

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

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

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

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

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

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

12

13

## 14 **Abstract**

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

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

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

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

## 36 **Data description**

### 37 **Introduction**

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

50

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

60 mobile elements and by segmental duplications, and only support highly fragmented genome  
61 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  
66 by hierarchically integrating PacBio long reads and BioNano for a hybrid assembly, which also  
67 uncovered novel tandem repeats and replicated the structural variants that were newly included  
68 in the updated hg38 human reference sequence. Another study [14] produced human genome  
69 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  
74 genome 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  
79 intensive to generate and sequence, this approach is not generalizable to personal genomes.  
80 The "Long Fragment Read" (LFR) approach [17], where a long fragment is sequenced at high  
81 depth via single-molecule fragmented amplification, reported promising personal genome  
82 assembly and variant phasing by attaching a barcode to the short reads derived from the same  
83 long fragment. However, because LFR is implemented in a 384 well plate, many long fragments  
84 would be 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  
88 the 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] ,  
98 and *de novo* assembly of human and other genomes [24-26], compared to standard Illumina  
99 short-fragment sequencing. While the uniformity of sequence coverage is not as good as with  
100 PCR-free Illumina libraries, 10x Linked-Read sequencing is a promising technology that  
101 combines low per-base error and good small-variant discovery with long-range information for  
102 much improved SV detection in mapping-based approaches [22, 27], and the possibility of long-  
103 range contiguity 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  
107 vary, along with other factors, to influence an interconnected array of parameters that are  
108 relevant to genome assembly and reconstruction. The parameters over which the experimenter  
109 has influence are (**Figure 1**): i).  $C_R$ : average **C**overage of short **R**eads per fragment; ii).  $C_F$ :  
110 average physical **C**overage of the genome by long DNA **F**ragments; iii).  $N_{FP}$ : **N**umber of  
111 **F**ragments per **P**artition; iv). Fragment length distribution, several parameters of which are used,

112 specifically  $\mu_{FL}$ : Average Unweighted DNA Fragment Length and  $W\mu_{FL}$ : Length-Weighted  
113 average of DNA Fragment Length. Note that several parameters depend on each other. For  
114 example, a greater amount of input DNA will increase  $N_{FP}$ ; shorter fragments increase  $N_{FP}$  at  
115 the same DNA input amount compared to longer fragments; less input DNA will (within practical  
116 constraints) increase  $C_R$  and decrease  $C_F$ , and their absolute values are set by how much total  
117 sequence coverage is 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  
123 parameter space, we first performed computer simulations with reasonably realistic synthetic  
124 data. The simulation results suggested certain parameter combinations that we then  
125 approximated in the generation of real, high-depth, sequence data on two human reference  
126 genome cell lines, NA12878 and NA24385. These simulated and real data sets were then used  
127 to produce *de novo* assemblies, with an emphasis on the performance of 10x's Supernova2 [24].  
128 We finally assessed the quality of the assemblies using standard metrics of contiguity and  
129 accuracy, facilitated by the existence of a gold standard (in the case of simulations) and  
130 comparisons to the reference genome (in the case of real data).

131

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

133 We made six DNA preparations that varied in fragment size distribution and amount of input  
134 DNA, three each from NA12878 and NA24385. From these, we prepared eight libraries, five  
135 from NA12878 and three from NA24385 (**Table S1**). To generate libraries  $L_{1L}$ ,  $L_{1M}$  and  $L_{1H}$  (the  
136 subscripts  $L$ ,  $M$  and  $H$  represent low, medium and high  $C_F$ , respectively), genomic DNA was

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

157

### 158 **Linked-Reads subsampling**

159 The high sequencing coverage in the libraries allowed subsampling to facilitate the matching of  
160 parameters among the different libraries, for purposes of comparability; these subsampled  
161 Linked-Read sets are denoted  $R_{id}$  (**Figure 1**). We aligned the 10x Linked-Reads to human

162 reference genome (hg38, GRCh38 Reference 2.1.0 from 10x website) followed by removing  
163 PCR duplication by barcode-aware analysis in Long Ranger[21]. Original input DNA fragments  
164 were inferred by collecting the read-pairs with the same barcode that were aligned in proximity  
165 to each other. A fragment was terminated if the distance between two consecutive reads with  
166 the identical barcode larger than 50kb. Fragments were required to have at least two read pairs  
167 with the same barcode and a length of at least 2 kb. Partitions with fewer than three fragments  
168 were removed. We subsampled short-reads for each fragment to satisfy the expected  $C_R$ .

169

### 170 **Generating 10x simulated libraries by LR TK-SIM**

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



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

202

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

204 The scaffolds were generated by the “pseudohap2” output of Supernova2, which explicitly  
205 generated two haploid scaffolds, simultaneously. Contigs were generated by breaking the  
206 scaffolds if at least 10 consecutive ‘N’s appeared, per definition by Supernova2. For the  
207 simulations of human chromosome 19, we used the scaffolds from the “megabubbles” output.  
208 Contig and scaffold N50 and NA50 were used to evaluate assembly quality. Contigs longer than  
209 500bp were aligned to hg38 by Minimap2[29]. We calculated contig NA50 on the basis of contig  
210 misassemblies reported by QUAST-LG [30]. For scaffolds (longer than 1kb), we calculated the  
211 NA50 following Assemblathon 1’s procedure [31] (**Supplementary Information**).

