

GigaScience

Entering the era of conservation genomics: Cost-effective assembly of the African wild dog genome using linked reads.

--Manuscript Draft--

Manuscript Number:	GIGA-D-17-00324	
Full Title:	Entering the era of conservation genomics: Cost-effective assembly of the African wild dog genome using linked reads.	
Article Type:	Research	
Funding Information:	John Stuelpnagel	Not applicable
Abstract:	<p>A high-quality reference genome assembly is a valuable tool for the study of non-model organisms across disciplines. Genomic techniques can provide important insights about past population sizes, local adaptation, and even aid in the development of breeding management plans. This information can be particularly important for fields like conservation genetics, where endangered species require critical and immediate attention. However, funding for genomic-based methods can be sparse for conservation projects, as costs for general species management can consume budgets. Here we report the generation of high-quality reference genomes for the African wild dog (<i>Lycaon pictus</i>) at a low cost, thereby facilitating future studies of this endangered canid. We generated assemblies for three individuals from whole blood samples using the linked-read 10x Genomics Chromium system. The most continuous assembly had a scaffold N50 of 21 Mb, a contig N50 of 83 Kb, and completely reconstructed 95% of conserved mammalian genes as reported by BUSCO v2, indicating a high assembly quality. Thus, we show that 10x Genomics Chromium data can be used to effectively generate high-quality genomes of mammal species from Illumina short-read data of intermediate coverage (~25-50x). Interestingly, the wild dog shows a much higher heterozygosity than other species of conservation concern, possibly as a result of its behavioral ecology. The availability of reference genomes for non-model organisms will facilitate better genetic monitoring of threatened species such as the African wild dog and help researchers and conservationists to better understand the ecology and adaptability of those species in a changing environment.</p>	
Corresponding Author:	Ellie Armstrong Stanford University UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Stanford University	
Corresponding Author's Secondary Institution:		
First Author:	Ellie Armstrong	
First Author Secondary Information:		
Order of Authors:	Ellie Armstrong	
	Ryan W Taylor	
	Stefan Prost	
	Peter Blinston	
	Esther van der Meer	
	Hillary Madzikanda	
	Olivia Mufute	
	Roseline Madisodza-Chikerema	
	John Stuelpnagel	

	Claudio Sillero-Zubiri
	Dmitri Petrov
Order of Authors Secondary Information:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</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 Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (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 Minimum</p>	Yes

[Standards Reporting Checklist?](#)

1 **Entering the era of conservation genomics: Cost-effective assembly of the**
2 **African wild dog genome using linked reads.**

3
4
5
6 4 Ellie E. Armstrong^{1*}, Ryan W. Taylor^{1*}, Stefan Prost^{1,2}, Peter Blinston³, Esther van der
7
8 5 Meer³, Hillary Madzikanda³, Olivia Mufute⁴, Roseline Mandisodza-Chikerema⁴, John
9
10 6 Stuelpnagel⁵, Claudio Sillero-Zubiri⁶, Dmitri Petrov¹

11
12
13
14
15 8 ¹Program for Conservation Genomics, Department of Biology, Stanford University,
16
17 9 Stanford, CA, USA

18
19
20 10 ²Department of Integrative Biology, University of California, Berkeley, CA, USA

21
22 11 ³Painted Dog Conservation, Dete, Zimbabwe

23
24 12 ⁴The Zimbabwe Parks & Wildlife Management Authority, Zimbabwe

25
26 13 ⁵10x Genomics, Inc., Pleasanton, CA

27
28 14 ⁶Wildlife Conservation Research Unit, Zoology, University of Oxford, The Recanati-
29
30 15 Kaplan Centre, Tubney, UK

31
32
33 16

34
35
36 17 * These authors contributed equally to this work.

37
38 18 Corresponding Author: Ellie E. Armstrong (elliea@stanford.edu)

39
40 19

41
42 20

43
44 21 **Abstract**

45
46 22

47
48 23 **Background**

49
50
51 24 A high-quality reference genome assembly is a valuable tool for the study of non-
52
53 25 model organisms across disciplines. Genomic techniques can provide important
54
55 26 insights about past population sizes, local adaptation, and even aid in the
56
57 27 development of breeding management plans. This information can be particularly
58
59 28 important for fields like conservation genetics, where endangered species require

60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29 critical and immediate attention. However, funding for genomic-based methods can
30 be sparse for conservation projects, as costs for general species management can
31 consume budgets.

