

1  
2  
3  
4 **A recurrence based approach for validating structural variation using long-read sequencing**  
5 **technology.**  
6  
7

8  
9 Xuefang Zhao<sup>1</sup>, Alexandra M. Weber<sup>1</sup>, and Ryan E. Mills<sup>1,2,#</sup>  
10

11  
12 <sup>1</sup>Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI,  
13 48109, USA, <sup>2</sup>Department of Neurology, University of Michigan, Ann Arbor, MI, 48109, USA,  
14

15  
16 <sup>3</sup>Veterans Affairs Medical Center, Ann Arbor, MI, 48105, USA, <sup>4</sup>Department of Human Genetics,  
17 University of Michigan, Ann Arbor, MI, 48109, USA  
18

19  
20  
21 #Corresponding author: remills@umich.edu  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
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

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
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

# **ABSTRACT**

## **Background**

Although numerous algorithms have been developed to identify structural variation (SVs) in genomic sequences, there is a dearth of approaches that can be used to evaluate their results. This is significant, as the accurate identification of structural variation is still an outstanding but important problem in genomics. The emergence of new sequencing technologies that generate longer sequence reads can, in theory, provide direct evidence for all types of SVs regardless of the length of region through which it spans. However, current efforts to use these data in this manner require the use of large computational resources to assemble these sequences as well as visual inspection of each region.

## **Results**

Here we present VaPoR, a highly efficient algorithm that autonomously validates large SV sets using long read sequencing data. We assessed of the performance of VaPoR on SVs in both simulated and real genomes and report a high-fidelity rate for overall accuracy across different levels of sequence depths. We show that VaPoR can interrogate a much larger range of SVs while still matching existing methods in terms of false positive validations and providing additional features considering breakpoint precision and predicted genotype. We further show that VaPoR can run quickly and efficiency without requiring a large processing or assembly pipeline.

## **Conclusions**

VaPoR serves as a high efficient long read based validation approach for genomic SVs that requires relatively low read depth and computing resources and thus will provide utility with targeted or low-pass sequencing coverage for accurate SV assessment.

**Keywords:** structural variation, copy number variation, sequence analysis

1  
2  
3  
4 **BACKGROUND**  
5  
6

7 Structural variants (SVs) are one of the major forms of genetic variation in humans and have been  
8 revealed to play important roles in numerous diseases including cancers and neurological disorders [1,  
9 2]. Various approaches have been developed and applied to paired-end sequencing to detect SVs in  
10 whole genomes [3-6], however individual algorithms often exhibit complementary strengths that  
11 sometimes lead to disagreements as to the precise structure of the underlying variant. The emergence of  
12 long read sequencing technology, such as Single Molecule Real-Time (SMRT) sequencing from Pacific  
13 Biosciences (PacBio) [7, 8], can deliver reads ranging from several to hundreds of kilobases and provide  
14 direct evidence for the presence of an SV. Current strategies make use of de novo assembly to create  
15 long contigs with minimized error rate [9-11], and then predict SVs, usually with single base resolution,  
16 through direct comparison of the assembly against the reference. Though such approaches are powerful,  
17 they require both a very high sequencing depth and significant computing power and are currently  
18 impracticable for many ongoing research studies.  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28

29  
30 The additional information obtained from using long reads can still be leveraged to improve variant  
31 calling, however. Indeed, such approaches have already been implemented to combine high depth  
32 Illumina sequencing with lower depth PacBio reads to improve error correction and variant calling in the  
33 context of *de novo* genome assembly [12]. With structural variation, the current state of the art is to use  
34 long reads to manually assess potential SVs using subsequent recurrence (dot) plots [13], where the  
35 sequences are compared against the reference through a fixed size sliding window (k-mer) and the  
36 matches are plotted for visual inspection. The k-mer method is of higher robustness compared against  
37 the direct sequences comparison [14], which is why these types of dot plots have been used for decades  
38 to examine the specific features of sequence alignments [15]. However, they require manual curation  
39 and, coupled with the computational costs of sequence assembly, are time-consuming and inefficient at  
40 scale for the high throughput validation of large sets of SVs.  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50

51  
52 Here, we present a high-speed long read based assessment tool, VaPoR, that scores each SV prediction  
53 by autonomously analyzing the recurrence of windows within a local read against the reference genome  
54 in both their original and rearranged format per the prediction. A positive score of each read on the  
55 altered reference, normalized against the score of the read on the original reference, supports the  
56 predicted structure. A baseline model is constructed as well by interrogating the reference sequence  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 against itself at the query location. We show that our approach can quickly and accurately distinguish  
5 true from false positive predictions of both simple and complex SVs as well as their underlying  
6 genotypes and is also able to assess the breakpoint accuracy of individual algorithms.  
7  
8  
9

## 10 11 12 13 14 **DATA DESCRIPTION**

### 15 16 17 *Simulated Data:*

18  
19  
20 Non-overlapping simple deletions, inversions, insertions and duplications as well as complex structural  
21 variants as previously categorized [5] were independently incorporated into GRCh38 in both  
22 heterozygous and homozygous states, excluding regions of the genome that are known to be difficult to  
23 assess as described from the ENCODE project [16]. Detailed descriptions of each simulated SV types  
24 simulated are summarized in Supplementary Tables 1- 2. We applied PBSIM [17] to simulate the  
25 modified reference sequences to different read depth ranging from 2X to 70X with a parameters  
26 difference-ratio of 5:75:20, length-mean 12000, accuracy-mean 0.85 and *model\_qc model\_qc\_clr*.  
27 Simulated data can be obtained from <https://umich.box.com/v/vapor>.  
28  
29  
30  
31  
32  
33  
34