212

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

214 We compared single nucleotide variants (SNVs) and structural variants (SVs) from the diploid  
215 regions of our assemblies with the ones from standard Illumina data and reference-based  
216 processing of our 10x data. The standard Illumina data were downloaded from Genome in a  
217 Bottle [32] and analyzed with SVABA [33] to generate SV calls, and with BWA [34] and  
218 FreeBayes [35] to generate SNV calls. Long ranger ([https://support.10xgenomics.com/genome-](https://support.10xgenomics.com/genome-exome/software/pipelines/latest/what-is-long-ranger)  
219 [exome/software/pipelines/latest/what-is-long-ranger](https://support.10xgenomics.com/genome-exome/software/pipelines/latest/what-is-long-ranger)) was used to generate SNV and SV (only  
220 deletions) calls for 10x reference-based analysis. We noted that  $R_9$  failed to be analyzed by  
221 Long Ranger due to its extremely large  $C_F$ . For SNVs, we benchmarked the calls from three  
222 strategies using the gold standard of NA12878 ([ftp://ftp-](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/)  
223 [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  
224 ([ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002\\_NA24385\\_](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002_NA24385_)  
225 [son/latest/GRCh38/](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/AshkenazimTrio/HG002_NA24385_)). For SVs, we compared three linked-read sets ( $R_9$ ,  $R_{10}$ ,  $R_{11}$ ) from HG002  
226 with the Tier 1 SV benchmark from Genome in a Bottle [36] and used VaPoR [37] to validate our  
227 SV calls based on PacBio CCS reads from NA24385 (Highly-accurate long-read sequencing  
228 improves variant detection and assembly of a human genome). We compared SNV and SV  
229 calls among the different approaches using vcfEval [38] and truvari [36], respectively.

230

231 **Performance of diploid assembly: influence of total coverage** Diploid assembly by Linked-  
232 Reads requires sufficient total read coverage ( $C=C_R \times C_F$ ) to generate long contigs and scaffolds.  
233 In this experiment, to explore the roles of both physical coverage ( $C_F$ ) and per-fragment read  
234 coverage ( $C_R$ ), we first generated eight simulated libraries whose total coverage  $C$  ranged from  
235 16x to 78x: four with  $C_R$  fixed and increasing  $C_F$  and four with fixed  $C_F$ , and increasing  $C_R$  (**Table**  
236 **S2**). Contig and scaffold N50s increased along with increasing either  $C_F$  or  $C_R$  (**Figure 2A** and  
237 **2B**). To investigate whether the trend was also present in the real datasets, we analyzed six real

238 libraries (three by varying  $C_F$ , and the other three by varying  $C_R$ ; **Figure 1**): as  $C$  increased, we  
239 varied  $C_F$  and  $C_R$  independently by fixing the other parameter. Contig and scaffold N50s also  
240 increased in these simulation (**Figure 2C** and **2D**) and real linked-read sets (**Figure 2E** and **2F**)  
241 as a function of total coverage  $C$ . Contig lengths did increase a little (621.4kb to 758.1kb for  
242 simulation; 110.7kb to 119.6kb for real data) when  $C$  was increased beyond 56X. Accuracy,  
243 which we define as the ratio between NA50 (N50 after breaking contigs or scaffolds at assembly  
244 errors) and N50 (**Figure 2C** and **2E**), changed 18% for simulation and 7% for real data (587.5kb  
245 to 713.3kb for simulation; 97.1kb to 104.5kb for real data). For scaffolds in the real data sets,  
246 when  $C$  increased from 48X ( $R_3$ ) to 67X ( $R_4$ ), both scaffold N50 and NA50 were significantly  
247 improved (N50: 13.4Mb to 30.6Mb; NA50: 6.3Mb to 12.0Mb), but the accuracy dropped slightly  
248 from 46.6% to 39.1%, which indicated that scaffold accuracy may be refractory to extremely  
249 high  $C$  (**Figure 2F**). These results indicated that assembly length and accuracy were  
250 comparable over a broad range of  $C_F$  and  $C_R$  at constant  $C$ , which implied that assembly quality  
251 was mainly determined by  $C$ .

252  
253 **Performance of diploid assembly: influence of fragment length and physical coverage.** To  
254 investigate if input weighted fragment length (as measured by  $W_{\mu_{FL}}$ ) influenced assembly  
255 quality, we generated four simulated libraries (**Table S3**) with fixed  $C_F$  and  $C_R$  and a range of  
256 fragment lengths (**Figure 3A**). Contig length decreased with increasing fragment length, a trend  
257 that was also seen in six real libraries (**Figure 3B**;  $C=56X$ ;  $R_6$  to  $R_{11}$  in **Figure 1**). We then  
258 simulated another six libraries with the same parameters as the real ones to explore the effects  
259 of physical coverage at constant  $C=56x$  (**Figure 3C**). Contig lengths decreased as a function of  
260 increasing physical coverage, a trend that is somewhat less clear in real data possibly due to  
261 confounding other parameters such as fragment length (**Figure 3D**). The two linked-read sets

262 with the worst contig qualities in NA12878 ( $R_7$ ) and NA24385 ( $R_{10}$ ) also showed a significant  
263 increase of the number of breakpoints (**Table S4**)

264  
265 **Performance of diploid assembly: nature of the source genome.** Assembly errors may  
266 occur because of heterozygosity, repetitive sequences, or sequencing error. To illuminate  
267 possible sources of assembly error, we performed simulations by generating 10x-like Linked-  
268 Reads as above from human chromosome 19, and then quantified assembly error against these  
269 synthetic gold standards. Removal of interspersed repeat sequences from the source genome  
270 resulted in better contigs with no loss of accuracy in experiments by varying  $C_F$ ,  $C_R$  and  
271  $\mu_{FL}$  (**Figure 4A, 4C and 4E**) and better scaffolds only if  $C_R$  was above 1X (**Figure 4D**). Removal  
272 of variation had little effect on contigs and only gave rise to longer scaffolds if  $C_R$  was above  
273 0.8X (**Figure S11**), which is difficult to achieve with real libraries. Finally, a 1% uniform  
274 sequencing error had no discernible effect (**Figure S12**).