32

33 **Findings**

34 Here we report the generation of high-quality reference genomes for the African wild
35 dog (*Lycaon pictus*) at a low cost, thereby facilitating future studies of this
36 endangered canid. We generated assemblies for three individuals from whole blood
37 samples using the linked-read 10x Genomics Chromium system. The most
38 continuous assembly had a scaffold N50 of 21 Mb, a contig N50 of 83 Kb, and
39 completely reconstructed 95% of conserved mammalian genes as reported by
40 BUSCO v2, indicating a high assembly quality.

41

42 **Conclusions**

43 We show that 10x Genomics Chromium data can be used to effectively generate
44 high-quality genomes of mammal species from Illumina short-read data of
45 intermediate coverage (~25-50x). Interestingly, the wild dog shows a much higher
46 heterozygosity than other species of conservation concern, possibly as a result of its
47 behavioral ecology. The availability of reference genomes for non-model organisms
48 will facilitate better genetic monitoring of threatened species such as the African wild
49 dog and help researchers and conservationists to better understand the ecology and
50 adaptability of those species in a changing environment.

51

52 **Keywords**

53 Conservation genomics, 10x Genomics Chromium, African wild dog, *Lycaon pictus*,
54 *de novo* Assembly

55

56 **Background**

61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57 Major population declines have been observed in vertebrate groups over the
58 past several hundred years, primarily due to anthropogenic change [1]. This decline
59 has resulted in extinction rates unprecedented in recent history [1, 2]. The
60 conservation of extant species will require major efforts in restoring and preserving
61 habitat, along with protection, management, and investment by local stakeholders.
62 Though, by definition, all species of conservation concern exist as small populations,
63 populations can still retain genetic variation that was generated and maintained a few
64 generations back, when population sizes were much larger. Within patterns of
65 historic genetic variation are signals of demographic history, gene flow, and natural
66 selection which can inform efforts towards the long-term survival of species. In
67 addition to signals of a species history, genetic information can be used to uncover
68 important contemporary or very recent events and processes. For example, Epstein,
69 Jones [3] identified genes that may confer facial tumor resistance in Tasmanian
70 devils, suggesting that the ability to artificially select for resistance in non-infected
71 populations may allow for a more robust population rescue and recovery. Genetic
72 markers can be used to track individual movement across landscapes either
73 indirectly by measuring relatedness, or directly by genotyping scat or hair left by an
74 individual as it moves. Additionally, the identification and assignment of individuals
75 through genotyping can be an important tool for law enforcement to assign
76 contraband and confiscated materials to their geographic origin. Conservationists
77 can also use fine grained measurements of reproductive success along with
78 genotypes and environmental variables to gather a detailed understanding of the
79 factors contributing to or limiting population growth, such as inbreeding depression.
80 Taken together genomic tools are poised to have a major contribution to
81 conservation [4, 5].

55
56
57
58
59
60
61
62
63
64
65
82 The African wild dog (*Lycaon pictus*) is a medium-sized (18-34kg),
83 endangered carnivore that lives in scattered populations in sub Saharan Africa (Fig.
84 1A). The species is the only surviving member of a lineage of wolf-like canids [6].

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
85 Wild dogs have been subject to intense recovery efforts across their range [7, 8], but
86 their global population is decreasing. It is estimated that only 6,600 adult wild dogs
87 remain in 39 subpopulations [9]. The primary reasons for the species' population
88 decline include habitat loss and fragmentation, as well as anthropogenic mortality
89 (e.g. snaring, persecution, road kills, exposure to infectious diseases from domestic
90 dogs) when they range beyond the borders of protected areas [7, 8, 10]. Due to their
91 large ranges and low population densities, African wild dogs are more susceptible to
92 these threats than most other carnivore species [8]. In addition, their complex social
93 system and susceptibility to Allee effects appears to increase the species extinction
94 risk [11, 12]. The dogs are obligate cooperative breeders which form packs
95 consisting of an alpha male and female, their adult siblings, and pups and subadults
96 from the dominant pair [13]. Subadults that have reached reproductive age disperse
97 in single sex groups and form new packs by joining dispersing groups from the
98 opposite sex [14]. Pack members rely on each other for hunting, breeding, and
99 defense against natural enemies and pack size has been found to be significant for
100 hunting and breeding success [13, 15, 16]. When pack size becomes critically low,
101 e.g. due to anthropogenic mortality, this dependence on helpers increases the risk of
102 pack extinction and reduces the number of successful dispersals (Courchamp,
103 Clutton-Brock [12], but see Creel and Creel [17]).

