
655 HighConf calls. (f) Average tumor purity fitting scores for the VAF of each SNV across the four
656 different confidence levels vs. the observed VAF in the tumor-normal titration series. The
657 formula for fitting scores is described in Eq. 1 in the Online Methods.

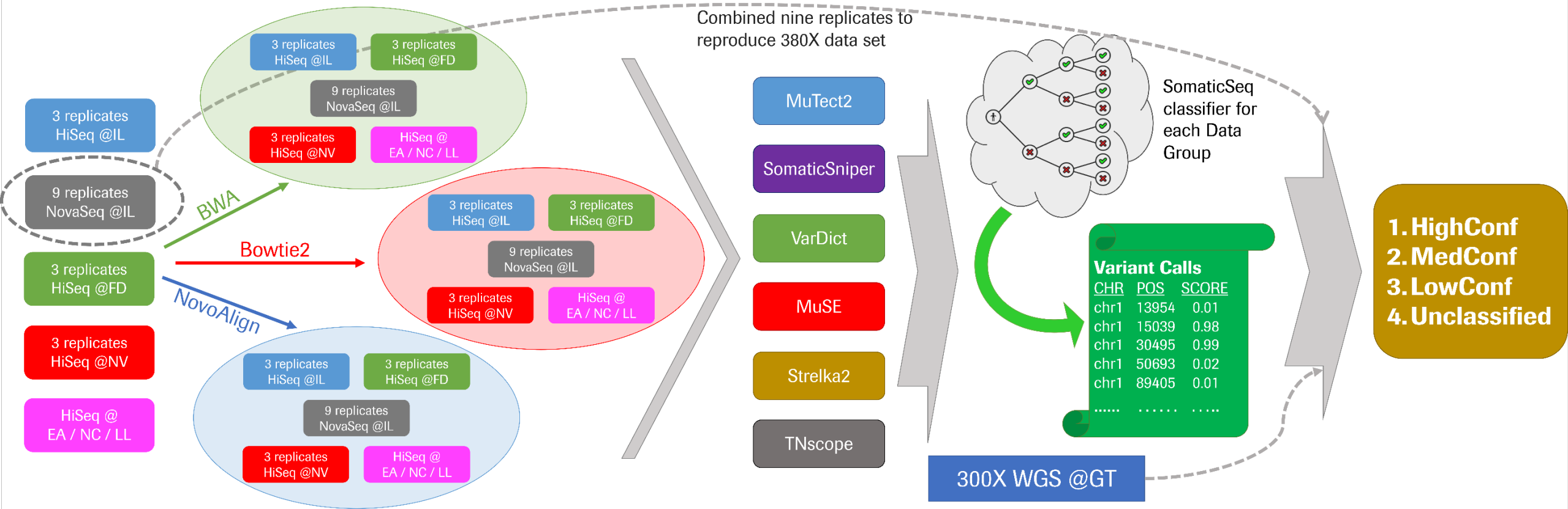
658

659 **Figure 3:** Initial definition of germline variants and validation. (a) Histogram of SNV call
660 probability for germline variants identified by four callers from 63 BAM files. (b) VAF scatter
661 plot of germline SNVs by the truth set and AmpliSeq. $R=0.986$ for SCP=1 calls. (c) VAF scatter
662 plot of germline SNVs by the truth set and Ion Torrent WES. $R=0.758$ for SCP=1 calls.

