1  **Extensive and deep sequencing of the Venter/HuRef genome for**
2  **developing and benchmarking genome analysis tools**
3
4  Bo Zhou[1,3], Joseph G. Arthur[2,3], Steve S. Ho[1], Reenal Pattni[1], Wing H. Wong[2],
5  Alexander E. Urban[1]
6
7  [1]Department of Psychiatry and Behavioral Sciences, Department of Genetics, Stanford
8  University School of Medicine, Stanford, California 94305, USA
9  [2]Department of Statistics, Department of Biomedical Data Science, Bio-X Program,
10 Stanford University, Stanford, California 94305, USA
11 [3]These authors contributed equally to this work.
12
13 corresponding author:
14
15 Alexander E. Urban (aeurban@stanford.edu)

16 **ABSTRACT**

17

18 We produced an extensive collection of deep re-sequencing datasets for the
19 Venter/HuRef genome using the Illumina massively-parallel DNA sequencing platform.
20 The original Venter genome sequence is a very-high quality phased assembly based on
21 Sanger sequencing. Therefore, researchers developing novel computational tools for
22 the analysis of human genome sequence variation for the dominant Illumina sequencing
23 technology can test and hone their algorithms by making variant calls from these
24 Venter/HuRef datasets and then immediately confirm the detected variants in the
25 Sanger assembly, freeing them of the need for further experimental validation. This
26 process also applies to implementing and benchmarking existing genome analysis
27 pipelines. We prepared and sequenced 200 bp and 350 bp short-insert whole-genome
28 sequencing libraries (sequenced to 100x and 40x genomic coverages respectively) as
29 well as 2 kb, 5 kb, and 12 kb mate-pair libraries (43x, 97x, and 122x physical coverages
30 respectively). Lastly, we produced a linked-read library (133x physical coverage) from
31 which we also performed haplotype phasing.

## BACKGROUND & SUMMARY

Almost two decades ago the extensive efforts of the Human Genome Project, backed up by work from Celera, resulted in the release of a draft of the first complete sequence of the human genome [1,2]. This catalyzed a new era of human whole-genome analysis where the now-available human genome sequence has been studied intensely to understand the functions of its parts and their interactions with each other and where a concurrent genome technology revolution has produced ever more powerful platforms to carry out such functional studies [3]. Since then, increasingly large numbers of human genomes have been sequenced, yielding insights into population-level genetic variation [4–6], structural genome variation [7–9], and mutational mechanisms [10]. Technological advances have progressively improved the information content and reduced the noise profile of sequencing data [11]. A large variety of methodologies for the routine analysis of sequencing data is now available [12]. "Whole-genome sequencing" is now a standing term that refers to the re-sequencing of a given sample of human genomic DNA using, typically, the dominant Illumina DNA sequencing platforms which can quickly produce several hundred million short sequencing reads at affordable costs. These reads are then aligned to the human reference genome and analyzed using various approaches [12–14], such as mismatch analysis, read-depth analysis, split-read analysis and discordant read-pairs analysis, producing an extensive catalog of sequence variants that are present in the DNA sample in question relative to the human reference sequence. The promise of human genome research is nothing short of a complete transformation of basic life science research, translational research, and eventually the way we diagnose, treat, and find cures for human disease.

It is clear, however, that current standard whole-genome sequence analysis leaves a rather large room for improvement. The standard genome analysis practices of today perform rather poorly in certain contexts, such as in repetitive regions (i.e. in around half the human genome), in the detection and resolution of complex structural variation, or in placing detected variants in their proper haplotypes. Although more advanced and novel computational algorithms that address these limitations are continuously being developed, one essential requirement during this process is that the detected variants are to be experimentally validated in order to establish false-positive rates and to make it possible to further tune and optimize the new algorithms. Experimental validation, especially of complex variants, during the tool development and testing phases is a very laborious and time-consuming process, but it can be circumvented by using a genome for which sufficiently large numbers of variants are already known, i.e. prevalidated. Several studies have been conducted with the goal of extensively characterizing the variants in a small number of human genomes using multiple sequencing technologies [15,16]. In some human genomes, variants have been carefully and extensively documented, providing a benchmark for other studies [9,17–20].

The Venter (HuRef) Genome, however, is especially distinguished for quality among the publicly-available human genome sequences as it is the only one for which its complete diploid assembly was generated from high-quality Sanger reads [17] and for which extensive catalogs of SNPs, indels, and structural variation are available [18,20]. To date,

78  no extensive Illumina sequencing datasets have been available for the Venter/HuRef
79  genome in contrast to other genomes that have been characterized for benchmarking
80  purposes [15,16].
81
82  To unlock the potential of the Venter/HuRef genome as the outstanding benchmark
83  genome, we have conducted deep whole-genome sequencing (WGS) using a variety of
84  sequencing strategies for the Illumina platform (**Table 1**). Specifically, we produced
85  short-insert paired-end WGS datasets at a combined sequence coverage of 140x,
86  linked-read data at 42x de-duplicated sequencing coverage (133x physical coverage),
87  and three long-insert (2 kb, 5 kb, and 12 kb) paired-end (i.e. mate-pair) WGS datasets
88  with physical coverages of 43x, 97x, and 122x, respectively (**Figure 1**]. These datasets
89  are of very high quality (**Figures 2-4**] and are complemented by the existing
90  Venter/HuRef assembly-quality Sanger reads [17] and long-read sequencing data, which
91  was produced using the Pacific Biosciences platform [21].
92
93  Researchers developing novel computational tools for analyzing whole-genome
94  sequencing data can now test their algorithms by processing the appropriate
95  Venter/HuRef Illumina datasets described here and then turn to the already-available
96  catalogs of sequence variants, or to the original Sanger reads [17], to confirm the
97  characterization of variants detected by their algorithms . Likewise, whenever a
98  laboratory implements a new computational pipeline for human genome analysis, it can
99  now use these Illumina Venter/HuRef datasets to confirm proper implementation and to
100  optimize proper settings for the pipeline.
101
102  **METHODS**
103
104  **Venter/HuRef DNA Sample**
105
106  The Venter/HuRef DNA sample as obtained as a 50 µg aliquot of LCL-extracted DNA
107  (NS12911) from the Coriell Institute for Medical Research where the iPSC (GM25430)
108  of the same subject is also available (https://catalog.coriell.org/1/HuRef).
109
110  **Illumina paired-end WGS**
111
112  *Library Preparation*
113
113  The library preparation was previously described in detail in Mu *et al* [20].  Briefly, 1 µg
114  of genomic DNA was fragmented using 2µL of NEBNext dsDNA fragmentase (New
115  England Biolabs, Ipswich, MA) in 1x fragmentation buffer and 1x BSA.  Reaction was
116  kept on ice for 5 minutes before adding the fragmentase and was incubated at 37□°C
117  for 20 minutes.  The reaction was stopped by addition of 5 µL of 0.5 M EDTA.  DNA
118  was purified from the reaction mixture using 0.9x by volume AMPure XP beads
119  (Beckman Coulter, Cat# A63880) and eluted in 50 µL of 10mM Tris-Acetate (pH 8.0)
120  buffer.  Six independent fragmentation reaction replicates were performed, and the
121  sizes of the DNA were analyzed using Agilent 2100 Bioanalyzer before library
122  preparation.