42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
104 Prior genetic studies on wild dogs using a combination of mitochondrial,
105 microsatellite, and MHC markers have resulted in varying estimates of the start of the
106 species decline on the African continent [18, 19]. Consistent with expectation, the
107 data shows strong structuring between populations due to habitat fragmentation and
108 isolation, as well as low genetic diversity within populations [19, 20]. For species that
109 are experiencing such rapid and alarming declines, estimates that are particularly
110 important for management decisions, such as local adaptation, effective population
111 size, and inbreeding are greatly improved by the use of whole-genome methods.
112 Recently, Campana and colleagues [21] sequenced low-coverage genomes of two

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

113 African wild dog individuals from Kenya and South Africa, respectively, to investigate
114 demographic history and signatures of selection of these two separate populations.
115 By mapping these data to the domestic dog genome, they discovered approximately
116 780,000 single nucleotide polymorphisms (SNPs) between their two individuals
117 which could be used to develop SNP typing for the two populations. However, given
118 the low coverage of their genomes (5.7-5.8x average coverage) and the small
119 number of individuals, additional sequencing will be needed to verify the authenticity
120 of those SNPs. Further, important structural variation can be overlooked when
121 mapping against a reference genome from a different genus, and mapping can be
122 hindered if the divergence is high between the sample and the reference (see e.g.
123 Shapiro and Hofreiter [22]). The groups containing the African wild dog and the
124 domestic dog are estimated to have split approximately 7.5-10 Mya and furthermore,
125 the domestic dog has undergone significant genomic selection in recent time [23].

126 Despite the ever-declining cost to sequence DNA, the routine use of genomic
127 approaches in conservation is still far from a reality. One of the major remaining
128 barriers is the lack of reference genomes for species of conservation concern.
129 Generating a *de novo* reference genome requires the sequencing and assembly of
130 the 100s of millions to billions of base-pairs that make up a genome. The first
131 mammalian genome (human) required a massive collaboration between hundreds of
132 scientists and nearly \$3 billion US dollars (1990-2001; [24, 25]). Fortunately, the cost
133 to sequence DNA is now low enough that every base-pair in a typical mammalian
134 genome can be sequenced to high coverage for a few thousand US dollars.
135 However, these low cost sequencing methods produce very short sequences of 150-
136 300 base-pairs in length (for a review on sequencing methods see Goodwin,
137 McPherson [26]). Because large proportions of typical mammal genomes consist of
138 repetitive sequences, it has been impossible to assemble highly-contiguous
139 genomes from only these short sequences. In order to achieve higher continuity,
140 more elaborate and expensive library preparation or alternative sequencing

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

141 technologies have to be used [26, 27]. Among others, these include mate-pair
142 libraries, chromatin folding based libraries, such as cHiCago [28] or HiC [29], and
143 long-read sequencing technologies, such as Pacific Biosciences and Oxford
144 Nanopore Technology. While the resulting genomes can show high continuity, those
145 methods substantially increase the costs of sequencing projects and thus can hinder
146 the generation of genomes for conservation biology purposes.

147 Here we report the use of the Chromium system developed by 10x Genomics
148 [30], a genomic library preparation technique that facilitates cost-effective (around
149 \$2,500) assemblies using short sequencing reads, to assemble three African wild
150 dog genomes. In brief, the 10x Genomics Chromium system is based on dilution of
151 high molecular weight (HMW) DNA. It uses as little as 1ng of input DNA, which is
152 well-suited for a variety of applications. During library preparation, gel beads, so-
153 called GEMs, are mixed with DNA and polymerase for whole-genome amplification.
154 Each gel bead has primer oligos (44nt long) attached to its surface. These contain a
155 priming site (22nt partial R1), a 16nt barcode region, and a 6nt N-mer region that
156 binds to different places on the original DNA fragment. The low amount of input DNA
157 ensures that each gel bead only binds a single (up to ~100kb) DNA fragment. In the
158 next step, amplification of short reads along the original DNA fragment is performed
159 within each gel bead. In most cases, this amplification results in spotted read
160 coverage along the fragment. However, all reads from a respective GEM contain
161 identical barcodes and can later be assigned to groups originating from the same
162 DNA molecule. The information about which molecule of DNA the sequence
163 originated from greatly increases the ability to identify the location of repetitive
164 sequences. The library is then sequenced on an Illumina platform and the raw read
165 data is assembled by the 10x Genomics Supernova assembler. This assembler is
166 very user-friendly and does not require any prior knowledge about input parameters
167 for the assembly.