663

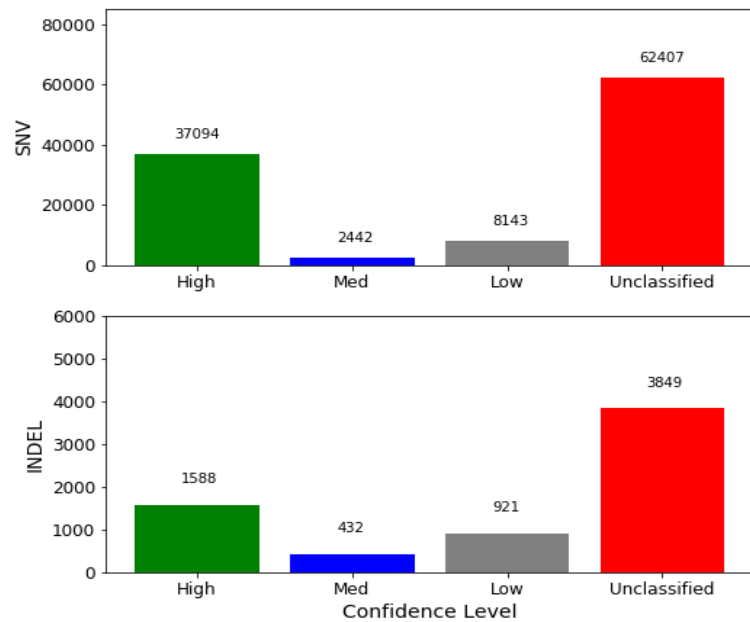
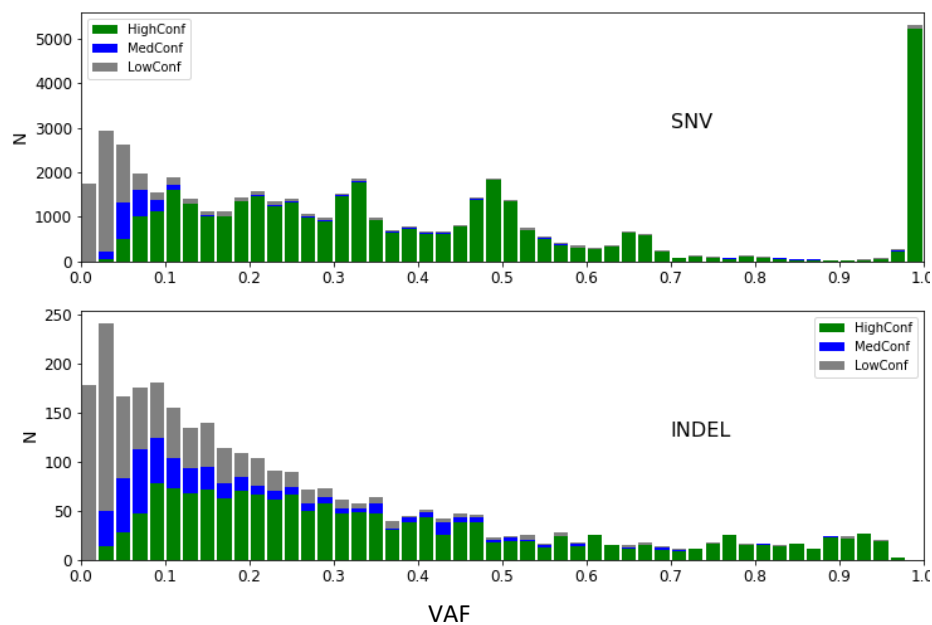
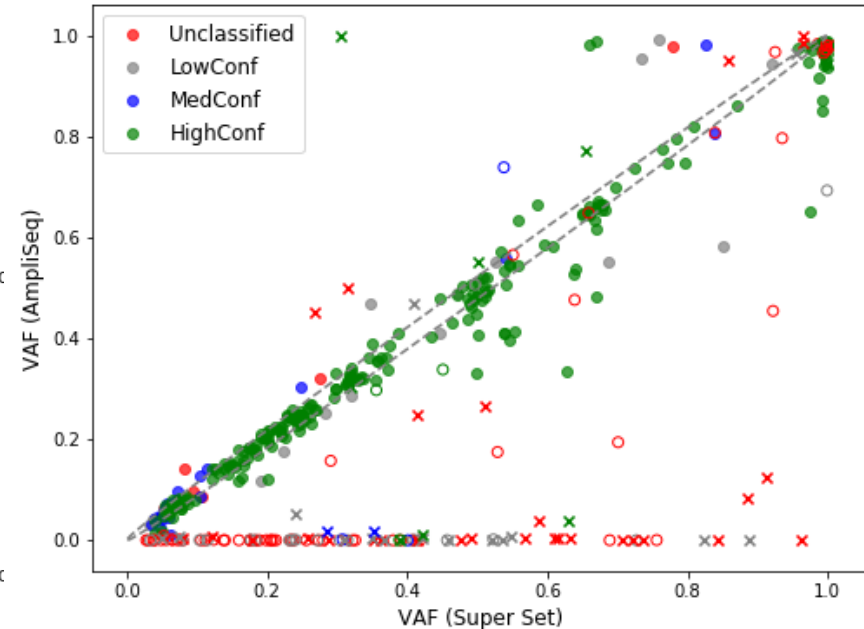
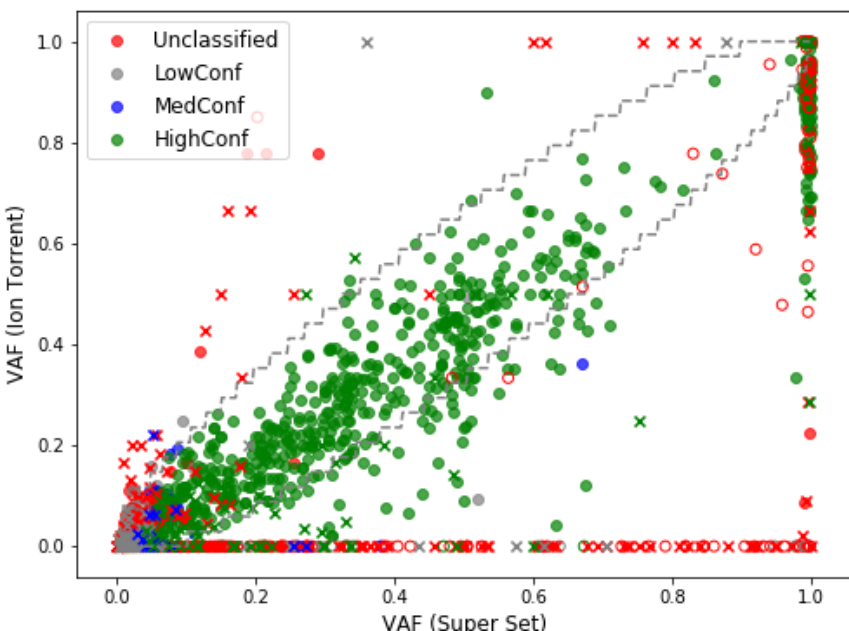
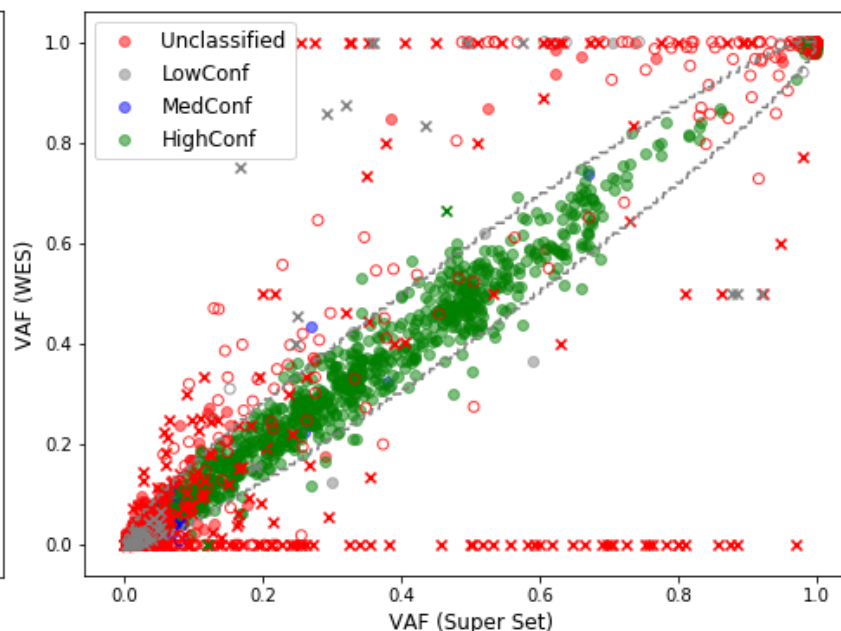