123  Library preparation was performed using the KAPA Library Preparation kit (KAPA
124  Biosystems, Wilington, MA) where 200 ng of fragmented DNA was used as input.
125  Library was constructed according to manufacture's protocol where the DNA was
126  end-repaired and A-tailed before adapter ligation with Illumina TruSeq Adapter (Index
127  1). DNA was then purified using 0.8x by volume AMPure XP beads and quantified
128  using the Qubit ds DNA High Sensitivity Assay Kit (Life Technologies, Cat# Q32851).
129  For PCR amplification, 50 ng of DNA was amplified using the KAPA HiFi DNA
130  Polymerase with the following thermocycling conditions:  98□°C/45□s, 5 cycles of
131  (98□°C/15□s, 60□°C/30□s, 72□°C/45□s), 72□°C/1□min, and 4□°C /hold.  Primers
132  from the KAPA Library Preparation kit was used for PCR amplification.  Afterwards,
133  DNA was purified from the PCR reaction using AMPure XP beads and eluted in 30 µL
134  of 10mM Tris-Acetate (pH 8.0) buffer.  Six independent experimental replicates were
135  performed, and the purified PCR amplified DNA fragments from each replicate was
136  pooled for size selection and gel-purified from 2% agarose gel.  Two size selections
137  were made at 200 bp and 350 bp.

138  *Sequencing*

139  Sequencing of the 200 bp and 350 bp insert-size libraries was described previously in
140  Mu et al [20].  The libraries were sequenced separately (2x100 bp) on an Illumina
141  HiSeq 2000 instrument in rapid run mode. For the 200 bp insert-size library, a total of
142  3,214,626,588 reads generated from 5 sequencing runs was pooled together to
143  obtain 100x genomic coverage. For the 350 bp insert-size library, a total of
144  1,280,576,580 reads generated from two sequencing runs was pooled together to
145  obtain 40x genomic coverage.

146  *Analysis*
147
148  Reads were trimmed at the 3' end to a uniform length of 100 bp using FASTX toolkit
149  (http://hannonlab.cshl.edu/fastx_toolkit/; version 0.0.13). The trimmed reads were
150  aligned by BWA-MEM (Li and Durbin 2009; version 0.7.17-r1188) using the hg38
151  reference with ALT alleles removed
152  (ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/001/405/GCA_000001405.15_GRCh38/
153  seqs_for_alignment_pipelines.ucsc_ids/), and the resulting alignment records were
154  sorted with Samtools (http://www.htslib.org/; version 1.7). Marking of PCR duplicates
155  and calculations of insert-size and coverage information was performed using Picard
156  (http://picard.sourceforge.net; version 2.17.10).
157
158  **Illumina mate-pair WGS**
159
160  *Library Preparation*
161
162  Mate Pair libraries at insert sizes 2 kb, 5 kb, and 12 kb were generated from
163  Venter/HuRef DNA using the Nextera Mate Pair Sample Preparation Kit (Illumina,
164  Cat# FC-132–1001) following standard manufacturer's instructions with the exception
165  of the shearing step (see below).  The Venter/HuRef DNA sample was first verified as