### 35 36 *Real Data*

37  
38  
39 We applied VaPoR to a set of diverse samples (HG00513 from CHS, HG00731 and HG00732 from  
40 PUR, NA19238 and NA19239 from YRI) that were initially sequenced by the 1000 Genomes Project  
41 and for which a high-quality set of SVs were reported in the final phase of the project [18]. These  
42 samples were recently re-sequenced using PacBio and therefore provides a platform for assessing  
43 VaPoR on known data. The 1000 Genomes Project (1KGP) Phase 3 data were obtained from [ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp/phase3/integrated\\_sv\\_map/](ftp://ftp-trace.ncbi.nih.gov/1000genomes/ftp/phase3/integrated_sv_map/) and lifted over to GRCh38. PacBio  
44 sequence data were obtained from  
45  
46  
47  
48  
49  
50  
51  
52 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data\\_collections/hgsv\\_sv\\_discovery/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/data_collections/hgsv_sv_discovery/).  
53

54  
55 We have also compared VaPoR against the long read validation approach developed by Layer et al. [3],  
56 which requires both PacBio and Moleculo long sequences for full evaluation of SVs. These comparisons  
57 made use of NA12878, one of few samples that have been sequenced with various technologies  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 including Illumina NGS, PacBio and Moleculo with a truth SV set included in the 1KGP Phase 3 report.  
5  
6 The software for the long-read validation approach was obtained from: [https://github.com/hall-lab/long-](https://github.com/hall-lab/long-read-validation)  
7 [read-validation](https://github.com/hall-lab/long-read-validation). The PacBio and the Moleculo sequences of this individual were obtained from :  
8  
9 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20131209\\_na12878\\_pacbio/si/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20131209_na12878_pacbio/si/) and  
10  
11 [http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated\\_sv\\_map/supporting/NA12878/moleculo/ali-](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated_sv_map/supporting/NA12878/moleculo/alignment/)  
12 [gnment/](http://ftp.1000genomes.ebi.ac.uk/vol1/ftp/phase3/integrated_sv_map/supporting/NA12878/moleculo/alignment/) respectively.  
13  
14  
15  
16  
17  
18

## 19 **RESULTS**

20  
21  
22 We assessed the performance of VaPoR on both simulated sequences and real genomes from the 1000  
23 Genomes Project to assess the following characteristics: sensitivity and false discovery rate on  
24 validating structural variants in simple and complex structures; sensitivity of VaPoR on validating  
25 different levels of predicted breakpoint efficacy; stratification of VaPoR scores by genotype; and time  
26 and computational cost of VaPoR.  
27  
28  
29  
30  
31

### 32 **VaPoR on Simulated Data**

33  
34  
35 We applied VaPoR to simulated simple deletions, inversions, insertions and duplications as well as  
36 complex structural variants and first assessed the proportion of SVs that VaPoR is capable of  
37 interrogating (i.e. passed VaPoR QC). We found that VaPoR can successfully evaluate >80% of  
38 insertions, >85% deletion-duplications and >90% SVs in all other categories when the read depth is 10X  
39 or higher. We then assessed the sensitivity and false discovery rate (FDR) at different VaPoR score  
40 cutoffs and found that a sensitivity >90% is achieved for most SV types across a wide range of read  
41 depths while maintaining a false discovery rate <10% at a VaPoR score cutoff of 0.15 (Supplemental  
42 Figures 1-2). We further observed that there were no significant changes of sensitivity or false discovery  
43 rate once the read depth was at or above 20X and is consistent across different SV types (Figure 2,  
44 Supplemental Figure 3-4, Supplemental Table 3).  
45  
46  
47  
48  
49  
50  
51  
52  
53

### 54 **VaPoR on 1000 Genomes Project Samples**

55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 We next examined SVs reported on chr1 of 5 individuals from the 1000 Genomes Project [19] to assess  
5 the sensitivity of VaPoR on real genomes (Table 1). We first observed that >95% of deletions and  
6 insertions could be successfully evaluated by VaPoR. For inversions, there were a limited number of  
7 events reported but at maximum only 1 event failed the VaPoR quality control per individual. A  
8 sensitivity of >90% was achieved for deletions (Figure 3a) and >80% for insertions (Figure 3b) at the  
9 recommended cutoff of 0.15. To examine the false validation rate of VaPoR, we modified reported  
10 events on chr2 to appear at the same coordinates on chr1 and assessed them as though they were real  
11 events using the same sequence data set. VaPoR validated very few deletions or inversion and <10% of  
12 insertions. We further assessed the performance of VaPoR to validate SVs with varying degrees of  
13 breakpoint accuracy. Real coordinates were artificially shifted each direction by -1000 to 1000 base  
14 pairs and re-assessed with VaPoR for both simulated and real samples. In both cases, VaPoR exhibited a  
15 robust validation score up to approximately 200bp overall, with some slight differences observed  
16 between different SV types (Figure 3c,d).

17  
18  
19 We also compared VaPoR against a long-read validation approach developed in conjunction with  
20 Lumpy [3] using SVs on chr1 of NA12878 reported by the 1000 Genomes Project Phase 3. VaPoR  
21 achieved a sensitivity of 72% for deletions and 86% for insertions, while the Lumpy-associated  
22 approach was only able to assess 11% and 0% respectively. Both approaches exhibited a very low false  
23 validation rate when synthetically assigning the variants to chr2, with 0 for all SV types by the Layer et  
24 al approach and varying between 0 and 2.5% for VaPoR (Supplemental Table 4).

### 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 **Discrimination of SV types and genotypes**

42  
43  
44 We identified a small number of SVs in the high quality 1000 Genomes set that did not validate with  
45 VaPoR. Previous studies have shown that complex rearrangements are often misclassified as simple  
46 structural changes [5, 13], and indeed upon manual inspection these appeared to consist of multiple  
47 connected rearrangements. For example, we observed a reported inversion in HG00513 and NA19239  
48 on chromosome 1 (chr1:239952707-239953529) that was invalidated by VaPoR; an investigation into  
49 the long-reads aligned in the region showed the signature of an inverted duplication (Figure 4a) which,  
50 when incorporated into a modified reference that location, matched almost exactly with the read  
51 sequence (Figure 4b).