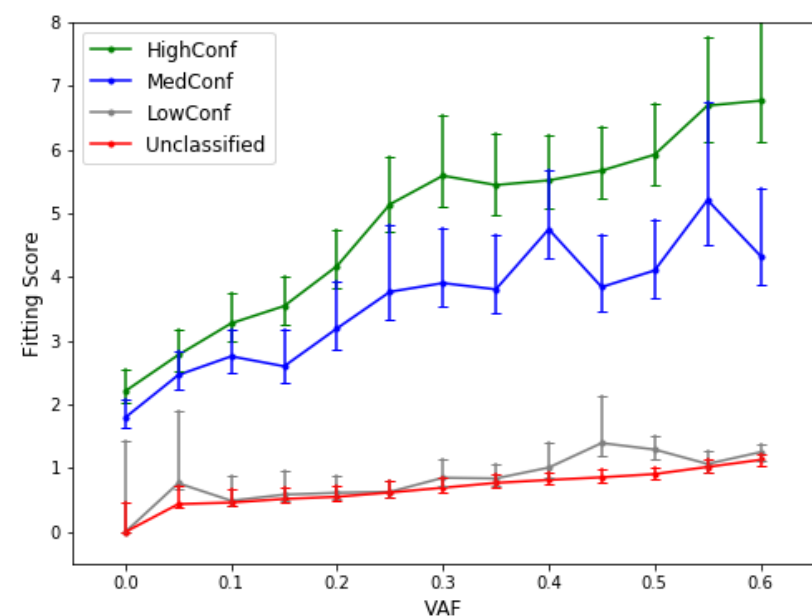
664 **Figure 4:** Genome coverage and high-confidence regions on reference genome GRCh38. a)
665 Genome coverage comparison between three technologies. Inner track: PacBio. Middle track:
