

<b>Manuscript Number:</b>	GIGA-D-19-00092R2	
<b>Full Title:</b>	Assessment of human diploid genome assembly with 10x Linked-Reads data	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	National Institute of Standards and Technology (na)	Not applicable
<b>Abstract:</b>	<p>Background: Producing cost-effective haplotype-resolved personal genomes remains challenging. 10x Linked-Read sequencing, with its high base quality and long-range information, has been demonstrated to facilitate de novo assembly of human genomes and variant detection. In this study, we investigate in depth how the parameter space of 10x library preparation and sequencing affects assembly quality, on the basis of both simulated and real libraries. Findings: We prepared and sequenced eight 10x libraries with a diverse set of parameters from standard cell lines NA12878 and NA24385 and performed whole genome assembly on the data. We also developed the simulator LRTK-SIM to follow the workflow of 10x data generation and produce realistic simulated Linked-Read data sets. We found that assembly quality could be improved by increasing the total sequencing coverage (C) and keeping physical coverage of DNA fragments (CF) or read coverage per fragment (CR) within broad ranges. The optimal physical coverage was between 332X and 823X and assembly quality worsened if it increased to greater than 1,000X for a given C. Long DNA fragments could significantly extend phase blocks, but decreased contig contiguity. The optimal length-weighted fragment length (<math>W\mu_{FL}</math>) was around 50 – 150kb. When broadly optimal parameters were used for library preparation and sequencing, ca. 80% of the genome was assembled in a diploid state. Conclusion: The Linked-Read libraries we generated and the parameter space we identified provide theoretical considerations and practical guidelines for personal genome assemblies based on 10x Linked-Read sequencing. Keywords: 10x Linked-Read sequencing, de novo assembly, diploid human genome, library preparation</p>	
<b>Corresponding Author:</b>	arend sidow  UNITED STATES	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>		
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Lu Zhang	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	Lu Zhang	
	Xin Zhou	
	Ziming Weng	
	arend sidow	
<b>Order of Authors Secondary Information:</b>		
<b>Response to Reviewers:</b>	<p>Reviewer reports: Reviewer #2: The authors improved the manuscript substantially and implemented many of the suggested changes. I wonder, however, whether there was a mixup of document versions because not all changes described in the response are reflected in the manuscript (including trivial ones like fixing the "_Alignment", now in line 283; also Luo et al. is still not cited). Maybe the authors can double check that they indeed</p>	

	<p>uploaded the latest version?</p> <p>Thank you for pointing out this oversight. We have double-checked everything and added a section to the appropriate place in the Methods where we explain the differences between Luo et al's method and ours. (The added text is highlighted in red.)</p> <p>Beyond that, the only concern left for me is the poor concordance of small variant calls. For the Illumina and 10x calls, my guess is that they went into the evaluation completely unfiltered, where FreeBayes (and the LongRanger pipeline which is based on FreeBayes) usually attain an acceptable precision only when the calls are filtered (e.g. for QUAL&gt;=10). Much more concerning is the observation that between a quarter and half of all calls are missed by the assembly strategy. How did the authors call variants from the assemblies? Given that the GIAB benchmark regions are (comparatively) easy genomic regions, I think that the authors should offer an explanation for the poor recall.</p> <p>We did not use any threshold to filter out low-quality variants from FreeBayes. To generate assembly-based calls, we aligned the two haploid contigs from Supernova to the reference genome (Mimimap2) independently and compared the two alleles of the corresponding coordinates (Paftools, mapQ&gt;20).</p> <p>For small SNV calls, we agree using Freebayes is a better choice since mapping-based algorithms have good base accuracy and assembly-based algorithms may lose sensitivity. The significant false negative rates of assembly-based calls likely come from two issues:</p> <ol style="list-style-type: none"> <li>1.Supernova cannot guarantee to generate diploid contigs (megabubbles) even for the “easy regions” from GIAB, because the diploid contigs would be influenced by SV also. As a result, in those regions we lose a large fraction of heterozygous variants.</li> <li>2.The single base variants in the de bruijn graph are represented as small bubbles, which would be flattened due to various reasons. The k-mer coverage is one of the critical thresholds, but the length of k-mer is much shorter than reads and the sequencing qualities are not taken into consideration. These may lead to miscount the coverage of variant alleles in the bubbles.</li> </ol> <p>We have added an explanatory sentence in the last section of the Results.</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
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
<b>Resources</b>	Yes

<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>	
<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>	<p>Yes</p>

[Click here to view linked References](#)

# 1 **Assessment of human diploid genome assembly with 10x**

## 2 **Linked-Reads data**

3

4 **Lu Zhang<sup>1,2,3,\*</sup>, Xin Zhou<sup>3,\*</sup>, Ziming Weng<sup>2</sup>, Arend Sidow<sup>2,4,†</sup>**

5 <sup>1</sup>Department of Computer Science, Hong Kong Baptist University

6 <sup>2</sup>Department of Pathology, Stanford University

7 <sup>3</sup>Department of Computer Science, Stanford University

8 <sup>4</sup>Department of Genetics, Stanford University

9 \*These authors contributed equally to this work. †Correspondence and requests for materials should be  
10 addressed to Arend Sidow (email: arend@stanford.edu)

11

12 ORCIDs:

13 Arend Sidow, 0000-0002-8287-331X;

14 Xin Zhou, 0000-0003-4015-4787;

15 Lu Zhang, 0000-0002-8157-9856

## 16 **Abstract**

17 **Background:** Producing cost-effective haplotype-resolved personal genomes remains  
18 challenging. 10x Linked-Read sequencing, with its high base quality and long-range information,  
19 has been demonstrated to facilitate *de novo* assembly of human genomes and variant detection.  
20 In this study, we investigate in depth how the parameter space of 10x library preparation and  
21 sequencing affects assembly quality, on the basis of both simulated and real libraries.

22 **Results:** We prepared and sequenced eight 10x libraries with a diverse set of parameters from  
23 standard cell lines NA12878 and NA24385 and performed whole genome assembly on the data.  
24 We also developed the simulator LRTK-SIM to follow the workflow of 10x data generation and  
25 produce realistic simulated Linked-Read data sets. We found that assembly quality could be  
26 improved by increasing the total sequencing coverage ( $C$ ) and keeping physical coverage of DNA  
27 fragments ( $C_F$ ) or read coverage per fragment ( $C_R$ ) within broad ranges. The optimal physical  
28 coverage was between 332X and 823X and assembly quality worsened if it increased to greater  
29 than 1,000X for a given  $C$ . Long DNA fragments could significantly extend phase blocks, but  
30 decreased contig contiguity. The optimal length-weighted fragment length ( $W_{\mu_{FL}}$ ) was around 50  
31 – 150kb. When broadly optimal parameters were used for library preparation and sequencing, ca.  
32 80% of the genome was assembled in a diploid state.

33 **Conclusions:** The Linked-Read libraries we generated and the parameter space we identified  
34 provide theoretical considerations and practical guidelines for personal genome assemblies  
35 based on 10x Linked-Read sequencing.

36 **Keywords:** 10x Linked-Read sequencing, *de novo* assembly, diploid human genome, library  
37 preparation

## 38 **Background**

39 The human genome holds the key for understanding the genetic basis of human evolution,  
40 hereditary illnesses and many phenotypes. Whole-genome reconstruction and variant discovery,  
41 accomplished by analysis of data from whole-genome sequencing experiments, are foundational  
42 for the study of human genomic variation and analysis of genotype-phenotype relationships. Over  
43 the past decades, cost-effective whole-genome sequencing has been revolutionized by short-  
44 fragment approaches, the most widespread of which have been the consistently improving  
45 generations of the original Solexa technology [1, 2], now referred to as Illumina sequencing.  
46 Illumina's strengths and weaknesses are inherent in the sample preparation and sequencing  
47 chemistry. Illumina generates short paired reads (2x150 base pairs for the highest-throughput  
48 platforms) from short fragments (usually 400-500 base pairs) [3]. Because many clonally amplified  
49 molecules generate a robust signal during the sequencing reaction, Illumina's average per-base  
50 error rates are very low.

51

52 The lack of long-range contiguity between end-sequenced short fragments limits their application  
53 for reconstructing personal genomes. Long-range contiguity is important for phasing variants and  
54 dealing with genomic complex regions. For haplotyping, variants can be phased by population-  
55 based methods [4, 5] or family-based recombination inference [6, 7]. However, such approaches  
56 are only feasible for common variants in single individuals or when a trio or larger pedigree is  
57 sequenced. Furthermore, highly polymorphic regions such as the HLA in which the reference  
58 sequence does not adequately capture the diversity segregating in the population are refractory  
59 to mapping-based approaches and require *de novo* assembly to reconstruct [8]. Short-read/short-  
60 fragment data are challenged by interspersed repetitive sequences from mobile elements and by  
61 segmental duplications, and only support highly fragmented genome reconstruction [9, 10].

62

63 In principle, many of these challenges can be overcome by long-read/long-fragment sequencing  
64 [11, 12]. Assembly of Pacific Biosciences (PacBio) or Oxford Nanopore (ONT) data can yield  
65 impressive contiguity of contigs and scaffolds. In one study [13], scaffold N50 reached 31.1Mb by  
66 hierarchically integrating PacBio long reads and BioNano for a hybrid assembly, which also  
67 uncovered novel tandem repeats and replicated the structural variants (SVs) that were newly  
68 included in the updated hg38 human reference sequence. Another study [14] produced human  
69 genome assemblies with ONT data, in which a contig N50 ~3Mb was achieved, and long contigs  
70 covered all class I HLA regions. A recent whole genome assembly of NA24385 [15] with high  
71 quality PacBio CCS reads generated contigs with an N50 of 15Mb. However, long-fragment  
72 sequencing suffers from extremely high cost (in the case of PacBio CCS), or low base quality (in  
73 the case of single-pass reads of either technology), hampering its usefulness for personal genome  
74 assembly.

75  
76 Hierarchical assembly pipelines in which multiple data types are used as another approach for  
77 genome assembly [16]. For example, in the reconstruction of an Asian personal genome, fosmid  
78 clone pools and Illumina data were merged, but because fosmid libraries are highly labor intensive  
79 to generate and sequence, this approach is not generalizable to personal genomes. The "Long  
80 Fragment Read" (LFR) approach [17], where a long fragment is sequenced at high depth via  
81 single-molecule fragmented amplification, reported promising personal genome assembly and  
82 variant phasing by attaching a barcode to the short reads derived from the same long fragment.  
83 However, because LFR is implemented in a 384 well plate, many long fragments would be  
84 labelled by the same barcodes, making it difficult for binning short-reads, and the great  
85 sequencing depth required rendered LFR not cost-effective.