168 We *de novo* assembled three African wild dog genomes using the 10x
169 Genomics Chromium platform in order to investigate whether this technology is
170 suitable for conservation genomic purposes. For any endangered species, a genome
171 can have large conservation impacts, but high-quality genomes have historically
172 been costly or impossible due to the sampling requirements and in addition,
173 downstream analyses can be challenging. Thus, in order for it to be useful for
174 conservation purposes the technology needs to be (a) cost-effective and (b) user-
175 friendly. Furthermore, we test the 10x Genomics Chromium based assemblies for
176 reproducibility, continuity, conserved gene completeness, and repetitive content, as
177 compared to the previously published domestic dog genome.

178

179 **Data Description & Analyses**

180

181 *Assembly of the African wild dog genome*

182 Using 10x Genomics Chromium technology, we generated DNA libraries for
183 three African wild dog individuals, two of which were collected from a wild pack in the
184 Hwange National Park, Zimbabwe and are presumed to be sisters (named Sister 1
185 and Sister 2), and a third unrelated individual from the Endangered Wolf Center,
186 Eureka, Missouri (named Eureka). A summary of the assembly statistics output by
187 the Supernova assembler can be found in Table 1 (detailed statistics for each
188 genome assembly can be found in Supporting Information Table 1). We generated
189 1,200 million paired-end reads for Sister 1, 801.56 million reads for Sister 2, and
190 427.6 million reads for Eureka. We then used the reads to assemble each genome
191 using the 10x Genomics Supernova assembler (as explained in
192 <https://support.10xgenomics.com/de-novo-assembly/software/overview/welcome>).
193 The mean input DNA molecule length reported by the Supernova assembler for
194 Sister 1 was 19.91kb, Sister 2 was 77.03kb, and Eureka was 52.00kb. All three
195 assemblies corroborate a genome size of approximately 2.3Gb, which is similar to

196 that of the domestic dog (2.4Gb). These three assemblies together constitute the first
197 reported *de novo* assemblies for the African wild dog species.

198 We then calculated the scaffold and contig N50 statistics, which are indicative
199 of assembly continuity. The Sister 1 assembly resulted in a contig and scaffold N50
200 of 61.34 kb and 7.91 Mb, respectively, the Sister 2 assembly achieved 83.47 kb
201 contig and 21.34 Mb scaffold N50s, and finally the Eureka assembly had 50.15 kb
202 contig and 15.31 Mb scaffold N50s (Table 1). While our contig and scaffold N50's are
203 smaller than the ones from the most recent dog genome (267kb and 45.9Mb,
204 respectively), they are still larger than most mammalian genomes assembled that
205 used only short read data (see e.g. Lok, Paton [31]).

207 *Conserved Genes*

208 The program BUSCO (Benchmarking Universal Copy Orthologs) uses highly
209 conserved single copy orthologous genes from a number of different taxa and groups
210 in order to test assemblies (both genomic and transcriptomic) for gene
211 completeness, fragmentation, or absence as an indicator of assembly quality. Using
212 BUSCO v2 on our assemblies, we found that the most continuous assembly, Sister
213 2, completely recovered 95.1% of conserved genes (Mammalia gene set; Table 2).
214 Sister 1 and Eureka recovered 95.4% and 93.3% of complete conserved genes,
215 respectively. Using the same analysis, we found 95.3% of complete conserved
216 genes in the latest dog assembly (canFam3.1). This indicates that although the
217 domestic dog assembly is more continuous overall, our assemblies recover nearly
218 the same or even higher number of conserved genes. Surprisingly, Sister 1 had the
219 least number of missing genes out of all the assemblies assessed, despite lower
220 continuity than Sister 2. We also ran BUSCO on the Hawaiian monk seal genome,
221 generated through the combination of 10x Genomics Chromium and Bionano
222 Genomics Irys data, and found it recovered 94.6% of conserved genes using
223 BUSCO. This suggests that using Bionano in addition to 10x does not greatly