275  
276 **Performance of diploid assembly: fraction of genome in diploid state.** While contiguity is  
277 an important parameter for any whole genome assembly, evaluation of diploid assemblies  
278 necessitates estimating the fraction of the genome in which the assembly recovered the diploid  
279 state. To this end, we divided the contigs generated by Supernova2 into “diploid contigs”, which  
280 were extracted from its megabubble structures, and “haploid contigs” from non-megabubble  
281 structures. Pairs of scaffolds were extracted as the two haplotypes from megabubble structures  
282 if they shared the same start and end nodes in the assembly graph. Diploid contigs were  
283 generated by breaking the candidate scaffolds at the sequences with least 10 consecutive ‘N’s  
284 and were aligned to human reference genome (hg38) by Minimap2. The genome was split into  
285 500bp windows and diploid regions were defined as the maximum extent of successive  
286 windows covered by two contigs, each from one haplotype. Alignment against the human  
287 reference genome revealed the overall genome coverages of the six assemblies to be around

288 91%. For most assemblies, 70%-80% of the genome was covered by two homologous contigs  
289 (**Table 1**), with  $R_6$  only reaching 58.9%, probably due to the short fragments of the DNA  
290 preparation ( $\mu_{FL}=24\text{kb}$ ). We also analyzed another seven assemblies produced by 10x  
291 Genomics, all of which had diploid fractions of about 80% as well (**Table S5**). In the male  
292 NA24385, non-pseudoautosomal regions of the X chromosome are hemizygous and should  
293 therefore be recovered as haploid regions. Between 79.9% and 87.6% of these regions were  
294 covered by one contig exactly depending on the assembled library. Library construction  
295 parameters other than fragment length appeared to have had little impact on the proportion of  
296 diploid regions (**Tables 1** and **Table S5**).

297  
298 Overlapping the diploid regions from the assemblies of the same individual revealed that 50.24%  
299 and 67.27% of the genome for NA12878 and NA24385 (**Figure S13**), respectively, were diploid  
300 in all the three assemblies. NA12878 was lower because of the low percentage of diploid  
301 regions in assembly  $R_6$  (**Table 1**). The overlaps were significantly greater than expected by  
302 chance (NA12878: 33.3%, p-value=0.0049; NA24385: 45.4%, p-value=0.0029. Chi square test).  
303 These observations were consistent with heterozygous variants being enriched in certain  
304 genomic segments, in which two haplotypes were more easily differentiated by Supernova2.  
305 Phase block lengths were mainly determined by total coverage  $C$  and increased in real data  
306 with increasing fragment length (**Figure S14, Table S6**).

307  
308 **Performance of diploid assembly: quality of variant calls.** The ultimate goal of human  
309 genome assembly is to accurately identify genomic variants. We compared the SNVs and SVs  
310 from our assemblies with the calls from referenced-based processing of standard Illumina and  
311 10x data, and benchmarked them using gold standard from Genome in a Bottle and PacBio  
312 CCS reads. We found the SNVs from referenced-based processing of standard Illumina and  
313 10x data were comparable and both of them were better than assembly-based calls (**Table S7**

314 and **S8**) For SVs, our assemblies generated many calls that were missed by the reference-  
315 based strategy (**Table S9-S12**) and even by the Tier 1 benchmark of Genome in a Bottle (**Table**  
316 **S13**), and half of the deletions and a majority of insertions could be validated by PacBio CCS  
317 reads (**Table S14**).

318

## 319 **Discussion**

320 In this study, we investigated human diploid assembly using 10x Linked-Read sequencing data  
321 on both simulated and real libraries. We developed the simulator LRTK-SIM to examine the  
322 likely impact of parameters in diploid assembly and compared results from simulated reads to  
323 those from real libraries. We thus determined the impact of key parameters ( $C_R$ ,  $C_F$ ,  $N_{FP}$  and  
324  $\mu_{FL}/W\mu_{FL}$ ) with respect to assembly continuity and accuracy. Our study provides a general  
325 strategy to evaluate assemblies of 10x data and may have implications for the evaluation of  
326 other barcode-based sequencing technologies such as CPTv2-seq [39] or stLRF [40] in the  
327 future.

328

## 329 **10x Practicalities**

330 For standard Illumina sequencing, library complexity is usually sufficient to generate  
331 tremendous numbers of reads from unique templates and read coverage can be increased  
332 simply by sequencing more. However, the 10x Chromium system performs amplification in each  
333 partition, and generally only about 20% to 40% of the original long fragment sequence can be  
334 captured as short fragments and eventually as reads, resulting in shallow sequencing coverage  
335 per fragment. Sequencing more deeply does not increase the per-fragment coverage much as  
336 most of the extra reads are from PCR duplicates. The solution is to sequence multiple 10x  
337 libraries constructed from the same DNA preparation and merge them for analysis. This means  
338 that  $C_R$  remains in the standard range where PCR duplicates are relatively rare, but  $C_F$

339 increases proportionally to the number of libraries used. A practical limitation to this approach is  
340 that Supernova2 limits the number of barcodes to 4.8 million.

