

# Evaluation on Efficient Detection of Structural Variants at Low Coverage by Long-Read Sequencing

Li Fang<sup>1</sup>, Jiang Hu<sup>1</sup>, Depeng Wang<sup>1</sup>, Kai Wang<sup>2,3,\*</sup>

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

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

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 on  
38 SVs 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 cancer have  
45 been associated with a large number of SVs in recent years [3-8]. Recent advances in  
46 next-generation sequencing (NGS) technologies have facilitated the analysis of  
47 variations such as SNPs and small Indels in unprecedented details, but the discovery of  
48 SVs using short reads still remains challenging [9]. Single-molecule, real-time (SMRT)  
49 sequencing developed by Pacific BioSciences (PacBio) offers a long read length,  
50 making it potentially well-suited for SV detection in personal genomes [9, 10]. Most  
51 recently, Merker et al. reported the application of low coverage whole genome PacBio  
52 sequencing to identify pathogenic structural variants from a patient with autosomal  
53 dominant Carney complex, for whom targeted clinical gene testing and whole genome  
54 short-read sequencing were negative [11].

55

56 Two SV software tools have been developed specifically for long-read sequencing:  
57 Pbhoney [12] and Sniffles [13]. Pbhoney identifies genomic variants via two algorithms,  
58 long-read discordance (Pbhoney-Spots) and interrupted mapping (Pbhoney-Tails).  
59 Sniffles is a SV caller written in C++ and it detects SVs using evidence from split-read  
60 alignments, high-mismatch regions, and coverage analysis. Due to the relative high cost  
61 of PacBio sequencing, users are often faced with issues such as what coverage is  
62 needed and how to get the best use of the available SV callers. In addition, it is unclear  
63 which software performs the best in low-coverage settings, and whether the  
64 combination of software tools can improve performance of SV calls. Finally, the  
65 execution of these software tools is often not straightforward and requires careful re-  
66 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 callers [14]. They also published sequencing data of seven human  
72 genomes, including PacBio data of an Ashkenazi Jewish family trio [15]. Previously, we  
73 sequenced a Chinese individual HX1 on the PacBio platform, and generated assembly-  
74 based SV call sets [16]. 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 BLASR [17] or BWA-MEM [18]. The bam files were  
93 down-sampled to different coverages using SAMtools (samtools view -s). The down-  
94 sampled 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) to increase sensitivity and discover SVs under low coverages (2~15X).

105

#### 106 ***SV detection using Sniffles***

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

112

#### 113 ***Comparing two SV call sets***

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

118

#### 119 ***Gold standard SV call set***

120 The gold standard SV call set for NA12878 was retrieved from the GIAB consortium [14],  
121 in which most of the calls were refined by experimental validation or other independent  
122 technologies. For the HX1 genome, we used the SV calls from a previously validated  
123 local assembly approach [10], as the initial high-quality calls. We also detected SVs on  
124 100X coverage PacBio data set of the HX1 genome using PBHoney-Tails, PBHoney-  
125 Spots and Sniffles. The initial high-quality calls that overlapped with one of the three  
126 100X call sets (PBHoney-Tails, PBHoney-Spots or Sniffles) were retained as final gold  
127 standard calls. SVs with length less than 200 bp were not considered. Number of SVs in  
128 the gold standard sets is shown in Table 2.

129

#### 130 ***Performance Evaluation of SV callers***

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

134 
$$\text{Precision} = \frac{TP}{TP+FP},$$

135 
$$\text{Recall} = \frac{TP}{TP+FN},$$

136 
$$F1 = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}},$$

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

141

## 142 **Results**

### 143 ***Performance of SV calling under different PacBio coverage***

144 To determine what sequencing coverage is needed for SV detection using PacBio data,  
145 we evaluated the performance of SV callers under several different coverages. We  
146 downloaded a recently published 22X PacBio data set of NA12878 [19] and down-  
147 sampled the data set to 2X, 4X, 6X, 8X, 10X, 12X, and 15X. SV calling was performed  
148 using PBHoney and Sniffles under each coverage. The resulting calls were compared  
149 with the gold standard SV set (including 2094 deletion calls and 68 insertion calls) from  
150 the Genome In A Bottle (GIAB) consortium [19].

