

1 **NextSV: a computational pipeline for structural variation analysis**
2 **from low-coverage long-read sequencing**

3
4
5
6
7
8
9

Li Fang^{1,2}, Jiang Hu¹, Depeng Wang¹, Kai Wang^{2,3,*}

10 1: Grandomics Biosciences, Beijing 102206, China
11 2: Institute for Genomic Medicine, Columbia University Medical Center, New York, NY 10032,
12 USA
13 3: Department of Biomedical Informatics, Columbia University Medical Center, New York, NY
14 10032, USA

15
16 *Please address correspondence to:

17 Kai Wang, Ph.D.
18 622 W 168th St,
19 Room PH20-313,
20 New York, NY 10032
21 Tel: 212-305-3234
22 Email: kaichop@gmail.com

23
24

25 Abbreviations: ADI, allele drop-in; AJ, Ashkenazi Jewish; SV, structural variants.

26 **Abstract**

27

28 Structural variants (SVs) in human genomes are implicated in a variety of human diseases. Long-
29 read sequencing delivers much longer read lengths than short-read sequencing and may greatly
30 improve SV detection. However, due to the relatively high cost of long-read sequencing, users are
31 often faced with issues such as what coverage is needed and how to optimally use the aligners and
32 SV callers. Here, we developed NextSV, a meta SV caller and a computational pipeline to perform
33 SV calling from low coverage long-read sequencing data. NextSV integrates three aligners and
34 three SV callers and generates two integrated call sets (sensitive / stringent) for different analysis
35 purpose. We evaluated SV calling performance of NextSV under different PacBio coverages on
36 two personal genomes, NA12878 and HX1. Our results showed that, compared with running any
37 single SV caller, NextSV stringent call set had higher precision and balanced accuracy (F1 value)
38 while NextSV sensitive call set had a higher recall. At 10X coverage, the recall of NextSV sensitive
39 call set was 93.5%~94.1% for deletions and 87.9%~93.2% for insertions, indicating that ~10X
40 coverage might be an optimal coverage to use in practice, considering the balance between the
41 sequencing costs and the recall rates. We further evaluated the Mendelian errors on an Ashkenazi
42 Jewish trio dataset. Our results provide useful guidelines for SV detection from low coverage
43 whole-genome PacBio data and we expect that NextSV will facilitate the analysis of SVs on long-
44 read sequencing data.

45

46 **Keywords**

47 long-read sequencing, structure variation, structural variants, low coverage, PacBio

48

49

50 **1. Introduction**

51

52 Structural variants (SVs), including large variations such as deletions, insertions, duplications,
53 inversions, and translocations, play important roles in human diversity and disease susceptibility
54 (Feuk et al., 2006; Pang et al., 2010). Many inherited diseases and cancers have been associated
55 with a large number of SVs in recent years (Carvalho and Lupski, 2016; Moncunill et al., 2014;
56 Stankiewicz and Lupski, 2010; Weischenfeldt et al., 2013; Yang et al., 2013; Zhang et al., 2009).
57 Recent advances in next-generation sequencing (NGS) technologies have facilitated the analysis
58 of variations such as SNPs and small indels in unprecedented details, but the discovery of SVs
59 using short-read sequencing still remains challenging (English et al., 2015). Single-molecule, real-
60 time (SMRT) sequencing developed by Pacific Biosciences (PacBio) offers a long read length,
61 making it potentially well-suited for SV detection in personal genomes (Chaisson et al., 2015;
62 English et al., 2015). Most recently, Merker et al. reported the application of low coverage whole
63 genome PacBio sequencing to identify pathogenic structural variants from a patient with
64 autosomal dominant Carney complex, for whom targeted clinical gene testing and whole genome
65 short-read sequencing were both negative (Merker et al., 2016). This represents an clear example
66 that long-read sequencing may solve some negative cases in clinical diagnostic settings.

67

68 Two popular SV software tools have been developed specifically for long-read sequencing:
69 PBHoney (English et al., 2014) and Sniffles (<https://github.com/fritzsedlazeck/Sniffles>). PBHoney
70 identifies genomic variants via two algorithms, long-read discordance (PBHoney-Spots) and
71 interrupted mapping (PBHoney-Tails). Sniffles is a SV caller written in C++ and it detects SVs
72 using evidence from split-read alignments, high-mismatch regions, and coverage analysis.
73 PBHoney uses bam files generated by BLASR (Chaisson and Tesler, 2012) as input while Sniffles
74 requires BAM files from BWA-MEM (Li, 2013) or NGMLR (Rescheneder et al., 2016), a new
75 long-read aligner. Due to the relative high cost of PacBio sequencing, users are often faced with
76 issues such as what coverage is needed and how to get the best use of the available aligners and
77 SV callers. In addition, it is unclear which software performs the best in low-coverage settings,
78 and whether the combination of software tools can improve performance of SV calls. Finally, the
79 execution of these software tools is often not straightforward and requires careful re-
80 parameterization given specific coverage of the source data.

81

82 To address these challenges, we developed NextSV, an automated SV detection pipeline
83 integrating multiple tools. NextSV automatically execute these software tools with optimized
84 parameters for the specific coverage that user specified, then integrates results of each caller and
85 generates a sensitive call set and a stringent call set, for different analysis purpose.

86

87 Recently, the Genome in a Bottle (GIAB) consortium and the 1000 Genome Project Consortium
88 released high-confidence SV calls for the NA12878 genome, an extensively sequenced genome
89 by different platforms, enabling benchmarking of SV callers (Parikh et al., 2016; Sudmant et al.,
90 2015). They also published sequencing data of seven human genomes, including PacBio data of
91 an Ashkenazi Jewish (AJ) family trio (Zook et al., 2016). Previously, we sequenced a Chinese
92 individual HX1 on the PacBio platform, and generated assembly-based SV call sets (Shi et al.,
93 2016). Using data sets of NA12878, HX1 and the AJ family trio, we evaluated the performance of
94 four aligner/SV caller combinations (BLASR / PBHoney-Spots, BLASR / PBHoney-Tails, BWA
95 / Sniffles and NGMLR / Sniffles) as well as NextSV under different PacBio coverages. We expect
96 that NextSV will facilitate the detection and analysis of SVs on long-read sequencing data.