341  
342 Our results showed that in practice,  $C_F$  should be between 335X and 823X, but no larger than  
343 1000X, given the optimal coverage of  $C=56X$  recommended by 10x and the requirement for  
344 sufficient per-fragment read coverage. Surprisingly, we observed that including more extremely  
345 long fragments was detrimental for assembly quality. This is possibly due to the loss of barcode  
346 specificity for fragments spanning repetitive sequences. From a computational perspective, too  
347 many long fragments are harmful to deconvolving the *de bruijn* graph, as more complex paths  
348 need to be picked out. In our experiments,  $W_{\mu_{FL}}$  between 50kb and 150kb is the best choice to  
349 generate reliable assemblies.

350

### 351 **Parameters driving assembly quality**

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

365

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

373

### 374 **Cost-benefit analysis**

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

383

### 384 **Comparison with hybrid assemblies**

385 Hybrid assembly strategies have been applied successfully to produce human genome  
386 assembly of long contiguity [13, 14, 41]. In these studies, long contigs are first produced by  
387 single-molecule long-reads, such as PacBio (NG50=1.1Mb; [13]) or Nanopore (NG50=3.21Mb;  
388 [14]) comparing favorably to our best results for Linked-Reads assemblies (NG50=236kb).  
389 Scaffolding is then performed with complementary technologies such as BioNano to capture



390 chromosomal level long-range information. It promoted the scaffold N50 of PacBio to 31.1Mb  
391 [13] and Illumina mate-pair sequencing with 10x data to 33.5Mb [25]. Using SuperNova2, the  
392 scaffold N50 from our studies reached  $\sim 27.86\text{Mb}$  ( $R_6$ ) on the basis of 10x data alone,  
393 suggesting that 10x technology gives broadly comparable results at a fraction of the price of  
394 long-read-based hybrid assemblies.  
395

## 396 **Availability of supporting data**

397 The raw sequencing data are deposited in the Sequence Read Archive and the corresponding  
398 BioProject accession number is PRJNA527321. Diploid assemblies and the codes for  
399 comparison are currently available at  
400 [http://mendel.stanford.edu/supplementarydata/zhang\\_SN2\\_2019](http://mendel.stanford.edu/supplementarydata/zhang_SN2_2019) and  
401 [https://github.com/zhanglu295/Evaluate\\_diploid\\_assembly](https://github.com/zhanglu295/Evaluate_diploid_assembly). LRTK-SIM is publicly available at  
402 <https://github.com/zhanglu295/LRTK-SIM>.

403

## 404 **Additional files**

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

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

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

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

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

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

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

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

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

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

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

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

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

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

432 **Figure S1. Basic statistics for  $L_{1L}$ .** The distributions of **A.** the number of fragments per  
433 partition; **B.** sequencing depth per fragment; **C.** probability density function of unweighted  
434 fragment lengths; **D.** cumulative density function of unweighted fragment lengths; **E.** reversed  
435 cumulative density function of unweighted fragment lengths; **F.** reversed cumulative density  
436 function of weighted fragment lengths.

437 **Figure S2. Basic statistics for  $L_{1M}$ .** The distributions of **A.** number of fragments per partition;  
438 **B.** sequencing depth per fragment; **C.** probability density function of unweighted fragment  
439 lengths; **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative  
440 density function of unweighted fragment lengths; **F.** reversed cumulative density function of  
441 weighted fragment lengths.

442 **Figure S3. Basic statistics for  $L_{1H}$ .** The distributions of **A.** number of fragments per partition; **B.**  
443 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
444 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
445 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
446 fragment lengths.

447 **Figure S4. Basic statistics for  $L_2$ .** The distributions of **A.** number of fragments per partition; **B.**  
448 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
449 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
450 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
451 fragment lengths.

452 **Figure S5. Basic statistics for  $L_3$ .** The distributions of **A.** number of fragments per partition; **B.**  
453 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
454 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
455 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
456 fragment lengths.

457 **Figure S6. Basic statistics for  $L_4$ .** The distributions of **A.** number of fragments per partition; **B.**  
458 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
459 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density

460 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
461 fragment lengths.

462 **Figure S7. Basic statistics for  $L_5$ .** The distributions of **A.** number of fragments per partition; **B.**  
463 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
464 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
465 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
466 fragment lengths.

467 **Figure S8. Basic statistics for  $L_6$ .** The distributions of **A.** number of fragments per partition; **B.**  
468 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
469 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
470 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
471 fragment lengths.

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

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

475 **Figure S11.** Comparison of assembly qualities from 10x data with and without single  
476 nucleotide 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  
477 to 300X in **C** and **D**;  $C_R$  was fixed 0.2X and  $C_F$  was fixed 300X in **E** and **F**.

478 **Figure S12.** Comparison of assembly qualities from 10x data with (1% uniform) and without  
479 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  
480 300X in **C** and **D**;  $C_R$  was fixed 0.2X and  $C_F$  was fixed 300X in **E** and **F**.

