

# GigaScience

## CNVcaller: High efficient and Widely Applicable Software for Detecting Copy Number Variations in large Populations --Manuscript Draft--

<b>Manuscript Number:</b>	GIGA-D-17-00119	
<b>Full Title:</b>	CNVcaller: High efficient and Widely Applicable Software for Detecting Copy Number Variations in large Populations	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	National Natural Science Foundation of China (31572381)	Prof. Yu Jiang
	National Thousand Youth Talents Plan	Prof. Yu Jiang
<b>Abstract:</b>	<p>Background: The increasing sequencing data of a wide variety of species offers an opportunity for copy number variation (CNV) detection at population level. However, the growing sample size and the divergent complexity of non-human genomes challenge the efficiency and robustness of the current human-oriented CNV detection methods.</p> <p>Result: Here we present CNVcaller, a read depth based method for CNVs discovering of the population sequencing data. By the statistics-based signal detection and population-level noise reduction algorithms, the detection for 232 goats with complicated genome assembly takes only 1.4 days on a single compute node. Besides, the false segmental duplications in reference genome assemblies can be mitigated by a simplified absolute copy number correction, which consumes only a few minutes and increases the sensitivity in CNV enriched regions. Multiple validations showed that CNVcaller achieved increased total sensitivity, high genotyping accuracy and low false discovery rate in human, livestock and crop populations.</p> <p>Conclusion: The fast and general detection algorithms of CNVcaller overcome prior computational barriers for detecting CNV from large scale sequencing data with complicated genome structure. These advantages will promote the population genetic analysis of functional CNVs of more species.</p>	
<b>Corresponding Author:</b>	Yu Jiang, Ph.D Northwest Agriculture and Forestry University Yangling, Shaanxi CHINA	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Northwest Agriculture and Forestry University	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Xihong Wang	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Xihong Wang	
	Zhuqing Zheng	
	Yudong Cai	
	Ting Chen	
	Chao Li	
	Weiwei Fu	
	Yu Jiang	
<b>Order of Authors Secondary Information:</b>		
<b>Opposed Reviewers:</b>		

Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	Yes

---

# CNVcaller: High Efficient and Widely Applicable Software for Detecting Copy Number Variations in large Populations

Xihong Wang<sup>1†</sup>, Zhuqing Zheng<sup>1†</sup>, Yudong Cai<sup>1</sup>, Ting Chen<sup>1</sup>, Chao Li<sup>1</sup>, Weiwei Fu<sup>1</sup>, Yu Jiang<sup>1\*</sup>

<sup>1</sup> College of Animal Science and Technology, Northwest A&F University, Yangling, Shaanxi,  
China

<sup>†</sup> These authors contributed equally to this work.

\*Correspondence should be addressed to Yu Jiang (yu.jiang@nwafu.edu.cn).

## Abstract

**Background:** The increasing sequencing data of a wide variety of species offers an opportunity for copy number variation (CNV) detection at population level. However, the growing sample size and the divergent complexity of non-human genomes challenge the efficiency and robustness of the current human-oriented CNV detection methods.

**Result:** Here we present CNVcaller, a read depth based method for CNV discovering of the population sequencing data. By the statistics-based signal detection and population-level noise reduction algorithms, the detection for 232 goats with complicated genome assembly takes only 1.4 days on a single compute node. Besides, the false segmental duplications in reference genome assemblies can be mitigated by a simplified absolute copy number correction, which consumes only a few minutes and increases the sensitivity in CNV enriched regions. Multiple validations showed that CNVcaller achieved increased total sensitivity, high genotyping accuracy and low false discovery rate in human, livestock and crop populations.

**Conclusion:** The fast and general detection algorithms of CNVcaller overcome prior

---

23 computational barriers for detecting CNVs from large scale sequencing data with complicated  
24 genome structure. These advantages will promote the population genetic analysis of functional  
25 CNVs of more species.

## 27 **Keywords**

28 copy number variation (CNV), next-generation sequencing (NGS), population genetics, segmental  
29 duplication, absolute copy number.

## 31 **Introduction**

32 Copy number variants (CNVs) are the prevalent and important source of genetic diversity [1],  
33 which are highly correlated with diseases [2, 3], evolutions [4] and other phenotypes [5-8] for all  
34 kinds of species. Over the development for decades, the large-scale sequencing projects have  
35 provided us with enormous amount of data across the tree of life. The geometric growing sample  
36 size enables the population genetics variant association studies using the CNV regions (CNVRs)  
37 integrated from multi-sample CNVs [9, 10]. However, the increasing data size aggravates  
38 computational burden and challenges the efficiency of the current CNV detectors. In addition, the  
39 complicated genome structure of many non-human species demands more robust signal detection  
40 and noise reduction algorithms.

41 Currently, several strategies are used to for CNV detecting of whole genome sequencing data:  
42 read-pair/split-read [11-14], local assembly [15-17] and read-depth (RD) [18-20]. Although  
43 employing multiple methods in one dataset can increase the total sensitivity [21], the efficiency  
44 and convenience would consequently become the subsequent concern. With the increasing release

---

1 45 of large-scale sequencing data, the population genetic information is applied to improve the  
2  
3 46 detecting accuracy [22]. A typical strategy is to simultaneously scan the genomes of multiple  
4  
5  
6 47 samples, then decompose all signals into true variations and noises by priori distributions.  
7  
8  
9 48 Genome STRiP [23] is shown to be one of the best population-level CNV detectors in 1000  
10  
11  
12 49 Genome Project [24].

13  
14 50 Current human-oriented CNV detectors leave some uncertain points for application to the  
15  
16  
17 51 other species. Firstly, gaps and unplaced scaffolds are riddled with reference genome assemblies  
18  
19  
20 52 of most non-model organism [25, 26], leading to the increase of the abnormal mapping and the  
21  
22  
23 53 false positive rate of the read-pair/split-read algorithms. In comparison, the RD algorithm deduces  
24  
25  
26 54 copy number from the number of reads aligned to of a particular region, which can efficiently  
27  
28  
29 55 screen out noises by statistical hypothesis [18, 27]. In the RD based methods, CNVnator [19]  
30  
31  
32 56 which provides multi-sample genotyping function was used in yak, chicken and fish cohorts [28-  
33  
34 57 30]. Secondly, the alternative alleles lead to high-proportioned erroneous segmental duplications  
35  
36  
37 58 (SDs) for the animal reference genomes [31, 32]. Therefore, intensive filtering of the duplicated  
38  
39  
40 59 regions on the reference genome is recommended by many CNV detectors. However, the SDs  
41  
42  
43 60 enrich CNVs 10 times than other area of the genome [4] and contribute to the evolutionary  
44  
45  
46 61 adaptive traits [33]. A more precise solution is deducing the absolute copy number from mrsFAST  
47  
48  
49 62 alignment which reports multiple hits of a single read [34]. However, precise realignment always  
50  
51  
52 63 requires enormous time, especially for the crop plant genomes which frequently contain large  
53  
54  
55 64 duplicated regions. Therefore, this strategy was hired in very few non-human CNV researches  
56  
57  
58 65 [35].

59 66 In this study, we introduce a super-fast and generally applicable method, CNVcaller, for

---

1 67 CNV discovering sequencing data of large populations. This software is based on the RD  
2  
3 68 algorithm, and implies robust signal detection and noise deduction methods to increase the  
4  
5  
6 69 computational efficiency in all kinds of genomes. We applied it to the population sequencing data  
7  
8  
9 70 of human, livestock and crop to demonstrate the utility and benchmarked against the widely used  
10  
11  
12 71 and best practice CNV detectors.

13  
14  
15 72

## 18 73 **Materials and Methods**

### 22 74 **Input data**

23  
24  
25  
26 75 The main input of CNVcaller are the alignment files in BAM format. The following data/samples  
27  
28  
29 76 were included in the validation: 30 BAM files of human from the 1000 Genome Project Phase 3  
30  
31  
32 77 [36], including 27 normal (~ 12 X) and three deeply sequenced samples (~ 50 X); 30 BAM files of  
33  
34  
35 78 10 families from the Genomes of Netherlands (GoNL) project [37] (~ 20 X); 70 FASTQ files of  
36  
37  
38 79 domestic sheep samples (~ 10 X) from the NCBI BioProject: PRJNA160933; two maize [38] and  
39  
40  
41 80 two soybeans [8] FASTQ files (each species contain one ~5 X and one ~10 X sample). An  
42  
43  
44 81 additional table showed the downloaded files in detail (**Supplementary Table1**).