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

224 improve the reconstruction of the gene regions. However, the Hawaiian monk seal
225 genome has a scaffold N50 of approximately 28Mb, so Bionano may improve the
226 overall assembly continuity compared to 10x Genomics alone. The low coverage
227 genomes from Campana et al. 2016 achieved a BUSCO score of 92.8% for the
228 individual from Kenya and 94.8% for the individual from South Africa. These scores
229 are similar to those from the dog assembly, which was the reference the reads were
230 mapped to initially.

231

232 *Repeat annotation*

233 We identified repetitive regions of the genome in order to discern how well
234 these complex areas were assembled by the 10x Genomics Chromium technology.
235 Using both RepeatMasker and RepeatModeler, we found that for all three wild dog
236 assemblies, total repeat content was evaluated to be within 3% of one another, which
237 indicates consistency among assemblies from a single species (Supporting
238 Information Table S2). No single repeat category was disproportionately affected
239 during repeat annotation of the three genomes, which suggests that assembly quality
240 was likely the most influential factor. Furthermore, repeat content of all wild dog
241 assemblies was qualitatively similar to canFam3.1. As repetitive regions tend to be
242 the most difficult regions to assemble, the similarity in repeat content between the
243 African wild dog compared to that of the domestic dog, highlights the value of using
244 10x Genomics Chromium technology to produce accurate and continuous
245 assemblies.

246

247 *Gene annotation*

248 The genome annotation pipeline Maker3 resulted in very similar numbers of
249 annotated genes between all three African wild dog individuals and the domestic
250 dog. Annotations ranged from 20,649 (Sister 2) to 20,946 (Sister 1) genes
251 (Supporting Information Table S3). Using proteinortho to detect orthologous genes

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

252 between individuals and paralogous genes within individuals, we found 12,617
253 one:one orthologs present in all three individuals and 6,462 one:one orthologs in two
254 out of the three individuals. We found 268 multi copy genes present in all three
255 individuals and 37 not present in one individual. Overall, the number of annotated
256 genes was comparable to those found in the dog genome (Supporting Information
257 Table S3).

258

259 *Variant rates*

260

261 We found a high number of heterozygous sites to be shared between all three
262 individuals (321k; here we report the heterozygous sites called using a posterior
263 probability cutoff of 0.99; Fig. 1B). As expected, Sister 1 and Sister 2 share more
264 heterozygous sites (344k) than either sister with Eureka (168k and 170k, for Sister 1
265 and Sister 2, respectively). Each individual shows a high number of singletons
266 (heterozygous sites only found in one individual), with Sister 2 showing the highest
267 number (1,100k), followed by Sister 1 (968k) and Eureka (825k). Even if we include
268 the two low coverage genomes from Campana, Parker [21], we find a high number of
269 shared heterozygous sites between all individuals (134k; Supporting Information
270 Figure S1). As expected, we see a higher number of singletons in these two
271 individuals, due to the lower reliability of the genotype calls caused by the low
272 coverage (false positives caused by sequencing errors). We estimated a per site
273 heterozygosity of 0.0008 to 0.0012 for Sister 1, 0.0009 to 0.0012 for Sister 2, and
274 0.0007 to 0.001 for Eureka using posterior cutoffs for genotype calls from 0.95 to 1 in
275 ANGSD (Supporting Information, Fig. S1C). As can be seen in Figure 1C, except for
276 a posterior probability cutoff of 1, where Sister 1 shows the highest heterozygosity,
277 Sister 2 always shows the highest, Sister 1 the second highest and Eureka the
278 lowest heterozygosity. Interestingly, Eureka shows a lower heterozygosity than the
279 other two assemblies, even though its parents originated from South Africa and

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

280 Botswana. Our estimates show that, while being heavily threatened, African Wild
281 dogs seem to still retain a relatively high within individual heterozygosity. We did not
282 see any major difference between heterozygosity estimates from repeat-masked and
283 unmasked genomes (data not shown). The Supernova software estimated a
284 heterozygous position every 2.6kb, 3.1kb, and 7.14kb for Sister 1, Sister2, and
285 Eureka, respectively (Supporting Information Table S4). On the contrary, estimates
286 based on genotype calls using ANGSD showed much more frequent heterozygous
287 positions (850bp - 1.2kb, 814bp - 1.1kb and 999bp - 1.5kb depending on the
288 posterior cutoff used; Supporting Information Table S4).

