

<b>Manuscript Number:</b>	GIGA-D-19-00092R1
<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>	
<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.</p> <p>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.</p> <p>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.</p> <p>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:</p> <p>Reviewer #1: Zhang et al. explore the parameter space of 10X libraries and the subsequent effects of those parameters on de novo assembly performance. They also developed an in silico simulator and that generates results similar to experimental findings. The manuscript is well written and easy to understand.</p>

We thank the reviewer for these positive comments and address each point below.

That said, I think there are some analyses missing that should be included:

1. I think you should variant call off of the de novo assemblies to see if there are any differences you are missing because you're only looking at things at a very high structural level.

We have now called SNVs and SVs from our de novo assemblies and from other methods. Please find our results in the responses to points 2-4 of reviewer2.

2. How is phasing affected? I don't see any data on that other than total diploid regions. You should include the changes to the phase block N50. It's mentioned in the abstract, but I don't see it anywhere else.

We have showed the trend of phased block N50 in different linked-read sets in Figure S14, now we also provided the values of phase block N50s in Table S6

3. Besides NA50 you should include assembly errors such as breakpoints, translocations, inversions, relocations, etc.....

You have a nice dataset here, you should try to get more out of it.

Thank you for the suggestions. We have re-run QUAST and generated several detailed statistics which are now shown in Table S4. These results are consistent with the contig N50s reported in Figure 3.

Minor comments:

58-66, Probably should add this reference for PacBio CCS sequencing, contig N50 is 15 mb, <https://www.biorxiv.org/content/10.1101/519025v2>

We have added this reference

65-66, I'd argue that this statement is a bit strong, cost is lowering, and throughput is increasing for these systems

This is now lines 70-72. We have rephrased the sentence and now write: "However, long-fragment sequencing suffers from extremely high cost (in the case of PacBio CCS), or low base quality (in the case of single-pass reads of either technology), hampering its usefulness for personal genome assembly."

68 Not a complete sentence

We fixed this

Ref 27 isn't our stLFR paper, the doi for that is 10.1101/gr.245126.118, and it is commercially available now in some parts of the world

We have added the new reference and deleted the confusing words in this sentence.

Reviewer #2: Zhang and co-authors present a parameter study for 10x linked-read sequencing experiments with the objective of evaluating the influence of experimentally controllable parameters on the final diploid assembly quality. The authors perform basic performance evaluation in terms of common metrics such as N50 values and provide technical recommendations for designing linked-read sequencing experiments. Additionally, Zhang et al. implemented a software tool for simulating linked-read sequencing data, which they use for parameter assessment given the known (simulated) truth.

While such studies that provide guidance to users of a sequencing technology are very valuable in principle, I have a number of concerns that should be addressed:

1. There is a closely related article by Luo et al. (2017, DOI: 10.1016/j.csbj.2017.10.002) that has been missed. The authors should clarify what the added value of their study is beyond the work by Luo et al. This comment applies to both aspects: guidance to users in terms of 10x sequencing experiments and the utility/features of their data simulation tool (note that Luo et al. also provide a simulator).

We appreciate and cite the work by Luo et al. However, our study provides (1) a more flexible simulation tool and (2) an extensive set of new sequence data.

Regarding (1)

A. We explicitly allow users to input CF, CR,  $W\mu_{FL}$  and  $\mu_{FL}$ , which have strong connections with library preparation and Illumina sequencing. For example, CF is driven by input DNA amount and  $\mu_{FL}$  by DNA preparation and potential size selection. LRSIM only lets the user set the total number of reads.

B. The usability of LRTK-SIM is better than LRSIM. LRSIM requires many third party packages and software to be installed first, such as Inline::C perl library, DWGSIM etc. It is not convenient for the users with insufficient computer experience. LRTK-SIM was written in Python and no third-party software was required. It can be installed and gotten started easily. LRTK-SIM can parallel simulate multiple libraries with a variety of parameters simultaneously. The users can compare the performance of different parameters in one run.

Regarding (2)

Luo et al. compared the influence of different parameters by simulation only, which does not always reflect the situation in real sequencing. In our study, we prepared six real libraries with different parameters and could validate our observations from simulation data.

2. The focus of this manuscript is on guiding researchers who are after a cost-effective characterization of individual human genomes. In my view, Zhang et al. should go the full distance and additionally compare to standard Illumina sequencing followed by mapping and variant calling as a baseline. The assembly metrics employed are not so very informative when it comes to the question of which variation (relative to the reference genome) is been missed/captured in standard approaches.

While human assembly is the focus, we believe that much of the interest in our work will come mainly from researchers who are interested in assembling novel genomes. We use human as an assembly model because assembly quality can be gauged by comparison to the reference sequence. Nonetheless ...

Beyond comparing to standard Illumina sequencing, including a detailed comparison to reference-based processing of 10x data (e.g. using LongRanger) would be interesting. In this way, this study would be much more helpful for planning sequencing studies.

... in response to this comment, we now systematically investigate SNV and SV calls from our assemblies. We compare with standard Illumina data and reference-based processing of our 10x data. The standard Illumina data were downloaded from Genome In A Bottle and analyzed with SVABA to generate SV calls, and with BWA and FreeBayes to generate SNV calls. Long ranger was used to generate SNVs and SVs (only deletions) for 10x reference-based analysis. We noted that R9 failed to be analyzed by Long Ranger due to its extremely large CF. We compared SNV and SV calls among the different approaches using vcfEval (<https://github.com/RealTimeGenomics/rtg-tools>) and truvari (<https://github.com/spiralgenetics/truvari>), respectively.

For SNVs, we compared the calls from three strategies to the gold standard of NA12878 ([ftp://ftp-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 ([ftp://ftp-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/)).

We found that SNVs from reference-based processing of Illumina and 10x data were comparable, and both of them were better than assembly-based SNV calls. For SVs, our assemblies generated many calls that were missed by the reference-based strategy.

We now provide several additional supplementary tables (Table S7-S12) to present these results.

3. The main reason (in my view) for pursuing de novo assembly of human genomes is to access structural variation that is missed otherwise. An evaluation on how much structural variation is (accurately) captured would be of interest to many readers. This is actually something that the authors point out in the Discussion themselves: "Arguably, the metric that matters most in the context of a personal genome is the discovery of variation that lower-cost approaches do not enable."

As implied by the quote, we agree with the reviewer's comment. Consequently, we now compare three linked-read sets from HG002 with the Tier 1 SV benchmark from Genome in a Bottle by using *truvari* (<https://github.com/spiralgenetics/truvari>). The results are summarized in Table S13.

4. PacBio CCS reads are available for HG002 (see Wenger et al., <http://dx.doi.org/10.1101/519025>). Mapping those CCS reads back to your diploid assemblies and calling variants provides an easy and powerful opportunity to assess the sequence quality from an independent technology.

These data became available while our manuscript was in review. We note that the PacBio CCS calls on HG002 are generally reasonably accurate but are not guaranteed to be correct in the absence of a gold standard. Therefore, we prefer to compare them in an overlap analysis with our calls, as opposed to implying that they are a gold standard by using the term "validation". We used *vapor* (<https://github.com/mills-lab/vapor>) to validate our SV calls based on PacBio CCS reads from HG002 and include Table S14 to show the validation rates.