86  
87 An alternative approach is offered by the 10x Genomics Chromium system, which distributes the  
88 DNA preparation into millions of partitions where partition-specific barcode sequences are

89 attached to short amplification products that are templated off the input fragments. Because of  
90 the limited reaction efficiency in each partition, the sequencing depth for each fragment is too  
91 shallow to reconstruct the original long-fragment, distinguishing this approach from LFR [18].  
92 However, to compensate for the low read coverage of each fragment, each genomic region is  
93 covered by hundreds of DNA fragments, giving overall sequence coverage that is in a range  
94 comparable to standard Illumina short-fragment sequencing while providing very high physical  
95 coverage. Novel computational approaches leveraging the special characteristics of 10x  
96 Genomics data have already generated significant advances in power and accuracy of  
97 haplotyping [19, 20], cancer genome reconstruction [21, 22], metagenomic assemblies [23], and  
98 *de novo* assembly of human and other genomes [24-26], compared to standard Illumina short-  
99 fragment sequencing. While the uniformity of sequence coverage is not as good as with PCR-  
100 free Illumina libraries, 10x Linked-Read sequencing is a promising technology that combines low  
101 per-base error and good small-variant discovery with long-range information for much improved  
102 SV detection in mapping-based approaches [22, 27], and the possibility of long-range contiguity  
103 in *de novo* assembly [24, 26, 28].

104

105 Practical advantages of the technology include the low DNA input mass requirement (1ng per  
106 library, or approximately 300 haploid human genome equivalents). Real input quantities can vary,  
107 along with other factors, to influence an interconnected array of parameters that are relevant to  
108 genome assembly and reconstruction. The parameters over which the experimenter has influence  
109 are (**Figure 1**): i).  $C_R$ : average **C**overage of short **R**eads per fragment; ii).  $C_F$ : average physical  
110 **C**overage of the genome by long DNA **F**ragments; iii).  $N_{FP}$ : **N**umber of **F**ragments per **P**artition;  
111 iv). Fragment length distribution, several parameters of which are used, specifically  $\mu_{FL}$ : Average  
112 Unweighted DNA **F**ragment **L**ength and  $W\mu_{FL}$ : Length-**W**eighted average of DNA **F**ragment  
113 **L**ength. Note that several parameters depend on each other. For example, a greater amount of  
114 input DNA will increase  $N_{FP}$ ; shorter fragments increase  $N_{FP}$  at the same DNA input amount



115 compared to longer fragments; less input DNA will (within practical constraints) increase  $C_R$  and  
116 decrease  $C_F$ ; and their absolute values are set by how much total sequence coverage is  
117 generated because  $C_R \times C_F = C$ .

118

119 Our goal in this study was to experimentally explore the 10x parameter space and evaluate the  
120 quality of *de novo* diploid assembly as a function of the parameter values. For example, we set  
121 out to ask whether longer input fragments produce better assemblies, or what the effect of  
122 sequencing vs. physical coverage is on contiguity of assembly. In order to constrain the parameter  
123 space, we first performed computer simulations with reasonably realistic synthetic data. The  
124 simulation results suggested certain parameter combinations that we then approximated in the  
125 generation of real, high-depth, sequence data on two human reference genome cell lines,  
126 NA12878 and NA24385. These simulated and real data sets were then used to produce *de novo*  
127 assemblies, with an emphasis on the performance of 10x's Supernova2 [24]. We finally assessed  
128 the quality of the assemblies using standard metrics of contiguity and accuracy, facilitated by the  
129 existence of a gold standard (in the case of simulations) and comparisons to the reference  
130 genome (in the case of real data).

131

## 132 **Methods**

133

### 134 **Library preparation, physical parameters and sequencing coverage**

135 We made six DNA preparations that varied in fragment size distribution and amount of input DNA,  
136 three each from NA12878 (Coriell Cat# GM12878, RRID:CVCL\_7526) and NA24385 (Coriell Cat#  
137 GM24385, RRID:CVCL\_1C78). From these, we prepared eight libraries, five from NA12878 and  
138 three from NA24385 (**Table S1**). To generate libraries  $L_{1L}$ ,  $L_{1M}$  and  $L_{1H}$  (the subscripts  $L$ ,  $M$  and  
139  $H$  represent low, medium and high  $C_F$ , respectively), genomic DNA was extracted from ca. 1

140 million cultured NA12878 cells using the Gentra Puregene Blood Kit following manufacturer's  
141 instructions (Qiagen, Cat. No 158467). The GEMs were divided into 3 tubes with 5%, 20%, and  
142 75% to generate libraries  $L_{1L}$ ,  $L_{1M}$  and  $L_{1H}$ , respectively (**Figure S1-S3**). For the other libraries,  
143 to generate longer DNA fragments ( $W_{\mu FL}=150\text{kb}$  and longer, **Figure S4-S8**), a modified protocol  
144 was applied. Two-hundred thousand NA12878 or NA24385 cells of fresh culture were added to  
145 1mL cold 1x PBS in a 1.5 ml tube and pelleted for 5 minutes at 300g. The cell pellets were  
146 completely resuspended in the residual supernatant by vortexing and then lysed by adding 200ul  
147 Cell Lysis Solution and 1ul of RNaseA Solution (Qiagen, Cat. No 158467), mixing by gentle  
148 inversion, and incubating at 37°C for 15-30 minutes. This cell lysis solution is used immediately  
149 as input for the 10x Chromium preparation (Chromium™ Genome Library & Gel Bead Kit v2,  
150 PN-120258; Chromium™ i7 Multiplex Kit, PN-120262). Fragment size of the input DNA can be  
151 controlled by gentle handling during lysis and DNA preparation for Chromium. The amount of  
152 input DNA (between 1.25 and 4 ng) was varied to achieve a wide range of physical coverage  
153 ( $C_F$ ).The Chromium Controller was operated and the GEM preparation was performed as  
154 instructed by the manufacturer. Individual libraries were then constructed by end repairing, A-  
155 tailing, adapter ligation and PCR amplification. All libraries were sequenced with three lanes of  
156 paired-end 150bp runs on the Illumina HiSeqX to obtain very high coverage ( $C=94\text{x}-192\text{x}$ ), though  
157 the two with the fewest number of gel beads ( $L_{1L}$  and  $L_{1M}$ ) exhibited high PCR duplication rates  
158 because of the reduced complexity of the libraries (**Table S1**).

159

### 160 **Linked-Reads subsampling**

161 The high sequencing coverage in the libraries allowed subsampling to facilitate the matching of  
162 parameters among the different libraries, for purposes of comparability; these subsampled  
163 Linked-Read sets are denoted  $R_{id}$  (**Figure 1**). We aligned the 10x Linked-Reads to human  
164 reference genome (hg38, GRCh38 Reference 2.1.0 from 10x website) followed by removing PCR

165 duplication by barcode-aware analysis in Long Ranger[21]. Original input DNA fragments were  
166 inferred by collecting the read-pairs with the same barcode that were aligned in proximity to each  
167 other. A fragment was terminated if the distance between two consecutive reads with the identical  
168 barcode larger than 50kb. Fragments were required to have at least two read pairs with the same  
169 barcode and a length of at least 2 kb. Partitions with fewer than three fragments were removed.  
170 We subsampled short-reads for each fragment to satisfy the expected  $C_R$ .

171

### 172 **Generating 10x simulated libraries by LRTK-SIM**

173 To compare the observations from real data with a known truth set, we developed LRTK-SIM, a  
174 simulator that follows the workflow of the 10x Chromium system and generates synthetic Linked-  
175 Reads like those produced by an Illumina HiSeqX machine (**Supplementary Information** and  
176 **Figure S9**). Based on the parameters commonly employed by 10x Genomics Linked-Read  
177 sequencing and the characteristics of our libraries, LRTK-SIM generated simulated datasets from  
178 the human reference (hg38), explicitly modeling the five key steps in real data generation.  
179 Parameters in parentheses are from the standard 10x Genomics protocol: 1. Shearing genomic  
180 DNA into long fragments ( $W_{\mu_{FL}}$  from 50kb to 100kb); 2. Loading DNA to the 10x Chromium  
181 instrument (~1.25ng DNA); 3. Allocating DNA fragments into partitions which are attached the  
182 unique barcodes (~10 fragments per partition); 4. Generating short fragments; 5. Generating  
183 Illumina paired-end short reads (800M~1200M reads). LRTK-SIM first generated a diploid  
184 reference genome as a template by duplicating the human reference genome (hg38) into two  
185 haplotypes and inserting single nucleotide variants (SNVs) from high-confidence regions in GIAB  
186 of [NA12878](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/HG001_GRCh38_GIA) ([ftp://ftp-](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/HG001_GRCh38_GIA)  
187 [trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878\\_HG001/latest/GRCh38/HG001\\_GRCh38\\_GIA](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/HG001_GRCh38_GIA)  
188 [B\\_highconf\\_CG-IIIIFB-IIIIGATKHC-Ion-10X-SOLID\\_CHROM1-](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/HG001_GRCh38_GIA)  
189 [X\\_v.3.3.2\\_highconf\\_nosomaticdel\\_noCENorHET7.bed](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/HG001_GRCh38_GIA)); For low-confidence regions we  
190 randomly simulated 1 SNV per 1 kb. The ratio was 2:1 for heterozygous and homozygous SNVs.

191 From this diploid reference genome, LRTK-SIM generated long DNA fragments by randomly  
192 shearing each haplotype with multiple copies into pieces whose lengths were sampled from an  
193 exponential distribution with mean of  $\mu_{FL}$ . These fragments were then allocated to pseudo-  
194 partitions, and all the fragments within each partition were assigned the same barcode. The  
195 number of fragments for each partition was randomly picked from a Poisson distribution with mean  
196 of  $N_{FP}$ . Finally, paired-end short reads were generated according to  $C_R$  and replaced the first 16bp  
197 of the reads from forward strand to the assigned barcodes followed by 7 Ns. More information  
198 about implementation can be found in **Supplementary Information**. From that diploid genome,  
199 Linked-Read datasets were generated that varied in  $C_R$ ,  $C_F$  and  $\mu_{FL}$  ( $W\mu_{FL}$ ) (**Table S2-S3**).  
200 Varying  $N_{FP}$  was only done for chromosome 19 because of the infeasibility of running Supernova2  
201 on whole genome assemblies with large  $N_{FP}$ ; within practically reasonable values,  $N_{FP}$  does not  
202 appear to influence assembly quality (**Figure S10**). In total, we generated 17 simulated Linked-  
203 Read datasets to explore the overall parameter space (**Table S2-S3**) and 11 to match the  
204 parameters of the abovementioned real libraries (**Figure 1**).

205

206 LRTK-SIM provides more flexible simulation parameters than another method for simulating  
207 linked-read data, LRSIM [29]. It explicitly allows users to input  $C_F$ ,  $C_R$ ,  $W\mu_{FL}$  and  $\mu_{FL}$ , which have  
208 strong connections with library preparation and Illumina sequencing, whereas LRSIM only lets the  
209 user set the total number of reads. For example,  $C_F$  is driven by input DNA amount, and  $\mu_{FL}$  by  
210 DNA preparation and potential size selection. Also, LRSIM requires many third party packages  
211 and software to be installed first, such as Inline::C perl library and DWGSIM [30]. By contrast,  
212 LRTK-SIM was written in Python and no third-party software is required to run it. LRTK-SIM can  
213 simulate multiple libraries with a variety of parameters simultaneously, and users can compare  
214 the performance of different parameters in one run.

