

# Supplement to: Resolving the Full Spectrum of Human Genome Variation using Linked-Reads

*Patrick Marks<sup>a</sup>, Sarah Garcia<sup>a</sup>, Alvaro Martinez Barrio<sup>a</sup>, Kamila Belhocine<sup>a</sup>, Jorge Bernate<sup>a</sup>, Rajiv Bharadwaj<sup>a</sup>, Keith Bjornson<sup>a</sup>, Claudia Catalanotti<sup>a</sup>, Josh Delaney<sup>a</sup>, Adrian Fehr<sup>a</sup>, Ian T. Fiddes<sup>a</sup>, Brendan Galvin<sup>a</sup>, Haynes Heaton<sup>a,e,f</sup>, Jill Herschleb<sup>a</sup>, Christopher Hindson<sup>a</sup>, Esty Holt<sup>b</sup>, Cassandra B. Jabara<sup>a,g</sup>, Susanna Jett<sup>a,h</sup>, Nikka Keivanfar<sup>a</sup>, Sofia Kyriazopoulou-Panagiotopoulou<sup>a,i</sup>, Monkol Lek<sup>c,d</sup>, Bill Lin<sup>a</sup>, Adam Lowe<sup>a</sup>, Shazia Mahamdallie<sup>b</sup>, Shamoni Maheshwari<sup>a</sup>, Tony Makarewicz<sup>a</sup>, Jamie Marshall<sup>d</sup>, Francesca Meschi<sup>a</sup>, Chris O'keefe<sup>a</sup>, Heather Ordonez<sup>a</sup>, Pranav Patel<sup>a</sup>, Andrew Price<sup>a</sup>, Ariel Royall<sup>a</sup>, Elise Ruark<sup>b</sup>, Sheila Seal<sup>b</sup>, Michael Schnall-Levin<sup>a</sup>, Preyas Shah<sup>a</sup>, David Stafford<sup>a</sup>, Stephen Williams<sup>a</sup>, Indira Wu<sup>a</sup>, Andrew Wei Xu<sup>a</sup>, Nazneen Rahman<sup>b</sup>, Daniel MacArthur<sup>c,d</sup>, Deanna M. Church<sup>a,j</sup>*

2019-01-24 23:07:26

**Author affiliations** a: 10x Genomics, 7068 Koll Center Parkway, Suite 401, Pleasanton, CA 94566; b: The Institute of Cancer Research, Division of Genetics & Epidemiology, 15 Cotswold Road, London, SM2 5NG, UK; c: Analytic and Translational Genetics Unit, Massachusetts General Hospital, Boston, MA; d: Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA; e: Current affiliation, Wellcome Trust Sanger Institute, Hinxton CB10 1SA, UK; f: Current affiliation, University of Cambridge, Cambridge, UK; g: Current affiliation, Purigen Biosystems, Inc., 5700 Stoneridge Drive, Suite 100, Pleasanton, CA 94588; h: Current affiliation, LevitasBio, Inc., 3283 25th Street, 3, San Francisco, CA 94110; i: Current affiliation,

22 Illumina, Inc., 499 Illinois Street, Suite 201, San Francisco, CA 94158; j: Current affiliation, Inscripta  
23 Inc., 5500 Central Avenue #220, Boulder, CO 80301

## 24 **Methods**

### 25 **1 Lariat aligner**

26 Lariat is an aligner for barcoded Linked-Reads. All the Linked-Reads for a single barcode are  
27 aligned simultaneously, with the prior knowledge that the reads arise from a small number of long  
28 (10 kb - 200 kb) molecules. Lariat is an implementation of the RFA method (Bishara et al. 2015).  
29 Briefly, we model the observed reads for one barcode as being generated by a hierarchical process  
30 which first selects a small number of loci on the genome corresponding to the long input  
31 molecules, covering on average 500 kb of the genome. Then short reads are sampled with a  
32 uniform distribution over the selected loci. The sequencing process that generates the observed  
33 read sequence from the genome is modeled by the standard Smith-Waterman scoring scheme used  
34 by e.g. BWA-MEM (Li 2013). For each read a set of feasible candidate alignments is generated with  
35 traditional short-read alignment methods. An alignment configuration is a choice of one alignment  
36 from the feasible set for each read. We search for an alignment configuration that maximizes the  
37 likelihood of the data under the RFA model. A MAPQ can be derived from the likelihood ratio of  
38 the optimal alignment configuration to the sum of suboptimal configuration that select a different  
39 alignment for the read. The molecule selection process induces a strong prior that the aligned  
40 positions of reads cluster together on the genome. Reads with near-identical alignments to  $>1$   
41 locus would be assigned  $\text{MAPQ} < 10$  in typical short-read data, and could not be used for variant  
42 calling. In the RFA method, confidently mapped reads flanking a duplicated region will anchor the  
43 molecule to the correct locus, and the molecular prior will strongly favor the alignments proximal  
44 to the confidently placed molecule, allowing the assignment of  $\text{MAPQ} > 40$  for reads with two  
45 identical alignments.

46 Lariat is written in the Go language and is available at <https://github.com/10xGenomics/lariat>.  
47 Upstream stages in the Long Ranger pipeline extract and correct the molecular barcodes, and  
48 prepare barcode sorted FASTQ-like inputs. Lariat generates candidate alignment positions by  
49 calling the BWA (Li 2013) API. It then performs RFA inference to select the final mapping position  
50 and MAPQ, and emits alignment records to BAM.

## 51 **1.1 Adversarial alignments**

52 We discovered a surprisingly high rate of degenerate alignments to segmental duplications with  
53 strong molecular evidence for one locus, but a better alignment score (typically by a single  
54 mismatch) at the other locus. In this case the RFA model will typically select the position supported  
55 by the molecular evidence, but with low MAPQ. Typically ~50% or 100% of reads mapped to the  
56 mismatch position supported the alternate allele suggesting the presence of a variant. We postulate  
57 that these cases are an expected feature of reference-based analysis of segmental duplications.  
58 Studies of copy-number in segmental duplication have cataloged singly unique nucleotides (SUNs),  
59 which are bases within one copy of a duplication that uniquely tag that copy (Sudmant et al. 2010).  
60 Typically SUNs alleles are introduced after the duplication event. Reads carrying the ancestral  
61 allele at a SUN position will be biased away from the SUN position to an ancestral copy. We term  
62 such reads ‘adversarial’ since the best alignment is not the correct mapping position.