Beyond this, your evaluation could be improved by also adding an assembly evaluation perspective that is more biologically motivated, e.g., number of recovered genes/disrupted genes or similar (this should be supported by Quast-LG/BUSCO).

We have added this analysis in Table S4.

#### Minor comments

- line 51: pedigree based phasing is quite powerful even for trios (where it is able to phase all variants that are homozygous in at least one individual), so I disagree to the statement that this is only feasible in large pedigrees.

We fixed this and removed confusing words.

- lines 60ff: it is unclear which study you are referring to here, please add the citation at the end of the sentence (N50 31.1Mb)

We included a new reference here.

- line 68: broken sentence; also, putting the citation at the end of the sentence increases readability

We fixed this issue.

- lines 71/72: again, unclear which study you are referring to ("Long Fragment Read")

We included a new reference here.

- lines 125ff: is there a specific reason why five and three? (And not, e.g., five and five?) Also, the meaning of L, M, and H in the subscript of L should be explained. Because we generated two additional libraries (L\_1L and L\_1M for NA12878) to evaluate the effects of CF and CR in assembly, and we believe the trend should be consistent in the two samples. L, M and H represent low, medium and high CF in the experiments. We have clarified this in the manuscript.

- line 129: percent of what?

The percent of GEM in 10x Chromium system.

- line 151: please be more specific about which version of hg38 was used (detail once if identical hg38 was used throughout the rest of the paper [lines 165, 171, 195 and so on...])

The reference was downloaded from 10x website with the version of GRCh38

	<p>Reference 2.1.0.</p> <p>- line 172: please provide an exact reference for the high confidence regions that you used (e.g., file URL) We have added the URL in the manuscript.</p> <p>- line 208: "in in" We fixed this.</p> <p>- line 208: this sentence is talking about real data, so the reference to Fig 2C and 2D does not match. We clarified this in the manuscript.</p> <p>- line 209: "...but not dramatically... [...] ...appreciably" - this is subjective language, please rephrase and be more fact-oriented (for instance by including the numbers you refer to in parentheses). We included the numbers and rephrased the sentence to be more fact-oriented.</p> <p>- line 250: "_Alignment" ? We fixed this.</p> <p>- line 251: what is the denominator for these 91% all bases that are not Ns in the reference genome? (Note that for this analysis, the version of hg38 matters, see comment above). "N"s do not contribute to the denominator.</p> <p>- The authors mention stLFR in line 278. There's a new preprint that's worth citing/discussing: <a href="http://dx.doi.org/10.1101/324392">http://dx.doi.org/10.1101/324392</a> We have cited their latest version.</p> <p>- line 296: "extremely long" please say what extremely long means here We defined "extremely long" as the DNA fragments longer than 200kb.</p> <p>- line 570: please be more specific what you mean by "in-house programs", and where the respective sources are available (is that the "Evaluate_diploid_assembly" github?) All the source codes for assembly evaluation are available in <a href="https://github.com/zhanglu295/Evaluate_diploid_assembly">https://github.com/zhanglu295/Evaluate_diploid_assembly</a>. We added this information in the sentence.</p> <p>- please add a - preferably open source - license file to your github repositories We added the license files in the GitHub.</p> <p>- "sample prep" is jargon and should be replaced by "sample preparation" (eg. line 41, but also elsewhere) We have updated all the "sample prep" to "sample preparation" in the manuscripts.</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<b>Experimental design and statistics</b>	Yes
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>Have you included all the information requested in your manuscript?</p>	
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</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

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)

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 DNA  
25 fragments ( $C_F$ ) or read coverage per fragment ( $C_R$ ) within broad ranges. The optimal physical  
26 coverage was between 332X and 823X and assembly quality worsened if it increased to greater  
27 than 1,000X for a given  $C$ . Long DNA fragments could significantly extend phase blocks, but  
28 decreased contig contiguity. The optimal length-weighted fragment length ( $W\mu_{FL}$ ) was around 50  
29 – 150kb. When broadly optimal parameters were used for library preparation and sequencing, ca.  
30 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 foundational  
41 for the study of human genomic variation and analysis of genotype-phenotype relationships. Over  
42 the past decades, cost-effective whole-genome sequencing has been revolutionized by short-  
43 fragment approaches, the most widespread of which have been the consistently improving  
44 generations of the original Solexa technology [1, 2], now referred to as Illumina sequencing.  
45 Illumina's strengths and weaknesses are inherent in the sample preparation and sequencing  
46 chemistry. Illumina generates short paired reads (2x150 base pairs for the highest-throughput  
47 platforms) from short fragments (usually 400-500 base pairs) [3]. Because many clonally amplified  
48 molecules generate a robust signal during the sequencing reaction, Illumina's average per-base  
49 error rates are very low.

50

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

61

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

73

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

84

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

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

102  
103 Practical advantages of the technology include the low DNA input mass requirement (1ng per  
104 library, or approximately 300 haploid human genome equivalents). Real input quantities can vary,  
105 along with other factors, to influence an interconnected array of parameters that are relevant to  
106 genome assembly and reconstruction. The parameters over which the experimenter has influence  
107 are (**Figure 1**): i).  $C_R$ : average **C**overage of short **R**eads per fragment; ii).  $C_F$ : average physical  
108 **C**overage of the genome by long DNA **F**ragments; iii).  $N_{FP}$ : **N**umber of **F**ragments per **P**artition;  
109 iv). Fragment length distribution, several parameters of which are used, specifically  $\mu_{FL}$ : Average  
110 Unweighted DNA **F**ragment **L**ength and  $W\mu_{FL}$ : Length-**W**eighted average of DNA **F**ragment  
111 **L**ength. Note that several parameters depend on each other. For example, a greater amount of  
112 input DNA will increase  $N_{FP}$ ; shorter fragments increase  $N_{FP}$  at the same DNA input amount

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

116

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

129

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

131 We made six DNA preparations that varied in fragment size distribution and amount of input DNA,  
132 three each from NA12878 and NA24385. From these, we prepared eight libraries, five from  
133 NA12878 and three from NA24385 (**Table S1**). To generate libraries  $L_{1L}$ ,  $L_{1M}$  and  $L_{1H}$  (the  
134 subscripts  $L$ ,  $M$  and  $H$  represent low, medium and high  $C_F$ , respectively), genomic DNA was  
135 extracted from ca. 1 million cultured NA12878 cells using the Gentra Puregene Blood Kit following  
136 manufacturer's instructions (Qiagen, Cat. No 158467). The GEMs were divided into 3 tubes with  
137 5%, 20%, and 75% to generate libraries  $L_{1L}$ ,  $L_{1M}$  and  $L_{1H}$ , respectively (**Figure S1-S3**). For the

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

154

### 155 **Linked-Reads subsampling**

156 The high sequencing coverage in the libraries allowed subsampling to facilitate the matching of  
157 parameters among the different libraries, for purposes of comparability; these subsampled  
158 Linked-Read sets are denoted  $R_{id}$  (**Figure 1**). We aligned the 10x Linked-Reads to human  
159 reference genome (hg38, GRCh38 Reference 2.1.0 from 10x website) followed by removing PCR  
160 duplication by barcode-aware analysis in Long Ranger[21]. Original input DNA fragments were  
161 inferred by collecting the read-pairs with the same barcode that were aligned in proximity to each  
162 other. A fragment was terminated if the distance between two consecutive reads with the identical

163 barcode larger than 50kb. Fragments were required to have at least two read pairs with the same  
164 barcode and a length of at least 2 kb. Partitions with fewer than three fragments were removed.  
165 We subsampled short-reads for each fragment to satisfy the expected  $C_R$ .

