

1 **Evaluation on Detection of Structural Variants by Low-**
2 **Coverage Long-Read Sequencing**

3
4
5
6 **Li Fang¹, Jiang Hu¹, Depeng Wang¹, Kai Wang^{2,3,*}**

7
8
9
10 1: Grandomics Biosciences, Beijing 102206, China

11 2: Institute for Genomic Medicine, Columbia University Medical Center, New York, NY
12 10032, USA

13 3: Department of Biomedical Informatics, Columbia University Medical Center, New
14 York, NY 10032, USA

15
16 *: Correspondence should be addressed to K.W. (kai@openbioinformatics.org)

17
18

19 **Abstract**

20

21 Structural variants (SVs) in human genome are implicated in a variety of human
22 diseases. Long-read sequencing (such as those from PacBio) delivers much longer
23 read lengths than short-read sequencing (such as those from Illumina) and may greatly
24 improve SV detection. However, due to the relatively high cost of long-read sequencing,
25 users are often faced with issues such as what coverage is needed and how to
26 optimally use the aligners and SV callers. Here, we evaluated SV calling performance of
27 three SV calling algorithms (PBHoney-Tails, PBHoney-Spots and Sniffles) under
28 different PacBio coverages on two personal genomes, NA12878 and HX1. Our results
29 showed that, at 10X coverage, 76% ~ 84% deletions and 80% ~ 92 % insertions in the
30 gold standard set can be detected by PBHoney-Spots. Combining both PBHoney-Spots
31 and Sniffles greatly increased sensitivity, especially under lower coverages such as 6X.
32 We further evaluated the Mendelian errors on an Ashkenazi Jewish trio dataset with
33 low-coverage whole-genome PacBio sequencing. In addition, to automate SV calling,
34 we developed a computational pipeline called NextSV, which integrates PBhoney and
35 Sniffles and generates the union (high sensitivity) or intersection (high specificity) call
36 sets. Our results provide useful guidelines for SV identification from low coverage
37 whole-genome PacBio data and we expect that NextSV will facilitate the analysis of SVs
38 on long-read sequencing data.

39

40 Introduction

41

42 Structural variants (SVs), including large variations such as deletions, insertions,
43 duplications, inversions, and translocations, make an important contribution to human
44 diversity and disease susceptibility^{1,2}. Many inherited diseases and cancers have been
45 associated with a large number of SVs in recent years³⁻⁸. Recent advances in next-
46 generation sequencing (NGS) technologies have facilitated the analysis of variations
47 such as SNPs and small Indels in unprecedented details, but the discovery of SVs using
48 short reads still remains challenging⁹. Single-molecule, real-time (SMRT) sequencing
49 developed by Pacific Biosciences (PacBio) offers a long read length, making it
50 potentially well-suited for SV detection in personal genomes^{9,10}. Most recently, Merker
51 et al. reported the application of low coverage whole genome PacBio sequencing to
52 identify pathogenic structural variants from a patient with autosomal dominant Carney
53 complex, for whom targeted clinical gene testing and whole genome short-read
54 sequencing were negative¹¹.

55

56 Two SV software tools have been developed specifically for long-read sequencing:
57 PBhoney¹² and Sniffles (<https://github.com/fritzsedlazeck/Sniffles>). PBhoney identifies
58 genomic variants via two algorithms, long-read discordance (PBhoney-Spots) and
59 interrupted mapping (PBhoney-Tails). Sniffles is a SV caller written in C++ and it detects
60 SVs using evidence from split-read alignments, high-mismatch regions, and coverage
61 analysis. Due to the relative high cost of PacBio sequencing, users are often faced with
62 issues such as what coverage is needed and how to get the best use of the available
63 SV callers. In addition, it is unclear which software performs the best in low-coverage
64 settings, and whether the combination of software tools can improve performance of SV
65 calls. Finally, the execution of these software tools is often not straightforward and
66 requires careful re-parameterization given specific coverage of the source data.