151

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

161  
162 Second, we examined the precision and F1 scores of callers under different coverage.  
163 We calculated precision as the fraction of detected SVs that matching the gold standard  
164 set. As shown in Figure 1C, Sniffles has higher precision than PBHoney-Spots and  
165 PBHoney-Tails. The precision of Sniffles for deletion calls was 70% at 6X coverage, and  
166 decreased slightly as the coverage increased. F1 score, the harmonic mean of precision  
167 and recall, increased before 10X and then kept stable at higher coverage (Figure 1D).  
168 Precision for insertion calls was not assessed because there were only 86 insertion  
169 calls in the GIAB gold standard set, which was one order of magnitude smaller than the  
170 number of deletion calls, with potentially high false negative rates.  
171  
172 To verify the performance of SV detection on different individuals, we also did  
173 evaluation on a Chinese genome HX1, which was sequenced by us recently [16] at  
174 103X PacBio coverage. The genome was sequenced using a newer version of chemical  
175 reagents and thus the mean read length of HX1 was 40% longer than that of NA12878  
176 (Table 1). The total data set was down-sampled to 6X, 10X and 15X coverage. For each  
177 coverage data set, SVs were called and compared to the gold standard set. The results  
178 were similar to those of the NA12878 data set (Figure 3). At 10X coverage, 84%  
179 deletions and 92% insertions in the gold standard set can be detected by PBHoney-  
180 Spots. The precisions at 10X coverage range from 54% ~ 60% for deletion calls and 31%  
181 ~ 43% for insertion calls. At 15X coverage, the recall increased slightly but precision  
182 decreased. Thus, 10X may be an optimal coverage to use in practice, considering the  
183 sequencing costs and the balance of recall and precision.  
184  
185 ***Performance of SV calling using a combination of PBHoney and Sniffles***  
186 Although PBHoney-Spots detected most of the variants, we examined whether we can  
187 improve the recall rates by running both PBHoney-Spots and Sniffles, especially under  
188 low fold coverages. As shown in Figure 2, at 6X coverage, the union set of both callers  
189 discovered 77% deletions in the NA12878 gold standard set, which was 23% more than  
190 running PBHoney-Spots alone at 6X coverage and comparable to running PBHoney-

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

193

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

199

#### 200 ***Evaluation on Mendelian Errors***

201 As the germline mutation rate is very low [20, 21], Mendelian errors are more likely a  
202 result of genotyping errors and can be used as a quality control criteria in genome  
203 sequencing [22]. Here, we evaluated the errors of allele drop-in (ADI), which means  
204 that an offspring presents an allele that does not appear in either parent, using a whole  
205 genome sequencing data set of an Ashkenazi Jewish (AJ) family trio released by NIST  
206 [15]. The sequencing data of AJ son, AJ father and AJ mother was down-sampled to  
207 10X coverage. SVs were called using PBHoney-Tails, PBHoney-Spots and Sniffles. The  
208 calls from AJ son were compared with calls from AJ father and AJ mother. ADI rate was  
209 calculated as the proportion of calls in offspring not matching any call from either parent.  
210 The result shows that PBHoney-Spots returns the most calls. For deletion calls,  
211 PBHoney-Spots gives us a lowest ADI rate (14.1%), while the ADI rates for insertion  
212 calls are considerable higher (31.8% ~ 41.8). Therefore, further validation or manual  
213 inspection of the calls is needed when analyzing SVs that may be associated with  
214 diseases with low coverage sequencing.

215

#### 216 ***Automated pipeline for SV calling using PBhoney and Sniffles***

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

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

237

### 238 ***Computational Performance of NextSV***

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

251

252

253 **Discussion**

254 Depth of coverage is often a key consideration in genomic analyses [23]. In this study,  
255 we evaluated SV calling performance of three SV calling algorithms, PBHoney-Tails,  
256 PBHoney-Spots and Sniffles, at various PacBio coverages of 2 ~ 15X. Our results  
257 showed that, at 10X coverage, 76% ~ 84% deletions and 80% ~ 92 % insertions were  
258 detected by running PBHoney-Spots. By running both PBHoney-Spots and Sniffles,  
259 comparable recall can be achieved at coverage as low as 6X. At more than 10X  
260 coverage, the recall slightly increased. Thus, 10X can be an optimal PacBio coverage  
261 for efficient SV detection, yet 6X may also be an economic choice under limited budget.  
262