215

## 216 **Human genome diploid assembly and evaluation**

217 The scaffolds were generated by the “pseudohap2” output of Supernova2, which explicitly  
218 generated two haploid scaffolds, simultaneously. Contigs were generated by breaking the  
219 scaffolds if at least 10 consecutive ‘N’s appeared, per definition by Supernova2. For the  
220 simulations of human chromosome 19, we used the scaffolds from the “megabubbles” output.  
221 Contig and scaffold N50 and NA50 were used to evaluate assembly quality. Contigs longer than  
222 500bp were aligned to hg38 by Minimap2[31]. We calculated contig NA50 on the basis of contig  
223 misassemblies reported by QUAST-LG [32]. For scaffolds (longer than 1kb), we calculated the  
224 NA50 following Assemblathon 1's procedure [33] (**Supplementary Information**).

225

## 226 **Genomic variant calls from diploid assembly**

227 We compare SNVs and SVs from the diploid regions of our assemblies with the ones from  
228 standard Illumina data and reference-based processing of our 10x data. The standard Illumina  
229 data were downloaded from Genome in a Bottle (GIAB) [34] and analyzed with SVABA [35] to  
230 generate SV calls, and with BWA (BWA, RRID:SCR\_010910) [36] and FreeBayes (FreeBayes,  
231 RRID:SCR\_010761) [37] to generate SNV calls. Long ranger  
232 (<https://support.10xgenomics.com/genome-exome/software/pipelines/latest/what-is-long-ranger>)  
233 was used to generate SNV and SV (only deletions) calls for 10x reference-based analysis. We  
234 note that R<sub>9</sub> failed to be analyzed by Long Ranger due to its extremely large C<sub>F</sub>. For SNVs, we  
235 compared the calls from three strategies using the benchmark of NA12878 ([ftp://ftp-  
236 trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878\\_HG001/latest/GRCh38/](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/)) and NA24385  
237 ([ftp://ftp-  
238 trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002\\_NA24385\\_son/latest/GRCh38/](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002_NA24385_son/latest/GRCh38/)).  
239 For SVs, we compared three linked-read sets (R<sub>9</sub>, R<sub>10</sub>, R<sub>11</sub>) from HG002 with the Tier 1 SV  
240 benchmark from GIAB [38] and used VaPoR [39] to validate our SV calls based on PacBio CCS

241 reads from NA24385 [40]. We compared SNV and SV calls among the different approaches using  
242 vcfEval [41] and truvari [38], respectively.

243

## 244 **Results**

245

### 246 **Performance of diploid assembly: influence of total coverage**

247 Diploid assembly by Linked-Reads requires sufficient total read coverage ( $C=C_R \times C_F$ ) to generate  
248 long contigs and scaffolds. In this experiment, to explore the roles of both physical coverage ( $C_F$ )  
249 and per-fragment read coverage ( $C_R$ ), we first generated eight simulated libraries whose total  
250 coverage  $C$  ranged from 16x to 78x: four with  $C_R$  fixed and increasing  $C_F$  and four with fixed  $C_F$ ,  
251 and increasing  $C_R$  (**Table S2**). Contig and scaffold N50s increased along with increasing either  
252  $C_F$  or  $C_R$  (**Figure 2A** and **2B**). To investigate whether the trend was also present in the real  
253 datasets, we analyzed six real libraries (three by varying  $C_F$ , and the other three by varying  $C_R$ ;  
254 **Figure 1**): as  $C$  increased, we varied  $C_F$  and  $C_R$  independently by fixing the other parameter.  
255 Contig and scaffold N50s also increased in these simulation (**Figure 2C** and **2D**) and real linked-  
256 read sets (**Figure 2E** and **2F**) as a function of total coverage  $C$ . Contig lengths did increase a little  
257 (621.4kb to 758.1kb for simulation; 110.7kb to 119.6kb for real data) when  $C$  was increased  
258 beyond 56X. Accuracy, which we define as the ratio between NA50 (N50 after breaking contigs  
259 or scaffolds at assembly errors) and N50 (**Figure 2C** and **2E**), changed 18% for simulation and  
260 7% for real data (587.5kb to 713.3kb for simulation; 97.1kb to 104.5kb for real data). For scaffolds  
261 in the real data sets, when  $C$  increased from 48X ( $R_3$ ) to 67X ( $R_4$ ), both scaffold N50 and NA50  
262 were significantly improved (N50: 13.4Mb to 30.6Mb; NA50: 6.3Mb to 12.0Mb), but the accuracy  
263 dropped slightly from 46.6% to 39.1%, which indicated that scaffold accuracy may be refractory  
264 to extremely high  $C$  (**Figure 2F**). These results indicated that assembly length and accuracy were

265 comparable over a broad range of  $C_F$  and  $C_R$  at constant  $C$ , which implied that assembly quality  
266 was mainly determined by  $C$ .

267

### 268 **Performance of diploid assembly: influence of fragment length and physical** 269 **coverage**

270 To investigate if input weighted fragment length (as measured by  $W\mu_{FL}$ ) influenced assembly  
271 quality, we generated four simulated libraries (**Table S3**) with fixed  $C_F$  and  $C_R$  and a range of  
272 fragment lengths (**Figure 3A**). Contig length decreased with increasing fragment length, a trend  
273 that was also seen in six real libraries (**Figure 3B**;  $C=56X$ ;  $R_6$  to  $R_{11}$  in **Figure 1**). We then  
274 simulated another six libraries with the same parameters as the real ones to explore the effects  
275 of physical coverage at constant  $C=56x$  (**Figure 3C**). Contig lengths decreased as a function of  
276 increasing physical coverage, a trend that is somewhat less clear in real data possibly due to  
277 confounding other parameters such as fragment length (**Figure 3D**). The two linked-read sets  
278 with the worst contig qualities in NA12878 ( $R_7$ ) and NA24385 ( $R_{10}$ ) also showed a significant  
279 increase of the number of breakpoints (**Table S4**)

280

### 281 **Performance of diploid assembly: nature of the source genome**

282 Assembly errors may occur because of heterozygosity, repetitive sequences, or sequencing error.  
283 To illuminate possible sources of assembly error, we performed simulations by generating 10x-  
284 like Linked-Reads as above from human chromosome 19, and then quantified assembly error  
285 against these synthetic gold standards. Removal of interspersed repeat sequences from the  
286 source genome resulted in better contigs with no loss of accuracy in experiments by varying  $C_F$ ,  
287  $C_R$  and  $\mu_{FL}$  (**Figure 4A, 4C and 4E**) and better scaffolds only if  $C_R$  was above 1X (**Figure 4D**).  
288 Removal of variation had little effect on contigs and only gave rise to longer scaffolds if  $C_R$  was

289 above 0.8X (**Figure S11**), which is difficult to achieve with real libraries. Finally, a 1% uniform  
290 sequencing error had no discernible effect (**Figure S12**).

291

### 292 **Performance of diploid assembly: fraction of genome in diploid state**

293 While contiguity is an important parameter for any whole genome assembly, evaluation of diploid  
294 assemblies necessitates estimating the fraction of the genome in which the assembly recovered  
295 the diploid state. To this end, we divided the contigs generated by Supernova2 into “diploid  
296 contigs”, which were extracted from its megabubble structures, and “haploid contigs” from non-  
297 megabubble structures. Pairs of scaffolds were extracted as the two haplotypes from  
298 megabubble structures if they shared the same start and end nodes in the assembly graph.  
299 Diploid contigs were generated by breaking the candidate scaffolds at the sequences with least  
300 10 consecutive ‘N’s and were aligned to human reference genome (hg38) by Minimap2. The  
301 genome was split into 500bp windows and diploid regions were defined as the maximum extent  
302 of successive windows covered by two contigs, each from one haplotype. Alignment against the  
303 human reference genome revealed the overall genome coverages of the six assemblies to be  
304 around 91%. For most assemblies, 70%-80% of the genome was covered by two homologous  
305 contigs (**Table 1**), with  $R_6$  only reaching 58.9%, probably due to the short fragments of the DNA  
306 preparation ( $\mu_{FL}=24\text{kb}$ ). We also analyzed another seven assemblies produced by 10x Genomics,  
307 all of which had diploid fractions of about 80% as well (**Table S5**). In the male NA24385, non-  
308 pseudoautosomal regions of the X chromosome are hemizygous and should therefore be  
309 recovered as haploid regions. Between 79.9% and 87.6% of these regions were covered by one  
310 contig exactly depending on the assembled library. Library construction parameters other than  
311 fragment length appeared to have had little impact on the proportion of diploid regions (**Tables 1**  
312 and **Table S5**).

313



314 Overlapping the diploid regions from the assemblies of the same individual revealed that 50.24%  
315 and 67.27% of the genome for NA12878 and NA24385 (**Figure S13**), respectively, were diploid  
316 in all the three assemblies. NA12878 was lower because of the low percentage of diploid regions  
317 in assembly  $R_6$  (**Table 1**). The overlaps were significantly greater than expected by chance  
318 (NA12878: 33.3%, p-value=0.0049; NA24385: 45.4%, p-value=0.0029. Chi square test). These  
319 observations were consistent with heterozygous variants being enriched in certain genomic  
320 segments, in which two haplotypes were more easily differentiated by Supernova2. Phase block  
321 lengths were mainly determined by total coverage  $C$  and increased in real data with increasing  
322 fragment length (**Figure S14, Table S6**).

323

### 324 **Performance of diploid assembly: quality of variant calls**

325 The ultimate goal of human genome assembly is to accurately identify genomic variants. We  
326 therefore compared the SNVs and SVs from our assemblies with the calls from referenced-based  
327 processing of standard Illumina and 10x data, and benchmarked them using gold standard from  
328 GIAB [38, 42] and PacBio CCS reads [40]. Accuracy of SNV calls from reference-based  
329 processing of standard Illumina and 10x data were comparable, but both were better than  
330 assembly-based calls (**Table S7 and S8**). The likely reason for the relatively poor performance of  
331 assembly-based SNV calls is that the assemblies contain only about 80% of the genome in a  
332 diploid state. For SVs, our assemblies generated many calls that were missed by the reference-  
333 based strategy (**Table S9-S12**) and even by the Tier 1 benchmark of GIAB (**Table S13**), and half  
334 of the deletions and a majority of insertions could be validated by PacBio CCS reads (**Table S14**).

335

### 336 **Discussion**

337 In this study, we investigated human diploid assembly using 10x Linked-Read sequencing data  
338 on both simulated and real libraries. We developed the simulator LRTK-SIM to examine the likely

339 impact of parameters in diploid assembly and compared results from simulated reads to those  
340 from real libraries. We thus determined the impact of key parameters ( $C_R$ ,  $C_F$ ,  $N_{F/P}$  and  $\mu_{FL}/N\mu_{FL}$ )  
341 with respect to assembly continuity and accuracy. Our study provides a general strategy to  
342 evaluate assemblies of 10x data and may have implications for the evaluation of other barcode-  
343 based sequencing technologies such as CPTv2-seq [43] or stLRF [44] in the future.