166  high molecular weight (>15 kb) by running 60 ng, quantified by using the Qubit
167  dsDNA HS Assay Kit (Life Technologies, Cat# Q32851), on 0.8% 1X TAE agarose
168  gel next to the 1 kb Plus DNA Ladder (ThermoFisher Cat# 10787018).  Afterwards,
169  for each insert size, 4 µg of the high molecular weight genomic DNA was tagmented
170  with biotinylated junction adapters and fragmented to about 7-8 kb on average in a
171  400 µL tagmentation reaction containing 12 µL of Tagmentase at 55 °C for 30 min.
172  The tagmented DNA fragments were purified by adding 2X the volume of DNA
173  Binding Buffer with Zymo Genomic DNA Clean & Concentrator Kit (Zymo Research,
174  Cat# D4010) and eluted in 30 uL of Elution Buffer after two washes with the provided
175  Wash Buffer. To fill in the gaps in the DNA adjacent to the junction adapters as a by-
176  product of Tagmentation, single-strand displacement reaction was performed in a 200
177  µL reaction by adding 132 µL of water, 20 µL of 10x Strand Displacement Buffer, 8 µL
178  of dNTPs, and 10 µL of Strand Displacement Polymerase to the 30 µL elution and at
179  20 °C for 30 min. DNA purification was then performed in 30 µL elution with 0.5x
180  volume of AMPure XP Beads (Beckman Coulter, Cat# A63880) and size-selected by
181  using BluePippin (Sage Science).  The 0.75% DF 3-10kb Marker S1 – Improved
182  Recovery and the 0.75% DF 10-18kb Marker U1 protocols were used for size
183  selection on the BluePippin for insert sizes 5 kb and 12 kb respectively, and 0.75%
184  DF 1-6kb Marker S1 protocol was used for insert size 2 kb.  The "Tight Selection"
185  option was used instead of "Range" for all size selections.  The size selected DNA
186  was then circularized overnight (12-16 hours) at 30 °C with Circularization Ligase in
187  a 300 µL reaction.
188
189  After overnight circularization, the uncirculated linear DNA was digested by adding 9
190  µL of Exonuclease and incubated at 30 °C for 30 minutes and heat inactivated at
191  70 °C for 30 minutes. Afterwards, 12 µL of Stop Ligation Buffer was added.
192  Circularized DNA was then transferred to T6 (6×32 mm) glass tube (Covaris, Part#
193  520031 and 520042) and sheared *twice* on the Covaris S2 machine (Intensity of 8,
194  Duty Cycle of 20%, Cycles Per Burst of 200, Time of 40 s, Temperature of 2–6 °C).
195  We find that shearing *twice* creates a tighter final library size distribution which leads
196  to a higher fraction of pass-filter clusters during the Illumina sequencing step.
197  The mate pair fragments within the sheared DNA fragments contain the biotinylated
198  junction adapter and were selected by binding to Dynabeads M-280 Streptavidin
199  Magnetic Beads (Invitrogen, Part# 112-05D) by adding an equal volume of the Bead
200  Bind Buffer (incubated at 20 °C for 15 minutes on shaking heat block at highest rpm
201  setting).  The non-biotinylated molecules in solution were washed away using the
202  Wash Buffer. All downstream reactions were carried out on streptavidin beads with
203  magnetic immobilization and washes with the Wash Buffer between successive
204  reactions (e.g. End Repair, A-Tailing, and Adapter Ligation.  The sheared DNA was
205  first End-repaired followed by A-Tailing and TruSeq indexed adapter ligation.
206  The adapter-ligated DNA was resuspended in 20 µL of Resuspension Buffer and then
207  PCR amplified in a 50 µL reaction with 25 µL of PCR 2X Master Mix and 5 µL of
208  Primers both provided in the Nextera Mate Pair Sample Preparation Kit (Illumina,
209  Cat# FC-132–1001) to generate the final library.  The thermocycling conditions are
210  98 °C/1 min, 15 cycles of (98 °C/10 s, 60 °C/30 s, 72 °C/30 s), 72 °C/5 min,
211  and 4 °C /hold.  The amplified library (supernatant) was purified using a 0.66x

6

212  volume of AMPure XP Beads (0.67x vol) and eluted in 20 µL of Resuspension Buffer.
213  The size distribution of the library was determined by Agilent Technologies 2100
214  Bioanalyzer (High Sensitivity Assay), and the indexed library concentration was
215  measured by the Qubit dsDNA HS Assay Kit (Life Technologies, Cat# Q32851).
216
217  *Sequencing*
218
219  The Mate-Pair libraries were sequenced on the Illumina NextSeq 500 using the
220  NextSeq 500/550 Mid Output v2 kit (300 cycles) (Illumina, Cat# FC-404-2003) to
221  generate 2×151□bp paired-end reads. The libraries were loaded onto the flowcell at a
222  final concentration of 1.8pM and 1% PhiX Control v3 (Illumina, Cat# FC-110-3001).
223  Additional rounds of sequencing also used a final library concentration of 1.8pM and
224  1% PhiX Control v3.
225
226  *Analysis*
227
228  Illumina Nextera Mate Pair junction adapter sequences were first trimmed using
229  NxTrim (O'Connell et al. 2015; version 0.4.3) with the "--aggressive --preserve-mp"
230  settings in order to maximize the number of long-insert pairs. Nxtrim outputs four sets
231  of reads, designated "Mate Pair", "Paired-End", "Singleton", and "Unknown." "Mate
232  Pair" reads have junction adapter sequence trimmed off from the 3' end of Read 1
233  and/or Read 2; "Paired-End" (short-insert) reads have junction adapter sequence
234  trimmed from the 5' end of Read 1 and/or Read 2; "Singleton" reads have junction
235  adapter sequence trimmed from the middle of either Read 1 or Read 2 rendering one
236  of the reads useless. "Unknown" reads have no junction adapter sequences detected.
237  This is most likely because the junction adapter sequence sits in the un-sequenced
238  portion of the template, thus whether reads are "Mate Pair" or "Paired-End" cannot be
239  discerned. Nonetheless, mate-pair reads are present in the "Unknown" fractions as
240  well as paired-end reads. The "Unknown" reads can be used for alignment and
241  analysis if more long-insert information is desired??? (O'Connell et al. 2015). Here,
242  the reads designated as "Mate Pair" and "Unknown" were combined, aligned with
243  BWA-MEM (Li and Durbin 2009) against the hg38 reference without ALT alleles
244  (ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/001/405/GCA_000001405.15_GRCh
245  38/seqs_for_alignment_pipelines.ucsc_ids/), and sorted using samtools
246  (http://www.htslib.org/; version 1.7). Marking of PCR duplicates and calculations of
247  insert-size and coverage information was performed using Picard
248  (http://picard.sourceforge.net; version 2.17.10).
249
250  **10X Genomics Chromium library for Illumina sequencing**
251
252  *Input genomic DNA preparation*
253
254  The Venter/HuRef DNA sample (obtained from the Coriell Institute for Medical
255  Research) was first verified as high molecular weight (>15 kb) by running 60 ng,
256  quantified by using the Qubit dsDNA HS Assay Kit (Life Technologies, Cat# Q32851),
257  on 0.8% 1X TAE agarose gel next to the 1 kb Plus DNA Ladder (ThermoFisher Cat#