289

290 **Discussion**

291

292 *Assembly continuity and quality*

293 All three African wild dog assemblies produced with 10x Genomics Chromium
294 data showed high continuity, high recovery rates of conserved genes, and expected
295 proportions of repetitive sequence; indicating that they are high-quality assemblies.
296 The Sister 2 assembly, which has the highest mean molecule length, is also the most
297 continuous (Contig N50: 83.47kb, Scaffold N50: 21.34Mb; Table 1). Interestingly, the
298 Sister 1 genome has a higher contig N50 (61.34kb) than Eureka (50.15kb), but a
299 lower scaffold N50 (7.91Mb and 15.31Mb, respectively). This may indicate that input
300 molecule length is a key factor for scaffolding, while coverage is a key factor for
301 contig assembly. Despite having the highest continuity of all three assemblies, Sister
302 2 did not show the highest BUSCO completeness scores (see Table 2), although the
303 differences were minor and likely not meaningful as they could lie well within the
304 uncertainty of the BUSCO analysis (with 95.1% complete BUSCOs compared to
305 95.4% for Sister 1). Sister 1 achieved the highest BUSCO scores, even compared to
306 the latest domestic dog genome assembly (CanFam3.1; 95.2%), which has three
307 times higher contig N50 and an almost six times higher scaffold N50. The high

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

308 scores are remarkable for the limited number of reads used for the assemblies (as
309 low as 25x coverage). As expected, Sister 2, which showed the highest continuity
310 also had the highest repeat content (see Supporting Information Table S2). However,
311 all three assemblies resulted in similar repeat contents in terms of repeat
312 composition as well as overall percentage (within 3% of each other), with the most
313 continuous assembly (Sister 2) showing the highest number of repeats. Repeat
314 composition in the African wild dog genomes was also similar to the domestic dog.

315 All assemblies yielded similar amounts of genes, with Sister 1 showing the
316 highest number (see Supporting Information Table S3), which reflects its BUSCO
317 scores. Closer investigations of one:one and one:many orthologs further showed a
318 very good agreement between annotations obtained from all three individuals. The
319 numbers of annotated genes for all three African wild dogs were similar to those
320 calculated for the latest domestic dog assembly.

321 322 *10x Genomics Chromium system: Feasibility and caveats*

323
324 Most mammal genomes published in the last several years use a mixture of
325 paired-end (PE) and multiple mate pair (MP) Illumina libraries (e.g. Lok, Paton [31]
326 and Liu, Lorenzen [32]). While often resulting in good continuity (e.g. Liu, Lorenzen
327 [32] or Huang, Zhao [33]), using different insert libraries considerably increases the
328 cost per genome. On the contrary, 10x Genomics Chromium allows for assembly of a
329 comparable or even more continuous genome using only a single library for a
330 fraction of the cost (see below). Furthermore, as we show here, this library
331 technology generates high-quality assemblies from as low as 25x coverage (see
332 Eureka assembly), while the recommended coverage for PE plus MP assemblies is
333 100x [34]. Recently, Mohr and colleagues [35] presented a highly continuous
334 assembly of the endangered Hawaiian Monk seal (~2.4Gb total genome assembly
335 length) using a combination of 10x Genomics Chromium and Bionano Genomics

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

336 optical mapping. Interestingly, their 10x Genomics Chromium assembly showed
337 similar N50 statistics to those reported here (scaffold N50 22.23Mb), showing that
338 10x Genomics Chromium technology alone enables the generation of high-quality
339 mammalian genome assemblies.