344

### 345 **10x Practicalities**

346 For standard Illumina sequencing, library complexity is usually sufficient to generate tremendous  
347 numbers of reads from unique templates and read coverage can be increased simply by  
348 sequencing more. However, the 10x Chromium system performs amplification in each partition,  
349 and generally only about 20% to 40% of the original long fragment sequence can be captured as  
350 short fragments and eventually as reads, resulting in shallow sequencing coverage per fragment.  
351 Sequencing more deeply does not increase the per-fragment coverage much as most of the extra  
352 reads are from PCR duplicates. The solution is to sequence multiple 10x libraries constructed  
353 from the same DNA preparation and merge them for analysis. This means that  $C_R$  remains in the  
354 standard range where PCR duplicates are relatively rare, but  $C_F$  increases proportionally to the  
355 number of libraries used. A practical limitation to this approach is that Supernova2 limits the  
356 number of barcodes to 4.8 million.

357

358 Our results showed that in practice,  $C_F$  should be between 335X and 823X, but no larger than  
359 1000X, given the optimal coverage of  $C=56X$  recommended by 10x and the requirement for  
360 sufficient per-fragment read coverage. Surprisingly, we observed that including more extremely  
361 long fragments was detrimental for assembly quality. This is possibly due to the loss of barcode  
362 specificity for fragments spanning repetitive sequences. From a computational perspective, too  
363 many long fragments are harmful to deconvolving the *de bruijn* graph, as more complex paths

364 need to be picked out. In our experiments,  $W_{\mu_{FL}}$  between 50kb and 150kb is the best choice to  
365 generate reliable assemblies.

366

### 367 **Parameters driving assembly quality**

368 Our results regarding assembly quality, and the 10x parameters that influence it, may be useful  
369 for efforts in which *de novo* assemblies are important for generation of an initial reference  
370 sequence. We show that maximization of N50 does not necessarily reflect assembly quality,  
371 which we were able to compare to NA50 because there exists a high-quality human reference  
372 genome. Contig and scaffold lengths mostly increased with ascending sequencing coverage, and  
373 at sufficient overall sequence coverage it did not matter much whether the increasing coverage  
374  $C$  was accomplished by increasing  $C_R$  or  $C_F$ . However, both contig and scaffold accuracy  
375 decreased with increasing  $C$ . We also found, counterintuitively, that contig and scaffold length  
376 mostly decreased with increasing fragment length, a phenomenon that may be due to the specific  
377 implementation; however, until there is another assembler that can be compared to Supernova2  
378 it will not be possible to reason about this effect. In addition, intrinsic properties of the genome  
379 matter greatly, as removal of repeats or lack of variation dramatically improves assembly quality.

380

381 Diploid assembly is the appropriate approach for assembly of genomes of diploid organisms that  
382 harbor variation. Therefore, an important metric to evaluate diploid assembly is the fraction of the  
383 genome that is assembled in a diploid state. The short input fragment length of  $R_6$  resulted in  
384 roughly 20% less of the genome in a diploid state (<60% vs <80%) compared to the other libraries  
385 of the same individual. This observation suggests that in addition to metrics such as N50,  
386 evaluation of assembly quality should also include the fraction of the genome (or the assembly)  
387 that is in a diploid state.

388

389 **Cost-benefit analysis**

390 Overall, we have attempted to give practical guidelines to assembly of 10x data with Supernova2  
391 and evaluate the performance across a wide range of metrics. Arguably, the metric that matters  
392 most in the context of a personal genome is the discovery of variation that lower-cost approaches  
393 do not enable. We estimate that the cost increase over standard Illumina sequencing is about 2x,  
394 given the 10X preparation cost and the higher level of sequence coverage required. There may  
395 be many applications for which this combination of excellent single nucleotide variant detection  
396 (via barcode-aware read mapping) and precise structural variant discovery (via assembly),  
397 achieved by the same data set, is worth the price.

398

399 **Comparison with hybrid assemblies**

400 Hybrid assembly strategies have been applied successfully to produce human genome assembly  
401 of long contiguity [13, 14, 45]. In these studies, long contigs are first produced by single-molecule  
402 long-reads, such as PacBio (NG50=1.1Mb; [13]) or Nanopore (NG50=3.21Mb; [14]) comparing  
403 favorably to our best results for Linked-Reads assemblies (NG50=236kb). Scaffolding is then  
404 performed with complementary technologies such as BioNano to capture chromosomal level long-  
405 range information. It promoted the scaffold N50 of PacBio to 31.1Mb [13] and Illumina mate-pair  
406 sequencing with 10x data to 33.5Mb [25]. Using SuperNova2, the scaffold N50 from our studies  
407 reached ~27.86Mb ( $R_6$ ) on the basis of 10x data alone, suggesting that 10x technology gives  
408 broadly comparable results at a fraction of the price of long-read-based hybrid assemblies.

409

## 410 **Availability of Supporting Data and Materials**

411 The raw sequencing data are deposited in the Sequence Read Archive and the corresponding  
412 BioProject accession number is PRJNA527321. Diploid assemblies and the codes for comparison  
413 are currently available at [http://mendel.stanford.edu/supplementarydata/zhang\\_SN2\\_2019](http://mendel.stanford.edu/supplementarydata/zhang_SN2_2019) and  
414 [https://github.com/zhanglu295/Evaluate\\_diploid\\_assembly](https://github.com/zhanglu295/Evaluate_diploid_assembly). LRTK-SIM is publicly available at  
415 <https://github.com/zhanglu295/LRTK-SIM>. Additional supporting data is available in the  
416 *GigaScience* GigaDB database [46].

417

## 418 **Abbreviation**

419 LFR: Long Fragment Read; ONT: Oxford Nanopore; PacBio: Pacific Biosciences; SNVs: single  
420 nucleotide variants; SVs: structural variants

421

## 422 **Additional files**

423 **Table S1.** Parameters of libraries prepared for NA12878 and NA24385.

424 **Table S2.** Parameters used to generate linked-read sets for evaluating the impact of  $C_F$  and  $C_R$   
425 on assemblies.

426 **Table S3.** Parameters used to generate linked-read sets for evaluating the impact of  $\mu_{FL}$  and  
427  $N_{FP}$  on assemblies.

428 **Table S4.** Contig misassemblies and recovered transcripts of the six assemblies.

429 **Table S5.** Genomic coverage and fraction of contigs in diploid state generated by Supernova2  
430 for the seven libraries prepared by 10x Genomics. Non-PAR: non-pseudoautosomal regions of  
431 X chromosome. WFU, YOR, YORM, PR are female; HGP, ASH and CHI are male.

432 **Table S6.** Phase block N50s of the six assemblies.

433 **Table S7.** Comparison SNV calls from standard Illumina data, 10x reference-based calls, and  
434 assembly-based calls for NA12878. All calls were compared to the Genome in a Bottle benchmark.

435 **Table S8.** Comparison SNV calls from standard Illumina data, 10x reference-based calls, and  
436 assembly-based calls for NA24385. All calls were compared to the Genome in a Bottle benchmark.

437 **Table S9.** Comparison of SV calls from standard Illumina data and 10x assembly-based calls for  
438 NA12878.

439 **Table S10.** Comparison of SV calls from standard Illumina data and 10x assembly-based calls  
440 for NA24385.

441 **Table S11.** Comparison of SV calls from 10x reference-based and assembly-based calls for  
442 NA12878.

443 **Table S12.** Comparison of SV calls from 10x reference-based and assembly-based calls for  
444 NA24385.

445 **Table S13.** Comparison of SV calls from our de novo assemblies with the Tier 1 SV benchmark  
446 from Genome in a Bottle.

447 **Table S14.** Proportion of assembly-based SV calls supported by PacBio CCS reads.

448 **Figure S1. Basic statistics for  $L_{1L}$ .** The distributions of **A.** the number of fragments per partition;  
449 **B.** sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
450 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
451 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
452 fragment lengths.

453 **Figure S2. Basic statistics for  $L_{1M}$ .** The distributions of **A.** number of fragments per partition; **B.**  
454 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
455 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
456 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
457 fragment lengths.

458 **Figure S3. Basic statistics for  $L_{1H}$ .** The distributions of **A.** number of fragments per partition; **B.**  
459 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
460 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
461 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
462 fragment lengths.

463 **Figure S4. Basic statistics for  $L_2$ .** The distributions of **A.** number of fragments per partition; **B.**  
464 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
465 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
466 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
467 fragment lengths.

468 **Figure S5. Basic statistics for  $L_3$ .** The distributions of **A.** number of fragments per partition; **B.**  
469 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
470 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
471 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
472 fragment lengths.

473 **Figure S6. Basic statistics for  $L_4$ .** The distributions of **A.** number of fragments per partition; **B.**  
474 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
475 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
476 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
477 fragment lengths.

478 **Figure S7. Basic statistics for  $L_5$ .** The distributions of **A.** number of fragments per partition; **B.**  
479 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
480 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
481 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
482 fragment lengths.

483 **Figure S8. Basic statistics for  $L_6$ .** The distributions of **A.** number of fragments per partition; **B.**  
484 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
485 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
486 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
487 fragment lengths.

488 **Figure S9.** The workflow of LRTK-SIM to simulate linked-reads

489 **Figure S10.** The effect of  $N_{FP}$  on human diploid assembly of chromosome 19 by Supernova2,  
490 where  $C$  ( $C=60X$ ;  $C_F=300X$  and  $C_R=0.2X$ ) and  $\mu_{FL}$  ( $\mu_{FL}=37kb$ ) are fixed.

491 **Figure S11.** Comparison of assembly qualities from 10x data with and without single nucleotide  
492 variants by changing  $C_F$ ,  $C_R$  and  $\mu_{FL}$ .  $C_R$  was fixed to 0.2X in **A** and **B**;  $C_F$  was fixed to 300X in  
493 **C** and **D**;  $C_R$  was fixed 0.2X and  $C_F$  was fixed 300X in **E** and **F**.

494 **Figure S12.** Comparison of assembly qualities from 10x data with (1% uniform) and without  
495 sequencing error by changing  $C_F$ ,  $C_R$  and  $\mu_{FL}$ .  $C_R$  was fixed to 0.2X in **A** and **B**;  $C_F$  was fixed to  
496 300X in **C** and **D**;  $C_R$  was fixed 0.2X and  $C_F$  was fixed 300X in **E** and **F**.

497 **Figure S13.** Overlaps of diploid regions for the three libraries from the same sample. Diploid  
498 regions for NA12878 (**A**) and NA24385 (**B**). The percentages denote the proportion of genome is  
499 diploid.

500 **Figure S14.** Phase block N50s as a function of different parameter combinations. **A.** simulated  
501 linked-reads with predefined parameters (**Table S5**) by changing  $C_F$  and  $C_R$ ; **B.** simulated linked-  
502 reads with matched parameters of real linked-read sets (**Table S2**) by changing  $C_F$  and  $C_R$ ; **C.**  
503 real linked-read sets (**Table S2**) by changing  $C_F$  and  $C_R$ ; **D.** simulated linked-read sets (**Table S3**)  
504 with different  $W_{\mu_{FL}}$ ; **E.** simulated linked-read sets with matched parameters (**Table S3**) with real  
505 linked-read sets as  $C=56X$ ; **F.** real linked-read sets with  $C=56X$  (**Table S3**).

506

507

508 **Competing interests**

509 Arend Sidow is a consultant and shareholder of DNAnexus, Inc.

510

511 **Authors' contributions**

512 AS conceived the study. LZ and XZ wrote LRTK-SIM and performed the analyses. ZMW prepared  
513 the genomic DNA and 10x libraries. LZ, XZ, ZMW and AS analyzed the results and wrote the  
514 paper. All authors read and approved the final manuscript.