45 82 Another 63 sheep (~10 X) and 232 goats (~12 X) data were sequenced using pair-end  
46  
47  
48 83 libraries on the Illumina HiSeq 4000 platform. The FASTQ files were aligned to their respective  
49  
50  
51 84 reference assemblies using BWA 0.7.13 to generate BAM files [35]. The version of reference  
52  
53  
54 85 genomes are: human GRCh37, maize B73 RefGen\_v3, soybean Glycine\_max\_v2.0, sheep  
55  
56  
57 86 OAR\_v3.1 and goat ARS1. The GATK v3.5 [39] pre-processing workflow is used to produce  
58  
59  
60 87 analysis-ready BAM files. After alignment, PCR duplications were marked by Picard 2.1 and the

---

1 88 realignment was performed by GATK. The reads with 0x504 flag (indicating unmapped,  
2  
3 89 secondary mapped or PCR duplication) were removed.  
4  
5  
6

7 **90 Individual RD processing**  
8  
9

10 91 *RD Estimation.* The reference genome is segmented into overlapping sliding windows. For 5-10  
11  
12  
13 92 X sequencing data, 800 bp windows with a 400 bp overlap is recommended. The sliding windows  
14  
15  
16 93 with gaps are excluded from the computation. The windows are indexed to form a reference database  
17  
18  
19 94 which will be used in all samples. The BAM file of each individual is parsed out using SAMtools  
20  
21  
22 95 v1.3 [40]. The raw RD signal is calculated for each window as the number of placed reads with  
23  
24  
25 96 centers within window boundaries. This step consumes less than 500 Mb max memory for one BAM  
26  
27  
28 97 file, so parallel submitting is recommended.  
29

30 98  
31  
32

33 99 *Absolute copy number correction.* The standard mapping only aligns one sequencing read to one  
34  
35  
36 100 best position of the genome. For the regions with more than one assembled copy in reference  
37  
38  
39 101 genome, the reads will be split among the copies. Therefore, the deduced copy numbers are  
40  
41  
42 102 dependent to the segment number in reference genome, which is called relative copy number.  
43  
44  
45 103 CNVcaller implies a simple correction to deduce the copy number independent to the reference  
46  
47 104 genome, which was called absolute copy number.  
48

49  
50 105 To perform the absolute copy number correction, the windows with >97% sequence similarity  
51  
52  
53 106 are linked together to form a duplicated window record file before correction. This file can be  
54  
55  
56 107 generated by splitting the reference genome into non-overlapping windows and aligning them  
57  
58  
59 108 onto the reference genome using the precise aligner, e.g. BLAT v. 36X1 [41]. The windows with  
60  
61  
62  
63  
64  
65

---

109 more than 20 hits are excluded to remove the low complexity regions. The record files of human,  
110 livestock and main crops can be downloaded from the CNVcaller website.

111 Based on the duplicated window record file, the raw RD located on similar windows are added  
112 together to generate the absolute RD for all the high similarity windows.

$$113 \quad RD_{absolute}^i = \sum_{j=1}^t RD_{raw}^{ij}$$

114 Where  $i$  is the index of the window to be corrected,  $t$  is the total number of the high similarity  
115 windows.  $RD_{raw}^{ij}$  is the raw RD of the window similar with the  $i$ -th window (including the  $i$ -th  
116 window itself), which is counted directly from the BWA alignment, and  $RD_{absolute}^i$  is the  
117 corrected RD of the  $i$ -th window which can be used to deduce the absolute copy number.

118  
119 *GC correction and normalization.* Since the resequencing samples may show various GC content  
120 distribution, the GC bias is corrected individually basically as CNVnator [19] except using the RD  
121 of the windows with 40% GC as standard:

$$122 \quad RD_{corrected}^i = \frac{\overline{RD}_{40}}{\overline{RD}_{gc}} RD_{absolute}^i$$

123 Where  $i$  is the window index,  $RD_{absolute}^i$  is the RD after absolute copy number correction,  
124  $RD_{corrected}^i$  is final corrected RD for the window,  $\overline{RD}_{40}$  is the mean RD of windows with 40%  
125 percent GC as standard, and  $\overline{RD}_{gc}$  is the mean RD over all windows that have the same GC  
126 content with the  $i$ -th window.

127 Assuming that the majority part of the genome is in normal copy number, the corrected RDs are  
128 divided by the global median RD to normalized to one. For the sex chromosomes, if the median  
129 RD of the homogametic sex chromosomes (X or Z) is about half of the median RD of autosome,



---

1 130 the RDs on the X or Z chromosome are doubled before normalization.

2  
3 131

4  
5  
6  
7 132 **CNVR detection by multiple criteria**

8  
9  
10  
11 133 *Individual candidate CNV window definition.* The individual candidate CNV windows are defined

12  
13  
14 134 using two criteria: (1) The normalized RD is significantly higher or lower than the normalized

15  
16 135 mean RD (deletions  $< 1 - 2 * STDEV$ ; duplications  $> 1 + 2 * STDEV$ ). (2) Considering the

17  
18 136 normalized RD of heterozygous deletions and duplications should be around 0.5 and 1.5

19  
20 137 respectively, an empirical standard for the normalized RD (deletions  $< 0.65$ ; duplications  $> 1.35$ )

21  
22 138 also need to be achieved. For some strictly self-bred species, such as soybean and wheat, this

23  
24 139 empirical standard should be raised to 0.25 or 1.75 for the normalized RD of the homozygous

25  
26 140 deletions or duplications respectively.

27  
28  
29  
30  
31  
32 141

33  
34  
35  
36 142 *Population-level candidate CNV window definition.* All individual RD files are piled up by the

37  
38 143 universal window index to a two-dimensional population RD file. The window showing high

39  
40 144 frequency of individually candidate CNV (allele frequency  $> 0.05$ ) or have at least three

41  
42 145 homozygous duplicated/deleted individuals in large population are selected. Then Pearson's

43  
44 146 product-moment correlation coefficients of the multi-sample RDs are calculated between the two

45  
46 147 adjacent non-overlapping windows. Only the windows with significant correlation ( $P < 0.01$  by T

47  
48 148 test) are selected and merged into one call.

49  
50  
51  
52 149

53  
54  
55  
56 150 *CNV region definition.* Initial calls are selected if more than four sequential 800 bp overlapped

---

1 151 windows (total length  $\geq 2,000$  bp) are defined as the population-level candidate windows. To  
2  
3 152 tolerant noises, at most one unselected window out of four continuous candidate windows is  
4  
5  
6 153 allowed to exist. Then the two adjacent initial calls are further merged if their copy numbers are  
7  
8  
9 154 highly correlated and the distance between them is less than a certain percent of their own length.

10  
11 155

## 12 156 **CNVR Genotyping**

13  
14  
15  
16  
17  
18  
19 157 The copy number of a specific sample is initially estimated by two times the median RD of all the  
20  
21  
22 158 candidate windows in this region. The deleted and biallelic duplicated CNVRs (average copy  
23  
24  
25 159 number  $\leq 4$ ) will be clustered by a constrained mixture Gaussian model embedded in CNVcaller.  
26  
27  
28 160 This model presets the average copy number of homozygous deletion, heterozygous deletion,  
29  
30  
31 161 normal, heterozygous deletion and homozygous deletion at zero to four respectively. For  
32  
33 162 multiallelic CNVRs (average copy number  $>4$ ) we provide a clustering process by unsupervised  
34  
35  
36 163 mixture Gaussian model (calling R package mclust 5.2 [42]). In a population, the calls with the  
37  
38  
39 164 same copy number in all samples are defined as SDs while the polymorphic calls are defined as  
40  
41  
42 165 CNVRs. The output CNVR genotyping file is analyzable by the population genetic algorithms.

43  
44 166

## 45 46 47 167 **Performance evaluation**

48  
49  
50  
51  
52 168 *Competing methods.* Most validations were based on the 30 human BAM files from the 1000  
53  
54  
55 169 Genome Project Phase 3 unless otherwise noted. The performance of CNVcaller was compared  
56  
57  
58 170 with two pipelines: CNVnator\_v0.3.3 [19] which was well-used in animal population CNVR

---

1 171 detection and Genome STRiP (included in svtoolkit\_2.00.1696) [23] which was the state-of-the-  
2  
3 172 art human population CNV detector. The recommended parameters and QC filters were used. For  
4  
5  
6 173 Genome STRiP, both the deletion and CNV pipelines were performed. The unplaced scaffolds  
7  
8  
9 174 were removed from the reference genome and the whole genome was separated by chromosome  
10  
11  
12 175 as recommended. The standard screens were applied to select passing sites and remove duplicated  
13  
14 176 calls. For CNVnator, the gap regions and calls with p values less than 0.01 were removed. We also  
15  
16  
17 177 used the q0 filter to remove any predictions with  $q_0 < 0.5$  (reads with multiple mapping locations)  
18  
19  
20 178 as recommend. The individual CNVs of all sample were merged in to the population CNVRs by  
21  
22  
23 179 the arbitrary standards: two calls have  $>50\%$  reciprocal overlapping with each other or  $>90\%$  of  
24  
25  
26 180 on one call is covered by another call. Then the CNVRs were genotyped by the built-in function of  
27  
28  
29 181 CNVnator. Because the three software have different limitations in CNV detection, only the  
30  
31  
32 182 CNVRs on autosomes with  $>2,000$  bp length and allele frequency  $\geq 0.05$  were used in the  
33  
34  
35 183 following validation.  
36  
37  
38  
39 184  
40  
41 185 *Sensitivity validation.* Sensitivity was calculated as the proportion of high-confident CNVR  
42  
43  
44 186 database overlapped by predicted CNVRs. Two previously published database including the same  
45  
46  
47 187 samples in the test data were used. One is the 1000 Genome Project CNVR map [24] included 26  
48  
49  
50 188 tested samples, the other is the array comparative genomic hybridization, (aCGH) based CNVR  
51  
52  
53 189 database [43] included 10 tested samples. The CNVRs of the specific samples were extracted from  
54  
55  
56 190 the database then screened by the same length and frequency as detected CNVRs (length  $>2,000$   
57  
58  
59 191 bp and alternative allele frequency  $\geq 0.05$ ). The intersected length of the predicted CNVRs and  
60  
61  
62 192 the high-confident CNVR database were calculated by the bedtools v2.25.0 [44].  
63  
64  
65