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

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

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

## 290 Acknowledgments

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

## 296 Competing interests

297 L.F. and D.W. are employees and K.W. is a consultant for Grandomics Biosciences.  
298

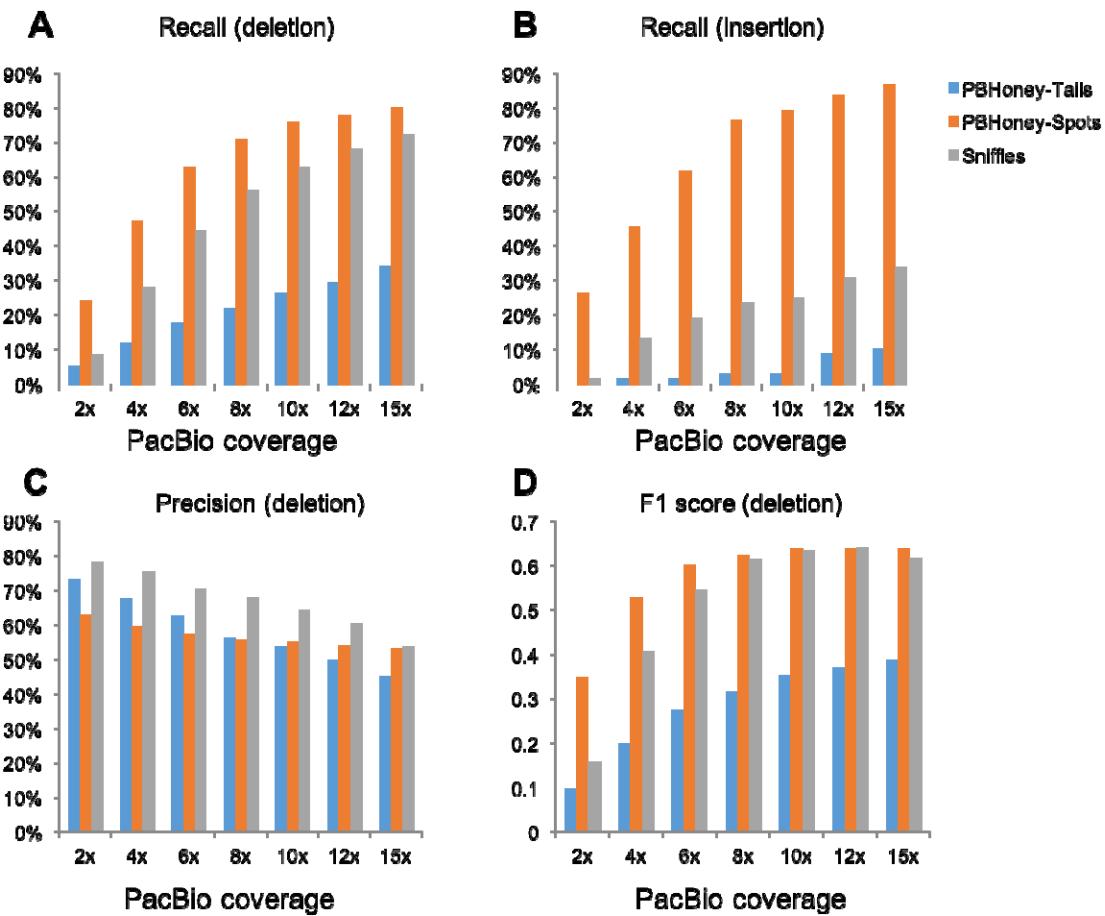
## 299 References

- 300 1. Feuk L, Carson AR, Scherer SW: **Structural variation in the human genome.** *Nature reviews Genetics* 2006, **7**(2):85-97.
- 301 2. Pang AW, MacDonald JR, Pinto D, Wei J, Rafiq MA, Conrad DF, Park H, Hurles ME, Lee C, Venter JC et al: **Towards a comprehensive structural variation map of an individual human genome.** *Genome biology* 2010, **11**(5):R52.
- 302 3. Stankiewicz P, Lupski JR: **Structural variation in the human genome and its role in disease.** *Annual review of medicine* 2010, **61**:437-455.
- 303 4. Weischenfeldt J, Symmons O, Spitz F, Korbel JO: **Phenotypic impact of genomic structural variation: insights from and for human disease.** *Nature reviews Genetics* 2013, **14**(2):125-138.
- 304 5. Yang L, Luquette LJ, Gehlenborg N, Xi R, Haseley PS, Hsieh CH, Zhang C, Ren X, Protopopov A, Chin L et al: **Diverse mechanisms of somatic structural variations in human cancer genomes.** *Cell* 2013, **153**(4):919-929.
- 305 6. Moncunill V, Gonzalez S, Bea S, Andrieux LO, Salaverria I, Royo C, Martinez L, Puiggros M, Segura-Wang M, Stutz AM et al: **Comprehensive characterization of complex structural variations in cancer by directly comparing genome sequence reads.** *Nature biotechnology* 2014, **32**(11):1106-1112.
- 306 7. Zhang F, Gu W, Hurles ME, Lupski JR: **Copy number variation in human health, disease, and evolution.** *Annual review of genomics and human genetics* 2009, **10**:451-481.
- 307 8. Carvalho CM, Lupski JR: **Mechanisms underlying structural variant formation in genomic disorders.** *Nature reviews Genetics* 2016, **17**(4):224-238.
- 308 9. English AC, Salerno WJ, Hampton OA, Gonzaga-Jauregui C, Ambreth S, Ritter DL, Beck CR, Davis CF, Dahdouli M, Ma S et al: **Assessing structural variation in a personal genome-towards a human reference diploid genome.** *BMC genomics* 2015, **16**:286.
- 309 10. Chaisson MJ, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, Antonacci F, Surti U, Sandstrom R, Boitano M et al: **Resolving the complexity of the human genome using single-molecule sequencing.** *Nature* 2015, **517**(7536):608-611.

- 326 11. Merker J, Wenger AM, Sneddon T, Grove M, Waggott D, Utiramerur S, Hou Y, Lambert CC, Eng  