1  
2  
3  
4 We further explored the distribution of VaPoR scores for this region and others across the sample set  
5 and observed clear delineations between allelic copy number when fitted with a Gaussian mixture model  
6 allowing for the generation of genotype likelihoods for each site (Figure 4c). These tracked with our  
7 expected genotypes for the inverted duplication on chr1 across the 5 individuals queried while showing  
8 no support for the originally predicted inversion (Figure 4d). This shows that VaPoR is not only able to  
9 accurately genotype variants but can also distinguish between similar but distinct SV predictions in the  
10 same region.  
11  
12  
13  
14  
15  
16  
17

### 18 **Runtime and efficiency**

19  
20  
21 The computation runtime of VaPoR was assessed using 2 Intel Xeon Intel Xeon E7-4860 processors  
22 with 4GB RAM each on both simulated and real genomes. The runtime of simulated event was observed  
23 to increase linearly with read depth (Supplemental Figure 7). For events sequenced up to 20X, VaPoR  
24 takes ~3 seconds to assess a simple SV and ~5s for a complex event. The assessment of real samples  
25 sequenced at 20X required ~1.4 seconds to assess a simple deletion or insertion and ~6 seconds for an  
26 inversion (Supplemental Table 5).  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36

## 37 **DISCUSSION**

38  
39 Here we present an automated assessment approach, named VaPoR, for exploring various features of  
40 predicted genomic structural variants using long read sequencing data. VaPoR directly compares the  
41 input reads with the reference sequences with relatively straightforward computational metrics, thus  
42 achieving high efficiency in both run time and computing cost. VaPoR exhibits high sensitivity and  
43 specificity in both simulated and real genomes, with the capability of discriminating partially resolved  
44 SVs either consisting of similar but incorrect SV types at the same location or correct SVs with offset  
45 breakpoints. Furthermore, we show that VaPoR performs well at low read depths (5-10X), thus  
46 providing the option of systematically assessing large-scale SVs with a lower sequencing cost.  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

## METHODS

### VaPoR Workflow

VaPoR takes in aligned sequence reads in BAM format and predicted SVs (>50bp) in various formats including VCF and BED. Evaluation of an SV is performed by comparing long reads that go through the event against reference sequences in two formats: (a) the original human reference to which the sample is aligned and (b) a modified reference sequence altered to match the predicted structural rearrangement. A recurrence matrix is then derived by sliding a fixed-size window with 1bp step through each read to mark positions where the read sequence and reference are identical. The matching patterns are then assessed as to the validity of the SV as described below and a validation score is reported. Given the large variance of SVs lengths, each SV is stratified into one of two groups: smaller SVs that can be completely encompassed within multiple (>10 by default) long sequences and larger events that are rarely covered by individual long reads, with different statistical model applied. The VaPoR workflow is briefly summarized in Figure 1.

#### *Small Variants Assessment:*

For an SV  $k$  in sample  $s$  that is covered by  $n$  reads, the recurrence matrix between each read and the reference sequences in original ( $R_o$ ) and altered ( $R_a$ ) format is calculated. The vertical distance between each record ( $x_{i,k,s,Rx}$ ,  $y_{i,k,s,Rx}$ ) in matrix  $x$  and the diagonal ( $x_{i,k,s,Rx}$ ,  $x_{i,k,s,Rx}$ ) line is calculated as  $d_{i,k,s,Rx} = \text{abs}(x_{i,k,s,Rx} - y_{i,k,s,Rx})$ , and the average distance of all records would be exported as the score of each matrix:

$$\text{Score}_{k,s,Rx} = \sum_{i=1}^m d_{i,k,s,Rx} / m,$$

where  $m$  is the total number of records in the matrix. Sequences that share higher identity with the read shall have a lower  $\text{Score}_{k,s,Rx}$ , such that the score of each read is normalized as:

$$\text{Score}_{k,s,R} = \text{Score}_{k,s,R_o} / \text{Score}_{k,s,R_a} - 1,$$



1  
2  
3  
4 where a positive  $Score_{k,s,R}$  represents the superiority of the predicted structure versus the original and  
5 vice versa for negative  $Score_{k,s,R}$ , with one exceptional case where there exists a duplicated structure in  
6 the predicted SV such that the predicted structure would show higher  $Score_{k,s,R}$  due to the multi-  
7 alignment of duplicated segments. To correct for duplications, VaPoR adopts the directed distance  
8  $d_{i,k,s,Rx} = x_{i,k,s,Rx} - y_{i,k,s,Rx}$  instead such that the distance contributed by centrosymmetric duplicated  
9 segments would offset each other.  
10  
11  
12  
13  
14

15  
16 *Large Variants Assessment:*  
17

18  
19 For larger SVs where there are few, if any, long reads that can transverse the predicted SV, VaPoR  
20 assesses the quality of each predicted junction instead using:  
21  
22

$$23$$

$$24$$

$$25 \quad Score_{k,s,Rx} = \frac{\sum_{i=1}^m I = \begin{cases} 1, & \text{if } abs(x_{i,k,s,Rx} - y_{i,k,s,Rx}) < 0.15 * x_{i,k,s,Rx} \\ 0, & \text{otherwise} \end{cases}}{m},$$

$$26$$

$$27$$

$$28$$

$$29$$

30 where a larger  $Score_{k,s,Rx}$  represents higher similarity between the read and the reference sequence. The  
31 normalized scores of each read is then defined as:  
32  
33

$$34$$

$$35 \quad Score_{k,s,R} = Score_{k,s,Ra} / Score_{k,s,Ro} - 1,$$

$$36$$

$$37$$

38 *VaPoR Score Calculation:*  
39

40  
41 With a score assigned to each read spanning through the predicted structural variants, the VaPoR score  
42 is summarized as:  
43  
44

$$45$$

$$46 \quad Score_{k,s} = \frac{\sum_{R=1}^n I = \begin{cases} 1, & \text{if } Score_{k,s,R} > 0 \\ 0, & \text{otherwise} \end{cases}}{n}$$

$$47$$

$$48$$

$$49$$

$$50$$

51 to represent the proportion of long reads supporting predicted structure.  
52  
53

54 The highest supportive score ( $\max(Score_{k,s,R})$ ) is also reported as a reference for users to meet the  
55 specific requirement of their study design, for which we recommend 0.1 as the cutoff.  
56  
57

58  
59 *Flexible window size:*  
60  
61  
62  
63  
64  
65

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
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

By default, VaPoR uses a window size of 10bp and requires an exact match between sequences, though these can be changed to user-defined parameters. However, many regions of the genome contain repetitive sequences resulting in an abundance of spurious matches in the recurrence matrix, thus introducing bias to the assessment. To address this, VaPoR adopts a quality control step by iteratively assessing the reference sequence against itself and tabulating the proportion of matches along the diagonal. The window size initially starts at 10bp and iteratively increases by 10bp until either (a) the proportion of matches on the diagonal exceeds 40% and the current window size is kept or (b) the window size exceeds 40bp whereby the event will be labeled as ‘non-assessable and excluded from the evaluation.