166

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

168 To compare the observations from real data with a known truth set, we developed LRTK-SIM, a  
169 simulator that follows the workflow of the 10x Chromium system and generates synthetic Linked-  
170 Reads like those produced by an Illumina HiSeqX machine (**Supplementary Information** and  
171 **Figure S9**). Based on the parameters commonly employed by 10x Genomics Linked-Read  
172 sequencing and the characteristics of our libraries, LRTK-SIM generated simulated datasets from  
173 the human reference (hg38), explicitly modeling the five key steps in real data generation.  
174 Parameters in parentheses are from the standard 10x Genomics protocol: 1. Shearing genomic  
175 DNA into long fragments ( $W_{\mu_{FL}}$  from 50kb to 100kb); 2. Loading DNA to the 10x Chromium  
176 instrument ( $\sim 1.25$ ng DNA); 3. Allocating DNA fragments into partitions which are attached the  
177 unique barcodes ( $\sim 10$  fragments per partition); 4. Generating short fragments; 5. Generating  
178 Illumina paired-end short reads (800M $\sim$ 1200M reads). LRTK-SIM first generated a diploid  
179 reference genome as a template by duplicating the human reference genome (hg38) into two  
180 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)  
181 [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)  
182 [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)  
183 [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  
184 randomly simulated 1 SNV per 1 kb. The ratio was 2:1 for heterozygous and homozygous SNVs.  
185 From this diploid reference genome, LRTK-SIM generated long DNA fragments by randomly  
186 shearing each haplotype with multiple copies into pieces whose lengths were sampled from an  
187 exponential distribution with mean of  $\mu_{FL}$ . These fragments were then allocated to pseudo-

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

199

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

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

209

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

211 We compared single nucleotide variants (SNVs) and structural variants (SVs) from the diploid  
212 regions of our assemblies with the ones from standard Illumina data and reference-based

213 processing of our 10x data. The standard Illumina data were downloaded from Genome in a Bottle  
214 [32] and analyzed with SVABA [33] to generate SV calls, and with BWA [34] and FreeBayes [35]  
215 to generate SNV calls. Long ranger ([https://support.10xgenomics.com/genome-](https://support.10xgenomics.com/genome-exome/software/pipelines/latest/what-is-long-ranger)  
216 [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  
217 deletions) calls for 10x reference-based analysis. We noted that R<sub>9</sub> failed to be analyzed by Long  
218 Ranger due to its extremely large C<sub>F</sub>. For SNVs, we benchmarked the calls from three strategies  
219 using the gold standard of NA12878 ([ftp://ftp-](ftp://ftp-trace.ncbi.nlm.nih.gov/giab/ftp/release/NA12878_HG001/latest/GRCh38/)  
220 [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  
221 ([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_)  
222 [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<sub>9</sub>, R<sub>10</sub>, R<sub>11</sub>) from HG002  
223 with the Tier 1 SV benchmark from Genome in a Bottle [36] and used VaPoR [37] to validate our  
224 SV calls based on PacBio CCS reads from NA24385 (Highly-accurate long-read sequencing  
225 improves variant detection and assembly of a human genome). We compared SNV and SV calls  
226 among the different approaches using vcfEval [38] and truvari [36], respectively.

227

228 **Performance of diploid assembly: influence of total coverage** Diploid assembly by Linked-  
229 Reads requires sufficient total read coverage ( $C=C_R \times C_F$ ) to generate long contigs and scaffolds.  
230 In this experiment, to explore the roles of both physical coverage ( $C_F$ ) and per-fragment read  
231 coverage ( $C_R$ ), we first generated eight simulated libraries whose total coverage  $C$  ranged from  
232 16x to 78x: four with  $C_R$  fixed and increasing  $C_F$  and four with fixed  $C_F$ , and increasing  $C_R$  (**Table**  
233 **S2**). Contig and scaffold N50s increased along with increasing either  $C_F$  or  $C_R$  (**Figure 2A** and  
234 **2B**). To investigate whether the trend was also present in the real datasets, we analyzed six real  
235 libraries (three by varying  $C_F$ , and the other three by varying  $C_R$ ; **Figure 1**): as  $C$  increased, we  
236 varied  $C_F$  and  $C_R$  independently by fixing the other parameter. Contig and scaffold N50s also  
237 increased in these simulation (**Figure 2C** and **2D**) and real linked-read sets (**Figure 2E** and **2F**)  
238 as a function of total coverage  $C$ . Contig lengths did increase a little (621.4kb to 758.1kb for



239 simulation; 110.7kb to 119.6kb for real data) when  $C$  was increased beyond 56X. Accuracy, which  
240 we define as the ratio between NA50 (N50 after breaking contigs or scaffolds at assembly errors)  
241 and N50 (**Figure 2C** and **2E**), changed 18% for simulation and 7% for real data (587.5kb to  
242 713.3kb for simulation; 97.1kb to 104.5kb for real data). For scaffolds in the real data sets, when  
243  $C$  increased from 48X ( $R_3$ ) to 67X ( $R_4$ ), both scaffold N50 and NA50 were significantly improved  
244 (N50: 13.4Mb to 30.6Mb; NA50: 6.3Mb to 12.0Mb), but the accuracy dropped slightly from 46.6%  
245 to 39.1%, which indicated that scaffold accuracy may be refractory to extremely high  $C$  (**Figure**  
246 **2F**). These results indicated that assembly length and accuracy were comparable over a broad  
247 range of  $C_F$  and  $C_R$  at constant  $C$ , which implied that assembly quality was mainly determined by  
248  $C$ .

249

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

261

262 **Performance of diploid assembly: nature of the source genome.** Assembly errors may occur  
263 because of heterozygosity, repetitive sequences, or sequencing error. To illuminate possible

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

272

273 **Performance of diploid assembly: fraction of genome in diploid state.** While contiguity is an  
274 important parameter for any whole genome assembly, evaluation of diploid assemblies  
275 necessitates estimating the fraction of the genome in which the assembly recovered the diploid  
276 state. To this end, we divided the contigs generated by Supernova2 into “diploid contigs”, which  
277 were extracted from its megabubble structures, and “haploid contigs” from non-megabubble  
278 structures. Pairs of scaffolds were extracted as the two haplotypes from megabubble structures  
279 if they shared the same start and end nodes in the assembly graph. Diploid contigs were  
280 generated by breaking the candidate scaffolds at the sequences with least 10 consecutive ‘N’s  
281 and were aligned to human reference genome (hg38) by Minimap2. The genome was split into  
282 500bp windows and diploid regions were defined as the maximum extent of successive windows  
283 covered by two contigs, each from one haplotype. Alignment against the human reference  
284 genome revealed the overall genome coverages of the six assemblies to be around 91%. For  
285 most assemblies, 70%-80% of the genome was covered by two homologous contigs (**Table 1**),  
286 with  $R_6$  only reaching 58.9%, probably due to the short fragments of the DNA preparation  
287 ( $\mu_{FL}=24\text{kb}$ ). We also analyzed another seven assemblies produced by 10x Genomics, all of which  
288 had diploid fractions of about 80% as well (**Table S5**). In the male NA24385, non-  
289 pseudoautosomal regions of the X chromosome are hemizygous and should therefore be

290 recovered as haploid regions. Between 79.9% and 87.6% of these regions were covered by one  
291 contig exactly depending on the assembled library. Library construction parameters other than  
292 fragment length appeared to have had little impact on the proportion of diploid regions (**Tables 1**  
293 and **Table S5**).