97

98 **2. Results**

99 **2.1 NextSV analysis pipeline**

100 As shown in Figure 1, NextSV currently supports four aligner / SV caller combinations: BLASR
101 / PBHoney-Spots, BLASR / PBHoney-Tails, BWA / Sniffles and NGMLR / Sniffles. Some
102 accessory programs (such as SAMtools) are included in NextSV. NextSV extracts FASTQ files
103 from PacBio raw data (.hdf5 or .bam) and performs QC according to users specified settings. Once
104 the aligner / SV caller combination is selected by user, NextSV automatically generates the scripts
105 for alignment, sorting, and SV calling with appropriate parameters. When the analysis is finished,
106 NextSV will format the raw result files (.tails, .spots, or .vcf files) into bed files. If multiple
107 aligner/SV caller combinations are selected, NextSV will integrate the calls to generate a sensitive
108 (by union) and a stringent (by intersection) call set. The output of NextSV is ANNOVAR-
109 compatible, so that users can easily perform downstream annotation using ANNOVAR (Wang et
110 al., 2010). In addition, NextSV also supports job submitting via Sun Grid Engine (SGE), a popular
111 batch-queuing system in cluster environment.

112

113 **2.2 Performance of SV calling on different coverages of the NA12878 Genome**

114 To determine the optimal coverage for SV detection on PacBio data, we evaluated the performance
115 of NextSV under several different coverages. We downloaded a recently published PacBio data
116 set of NA12878 (Pendleton et al., 2015) and down-sampled the data set to 2X, 4X, 6X, 8X, 10X,
117 12X, and 15X. SV calling was performed using NextSV under each coverage. All supported
118 aligner/SV caller combinations were run. At least two supporting reads is required for all SV calls.
119 The resulting calls were compared with the gold standard SV set (including 2094 deletion calls
120 and 1114 insertion calls) described in method section.

121

122 First, we examined how many calls in the gold set can be discovered. As shown in Figure 1, the
123 recall increased rapidly before 10X coverage but the slope of increase slowed down after 10X.
124 Among the four aligner / SV caller combinations, BLASR / PBHoney-Spots had the highest recall
125 for insertions while NGMLR / Sniffles had the highest recall for deletions. At 10X coverage,
126 BLASR / PBHoney-Spots detected 76.9% of deletions and 81.8% insertions in the gold standard
127 set; NGMLR / Sniffles discovered 90.9% deletions and 75.1% insertions in the gold standard set.
128 BWA / Sniffles had a lower recall for deletions (72.5%) and insertions (51.3%) than NGMLR /
129 Sniffles, indicating NGMLR is a better aligner for Sniffles. PBHoney-Tails only detected 26.6%
130 deletions and 0.09% insertions. NextSV sensitive call set, which was generated by the union call
131 set of BLASR / PBHoney-Spots, BLASR / PBHoney-Tails, and NGMLR / Sniffles, had the highest
132 recall. At 10X coverage, the recall of NextSV sensitive call set is 93.5~94.1% for deletions and
133 87.9~93.2% for insertions. At 15X coverage, the recall of NextSV sensitive call set increased
134 slightly. Therefore, 10X coverage might be an optimal coverage to use in practice, considering the
135 relatively high sequencing costs and the generally high recall rates.

136

137 Second, we examined the precision and balanced accuracy (F1 scores) under different coverages
138 (Figure 2). The precision is calculated as the fraction of detected SVs that matching the gold
139 standard set. For deletions calls, NextSV stringent call set had the second highest precision and
140 highest F1 value. For insertion calls, NextSV stringent call set had the highest precision and F1
141 value at each coverage. Therefore, NextSV stringent call set performs the best, considering the
142 balance between recall and precision.

143

144 **2.3 Performance of SV calling on different coverages on the HX1 Genome**

145 To verify the performance of SV detection on different individuals, we also performed evaluation
146 on a Chinese genome HX1, which was sequenced by us recently (Shi et al., 2016) at 103X PacBio
147 coverage. The genome was sequenced using a newer version of chemical reagents and thus the
148 mean read length of HX1 was 40% longer than that of NA12878 (Table 1). The total data set was
149 down-sampled to three representative coverages (6X, 10X and 15X). For each coverage, SVs were
150 called using the four pipelines described above and compared to the gold standard set. The results
151 were similar to those of the NA12878 data set (Figure 3). At 10X coverage, NextSV sensitive call
152 set had a recall of 94.1% for deletions and 93.2% for insertions, highest among all the call sets.
153 NextSV stringent call set had the highest precisions and F1 values. Among the four aligner / SV
154 caller combinations, NGMLR / Sniffles discovered the most deletions (91.5%) and BLASR /
155 PBHoney-Spots discovered the most insertions (81.7%) at 10X coverage. BWA / Sniffles had a
156 higher precision but a lower recall and F1 value than NGMLR / Sniffles.

157

158 **2.4 Evaluation on Mendelian Errors**

159 As the germline mutation rate is very low (Kong et al., 2012; Veltman and Brunner, 2012),
160 Mendelian errors are more likely a result of genotyping errors and can be used as a quality control
161 criteria in genome sequencing (Pilipenko et al., 2014). Due to the lack of gold standard call sets,
162 here, we evaluated the errors of allele drop-in (ADI), which means that the presence of an alleles
163 in offspring that does not appear in either parent. We used a whole genome sequencing data set of
164 an AJ family trio released by NIST (Zook et al., 2016) to do the evaluation. The sequencing data
165 of AJ son, AJ father and AJ mother was down-sampled to 10X coverage. SV detection was
166 performed using NextSV with all supported aligners and SV callers enabled. The calls from AJ
167 son were compared with calls from AJ father and AJ mother. The results showed that, NextSV
168 stringent call set had the lowest ADI rate for both deletions (10.4%) and insertions (23.5%).
169 Among the four aligner/SV caller combinations, NGMLR / Sniffles was the best for both deletions
170 and insertions.