1  
2  
3  
4 **FIGURE LEGENDS**  
5  
6

7 **Figure 1. Flowchart describing the VaPoR algorithm.** As input, the algorithm requires a set of  
8 structural variants in either VCF or BED format, a series of long reads and/or sequence contigs in BAM  
9 format, and the corresponding reference sequence. VaPoR then interrogates each variant individually at  
10 its corresponding reference location, assesses the quality of the region and assigns a score.  
11  
12  
13  
14

15 **Figure 2. Accuracy of VaPoR on simulated heterozygous and homozygous SVs at varying degrees**  
16 **of sequence coverage and VaPoR score cut-offs.** Receiver operator curves (ROC) are shown for  
17 simple deletions, duplications and inversions (a,b) as well as complex rearrangements including inverted  
18 duplications and deletion-inversion rearrangements (c,d).  
19  
20  
21  
22  
23

24 **Figure 3. Validation rate and breakpoint accuracy of VaPoR on the 1000 Genomes Projects phase**  
25 **3 calls.** VaPoR was applied on 5 individuals with reported SVs as a truth set: HG00513, HG00731,  
26 HG00732, NA19238, NA19239. The validation rate of deletions (a) and insertions (b) are shown here  
27 across different cutoff scores for VaPoR. Robustness to breakpoint accuracy was assessed by deviating  
28 breakpoints from their actual positions across varying distances for deletions (c) and insertions (d).  
29  
30  
31  
32  
33

34 **Figure 4. Validation and genotyping of assessed regions using VaPoR.** (a) Recurrence plot of  
35 reference genome (GRCh38) to an aligned long read in NA19239  
36 (m150208\_160301\_42225\_c100732022550000001823141405141504\_s1\_p0/3831/0\_12148) for a  
37 reported inversion at position chr1:239952707-239953529. The signature is consistent with an inverted  
38 duplication structure. (b) Recurrence plot of a different read  
39 (m150216\_212941\_42225\_c100729442550000001823151505141565\_s1\_p0/106403/0\_13205) against  
40 the same location, consistent with a non-variant (reference) structure. (c) Distribution of VaPoR scores  
41 on all reported SVs on chr1 in samples HG00513, HG00731, HG00732, NA19238, NA19239, stratified  
42 by color (solid) and modeled with a Gaussian mixture model (dashed). (d) VaPoR scores of SV above  
43 now stratified by color as indicated in (c) for both reported inversion (red) and predicted inverted  
44 duplication (blue).  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 **TABLES**  
5  
6  
7

8 **Table 1.** Sensitivity and false discovery rate of different SV types  
9

	deletion	insertion	inversion
Sample	Sens/FDR	Sens/FDR	Sens/FDR
HG00513	0.96/0.00 (0.94 <sup>1</sup> )	0.80/0.05 (0.93)	0.50/0.00 (0.71)
HG00731	0.94/0.00 (0.96)	0.85/0.07 (0.97)	0.60/0.00 (1.00)
HG00732	0.92/0.00 (0.98)	0.92/0.08 (0.96)	0.33/0.00 (0.86)
NA19238	0.90/0.00 (0.93)	0.88/0.10 (0.96)	1.00/0.00 (1.00)
NA19239	0.87/0.02 (0.95)	0.73/0.09 (0.96)	0.33/0.00 (1.00)

10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28 <sup>1</sup>Proportion of SVs passed VaPoR QC, as listed in brackets, are counted  
29 for events on chr1 and chr2 together.  
30  
31  
32  
33  
34  
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

1  
2  
3  
4 **Availability and Requirements**  
5

6  
7 Project name: VaPoR  
8

9  
10 Project home page: <https://github.com/millslab/vapor>  
11

12 Operating systems: Linux, OS X  
13

14 Programming languages: Python, R  
15

16  
17 Other requirements: Python v2.7.8+, rpy2, HTSeq, samtools v0.19+, pyfasta v0.5.2+, and pysam  
18 0.9.1.4+.  
19

20  
21 **Acknowledgements**  
22

23  
24 We thank the Human Genome Structural Variation Consortium (HGSVC) for generating and providing  
25 the deep PacBio sequencing. We also thank Yuanfang Guan and Kerby Shedden for discussions over  
26 specific statistical considerations.  
27  
28

29  
30 **Funding**  
31

32  
33 This work was supported by the National Institutes of Health [R01HG007068]. AMW was supported by  
34 the Genome Science Training Program at the University of Michigan [5T32HG000040]  
35  
36  
37

38  
39 **Authors' contributions**  
40

41 XZ designed the algorithm, wrote the program, comparatively benchmarked the different algorithms,  
42 and wrote the manuscript. AMW generated simulated data, aided in assessment testing, and revised the  
43 manuscript. REM conceived the study, modified the algorithm, and revised the manuscript. All authors  
44 read and approved the final manuscript.  
45  
46  
47  
48

49  
50 **Competing interests**  
51

52  
53 None declared.  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 **REFERENCES**  
5  
6