294

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

304

305 **Performance of diploid assembly: quality of variant calls.** The ultimate goal of human genome  
306 assembly is to accurately identify genomic variants. We compared the SNVs and SVs from our  
307 assemblies with the calls from referenced-based processing of standard Illumina and 10x data,  
308 and benchmarked them using gold standard from Genome in a Bottle and PacBio CCS reads.  
309 We found the SNVs from referenced-based processing of standard Illumina and 10x data were  
310 comparable and both of them were better than assembly-based calls (**Table S7** and **S8**) For SVs,  
311 our assemblies generated many calls that were missed by the reference-based strategy (**Table**  
312 **S9-S12**) and even by the Tier 1 benchmark of Genome in a Bottle (**Table S13**), and half of the  
313 deletions and a majority of insertions could be validated by PacBio CCS reads (**Table S14**).

314

## 315 **Discussion**

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

323

## 324 **10x Practicalities**

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

336

337 Our results showed that in practice,  $C_F$  should be between 335X and 823X, but no larger than  
338 1000X, given the optimal coverage of  $C=56X$  recommended by 10x and the requirement for  
339 sufficient per-fragment read coverage. Surprisingly, we observed that including more extremely

340 long fragments was detrimental for assembly quality. This is possibly due to the loss of barcode  
341 specificity for fragments spanning repetitive sequences. From a computational perspective, too  
342 many long fragments are harmful to deconvolving the *de bruijn* graph, as more complex paths  
343 need to be picked out. In our experiments,  $W_{\mu_{FL}}$  between 50kb and 150kb is the best choice to  
344 generate reliable assemblies.

345

### 346 **Parameters driving assembly quality**

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

359

360 Diploid assembly is the appropriate approach for assembly of genomes of diploid organisms that  
361 harbor variation. Therefore, an important metric to evaluate diploid assembly is the fraction of the  
362 genome that is assembled in a diploid state. The short input fragment length of  $R_6$  resulted in  
363 roughly 20% less of the genome in a diploid state (<60% vs <80%) compared to the other libraries  
364 of the same individual. This observation suggests that in addition to metrics such as N50,

365 evaluation of assembly quality should also include the fraction of the genome (or the assembly)  
366 that is in a diploid state.

367

### 368 **Cost-benefit analysis**

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

377

### 378 **Comparison with hybrid assemblies**

379 Hybrid assembly strategies have been applied successfully to produce human genome assembly  
380 of long contiguity [13, 14, 41]. In these studies, long contigs are first produced by single-molecule  
381 long-reads, such as PacBio (NG50=1.1Mb; [13]) or Nanopore (NG50=3.21Mb; [14]) comparing  
382 favorably to our best results for Linked-Reads assemblies (NG50=236kb). Scaffolding is then  
383 performed with complementary technologies such as BioNano to capture chromosomal level long-  
384 range information. It promoted the scaffold N50 of PacBio to 31.1Mb [13] and Illumina mate-pair  
385 sequencing with 10x data to 33.5Mb [25]. Using SuperNova2, the scaffold N50 from our studies  
386 reached ~27.86Mb ( $R_6$ ) on the basis of 10x data alone, suggesting that 10x technology gives  
387 broadly comparable results at a fraction of the price of long-read-based hybrid assemblies.

388

## 389 **Availability of supporting data**

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

395

## 396 **Additional files**

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

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

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

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

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

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

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

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

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

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

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

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

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

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

422 **Figure S1. Basic statistics for  $L_{1L}$ .** The distributions of **A.** the number of fragments per partition;  
423 **B.** sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
424 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
425 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
426 fragment lengths.

427 **Figure S2. Basic statistics for  $L_{1M}$ .** The distributions of **A.** number of fragments per partition; **B.**  
428 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
429 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
430 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
431 fragment lengths.

432 **Figure S3. Basic statistics for  $L_{1H}$ .** The distributions of **A.** number of fragments per partition; **B.**  
433 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
434 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
435 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
436 fragment lengths.

437 **Figure S4. Basic statistics for  $L_2$ .** The distributions of **A.** number of fragments per partition; **B.**  
438 sequencing depth per fragment; **C.** probability density function of unweighted fragment lengths;  
439 **D.** cumulative density function of unweighted fragment lengths; **E.** reversed cumulative density  
440 function of unweighted fragment lengths; **F.** reversed cumulative density function of weighted  
441 fragment lengths.

442 **Figure S5. Basic statistics for  $L_3$ .** 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 S6. Basic statistics for  $L_4$ .** 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 S7. Basic statistics for  $L_5$ .** 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 S8. Basic statistics for  $L_6$ .** 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 S9.** The workflow of LRTK-SIM to simulate linked-reads

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

465 **Figure S11.** Comparison of assembly qualities from 10x data with and without single nucleotide  
466 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  
467 **C** and **D**;  $C_R$  was fixed 0.2X and  $C_F$  was fixed 300X in **E** and **F**.

468 **Figure S12.** Comparison of assembly qualities from 10x data with (1% uniform) and without  
469 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  
470 300X in **C** and **D**;  $C_R$  was fixed 0.2X and  $C_F$  was fixed 300X in **E** and **F**.

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

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

480

481

482 **Competing interest**

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

484

485 **Author Contributions**

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

489

490 **Acknowledgements**

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

494

495

496 **Table**

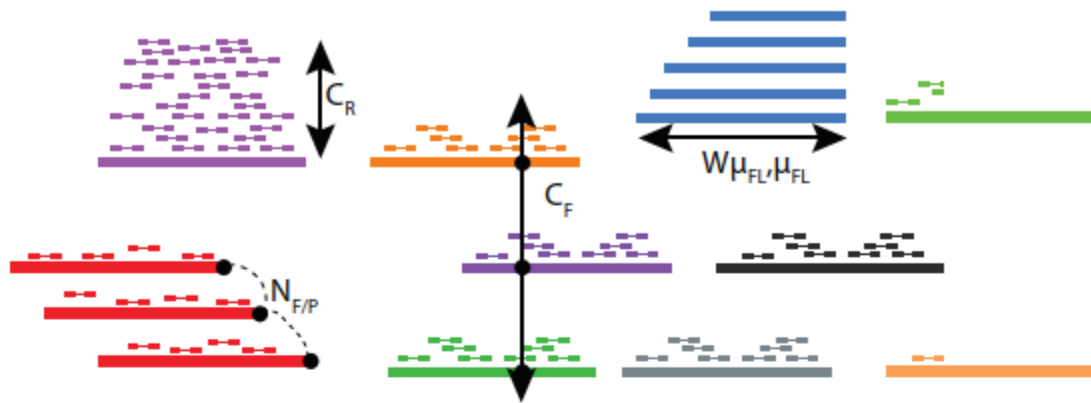
Linked- reads set	Overall (%)	Diploid regions (%)	Haploid regions (%)	Non-PAR (%)	Total contig length (contig>500bp)	Length of contigs from megabubble (contig>500bp)	Percentage (%)
$R_6$	91.9	58.9	27.7	-	5,632,483,053	3,758,345,846	66.73
$R_7$	91.1	73.3	11.3	-	5,613,140,437	4,668,186,478	83.17
$R_8$	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