63 We adopt a proposal from the RFA authors (A. Bishara, Y Lui, S. Batzoglou, private  
64 communication), to allow a collection of mapped reads from multiple barcodes to overturn  
65 reference alleles and realize a MAPQ improvement. Lariat implements a limited form of this  
66 approach. Reads are initially mapped in independent groups for each barcode. The second best  
67 alignment score and log-likelihood ratio of the molecular positioning analysis is stored in the BAM  
68 record. After mapping, Lariat scans the read pileups looking for sites with  $\geq 3$  alternate alleles  
69 among reads with molecule support for this location. In this case, the best explanation for the data  
70 is the presence of variant in the sample at this locus, rather than independent sequencing errors on

71 each read. We recompute the MAPQ for each read containing an alternate allele, but divide the  
72 mismatch error induced by the putative variant among all the reads supporting the variant. For  
73 isolated adversarial SNPs, MAPQ is typically increased from MAPQ=3 to MAPQ=40, leading to  
74 ~15,000 additional variant calls in degenerate regions.

## 75 **1.2 Additional genome coverage gained with Linked-Reads and Lariat**

76 To further investigate the properties of the parts of the genome with alignment coverage that is  
77 unique to a method, either Linked-Reads + Long Ranger (CrG) or PCR- short reads + BWA  
78 (TruSeq), we first looked at the distribution of these regions across the genome (Supplemental  
79 Figure S2). In the TruSeq data, the decoy sequence (hs37d5) has the greatest amount of unique  
80 sequence alignment in the two female samples, with the rest of the regions distributed roughly  
81 equally among the other chromosomes. This pattern is different for the one male sample  
82 (NA24385), where we see the largest sequence gain on the X chromosome.

83 In the CrG data, there is a completely different distribution of regions unique alignments. In all  
84 samples, chr9 shows the largest gain. This is driven by the ability to align sequences around the  
85 repetitive pericentromeric regions (Supplemental Figure S2). The ability of Lariat to resolve  
86 multi-mapping reads is a function of genome structure; the repeats need to be far enough apart that  
87 they are unlikely to share barcodes. This is reflected in the pattern of uniquely aligning sequence  
88 regions in CrG. Of note, there is a substantial gain of aligned regions on chrY in NA24385.

89 We then compared the uniquely aligning regions to exon and segmental duplication annotations.  
90 For both the TruSeq and CrG samples, we see roughly 5% of the uniquely aligning regions  
91 correspond to exon annotations using bedtools (Supplemental Table S1)(Quinlan and Hall 2010).  
92 For the TruSeq regions, we see a range of 17-30% of regions overlapping segmental duplications,  
93 and 13-28% aligning to the decoy. For the CrG samples, roughly 77% of the uniquely aligning  
94 regions correspond to segmental duplication annotations, and 12-14% to decoy alignments.

### 95 **1.3 Comparison to PacBio**

96 Raw PacBio FASTQs were aligned to the reference using BWA-MEM -x pacbio (Li 2013). To test a  
97 variant, we fetch all PacBio reads covering the variant position, and retain the substring aligned  
98 within 50 bp of the variant on the reference. We re-align the PacBio read sequence to the +/-50 bp  
99 interval of the reference, and the same interval with the alternate allele applied. A read is  
100 considered to support the alternate allele if the alignment score to the alt-edited template exceeds  
101 the alignment score of the reference template. A variant was considered to be validated if at least 2  
102 PacBio reads supported the alt allele, at least 10 PacBio reads covered the locus, and the overall  
103 alternate allele fraction seen in the PacBio reads was at least 25%.

## 104 **2 Variant Phasing**

105 A variety of methods for phasing haplotypes have been proposed, which optimize a variety of  
106 different objective functions (Bansal et al. 2008; Bansal and Bafna 2008). The HASH (Bansal et al.  
107 2008) phasing method optimizes the likelihood of generative probabilistic objective function that is  
108 a natural model of reads generated from haplotypes. The method was designed for ~800 bp Sanger  
109 reads as input fragments. HASH uses a MCMC approach to optimize the objective function, which  
110 may lead to long running times.

111 We build on this basic approach, and extend the probabilistic model to be robust to mixed  
112 fragments which contain alleles from both haplotypes. We add a new variant hypothesis we term  
113 ‘non-heterozygous’ that allows the model to identify variants that were initially called as  
114 heterozygous but whose alleles do not cleanly segregate onto the local haplotypes, which may be  
115 false-positive calls, or homozygous variants incorrectly called as heterozygous. We use local  
116 realignment to both allele sequences to carefully quantify the allele supported by a read, which is  
117 critical for good phasing performance of indel variants. Finally we develop an efficient search  
118 heuristic that combines direct beam-search to phase chunks of ~50 variants, a greedy stitching pass

119 to phase chunks to each other, and final polishing pass to correct any local errors.

## 120 **2.1 Introduction**

121 We take as input a pre-determined set of of biallelic variants. We label the alleles  $A_{i,p}$  where  
122  $i \in 1, \dots, N$  indexes the variant, and  $p \in 0, 1$  is an arbitrary label for the two alleles of the variant.

123 The set of alleles that come from the same parent chromosome is referred to as a haplotype, and  
124 are arbitrarily labeled  $H_0$  and  $H_1$ . The goal of the phasing algorithm is to determine which allele  
125 from each variant came from each parent chromosome. The phasing result can be described by a  
126 binary variable for each variant  $X_i \in 0, 1$  where  $X_i = 0$  indicates the  $A_{i,0} \in H_0$  and  $A_{i,1} \in H_1$   
127 and  $X_i = 1$  indicates that  $A_{i,0} \in H_1$  and  $A_{i,1} \in H_0$ .