- 7 1. Brand H, Pillalamarri V, Collins RL, Eggert S, O'Dushlaine C, Braaten EB, et al. Cryptic and  
8 complex chromosomal aberrations in early-onset neuropsychiatric disorders. *Am J Hum*  
9 *Genet.* 2014;95 4:454-61. doi:10.1016/j.ajhg.2014.09.005.
- 10 2. Chiang C, Jacobsen JC, Ernst C, Hanscom C, Heilbut A, Blumenthal I, et al. Complex  
11 reorganization and predominant non-homologous repair following chromosomal breakage  
12 in karyotypically balanced germline rearrangements and transgenic integration. *Nature*  
13 *Genetics.* 2012;44 4:390-U195. doi:10.1038/ng.2202.
- 14 3. Layer RM, Chiang C, Quinlan AR and Hall IM. LUMPY: a probabilistic framework for  
15 structural variant discovery. *Genome Biol.* 2014;15 6:R84. doi:10.1186/gb-2014-15-6-r84.
- 16 4. Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V and Korbel JO. DELLY: structural variant  
17 discovery by integrated paired-end and split-read analysis. *Bioinformatics.* 2012;28  
18 18:i333-i9. doi:10.1093/bioinformatics/bts378  
19  
20  
21  
22 bts378 [pii].
- 23 5. Zhao X, Emery SB, Myers B, Kidd JM and Mills RE. Resolving complex structural genomic  
24 rearrangements using a randomized approach. *Genome Biology.* 2016;17  
25 doi:10.1186/s13059-016-0993-1.
- 26 6. Chong Z, Ruan J, Gao M, Zhou W, Chen T, Fan X, et al. novoBreak: local assembly for  
27 breakpoint detection in cancer genomes. *Nat Methods.* 2017;14 1:65-7.  
28 doi:10.1038/nmeth.4084.
- 29 7. Rhoads A and Au KF. PacBio Sequencing and Its Applications. *Genomics Proteomics*  
30 *Bioinformatics.* 2015;13 5:278-89. doi:10.1016/j.gpb.2015.08.002.
- 31 8. Travers KJ, Chin CS, Rank DR, Eid JS and Turner SW. A flexible and efficient template format  
32 for circular consensus sequencing and SNP detection. *Nucleic Acids Res.* 2010;38 15:e159.  
33 doi:10.1093/nar/gkq543.
- 34 9. Chaisson MJ, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, et al.  
35 Resolving the complexity of the human genome using single-molecule sequencing. *Nature.*  
36 2015;517 7536:608-11. doi:10.1038/nature13907.
- 37 10. Pendleton M, Sebra R, Pang AW, Ummat A, Franzen O, Rausch T, et al. Assembly and diploid  
38 architecture of an individual human genome via single-molecule technologies. *Nat*  
39 *Methods.* 2015; doi:10.1038/nmeth.3454.
- 40 11. Shi L, Guo Y, Dong C, Huddleston J, Yang H, Han X, et al. Long-read sequencing and de novo  
41 assembly of a Chinese genome. *Nature communications.* 2016;7:12065.  
42 doi:10.1038/ncomms12065.
- 43 12. Koren S, Schatz MC, Walenz BP, Martin J, Howard JT, Ganapathy G, et al. Hybrid error  
44 correction and de novo assembly of single-molecule sequencing reads. *Nat Biotechnol.*  
45 2012;30 7:693-700. doi:10.1038/nbt.2280.
- 46 13. Huddleston J, Chaisson MJ, Meltz Steinberg K, Warren W, Hoekzema K, Gordon DS, et al.  
47 Discovery and genotyping of structural variation from long-read haploid genome sequence  
48 data. *Genome Res.* 2016; doi:10.1101/gr.214007.116.
- 49 14. Carvalho AB, Dupim EG and Goldstein G. Improved assembly of noisy long reads by k-mer  
50 validation. *Genome Res.* 2016;26 12:1710-20. doi:10.1101/gr.209247.116.
- 51 15. Gibbs AJ and McIntyre GA. The diagram, a method for comparing sequences. Its use with  
52 amino acid and nucleotide sequences. *European journal of biochemistry.* 1970;16 1:1-11.  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

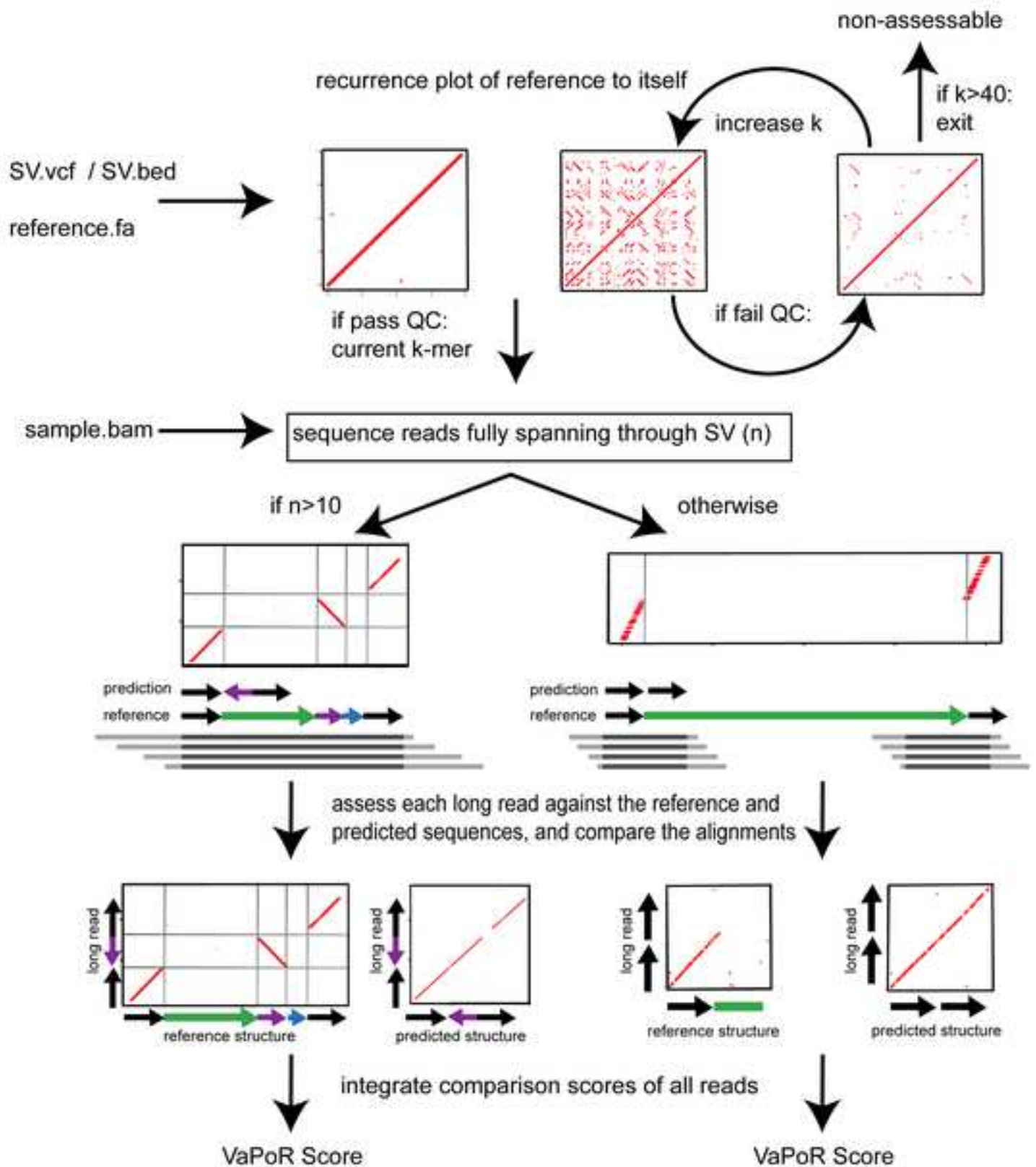
1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26  
27  
28  
29  
30  
31  
32  
33  
34  
35  
36  
37  
38  
39  
40  
41  
42  
43  
44  
45  
46  
47  
48  
49  
50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