171

172 **2.5 Computational Performance of NextSV**

173 To evaluate the computational resources consumed by NextSV, we used the whole genome
174 sequencing data set of HX1 (10X coverage) for benchmarking. All aligners and SV callers in
175 NextSV were tested using a machine equipped with 12-core Intel Xeon 2.66 GHz CPU and 48
176 Gigabytes of memory. As shown in Table 5, mapping is the most time-consuming step. BLASR
177 takes about 80 hours to map the reads, whereas NGMLR needs 11.2 hours, which is the fastest
178 among the three aligners. The SV calling step is much faster. PBHoney-Spots and Sniffles take
179 about 1 hour, while PBHoney-Tails needs 0.27 hour. In total, the BLASR / PBHoney combination
180 takes 80.8 hours while the NGMLR / Sniffles combination takes 12.5 hours, 84.5% less than the
181 former one. Since BLASR/PBHoney-Spots and NGMLR / Sniffles have good performance on SV
182 calling and running PBHoney-Tails is very fast given the BLASR output, the NextSV pipeline will
183 execute the three methods by default for generating the final results.

184

185 **3. Discussion**

186 Long-read sequencing such as PacBio sequencing has clear advantages over short-read sequencing
187 on SV discovery (English et al., 2015). However, its application in real-world setting is often
188 limited due to the relatively high sequencing cost and hence the relatively low sequencing coverage.
189 In this study, we developed NextSV, a computational pipeline integrating multiple aligners and
190 SV callers to improve SV discovery on low-coverage PacBio data sets. Our results showed that,
191 NextSV stringent call set had the highest precisions and F1 values while NextSV sensitive call set
192 had the highest recall. At 10X coverage, the recall of NextSV sensitive call set was 93.5%~94.1%
193 for deletions and 87.9%~93.2% for insertions. At 15X coverage, there is only a slight increase in
194 recall. Therefore, ~10X coverage can be an optimal coverage to use in practice, considering the
195 balance between the sequencing costs and the recall rates.

196

197 There is often a trade-off between recall and precision. NextSV generates a sensitive call set and
198 a stringent call set, for different purposes. NextSV sensitive call set is suitable for users who
199 consider recall more important than precision and who can afford extensive downstream analysis
200 (such as Sanger sequencing) to validate the candidate variants. This is often the case when doing
201 disease-casual variant discovery on personal genomes. NextSV stringent call set has the highest
202 precision, F1 value and Mendelian error. It is suitable for users who aim to perform genome-wide
203 analysis of SVs on a collection of samples, with limited downstream validation.

204

205 The performance of SV callers are affected by the parameter settings. By default, PBHoney
206 requires a minimal read support of 3 for an SV event and Sniffles requires a minimal read support
207 of 10 for an SV event. However, this may be too high for low coverage data set. In our evaluation
208 of recall and precision, we changed this setting to require a minimal read support of 2. This allows
209 detection SVs from very low coverage regions, with an acceptable precision. This result in
210 substantially higher number of true positives and less variants of interest would be missed. The
211 increased false positive calls can be removed by downstream validation or using call sets of two
212 SV callers (e.g. using the NextSV stringent call set).

213

214 In addition to test recalls and precisions, we examined the allele drop-in errors, which represent
215 the SV calls that in the offspring but not appear in either parent. The allele drop-in errors can come
216 from two sources: false positive calls of the offspring or false negatives in the parents, though in
217 very rare cases it could be due to de novo mutations. So this measure is related to both recall and
218 precision. Since we consider 10X coverage as a good choice, we did the evaluation on a family
219 trio data set with ~10X coverage. In our results, NextSV stringent call set has the lowest allele
220 drop-in error, which is consistent with the results that it has the highest F1 value.

221

222 NextSV currently supports four aligner / SV caller combinations: BLASR / PBHoney-Spots,
223 BLASR / PBHoney-Tails, BWA / Sniffles, NGMLR / Sniffles, but we expect to continuously
224 expand the support for other aligner / caller combinations. Users can choose to run any of them.
225 By default, NextSV will enable BLASR / PBHoney-Spots, BLASR / PBHoney-Tails and NGMLR
226 / Sniffles and integrate the results to generate the sensitive calls and stringent calls. We do not
227 enable BWA / Sniffles by default because Sniffles works better with NGMLR in our evaluation
228 and alignment is a time consuming step. SVs that are shorter than reads may result in intra-read
229 discordances while larger SVs may result in soft-clipped tails of long reads. We suggest running
230 both PBHoney-Spots and PBHoney-Tails because they are two complementary algorithms
231 designed to detect intra-read discordances and soft-clipped tails, respectively. Sniffles uses
232 multiple evidences to detect SV so it should be suitable for both small and large SVs.

233

234 In this study, we only evaluated the performance for insertions and deletions because we only have
235 the gold standard calls of insertions and deletions. This is another limitation of the study. We will
236 evaluate the performance on other types of SVs in the future when more gold standard SV calls
237 are available. Nonetheless, NextSV generates SV calls of all types. The output of NextSV is in
238 ANNOVAR-compatible bed format. Users can easily perform downstream annotation using
239 ANNOVAR and disease gene discovery using Phenolyzer (Yang et al., 2015). NextSV is available
240 at <http://github.com/Nextomics/NextSV> and can be installed by one simple command. We believe
241 that NextSV will facilitate the detection of structural variants from low coverage long-read
242 sequencing data.

243

244 **4. Materials and Methods**

245 **4.1 PacBio data sets used for this study**

246 Five whole-genome PacBio sequencing data sets were used to test the performance of SV calling
247 pipelines (Table 1). Data sets of NA12878 and HX1 genome were downloaded from NCBI SRA
248 database. Data sets of the AJ family trio were downloaded from ftp site of NIST ([ftp://ftp-
249 trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/)). After we obtained raw data, we extracted
250 subreads (reads that can be used for analysis) using the SMRT Portal software (Pacific Biosciences,
251 Menlo Park, CA) with filtering parameters (minReadScore=0.75, minLength=500). The subreads
252 were mapped to the reference genome using BLASR (Chaisson and Tesler, 2012), BWA-MEM
253 (Li, 2013) or NGMLR (Rescheneder et al., 2016). The BAM files were down-sampled to different
254 coverages using SAMtools (samtools view -s). The down-sampled coverages and mean read
255 lengths of the data sets are shown in Table 1.