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

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

490

491

## 492 **Competing interest**

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

494

## 495 **Author Contributions**

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

499

## 500 **Acknowledgements**

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

504

505

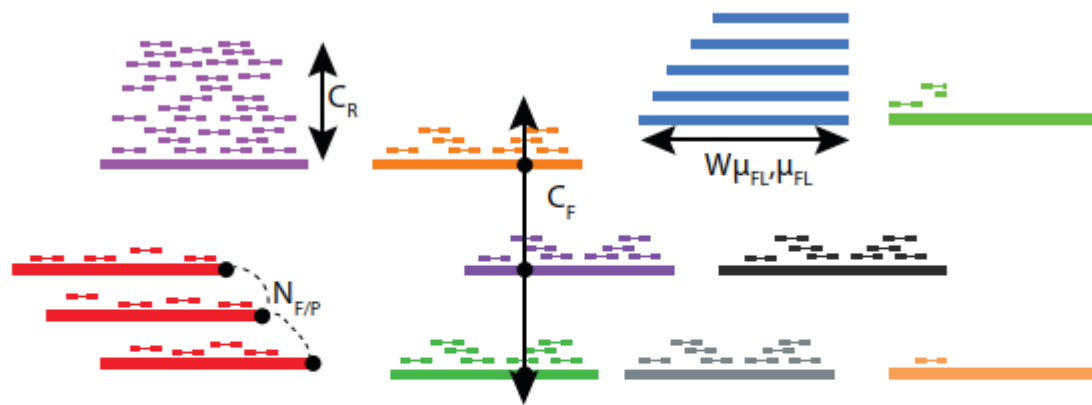
506 **Table**

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

507

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

511



**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**

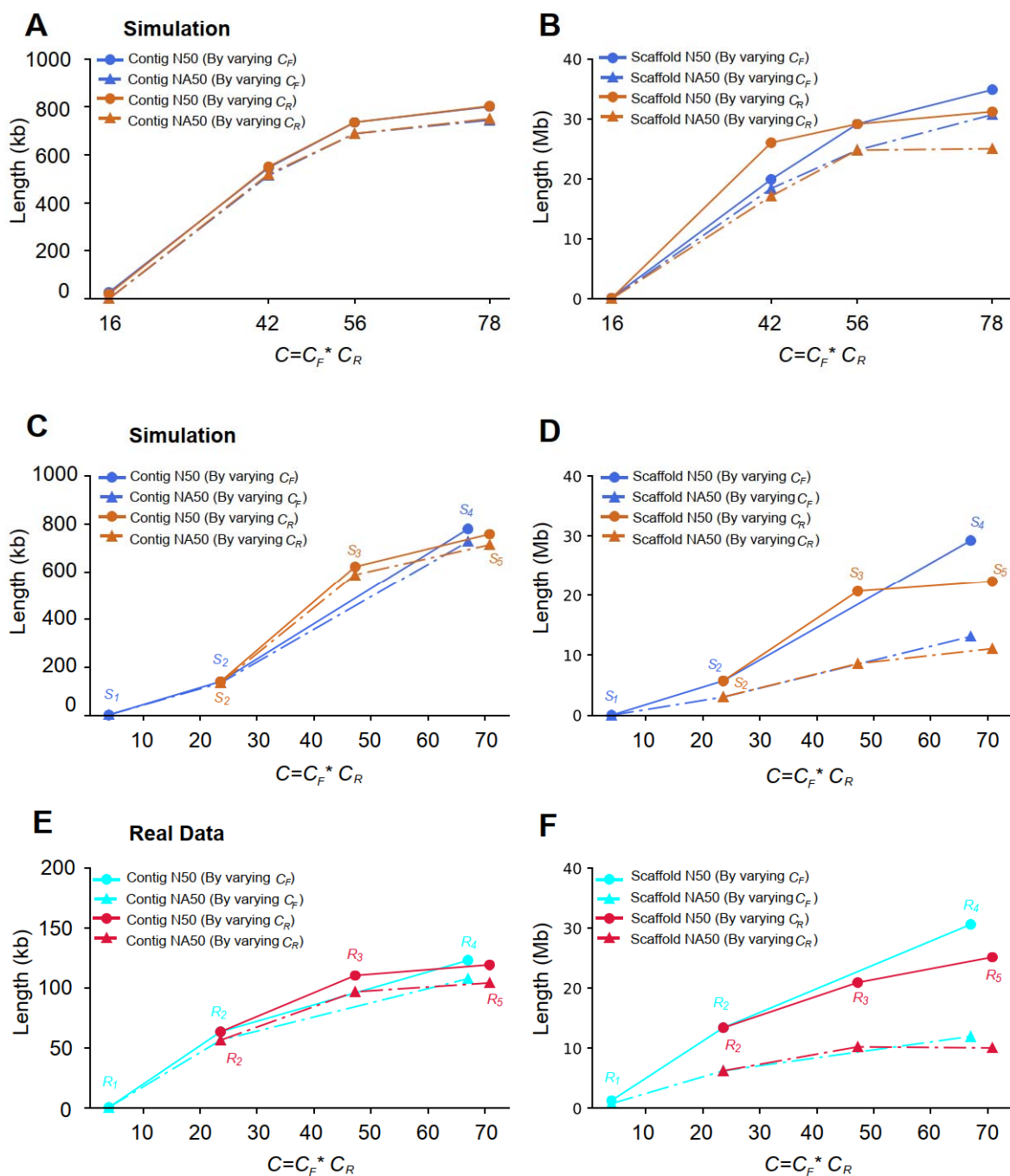
10 - 100  
 $\mu_{FL}$  = 10-100kb  
 $W\mu_{FL}$  = 20-400kb  
 $C_R$  = 0.1x - 0.4x  
 $C_F$  = 200x - 1000x  
 $C = C_R * C_F = 40x - 80x$

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	214.5/188.3	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