67

68 Recently, the Genome in a Bottle (GIAB) consortium hosted by National Institute of
69 Standards and Technology (NIST) distributed a set of high-confidence SV calls for the
70 NA12878 genome, an extensively sequenced genome by different platforms, enabling

71 benchmarking of SV callers¹³. They also published sequencing data of seven human
72 genomes, including PacBio data of an Ashkenazi Jewish family trio¹⁴. Previously, we
73 sequenced a Chinese individual HX1 on the PacBio platform, and generated assembly-
74 based SV call sets¹⁵. Using data sets of NA12878, HX1 and the AJ trio, we compared
75 the performance of PBhoney-Spots, PBhoney-Tails, Sniffles and their combination
76 under different PacBio coverages. In addition, we provided NextSV, an automated SV
77 calling pipeline using PBHoney-Spots, PBHoney-Tails and Sniffles. NextSV
78 automatically execute these three other software tools with optimized parameters for the
79 specific coverage that user specified, then integrates results of each caller and
80 generates the union (high sensitivity) or intersection (high specificity) call sets. We
81 expect that NextSV will facilitate the detection and analysis of SVs on long-read
82 sequencing data.

83

84 **Materials and Methods**

85 ***PacBio data sets used for this study***

86 Five whole-genome PacBio sequencing data sets were used to test the performance of
87 SV calling pipelines (Table 1). Data sets of NA12878 and HX1 genome were obtained
88 from NCBI SRA database. Data sets of the Ashkenazi Jewish (AJ) family trio were
89 downloaded from ftp site of NIST (<ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/>).
90 After we obtained raw data, we extracted subreads using the SMRT Portal software
91 (Pacific Biosciences, Menlo Park, CA) with default settings. The subreads were mapped
92 to the reference genome using BLASR¹⁶ or BWA-MEM¹⁷. The bam files were down-
93 sampled to different coverages using SAMtools (`samtools view -s`). The down-sampled
94 coverages and mean read lengths of the data sets are shown in Table 1.

95

96 ***SV detection using PBHoney***

97 PacBio subreads were iteratively aligned with the human reference genome (GRCh38
98 for HX1, GRCh37 for NA12878 and AJ trio genomes, depending on the reference of
99 gold standard set) using the BLASR aligner (parameter: `-bestn 1`). Each read's single
100 best alignment was stored in the SAM output. Unmapped portions of each read were
101 extracted from the alignments and remapped to the reference genome. The alignments

102 in SAM format were converted to BAM format and sorted by SAMtools. PBHoney-Tails
103 and PBHoney-Spots were run with slightly modified parameters (minimal read support 2,
104 instead of 3 and consensus polishing disabled) to increase sensitivity and discover SVs
105 under low coverages (2~15X).

106

107 ***SV detection using Sniffles***

108 PacBio subreads were aligned to the reference genome, using BWA-MEM with
109 parameters modified for PacBio reads (bwa mem -M -x pacbio), to generate the BAM
110 file. The BAM file was used as input of Sniffles. Sniffles was run with slightly modified
111 parameters (minimal read support 2, instead of 10) to increase sensitivity and discover
112 SVs under low fold of coverages (2~15X).

113

114 ***Comparing two SV call sets***

115 Calls which reciprocally overlapped by more than 50% (bedtools intersect -f 0.5 -r) were
116 considered to be the same SV and merged into a single call. For insertion calls, a
117 padding of 500 bp was added before intersection. When merging two SVs, the average
118 start and end positions were used.

119

120 ***Gold standard SV call set***

121 The gold standard SV call set for NA12878 was retrieved from the GIAB consortium ¹³,
122 in which most of the calls were refined by experimental validation or other independent
123 technologies. For the HX1 genome, we used the SV calls from a previously validated
124 local assembly approach ¹⁰, as the initial high-quality calls. We also detected SVs on
125 100X coverage PacBio data set of the HX1 genome using PBHoney-Tails, PBHoney-
126 Spots and Sniffles. The initial high-quality calls that overlapped with one of the three
127 100X call sets (PBHoney-Tails, PBHoney-Spots or Sniffles) were retained as final gold
128 standard calls. SVs with length less than 200 bp were not considered. Number of SVs in
129 the gold standard sets is shown in Table 2.

130

131 ***Performance Evaluation of SV callers***

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

$$\begin{aligned} \text{Precision} &= \frac{TP}{TP+FP}, \\ \text{Recall} &= \frac{TP}{TP+FN}, \\ \text{F1} &= 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}, \end{aligned}$$

138 where TP is the number of true positives (variants called by a variant caller and
139 matching the gold standard set), FP is the number of false positives (variants called by
140 a variant caller but not in the gold standard set), and FN is the number of false
141 negatives (variants in the gold standard set but not called by a variant caller).