---

1 193

2  
3 194 *Accuracy validation using human database.* The intensity rank-sum (IRS) test (included in the  
4  
5  
6 195 svtoolkit\_2.00.1696) was performed based on the intensity data of the Affymetrix SNP 6.0 array  
7  
8  
9 196 including 26 test samples. SD regions were removed as [23] because the probe design does not  
10  
11  
12 197 cover the high similarity regions. The genotyping accuracy were calculated based on the aCGH  
13  
14 198 CNVR database [43]. For a detected CNVR has >90% overlap with the database, the predicted  
15  
16  
17 199 copy number showed exact agreement with the integer genotyping from aCGH database were  
18  
19  
20 200 defined as correct. The Mendelian inconsistencies were calculated from the deleted and biallelic  
21  
22  
23 201 duplicated CNVRs (average copy number < 4) in the Dutch families and sheep trios.  
24

25 202

26  
27  
28 203 *Sheep genotyping validation by CNVplex assay.* A total of 73 sheep including Merino, Texel,  
29  
30  
31 204 Mongolia and Tibetan sheep were used for genotyping validation. Genomic DNA was extracted  
32  
33  
34 205 from the peripheral blood using the QIAamp DNA blood mini kit (Qiagen, Germany). ~10 X  
35  
36 206 resequencing was performed for each sheep and the CNVRs were detected by CNVcaller as  
37  
38  
39 207 described above. The predicted CNVRs with high variation frequency were selected for the  
40  
41  
42 208 validation. The copy numbers were validated by CNVplex® (Genesky Biotechnologies Inc.,  
43  
44  
45 209 Shanghai, China), which is based on double ligation and multiplex fluorescence PCR [45]. The  
46  
47  
48 210 probes were designed to target the candidate windows of the target CNVR. The sizes of the PCR  
49  
50  
51 211 fragments and target loci sequences in each reaction are listed in **Supplementary Table 2.**  
52  
53 212 Amplified probes were detected as fluorescent signals and peak areas were compared and  
54  
55  
56 213 normalized to determine the dosage of each target.  
57

58 214

---

1     **215 Absolute copy number validation**

2  
3  
4     **216 Putative X-linked scaffolds.** All the scaffolds of OAR v3.1 were mapped to the X chromosome of  
5  
6  
7     **217** sheep reference genome OAR v4.0, goat reference genome ARS1, and cattle reference genome  
8  
9  
10    **218** UMD 3.1 using BLASR [46]. If the best hit of a scaffold had a coverage >50% with >90% identity  
11  
12  
13    **219** and >3 Kb length, this scaffold was defined as the putative X-linked scaffold.  
14

15     **220**

16  
17  
18    **221 mrsFAST alignment.** The pair-end reads with multiple hits indicated by the XA tag in BWA  
19  
20  
21    **222** alignment were selected to realign by mrsFAST\_v3.3.10 [34]. The mrsFAST alignment was  
22  
23  
24    **223** performed basically as previously described [47]. Longer reads were trimmed into 40 bp to reduce  
25  
26  
27    **224** the read length heterogeneity prior to sequence alignments. After alignment, the reads with more  
28  
29  
30    **225** than 20 mapped hits were excluded to remove the low complexity regions.  
31

32     **226**

33  
34  
35    **227 Results**

36  
37  
38    **228 Overview of CNVcaller algorithm**

39  
40  
41  
42    **229** CNVcaller pipeline includes three main steps (**Figure 1**). First, considering the population  
43  
44  
45    **230** sequencing data may come from different platforms, the RD of each sample is counted and  
46  
47  
48    **231** corrected individually. An original absolute copy number correction is used to modify the standard  
49  
50  
51    **232** read alignments generated by BWA software to multi-hit alignments, as similar to mrsFAST  
52  
53  
54    **233** format (**Supplementary Figure 1**). This correction takes only 0.06 core-hour for a mammalian  
55  
56  
57    **234** genome with 10 X sequencing coverage, while 10 core-hours are needed for remapping the reads  
58  
59  
60    **235** by mrsFAST. After corrections and normalization, the comparable RDs of each sample is  
61  
62  
63  
64  
65

---

1 236 concentrated to a ~100 Mb intermediate file and output. This design avoids repeat calculation of a  
2  
3 237 same individual in different populations, and save much time since the individual step consumes  
4  
5  
6 238 more than 80% of the total running time.  
7

8  
9 239 In the second CNVR detection step, the RD files of all samples are piled up into a two-  
10  
11 240 dimensional population RD file. Multi-criteria are implied to remove the high-proportional noise  
12  
13 241 caused by low sequencing quality or assembly bias. Individually, the RD of the candidate CNV  
14  
15  
16 242 window should significantly deviates from average. The piled-up candidate windows should also  
17  
18  
19 243 meet two population-level criteria: CNV allele frequency > 5% and the multi-sample RDs of  
20  
21  
22 244 adjacent windows are significantly correlated (**Figure 1**). Compared with intensified individual  
23  
24  
25 245 RD screening, the multi-criteria filtering preserves heterozygous CNVs with half RD value of the  
26  
27  
28 246 homozygous CNVs.  
29

30  
31 247 After merging the candidate CNV windows into a CNVR, the RDs of all samples in each  
32  
33 248 CNVR are clustered by the mixture Gaussian model and deducing the integer copy number of  
34  
35  
36 249 each individual. This step is called genotyping as used in SNP detection. The final output is  
37  
38  
39 250 compatible with most SNP based population genetic algorithm.  
40

### 41 42 43 251 **Computational cost**

44  
45  
46 252 The robustness of CNVcaller was validated by the real sequencing data of different genomes. The  
47  
48  
49 253 individual RD processing step of CNVcaller was compared against CNVnator, which also detects  
50  
51  
52 254 CNVs individually. The processing time of CNVcaller was linear related to the genome size and  
53  
54  
55 255 sequencing coverage: 20-40 minutes for a 3 Gb genome with 10 X coverage (**Supplementary**  
56  
57  
58 256 **Table 3**). However, the processing time of CNVnator rose exponentially with the scaffold number,  
59  
60

---

1 257 which became the only index of time consuming when the scaffold number exceeds one thousand  
2  
3 258 **(Figure 2A)**. Consequently, CNVcaller achieved 145 fold speed increasing over CNVnator for  
4  
5  
6 259 goat CNV detection. Noteworthy, the goat reference genome ARS1 which contains nearly 30  
7  
8  
9 260 thousand scaffolds was newly assembled by single-molecule sequencing platform [48].  
10

11 261 The memory requirement of CNVcaller is extremely low and mainly related to the genome size:  
12  
13  
14 262 only about 500 Mb for a mammalian genome, which was less than one twentieth of CNVnator  
15  
16  
17 263 **(Figure 2B)**. Therefore, in multi-sample CNV detection, more than 20 missions of the individual  
18  
19  
20 264 RD processing step can be run in parallel on one node to further increase the population CNVR  
21  
22  
23 265 detection efficiency. The population-level performance of CNVcaller was evaluated and  
24  
25  
26 266 benchmarked by Genome STRiP which also detects CNVRs at population level. After removing  
27  
28 267 the unplaced scaffolds, CNVcaller was still 3.5-7.8 times faster than Genome STRiP **(Figure 2C)**,  
29  
30  
31 268 with 70% ~86% reduction in memory requirement **(Figure 2D)**. The CNV detection of 232 goats  
32  
33  
34 269 with mean 12 X coverage can be completed in 1.4 days by CNVcaller on one node.  
35  
36

## 37 270 **Sensitivity and accuracy**

38  
39  
40  
41 271 A total of 1,058 CNVRs with a total length of 24.5 Mb were detected by CNVcaller from a 30  
42  
43  
44 272 human cohort, 20% longer than CNVnator, and twice of Genome STRiP. CNVcaller covered 43%  
45  
46  
47 273 of the CNVRs detected by CNVnator, 65% of Genome STRiP and 76% of their intersection  
48  
49  
50 274 **(Figure 3A)**, indicating CNVcaller has higher sensitivity for the cross validated CNVRs. We also  
51  
52  
53 275 compared the CNVRs identified by CNVcaller from worldwide 133 sheep of 44 breeds with the  
54  
55  
56 276 other two recently released large scale sheep CNVR datasets. One is derived from a small  
57  
58 277 pedigrees using multiple platforms including aCGH, SNP chip and whole genome sequence [49],  
59  
60