16. Consortium EP. An integrated encyclopedia of DNA elements in the human genome. *Nature*. 2012;489 7414:57-74. doi:10.1038/nature11247.

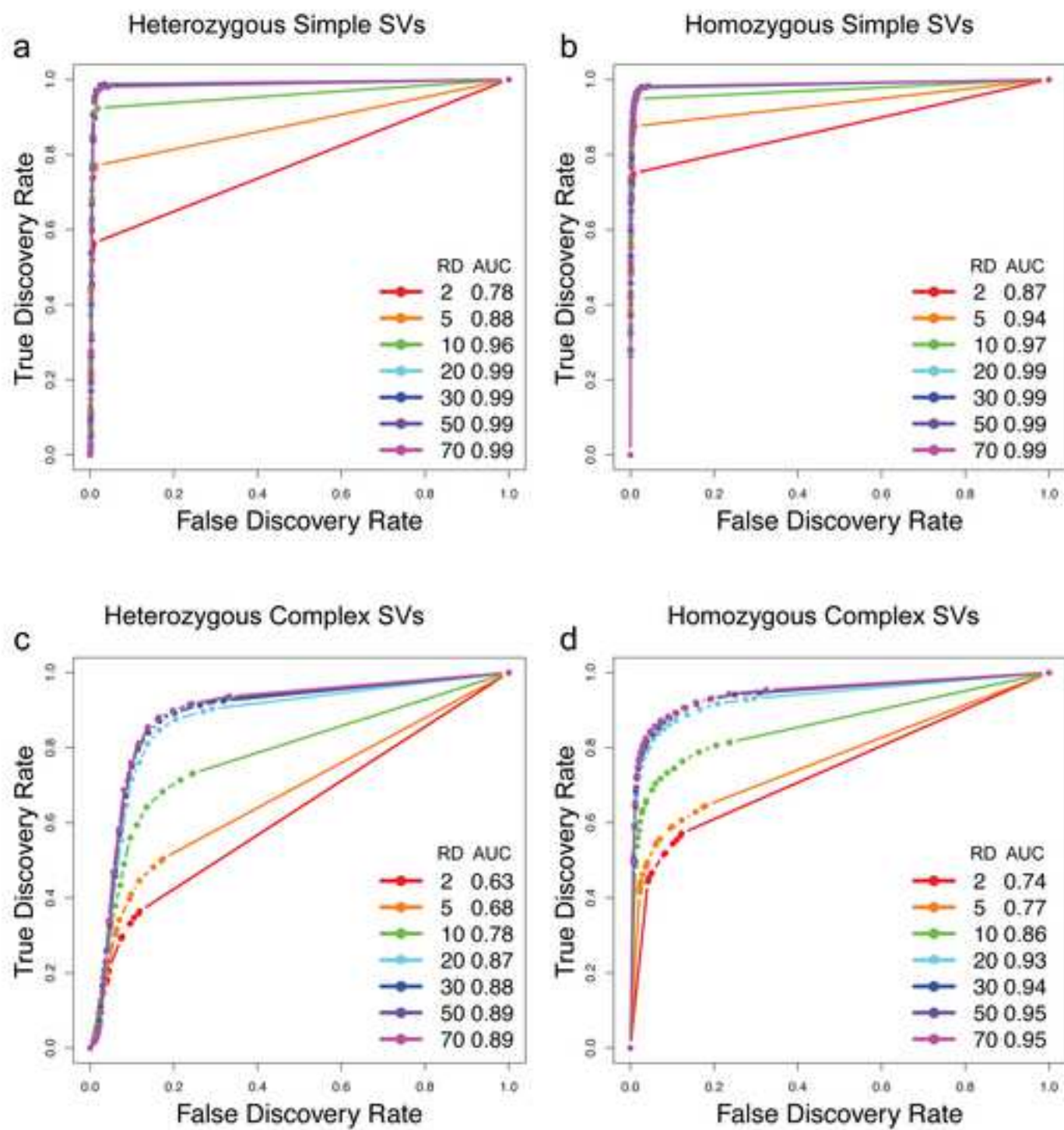
17. Ono Y, Asai K and Hamada M. PBSIM: PacBio reads simulator--toward accurate genome assembly. *Bioinformatics*. 2013;29 1:119-21. doi:10.1093/bioinformatics/bts649.