512 **Figures**

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

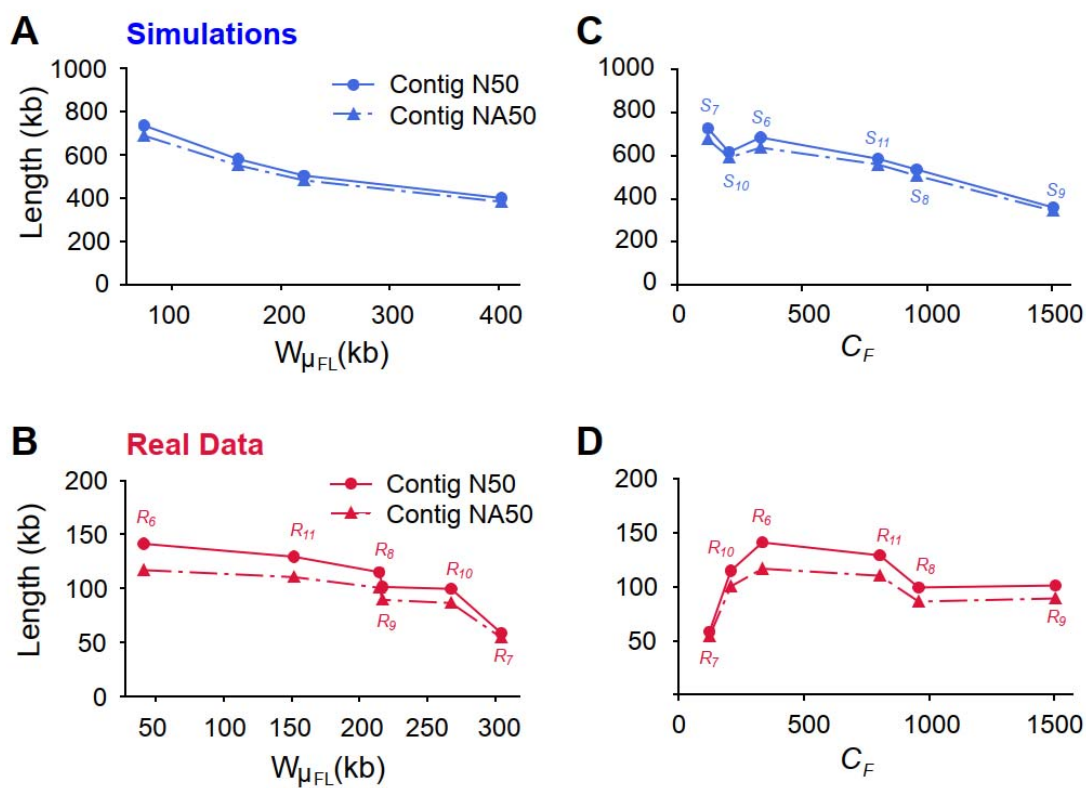
514 human diploid assembly.



515

516 **Figure 2.** Contig and scaffold lengths (N50 and NA50) as a function of  $C_F$  or  $C_R$ . **A and B:**  
 517 Simulated Linked-Reads with predefined parameters (**Table S2**); **C and D:** Simulated Linked-  
 518 reads with matched parameters of real Linked-Read data sets (**Figure 1**); **E and F:** Real linked-  
 519 read sets (**Figure 1**).