128 Neighboring variants on the genome are often separated by distances longer than the read-pair  
129 length, causing very short phase blocks. Long input fragments covering a small fraction ( 0.001) of  
130 the genome are exposed to each barcode, so the probability that a barcode contains reads from  
131 both haplotypes is small.

132 We cast the solution to the phasing problem as a search for the maximum likelihood phasing parity  
133 vector:

$$\hat{\mathbf{X}} = \underset{\mathbf{X}}{\operatorname{argmax}} P(\mathbf{O}|\mathbf{X})$$

134 where  $\mathbf{O}$  denotes the sets of barcoded reads observed, and  $\mathbf{X}$  is the phasing result we wish to infer.

135 Read pairs are aligned to the genome as usual. Reads are grouped by the attached barcode  
136 sequences. Reads with common barcodes are partitioned into groups that are likely to have  
137 originated from a single genomic input fragment, and thus provide evidence that the alleles  
138 covered by the reads came from the same haplotype.

139 We compute the probability of the observed reads covering variant  $i$  from fragment  $f$  as:

$$\log P(O_{i,f}|A_{i,p}) = \sum_{r \in O_{i,f}} \mathbf{1}(S_r = A_{i,p})(1 - 10^{-Q_r/10}) + \mathbf{1}(S_r \neq A_{i,p})(10^{-Q_r/10})$$

140 where  $r$  sums over reads,  $\mathbf{1}(S_r = A_{i,p})$  is the indicator function testing if the  $r$ th sequence  $S_r$   
 141 match allele  $A_{i,p}$ . The probability assigned is derived from the inverse-Phred transformed quality  
 142 value of relevant read base  $Q_r$ .

143 The data from a fragment  $f$  comes from one of three cases. First two cases are that the alleles  
 144 present are only from  $H_0$  or only from  $H_1$ . These cases are the typical case and have a high prior  
 145 probability, governed by the fraction of the genome present in each partition. The third case is that  
 146 input DNA from both haplotypes was present at the locus, so both either allele is equally likely to  
 147 be observed:

$$P(O_{1,f}, \dots, O_{N,f}|\mathbf{X}, H_f = 0) = \prod_i P(O_{i,f}|A_{i,X_i})$$

$$148 \quad P(O_{1,f}, \dots, O_{N,f}|\mathbf{X}, H_f = 1) = \prod_i P(O_{i,f}|A_{i,1-X_i})$$

$$149 \quad P(O_{1,f}, \dots, O_{N,f}|\mathbf{X}, H_f = M) = \prod_i 0.5$$

150 These equations give the probability of the observed reads from fragment  $f$  at variant location  $i$ ,  
 151  $X_i$ , and fragment haplotype  $H_f$ . Observations are independent given the variant party and  
 152 fragment haplotype. The prior probability of third case is  $\alpha$  – the probability that a partition  
 153 contains both haplotypes at a locus. We can then compute the overall likelihood by summing over  
 154 the three cases:

$$P(O_{1,f}, \dots, O_{N,f}|\mathbf{X}) = \frac{(1 - \alpha)}{2} \left( \prod_i P(O_{i,f}|A_{i,X_i}) + \prod_i P(O_{i,f}|A_{i,1-X_i}) \right) + \alpha \prod_i 0.5$$

155 Fragments are independent given the variant party  $X_i$  letting us form the overall objective

156 function as:

$$P(\mathbf{O}|\mathbf{X}) = \prod_f P(O_{1,f}, \dots, O_{N,f}|\mathbf{X})$$

## 157 2.2 Optimization

158 We optimize the overall objective function using a hierarchical search over the phasing vector  $\mathbf{X}$ .

159 Initially we break up  $\mathbf{X}$  into local chunks of  $n \approx 40$  variants and determine the relative phasing of  
160 the block using beam search over the assignments of  $X_k, X_{k+1}, \dots, X_{k+n}$ . Where  $k$  is the first  
161 variant in the local block. Beam search is a standard method that has existed for a long time (see  
162 [http://en.wikipedia.org/wiki/Beam\\_search](http://en.wikipedia.org/wiki/Beam_search)).

163 The relative phasing of neighboring blocks is found greedily, yielding a candidate phasing vector  
164  $\mathbf{X}$ . Finally  $\mathbf{X}$  is iteratively refined by swapping the phase of individual variants. When refinement  
165 converges we are left with our estimate of the optimal phasing configuration  $\hat{\mathbf{X}}$ .

## 166 2.3 QV Testing

167 We can compute estimates of the accuracy of the phasing configuration by computing the  
168 likelihood ratio between the optimal configuration  $\hat{\mathbf{X}}$  and some alternate configuration  $\mathbf{X}_{alt}$  by  
169 computing the likelihood ratio between the hypotheses. The confidence is then reported as a  
170 Phred-scaled quality value:

$$Q(\mathbf{X}_{alt}) = -10 \log_{10} \left( \frac{P(\mathbf{O}|\mathbf{X}_{alt})}{P(\mathbf{O}|\hat{\mathbf{X}})} \right)$$

171 There are two classes of errors we consider: short switch errors and long switch errors. Short  
172 switch errors are single variants that are assigned the wrong phasing in an otherwise correctly  
173 phased region - to measure the short switch confidence of variant  $i$ , we flip  $X_i$  to form  $\mathbf{X}_{alt}$ . When  
174 the short switch confidence is low, the variant is marked as not phased in the output, rather than  
175 reporting a phasing call likely to be erroneous.

176 Long switch errors occur when two neighboring blocks of variants ...,  $X_{i-2}$ ,  $X_{i-1}$  and  $X_i$ ,  $X_{i+1}$ , ...  
177 are correctly phased internally, but have the wrong relative phasing between the two blocks. In  
178 this case we say a long switch error occurred at position  $i$ . We test the long switch confidence at  
179 position  $i$  by inverting the phase of  $X_j$  for all  $i \geq j$ . When the long switch confidence falls below  
180 a threshold we start a new phase block – variants in different phase blocks are not called as phased  
181 with respect to one another.