258    10787018).  Afterwards, 4□µg of the high molecular weight genomic DNA was loaded
259    on a BluePippin (Sage Science) instrument to select for DNA fragments 30 kb to 80
260    kb using the "0.75%DF Marker U1 high-pass 30- 40 kb vs3" protocol.  The
261    concentration of the selected DNA fragments was then quantified by using the Qubit
262    dsDNA HS Assay Kit (Life Technologies, Cat# Q32851) and diluted to 1 ng/ µL.  The
263    final dilution concentration of 1 ng/ µL was verified again by performing three
264    technical replicates of Qubit dsDNA HS Assay with 5 µL of the DNA dilution as input.
265    *Chromium whole-genome linked-read library preparation and sequencing*
266    The linked-read whole-genome library was prepared using the Chromium Genome kit
267    and reagent delivery system (10X Genomics, Pleasanton, CA). The linked-read
268    library was made following standard manufacturer's protocol with 10 cycles of PCR
269    amplification. Briefly 1□ng of DNA (~300 genome equivalents) of size-selected high
270    molecular DNA was partitioned into ~1.5 million oil droplets in emulsion, tagged with a
271    unique 16 bp barcode within each droplet, and subjected isothermal amplification
272    (30□°C for 3 hours; 65□°C for 30□minutes) by random priming within each droplet.
273    Amplified (isothermal) DNA was then purified from the droplet emulsion following the
274    manufacturer's protocol using SPRI beads.  The purified DNA was then End-Repaired
275    and A-tailed followed by adapter ligation of adapter in the same reaction mixture.
276    DNA was purified from the was the reaction mixture using SPRI beads and eluted in
277    40 uL.  Sample Index PCR amplification (primers and 2x master mix provided in the
278    Chromium Genome kit) was then performed on the eluted DNA in a toal volume of
279    100 uL with the following thermocycling conditions:  98□°C/45□s, 10 cycles of
280    (98□°C/20□s, 54□°C/30□s, 72□°C/20□s), 72□°C/1□min, and 4□°C /hold.  Primer
281    index SI-GA-A6 was used.  DNA (final linked-read library) was purified from the PCR
282    reaction with SPRI bead size selection following manufacturer's protocol.
283    The final purified library was quantified by qPCR (KAPA Library Quantification Kit for
284    Illumina platforms, Kapa Biosystems, Wilmington, MA) using the following
285    thermocycling conditions: 95□°C/3 min, 30 cycles of (95□°C/5□s, 67□°C/30□s).  The
286    library concentration was calculated in nanomolar (nM) concentration and then diluted
287    to 5 nM.  Sequencing (2x151bp, 8 cycles of single indexing) on two lanes of Illumina
288    HiSeq X was performed at Macrogen (Rockville, MD).
289
290    *Sequencing*
291
292    The final purified library was quantified by qPCR (KAPA Library Quantification Kit for
293    Illumina platforms, Kapa Biosystems, Wilmington, MA) using the following
294    thermocycling conditions: 95□°C/3 min, 30 cycles of (95□°C/5□s, 67□°C/30□s).  The
295    library concentration was calculated in nanomolar (nM) concentration and then diluted
296    to 5 nM.  Sequencing (2x151bp, 8 cycles of single indexing) on two lanes of Illumina
297    HiSeq X (flowcell ID:  H3MHGALXX, lanes #4 and #5) was performed at Macrogen
298    (Rockville, MD) resulting in a total of 789,239,544 paired reads (**Table 1**).
299
300    *Analysis*
301
302    FASTQ files were generated raw BCL files using "*mkfastq*" mode in the Long Ranger
303    software (version 2.1.3) from 10X Genomics (Pleasanton, CA).  10X Genomics

8

304   Chromium library index "SI-GA-A6" was specified in the required sample sheet file for
305   "*mkfastq*".  Before alignment, the hg38 genome files were downloaded from
306   ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCA/000/001/405/GCA_000001405.15_GRCh38/
307   seqs_for_alignment_pipelines.ucsc_ids/GCA_000001405.15_GRCh38_no_alt_analysis
308   _set.fna.gz and indexed using the "*mkref*" mode in Long Ranger.  Sequencing
309   alignment and haplotype phasing were performed using the "*wgs*" mode in Long Ranger,
310   and the options *"--sex=male"* and *"--vcmode=freebayes"* were specified.  Only "PASS"
311   SNPs and Indels 50 bp or smaller were included in the final phased variant vcf.
312
313   **DATA RECORDS**
314
315   The Venter/HuRef genome sequenced is publicly available through The Coriell Institute
316   for Medical Research (Camden, NJ, USA) both as genomic DNA (catalog ID:  NS12911)
317   extracted from lymphoblastoid cell line (LCL) or as retroviral reprogrammed induced
318   pluriplotent stem cell culture (catalog ID:  GM25430).  As described in the Methods,
319   Venter/HuRef LCL DNA (NCBI SRA biosample accession SAMN03491120) was used
320   for sequencing library preparation in this work.
321
322   **Illumina short-insert WGS**
323
324   Approximately 100x sequencing coverage 2x100bp Illumina short-insert (200 bp) WGS
325   data generated from the Illumina HiSeq 2000 is available through NCBI SRA accession
326   SRR7097858 [Data Citation 1:  NCBI SRA SRR7097858].  Approximately 40x
327   sequencing coverage 2x100bp Illumina short-insert (350 bp) WGS data generated from
328   the Illumina HiSeq 2000 platform is available through NCBI SRA accession
329   SRR7097859 [Data Citation 2:  NCBI SRA SRR7097859].
330
331   **Illumina mate-pair WGS**
332
333   Illumina mate-pair data sequenced (2x150 bp) on the Illumina NextSeq 500 are
334   available through NCBI SRA accessions SRR6951312 [Data Citation 3:  NCBI SRA
335   SRR6951312], SRR6951313 [Data Citation 4:  NCBI SRA SRR6951313], and
336   SRR6951310 [Data Citation 5:  NCBI SRA SRR6951310] for insert sizes 2 kb, 5 kb, and
337   12 kb respectively.
338
339   **10X Genomics Chromium linked-read Library**
340
341   10X Genomics Chromium linked-read data sequenced (2x150 bp) on two lanes of the
342   Illumina HiSeq X Ten is available through NCBI SRA accession SRR6951311 [Data
343   Citation 6:  NCBI SRA SRR6951311].  The phased variants of the Venter/HuRef
344   genome obtained through the analysis linked reads is available through dbSNP
345   NCBI_ss# 2137543904 to 3651364986 (For phasing information, request for original
346   submitted vcf file through NCBI dbSNP.) [Data Citation 7:  NCBI dbSNP NCBI_ss#
347   2137543904-3651364986].