142

143 **Results**

144 ***Performance of SV calling under different PacBio coverage***

145 To determine what sequencing coverage is needed for SV detection using PacBio data,
146 we evaluated the performance of SV callers under several different coverages. We
147 downloaded a recently published 22X PacBio data set of NA12878¹⁸ and down-
148 sampled the data set to 2X, 4X, 6X, 8X, 10X, 12X, and 15X. SV calling was performed
149 using PBHoney and Sniffles under each coverage. The resulting calls were compared
150 with the gold standard SV set (including 2094 deletion calls and 68 insertion calls) from
151 the Genome In A Bottle (GIAB) consortium¹⁸.

152

153 First, we examined how many calls in the gold set can be discovered. As shown in
154 Figure 1A and 1B, the recall increased rapidly before 6X coverage but the slope of
155 increase slowed down after 10X. Among the three callers, PBHoney-Spots discovered
156 more SV calls than Sniffles and PBHoney-Tails. At 10X coverage, PBHoney-Spots
157 detected 76% of deletions and 80% insertions in the gold standard set; Sniffles
158 discovered 63% deletions and 25% insertions in the gold standard set; PBHoney-Tails
159 recalled 26% deletions and 3% insertions. At 15X coverage, the recall of PBHoney-
160 Spots was 80% for deletion calls and 87% for insertion calls, which is only 6% ~ 9%
161 higher than the recall at 10X.

162
163 Second, we examined the precision and F1 scores of callers under different coverage.
164 We calculated precision as the fraction of detected SVs that matching the gold standard
165 set. As shown in Figure 1C, Sniffles has higher precision than PBHoney-Spots and
166 PBHoney-Tails. The precision of Sniffles for deletion calls was 70% at 6X coverage, and
167 decreased slightly as the coverage increased. F1 score, the harmonic mean of precision
168 and recall, increased before 10X and then kept stable at higher coverage (Figure 1D).
169 Precision for insertion calls was not assessed because there were only 86 insertion
170 calls in the GIAB gold standard set, which was one order of magnitude smaller than the
171 number of deletion calls, with potentially high false negative rates.

172
173 To verify the performance of SV detection on different individuals, we also did
174 evaluation on a Chinese genome HX1, which was sequenced by us recently¹⁵ at 103X
175 PacBio coverage. The genome was sequenced using a newer version of chemical
176 reagents and thus the mean read length of HX1 was 40% longer than that of NA12878
177 (Table 1). The total data set was down-sampled to 6X, 10X and 15X coverage. For each
178 coverage data set, SVs were called and compared to the gold standard set. The results
179 (Figure 3) were similar to those of the NA12878 data set. At 10X coverage, 84%
180 deletions and 92% insertions in the gold standard set can be detected by PBHoney-
181 Spots. The precisions at 10X coverage range from 54% ~ 60% for deletion calls and 31%
182 ~ 43% for insertion calls. At 15X coverage, the recall increased slightly but precision
183 decreased. Thus, 10X may be an optimal coverage to use in practice, considering the
184 sequencing costs and the balance of recall and precision.

185

186 ***Performance of SV calling using a combination of PBHoney and Sniffles***

187 Although PBHoney-Spots detected most of the variants, we examined whether we can
188 improve the recall rates by running both PBHoney-Spots and Sniffles, especially under
189 low fold coverages. As shown in Figure 2, at 6X coverage, the union set of both callers
190 discovered 77% deletions in the NA12878 gold standard set, which was 23% more than
191 running PBHoney-Spots alone at 6X coverage and comparable to running PBHoney-

192 Spots alone at 10X. At 15X coverage, the union set recalled 93% deletions and 88%
193 insertions.

194

195 In addition, we tested whether we can get high confidence calls by running both callers.
196 We evaluated precision of the intersection call sets of both callers on 6X, 10X and 15X
197 data sets of the HX1 genome (Figure 3 B, D). The precision of the intersection sets was
198 87% ~ 90% for deletion calls and 64% ~ 73% for insertion calls, which was half to one-
199 fold higher than that of PBHoney-Spots only.

200

201 ***Evaluation on Mendelian Errors***

202 As the germline mutation rate is very low^{19,20}, Mendelian errors are more likely a result
203 of genotyping errors and can be used as a quality control criteria in genome sequencing
204²¹. Here, we evaluated the errors of allele drop-in (ADI), which means that an offspring
205 presents an allele that does not appear in either parent, using a whole genome
206 sequencing data set of an Ashkenazi Jewish (AJ) family trio released by NIST¹⁴. The
207 sequencing data of AJ son, AJ father and AJ mother was down-sampled to 10X
208 coverage. SVs were called using PBHoney-Tails, PBHoney-Spots and Sniffles. The
209 calls from AJ son were compared with calls from AJ father and AJ mother. ADI rate was
210 calculated as the proportion of calls in offspring not matching any call from either parent.
211 The result shows that PBHoney-Spots returns the most calls. For deletion calls,
212 PBHoney-Spots gives us a lowest ADI rate (14.1%), while the ADI rates for insertion
213 calls are considerable higher (31.8% ~ 41.8). Therefore, further validation or manual
214 inspection of the calls is needed when analyzing SVs that may be associated with
215 diseases with low coverage sequencing.