515

516 **Acknowledgements**

517 This research was supported by training and research grants from the National Institute of  
518 Standards and Technology (NIST). We would like to thank Justin Zook, Marc Salit, Alex Bishara,  
519 Noah Spies, Nancy Hansen, David Jaffe, and Deanna Church for informative discussions.

520

521



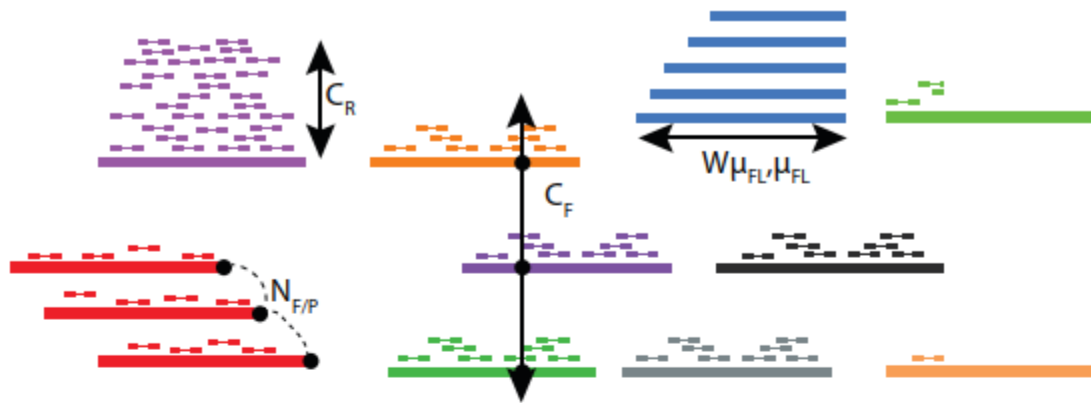
522 **Table**

Linked- reads set	Overall (%)	Diploid regions (%)	Haploid regions (%)	Non-PAR (%)	Total contig length (contig>500bp)	Length of contigs from megabubble (contig>500bp)	Percentage (%)
$R_6$	91.9	58.9	27.7	-	5,632,483,053	3,758,345,846	66.73
$R_7$	91.1	73.3	11.3	-	5,613,140,437	4,668,186,478	83.17
$R_8$	91.7	77.2	9.2	-	5,635,127,471	4,896,821,850	86.90
$R_9$	91.3	73.4	12.2	85.9	5,637,615,919	4,438,175,621	78.72
$R_{10}$	91.7	79.2	5.8	79.9	5,749,001,471	4,793,226,150	83.37
$R_{11}$	91.7	78.1	7.9	87.6	5,677,566,094	4,723,083,367	83.19

523

524 **Table 1.** Genomic coverage of contigs generated by Supernova2. Non-PAR: non-  
525 pseudoautosomal regions of X chromosome.  $R_6$ ,  $R_7$  and  $R_8$  are female;  $R_9$ ,  $R_{10}$  and  $R_{11}$  are male.

526



**Parameter**

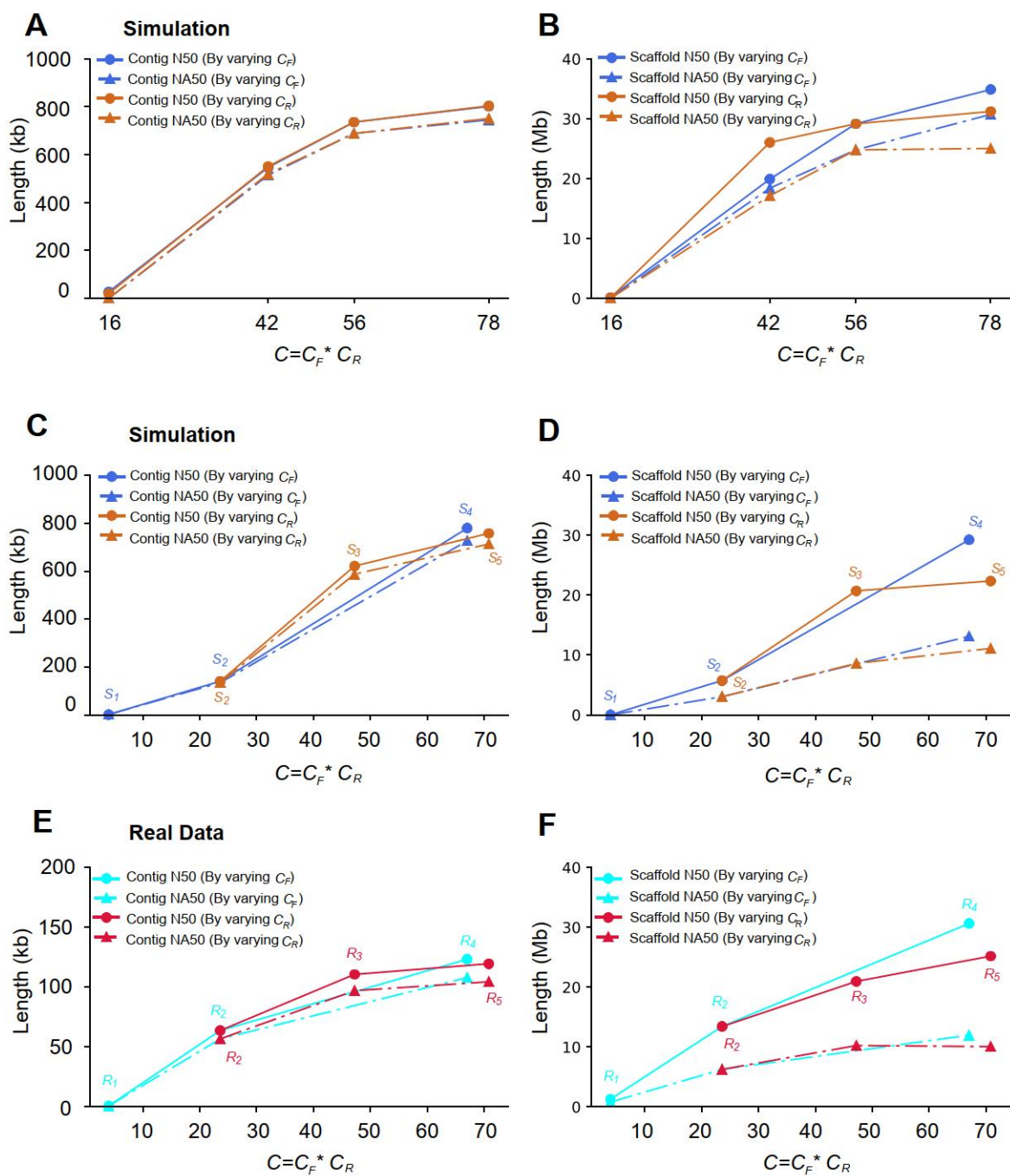
$N_{F/P}$  = Number of fragments per partition  
 $\mu_{FL}$  = Mean fragment length  
 $W\mu_{FL}$  = Weighted mean fragment length  
 $C_R$  = Read coverage per fragment  
 $C_F$  = Physical (fragment) coverage  
 $C$  = total coverage

**Typical values**

10 - 100  
 $\mu_{FL}$  = 10-100kb  
 $W\mu_{FL}$  = 20-400kb  
 $C_R$  = 0.1x - 0.4x  
 $C_F$  = 200x - 1000x  
 $C = C_R * C_F = 40x - 80x$

Linked-read set R (Real) / S (Simulated)	Sequenced Library	$\mu_{FL}$ (kb)	$W\mu_{FL}$ (kb)	$C_F$ (X)	$C_R$ (X)	$C$ (X)
$R_1 / S_1$	$L_{1L}$	21.6	38.6/35.7	19	0.2	4
$R_2 / S_2$	$L_{1M}$	22.4	39.7/37.4	117	0.2	24
$R_3 / S_3$	$L_{1M}$	22.4	39.7/36.8	117	0.4	48
$R_4 / S_4$	$L_{1H}$	24.0	41.1/40.7	334	0.2	67
$R_5 / S_5$	$L_{1M}$	22.4	39.7/36.8	117	0.6	72
$R_6 / S_6$	$L_{1H}$	24.0	41.1/40.6	334	0.17	56
$R_7 / S_7$	$L_2$	79.0	304.3/131.8	123	0.45	56
$R_8 / S_8$	$L_3$	99.2	214.5/168.3	958	0.058	56
$R_9 / S_9$	$L_4$	92.1	216.9/154.1	1504	0.036	56
$R_{10} / S_{10}$	$L_5$	120.8	267.4/203.7	208	0.27	56
$R_{11} / S_{11}$	$L_6$	64.2	151.7/107.6	803	0.07	56

528 **Figure 1.** The linked-read sets prepared to evaluate the impact of  $C_F$ ,  $C_R$ ,  $\mu_{FL}$  and  $W\mu_{FL}$  on  
 529 human diploid assembly.



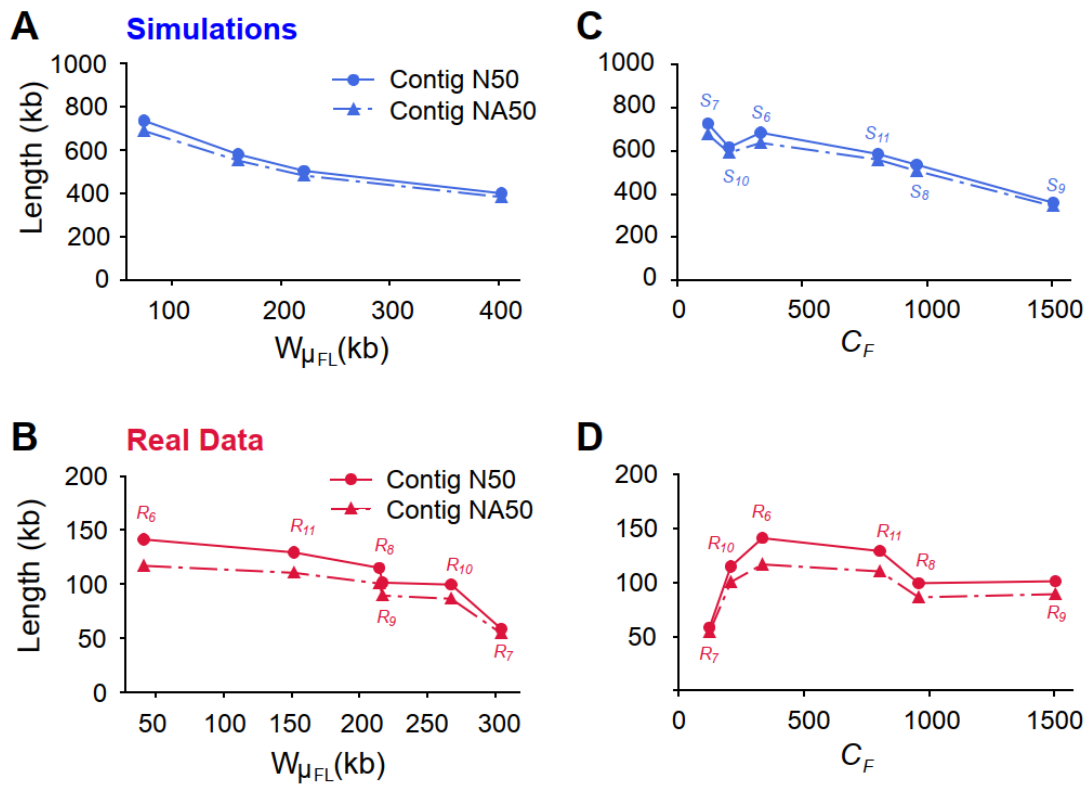
530

531 **Figure 2.** Contig and scaffold lengths (N50 and NA50) as a function of  $C_F$  or  $C_R$ . **A and B:**