---

1 278 the other is based on three Chinese sheep breeds and 600K SNP array [50]. Although based on  
2  
3 279 different technologies and breeds, CNVcaller still covers 51% of their intersection (**Figure 3B**).  
4  
5  
6 280 The 1000 Genome Project samples with experimental validated CNV database was mainly used  
7  
8  
9 281 to evaluate the sensitivity and accuracy of CNVcaller and other two methods. The sensitivity was  
10  
11 282 estimated as the proportion of high-confident CNV database overlapped by predicted CNVs  
12  
13 283 (**Table 1**). Based on the aCGH database [43], the sensitivity of CNVcaller was 13%-18% higher  
14  
15 284 than the other two methods. Even based on the 1000 Genome Project CNV maps which was  
16  
17 285 constructed by multiple methods including the Genome STRiP and CNVnator [24], CNVcaller  
18  
19 286 still achieved higher sensitivity than other software.  
20  
21  
22 287 False discovery rate (FDR) on human genome was estimated by multiple strict sample specific  
23  
24 288 methods (**Table 1**). (1) IRS test based on the intensity data of the SNP array; (2) the integer copy  
25  
26 289 numbers in aCGH database; (3) the Mendelian inconsistencies from 10 Dutch families. For the  
27  
28 290 three CNV detectors, Genome STRiP achieved the best accuracy (0.8% - 3.9%) in all the human  
29  
30 291 based validations, and a little higher Mendelian inconsistencies (5.2%) in the three sheep trios.  
31  
32 292 CNVcaller had median FDR (2.8% - 5.4%) in human validations, however, it achieved the best  
33  
34 293 accuracy (2.4% Mendelian inconsistencies) in the three sheep trios, indicating its superiority in  
35  
36 294 non-human genomes. To evaluate the genotyping accuracy, we turned to a recently developed  
37  
38 295 molecular biology technique, CNVplex, which counts the copy number of a genomic sequence  
39  
40 296 based on the multiplex ligation-dependent probe amplification (MLPA) method [45]. When we  
41  
42 297 compared the copy numbers predicted by CNVcaller from sequencing data and the CNVplex  
43  
44 298 result, the Pearson's product-moment correlation coefficients were higher than 0.95, and the  
45  
46 299 integer genotype concordance was 98% (**Figure 4**).



---

1 300 **The absolute copy number correction in duplicated region of the reference**  
2  
3 301 **genome**  
4  
5  
6  
7 302 Compared with human, the sheep sample had much lower copy numbers in putative duplicated  
8  
9  
10 303 regions than expected (**Figure 5A**), indicating the sequencing reads were split among the false  
11  
12  
13 304 duplications as previous reported on other animal genome assemblies [31]. This bias led the  
14  
15  
16 305 raw copy number distribution of the putative two-copy segment duplications peak at one (**Figure**  
17  
18 306 **5B**), thus most of these windows were likely to be mistaken for heterozygote deletions. After  
19  
20  
21 307 absolute copy number correction (**Figure 1**), the copy number distribution was more reasonable:  
22  
23  
24 308 the main peak shift to around two (normal biallelic copy number), and the smaller peaks around  
25  
26  
27 309 one and four indicated the detectable heterozygous deletions and duplications respectively (**Figure**  
28  
29 310 **5B**). Because this correction preserves the complicated regions and the multi-hit reads, CNVcaller  
30  
31  
32 311 increased the sensitivity of SD region about six time than the other two methods. Moreover, it is  
33  
34  
35 312 more reasonable for most CNVs in SD regions were genotyped as duplications instead of deletions  
36  
37  
38 313 (**Figure 5C**).

39  
40 314 The detection and genotyping accuracy in SD region were further estimated by the sex  
41  
42  
43 315 information of 133 sheep. We first defined 138 unplaced scaffolds with high sequence similarity to  
44  
45  
46 316 X chromosome as X-linked scaffolds. In theory, all these scaffolds were expected to be detected as  
47  
48  
49 317 high frequency CNVRs because the RDs of unplaced scaffolds were not corrected by sex.  
50  
51  
52 318 CNVcaller detected 101 out of the 138 X-linked scaffolds with a sensitivity of 73%, while  
53  
54  
55 319 CNVnator did not catch any of these regions. The corrected copy number of these scaffold  
56  
57 320 centralized at one and two in rams and doubled in ewes indicating the unique and duplicated X-

---

1 321 linked regions. However, the peaks of raw copy numbers were ambiguous and not at integer  
2  
3 322 (**Figure 5D**). Further examination of the duplicated regions showed the higher divergence was  
4  
5  
6 323 caused by splitting of raw RDs among the mistaken-assembly segments (**Supplementary Figure**  
7  
8  
9 324 **2**).

10  
11 325

## 14 326 **Discussion**

17  
18 327 CNVcaller was designed to detect the CNVRs from large scale resequencing data of all kinds of  
19  
20  
21 328 genomes. It takes full consideration of the complexity of genome, and implies several general  
22  
23  
24 329 applicable signal detection and noise deduction algorithms to increase the computational  
25  
26  
27 330 efficiency. The detection of 232 individuals can be complete on one compute node within two  
28  
29  
30 331 days, and this speed up does not compromise the accuracy. Meanwhile, validated through multiple  
31  
32 332 rigorous assessment, CNVcaller increased the sensitivity by 13%-18% with a low FDR in human  
33  
34  
35 333 and non-human species.

36  
37  
38 334 The statistics-based detection algorithm of CNVcaller assumes a true CNV signal can be  
39  
40  
41 335 interrupted with high ratio of noise, while the segmentation algorithms assume that the RD signal  
42  
43  
44 336 is basically piecewise constant in the genome sequence. The formal model is universally  
45  
46  
47 337 applicable to high quality as well as fragmented genomes. Therefore, the speed of CNVcaller is  
48  
49  
50 338 still fast with thousands of scaffolds, which is common in non-human reference genome  
51  
52  
53 339 assemblies. The robustness of CNVcaller reduces the restrictions of the reference genome, which  
54  
55  
56 340 will promote CNV research of the species with only scaffold level reference genome assemblies.  
57  
58 341 More importantly, this feature enables the comprehensive variation discoveries using multiple  
59  
60  
61  
62  
63  
64  
65

---

1 342 assemblies or pan-genomes. Defined as the entire set of genes possessed by all members of a  
2  
3 343 particular species, pan-genomes reveal numerous functional important genes unplaced on one  
4  
5  
6 344 single reference genome [50-52]. However, its complexity and diversity hinder the application of  
7  
8  
9 345 almost all CNV detectors. Up on our unpublished data, CNVcaller is efficient and friendly to  
10  
11  
12 346 detect the present/absent variations for pan-genomes.

13  
14  
15 347 Another optimization of CNVcaller is the simplified absolute copy number correction.  
16  
17  
18 348 Although absolute copy number can reduce the bias of misassembled duplications in non-human  
19  
20  
21 349 genomes, conventional solution requires mapping the sequencing reads of each individual by  
22  
23  
24 350 time-consuming precise aligner. CNVcaller simplifies the calculation by generating a duplicated  
25  
26  
27 351 window record file through one precise alignment of the reference genome (the duplicated  
28  
29  
30 352 window record file of the latest reference genome of human, livestock and main crops can be  
31  
32 353 downloaded from website), then the standard alignment of all samples can be corrected with high  
33  
34  
35 354 speed. In this way, the misassembled duplications can be mitigated with a great reduction of  
36  
37  
38 355 running time. Validation of sheep genomes showed this correction multiplied the detection  
39  
40  
41 356 efficiency of SD regions, and deduced the more reasonable copy numbers.

42  
43 357 Several limitations still exist in CNVcaller pipeline. First, population level criteria are used to  
44  
45  
46 358 screen out the low frequency and uncorrelated windows. Therefore, CNVcaller is not suitable for  
47  
48  
49 359 detection of rare CNVs, however, the influence is much less on GWAS which also has low power  
50  
51  
52 360 to capture rare functional variants [51]. Besides, the RD algorithm has disadvantage in short CNV  
53  
54  
55 361 detection and breakpoint definition. The visual examination for the specific interval using IGV  
56  
57  
58 362 [52] or combined with other CNV detection methods can improve the results. In the latter  
59  
60  
61 363 situation, CNVcaller can provide high-confidence RD information as solid prior for the read

---

1 364 pair/split read pipeline [53].  
2  
3

4 365 In summary, CNVcaller offers a fast and robust pipeline to detect CNVRs from population-  
5  
6  
7 366 scale resequencing data. The high computational efficiency reduces the hardware requirements  
8  
9  
10 367 and facilitates the CNVR detection of large populations. The general applicable detection and  
11  
12  
13 368 correction algorithms have greatly increased the sensitivity in non-model species and enabled the  
14  
15  
16 369 CNV detection for a wide range of species. The rapid and reliable population-level CNV detection  
17  
18  
19 370 will promote the discovery of the missing heritability of complex traits and accurately  
20  
21  
22 371 determination of the causative mutations for more species.  
23  
24  
25 372  
26  
27

## 28 373 **Availability and requirements**

31 374 Project name: CNVcaller

34 375 Project home page: <http://animal.nwsuaf.edu.cn/software>

37 376 <https://github.com/JiangYuLab/CNVcaller>

40 377 Operating system(s): platform independent

43 378 Programming language: Perl, C++

46 379 Other requirements: Samtools 1.3 (using htlib 1.3)

49 380 License: GNU General Public License, version 3.0 (GPL-3.0)

52 381

## 55 382 **Conflict of interest**

58 383 The authors declare that they have no competing interests

---

1 384

2  
3  
4 385 **Author contributions**

5  
6  
7 386 WXH and JY designed this software; ZZQ and CT wrote the code; WXH and ZZQ improved the  
8  
9  
10 387 pipeline structures; ZZQ and CYD tested the software prototype; LC and FWW contributed to the  
11  
12  
13 388 data organization; WXH and JY drafted the manuscript. All authors read and approved the final  
14  
15  
16 389 manuscript.