### 182 **3 SV calling from linked-reads**

#### 183 **3.1 Finding candidate regions with a lot of barcode overlap**

184 The goal of this step is to get a high-sensitivity/low-specificity list of potential SV candidates.  
185 Given two loci, we want a quick way to decide whether they share a significant number of  
186 common barcodes. The list of these loci will go into the next step of the algorithm, which uses a  
187 probabilistic calculation to make a more accurate prediction as to whether the observed barcode  
188 overlap is consistent with the presence of a structural variant.

189 **Expected barcode overlap between distant loci** If the two loci are on different chromosomes  
190 or the distance between them is much larger than the average molecule length, then we can use a  
191 binomial test to determine if the observed barcode overlap between the loci is larger than expected  
192 by chance. Let  $N_1$ ,  $N_2$ , and  $N$  be the observed number of barcodes at the first locus, the observed  
193 number of barcodes at the second locus, and the barcode diversity respectively. Then, the  
194 probability of observing  $n$  common barcodes between the two loci is governed by the binomial  
195 distribution:

$$Binom(n; N_1, N_2/N)$$

196 Therefore, we can pick a p-value cutoff and select all pairs of loci for which the above probability is  
197 less than our cutoff. These loci will serve as candidates for distal SVs.

198 **Expected barcode overlap between not so distant loci** The binomial test above assumes that

199 the two loci under consideration are independent in that no molecule can span both loci. This  
200 assumption clearly does not hold when the distance  $d$  between two loci is in the order of the  
201 molecule length.

202 Given the count of barcodes on each of the loci and the distance between them, we want to  
203 compute the expected number of common barcodes between the two loci. We start by computing  
204 the probability that a molecule with barcode  $b$  present at locus  $X$  will reach locus  $X + d$ .

$$f_b(d) = P(b \text{ present at } X + d | b \text{ present at } X) = \\ \sum_{m: L(m) > d} (P(\text{molecule at } X \text{ is } m) P(m \text{ present at } X + d | m \text{ present at } X))$$

205 Here the sum is over molecules  $m$  from barcode  $b$  with length  $L(m) > d$ . The first probability  
206 above is  $L(m) / \sum_{m'} L(m')$ . The second is  $(L(m) - d) / L(m)$ . So after simplifying we get  
207  $\sum_{m: L(m) > d} (L(m) - d) / \sum_{m'} (L(m'))$ . In practice, we get good results when simplifying this  
208 to  $\sum_{m: L(m) > d} L(m) / \sum_{m'} L(m')$ .

209 Given two loci at distance  $d$  with  $N_1$  and  $N_2$  barcodes respectively, we estimate the expected  
210 barcode overlap between them as

$$\min(N_1, N_2) \times f(d)$$

211 where  $f(d) = \text{avg}_b f_b(d)$

212 In practice, we pre-compute  $f(d)$  for a range of values of  $d$ . We can also further reduce the time  
213 required to compute  $f(d)$  by sampling a large number of barcodes instead of using all of them to  
214 compute the above average.

## 215 3.2 Probabilistic model

### 216 3.2.1 Setting up the maximum-likelihood problem

217 Given two candidate loci for structural variation, we want to determine whether the observed  
218 reads in the two loci are more consistent with the presence or the absence of an SV. In particular,  
219 we want to find the model that maximizes the data (log-)likelihood

$$\log P(D; m) = \sum_b \log P(D_b; m)$$

220 Here,  $D_b$  is the observed data from barcode  $b$  (at the loci of interest - the presence of the barcode at  
221 very distant loci is considered irrelevant). Data from different barcodes are independent  
222 (conditioning on the model).  $m$  is the model and comes from a discrete set of models:

- 223 1. no SV (no-SV or reference model),
- 224 2. homozygous SV at loci  $x$  and  $y$ , or
- 225 3. SV at loci  $x$  and  $y$  on haplotypes  $i$  and  $j$  respectively.

226 Here,  $x$  and  $y$  could be any pair of loci of the genome, but in practice we only consider a relatively  
227 small list of loci pairs, based for example on barcode overlaps or read-pair support.  $i$  and  $j$  are in  
228  $\{0, 1\}$  and denote the haplotype assignment of the breakpoints  $x$  and  $y$ . We further assume that if  
229  $x$  and  $y$  are on the same phase block, then  $i$  and  $j$  must be equal (i.e. the SV-calling cannot redefine  
230 phase blocks). We can further refine this set of SV models based on the type of the SV (more on  
231 this later).

232 There are two sets of latent variables,  $H_b^{x,y}$  the haplotype assignment of barcode  $b$  at loci  $x$  and  $y$ ,  
233 and  $M_b$ , the number of molecules from which the reads with barcode  $b$  were generated. For  
234 simplicity, we assume that  $M_b$  can be at most 2, since it is extremely unlikely that there are more  
235 than two molecules from the same locus in the same partition (or that we had multiple partitions  
236 with the same barcode).

237 Below is summary of notation:

- 238 1.  $D$  is the observed data (positions of reads, their barcodes, and the ph) in the loci under  
239 consideration.
- 240 2.  $D_b$  is the data (i.e. read positions) from barcode  $b$ .
- 241 3.  $D_{b_{1\dots k}}$  is a subset of  $D_b$  comprising the first  $k$  reads from barcode  $b$ .
- 242 4.  $R_b$  is the event that there is no SV on barcode  $b$  (or that  $b$  was generated from the reference).
- 243 5.  $SV_b^{x,y}$  is the event that there is an SV between positions  $x$  and  $y$  on the haplotype that  
244 generated  $b$ .
- 245 6.  $SV_{ij}^{x,y}$  is the event that there is an SV at positions  $x, y$  on haplotypes  $i$  and  $j$  respectively,  
246 where  $i, j \in \{0, 1\}$ .
- 247 7. We assume that reads are generated from a Poisson distribution with rate  $\alpha$  (uniform across  
248 the genome). That is,  $\alpha$  is the expected number of reads per basepair.
- 249 8.  $P_L(\ell)$  is the probability of having a molecule of length  $\ell$ . In practice, we use the empirical  
250 molecule length distribution.
- 251 9.  $L_{\max}$  is the maximum possible length of an input molecule.