532 Simulated Linked-Reads with predefined parameters (**Table S2**); **C and D:** Simulated Linked-

533 reads with matched parameters of real Linked-Read data sets (**Figure 1**); **E and F:** Real linked-

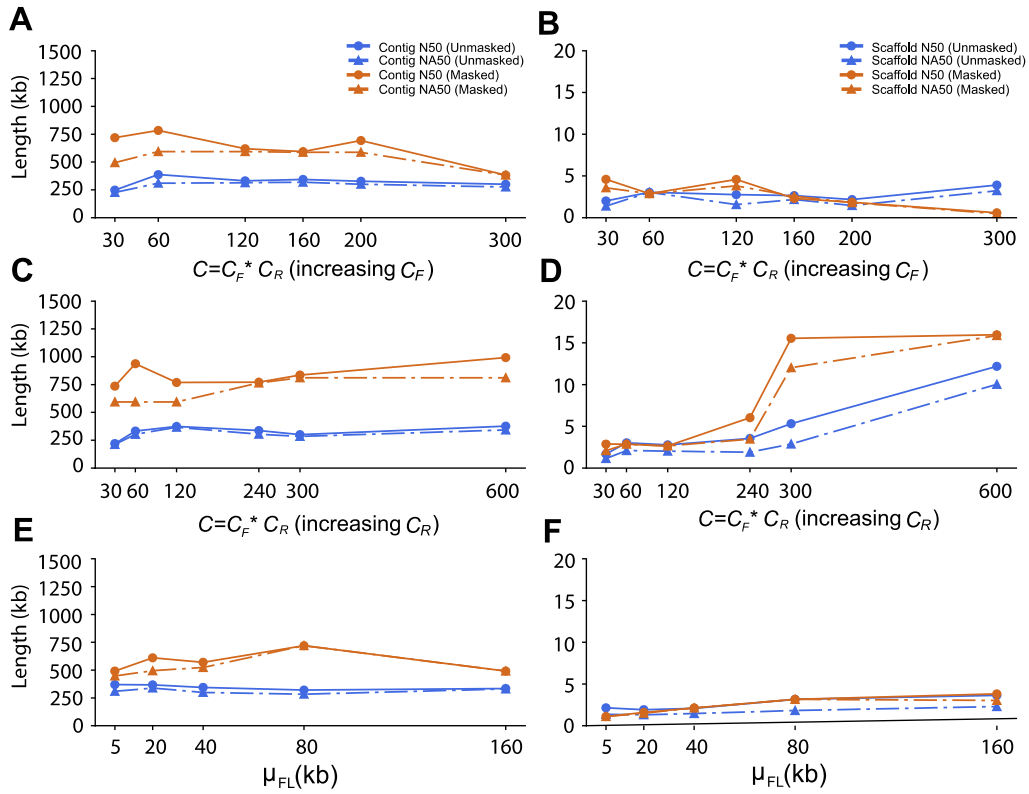
534 read sets (**Figure 1**).



535

536 **Figure 3.** Contig qualities (N50 and NA50) as a function of fragment length  $W_{\mu_{FL}}$  or physical

537 coverage  $C_F$ , at  $C=56X$ . **A** and **C**, results from simulations; **B** and **D**, results from real data.



538

539 **Figure 4.** Comparison of contig and scaffold lengths from 10x data with masked and unmasked

540 repetitive sequences by changing  $C_F$ ,  $C_R$  and  $\mu_{FL}$ .  $C_R$  was fixed to 0.2X in **A** and **B**;  $C_F$  was fixed

541 to 300X in **C** and **D**;  $C_R$  was fixed to 0.2X and  $C_F$  was fixed to 300X in **E** and **F**.

542

543 **References**

544 1. Metzker ML. Sequencing technologies - the next generation. *Nat Rev Genet.* 2010;11 1:31-  
545 46. doi:10.1038/nrg2626.

546 2. Shendure J, Balasubramanian S, Church GM, Gilbert W, Rogers J, Schloss JA, et al. DNA  
547 sequencing at 40: past, present and future. *Nature.* 2017;550 7676:345-53.  
548 doi:10.1038/nature24286.

549 3. Head SR, Komori HK, LaMere SA, Whisenant T, Van Nieuwerburgh F, Salomon DR, et  
550 al. Library construction for next-generation sequencing: overviews and challenges.  
551 *Biotechniques.* 2014;56 2:61-4, 6, 8, passim. doi:10.2144/000114133.

552 4. O'Connell J, Sharp K, Shrine N, Wain L, Hall I, Tobin M, et al. Haplotype estimation for  
553 biobank-scale data sets. *Nat Genet.* 2016;48 7:817-20. doi:10.1038/ng.3583.

554 5. Delaneau O, Zagury JF and Marchini J. Improved whole-chromosome phasing for disease  
555 and population genetic studies. *Nat Methods.* 2013;10 1:5-6. doi:10.1038/nmeth.2307.

556 6. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, et al. A general  
557 approach for haplotype phasing across the full spectrum of relatedness. *PLoS Genet.*  
558 2014;10 4:e1004234. doi:10.1371/journal.pgen.1004234.

559 7. Roach JC, Glusman G, Hubley R, Montsaroff SZ, Holloway AK, Mauldin DE, et al.  
560 Chromosomal haplotypes by genetic phasing of human families. *Am J Hum Genet.*  
561 2011;89 3:382-97. doi:10.1016/j.ajhg.2011.07.023.

562 8. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al. Efficient de  
563 novo assembly of highly heterozygous genomes from whole-genome shotgun short reads.  
564 *Genome Res.* 2014;24 8:1384-95. doi:10.1101/gr.170720.113.

565 9. Alkan C, Sajjadian S and Eichler EE. Limitations of next-generation genome sequence  
566 assembly. *Nat Methods.* 2011;8 1:61-5. doi:10.1038/nmeth.1527.

567 10. Treangen TJ and Salzberg SL. Repetitive DNA and next-generation sequencing:  
568 computational challenges and solutions. *Nat Rev Genet.* 2011;13 1:36-46.  
569 doi:10.1038/nrg3117.

570 11. Huddleston J, Ranade S, Malig M, Antonacci F, Chaisson M, Hon L, et al. Reconstructing  
571 complex regions of genomes using long-read sequencing technology. *Genome Res.*  
572 2014;24 4:688-96. doi:10.1101/gr.168450.113.

573 12. Lu H, Giordano F and Ning Z. Oxford Nanopore MinION Sequencing and Genome  
574 Assembly. *Genomics Proteomics Bioinformatics.* 2016;14 5:265-79.  
575 doi:10.1016/j.gpb.2016.05.004.

576 13. Pendleton M, Sebra R, Pang AW, Ummat A, Franzen O, Rausch T, et al. Assembly and  
577 diploid architecture of an individual human genome via single-molecule technologies. *Nat*  
578 *Methods.* 2015;12 8:780-6. doi:10.1038/nmeth.3454.

579 14. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore sequencing and  
580 assembly of a human genome with ultra-long reads. *Nat Biotechnol.* 2018;36 4:338-45.  
581 doi:10.1038/nbt.4060.

582 15. Wenger AM, Peluso P, Rowell WJ, Chang P-C, Hall RJ, Concepcion GT, et al. Accurate  
583 circular consensus long-read sequencing improves variant detection and  
584 assembly of a human genome. *Nat Biotechnol.* 2019 Oct;37(10):1155-1162. doi:  
585 10.1038/s41587-019-0217-9.

586 16. Cao H, Wu H, Luo R, Huang S, Sun Y, Tong X, et al. De novo assembly of a haplotype-  
587 resolved human genome. *Nat Biotechnol.* 2015;33 6:617-22. doi:10.1038/nbt.3200.

- 588 17. Kuleshov V, Xie D, Chen R, Pushkarev D, Ma Z, Blauwkamp T, et al. Whole-genome  
589 haplotyping using long reads and statistical methods. *Nat Biotechnol.* 2014;32 3:261-6.  
590 doi:10.1038/nbt.2833.
- 591 18. Peters BA, Kermani BG, Sparks AB, Alferov O, Hong P, Alexeev A, et al. Accurate whole-  
592 genome sequencing and haplotyping from 10 to 20 human cells. *Nature.* 2012;487  
593 7406:190-5. doi:10.1038/nature11236.
- 594 19. Edge P, Bafna V and Bansal V. HapCUT2: robust and accurate haplotype assembly for  
595 diverse sequencing technologies. *Genome Res.* 2017;27 5:801-12.  
596 doi:10.1101/gr.213462.116.
- 597 20. Patterson M, Marschall T, Pisanti N, van Iersel L, Stougie L, Klau GW, et al. WhatsHap:  
598 Weighted Haplotype Assembly for Future-Generation Sequencing Reads. *J Comput Biol.*  
599 2015;22 6:498-509. doi:10.1089/cmb.2014.0157.
- 600 21. Zheng GX, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, et al. Haplotyping  
601 germline and cancer genomes with high-throughput linked-read sequencing. *Nat*  
602 *Biotechnol.* 2016;34 3:303-11. doi:10.1038/nbt.3432.
- 603 22. Spies N, Weng Z, Bishara A, McDaniel J, Catoe D, Zook JM, et al. Genome-wide  
604 reconstruction of complex structural variants using read clouds. *Nat Methods.* 2017;14  
605 9:915-20. doi:10.1038/nmeth.4366.
- 606 23. Bishara A, Moss EL, Kolmogorov M, Parada AE, Weng Z, Sidow A, et al. High-quality  
607 genome sequences of uncultured microbes by assembly of read clouds. *Nat Biotechnol.*  
608 2018; doi:10.1038/nbt.4266.
- 609 24. Weisenfeld NI, Kumar V, Shah P, Church DM and Jaffe DB. Direct determination of  
610 diploid genome sequences. *Genome Res.* 2017;27 5:757-67. doi:10.1101/gr.214874.116.
- 611 25. Mostovoy Y, Levy-Sakin M, Lam J, Lam ET, Hastie AR, Marks P, et al. A hybrid approach  
612 for de novo human genome sequence assembly and phasing. *Nat Methods.* 2016;13 7:587-  
613 90. doi:10.1038/nmeth.3865.
- 614 26. Hulse-Kemp AM, Maheshwari S, Stoffel K, Hill TA, Jaffe D, Williams SR, et al. Reference  
615 quality assembly of the 3.5-Gb genome of *Capsicum annuum* from a single linked-read  
616 library. *Hortic Res.* 2018;5:4. doi:10.1038/s41438-017-0011-0.
- 617 27. Elyanow R, Wu HT and Raphael BJ. Identifying structural variants using linked-read  
618 sequencing data. *Bioinformatics.* 2017; doi:10.1093/bioinformatics/btx712.
- 619 28. Jones SJ, Haulena M, Taylor GA, Chan S, Bilobram S, Warren RL, et al. The Genome of  
620 the Northern Sea Otter (*Enhydra lutris kenyoni*). *Genes (Basel).* 2017;8 12  
621 doi:10.3390/genes8120379.
- 622 29. Luo R, Sedlazeck FJ, Darby CA, Kelly SM and Schatz MC. LRSim: A Linked-Reads  
623 Simulator Generating Insights for Better Genome Partitioning. *Comput Struct Biotechnol*  
624 *J.* 2017;15:478-84. doi:10.1016/j.csbj.2017.10.002.
- 625 30. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence  
626 Alignment/Map format and SAMtools. *Bioinformatics.* 2009;25 16:2078-9.  
627 doi:10.1093/bioinformatics/btp352.
- 628 31. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34  
629 18:3094-100. doi:10.1093/bioinformatics/bty191.
- 630 32. Mikheenko A, Prjibelski A, Saveliev V, Antipov D and Gurevich A. Versatile genome  
631 assembly evaluation with QUAST-LG. *Bioinformatics.* 2018;34 13:i142-i50.  
632 doi:10.1093/bioinformatics/bty266.