17  
18  
19 390

20  
21 391 **Acknowledgements**

22  
23  
24 392 This work is supported by grants from National Natural Science Foundation of China (31572381),  
25  
26  
27 393 and the National Thousand Youth Talents Plan. We thank the International Sheep Genomics  
28  
29  
30 394 Consortium (ISGC) for access to the unpublished sheep sequencing data provided under the  
31  
32  
33 395 Toronto guidelines for data users. We thank for the Genomes of Netherlands (GoNL) project for  
34  
35  
36 396 the human family data.

37  
38 397

39  
40  
41 398 **References**

- 42  
43 399 1. Chiang C, Scott AJ, Davis JR, Tsang EK, Li X, Kim Y, et al. The impact of  
44  
45  
46 400 structural variation on human gene expression. *Nature genetics*. 2017.
- 47  
48  
49 401 2. Stefansson H, Meyer-Lindenberg A, Steinberg S, Magnusdottir B, Morgen K,  
50  
51  
52 402 Arnarsdottir S, et al. CNVs conferring risk of autism or schizophrenia affect cognition  
53  
54  
55 403 in controls. *Nature*. 2014;505 7483:361-6.
- 56  
57 404 3. Khurana E, Fu Y, Colonna V, Mu XJ, Kang HM, Lappalainen T, et al.

---

1 405 Integrative annotation of variants from 1092 humans: application to cancer genomics.  
2  
3 406 Science. 2013;342 6154:1235587.  
4  
5  
6 407 4. Sudmant PH, Mallick S, Nelson BJ, Hormozdiari F, Krumm N, Huddleston J, et  
7  
8  
9 408 al. Global diversity, population stratification, and selection of human copy-number  
10  
11 409 variation. Science. 2015;349 6253:aab3761.  
12  
13  
14 410 5. Norris BJ and Whan VA. A gene duplication affecting expression of the ovine  
15  
16  
17 411 ASIP gene is responsible for white and black sheep. Genome research. 2008;18  
18  
19 412 8:1282-93.  
20  
21  
22 413 6. Wright D, Boije H, Meadows JR, Bed'Hom B, Gourichon D, Vieaud A, et al.  
23  
24  
25 414 Copy number variation in intron 1 of SOX5 causes the Pea-comb phenotype in  
26  
27 415 chickens. PLoS Genet. 2009;5 6:e1000512.  
28  
29  
30 416 7. Seo B-Y, Park E-W, Ahn S-J, Lee S-H, Kim J-H, Im H-T, et al. An accurate  
31  
32  
33 417 method for quantifying and analyzing copy number variation in porcine KIT by an  
34  
35  
36 418 oligonucleotide ligation assay. BMC genetics. 2007;8 1:81.  
37  
38  
39 419 8. Zhou Z, Jiang Y, Wang Z, Gou Z, Lyu J, Li W, et al. Resequencing 302 wild  
40  
41  
42 420 and cultivated accessions identifies genes related to domestication and improvement  
43  
44  
45 421 in soybean. Nature biotechnology. 2015;33 4:408-14.  
46  
47  
48 422 9. Craddock N, Hurles ME, Cardin N, Pearson RD, Plagnol V, Robson S, et al.  
49  
50  
51 423 Genome-wide association study of CNVs in 16,000 cases of eight common diseases  
52  
53 424 and 3,000 shared controls. Nature. 2010;464 7289:713-20. doi:10.1038/nature08979.  
54  
55  
56 425 10. Xu L, Cole JB, Bickhart DM, Hou Y, Song J, VanRaden PM, et al. Genome  
57  
58 426 wide CNV analysis reveals additional variants associated with milk production traits in  
59  
60  
61  
62  
63  
64  
65

---

1 427 Holsteins. BMC genomics. 2014;15:683. doi:10.1186/1471-2164-15-683.  
2  
3  
4 428 11. Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, et al.  
5  
6 429 BreakDancer: an algorithm for high-resolution mapping of genomic structural  
7  
8  
9 430 variation. Nat Meth. 2009;6 9:677-81.  
10  
11 431 12. Hormozdiari F, Hajirasouliha I, Dao P, Hach F, Yorukoglu D, Alkan C, et al.  
12  
13  
14 432 Next-generation VariationHunter: combinatorial algorithms for transposon insertion  
15  
16  
17 433 discovery. Bioinformatics. 2010;26 12:i350-i7. doi:10.1093/bioinformatics/btq216.  
18  
19  
20 434 13. Ye K, Schulz MH, Long Q, Apweiler R and Ning Z. Pindel: a pattern growth  
21  
22  
23 435 approach to detect break points of large deletions and medium sized insertions from  
24  
25  
26 436 paired-end short reads. Bioinformatics. 2009;25 21:2865-71.  
27  
28  
29 437 doi:10.1093/bioinformatics/btp394.  
30  
31 438 14. Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V and Korbel JO. DELLY:  
32  
33  
34 439 structural variant discovery by integrated paired-end and split-read analysis.  
35  
36  
37 440 Bioinformatics. 2012;28 18:i333-i9.  
38  
39  
40 441 15. Chen K, Chen L, Fan X, Wallis J, Ding L and Weinstock G. TIGRA: A targeted  
41  
42  
43 442 iterative graph routing assembler for breakpoint assembly. Genome research.  
44  
45  
46 443 2014;24 2:310-7. doi:10.1101/gr.162883.113.  
47  
48  
49 444 16. Zerbino DR and Birney E. Velvet: Algorithms for de novo short read assembly  
50  
51  
52 445 using de Bruijn graphs. Genome research. 2008;18 5:821-9.  
53  
54  
55 446 doi:10.1101/gr.074492.107.  
56  
57  
58 447 17. Zhuang J and Weng Z. Local sequence assembly reveals a high-resolution  
59  
60  
61 448 profile of somatic structural variations in 97 cancer genomes. Nucleic Acids Res.

---

1 449 2015;43 17:8146-56. doi:10.1093/nar/gkv831.

2  
3 450 18. Yoon S, Xuan Z, Makarov V, Ye K and Sebat J. Sensitive and accurate  
4  
5  
6 451 detection of copy number variants using read depth of coverage. *Genome research*.  
7  
8  
9 452 2009;19 9:1586-92.

10  
11 453 19. Abyzov A, Urban AE, Snyder M and Gerstein M. CNVnator: an approach to  
12  
13  
14 454 discover, genotype, and characterize typical and atypical CNVs from family and  
15  
16  
17 455 population genome sequencing. *Genome research*. 2011;21 6:974-84.

18  
19  
20 456 20. Szatkiewicz JP, Wang W, Sullivan PF, Wang W and Sun W. Improving  
21  
22  
23 457 detection of copy-number variation by simultaneous bias correction and read-depth  
24  
25  
26 458 segmentation. *Nucleic acids research*. 2013;41 3:1519-32.

27  
28 459 21. Wang M, Tu L, Lin M, Lin Z, Wang P, Yang Q, et al. Asymmetric subgenome  
29  
30  
31 460 selection and cis-regulatory divergence during cotton domestication. *Nature genetics*.  
32  
33  
34 461 2017.

35  
36 462 22. Klambauer G, Schwarzbauer K, Mayr A, Clevert D-A, Mitterecker A,  
37  
38  
39 463 Bodenhofer U, et al. cn. MOPS: mixture of Poissons for discovering copy number  
40  
41  
42 464 variations in next-generation sequencing data with a low false discovery rate. *Nucleic*  
43  
44  
45 465 *acids research*. 2012:gks003.

46  
47 466 23. Handsaker RE, Van Doren V, Berman JR, Genovese G, Kashin S, Boettger  
48  
49  
50 467 LM, et al. Large multiallelic copy number variations in humans. *Nature genetics*.  
51  
52  
53 468 2015;47 3:296-303.

54  
55  
56 469 24. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J,  
57  
58  
59 470 et al. An integrated map of structural variation in 2,504 human genomes. *Nature*.