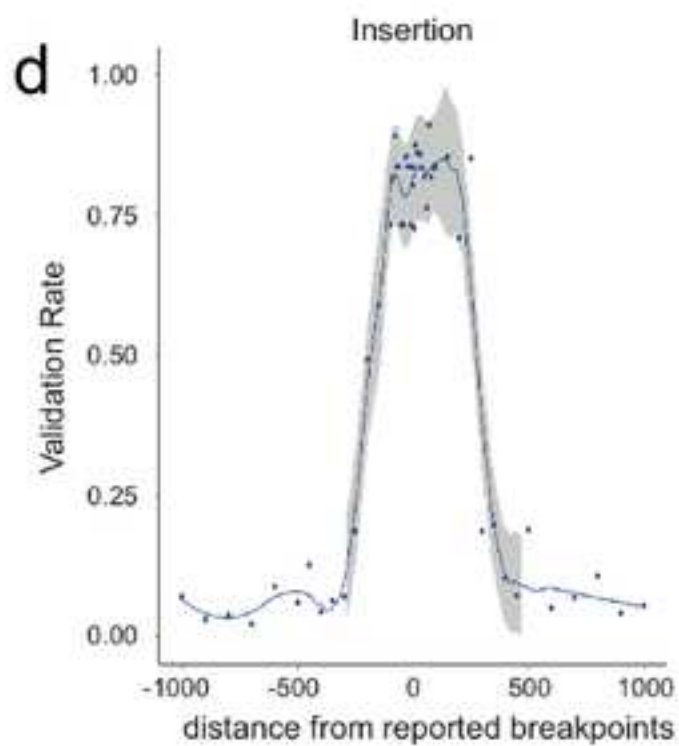
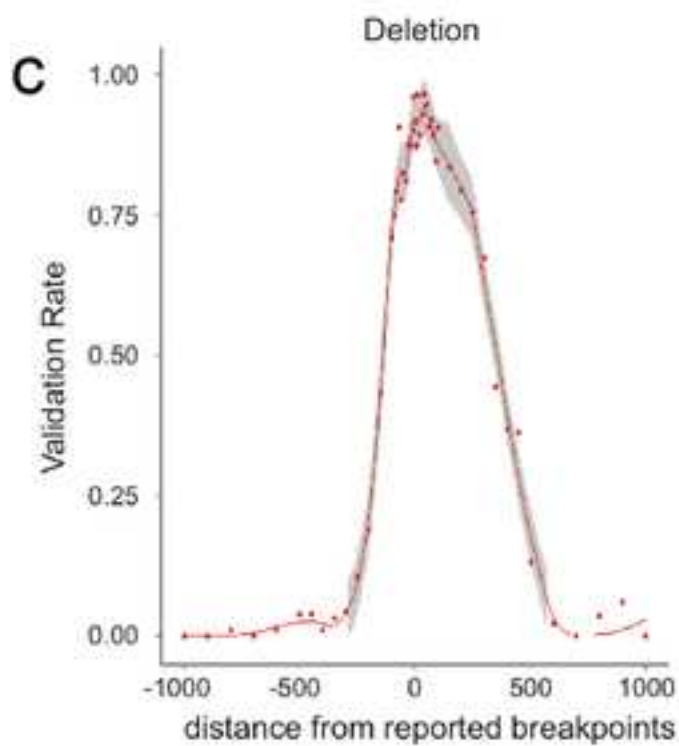
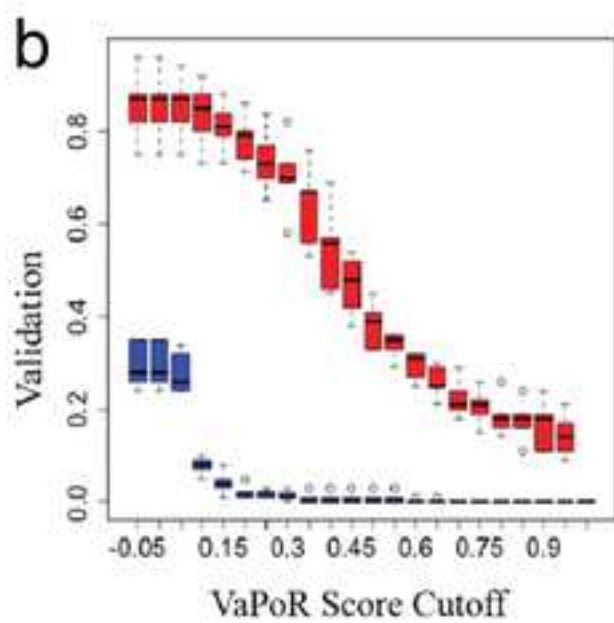
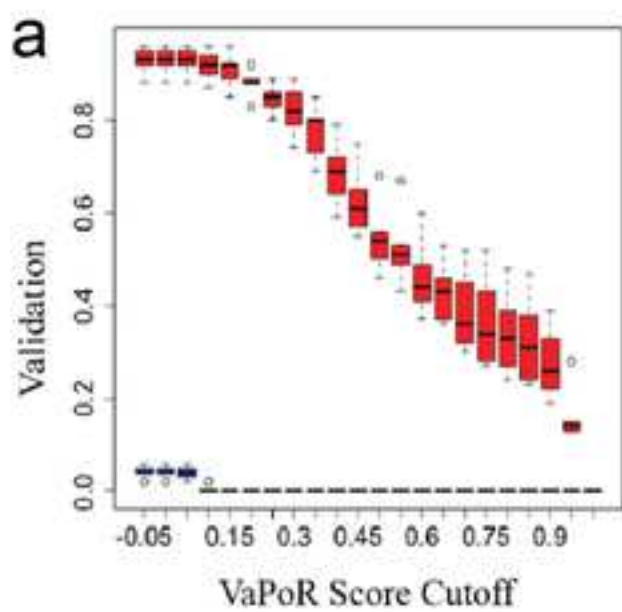
18. Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al. An integrated map of structural variation in 2,504 human genomes. *Nature*. 2015;526 7571:75-+. doi:10.1038/nature15394.