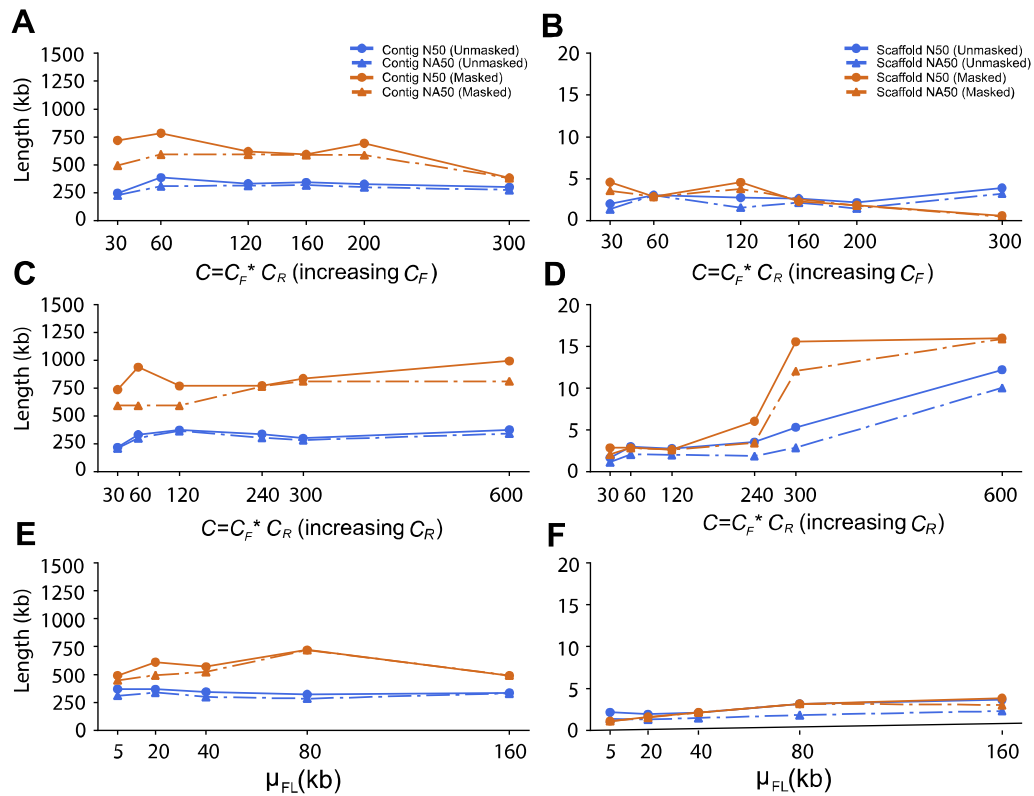




520

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

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



523

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

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

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

527

## 528 References

- 529 1. Metzker ML. Sequencing technologies - the next generation. *Nat Rev Genet.* 2010;11  
530 1:31-46. doi:10.1038/nrg2626.
- 531 2. Shendure J, Balasubramanian S, Church GM, Gilbert W, Rogers J, Schloss JA, et al.  
532 DNA sequencing at 40: past, present and future. *Nature.* 2017;550 7676:345-53.  
533 doi:10.1038/nature24286.
- 534 3. Head SR, Komori HK, LaMere SA, Whisenant T, Van Nieuwerburgh F, Salomon DR, et  
535 al. Library construction for next-generation sequencing: overviews and challenges.  
536 *Biotechniques.* 2014;56 2:61-4, 6, 8, passim. doi:10.2144/000114133.
- 537 4. O'Connell J, Sharp K, Shrine N, Wain L, Hall I, Tobin M, et al. Haplotype estimation for  
538 biobank-scale data sets. *Nat Genet.* 2016;48 7:817-20. doi:10.1038/ng.3583.
- 539 5. Delaneau O, Zagury JF and Marchini J. Improved whole-chromosome phasing for  
540 disease and population genetic studies. *Nat Methods.* 2013;10 1:5-6.  
541 doi:10.1038/nmeth.2307.
- 542 6. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, et al. A general  
543 approach for haplotype phasing across the full spectrum of relatedness. *PLoS Genet.*  
544 2014;10 4:e1004234. doi:10.1371/journal.pgen.1004234.
- 545 7. Roach JC, Glusman G, Hubley R, Montsaroff SZ, Holloway AK, Mauldin DE, et al.  
546 Chromosomal haplotypes by genetic phasing of human families. *Am J Hum Genet.*  
547 2011;89 3:382-97. doi:10.1016/j.ajhg.2011.07.023.
- 548 8. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, et al. Efficient de  
549 novo assembly of highly heterozygous genomes from whole-genome shotgun short reads.  
550 *Genome Res.* 2014;24 8:1384-95. doi:10.1101/gr.170720.113.
- 551 9. Alkan C, Sajjadian S and Eichler EE. Limitations of next-generation genome sequence  
552 assembly. *Nat Methods.* 2011;8 1:61-5. doi:10.1038/nmeth.1527.
- 553 10. Treangen TJ and Salzberg SL. Repetitive DNA and next-generation sequencing:  
554 computational challenges and solutions. *Nat Rev Genet.* 2011;13 1:36-46.  
555 doi:10.1038/nrg3117.
- 556 11. Huddleston J, Ranade S, Malig M, Antonacci F, Chaisson M, Hon L, et al.  
557 Reconstructing complex regions of genomes using long-read sequencing technology.  
558 *Genome Res.* 2014;24 4:688-96. doi:10.1101/gr.168450.113.
- 559 12. Lu H, Giordano F and Ning Z. Oxford Nanopore MinION Sequencing and Genome  
560 Assembly. *Genomics Proteomics Bioinformatics.* 2016;14 5:265-79.  
561 doi:10.1016/j.gpb.2016.05.004.
- 562 13. Pendleton M, Sebra R, Pang AW, Ummat A, Franzen O, Rausch T, et al. Assembly and  
563 diploid architecture of an individual human genome via single-molecule technologies.  
564 *Nat Methods.* 2015;12 8:780-6. doi:10.1038/nmeth.3454.
- 565 14. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, et al. Nanopore sequencing  
566 and assembly of a human genome with ultra-long reads. *Nat Biotechnol.* 2018;36 4:338-  
567 45. doi:10.1038/nbt.4060.
- 568 15. Wenger AM, Peluso P, Rowell WJ, Chang P-C, Hall RJ, Concepcion GT, et al. Highly-  
569 accurate long-read sequencing improves variant detection and assembly of a human  
570 genome. *bioRxiv.* 2019.
- 571 16. Cao H, Wu H, Luo R, Huang S, Sun Y, Tong X, et al. De novo assembly of a haplotype-  
572 resolved human genome. *Nat Biotechnol.* 2015;33 6:617-22. doi:10.1038/nbt.3200.