497

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

500

501 **Figures**



**Parameter**

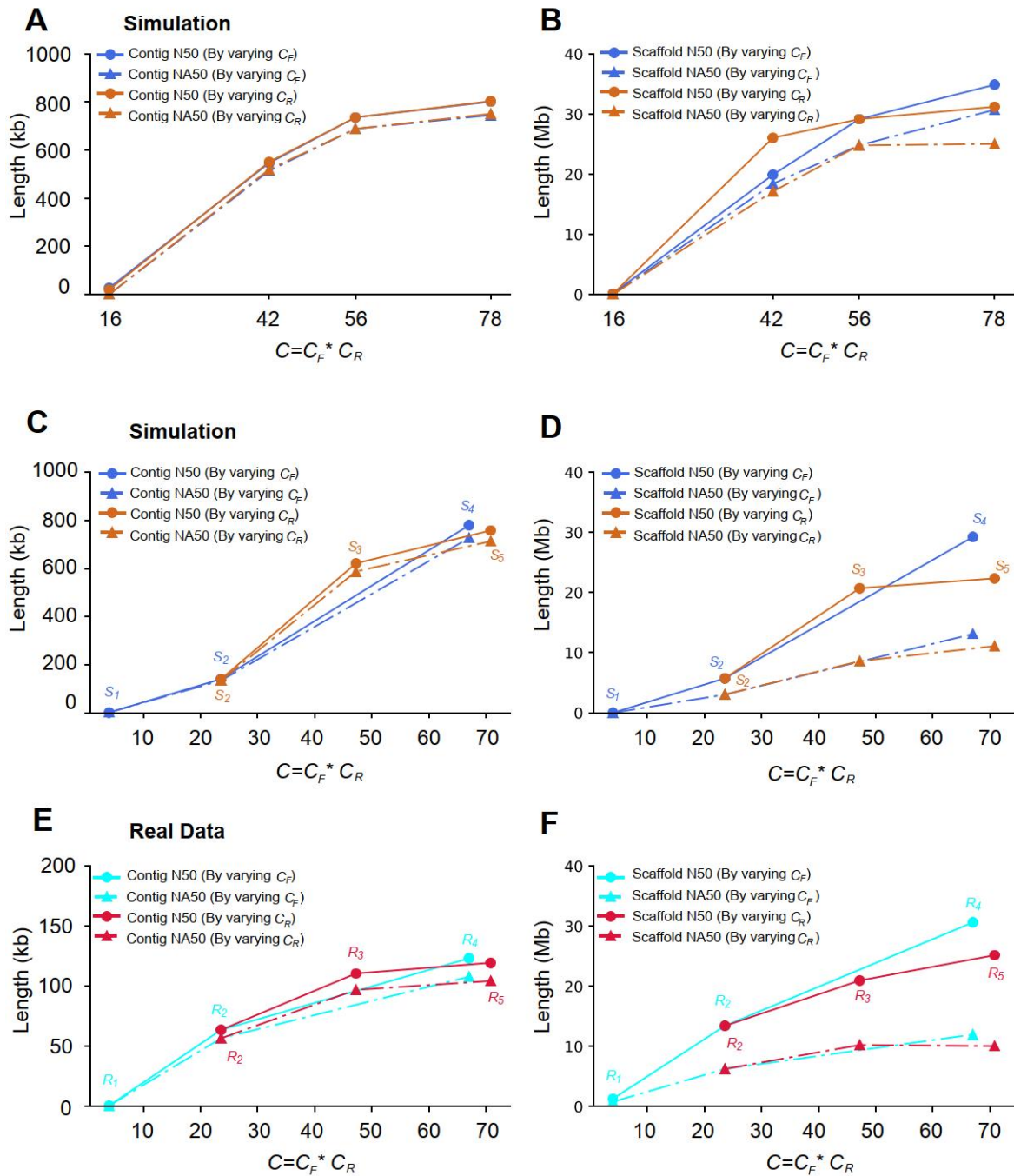
$N_{F/P}$  = Number of fragments per partition  
 $\mu_{FL}$  = Mean fragment length  
 $W\mu_{FL}$  = Weighted mean fragment length  
 $C_R$  = Read coverage per fragment  
 $C_F$  = Physical (fragment) coverage  
 $C$  = total coverage