633 33. Earl D, Bradnam K, St John J, Darling A, Lin D, Fass J, et al. Assemblathon 1: a  
634 competitive assessment of de novo short read assembly methods. *Genome Res.* 2011;21  
635 12:2224-41. doi:10.1101/gr.126599.111.

636 34. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, et al. Extensive sequencing of  
637 seven human genomes to characterize benchmark reference materials. *Sci Data.*  
638 2016;3:160025. doi:10.1038/sdata.2016.25.

639 35. Wala JA, Bandopadhyay P, Greenwald NF, O'Rourke R, Sharpe T, Stewart C, et al.  
640 SvABA: genome-wide detection of structural variants and indels by local assembly.  
641 *Genome Res.* 2018;28 4:581-91. doi:10.1101/gr.221028.117.

642 36. Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler  
643 transform. *Bioinformatics.* 2009;25 14:1754-60. doi:10.1093/bioinformatics/btp324.

644 37. Garrison E and Marth G. Haplotype-based variant detection from short-read sequencing.  
645 arXiv e-prints. 2012.

646 38. Zook JM, Hansen NF, Olson ND, Chapman LM, Mullikin JC, Xiao C, et al. A robust  
647 benchmark for germline structural variant detection. *bioRxiv.* 2019.

648 39. Zhao X, Weber AM and Mills RE. A recurrence-based approach for validating structural  
649 variation using long-read sequencing technology. *GigaScience.* 2017;6 8:1-9.  
650 doi:10.1093/gigascience/gix061.

651 40. Wenger AM, Peluso P, Rowell WJ, Chang PC, Hall RJ, Concepcion GT, et al. Accurate  
652 circular consensus long-read sequencing improves variant detection and assembly of a  
653 human genome. *Nat Biotechnol.* 2019;37 10:1155-62. doi:10.1038/s41587-019-0217-9.

654 41. Krusche P, Trigg L, Boutros PC, Mason CE, De La Vega FM, Moore BL, et al. Best  
655 practices for benchmarking germline small-variant calls in human genomes. *Nat*  
656 *Biotechnol.* 2019;37 5:555-60. doi:10.1038/s41587-019-0054-x.

657 42. Zook JM, McDaniel J, Olson ND, Wagner J, Parikh H, Heaton H, et al. An open resource  
658 for accurately benchmarking small variant and reference calls. *Nat Biotechnol.* 2019;37  
659 5:561-6. doi:10.1038/s41587-019-0074-6.

660 43. Zhang F, Christiansen L, Thomas J, Pokholok D, Jackson R, Morrell N, et al. Haplotype  
661 phasing of whole human genomes using bead-based barcode partitioning in a single tube.  
662 *Nat Biotechnol.* 2017;35 9:852-7. doi:10.1038/nbt.3897.

663 44. Wang O, Chin R, Cheng X, Wu MKY, Mao Q, Tang J, et al. Efficient and unique  
664 cobarcoding of second-generation sequencing reads from long DNA molecules enabling  
665 cost-effective and accurate sequencing, haplotyping, and de novo assembly. *Genome Res.*  
666 2019;29 5:798-808. doi:10.1101/gr.245126.118.

667 45. Ma ZS, Li L, Ye C, Peng M and Zhang YP. Hybrid assembly of ultra-long Nanopore reads  
668 augmented with 10x-Genomics contigs: Demonstrated with a human genome. *Genomics.*  
669 2018; doi:10.1016/j.ygeno.2018.12.013.

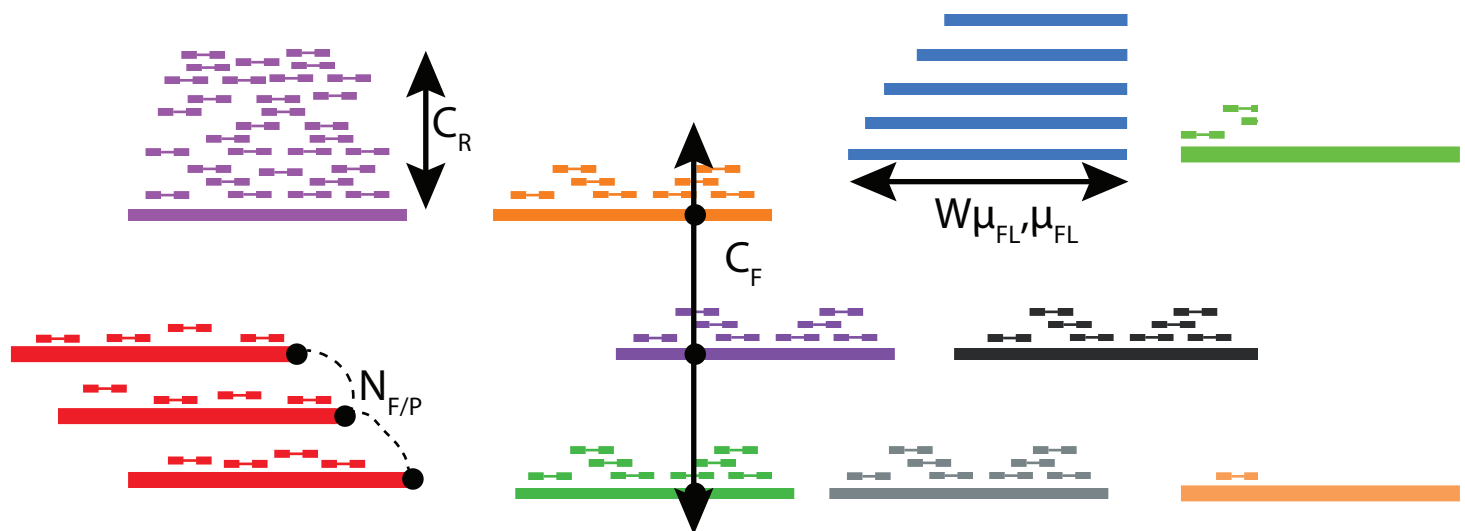
670 46. Zhang L; Zhou X; Weng Z; Sidow A (2019): Supporting data for "Assessment of human  
671 diploid genome assembly with 10x Linked-Reads data" *GigaScience Database.*  
672 <http://dx.doi.org/10.5524/100668>  
673



**Table**

Linked-reads set	Overall (%)	Diploid regions (%)	Haploid regions (%)	Non-PAR (%)	Total contig length (contig>500bp)	Length of contigs from megabubble (contig>500bp)	Percentage (%)
$R_6$	91.9	58.9	27.7	-	5,632,483,053	3,758,345,846	66.73
$R_7$	91.1	73.3	11.3	-	5,613,140,437	4,668,186,478	83.17
$R_8$	91.7	77.2	9.2	-	5,635,127,471	4,896,821,850	86.90
$R_9$	91.3	73.4	12.2	85.9	5,637,615,919	4,438,175,621	78.72
$R_{10}$	91.7	79.2	5.8	79.9	5,749,001,471	4,793,226,150	83.37
$R_{11}$	91.7	78.1	7.9	87.6	5,677,566,094	4,723,083,367	83.19

**Table 1.** Genomic coverage of contigs generated by Supernova2. Non-PAR: non-pseudoautosomal regions of X chromosome.  $R_6$ ,  $R_7$  and  $R_8$  are female;  $R_9$ ,  $R_{10}$  and  $R_{11}$  are male.