### 252 3.2.2 Some useful probabilities

253 **Probability of a molecule** Let  $x_{b_1} \leq x_{b_2} \leq \dots \leq x_{b_n}$  be the positions of the reads from a single  
254 molecule with barcode  $b$ . We assume that the reads are generated from a single molecule with  
255 hidden length  $\ell$ . The distances  $x_{b_{i+1}} - x_{b_i}$  are the waiting times between events of a Poisson  
256 process. The log-probability of observing the molecule is:

$$\begin{aligned}
& \log P_m(x_{b_1}, x_{b_2}, \dots, x_{b_n}) = \\
& \log \left[ \sum_{\ell \geq x_{b_n} - x_{b_1}} P_L(\ell) \alpha e^{-\alpha(\ell - (x_{b_n} - x_{b_1}))} \prod_{i=1}^{n-1} \alpha e^{-\alpha(x_{b_{i+1}} - x_{b_i})} \right] = \\
& \log \left[ \alpha^n \prod_{i=1}^{n-1} e^{-\alpha(x_{b_{i+1}} - x_{b_i})} \sum_{\ell \geq x_{b_n} - x_{b_1}} P_L(\ell) e^{-\alpha(\ell - (x_{b_n} - x_{b_1}))} \right] = \\
& \log \left[ \alpha^n e^{-\alpha(x_{b_n} - x_{b_1})} \sum_{\ell \geq x_{b_n} - x_{b_1}} P_L(\ell) e^{-\alpha(\ell - (x_{b_n} - x_{b_1}))} \right] = \\
& \log \left[ \alpha^n \sum_{\ell \geq x_{b_n} - x_{b_1}} P_L(\ell) e^{-\alpha \ell} \right] = \\
& n \log \alpha + \text{logaddexp}_{\ell \geq x_{b_n} - x_{b_1}} [\log P_L(\ell) - \alpha \ell]
\end{aligned}$$

257 where logaddexp is the log of the sum of the exponentials of the arguments. Intuitively, the  
258 probability of observing the molecule is the product of the following probabilities:

- 259 1. The probability of getting a molecule of length  $\ell$  given that the molecule length was greater  
260 than  $x_{b_n} - x_{b_1}$ .
- 261 2. The probability of observing waiting times  $x_{b_{i+1}} - x_{b_i}$ .
- 262 3. The probability of observing no reads in a length  $\ell - (x_{b_n} - x_{b_1})$ .

263 These probabilities are then summed over all possible lengths  $\ell \geq x_{b_n} - x_{b_1}$ . Since  $P_m$  only  
264 depends on the observed length  $d = x_{b_n} - x_{b_1}$  and the number of reads  $n$ , below we will also use  
265 the (overloaded) notation  $P_m(n, d)$ .

266 **Barcode likelihood assuming no SV** The likelihood of the data from barcode  $b$  assuming that  
267 all of the data from barcode  $b$  were generated from a single molecule from the reference is:

$$P(D_b | M_b = 1; R_b) = P_m(n, d)$$

268 if  $x_{b_1}, \dots, x_{b_n}$  are all on the same chromosome and  $x_{b_n} - x_{b_1} < L_{\max}$ . Otherwise,

269  $P(D_b|M_b = 1; R_b) = \epsilon.$

270 Similarly:

$$P(D_b|M_b = 2; R_b) = \sum_{k=2}^{n-1} P(D_{b_{1\dots k}}|M_{b_{1\dots k}} = 1; R_{b_{1\dots k}})P(D_{b_{k+1\dots n}}|M_{b_{k+1\dots n}} = 1; R_{b_{k+1\dots n}})$$

271 More accurately, we would sum over all possible splits into two disjoint subsets. However, this  
 272 would add too much complexity (especially given how unlikely barcode collisions are and how few  
 273 molecules are typically within a partition), so we assume that molecules cannot overlap but can  
 274 “touch”.

275 **Barcode likelihood assuming a homozygous SV** The likelihood assuming that the data from  
 276 barcode  $b$  were generated from an SV haplotype  $P(D_b|M_b = 1; SV_b^{x,y})$  depends on the type of the  
 277 SV.

278 **Deletions** Assume that the SV is a deletion between  $x$  and  $y$  ( $x < y$ ) and that  $x_{b_i} < x \leq x_{b_{i+1}}$  and  
 279  $x_{b_j} < y \leq x_{b_{j+1}}$ .

280 1. If  $x > x_{b_n}$  or  $y < x_{b_1}$  then  $P(D_b|M_b = 1, SV_b^{x,y}) = P(D_b|M_b = 1; R_b)$ . We assume that SVs  
 281 are independent from each other, in that we can only have at most one SV within the length  
 282 of a molecule.

283 2. If  $i \neq j$ , this means that the molecule has reads inside the deletion, so

284  $P(D_b|M_b = 1, SV_b^{x,y}) = \epsilon.$

285 3. If none of the above holds, we have  $x_{b_1} \leq x_{b_2} \leq \dots \leq x_{b_i} < x < y \leq x_{b_{i+1}} \leq \dots \leq x_{b_n}$ . Let  
 286  $d = y - x$  be the length of the deleted sequence. Then

287  $P(D_b|M_b = 1, SV_b^{x,y}) = P_m(x_{b_1}, x_{b_2}, \dots, x_{b_i}, x_{b_{i+1}} - d, \dots, x_{b_n} - d) = P_m(n, x_{b_n} - x_{b_1} - d).$

288 To compute  $P(D_b|M_b = 2, SV_b^{x,y})$  we again need to consider all splits of the reads from barcode  $b$   
 289 into two chunks. Like before we simplify by only considering non-overlapping chunks.

$$P(D_b|M_b = 2; SV_b^{x,y}) = \sum_{k=2}^{n-1} P(D_{b_1\dots k}|M_b = 1; SV_b^{x,y})P(D_{b_{k+1}\dots n}|M_b = 1; SV_b^{x,y})$$

290 Depending on where  $x_k$  is with respect to  $x$  and  $y$  each of the probabilities above is equal to the  
291 probability under either the reference or the SV model.