**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

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



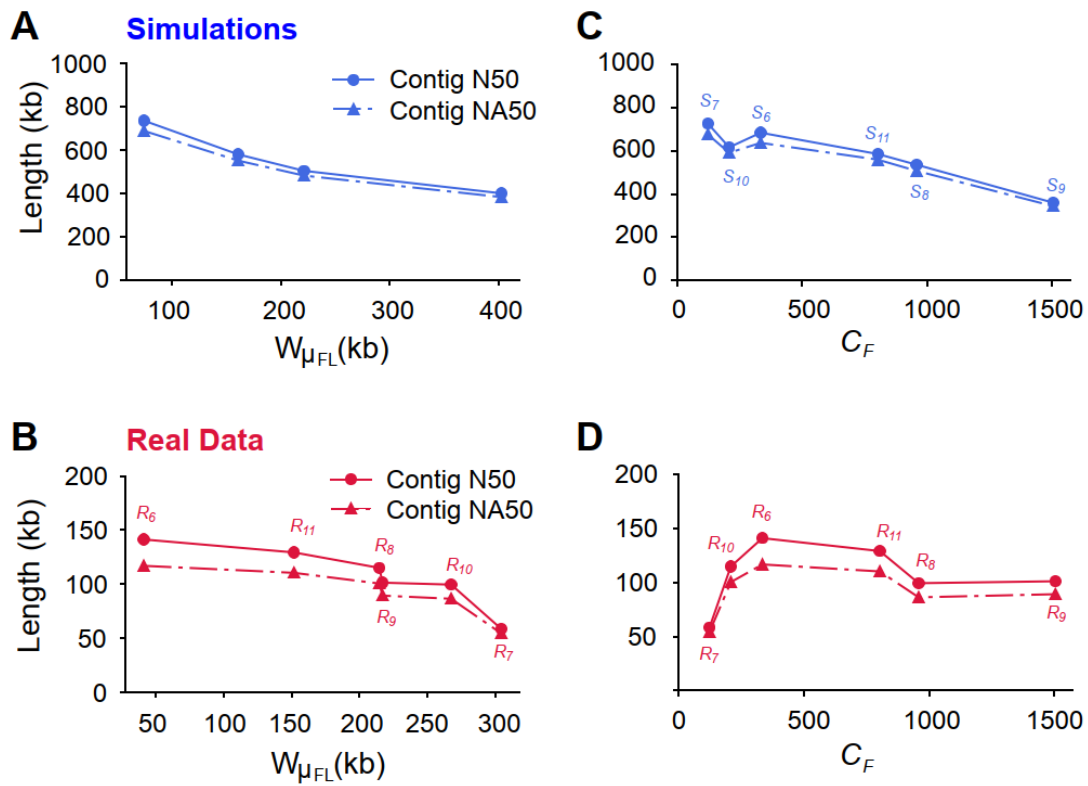
504

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

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

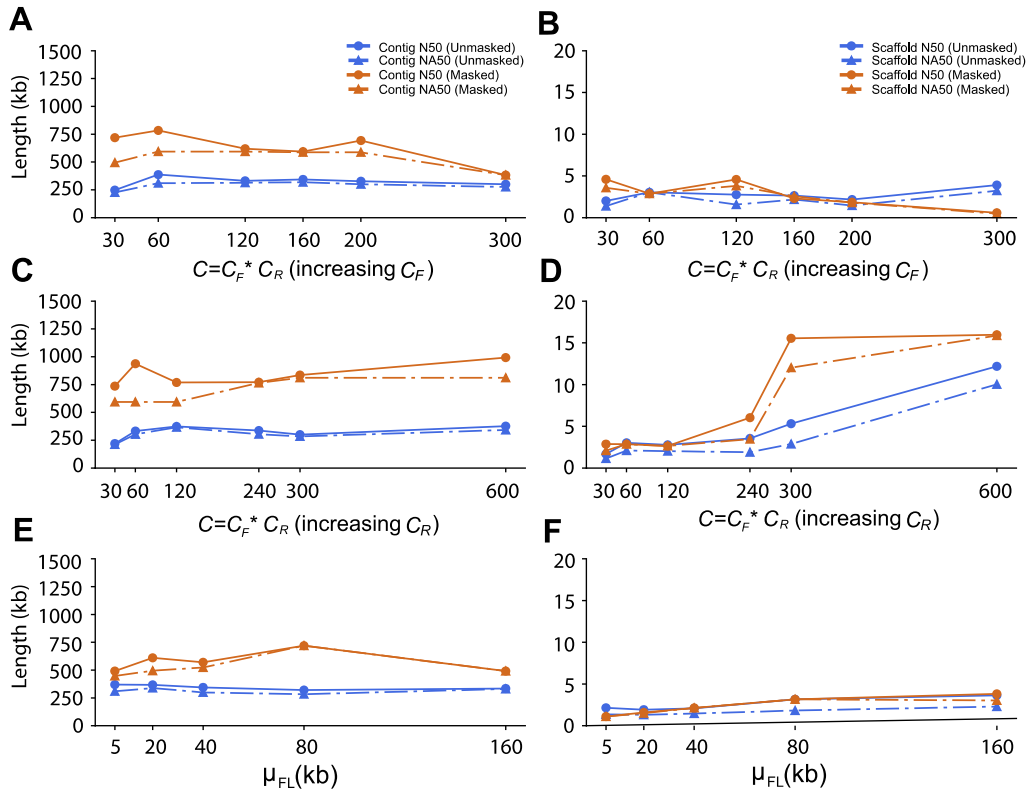
507 reads with matched parameters of real Linked-Read data sets (**Figure 1**); **E and F:** Real linked-

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



509

510 **Figure 3.** Contig qualities (N50 and NA50) as a function of fragment length  $W_{\mu_{FL}}$  or physical  
 511 coverage  $C_F$ , at  $C=56X$ . **A** and **C**, results from simulations; **B** and **D**, results from real data.



512

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

514 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

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

516

517 **References**

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



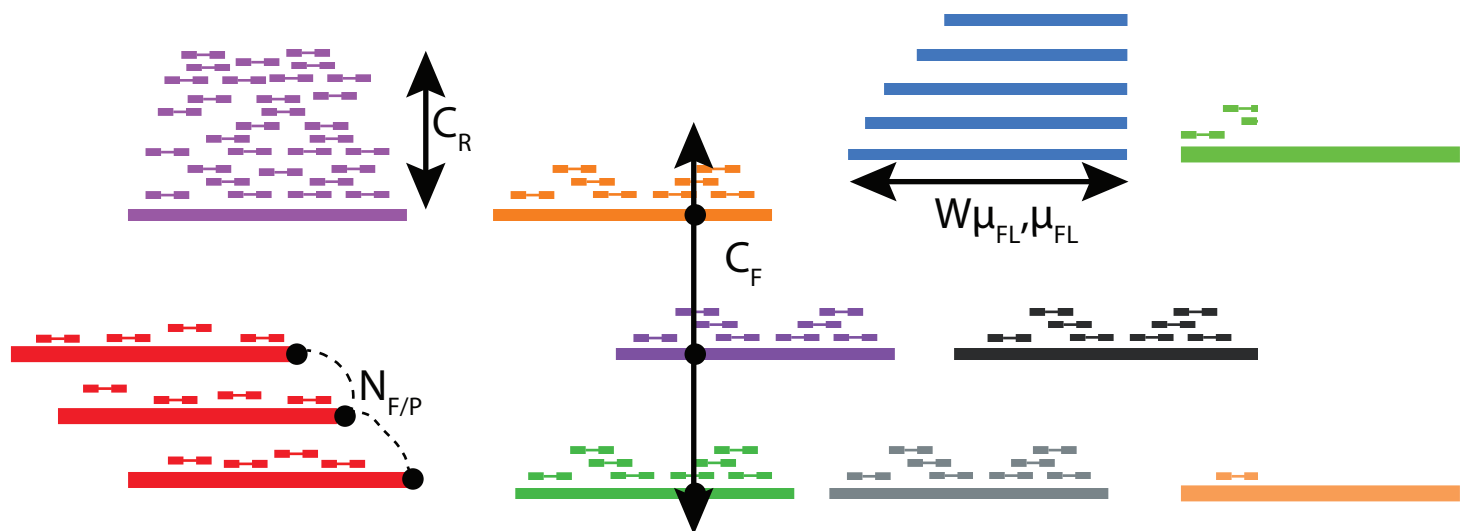
- 561 17. Kuleshov V, Xie D, Chen R, Pushkarev D, Ma Z, Blauwkamp T, et al. Whole-genome  
562 haplotyping using long reads and statistical methods. *Nat Biotechnol.* 2014;32 3:261-6.  
563 doi:10.1038/nbt.2833.
- 564 18. Peters BA, Kermani BG, Sparks AB, Alferov O, Hong P, Alexeev A, et al. Accurate whole-  
565 genome sequencing and haplotyping from 10 to 20 human cells. *Nature.* 2012;487  
566 7406:190-5. doi:10.1038/nature11236.
- 567 19. Edge P, Bafna V and Bansal V. HapCUT2: robust and accurate haplotype assembly for  
568 diverse sequencing technologies. *Genome Res.* 2017;27 5:801-12.  
569 doi:10.1101/gr.213462.116.
- 570 20. Patterson M, Marschall T, Pisanti N, van Iersel L, Stougie L, Klau GW, et al. WhatsHap:  
571 Weighted Haplotype Assembly for Future-Generation Sequencing Reads. *J Comput Biol.*  
572 2015;22 6:498-509. doi:10.1089/cmb.2014.0157.
- 573 21. Zheng GX, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, et al. Haplotyping  
574 germline and cancer genomes with high-throughput linked-read sequencing. *Nat*  
575 *Biotechnol.* 2016;34 3:303-11. doi:10.1038/nbt.3432.
- 576 22. Spies N, Weng Z, Bishara A, McDaniel J, Catoe D, Zook JM, et al. Genome-wide  