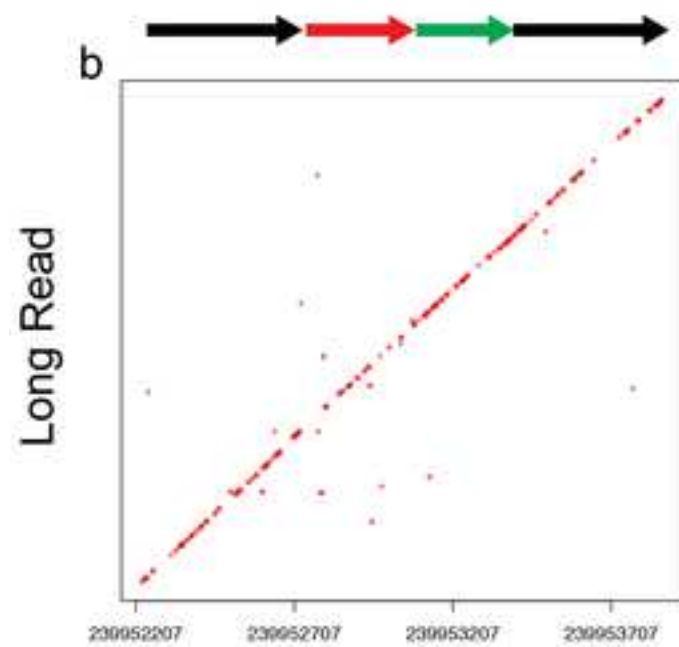
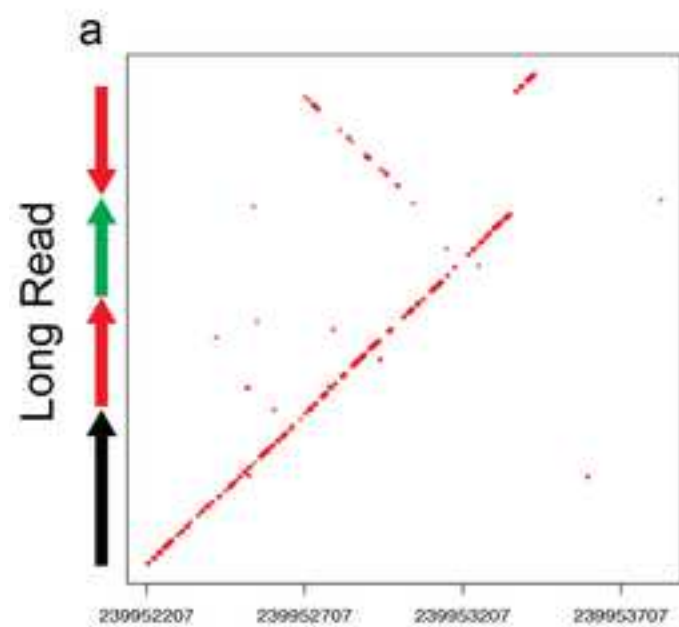
19. Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, et al. A global reference for human genetic variation. *Nature*. 2015;526 7571:68-74. doi:10.1038/nature15393.



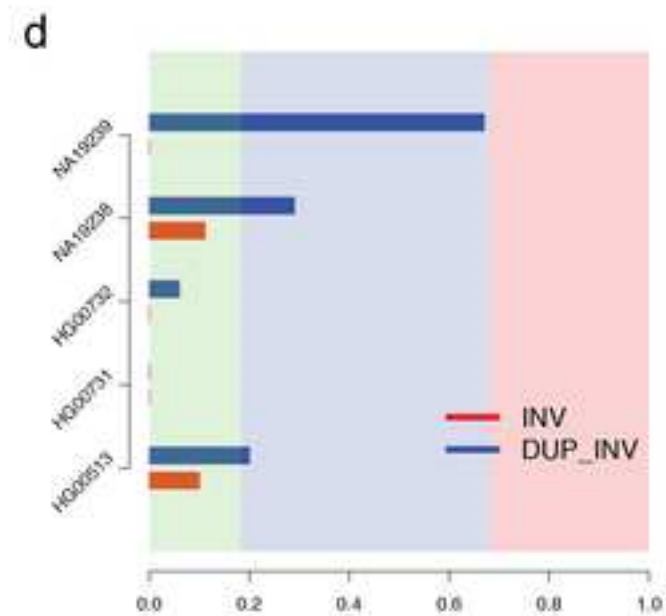
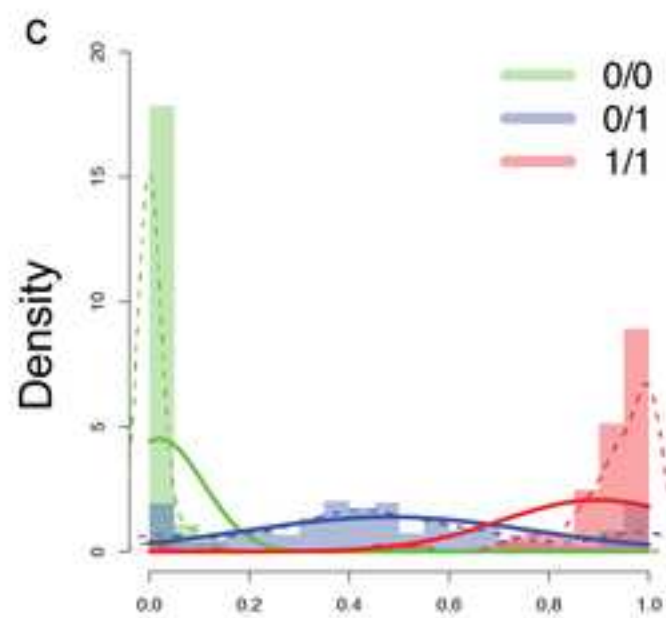








Chromosome 1



VaPoR Score



Click here to access/download  
**Supplementary Material**  
supplementary\_figures.docx





Click here to access/download  
**Supplementary Material**  
supplementary\_tables.xlsx