216

217 ***Automated pipeline for SV calling using PBhoney and Sniffles***

218 Although we can get highly confident calls at low PacBio coverage using PBhoney and
219 Sniffles, there are still challenges for installation, execution and integration of the
220 aligners and SV callers for average users. Therefore, we developed NextSV, an
221 automated computational pipeline that allows SV calling from PacBio sequencing data
222 using PBhoney and Sniffles. The workflow of NextSV is shown in Figure 4. Two

223 mapping tools (BWA-MEM, BLASR), three SV callers (PBHoney-Tails, PBHoney-Spots
224 and Sniffles) and some accessory programs (such as SAMtools, BEDtools) were
225 included in NextSV. NextSV takes FASTA or FASTQ files as input. Once the SV caller
226 is selected, NextSV automatically chooses the compatible aligner and performs
227 mapping. The alignments will be automatically sorted and then presented to the SV
228 caller with appropriate parameters. When the analysis is finished, NextSV will examine
229 the FASTA/FASTQ, BAM, and result files and generate a report showing various
230 statistics. If more than one caller is selected, NextSV will format the raw result files
231 (.tails, .spots, or .vcf files) into bed files and generate the intersection or union call set
232 for the purpose of higher accuracy or sensitivity. In addition, NextSV also supports
233 analyzing high coverage samples via Sun Grid Engine (SGE), a popular batch-queuing
234 system in cluster environment. NextSV splits the input FASTA/FASTQ file into several
235 files of equal sizes and generates mapping task for each file. The mapping tasks are
236 then submitted to the queue. After mapping is done, the alignments are automatically
237 merged and subjected to the caller.

238

239 ***Computational Performance of NextSV***

240 To evaluate the computational resources consumed by NextSV, we used the whole
241 genome sequencing data set of HX1 (10X coverage) for benchmarking. All aligners and
242 SV callers in NextSV were tested using a machine equipped with 12-core Intel Xeon
243 2.66 GHz CPU and 48 Gigabytes of memory. As shown in Table 5, mapping is the most
244 time-consuming step. BLASR takes about 80 hours to map the reads, whereas BWA-
245 MEM needs 27 hours. The SV calling step is much faster. PBHoney-Spots and Sniffles
246 take about 1 hour, while PBHoney-Tails needs 0.27 hours. In total, the BLASR /
247 PBHoney-Spots pipeline takes 80.8 hours while the BWA-MEM / Sniffles pipeline takes
248 28.1 hours, two thirds less than the former one. Since the BLASR/PBHoney-Spots
249 pipeline has improved performance on SV calling and the BWA-MEM/Sniffles pipeline is
250 faster and complementary of PBHoney, we suggest running both to get the best results
251 in practice.

252

253

254 **Discussion**

255 Depth of coverage is often a key consideration in genomic analyses²². In this study, we
256 evaluated SV calling performance of three SV calling algorithms, PBHoney-Tails,
257 PBHoney-Spots and Sniffles, at various PacBio coverages of 2 ~ 15X. Our results
258 showed that, at 10X coverage, 76% ~ 84% deletions and 80% ~ 92 % insertions were
259 detected by running PBHoney-Spots. By running both PBHoney-Spots and Sniffles,
260 comparable recall can be achieved at coverage as low as 6X. At more than 10X
261 coverage, the recall slightly increased. Thus, 10X can be an optimal PacBio coverage
262 for efficient SV detection, yet 6X may also be an economic choice under limited budget.

263

264 Given the long read length, structural variants can be spanned by reads. In our results,
265 the “Spots” algorithm of PBHoney, which was specifically designed for detection of intra-
266 read SV events, uncovered the most calls among the three algorithms. Sniffles was a
267 newly designed SV caller, and its pre-publication release version was tested in our
268 study. There are several advantages of running both PBHoney and Sniffles. First, the
269 overlapping calls are more accurate. In our results, the precisions of the intersection
270 sets were half to one-fold higher than those of PBHoney-Spots only. The recall of the
271 intersection set was 45% at 10X coverage, meaning that 45% calls can be detected at a
272 very high accuracy. Second, more calls can be discovered by running both, especially
273 for deletion calls. In our results, under 6X coverage, the union call set of two callers
274 covered 77% deletions in the NA12878 gold standard set, which was 23% more than
275 the call set of PBHoney-Spots alone. In addition, by running both BLASR/PBHoney and
276 BWA-MEM/Sniffles, we can have two BAM files for necessary manual inspection,
277 potentially eliminating the mapping artifacts that are specific to one aligner.