577 reconstruction of complex structural variants using read clouds. *Nat Methods.* 2017;14  
578 9:915-20. doi:10.1038/nmeth.4366.
- 579 23. Bishara A, Moss EL, Kolmogorov M, Parada AE, Weng Z, Sidow A, et al. High-quality  
580 genome sequences of uncultured microbes by assembly of read clouds. *Nat Biotechnol.*  
581 2018; doi:10.1038/nbt.4266.
- 582 24. Weisenfeld NI, Kumar V, Shah P, Church DM and Jaffe DB. Direct determination of  
583 diploid genome sequences. *Genome Res.* 2017;27 5:757-67. doi:10.1101/gr.214874.116.
- 584 25. Mostovoy Y, Levy-Sakin M, Lam J, Lam ET, Hastie AR, Marks P, et al. A hybrid approach  
585 for de novo human genome sequence assembly and phasing. *Nat Methods.* 2016;13 7:587-  
586 90. doi:10.1038/nmeth.3865.
- 587 26. Hulse-Kemp AM, Maheshwari S, Stoffel K, Hill TA, Jaffe D, Williams SR, et al. Reference  
588 quality assembly of the 3.5-Gb genome of *Capsicum annuum* from a single linked-read  
589 library. *Hortic Res.* 2018;5:4. doi:10.1038/s41438-017-0011-0.
- 590 27. Elyanow R, Wu HT and Raphael BJ. Identifying structural variants using linked-read  
591 sequencing data. *Bioinformatics.* 2017; doi:10.1093/bioinformatics/btx712.
- 592 28. Jones SJ, Haulena M, Taylor GA, Chan S, Bilobram S, Warren RL, et al. The Genome of  
593 the Northern Sea Otter (*Enhydra lutris kenyoni*). *Genes (Basel).* 2017;8 12  
594 doi:10.3390/genes8120379.
- 595 29. Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics.* 2018;34  
596 18:3094-100. doi:10.1093/bioinformatics/bty191.
- 597 30. Mikheenko A, Prjibelski A, Saveliev V, Antipov D and Gurevich A. Versatile genome  
598 assembly evaluation with QUAST-LG. *Bioinformatics.* 2018;34 13:i142-i50.  
599 doi:10.1093/bioinformatics/bty266.
- 600 31. Earl D, Bradnam K, St John J, Darling A, Lin D, Fass J, et al. Assemblathon 1: a  
601 competitive assessment of de novo short read assembly methods. *Genome Res.* 2011;21  
602 12:2224-41. doi:10.1101/gr.126599.111.
- 603 32. Zook JM, Catoe D, McDaniel J, Vang L, Spies N, Sidow A, et al. Extensive sequencing of  
604 seven human genomes to characterize benchmark reference materials. *Sci Data.*  
605 2016;3:160025. doi:10.1038/sdata.2016.25.