666 10X Genomics. Outer track: Illumina. Red line: HCC1395. Green line: HCC1395BL. b) Genome
667 regions coverage by Illumina short reads in comparison to NA12878. Inner track: NA12878.
668 Middle track: the callable regions in HCC1395 and HCC1395BL. Outer track: gene density

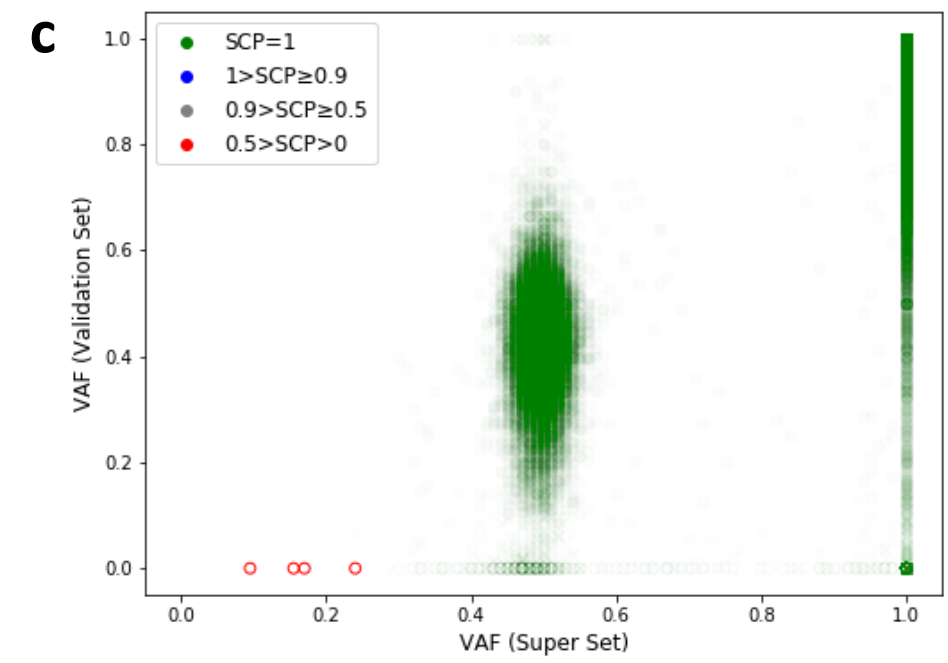
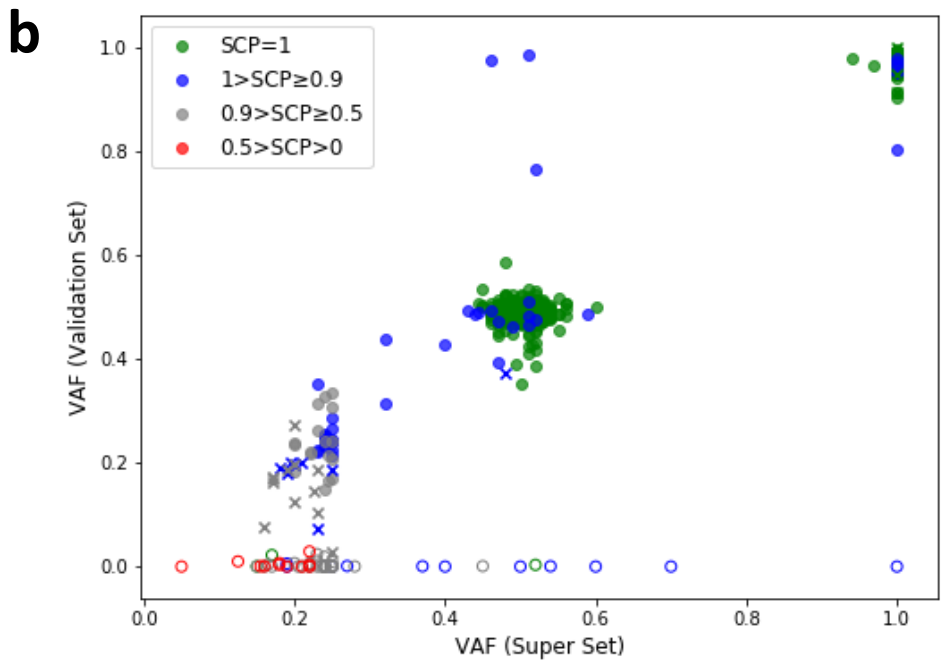
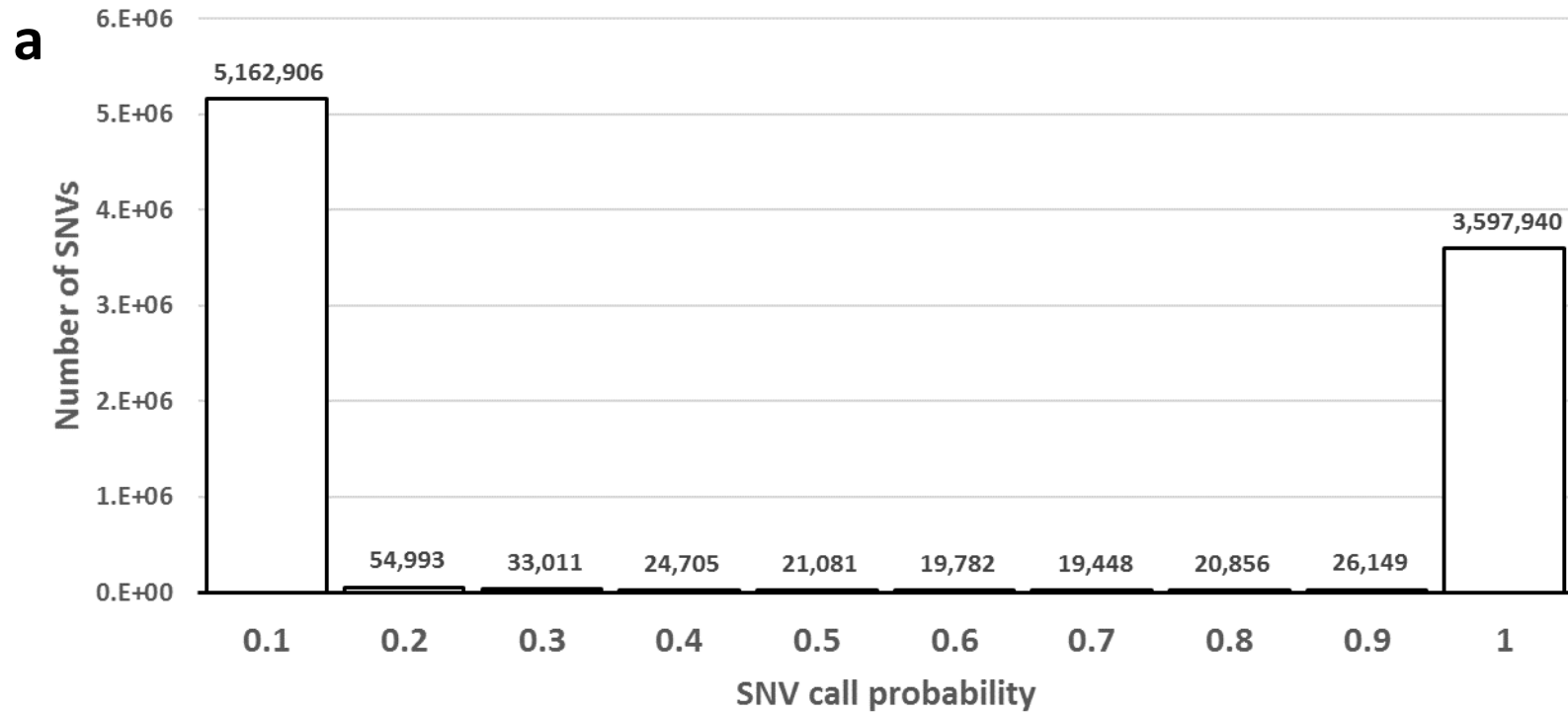