278

279 Besides installation of the aligners and callers, several steps are required to perform SV
280 detection using the combination of PBHoney and Sniffles, including quality check,
281 mapping, sorting, SV calling, generating union/intersection call set, and generating
282 summary statistics. In addition, several issues need to be considered during analysis.
283 PBHoney typically takes alignments from BLASR as input but Sniffles requires output
284 from BWA-MEM. The output files of PBHoney in tails or spots format should be

285 converted to standard format (such as bed or vcf) for the convenience of further
286 analysis. When two calls are merged, original information from each caller should be
287 retained. Therefore, we developed NextSV, a comprehensive solution to address this.
288 NextSV is available at <http://github.com/Nextomics/NextSV>. We believe that NextSV will
289 facilitate the detection of structural variants from low fold of PacBio sequencing data.

290

291 **Acknowledgments**

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

296

297 **Author Contributions**

298 L.F. performed the evaluation and wrote the NextSV software. J.H. and D.W. tested the
299 software and advised on the study. K.W. conceived of and supervised the study.

300

301 **Competing Interests**

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

304

305 **References**

- 306 1 Feuk, L., Carson, A. R. & Scherer, S. W. Structural variation in the human genome. *Nat Rev*
307 *Genet* **7**, 85-97 (2006).
308 2 Pang, A. W. *et al.* Towards a comprehensive structural variation map of an individual human
309 genome. *Genome Biol* **11**, R52 (2010).
310 3 Stankiewicz, P. & Lupski, J. R. Structural variation in the human genome and its role in disease.
311 *Annu Rev Med* **61**, 437-455 (2010).
312 4 Weischenfeldt, J., Symmons, O., Spitz, F. & Korbel, J. O. Phenotypic impact of genomic
313 structural variation: insights from and for human disease. *Nat Rev Genet* **14**, 125-138 (2013).
314 5 Yang, L. *et al.* Diverse mechanisms of somatic structural variations in human cancer genomes.
315 *Cell* **153**, 919-929 (2013).
316 6 Moncunill, V. *et al.* Comprehensive characterization of complex structural variations in cancer by
317 directly comparing genome sequence reads. *Nat Biotechnol* **32**, 1106-1112 (2014).
318 7 Zhang, F., Gu, W., Hurles, M. E. & Lupski, J. R. Copy number variation in human health, disease,
319 and evolution. *Annu Rev Genomics Hum Genet* **10**, 451-481 (2009).
320 8 Carvalho, C. M. & Lupski, J. R. Mechanisms underlying structural variant formation in genomic
321 disorders. *Nat Rev Genet* **17**, 224-238 (2016).