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

**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

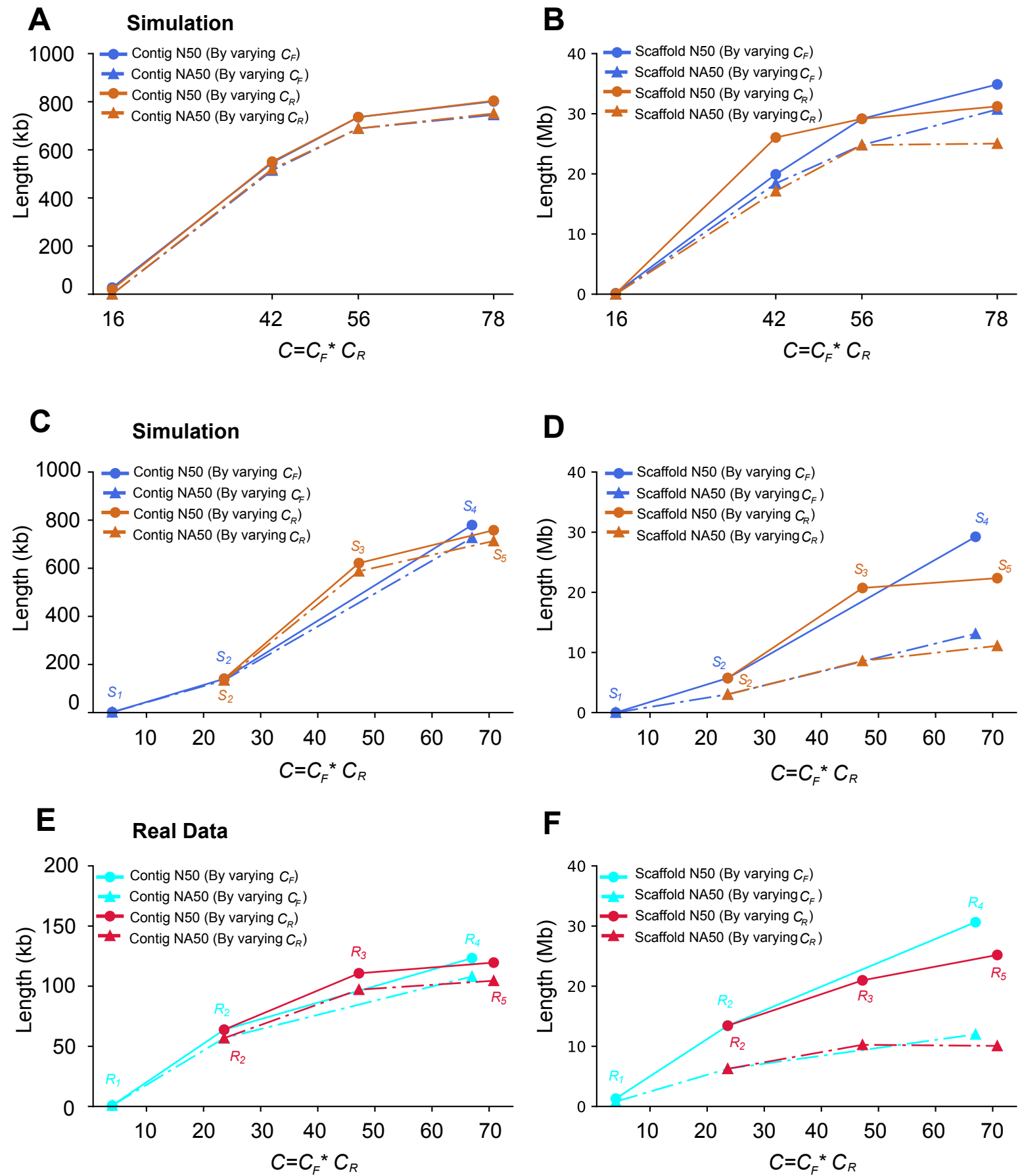
$W\mu_{FL}$  = 20-400kb

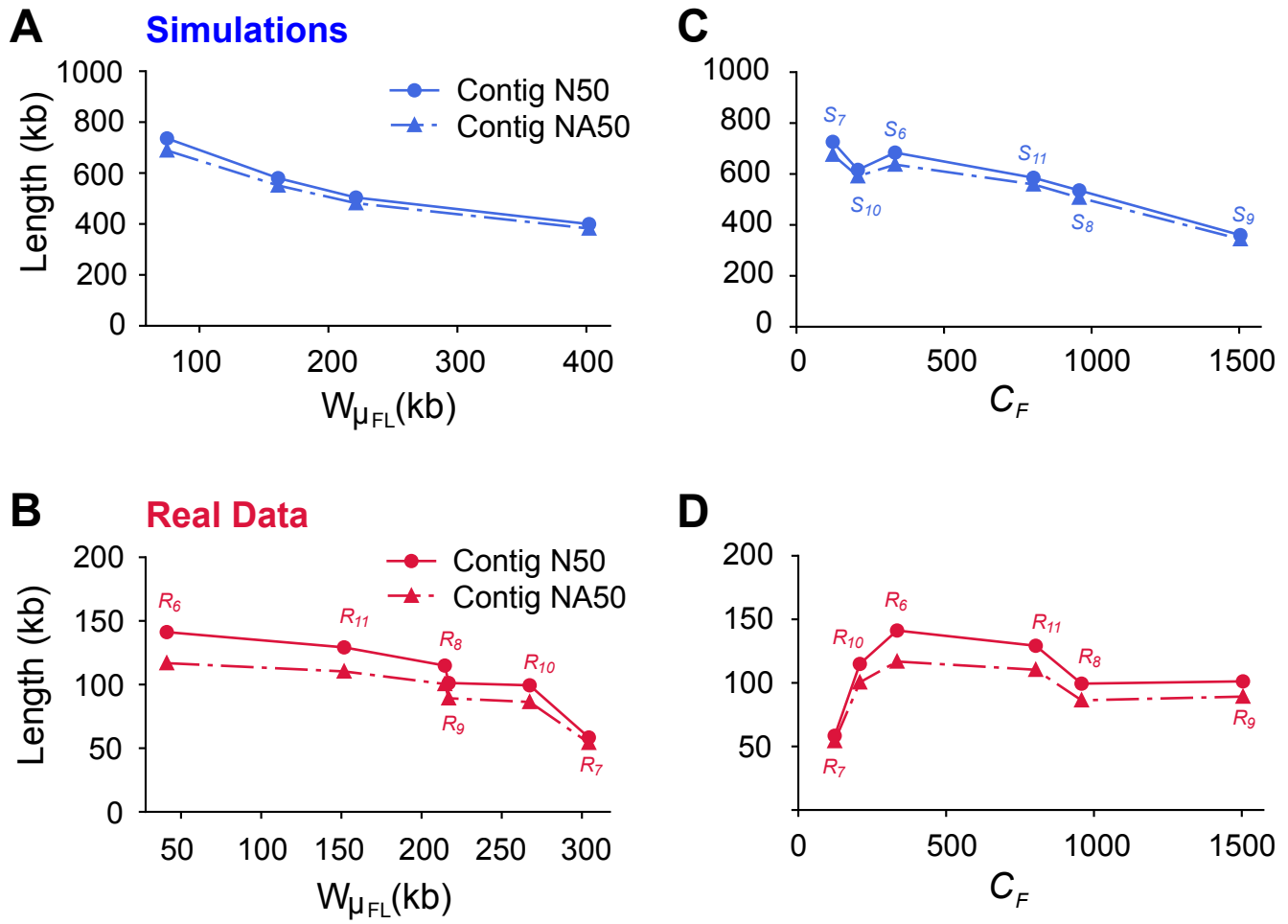
$C_R$  = 0.1x - 0.4x

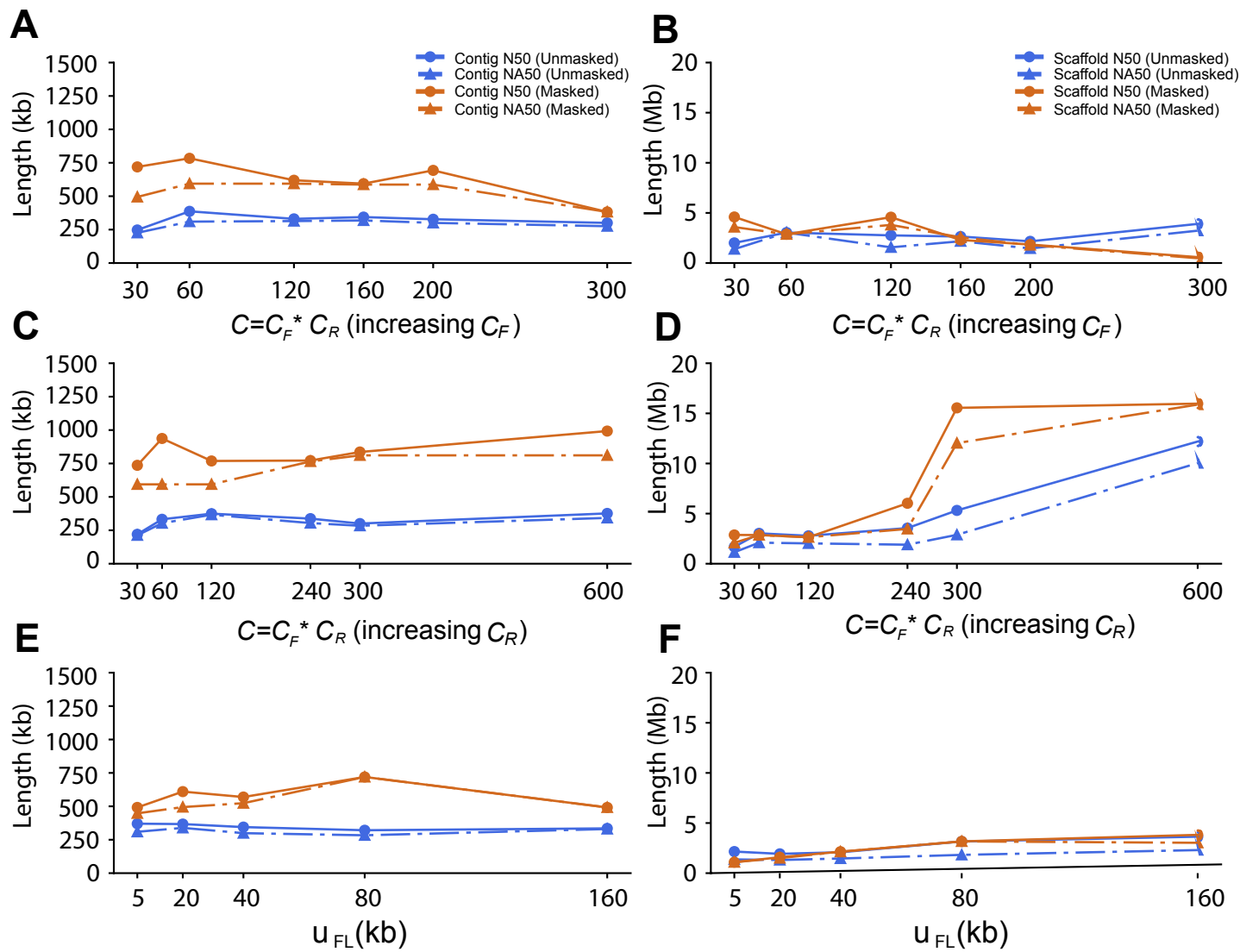
$C_F$  = 200x - 1000x

$C = C_R * C_F = 40x - 80x$

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









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, August 7, 2019

Dr. Hongling Zhou  
Editor  
GigaScience

Dear Dr. Zhou,

It is my pleasure to resubmit our revised, significantly improved and extended, 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 all of the reviewers' comments, which are addressed point by point in our response, and hope that you will be able to reach a positive decision.

Sincerely,



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>