256

257 **4.2 SV detection using BLASR / PBHoney-Spots and BLASR / PBHoney-Tails**

258 PacBio subreads were iteratively aligned with the human reference genome (GRCh38 for HX1,
259 GRCh37 for NA12878 and AJ trio genomes, depending on the reference of gold standard set)
260 using the BLASR aligner (parameter: -bestn 1). Each read's single best alignment was stored in
261 the SAM output. Unmapped portions of each read were extracted from the alignments and
262 remapped to the reference genome. The alignments in SAM format were converted to BAM format
263 and sorted by SAMtools. PBHoney-Tails and PBHoney-Spots (from PBSuite-15.8.24) were run

264 with slightly modified parameters (minimal read support 2, instead of 3 and consensus polishing
265 disabled) to increase sensitivity and to discover SVs under low coverages (2~15X).

266

267 **4.3 SV detection using BWA / Sniffles and NGMLR / Sniffles**

268 PacBio subreads were aligned to the reference genome, using BWA-MEM (bwa mem -M -x pacbio)
269 or NGMLR (default parameters) to generate the BAM file. The BAM file was sorted by SAMtools,
270 then used as input of Sniffles (version 1.0.5). Sniffles was run with slightly modified parameters
271 (minimal read support 2, instead of 10) to increase sensitivity and discover SVs under low fold of
272 coverages (2~15X).

273

274 **4.4 NextSV sensitive call set and NextSV stringent call set**

275 NextSV sensitive call set is generated as

$$276 \quad \text{SNIF} \cup (\text{SPOT} \cup \text{TAIL}),$$

277 and NextSV stringent call set is generated as

$$278 \quad \text{SNIF} \cap (\text{SPOT} \cup \text{TAIL}),$$

279 where SNI denotes the call set of NGMLR / Sniffles, SPOT denotes the call set of BLASR /
280 PBHoney-Spots and TAIL denotes the call set of BLASR / PBHoney-Tails.

281

282 **4.5 Comparing two SV call sets**

283 Calls which reciprocally overlapped by more than 50% (bedtools intersect -f 0.5 -F 0.5) were
284 considered to be the concordant SV calls and were merged into a single call. For insertion calls, a
285 padding of 500 bp was added before intersection. When merging two SVs, the average start and
286 end positions were taken.

287

288 **4.6 Gold standard SV call set**

289 The gold standard deletion call set of the NA12878 genome was release by the Genome In A Bottle
290 (GIAB) consortium (Parikh et al., 2016), in which most of the calls were refined by experimental
291 validation or other independent technologies. The gold standard insertion call set of the NA12878
292 genome was obtained by merging the high-confidence insertion calls of 1000 Genome phase 3
293 (Sudmant et al., 2015) and high-confidence insertion calls from GIAB. For the HX1 genome, due
294 to the availability of high-coverage (>100X) data, we used the SV calls from a previously validated

295 local assembly-based approach (Chaisson et al., 2015) as the initial high-quality calls. We also
296 detected SVs on 100X coverage PacBio data set of the HX1 genome using BLASR / PBHoney-
297 Spots, BLASR / PBHoney-Tails, BWA / Sniffles and NGMLR / Sniffles (minimal read support=20
298 for each SV caller). The initial high-quality calls that overlapped with one of the four 103X call
299 sets were retained as final gold standard calls. SVs with length less than 200 bp were not considered.
300 Number of SVs in the gold standard sets is shown in Table 2.

301

302 **4.7 Performance Evaluation of SV callers**

303 The SV calls of each caller were compared with the gold standard SV set. Precision, recall, and F1
304 score were used to evaluate the performance of the callers. Precision, recall, and F1 were calculated
305 as

$$306 \text{ Precision} = \frac{TP}{TP+FP},$$

$$307 \text{ Recall} = \frac{TP}{TP+FN},$$

$$308 \text{ F1} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}},$$

309 where TP is the number of true positives (variants called by a variant caller and matching the gold
310 standard set), FP is the number of false positives (variants called by a variant caller but not in the
311 gold standard set), and FN is the number of false negatives (variants in the gold standard set but
312 not called by a variant caller).

313

314 **Acknowledgments**

315 The authors wish to thank the National Institute of Standards and Technology and Genome in a
316 Bottle Consortium for making the reference data on PacBio sequencing available to benchmark
317 bioinformatics software tools. We also thank members of Grandomics to test the software tools
318 and offering valuable feedback.

319

320 **Author Contributions**

321 L.F. performed the evaluation and wrote the software. J.H. and D.W. tested the software and
322 advised on the study. K.W. conceived and supervised the study, and revised the manuscript.