292 **Inversions** Assume that the SV is an inversion between  $x$  and  $y$  ( $x < y$ ) and that  $x_{b_i} < x \leq x_{b_{i+1}}$   
293 and  $x_{b_j} < y \leq x_{b_{j+1}}$ .

294 1. If  $x_{b_1} \leq x_{b_2} \leq \dots \leq x_{b_i} < x \leq x_{b_{i+1}} \leq \dots \leq x_{b_n} < y$  (reads span  $x$  but end before  $y$ ) or

295  $x < x_{b_1} \leq \dots \leq x_{b_i} < y \leq \dots \leq x_{b_n}$  (reads start after  $x$  and span  $y$ ). In the first case,

296  $P(D_b|M_b = 1, SV_b^{x,y}) = P_m(x_{b_1}, x_{b_2}, \dots, x_{b_i}, d - x_{b_n}, d - x_{b_{n-1}}, \dots, d - x_{b_{i+1}}) =$

297  $P_m(n, x - x_{b_1} + y - x_{b_{i+1}}) = P_m(n, d - x_{b_1} - x_{b_{i+1}})$ , where  $d = x + y$ . The second case is  
298 similar.

299 2. In all other cases (reads entirely before  $x$ , reads entirely after  $y$ , reads entirely between  $x$  and  
300  $y$ , or reads spanning across  $x$  and  $y$ ),  $P(D_b|M_b = 1, SV_b^{x,y}) = P(D_b|M_b = 1; R_b)$ .

301 **Duplications** Assume that the SV is a duplication between  $x$  and  $y$  ( $x < y$ ) and that

302  $x_{b_i} < x \leq x_{b_{i+1}}$  and  $x_{b_j} < y \leq x_{b_{j+1}}$ .

303 1. If  $x < x_{b_1}$  and  $y > x_{b_n}$ , then the reads span the duplication and

304  $P(D_b|M_b = 1, SV_b^{x,y}) = P_m(n, d + y - x)$ .

305 2. If  $x < x_{b_1}$  and  $y > x_{b_n}$  (reads entirely within the duplication), then

306  $P(D_b|M_b = 1, SV_b^{x,y}) = \max\left(P_m(n, x_{b_n} - x_{b_1}), \max_j P_m(n, y - x - x_{b_{j+1}} + x_{b_j})\right)$ .

307 3. Otherwise,  $P(D_b|M_b = 1, SV_b^{x,y}) = P(D_b|M_b = 1; R_b)$ .

308 **Large-scale translocations** We only consider this case if  $x_{b_1}, \dots, x_{b_n}$  are generated from two

309 different chromosomes or  $x_{b_n} - x_{b_1} > L_{\max}$ . We can then split the reads into two groups

310  $x'_{b_1}, \dots, x'_{b'_n}, x''_{b_1}, \dots, x''_{b''_n}$  such that  $n' + n'' = n$ . Each group contains the subset of reads closer to

311  $x$  and  $y$  respectively.

312 1. If any of the two sets of reads above are empty then

313 
$$P(D_b|M_b = 1, SV_b^{x,y}) = P(D_b|M_b = 1; R_b).$$

314 2. If  $x'_{b'_n} < x$  and  $x''_{b''_1} > y$ , then  $P(D_b|M_b = 1, SV_b^{x,y}) = P_m(n, x - x'_{b'_1} + x''_{b''_n} - y)$ . All cases  
315 where all reads from the first set are on the same side of  $x$  and all reads from the second set  
316 are on the same side of  $y$  are similar.

317 3. Otherwise,  $P(D_b|M_b = 1, SV_b^{x,y}) = \epsilon$ .

318 **Unknowns** Since Long Ranger identifies event types by matching to simple models of deletions,  
319 duplications and inversions, there are additional events where Long Ranger identifies clear  
320 evidence for anomalous barcode overlap or coverage, but is unable to match the event to one of the  
321 pre-defined models, these are labeled “unknown”.

### 322 3.3 EM

323 We can use an EM approach to maximize the likelihood. This involves repeatedly conditioning on  
324 the latent variables to compute the maximum likelihood model and then getting a posterior  
325 estimate of the latent variables.

#### 326 3.3.1 M-step: Likelihood conditioning on the latent variables

327 **Homozygous reference** The likelihood of the data under the homozygous reference model is:

$$\prod_d \sum_{c=1}^2 P(D_b|M_b = c, R_b)P(M_b = c)$$

328 **Homozygous SV** The likelihood of the data under the homozygous SV model is:

$$\prod_d \sum_{c=1}^2 P(D_d|M_b = c; SV_b^{x,y})P(M_b = c)$$

329 We need to compute this for every type of SV.

330 **Heterozygous SV**

$$P(D_b; m) = \sum_{i,j \in [0,1]^2} \sum_{c=1}^2 P(D_b | H_b^{x,y} = (i,j), M_b = c; m) P(H_b^{x,y} = (i,j), M_b = c; m)$$

331 where  $m$  is the model (reference or SV).

$$P(D_b | H_b^{x,y} = (i,j), M_b = 1; SV_{i,j}^{x,y}) = P(D_b | SV_b^{x,y}, M_b = 1)$$

$$P(D_b | H_b^{x,y} \neq (i,j), M_b = 1; SV_{i,j}^{x,y}) = P(D_b | R_b, M_b = 1)$$

332 To compute  $P(D_b | H_b^{x,y} = (i,j), M_b = 2; SV_{i,j}^{x,y})$ , we start with the case where  $x$  and  $y$  are on the  
333 same phase block, so  $i$  and  $j$  must be equal.