9

348
349 **TECHNICAL VALIDATION**
350
351 **Illumina short-insert WGS**
352
353 Sequencing quality of the WGS mate-pair libraries were assessed using FastQC
354 (Supplementary Information).  Insert-size, coverage, GC-bias, alignment, and
355 duplication metrics were analyzed using Picard tools
356 (http://broadinstitute.github.io/picard/).  These statistics are summarized in Table 1,
357 Table 2 and Figure 2A.
358
359 **Illumina mate-pair WGS**
360
361 Sequencing quality of the WGS mate-pair libraries was assessed using FastQC
362 (Supplementary Information).  Insert-size, coverage, GC-bias, alignment, and
363 duplication metrics were analyzed using Picard tools
364 (http://broadinstitute.github.io/picard/).  These statistics are summarized in **Table 1,**
365 **Table 2 and Figure 2C-J**.  Read fractions that were designated by NxTrim [23] as "Mate
366 Pair", "Paired-End", "Singletons", and "Unknown" are summarized in **Table 3**.  The
367 "Mate Pair" fraction for all libraries fall within the expected range (~40-60%).  The
368 relatively high rates of PCR duplication (expected for mate-pair libraries) result in
369 significant decreases in sequence coverage (3x to 7x) (**Table 1, Table 2, Figure 2,**
370 **Supplementary Information**). However, the more useful metric for mate-pair
371 sequencing is high physical coverage [15].  The mean insert sizes for the mate pair
372 libraries are 1.8 kb, 4.8 kb, and 12.2 kb (**Table 2, Figure 2**), which results in physical
373 genomic coverage values of 62x, 136x, and 162x respectively.
374
375 **10X Genomics Chromium Library**
376
377 Sequencing quality of the WGS mate-pair libraries were assessed using FastQC
378 (Supplementary Information).  Input molecule length, coverage, alignment, duplication,
379 droplet barcode, and phasing metrics were analyzed using Long Ranger software
380 version 2.1.5 from 10X Genomics (Pleasanton, CA, USA).  These statistics are
381 summarized in **Table 1, Table 2, Table 4 and Figure 3**.  Overall, 2.4 million and 1.5
382 million, 0.42 million and 0.29 million heterozygous and homozygous SNVs and indels
383 respectively were called (**Table 4**).  Of which, 96.7% and 93.85% of heterozygous SNVs
384 and Indels respectively were successfully phased in the Venter/HuRef Genome in a
385 total of 8882 haplotype blocks (N50 ~ 0.9 Mbp, longest phase block ~ 6.5 Mbp) (**Table**
386 **4**).  Phase blocks for each chromosome are shown in **Figure 3**.  Similar to mate-pair
387 libraries, the physical coverage of the linked read library is calculated to be 133x from
388 the mean input DNA molecule length of 32kb.
389
390 **USAGE NOTES (optional)**
391
392 The Venter/HuRef genome sequenced in this work is publicly available as both cell line
393 and DNA from Coriell Institute for Medical Research.  The mate-pair and linked-read

394 sequencing data used the same DNA sample/extraction as input. It is possible that
395 small differences may exist when compared to the short-insert datasets since the
396 input DNA came from different cell passages and extractions. Researchers are
397 especially encouraged to use the sequencing data in this work in combination with
398 diploid Sanger sequencing data available for the Venter/HuRef genome published in
399 Levy et al.
400
401 **ACKNOWLEDGEMENTS**
402
408
409 **AUTHOR CONTRIBUTIONS**
410
411 BZ and RP performed experiments. JGA, BZ, SSH performed data analysis. BZ, JGA,
412 and AEU wrote the manuscript.
413
414 **COMPETING INTERESTS**
415
416 The authors declare no conflict of interest.
417
418 **REFERENCES**
419
420 1.  Lander, E. S. *et al.* Initial sequencing and analysis of the human genome. *Nature*
421     **409,** 860–921 (2001).
422 2.  Venter, J. C. *et al.* The sequence of the human genome. *Science (80-. ).* **291,**
423     1304–51 (2001).
424 3.  Reuter, J. A., Spacek, D. V. & Snyder, M. P. High-Throughput Sequencing
425     Technologies. *Mol. Cell* **58,** 586–597 (2015).
426 4.  1000 Genomes Project Consortium *et al.* A map of human genome variation from
427     population-scale sequencing. *Nature* **467,** 1061–73 (2010).
428 5.  1000 Genomes Project Consortium *et al.* An integrated map of genetic variation
429     from 1,092 human genomes. *Nature* **491,** 56–65 (2012).
430 6.  1000 Genomes Project Consortium *et al.* A global reference for human genetic
431     variation. *Nature* **526,** 68–74 (2015).
432 7.  Korbel, J. O. *et al.* Paired-end mapping reveals extensive structural variation in
433     the human genome. *Science* **318,** 420–6 (2007).
434 8.  Sudmant, P. H. *et al.* An integrated map of structural variation in 2,504 human
435     genomes. *Nature* **526,** 75–81 (2015).
436 9.  Chaisson, M. J. P. *et al.* Multi-platform discovery of haplotype-resolved
437     structural variation in human genomes. *bioRxiv* 193144 (2017).
438     doi:10.1101/193144
439 10. Jónsson, H. *et al.* Parental influence on human germline de novo mutations in

440    1,548 trios from Iceland. *Nature* **549,** 519–522 (2017).

441 11. Kumar, V. *et al.* Uniform, optimal signal processing of mapped deep-sequencing
442    data. *Nat. Biotechnol.* **31,** 615–22 (2013).

443 12. Pabinger, S. *et al.* A survey of tools for variant analysis of next-generation
444    genome sequencing data. *Brief. Bioinform.* **15,** 256–278 (2014).

445 13. Alkan, C., Coe, B. P. & Eichler, E. E. Genome structural variation discovery and
446    genotyping. *Nat. Rev. Genet.* **12,** 363–376 (2011).