327 KS, Hickey L et al: **Long-read whole genome sequencing identifies causal structural variation in a**  
328 **Mendelian disease.** *bioRxiv* 2016.
- 329 12. English AC, Salerno WJ, Reid JG: **PBHoney: identifying genomic variants via long-read**  
330 **discordance and interrupted mapping.** *BMC Bioinformatics* 2014, **15**:180.
- 331 13. <https://github.com/fritzsedlazeck/Sniffles>
- 332 14. Parikh H, Mohiyuddin M, Lam HY, Iyer H, Chen D, Pratt M, Bartha G, Spies N, Losert W, Zook  
333 JM et al: **svclassify: a method to establish benchmark structural variant calls.** *BMC genomics* 2016,  
334 **17**:64.
- 335 15. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, Weng Z, Liu Y, Mason CE, Alexander  
336 N et al: **Extensive sequencing of seven human genomes to characterize benchmark reference**  
337 **materials.** *Scientific data* 2016, **3**:160025.
- 338 16. Shi L, Guo Y, Dong C, Huddleston J, Yang H, Han X, Fu A, Li Q, Li N, Gong S et al: **Long-read**  
339 **sequencing and de novo assembly of a Chinese genome.** *Nature communications* 2016, **7**:12065.
- 340 17. Chaisson MJ, Tesler G: **Mapping single molecule sequencing reads using basic local**  
341 **alignment with successive refinement (BLASR): application and theory.** *BMC Bioinformatics* 2012,  
342 **13**(1):238.
- 343 18. Li H: **Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.**  
344 *arXiv* 2013, **1303.3997v1** [q-bio.GN].
- 345 19. Pendleton M, Sebra R, Pang AW, Ummat A, Franzen O, Rausch T, Stutz AM, Stedman W,  
346 Anantharaman T, Hastie A et al: **Assembly and diploid architecture of an individual human genome**  
347 **via single-molecule technologies.** *Nature methods* 2015, **12**(8):780-786.
- 348 20. Kong A, Frigge ML, Masson G, Besenbacher S, Sulem P, Magnusson G, Gudjonsson SA,  
349 Sigurdsson A, Jonasdottir A, Jonasdottir A et al: **Rate of de novo mutations and the importance of**  
350 **father's age to disease risk.** *Nature* 2012, **488**(7412):471-475.
- 351 21. Veltman JA, Brunner HG: **De novo mutations in human genetic disease.** *Nature reviews*  
352 *Genetics* 2012, **13**(8):565-575.
- 353 22. Pilipenko VV, He H, Kurowski BG, Alexander ES, Zhang X, Ding L, Mersha TB, Kottyan L, Fardo  
354 DW, Martin LJ: **Using Mendelian inheritance errors as quality control criteria in whole genome**  
355 **sequencing data set.** *BMC proceedings* 2014, **8**(Suppl 1 Genetic Analysis Workshop 18Vanessa  
356 Olmo):S21.
- 357 23. Sims D, Sudbery I, Ilott NE, Heger A, Ponting CP: **Sequencing depth and coverage: key**  
358 **considerations in genomic analyses.** *Nature reviews Genetics* 2014, **15**(2):121-132.
- 359
- 360
- 361