$$P(D_b | H_b^{x,y} = (i,i), M_b = 2; SV_{i,i}^{x,y}) = \sum_{k=2}^{n-1} P(D_{b_{1\dots k}} | H_{b_{1\dots k}}^{x,y} = (i,i), M_b = 1; SV_{i,i}^{x,y}) P(D_{b_{k+1\dots n}} | H_{b_{k+1\dots n}}^{x,y} = (i,i), M_b = 1; SV_{i,i}^{x,y})$$

334 Here the sum is taken over all ways of splitting the reads from  $b, x_1, x_2, \dots, x_n$  into two (non  
335 empty) sequences  $x_1, \dots, x_k$  and  $x_{k+1}, \dots, x_n$ .  $D_{b_{1\dots k}}$  and  $D_{b_{k+1\dots n}}$  are the sets of reads resulting  
336 from such a split. Depending where  $x_k$  is with respect to  $x$   
337  $P(D_{b_{1\dots k}} | H_{b_{1\dots k}}^{x,y} = (i,i), M_b = 1; SV_{i,i}^{x,y})$  is either  $P(D_{b_{1\dots k}} | R_{b_{1\dots k}}, M_{b_{1\dots k}} = 1)$  or  
338  $P(D_{b_{1\dots k}} | SV_{b_{1\dots k}}^{x,y}, M_{b_{1\dots k}} = 1)$ . The sum is taken similarly for the likelihood of the second chunk of  
339 data.

340 If  $x$  and  $y$  are on different phase blocks, then  $i$  and  $j$  can be different. To simplify things a bit, we  
341 assume that the only valid split is the one that assigns the points closer to  $x$  to haplotype  $i$  and the

342 points closer to  $y$  to haplotype  $j$ . The computation is then similar to the case above.

### 343 3.3.2 E-step: Posterior of the latent variables

$$P(H_b^{x,y} = (i, j), M_b = c | D_b; m) \propto P(D_b | H_b^{x,y} = (i, j), M_b = c; m) P(H_b^{x,y} = (i, j), M_b = c)$$

344 All we need is a prior on the latent variables. First, we assume that

$$P(H_b^{x,y} = (i, j), M_b = c) = P(H_b^{x,y} = (i, j)) P(M_b = c)$$

345 To compute  $P(H_b^{x,y} = (i, j))$ , let  $p_b^x(0), p_b^x(1)$  be the probability that barcode  $b$  at locus  $x$  is phased  
346 on haplotype 0 or 1 respectively. We assume that these probabilities are pre-computed during SNP  
347 phasing. If  $b$  is un-phased at  $x$ , then we can set  $p_b^x(0)$  to 0.5 or to the fraction of barcodes that are  
348 phased to haplotype 0 at locus  $x$ . If  $x$  and  $y$  are in the same phase set, then

349  $P(H_b^{x,y} = (i, j)) = p_b^x(i)$  if  $i == j$  and  $P(H_b^{x,y} = (i, j)) = 0$  otherwise. If  $x$  and  $y$  are on different  
350 phase blocks, then  $P(H_b^{x,y} = (i, j)) = p_b^x(i)p_b^y(j)$ .

351 To compute  $P(M_b = c)$ , let  $p_{ov}$  be the probability of having two overlapping molecules in the same  
352 partition. The probability that the reads with barcode  $b$  came from a single molecule is the product  
353 of the probability of generating a molecule greater than the observed length and the probability  
354 that there is no molecule overlap:  $P(M_b = 1) = \sum_{\ell \geq x_{b_n} - x_{b_1}} P_L(\ell)(1 - p_{ov})$  and  
355  $P(M_b = 2) = 1 - P(M_b = 1)$ .

### 356 3.4 Computing SV phasing scores

357 We can assign a score to the haplotype assignment of the SV as follows:

$$\frac{P(D; SV_{i,j}^{x,y})}{\sum_{(i,j) \in [0,1]^2} P(D; SV_{i,j}^{x,y})}$$

### 358 **3.5 Barcode coverage**

359 The genomic extent of long input molecules is inferred by ‘linking’ successive read-pairs with the  
360 same barcode if there are separated by <60 kb. A barcode coverage track is computed by counting  
361 the number of inferred molecules that span each position in the genome.

## 362 **4 SV calling details**

### 363 **4.1 Large SVs (>30 kb)**

#### 364 **Inherency consistency analysis for NA12878**

365 To evaluate 1) calls in the svclassify set that were not called by Long Ranger and 2) PASS SV calls  
366 contained in NA12878 that were not in the svclassify set, we looked to Long Ranger analyses of the  
367 parental samples for evidence of these events as well as manual review of the barcode overlap and  
368 coverage data.

369 With regard to svclassify events not called by Long Ranger- one event (chr12:8,558,486-8,590,846)  
370 is well-supported in the Linked-Read data by barcode overlap. For this event, Long Ranger calls a  
371 10 kb small deletion with a consistent 5’ breakpoint to the svclassify event, but prematurely closes  
372 the event, missing 22 kb of the deletion. The event is called correctly in both parents. A second  
373 event (chr22:24,274,143-24,311,297) is also well-supported in all three individuals and is called by  
374 Long Ranger but is filtered out from PASS and moved to Filtered as it overlaps with a segmental  
375 duplication and is thus an error of annotation. There is no support for the last missing call  
376 (chr14:37,631,608-37,771,227). Further investigation of this call reveals that it is genotyped as  
377 homozygous reference in NA12878 in the 1000G SV set (Sudmant et al. 2015), and represents an

378 error in the svclassify set relative to GRCh37.p13.

379 The three Long Ranger-only events that did not show inheritance consistency with breakpoint  
380 consistency all overlap more complex events/regions of the genome. The first,  
381 chr1:189690000-189790000, entirely contains an event that does show inheritance and breakpoint  
382 consistency (chr1:189704517-189783347). The second, chr11:55360000-55490000, overlaps with an  
383 event annotated as UNK (chr11:55365428-55445878) that is inheritance and breakpoint consistent  
384 with NA12891 and thus represents a more complex event than just a simple deletion. The final  
385 event, chr2:242900000-243080000, overlaps with four known GRCh37 assembly issues (HG-1616,  
386 HG-1709, HG-1714, and HG-1911), and is immediately upstream of a known assembly gap. We see  
387 a drop in phased coverage on haplotype 2 in NA12878 that is not seen in either parent, and thus  
388 this likely represents a false positive call in a complex assembly region (Supplemental Table S7).