340 A limitation of 10x Genomics Chromium technology is the requirement of
341 fresh tissue samples for the isolation of HMW DNA. This can be difficult or
342 impossible to obtain from some endangered species. Fortunately, small amounts of
343 mammalian blood yield sufficient amounts of HMW DNA when properly stored.
344 Additionally, DNA extraction kits such as the Qiagen MagAttract kit can extract
345 sufficient amounts of HMW DNA from as little as 200µl. For museum samples, or
346 tissues stored for extended periods of time, reference-based mapping might be the
347 only option to extract long-range genomic information. However, for extant
348 endangered species, especially those with individuals in captivity, 10x Genomics
349 Chromium offers a cost-effective approach to sequence genomes. For species with
350 genome sizes <1Gb and between ~3Gb and 5.8Gb special data processing will need
351 to be applied (see [https://support.10xgenomics.com/de-novo-assembly/sample-](https://support.10xgenomics.com/de-novo-assembly/sample-prep/doc/technical-note-supernova-guidance)
352 [prep/doc/technical-note-supernova-guidance](https://support.10xgenomics.com/de-novo-assembly/sample-prep/doc/technical-note-supernova-guidance)). In addition, the amplification primers
353 for the 10x Chromium library preparation are designed for GC contents similar to
354 human (~41%), implying that the method might not work as well for genomes that
355 strongly divert from this GC content (e.g. for some invertebrates).

356 357 *Cost effectiveness*

358
359 Sequencing costs are steadily dropping. At the time the sequencing for this
360 project was carried out a lane on the Illumina HiSeqX cost approximately \$1,500 -
361 \$2,000 and a 10x Genomics library ranged from \$450 to \$1000, thus allowing the
362 generation of high quality *de novo* genomes for less than \$3,000 total (prices
363 obtained from US sequencing facilities). Furthermore, prices are likely to decline as

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

364 technology improves. Even more so, independent of sequencing lane costs, this
365 method only requires a single library to be sequenced to an average coverage of 25 -
366 75x, unlike other methods which require multiple libraries at higher coverage. As we
367 have shown here, a continuous assembly can be generated from as little as 25x.
368 Furthermore, computational resources required to assemble the genome are very
369 low. The current version of Supernova 1.2 only requires a minimum of 16 CPU cores
370 and 244Gb of memory (for a human genome at 56x coverage;
371 <https://www.10xgenomics.com/>), and the assembly can be carried out in only few
372 days (depending on the number of available CPU cores). This is about a reduction of
373 five times the memory requirement compared to the first version of Supernova. Even
374 more so, Supernova does not require parameter input or tuning, thus allowing even
375 novices to easily assemble 10x Genomics Chromium based genomes.

376

377 *Applications in conservation*

378

379 Traditionally, conservation biologists have obtained a great deal of genetic
380 information from a few microsatellite markers and/or nuclear and mitochondrial loci.
381 The analysis of microsatellite markers can provide a snapshot into contemporary
382 population structure, but this method risks providing incomplete information on
383 selection and migration and can be an unreliable way to identify individuals from
384 degraded low-quality DNA samples (such as scat) due to the stochastic behavior of
385 marker amplification (allelic dropout; Frantzen, Silk [36]; Taberlet and Luikart [37] ;
386 Morin, Luikart [38]). Moreover, microsatellites can be difficult to successfully design
387 and develop, which can quickly increase costs for species that have little to no
388 genetic information available. The ability to rapidly and cost-effectively generate full
389 genomes will allow conservation biologists to bridge this gap and harvest crucial fine-
390 scale population information for population parameters such as inbreeding (e.g.
391 Vieira, Fumagalli [39]), load of deleterious mutations (e.g. Robinson, Ortega-Del

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

392 Vecchy [40]), gene flow (e.g. Pazmiño, Maes [41]) and population structure (e.g.
393 Hampton, Spencer [42]). Once a reference genome has been assembled, optional
394 (low coverage) re-sequencing data from several individuals allows for the typing of
395 genome-wide information such as single-nucleotide polymorphisms (SNPs),
396 potentially neutral microsatellite loci, and other genomic regions of interest. These
397 data can then be used to investigate the abovementioned population parameters, but
398 also further yield insights into adaptive genetic variation and perhaps the adaptive
399 potential of different populations or species.