362 **Figure and Tables**

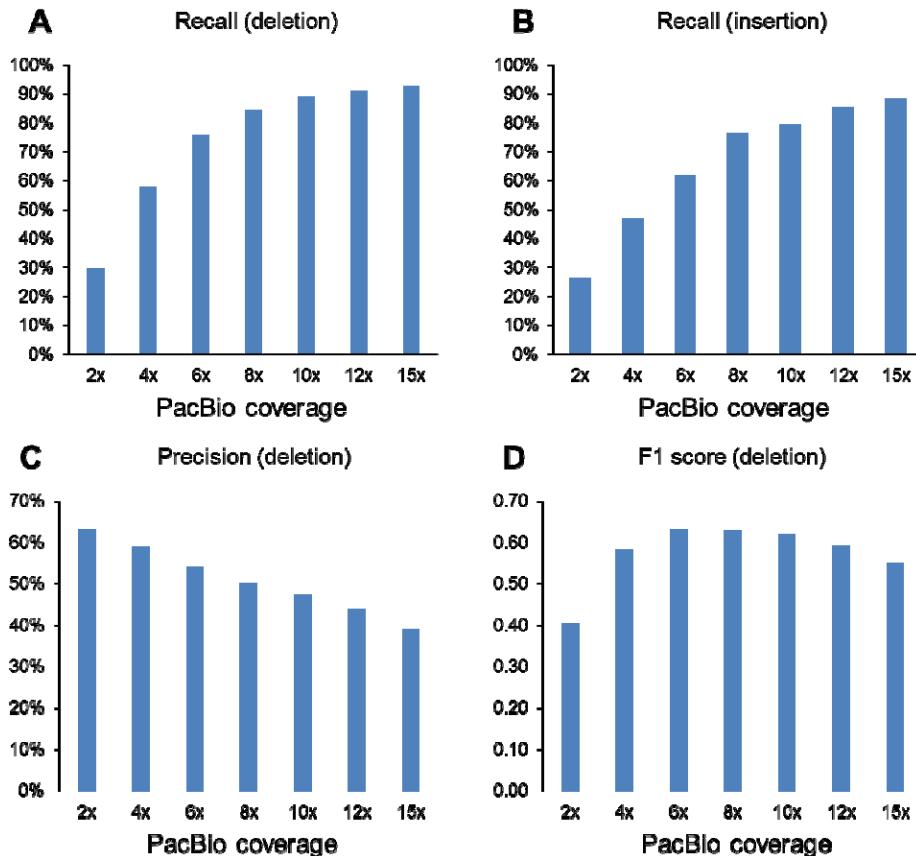
363



364

365

366 Figure 1. SV calling performance for each SV caller under different coverage on the  
367 NA12878 genome.

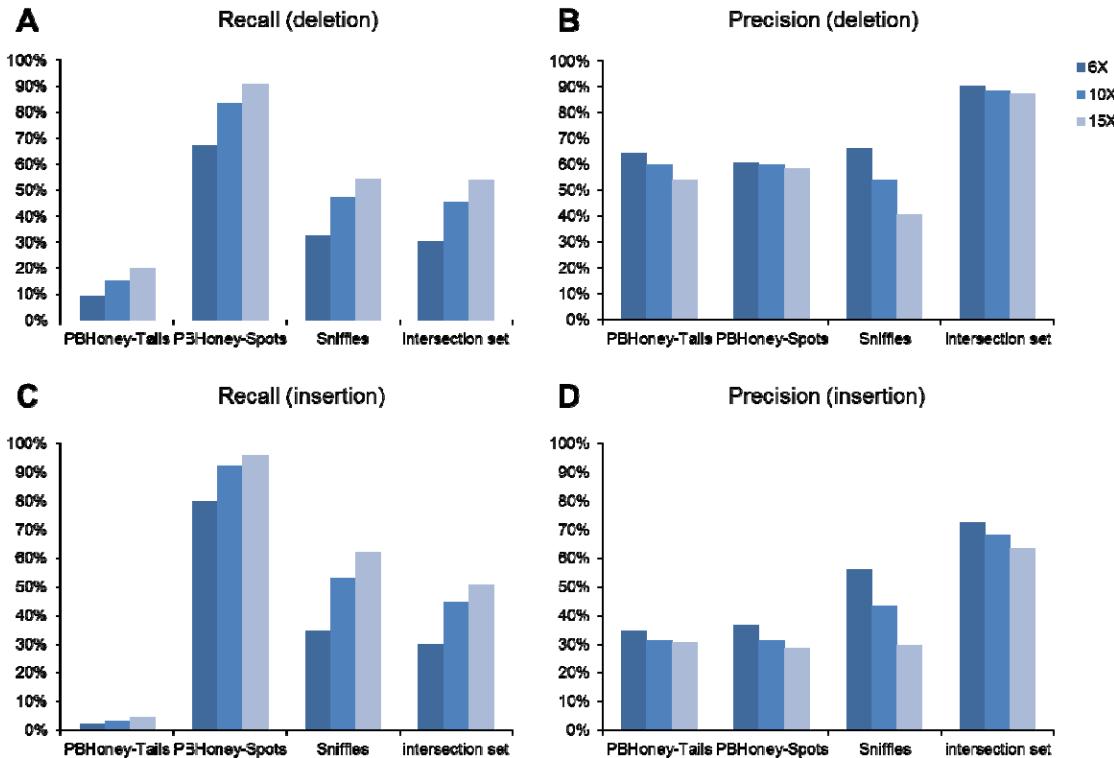


368

369

370     Figure 2. SV calling performance for the union call set of PBHoney-Spots, PBHoney-Tails and Sniffles under different coverage on the NA12878 genome.  
371

372



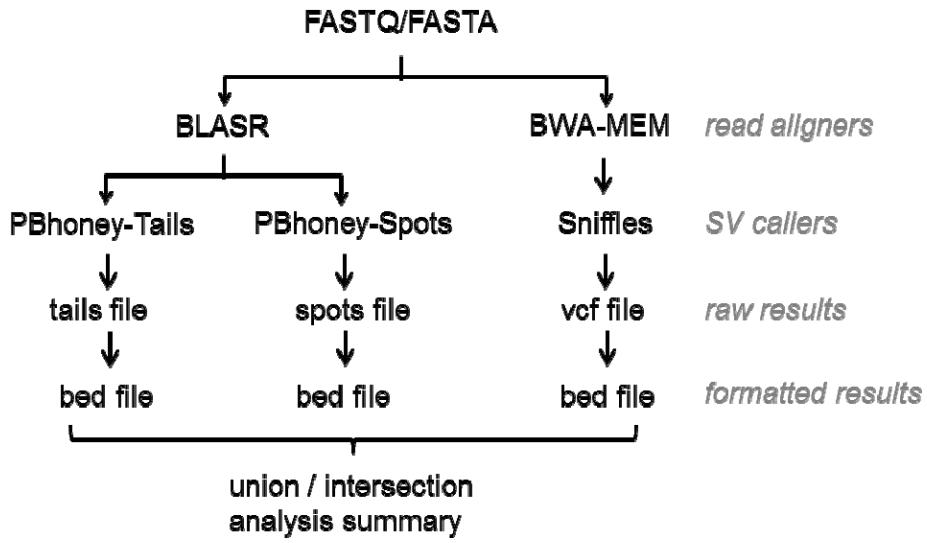
373

374

375

Figure 3. SV calling performance on the HX1 genome.

376



377

378

Figure 4. Scheme of NextSV workflow.

379

380

381

382

383

**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	[19]
SRX1424851	HX1	6~15X	7.0 kb	[16]
NIST	AJ son	10X	8.0 kb	[15]
NIST	AJ father	10X	7.3 kb	[15]
NIST	AJ mother	10X	7.8 kb	[15]

384

385

386

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

Genome	Platform	Number of Deletions ( $\geq$ 200bp)	Number of Insertions ( $\geq$ 200bp)	Reference
NA12878	Illumina	2094	68	[14]
HX1	PacBio	2976	2944	[16]

387

388

389

**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%

390

391

392

**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

393

394

395

396 **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

397

398