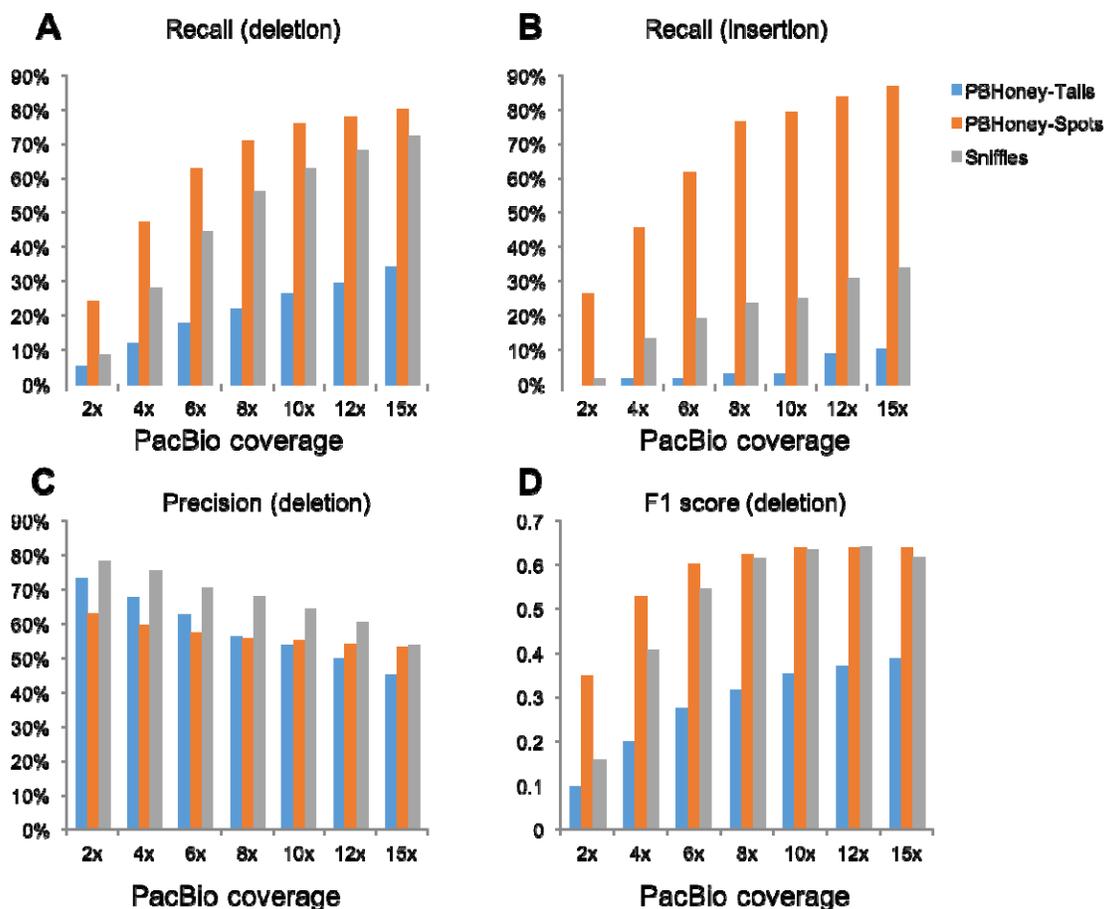
- 322 9 English, A. C. *et al.* Assessing structural variation in a personal genome-towards a human
323 reference diploid genome. *BMC Genomics* **16**, 286 (2015).
- 324 10 Chaisson, M. J. *et al.* Resolving the complexity of the human genome using single-molecule
325 sequencing. *Nature* **517**, 608-611 (2015).
- 326 11 Merker, J. *et al.* Long-read whole genome sequencing identifies causal structural variation in a
327 Mendelian disease. *bioRxiv* doi:10.1101/090985 (2016).
- 328 12 English, A. C., Salerno, W. J. & Reid, J. G. PBHoney: identifying genomic variants via long-read
329 discordance and interrupted mapping. *BMC Bioinformatics* **15**, 180 (2014).
- 330 13 Parikh, H. *et al.* svclassify: a method to establish benchmark structural variant calls. *BMC*
331 *Genomics* **17**, 64 (2016).
- 332 14 Zook, J. M. *et al.* Extensive sequencing of seven human genomes to characterize benchmark
333 reference materials. *Sci Data* **3**, 160025 (2016).
- 334 15 Shi, L. *et al.* Long-read sequencing and de novo assembly of a Chinese genome. *Nat Commun* **7**,
335 12065 (2016).
- 336 16 Chaisson, M. J. & Tesler, G. Mapping single molecule sequencing reads using basic local
337 alignment with successive refinement (BLASR): application and theory. *BMC Bioinformatics* **13**, 238
338 (2012).
- 339 17 Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*
340 1303.3997v2 [q-bio.GN] (2013).
- 341 18 Pendleton, M. *et al.* Assembly and diploid architecture of an individual human genome via single-
342 molecule technologies. *Nat Methods* **12**, 780-786 (2015).
- 343 19 Kong, A. *et al.* Rate of de novo mutations and the importance of father's age to disease risk.
344 *Nature* **488**, 471-475 (2012).
- 345 20 Veltman, J. A. & Brunner, H. G. De novo mutations in human genetic disease. *Nat Rev Genet* **13**,
346 565-575 (2012).
- 347 21 Pilipenko, V. V. *et al.* Using Mendelian inheritance errors as quality control criteria in whole
348 genome sequencing data set. *BMC Proc* **8**, S21 (2014).
- 349 22 Sims, D., Sudbery, I., Illott, N. E., Heger, A. & Ponting, C. P. Sequencing depth and coverage: key
350 considerations in genomic analyses. *Nat Rev Genet* **15**, 121-132 (2014).