---

1 471 2015;526 7571:75-81.  
2  
3 472 25. Claros MG, Bautista R, Guerrero-Fernández D, Benzerki H, Seoane P and  
4  
5  
6 473 Fernández-Pozo N. Why assembling plant genome sequences is so challenging.  
7  
8  
9 474 Biology. 2012;1 2:439-59.  
10  
11 475 26. Warr A, Hume D, Archibald AL, Deeb N and Watson M. Identification of low-  
12  
13 476 confidence regions in the pig reference genome (*Sscrofa10. 2*). *Frontiers in genetics*.  
14  
15 477 2015;6:338.  
16  
17 478 27. Xie C and Tammi MT. CNV-seq, a new method to detect copy number variation  
18  
19 479 using high-throughput sequencing. *BMC bioinformatics*. 2009;10 1:80.  
20  
21 480 28. Zhang X, Wang K, Wang L, Yang Y, Ni Z, Xie X, et al. Genome-wide patterns  
22  
23 481 of copy number variation in the Chinese yak genome. *BMC genomics*. 2016;17 1:1.  
24  
25 482 29. Chain FJ, Feulner PG, Panchal M, Eizaguirre C, Samonte IE, Kalbe M, et al.  
26  
27 483 Extensive copy-number variation of young genes across stickleback populations.  
28  
29 484 *PLoS Genet*. 2014;10 12:e1004830. doi:10.1371/journal.pgen.1004830.  
30  
31 485 30. Yi G, Qu L, Liu J, Yan Y, Xu G and Yang N. Genome-wide patterns of copy  
32  
33 486 number variation in the diversified chicken genomes using next-generation  
34  
35 487 sequencing. *BMC genomics*. 2014;15:962. doi:10.1186/1471-2164-15-962.  
36  
37 488 31. Kelley DR and Salzberg SL. Detection and correction of false segmental  
38  
39 489 duplications caused by genome mis-assembly. *Genome biology*. 2010;11 3:1.  
40  
41 490 32. Zimin AV, Kelley DR, Roberts M, Marçais G, Salzberg SL and Yorke JA. Mis-  
42  
43 491 Assembled “Segmental Duplications” in Two Versions of the *Bos*  
44  
45 492 *taurus* Genome. *PLoS One*. 2012;7 8:e42680.  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

---

1 493 doi:10.1371/journal.pone.0042680.  
2  
3 494 33. Dennis MY, Harshman L, Nelson BJ, Penn O, Cantsilieris S, Huddleston J, et  
4  
5  
6 495 al. The evolution and population diversity of human-specific segmental duplications.  
7  
8  
9 496 Nature Ecology & Evolution. 2017;1:0069.  
10  
11 497 34. Hach F, Hormozdiari F, Alkan C, Hormozdiari F, Birol I, Eichler EE, et al.  
12  
13  
14 498 mrsFAST: a cache-oblivious algorithm for short-read mapping. Nature methods.  
15  
16  
17 499 2010;7 8:576-7.  
18  
19  
20 500 35. Liu GE, Hou Y, Zhu B, Cardone MF, Jiang L, Cellamare A, et al. Analysis of  
21  
22  
23 501 copy number variations among diverse cattle breeds. Genome research. 2010;20  
24  
25 502 5:693-703.  
26  
27  
28 503 36. Zarrei M, MacDonald JR, Merico D and Scherer SW. A copy number variation  
29  
30  
31 504 map of the human genome. Nature Reviews Genetics. 2015.  
32  
33  
34 505 37. Consortium GotN. Whole-genome sequence variation, population structure and  
35  
36  
37 506 demographic history of the Dutch population. Nature genetics. 2014;46 8:818-25.  
38  
39  
40 507 38. Diez CM, Meca E, Tenaillon MI and Gaut BS. Three groups of transposable  
41  
42  
43 508 elements with contrasting copy number dynamics and host responses in the maize  
44  
45 509 (*Zea mays ssp. mays*) genome. PLoS Genet. 2014;10 4:e1004298.  
46  
47  
48 510 39. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al.  
49  
50  
51 511 The Genome Analysis Toolkit: a MapReduce framework for analyzing next-  
52  
53 512 generation DNA sequencing data. Genome Res. 2010;20 9:1297-303.  
54  
55  
56 513 doi:10.1101/gr.107524.110.  
57  
58  
59 514 40. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The

---

1 515 Sequence Alignment/Map format and SAMtools. *Bioinformatics*. 2009;25 16:2078-9.  
2  
3 516 doi:10.1093/bioinformatics/btp352.  
4  
5  
6 517 41. Kent WJ. BLAT—the BLAST-like alignment tool. *Genome research*. 2002;12  
7  
8  
9 518 4:656-64.  
10  
11  
12 519 42. Scrucca L, Fop M, Murphy TB and Raftery AE. mclust 5: Clustering,  
13  
14 520 classification and density estimation using gaussian finite mixture models. *The R*  
15  
16 521 *Journal*. 2016;8 1:289.  
17  
18  
19 522 43. Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, et al. Origins  
20  
21 523 and functional impact of copy number variation in the human genome. *Nature*.  
22  
23 524 2010;464 7289:704-12.  
24  
25  
26 525 44. Quinlan AR and Hall IM. BEDTools: a flexible suite of utilities for comparing  
27  
28 526 genomic features. *Bioinformatics*. 2010;26 6:841-2.  
29  
30  
31 527 45. Zhang X, Xu Y, Liu D, Geng J, Chen S, Jiang Z, et al. A modified multiplex  
32  
33 528 ligation-dependent probe amplification method for the detection of 22q11. 2 copy  
34  
35 529 number variations in patients with congenital heart disease. *BMC genomics*. 2015;16  
36  
37 530 1:364.  
38  
39  
40 531 46. Chaisson MJ and Tesler G. Mapping single molecule sequencing reads using  
41  
42 532 basic local alignment with successive refinement (BLASR): application and theory.  
43  
44 533 *BMC bioinformatics*. 2012;13 1:238.  
45  
46  
47 534 47. Bickhart DM, Hou Y, Schroeder SG, Alkan C, Cardone MF, Matukumalli LK, et  
48  
49 535 al. Copy number variation of individual cattle genomes using next-generation  
50  
51 536 sequencing. *Genome research*. 2012;22 4:778-90.  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

---

1 537 48. Bickhart DM, Rosen BD, Koren S, Sayre BL, Hastie AR, Chan S, et al. Single-  
2  
3 538 molecule sequencing and chromatin conformation capture enable de novo reference  
4  
5  
6 539 assembly of the domestic goat genome. *Nature Genetics*. 2017;49 4:643-50.  
7  
8  
9 540 49. Jenkins GM, Goddard ME, Black MA, Brauning R, Auvray B, Dodds KG, et al.  
10  
11 541 Copy number variants in the sheep genome detected using multiple approaches.  
12  
13  
14 542 *BMC genomics*. 2016;17 1:1.  
15  
16  
17 543 50. Zhu C, Fan H, Yuan Z, Hu S, Ma X, Xuan J, et al. Genome-wide detection of  
18  
19  
20 544 CNVs in Chinese indigenous sheep with different types of tails using ovine high-  
21  
22  
23 545 density 600K SNP arrays. *Scientific reports*. 2016;6.  
24  
25  
26 546 51. Robinson MR, Wray NR and Visscher PM. Explaining additional genetic  
27  
28 547 variation in complex traits. *Trends in Genetics*. 2014;30 4:124-32.  
29  
30  
31 548 52. Thorvaldsdóttir H, Robinson JT and Mesirov JP. Integrative Genomics Viewer  
32  
33  
34 549 (IGV): high-performance genomics data visualization and exploration. *Briefings in*  
35  
36  
37 550 *bioinformatics*. 2013;14 2:178-92.  
38  
39  
40 551 53. Layer RM, Chiang C, Quinlan AR and Hall IM. LUMPY: a probabilistic  
41  
42 552 framework for structural variant discovery. *Genome biology*. 2014;15 6:R84.  
43  
44  
45 553

46  
47 **Figure Legends**  
48

49  
50 **Figure 1** CNVcaller algorithm flowchart (left) and the key algorithms of each step (right). (1)  
51  
52  
53 556 Individually RD processing. In the absolute copy number correction, the RDs of high similar  
54  
55  
56 557 windows are added together to deduce the absolute copy number. (2) Multi-criteria CNVR  
57  
58  
59 558 selection. Curves show copy numbers in a specific region for multiple samples. Blue transverse  
60  
61  
62  
63  
64  
65

---

1 559 boxes mark the windows with significantly distinguish copy number from the average (individual  
2  
3 560 criterion). Green vertical boxes indicate regions with the CNV allele frequency >5% in a specific  
4  
5  
6 561 region, and red frame indicates the RDs between two adjacent windows are significantly  
7  
8  
9 562 correlated (population criteria). Only the region with continuous high CNV allele frequency and  
10  
11 563 high correlation (the forth bar from the left) are selected as the CNVRs. (3) Genotyping: The copy  
12  
13 564 numbers in each CNVR are clustered by mixture Gaussian model to distinguish the normal,  
14  
15  
16  
17 565 heterozygous and homozygous samples.  
18  
19

20 566

21  
22 567 **Figure 2** Computational performance of CNVcaller, CNVnator and Genome STRiP. All the  
23  
24  
25 568 programs were executed on one node with two 2.40-GHz Intel Xeon E5-2620 v3 processors. (A,  
26  
27  
28 569 B) Log plots of processing time (A) and max memory (B) for one individual. The numbers of  
29  
30  
31 570 unplaced scaffolds of the reference genome are indicated in brackets. The processing time was  
32  
33  
34 571 normalized by genome size and sequencing coverage to simulate a 3 Gb genome with 5 X or 10 X  
35  
36 572 sequencing coverage. (C, D) Log plots of total running time (C) and max memory (D) of  
37  
38  
39 573 population CNVR detection. The test cohorts are: 8 sheep, 30 humans and 232 goats with 19 X, 16  
40  
41  
42 574 X and 12 X average sequencing coverage respectively. In Genome STRiP running, the unplaced  
43  
44  
45 575 scaffolds were removed from the reference genome.  
46

47 576

48  
49  
50 577 **Figure 3** Overlap of the CNVRs detected by CNVcaller and other approaches/platforms (A)  
51  
52  
53 578 Intersection of the CNVRs detected by CNVcaller, CNVnator and Genome STRiP based on the  
54  
55  
56 579 same 30 human data. (B) Intersection of the CNVRs detected by CNVcaller and two other large  
57  
58  
59 580 scale sheep CNVR studies.  
60

---

1 581

2  
3 582 **Figure 4** Evaluation of CNV genotypes by CNVplex. Two duplicated (A, B) and two deleted (C,  
4  
5  
6 583 D) CNVRs with high variation frequency were typed in CNVplex in 73 sheep samples. The copy  
7  
8  
9 584 number genotypes predicted by CNVcaller from sequencing data were plotted against the  
10  
11 585 measurements from CNVplex of the same animal.