323

324 **Competing Interests**

325 L.F., J.H. and D.W. are employees and K.W. is a consultant for Grandomics Biosciences.

326

327 **References**

- 328 Carvalho, C.M., Lupski, J.R., 2016. Mechanisms underlying structural variant formation in genomic
329 disorders. *Nat Rev Genet* 17, 224-238.
- 330 Chaisson, M.J., Huddleston, J., Dennis, M.Y., Sudmant, P.H., Malig, M., Hormozdiari, F., Antonacci, F.,
331 Surti, U., Sandstrom, R., Boitano, M., Landolin, J.M., Stamatoyannopoulos, J.A., Hunkapiller, M.W.,
332 Korfach, J., Eichler, E.E., 2015. Resolving the complexity of the human genome using single-molecule
333 sequencing. *Nature* 517, 608-611.
- 334 Chaisson, M.J., Tesler, G., 2012. Mapping single molecule sequencing reads using basic local alignment
335 with successive refinement (BLASR): application and theory. *BMC Bioinformatics* 13, 238.
- 336 English, A.C., Salerno, W.J., Hampton, O.A., Gonzaga-Jauregui, C., Ambreth, S., Ritter, D.I., Beck, C.R.,
337 Davis, C.F., Dahdouli, M., Ma, S., Carroll, A., Veeraraghavan, N., Bruestle, J., Drees, B., Hastie, A., Lam,
338 E.T., White, S., Mishra, P., Wang, M., Han, Y., Zhang, F., Stankiewicz, P., Wheeler, D.A., Reid, J.G.,
339 Muzny, D.M., Rogers, J., Sabo, A., Worley, K.C., Lupski, J.R., Boerwinkle, E., Gibbs, R.A., 2015.
340 Assessing structural variation in a personal genome-towards a human reference diploid genome. *BMC*
341 *Genomics* 16, 286.
- 342 English, A.C., Salerno, W.J., Reid, J.G., 2014. PBHoney: identifying genomic variants via long-read
343 discordance and interrupted mapping. *BMC Bioinformatics* 15, 180.
- 344 Feuk, L., Carson, A.R., Scherer, S.W., 2006. Structural variation in the human genome. *Nat Rev Genet* 7,
345 85-97.
- 346 Kong, A., Frigge, M.L., Masson, G., Besenbacher, S., Sulem, P., Magnusson, G., Gudjonsson, S.A.,
347 Sigurdsson, A., Jonasdottir, A., Jonasdottir, A., Wong, W.S., Sigurdsson, G., Walters, G.B., Steinberg, S.,
348 Helgason, H., Thorleifsson, G., Gudbjartsson, D.F., Helgason, A., Magnusson, O.T., Thorsteinsdottir, U.,
349 Stefansson, K., 2012. Rate of de novo mutations and the importance of father's age to disease risk.
350 *Nature* 488, 471-475.
- 351 Li, H., 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv
352 1303.3997v2 [q-bio.GN].
- 353 Merker, J., Wenger, A.M., Sneddon, T., Grove, M., Waggott, D., Utiramerur, S., Hou, Y., Lambert, C.C.,
354 Eng, K.S., Hickey, L., Korfach, J., Ford, J., Ashley, E.A., 2016. Long-read whole genome sequencing
355 identifies causal structural variation in a Mendelian disease. bioRxiv doi:10.1101/090985.
- 356 Moncunill, V., Gonzalez, S., Bea, S., Andrieux, L.O., Salaverria, I., Royo, C., Martinez, L., Puiggros, M.,
357 Segura-Wang, M., Stutz, A.M., Navarro, A., Royo, R., Gelpi, J.L., Gut, I.G., Lopez-Otin, C., Orozco, M.,
358 Korb, J.O., Campo, E., Puente, X.S., Torrents, D., 2014. Comprehensive characterization of complex
359 structural variations in cancer by directly comparing genome sequence reads. *Nat Biotechnol* 32, 1106-
360 1112.
- 361 Pang, A.W., MacDonald, J.R., Pinto, D., Wei, J., Rafiq, M.A., Conrad, D.F., Park, H., Hurler, M.E., Lee,
362 C., Venter, J.C., Kirkness, E.F., Levy, S., Feuk, L., Scherer, S.W., 2010. Towards a comprehensive
363 structural variation map of an individual human genome. *Genome Biol* 11, R52.
- 364 Parikh, H., Mohiyuddin, M., Lam, H.Y., Iyer, H., Chen, D., Pratt, M., Bartha, G., Spies, N., Losert, W.,
365 Zook, J.M., Salit, M., 2016. svclassify: a method to establish benchmark structural variant calls. *BMC*
366 *Genomics* 17, 64.
- 367 Pendleton, M., Sebra, R., Pang, A.W., Ummat, A., Franzen, O., Rausch, T., Stutz, A.M., Stedman, W.,
368 Anantharaman, T., Hastie, A., Dai, H., Fritz, M.H., Cao, H., Cohain, A., Deikus, G., Durrett, R.E.,
369 Blanchard, S.C., Altman, R., Chin, C.S., Guo, Y., Paxinos, E.E., Korb, J.O., Darnell, R.B., McCombie,
370 W.R., Kwok, P.Y., Mason, C.E., Schadt, E.E., Bashir, A., 2015. Assembly and diploid architecture of an
371 individual human genome via single-molecule technologies. *Nat Methods* 12, 780-786.
- 372 Pilipenko, V.V., He, H., Kurowski, B.G., Alexander, E.S., Zhang, X., Ding, L., Mersha, T.B., Kottyan, L.,
373 Fardo, D.W., Martin, L.J., 2014. Using Mendelian inheritance errors as quality control criteria in whole
374 genome sequencing data set. *BMC Proc* 8, S21.
- 375 Rescheneder, P., Sedlazeck, F.J., Haeseler, A.V., Schatz, M.C., 2016. NGMLR: Highly accurate read
376 mapping of third generation sequencing reads for improved structural variation analysis, *Genome*
377 *Informatics* 2016, Wellcome Genome Campus Conference Centre, Hinxton, Cambridge, UK.