351
352

353

354 **Figure and Tables**

355



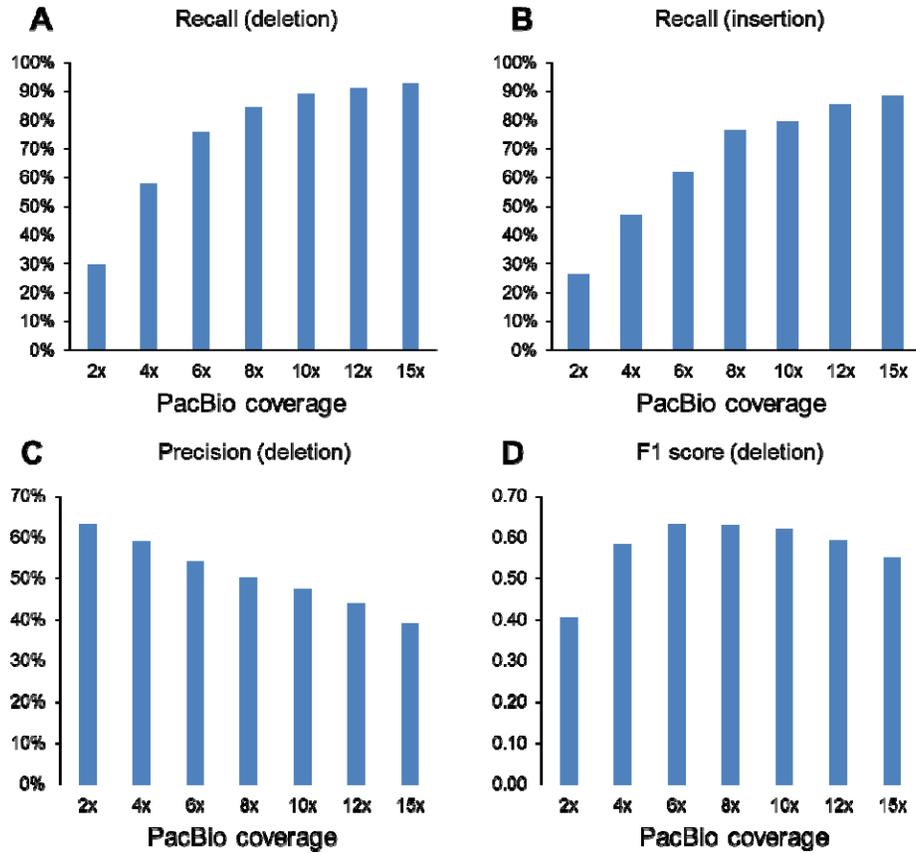
356

357

358 Figure 1. SV calling performance for each SV caller under different coverage on the

359

NA12878 genome.