12  
13  
14 586

15  
16  
17 587 **Figure 5** Absolute copy number correction in sheep genome. (A) The copy numbers of all  
18  
19 588 windows with no more than six repeats were plotted against the repeat numbers in reference  
20  
21 589 genome. (B) Distribution of copy numbers of two-copy loci in sheep genome before and after  
22  
23 590 absolute copy number correction. (C) The detected CNVRs resided in SD regions. The sheep SD  
24  
25 591 regions include the regions longer than 2 Kb with >97% identity. The CNVRs resided in SD  
26  
27 592 regions were defined if more than 50% of this CNVR was overlapped with the SD regions. (D)  
28  
29 593 The raw and corrected copy numbers of all X-linked scaffolds of 133 sheep.

30  
31  
32  
33  
34 594  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

**Table 1.** Sensitivity and FDR of CNVcaller, CNVnator and Genome STRiP.

Methods	Estimated sensitivity		Estimated FDR			
	aCGH	1000 GP	IRS	CGH Genotype	Mendelian error Human*	Mendelian error Sheep*
CNVcaller	45.4%	56.1%	4.1%	5.4%	2.8%	2.4%
CNVnator	32.6%	51.7%	11.4%	5.4%	5.5%	3.7%
Genome STRiP	27.2%	50.4%	3.9%	2.2%	0.8%	5.2%

\*The Mendelian errors in human and sheep were calculated based on 10 Dutch families and three sheep trios respectively. Other evaluations were based on 30 human BAM files downloaded from 1000 genome project.

### (1) Individual RD processing

Count the RD of each sliding window across genome



Absolute copy number correction, GC correction and normalization

Pile up corrected RD of all samples

### (2) Multi-criteria CNVR selection

Individual RD higher or lower than global average



CNV allele frequency > 5%



RDs of the adjacent windows significant correlated

Define CNVR boundary

### (3) Genotyping

Report integer copy numbers by Gaussian Mixture Model

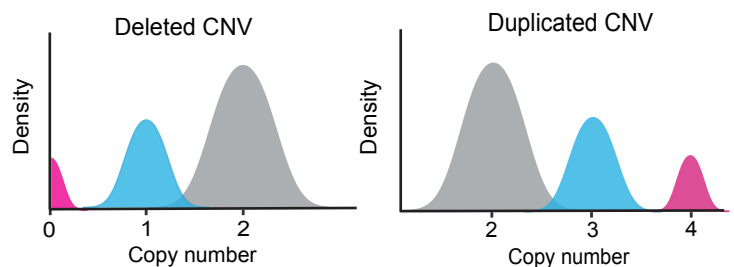
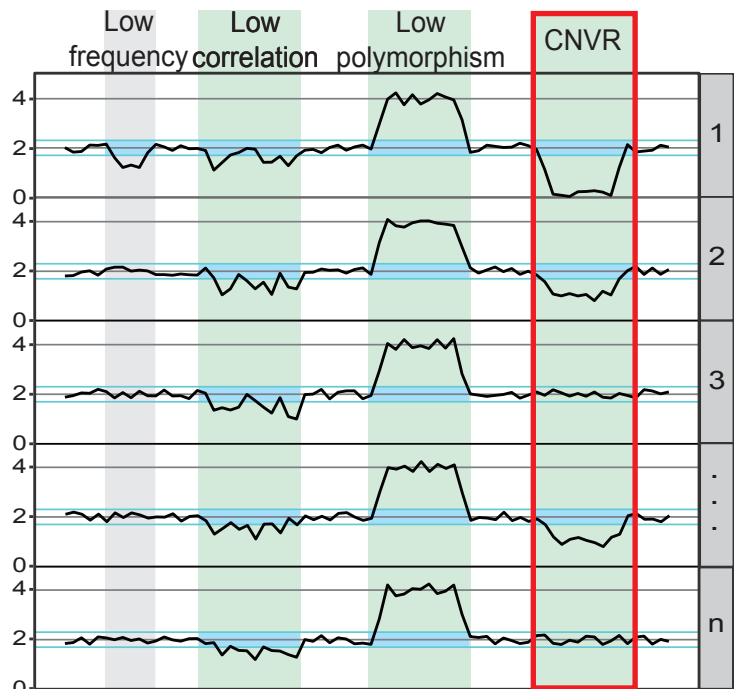
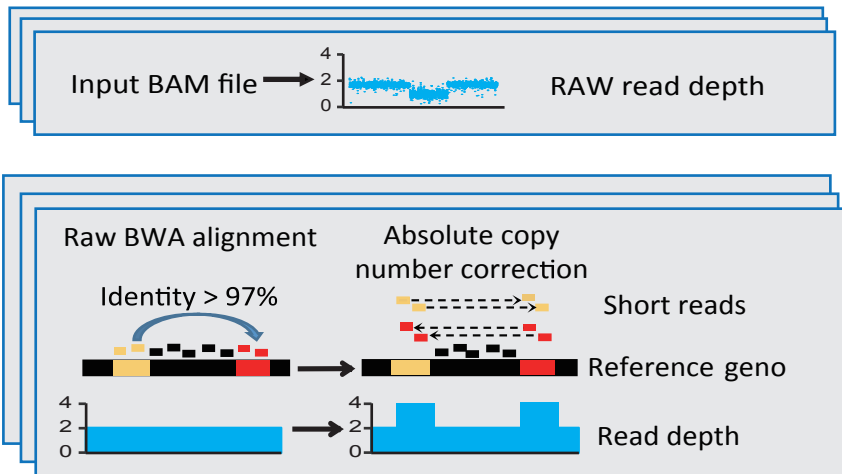
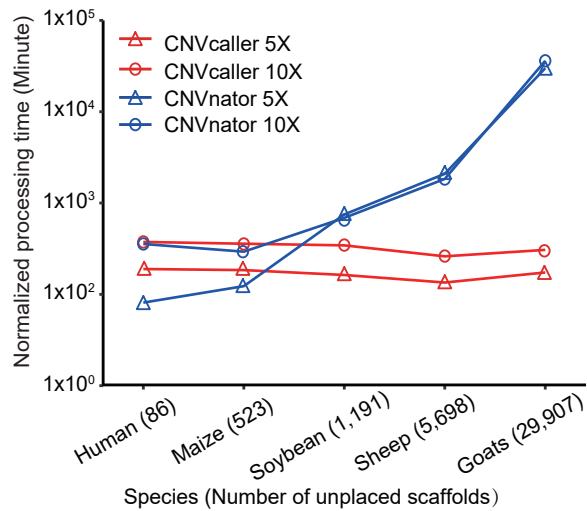
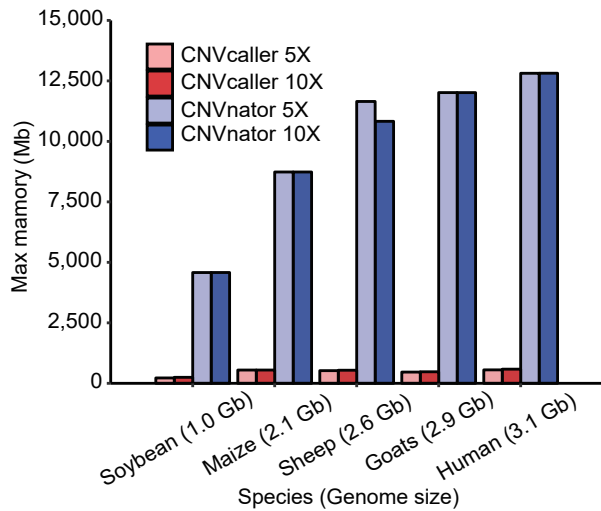
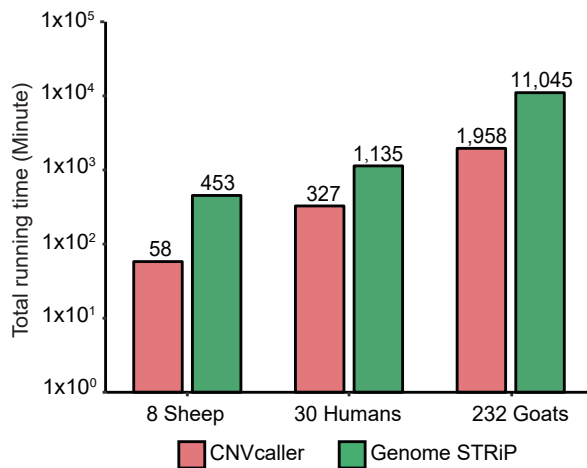