- 573 17. Kuleshov V, Xie D, Chen R, Pushkarev D, Ma Z, Blauwkamp T, et al. Whole-genome  
574 haplotyping using long reads and statistical methods. *Nat Biotechnol.* 2014;32 3:261-6.  
575 doi:10.1038/nbt.2833.
- 576 18. Peters BA, Kermani BG, Sparks AB, Alferov O, Hong P, Alexeev A, et al. Accurate  
577 whole-genome sequencing and haplotyping from 10 to 20 human cells. *Nature.* 2012;487  
578 7406:190-5. doi:10.1038/nature11236.
- 579 19. Edge P, Bafna V and Bansal V. HapCUT2: robust and accurate haplotype assembly for  
580 diverse sequencing technologies. *Genome Res.* 2017;27 5:801-12.  
581 doi:10.1101/gr.213462.116.
- 582 20. Patterson M, Marschall T, Pisanti N, van Iersel L, Stougie L, Klau GW, et al. WhatsHap:  
583 Weighted Haplotype Assembly for Future-Generation Sequencing Reads. *J Comput Biol.*  
584 2015;22 6:498-509. doi:10.1089/cmb.2014.0157.
- 585 21. Zheng GX, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, et al.  
586 Haplotyping germline and cancer genomes with high-throughput linked-read sequencing.  
587 *Nat Biotechnol.* 2016;34 3:303-11. doi:10.1038/nbt.3432.
- 588 22. Spies N, Weng Z, Bishara A, McDaniel J, Catoe D, Zook JM, et al. Genome-wide  
589 reconstruction of complex structural variants using read clouds. *Nat Methods.* 2017;14  
590 9:915-20. doi:10.1038/nmeth.4366.
- 591 23. Bishara A, Moss EL, Kolmogorov M, Parada AE, Weng Z, Sidow A, et al. High-quality  
592 genome sequences of uncultured microbes by assembly of read clouds. *Nat Biotechnol.*  
593 2018; doi:10.1038/nbt.4266.
- 594 24. Weisenfeld NI, Kumar V, Shah P, Church DM and Jaffe DB. Direct determination of  
595 diploid genome sequences. *Genome Res.* 2017;27 5:757-67. doi:10.1101/gr.214874.116.
- 596 25. Mostovoy Y, Levy-Sakin M, Lam J, Lam ET, Hastie AR, Marks P, et al. A hybrid  
597 approach for de novo human genome sequence assembly and phasing. *Nat Methods.*  
598 2016;13 7:587-90. doi:10.1038/nmeth.3865.
- 599 26. Hulse-Kemp AM, Maheshwari S, Stoffel K, Hill TA, Jaffe D, Williams SR, et al.  
600 Reference quality assembly of the 3.5-Gb genome of *Capsicum annuum* from a single  
601 linked-read library. *Hortic Res.* 2018;5:4. doi:10.1038/s41438-017-0011-0.
- 602 27. Elyanow R, Wu HT and Raphael BJ. Identifying structural variants using linked-read  
603 sequencing data. *Bioinformatics.* 2017; doi:10.1093/bioinformatics/btx712.
- 604 28. Jones SJ, Haulena M, Taylor GA, Chan S, Bilobram S, Warren RL, et al. The Genome of  
605 the Northern Sea Otter (*Enhydra lutris kenyoni*). *Genes (Basel).* 2017;8 12  
606 doi:10.3390/genes8120379.
- 607 29. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34  
608 18:3094-100. doi:10.1093/bioinformatics/bty191.
- 609 30. Mikheenko A, Prjibelski A, Saveliev V, Antipov D and Gurevich A. Versatile genome  
610 assembly evaluation with QUAST-LG. *Bioinformatics.* 2018;34 13:i142-i50.  
611 doi:10.1093/bioinformatics/bty266.
- 612 31. Earl D, Bradnam K, St John J, Darling A, Lin D, Fass J, et al. Assemblathon 1: a  
613 competitive assessment of de novo short read assembly methods. *Genome Res.* 2011;21  
614 12:2224-41. doi:10.1101/gr.126599.111.
- 615 32. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, et al. Extensive sequencing  
616 of seven human genomes to characterize benchmark reference materials. *Sci Data.*  
617 2016;3:160025. doi:10.1038/sdata.2016.25.

- 618 33. Wala JA, Bandopadhyay P, Greenwald NF, O'Rourke R, Sharpe T, Stewart C, et al.  
619 SvABA: genome-wide detection of structural variants and indels by local assembly.  
620 *Genome Res.* 2018;28 4:581-91. doi:10.1101/gr.221028.117.
- 621 34. Li H and Durbin R. Fast and accurate short read alignment with Burrows-Wheeler  
622 transform. *Bioinformatics.* 2009;25 14:1754-60. doi:10.1093/bioinformatics/btp324.
- 623 35. Garrison E and Marth G. Haplotype-based variant detection from short-read sequencing.  
624 arXiv e-prints. 2012.
- 625 36. Zook JM, Hansen NF, Olson ND, Chapman LM, Mullikin JC, Xiao C, et al. A robust  
626 benchmark for germline structural variant detection. *bioRxiv.* 2019.
- 627 37. Zhao X, Weber AM and Mills RE. A recurrence-based approach for validating structural  
628 variation using long-read sequencing technology. *Gigascience.* 2017;6 8:1-9.  
629 doi:10.1093/gigascience/gix061.
- 630 38. Krusche P, Trigg L, Boutros PC, Mason CE, De La Vega FM, Moore BL, et al. Best  
631 practices for benchmarking germline small-variant calls in human genomes. *Nat*  
632 *Biotechnol.* 2019;37 5:555-60. doi:10.1038/s41587-019-0054-x.
- 633 39. Zhang F, Christiansen L, Thomas J, Pokholok D, Jackson R, Morrell N, et al. Haplotype  
634 phasing of whole human genomes using bead-based barcode partitioning in a single tube.  
635 *Nat Biotechnol.* 2017;35 9:852-7. doi:10.1038/nbt.3897.
- 636 40. Wang O, Chin R, Cheng X, Wu MKY, Mao Q, Tang J, et al. Efficient and unique  
637 cobarcodeing of second-generation sequencing reads from long DNA molecules enabling  
638 cost-effective and accurate sequencing, haplotyping, and de novo assembly. *Genome Res.*  
639 2019;29 5:798-808. doi:10.1101/gr.245126.118.
- 640 41. Ma ZS, Li L, Ye C, Peng M and Zhang YP. Hybrid assembly of ultra-long Nanopore  
641 reads augmented with 10x-Genomics contigs: Demonstrated with a human genome.  
642 *Genomics.* 2018; doi:10.1016/j.ygeno.2018.12.013.  
643