447 14. DePristo, M. a *et al.* A framework for variation discovery and genotyping using
448    next-generation DNA sequencing data. *Nat Genet* **43,** 491–8 (2011).

449 15. Zook, J. M. *et al.* Extensive sequencing of seven human genomes to characterize
450    benchmark reference materials. *Sci. data* **3,** 160025 (2016).

451 16. Eberle, M. A. *et al.* A reference data set of 5.4 million phased human variants
452    validated by genetic inheritance from sequencing a three-generation 17-member
453    pedigree. *Genome Res.* **27,** 157–164 (2017).

454 17. Levy, S. *et al.* The Diploid Genome Sequence of an Individual Human. *PLoS Biol.*
455    **5,** e254 (2007).

456 18. Pang, A. W. *et al.* Towards a comprehensive structural variation map of an
457    individual human genome. *Genome Biol.* **11,** R52 (2010).

458 19. Parikh, H. *et al.* svclassify: a method to establish benchmark structural variant
459    calls. *BMC Genomics* **17,** 64 (2016).

460 20. Mu, J. C. *et al.* Leveraging long read sequencing from a single individual to
461    provide a comprehensive resource for benchmarking variant calling methods. *Sci.*
462    *Rep.* **5,** 14493 (2015).

463 21. Lin, M. Comparing de novo assemblies of J. Craig Venter's genome. (2015).
464    doi:10.6084/m9.figshare.1319564.v1

465 22. Arthur, J. G., Chen, X., Zhou, B. & Urban, A. E. Detection of complex structural
466    variation from paired-end sequencing data. *bioRxiv* 1–32 (2017).
467    doi:10.1101/200170

468 23. O'Connell, J. *et al.* NxTrim: optimized trimming of Illumina mate pair reads.
469    *Bioinformatics* **31,** 2035–2037 (2015).

470

**DATA CITATIONS**

1. **Arthur, J. G. NCBI SRA SRR7097858 (2015)**
2. **Arthur, J. G. NCBI SRA SRR7097859 (2015)**
3. **Zhou, B. NCBI SRA SRR6951312 (2018)**
4. **Zhou, B. NCBI SRA SRR6951313 (2018)**
5. **Zhou, B. NCBI SRA SRR6951310 (2018)**
6. **Zhou, B. NCBI SRA SRR6951311 (2018)**
7. **Zhou, B. NCBI dbSNP NCBI_ss# 2137543904-3651364986 (2018)***

*For phasing information, request the originally submitted VCF file through NCBI dbSNP.

**FIGURE & TABLE LEGENDS**

**Figure 1. (A)** Schematic diagram of the study. Venter/HuRef genomic DNA was used to generate short-insert (200 bp and 350 bp), mate-pair (2 kb, 5 kb, and 12 kb), and linked-read libraries. (**B**) Detailed overview of data generation including bio-sample used, types of Illumina WGS libraries constructed, sequencing instrument platforms, types of sequencing runs, and subsequent analysis of data.

**Figure 2**. Normalized coverage, GC (%) content windows, base quality at GC (%), and corresponding insert-size histograms for all WGS libraries: 200 bp short-insert (**A,B**), 350 bp short-insert (**C,D**), 2kb-mate-pair (**E,F**), 5kb-mate-pair (**G,H**), 12kb-mate-pair (**I,J**).

**Figure 3**. Coverage (deduplicated) histograms of (**A,B**) short-insert, (**C,D,E**) 2 kb, 5 kb, and 12 kb mate-pair, and (**F**) linked-read libraries. Only reads with mapping score > 20 were used.

**Figure 4**. Violin plot of sizes of haplotype blocks constructed using linked-read sequencing (133x physical coverage) for HuRef/Venter Genome for all chromosomes.

**Table 1**. Summary of library construction and sequencing for short-insert, mate-pair, and linked-read HuRef/Venter WGS libraries.