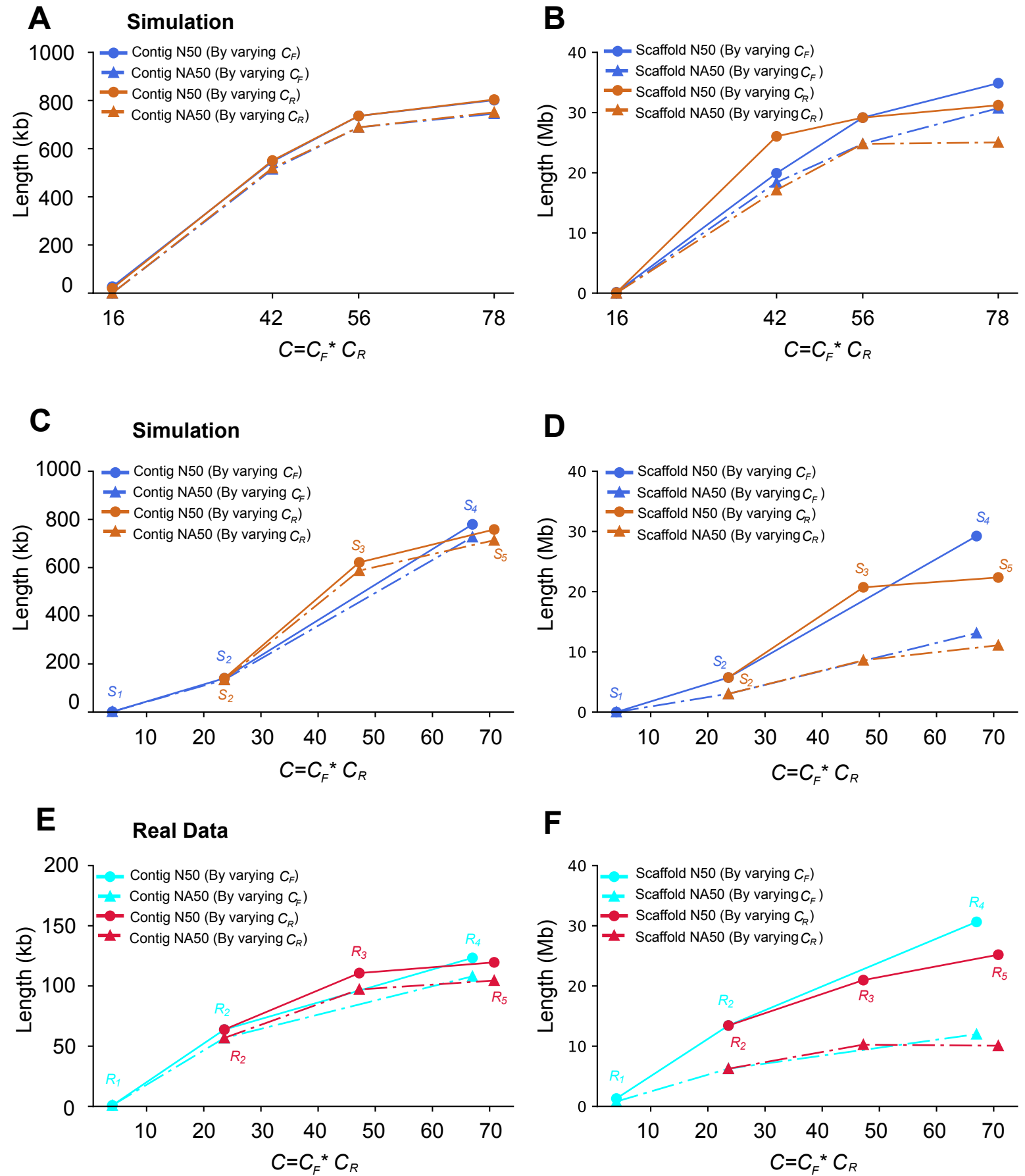
### Parameter

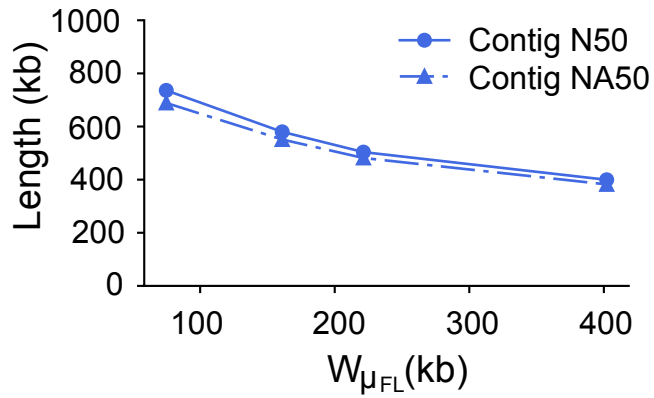
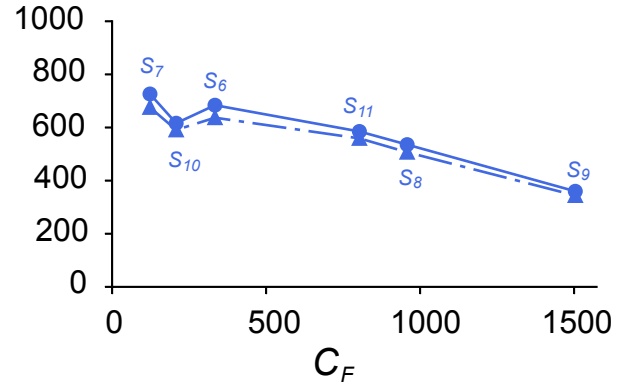
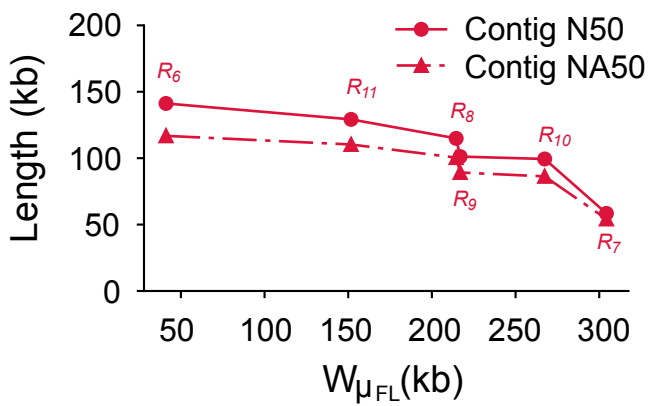
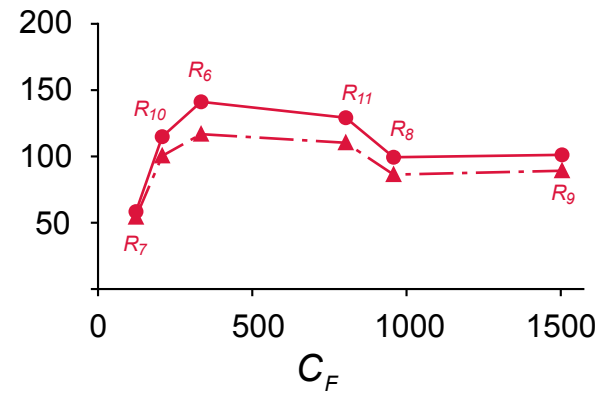
$N_{F/P}$  = Number of fragments per partition  
 $\mu_{FL}$  = Mean fragment length  
 $W\mu_{FL}$  = Weighted mean fragment length  
 $C_R$  = Read coverage per fragment  
 $C_F$  = Physical (fragment) coverage  
 $C$  = total coverage

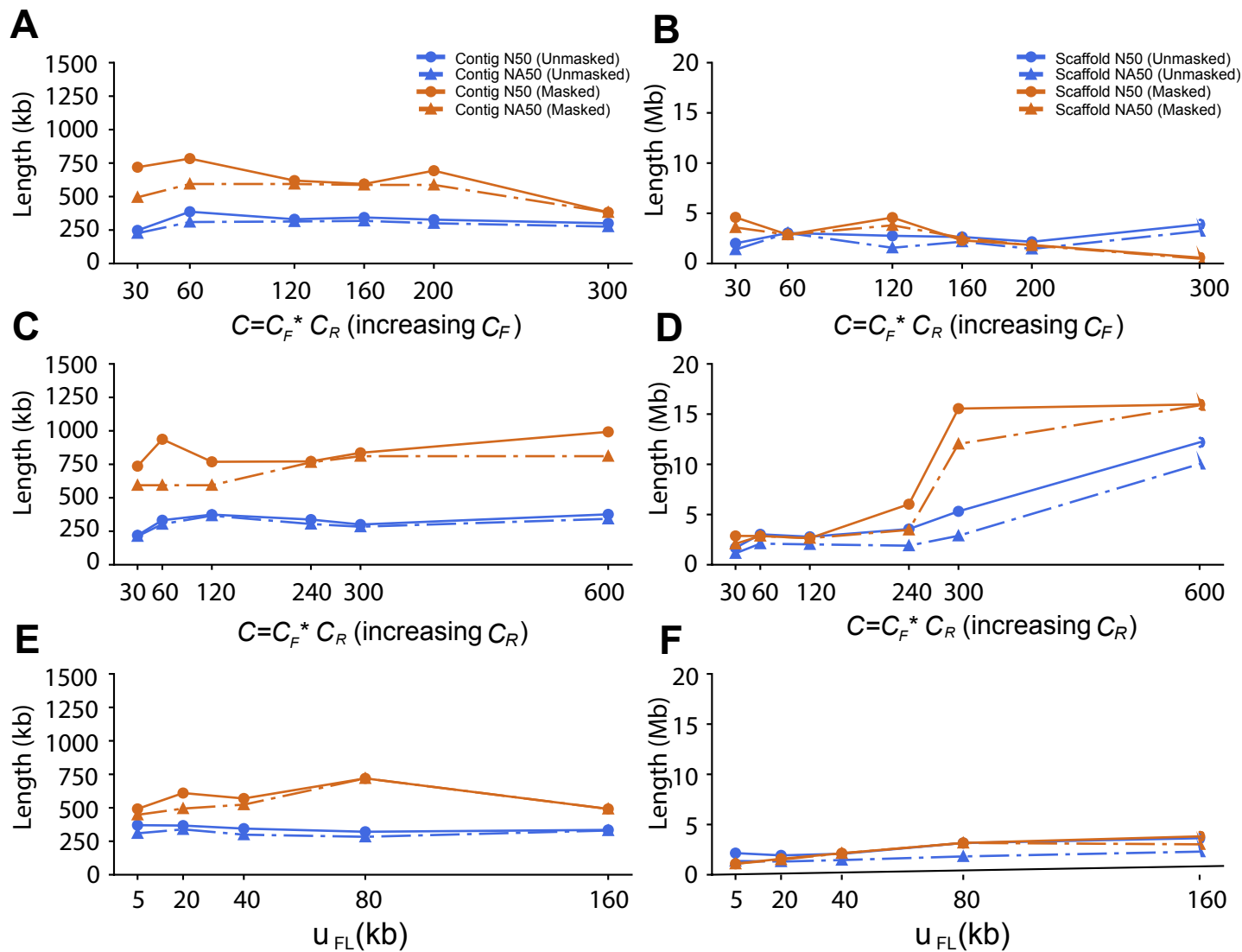
### Typical values

10 - 100  
 $\mu_{FL}$  = 10-100kb  
 $W\mu_{FL}$  = 20-400kb  
 $C_R$  = 0.1x - 0.4x  
 $C_F$  = 200x - 1000x  
 $C = C_R * C_F = 40x - 80x$

Linked-read set R (Real) / S (Simulated)	Sequenced Library	$\mu_{FL}$ (kb)	$W\mu_{FL}$ (kb)	$C_F$ (X)	$C_R$ (X)	$C$ (X)
$R_1 / S_1$	$L_{1L}$	21.6	38.6/35.7	19	0.2	4
$R_2 / S_2$	$L_{1M}$	22.4	39.7/37.4	117	0.2	24
$R_3 / S_3$	$L_{1M}$	22.4	39.7/36.8	117	0.4	48
$R_4 / S_4$	$L_{1H}$	24.0	41.1/40.7	334	0.2	67
$R_5 / S_5$	$L_{1M}$	22.4	39.7/36.8	117	0.6	72
$R_6 / S_6$	$L_{1H}$	24.0	41.1/40.6	334	0.17	56
$R_7 / S_7$	$L_2$	79.0	304.3/131.8	123	0.45	56
$R_8 / S_8$	$L_3$	99.2	214.5/168.3	958	0.058	56
$R_9 / S_9$	$L_4$	92.1	216.9/154.1	1504	0.036	56
$R_{10} / S_{10}$	$L_5$	120.8	267.4/203.7	208	0.27	56
$R_{11} / S_{11}$	$L_6$	64.2	151.7/107.6	803	0.07	56



**A Simulations****C****B Real Data****D**





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





DEPARTMENT OF PATHOLOGY  
DEPARTMENT OF GENETICS  
STANFORD UNIVERSITY SCHOOL OF MEDICINE  
STANFORD, CA 94305-5324

Stanford, October 22, 2019

Dr. Hongling Zhou  
Editor  
GigaScience

Dear Dr. Zhou,

Thank you for letting us submit a final revised version of our manuscript "Assessment of human diploid genome assembly with 10x Linked-Reads data" for your further consideration for publication in GigaScience. We were able to address the remaining comments, which are addressed point by point in our response, and hope that you will be able to reach a positive decision.

Sincerely,

A handwritten signature in blue ink that reads "Arend Sidow".

Arend Sidow, Ph.D.  
Professor of Pathology and of Genetics  
SUMC R353  
Stanford, CA 94305-5324

arend@stanford.edu  
+1-650-498-7024  
<http://www.sidowlab.org>  
<http://jimb.stanford.edu>