378 Shi, L., Guo, Y., Dong, C., Huddleston, J., Yang, H., Han, X., Fu, A., Li, Q., Li, N., Gong, S., Lintner, K.E.,
379 Ding, Q., Wang, Z., Hu, J., Wang, D., Wang, F., Wang, L., Lyon, G.J., Guan, Y., Shen, Y., Evgrafov, O.V.,
380 Knowles, J.A., Thibaud-Nissen, F., Schneider, V., Yu, C.Y., Zhou, L., Eichler, E.E., So, K.F., Wang, K.,
381 2016. Long-read sequencing and de novo assembly of a Chinese genome. *Nat Commun* 7, 12065.
382 Stankiewicz, P., Lupski, J.R., 2010. Structural variation in the human genome and its role in disease.
383 *Annu Rev Med* 61, 437-455.
384 Sudmant, P.H., Rausch, T., Gardner, E.J., Handsaker, R.E., Abyzov, A., Huddleston, J., Zhang, Y., Ye,
385 K., Jun, G., Hsi-Yang Fritz, M., Konkol, M.K., Malhotra, A., Stutz, A.M., Shi, X., Paolo Casale, F., Chen,
386 J., Hormozdiari, F., Dayama, G., Chen, K., Malig, M., Chaisson, M.J., Walter, K., Meiers, S., Kashin, S.,
387 Garrison, E., Auton, A., Lam, H.Y., Jasmine Mu, X., Alkan, C., Antaki, D., Bae, T., Cerveira, E., Chines,
388 P., Chong, Z., Clarke, L., Dal, E., Ding, L., Emery, S., Fan, X., Gujral, M., Kahveci, F., Kidd, J.M., Kong,
389 Y., Lameijer, E.W., McCarthy, S., Flicek, P., Gibbs, R.A., Marth, G., Mason, C.E., Menelaou, A., Muzny,
390 D.M., Nelson, B.J., Noor, A., Parrish, N.F., Pendleton, M., Quitadamo, A., Raeder, B., Schadt, E.E.,
391 Romanovitch, M., Schlattl, A., Sebra, R., Shabalina, A.A., Untergasser, A., Walker, J.A., Wang, M., Yu, F.,
392 Zhang, C., Zhang, J., Zheng-Bradley, X., Zhou, W., Zichner, T., Sebati, J., Batzer, M.A., McCarroll, S.A.,
393 Genomes Project, C., Mills, R.E., Gerstein, M.B., Bashir, A., Stegle, O., Devine, S.E., Lee, C., Eichler,
394 E.E., Korb, J.O., 2015. An integrated map of structural variation in 2,504 human genomes. *Nature* 526,
395 75-81.
396 Veltman, J.A., Brunner, H.G., 2012. De novo mutations in human genetic disease. *Nat Rev Genet* 13,
397 565-575.
398 Wang, K., Li, M., Hakonarson, H., 2010. ANNOVAR: functional annotation of genetic variants from high-
399 throughput sequencing data. *Nucleic Acids Res* 38, e164.
400 Weischenfeldt, J., Symmons, O., Spitz, F., Korb, J.O., 2013. Phenotypic impact of genomic structural
401 variation: insights from and for human disease. *Nat Rev Genet* 14, 125-138.
402 Yang, H., Robinson, P.N., Wang, K., 2015. Phenolyzer: phenotype-based prioritization of candidate
403 genes for human diseases. *Nat Methods* 12, 841-843.
404 Yang, L., Luquette, L.J., Gehlenborg, N., Xi, R., Haseley, P.S., Hsieh, C.H., Zhang, C., Ren, X.,
405 Protopopov, A., Chin, L., Kucherlapati, R., Lee, C., Park, P.J., 2013. Diverse mechanisms of somatic
406 structural variations in human cancer genomes. *Cell* 153, 919-929.
407 Zhang, F., Gu, W., Hurles, M.E., Lupski, J.R., 2009. Copy number variation in human health, disease,
408 and evolution. *Annu Rev Genomics Hum Genet* 10, 451-481.
409 Zook, J.M., Catoe, D., McDaniel, J., Vang, L., Spies, N., Sidow, A., Weng, Z., Liu, Y., Mason, C.E.,
410 Alexander, N., Henaff, E., McIntyre, A.B., Chandramohan, D., Chen, F., Jaeger, E., Moshrefi, A., Pham,
411 K., Stedman, W., Liang, T., Saghbini, M., Dzakula, Z., Hastie, A., Cao, H., Deikus, G., Schadt, E., Sebra,
412 R., Bashir, A., Truty, R.M., Chang, C.C., Gulbahce, N., Zhao, K., Ghosh, S., Hyland, F., Fu, Y., Chaisson,
413 M., Xiao, C., Trow, J., Sherry, S.T., Zaranek, A.W., Ball, M., Bobe, J., Estep, P., Church, G.M., Marks, P.,
414 Kyriazopoulou-Panagiotopoulou, S., Zheng, G.X., Schnall-Levin, M., Ordonez, H.S., Mudivarti, P.A.,
415 Giorda, K., Sheng, Y., Rypdal, K.B., Salit, M., 2016. Extensive sequencing of seven human genomes to
416 characterize benchmark reference materials. *Sci Data* 3, 160025.
417

418

419 **Tables**

420 **Table 1. Description of PacBio data sets used for this study.**

Data Source / Accession	Genome	Down-sampled Coverage	Mean Read Length	Reference
SRX627421	NA12878	2~15X	4.9 kb	(Pendleton et al., 2015)
SRX1424851	HX1	6~15X	7.0 kb	(Shi et al., 2016)
NIST	AJ son	10X	8.0 kb	(Zook et al., 2016)
NIST	AJ father	10X	7.3 kb	(Zook et al., 2016)
NIST	AJ mother	10X	7.8 kb	(Zook et al., 2016)

421

422

423 **Table 2. Number of calls in gold standard SV set**

Genome	Platform	Number of Deletions (≥ 200 bp)	Number of Insertions (≥ 200 bp)	Reference
NA12878	Illumina	2094	1114	(Parikh et al., 2016; Sudmant et al., 2015)
HX1	PacBio	2387	2937	(Shi et al., 2016)

424

425

426 **Table 3. Mendelian error of deletion calls under 10X coverage**

	BLASR / PBHoney- Spots	BLASR / PBHoney- Tails	BWA / Sniffles	NGMLR / Sniffles	NextSV Sensitive Calls	NextSV Stringent Calls
No. of calls (AJ father)	2943	775	2173	3109	4172	2342
No. of calls (AJ mother)	3090	789	2008	3169	4299	2399
No. of calls (AJ son)	3120	727	1976	3182	4246	2444
No. of calls inherited from father	2047	306	1238	2151	2812	1684
No. of calls inherited from mother	2166	295	1235	2232	2929	1747
No. of ADI	447	296	335	376	600	253
ADI rate	14.3%	40.7%	17.0%	11.8%	14.1%	10.4%

427

428

429

430

431

432

433

434

Table 4. Mendelian error of insertion calls under 10X coverage

	BLASR / PBHoney- Spots	BLASR / PBHoney- Tails	BWA / Sniffles	NGMLR / Sniffles	NextSV Sensitive Calls	NextSV Stringent Calls
No. of calls (AJ father)	4817	18	764	2601	5326	2103
No. of calls (AJ mother)	5151	20	855	2781	5708	2237
No. of calls (AJ son)	5341	17	903	2708	5815	2243
No. of calls inherited from father	2837	5	255	1393	3142	1182
No. of calls inherited from mother	2907	6	247	1458	3228	1238
No. of ADI	1625	10	520	711	1756	528
ADI rate	30.4%	58.8%	57.6%	26.3%	30.2%	23.5%