**Table 2**. Summary of post sequencing QC, alignment, duplication, coverage and insert-size analysis for all libraries.
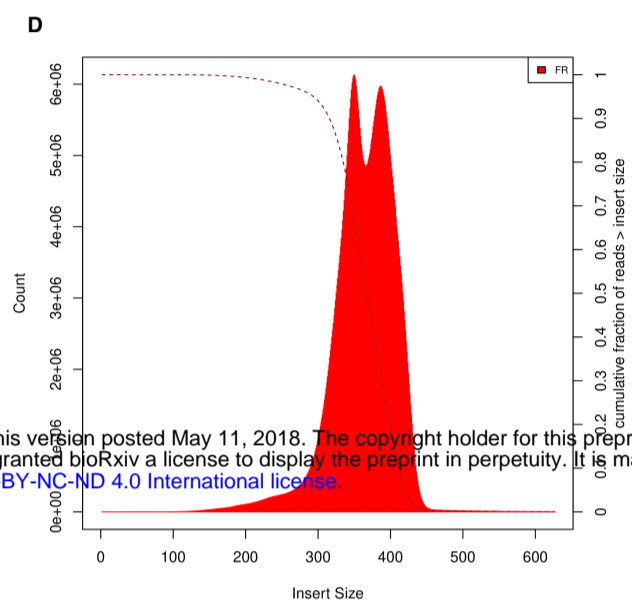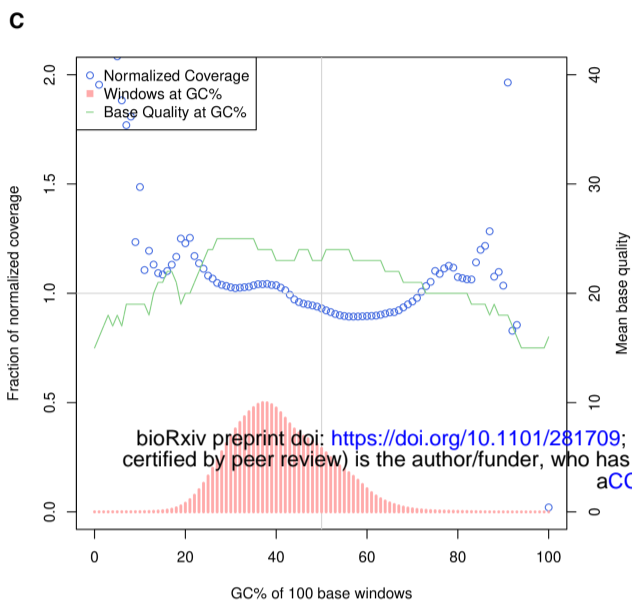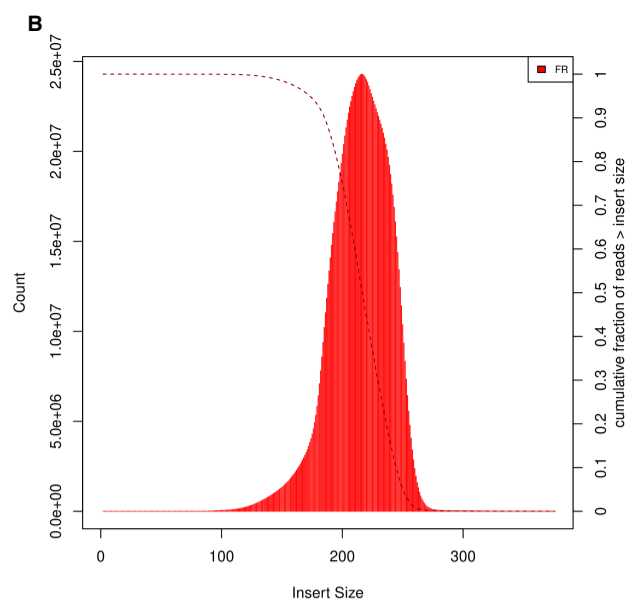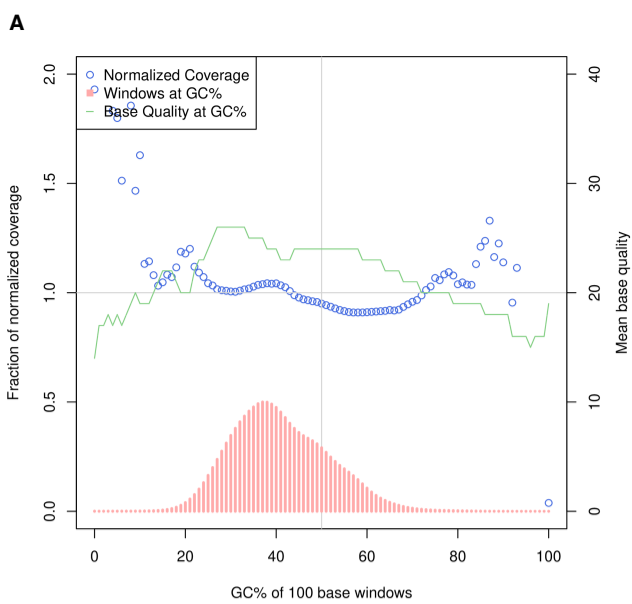
**Table 3**. Statistics for trimming of Nexera junction adapter sequence using NxTrim [23] for all mate-pair libraries.
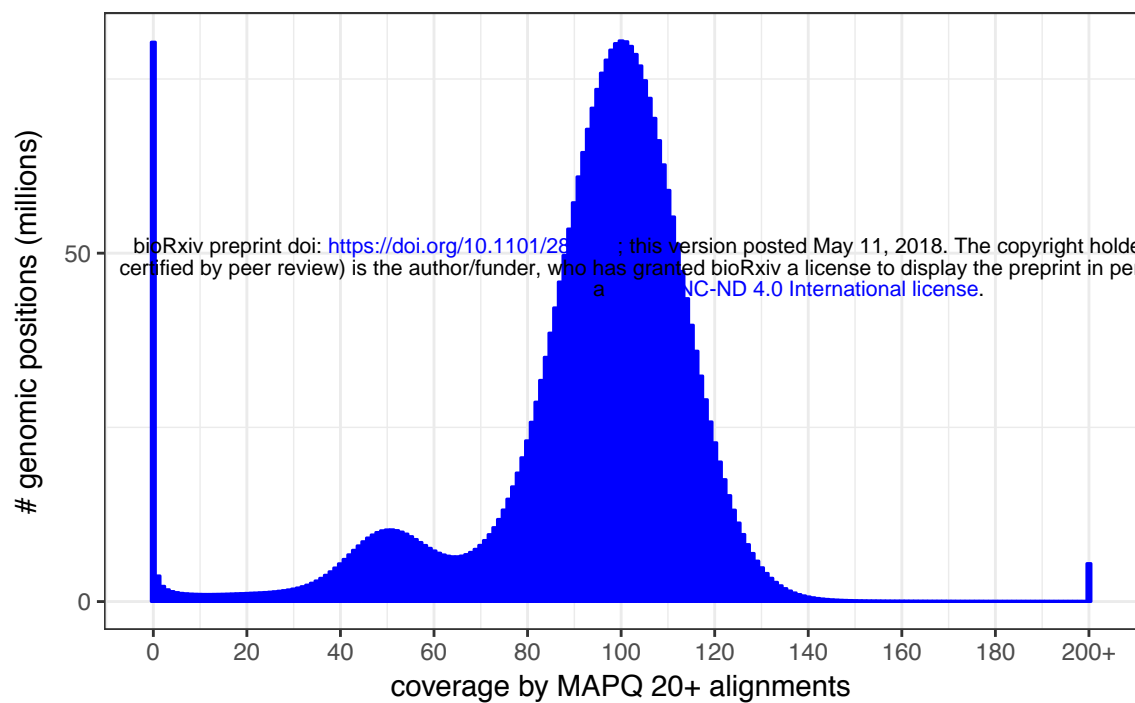
**Table 4**. Summary of metrics for linked-read sequencing and phasing of the HuRef/Venter genome.