669

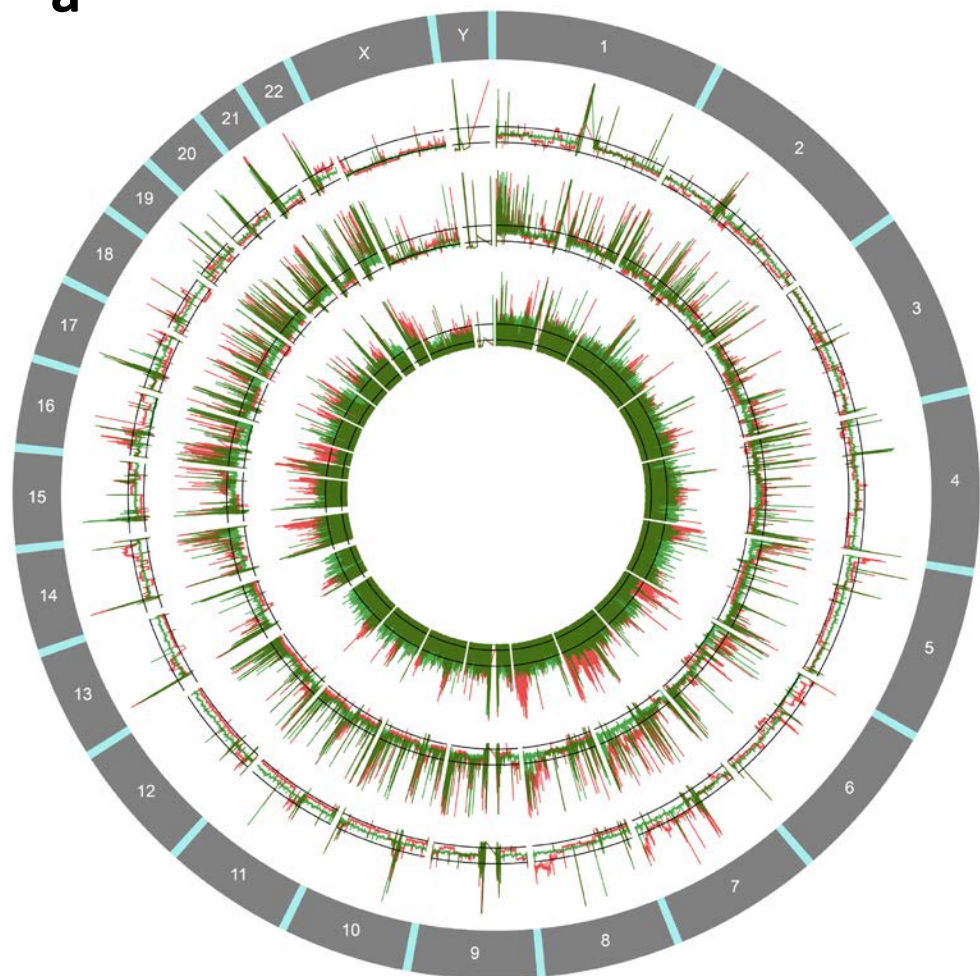


Variant Calls

CHR	POS	SCORE
chr1	13954	0.01
chr1	15039	0.98
chr1	30495	0.99
chr1	50693	0.02
chr1	89405	0.01
.....

a**b****c****d****e****f**



a**b**