435

436

437

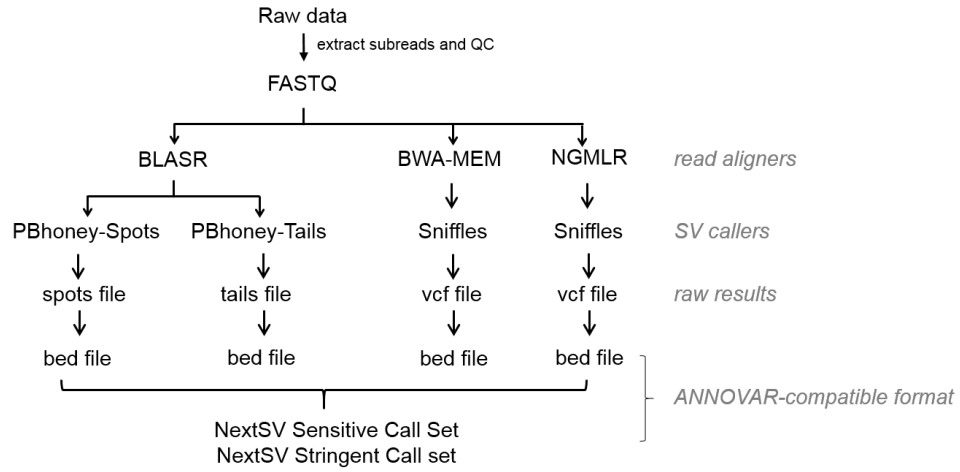
Table 5. Time consumption for each steps in the NextSV pipeline for 10X PacBio data set

SV caller	Aligner	CPU (number of threads)	Alignment time (hour)	SV calling time (hour)	Total Time (hour)
PBHoney	BLASR	12	79.6	0.27 (Tails) 0.96 (Spots)	80.8
Sniffles	BWA- MEM	12	27.0	1.1	28.1
Sniffles	NGMLR	12	11.2	1.3	12.5

438

439

440 **Figures**



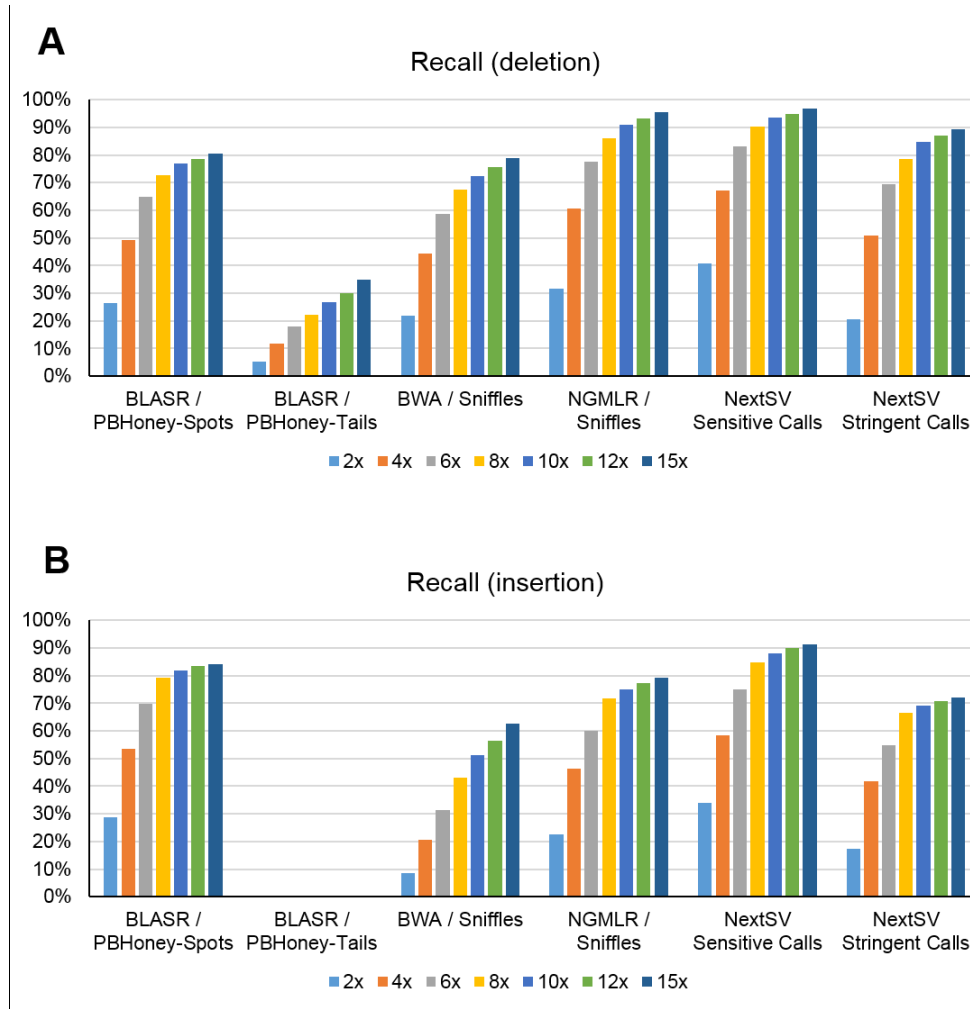
441

442

443

Figure 1. Scheme of NextSV workflow.