389 **Details on GeT-RM CNVPanel samples** To assess for performance of Linked-Reads on  
390 clinically-relevant variant types typically assessed by aCGH or karyotype we performed  $30\times$   
391 Linked-Read genome sequencing on a set of 23 samples with known balanced or unbalanced SVs.  
392 These 23 samples are derived from 1) the GeT-Rm CNVPanel, a collection of unbalanced events  
393 including large deletions, duplications, inversions, balanced translocations and unbalanced  
394 translocations designed to assess performance of clinical aCGH or 2) the Coriell general Cell  
395 Repository (balanced events). These cell lines have multiple, orthogonal assays confirming the  
396 presence of their described structural variants. We assessed the performance of each of the  
397 barcode-based detection methods (barcode coverage and barcode overlap) individually in our panel  
398 of samples covering a wide range of structural variant types (Supplemental Table S1). Together, the  
399 barcode overlap and barcode coverage methods detected 27 of the 29 structural variants, correctly  
400 characterizing 22 of the 23 samples tested. The barcode coverage method called all copy number  
401 variants except one; a chr7 deletion filtered to a secondary list of candidate structural variants  
402 because the breakpoints overlap a segmental duplication, a feature known to complicate genome  
403 analysis. However, this event was detected with the barcode overlap method, demonstrating the  
404 advantage of utilizing multiple detection methods for variant calling. The barcode overlap method

405 is unable to detect terminal events because it requires the ability to examine barcode sharing  
406 patterns on both sides of a breakpoint. This algorithm was able to call all non-terminal copy  
407 number events and all balanced events, except one. This undetected event is a balanced  
408 translocation with a breakpoint in a heterochromatic region of chromosome 16. This region is  
409 represented by Ns in the reference assembly and will be invisible to any sequence-based method  
410 relying on the reference genome (Schneider et al. 2017).

411 Downsampling of sequence data was performed in silico to determine the minimum sequence  
412 depth required to detect each CNV. Due to the evenness provided by barcode coverage, the  
413 deletion and duplication signals for these two samples are detectable even with coverage as low as  
414 5 Gb ( $\sim 1 \times$  genomic read coverage) (Supplemental Figure S8 C and D). The barcode overlap method  
415 was less robust with reduced sequencing and did not call CNVs with less than 50 Gb of sequence.  
416 This result is expected, given that the algorithm was designed for use with full-depth data.  
417 However, there was an observable signal in the barcode overlap data with as little as 5-10 Gb for  
418 many of the samples, indicating that the algorithm is likely extensible to lower depth data.  
419 (Supplemental Figure S9).

## 420 **5 TruSeq PCR-free library preparation**

421 350-800 ng of genomic DNA was sheared to a size of  $\sim 385$  bp using a CovarisM220 Focused  
422 Ultrasonicator using the following shearing parameters: Duty factor = 20%, cycles per burst = 200,  
423 time = 90 seconds, Peak power 50. Fragmented DNA was then cleaned up with  $0.8 \times$  SPRI beads  
424 and left bound to the beads. Then, using the KAPA Library Preparation Kit reagents (KAPA  
425 Biosystems, Catalog # KK8223), DNA fragments bound to the SPRI beads were subjected to end  
426 repair, A-base tailing and Illumina 'PCR-free' TruSeq adapter ligation (1.5 M final concentration of  
427 adapter was used). Following adapter ligation, two consecutive SPRI cleanup steps ( $1.0 \times$  and  $0.7 \times$ )  
428 were performed to remove adapter dimers and library fragments below  $\sim 150$  bp in size. No library  
429 PCR amplification enrichment was performed. Libraries were then eluted off the SPRI beads in 25

<sup>430</sup>  $\mu$ l elution buffer and quantified with quantitative PCR using KAPA Library Quant kit (KAPA  
<sup>431</sup> Biosystems, Catalog # KK4824) and an Agilent Bioanalyzer High Sensitivity Chip (Agilent  
<sup>432</sup> Technologies) following the manufacturer's recommendations.

## References

- 433
- 434 Bansal V, Bafna V. 2008. HapCUT: An efficient and accurate algorithm for the haplotype assembly  
435 problem. *Bioinformatics* **24**: i153–i159.
- 436 Bansal V, Halpern AL, Axelrod N, Bafna V. 2008. An MCMC algorithm for haplotype assembly  
437 from whole-genome sequence data. *Genome Res* **18**: 1336–1346.
- 438 Bishara A, Liu Y, Weng Z, Kashef-Haghighi D, Newburger DE, West R, Sidow A, Batzoglou S. 2015.  
439 Read clouds uncover variation in complex regions of the human genome. *Genome Res* **25**:  
440 1570–1580.
- 441 Li H. 2013. Aligning sequence reads , clone sequences and assembly contigs with BWA-MEM.  
442 *ArXiv* **00**: 1–2.
- 443 Quinlan AR, Hall IM. 2010. BEDTools: A flexible suite of utilities for comparing genomic features.  
444 *Bioinformatics* **26**: 841–842.
- 445 Schneider VA, Graves-Lindsay T, Howe K, Bouk N, Chen H-C, Kitts PA, Murphy TD, Pruitt KD,  
446 Thibaud-Nissen F, Albracht D, et al. 2017. Evaluation of grch38 and de novo haploid genome  
447 assemblies demonstrates the enduring quality of the reference assembly. *Genome Res* **27**: 849–864.
- 448 Sudmant PH, Kitzman JO, Antonacci F, Alkan C, Malig M, Tsalenko A, Sampas N, Bruhn L,  
449 Shendure J, 1000 Genomes Project, et al. 2010. Diversity of human copy number variation and  
450 multicopy genes. *Science* **330**: 641–646.