14

**a**

350bp

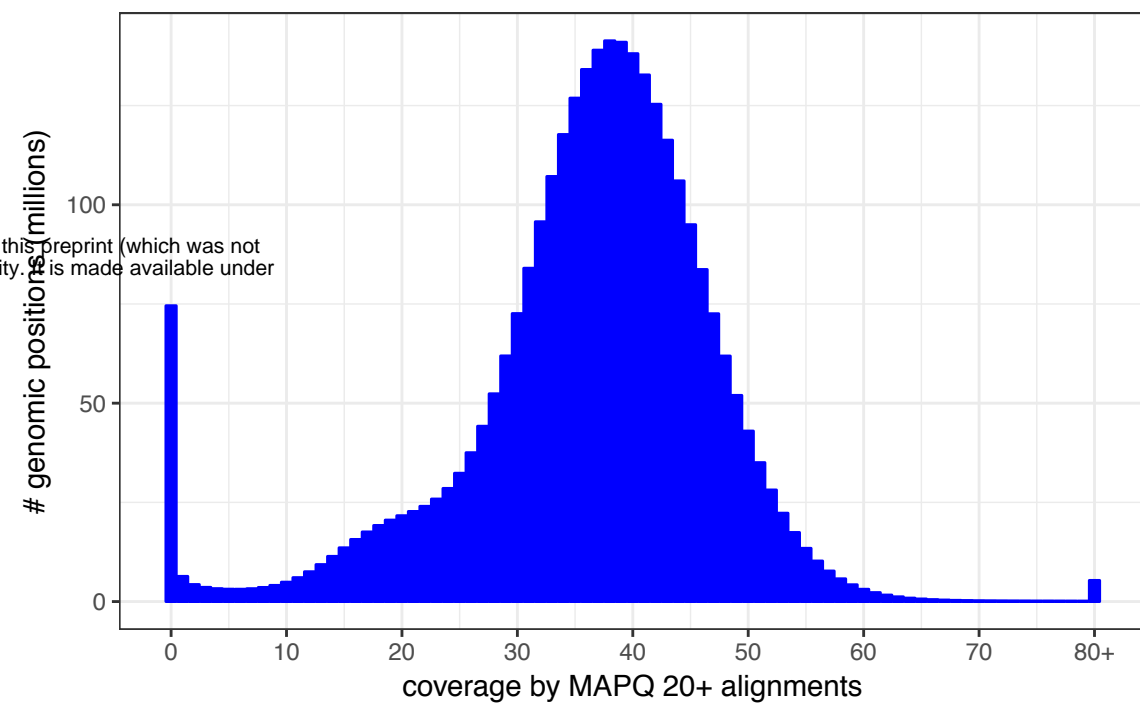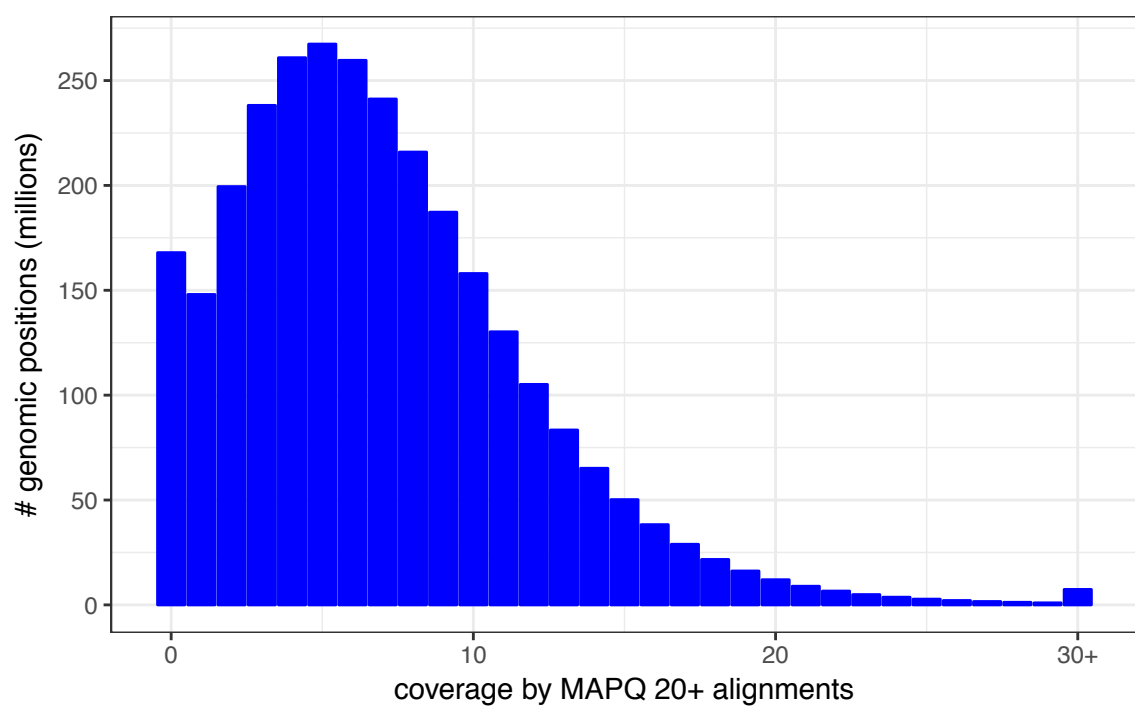500bp

B    B    B

2kb    5kb    12kb

**b**

| Sample | WGS Library | Sequencing | Analysis |
|---|---|---|---|
| Venter/HuRef (LCL DNA) | Short-insert (350bp) | HiSeq 2000 (2x100bp) | BWA, Picard |
| | Short-insert (500bp) | HiSeq 2000 (2x100bp) | BWA, Picard |
| | Mate-Pair (2kb) | NextSeq 500 (2x151bp) | NxTrim, BWA, Picard |
| | Mate-Pair (5kb) | NextSeq 500 (2x151bp) | NxTrim, BWA, Picard |
| | Mate-Pair (12kb) | NextSeq 500 (2x151bp) | NxTrim, BWA, Picard |
| | Linked-Read | HiSeq X (2x151bp) | Long Ranger |