444



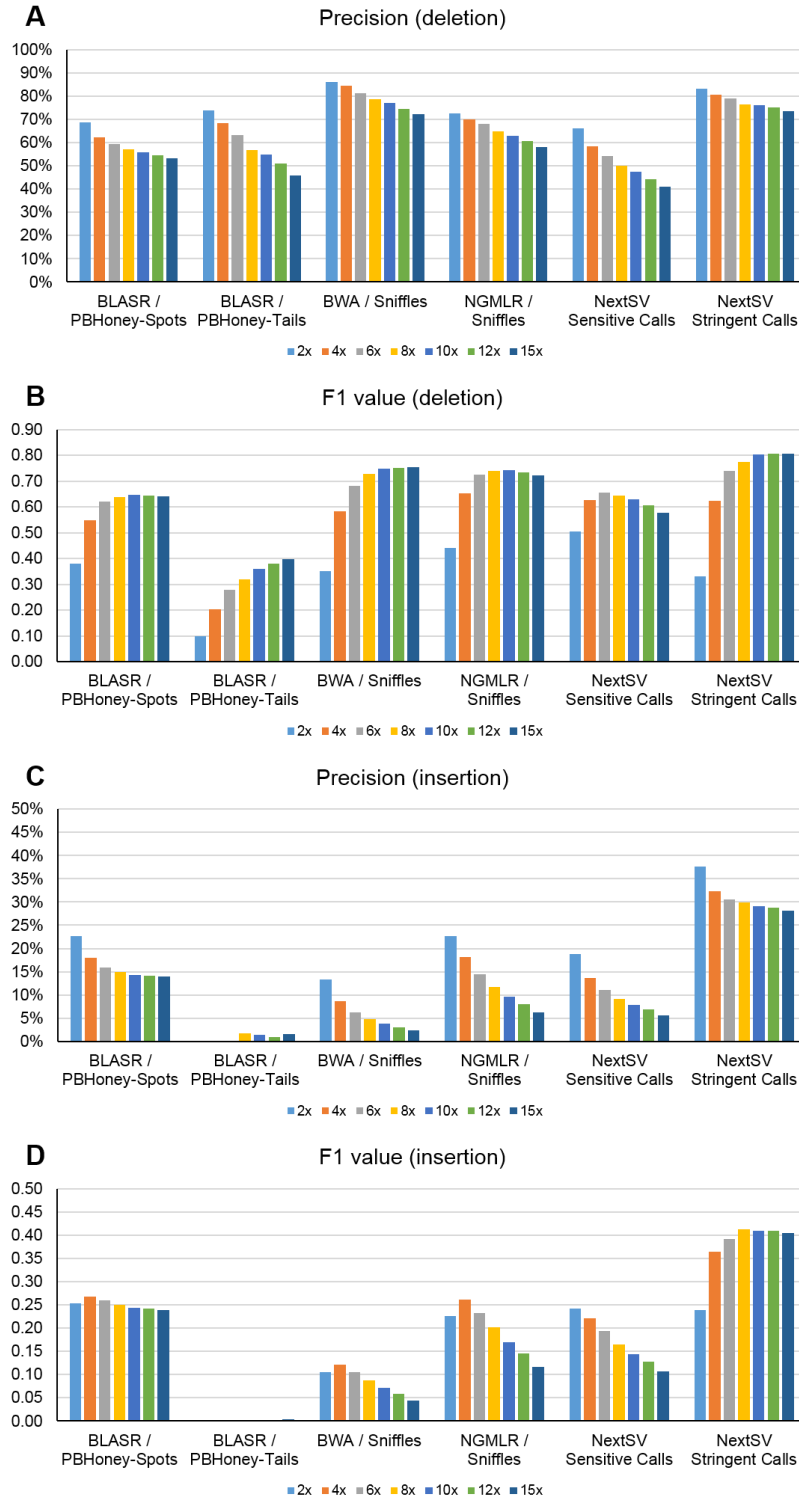
445

446

447

Figure 2. Evaluation of recall rates under different coverages on the NA12878 genome.

448



449

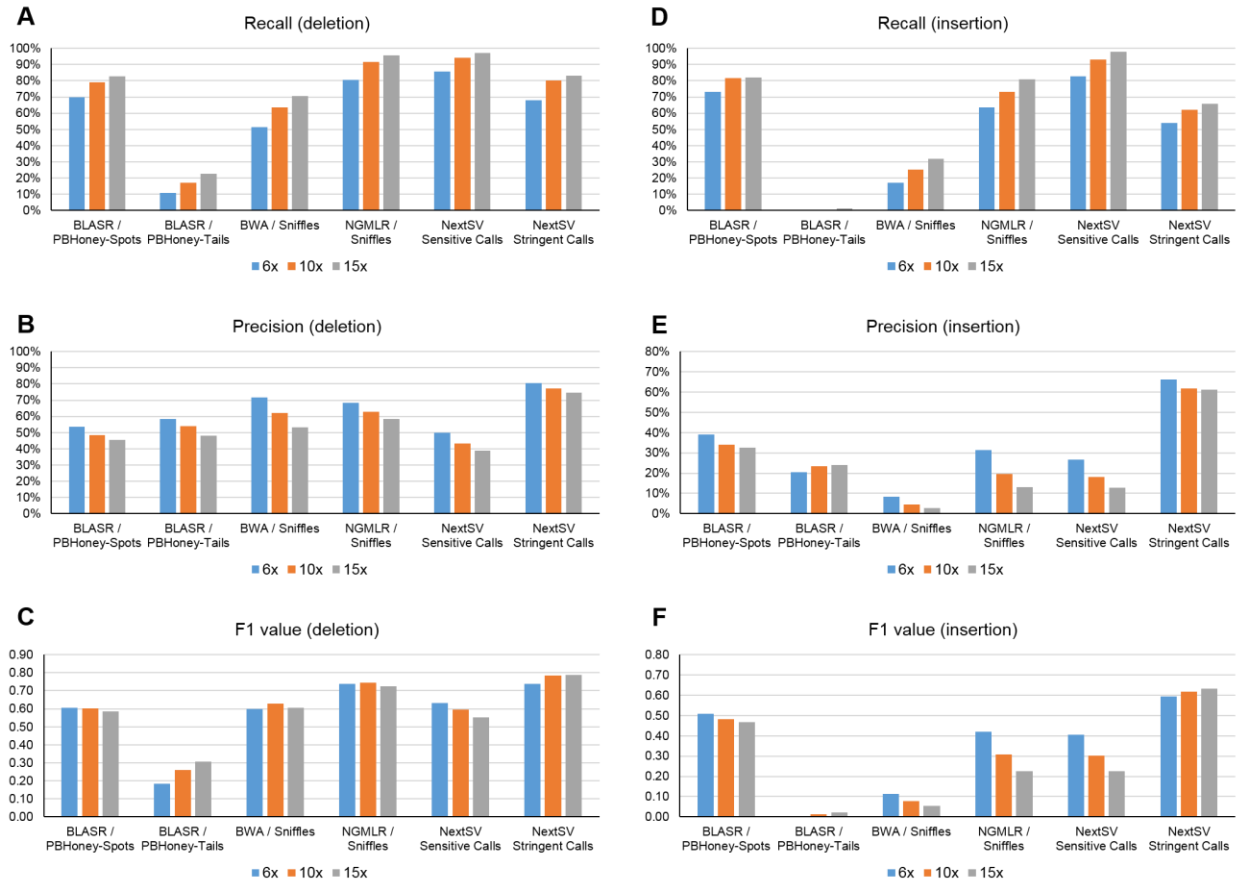
450

451

452

Figure 3. Evaluation of precisions and F1 values under different coverages on the NA12878 genome.

453



454

455

Figure 4. SV calling performance on the HX1 genome.