360

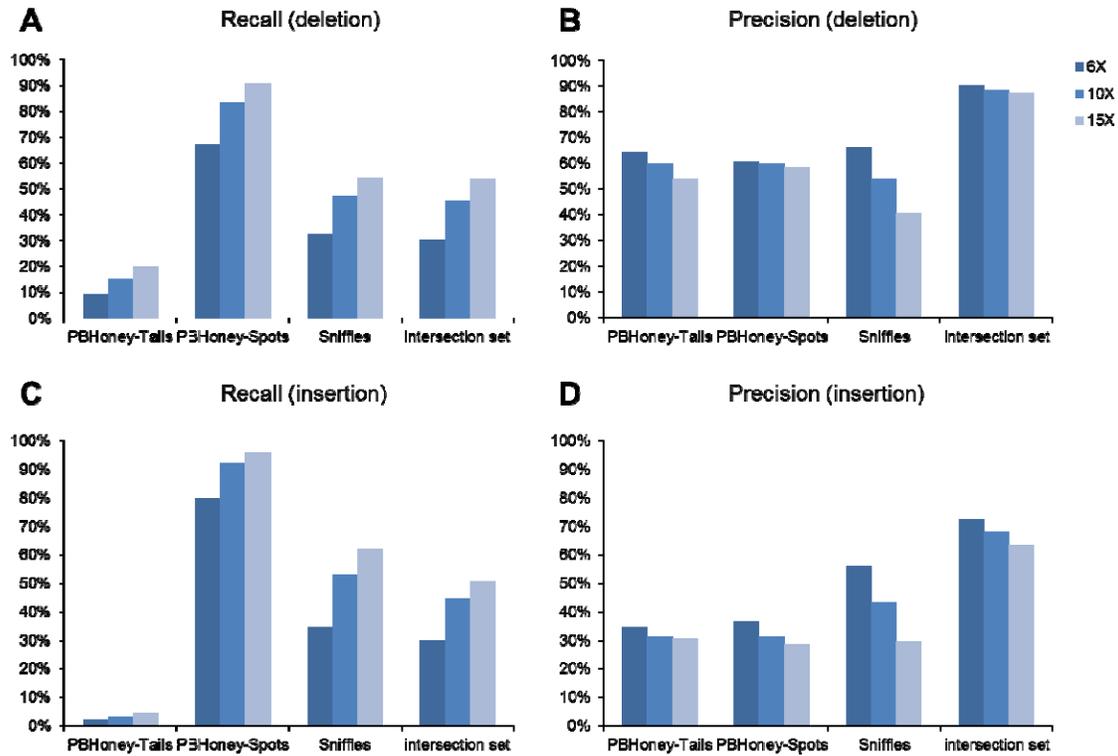
361

362 Figure 2. SV calling performance for the union call set of PBHoney-Spots and Sniffles

363

under different coverage on the NA12878 genome.

364



365

366

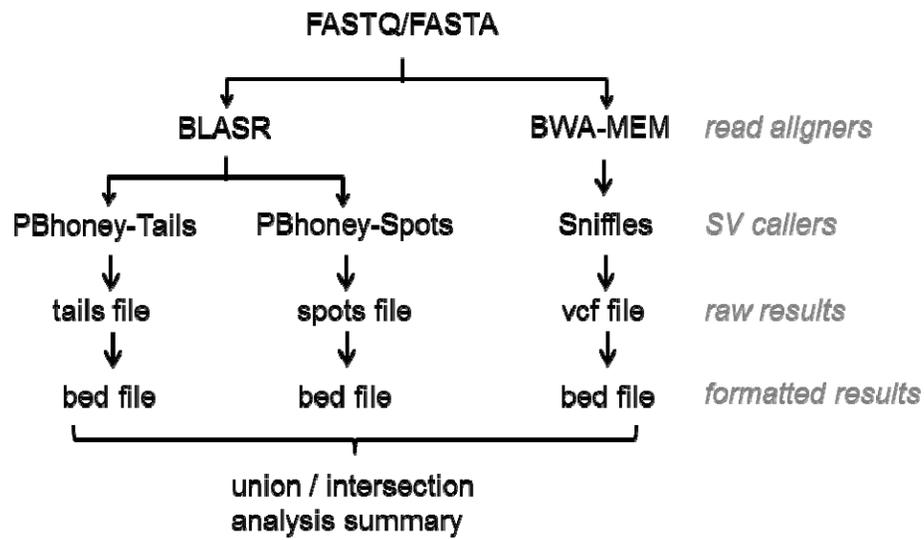
367

368

369

Figure 3. SV calling performance on the HX1 genome. Recalls and precisions of PBHoney-Tails, PBHoney-Spots, Sniffles and the intersection set of the latter two are shown.

370



371

372

Figure 4. Scheme of NextSV workflow.

373

374

375

376

377

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	18
SRX1424851	HX1	6~15X	7.0 kb	15
NIST	AJ son	10X	8.0 kb	14
NIST	AJ father	10X	7.3 kb	14
NIST	AJ mother	10X	7.8 kb	14

378

379

380

Table 2. Number of calls in gold standard SV set

Genome	Platform	Number of Deletions (≥ 200bp)	Number of Insertions (≥ 200bp)	Reference
NA12878	Illumina	2094	68	13
HX1	PacBio	2976	2944	15

381

382

383

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

	PBhoney-Tails	PBhoney-Spots	Sniffles	Union set
No. of calls (AJ father)	775	2944	2206	4020
No. of calls (AJ mother)	789	3091	2178	4165
No. of calls (AJ son)	728	3121	2198	4090
No. of calls inherited from father	370	1867	1006	2356
No. of calls inherited from mother	375	2095	987	2539
No. of ADI	282	441	814	937
ADI rate	38.6%	14.1%	37.0%	22.9%

384

385

386

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

	PBhoney-Tails	PBhoney-Spots	Sniffles	Union set
No. of calls (AJ father)	168	6691	1096	6952
No. of calls (AJ mother)	148	7183	1181	7476
No. of calls (AJ son)	151	7522	1148	7778
No. of calls inherited from father	104	2952	452	3897
No. of calls inherited from mother	87	3541	476	3986
No. of ADI	49	2721	479	2911
ADI rate	31.8%	36.2%	41.8%	37.4

387

388

389

390 **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	1.08	28.1

391

392