Figure 2

B [Click here to download Figure Figure2.pdf](#)

C



D

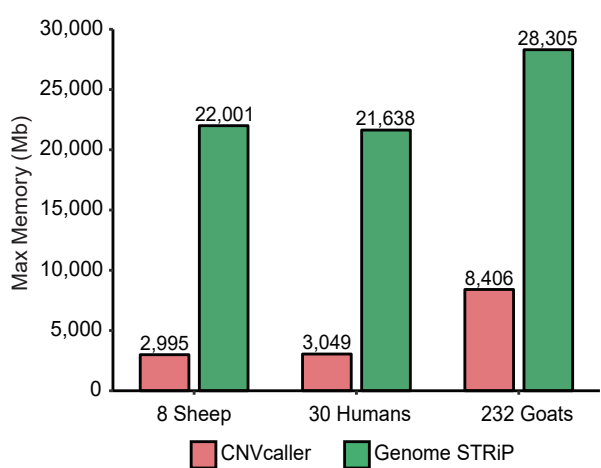
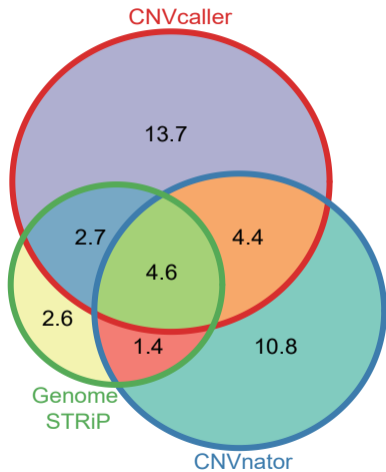


Figure3

A



[Click here to download Figure Figure3.pdf](#)

B

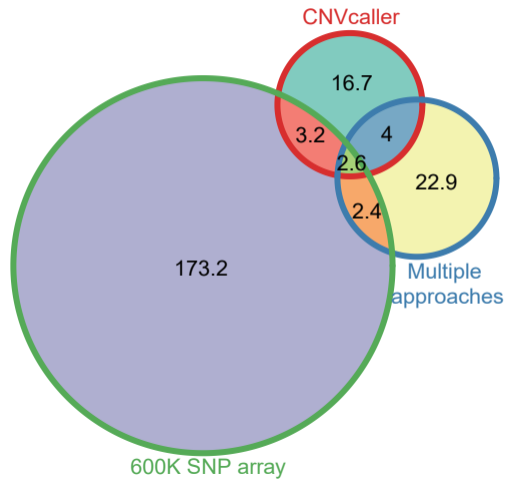
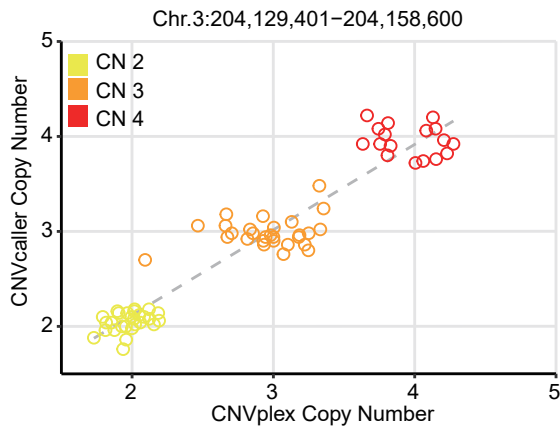


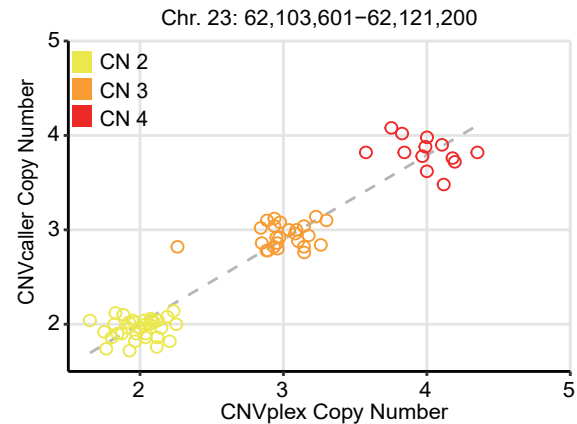
Figure4

[Click here to download Figure Figure4.pdf](#)

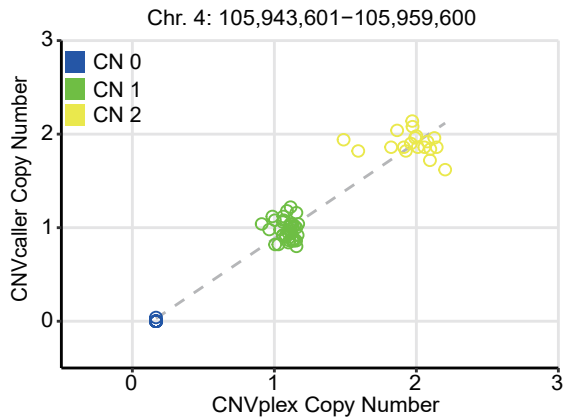
A



B



C



D

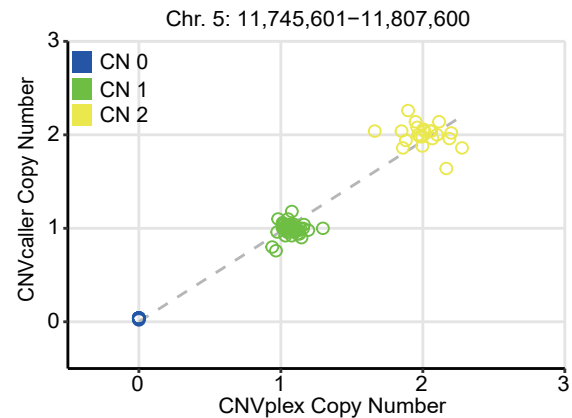
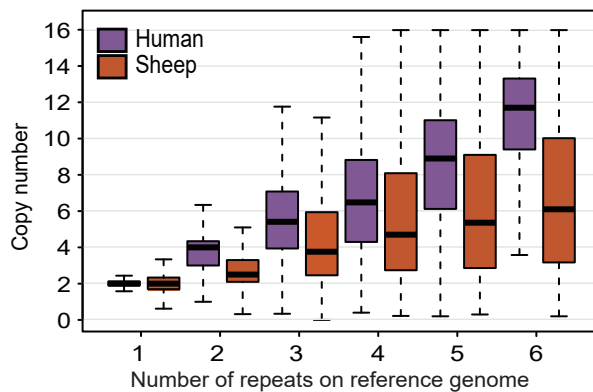
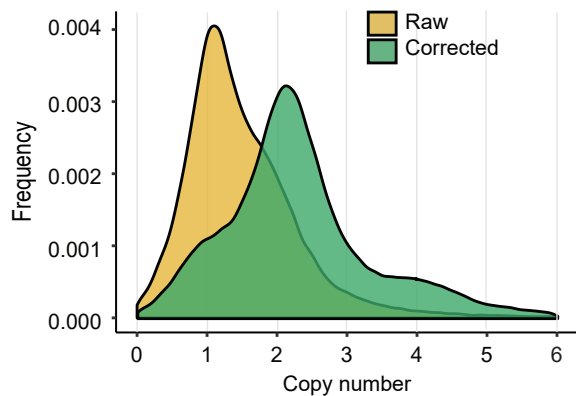


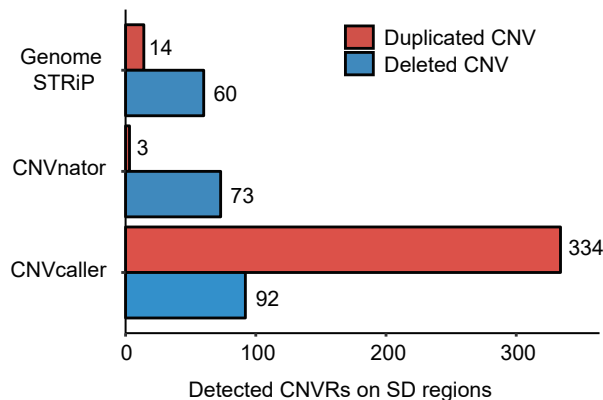
Figure 5  
A

B

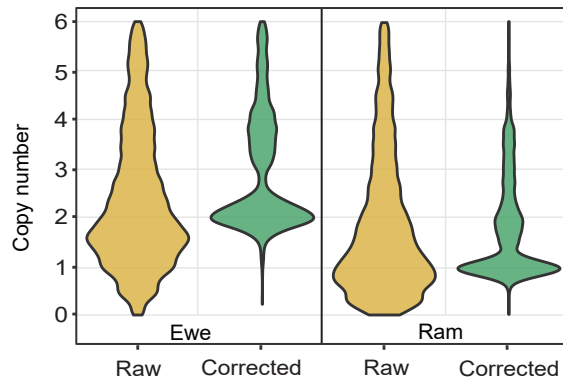
[Click here to download Figure Figure5.pdf](#)



C



D





Click here to access/download  
**Supplementary Material**  
Supplementary Materials.docx