400

401 *Heterozygosity within African wild dog individuals*

402

403 A high number of heterozygous sites were shared between all three
404 individuals in this study, with Sister 1 and Sister 2 sharing more heterozygous sites
405 than either with Eureka. Each of the individuals further shows a high number of
406 singletons (heterozygous sites only found in one individual). Even when compared to
407 the two low coverage genomes from Campana et al. (2016) we find a high number of
408 shared sites. As expected, we see a much higher rate of singletons in these two
409 individuals. Due to the low coverage (5.7 - 5.8x average coverage) we predict a
410 higher proportion of the called heterozygous sites to be false positives due to
411 sequencing errors. Heterozygosity per site estimates indicate a high within individual
412 diversity. Estimates ranged from 0.0007 - 0.001 for Eureka to 0.0009 - 0.0012 for
413 Sister 2, which are similar to those obtained for lions (0.00074 – 0.00148) and tigers
414 (0.00087 – 0.00104) [45]. Intriguingly, other threatened large bodied carnivores, such
415 as the Iberian lynx (*Lynx pardinus*), the cheetah (*Acinonyx jubatus*), and the island
416 fox (*Urocyon littoralis*) show nearly 10 fold lower heterozygosity (0.0001 [43], 0.0002
417 [44] and 0.000014 - 0.0004 [40], respectively). The high within-individual
418 heterozygosity could be a result of their social structure, as only unrelated individuals
419 come together to form new packs through dispersal. However, Hwange National

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

420 Park is considered to be a part of the most continuous population of African wild
421 dogs which may explain the high heterozygosity of Sister 1 and Sister 2 (Girman et
422 al. 2001). Further sequencing of other populations will be needed to assess whether
423 the high within-individual heterozygosity is a range-wide phenomenon in African wild
424 dogs. If true, this could be very good news for the survival of these species if external
425 pressures (such as hunting, habitat fragmentation, etc.) can be reduced.

426 The Supernova software reports distance between heterozygous site
427 estimates (see Supporting Information Table S1). Interestingly, those estimates were
428 much lower than the ones obtained based on the genotype calls produced with
429 ANGSD. While Supernova estimated this distance to be 2.6kb in Sister 1, 3.1kb in
430 Sister 2 and 7.1kb in Eureka, the ANGSD based estimates range from 850bp - 1.2kb
431 for Sister 1, 814bp - 1.1kb for Sister 2 and 999bp - 1.5kb for Eureka, depending on
432 the posterior cutoff used. Supernova calculates the distance between heterozygous
433 sites as part of the assembly process, however, when the fasta consensus sequence
434 is called part of the variation can get flattened (see Weisenfeld, Kumar [30]). This
435 phenomenon is typically seen in regions between megabubbles, which are nominally
436 homozygous, but could actually have some variation that cannot be phased by
437 Supernova. We also note that heterozygosity values obtained using genotype calls in
438 ANGSD could also be biased, as they are based on the nominal and not the effective
439 coverage. The nominal coverage is the total number of reads that cover a site in the
440 assembly, whereas for the effective coverage only reads from different barcodes are
441 included in the estimation. If individual barcoded regions amplified with different
442 efficiency during the library preparation step, then heterozygosity estimates could be
443 unreliable. However, this should not strongly affect genome-wide heterozygosity
444 estimates, as we expect this issue to be rare. Heterozygosity

445

446 **Potential Implications**

447

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

448 We find that the 10x Genomics Chromium system can be used to assemble
449 highly continuous and accurate mammalian genome assemblies for less than \$3,000
450 US dollars per genome (sequenced 2016 and 2017). The method can be easily
451 applied to species of conservation concern for which genomic methods could greatly
452 benefit their management and monitoring programs. For the African wild dog, these
453 genomes will facilitate more reliable and cost-effective conservation efforts through
454 the use of re-sequencing and SNP-typing methods. Compared to other species of
455 conservation concern, the African wild dog has a relatively high heterozygosity. More
456 studies are required to understand how both the social biology and recent precipitous
457 population declines have impacted the population genomic structure of African wild
458 dogs, and how management might use this information for the benefit and longevity
459 of the species.

460

461 **Methods**

462

463 Detailed Methods can be found in Supporting Information.

464

465 *Samples*

466 Blood samples from two individuals belonging to the same pack in Hwange
467 National Park, Zimbabwe were provided by Painted Dog Conservation (CITES
468 Export permit: ZW/0842/2015, ESA import permit: MA66259B-0, Research Council of
469 Zimbabwe permit: 02553). These individuals were presumed to be sisters from direct
470 observation of their litter at the den (here, named Sister 1 and Sister 2). DNA was
471 extracted two weeks after storage at -80°C. The third sample was provided by the
472 Endangered Wolf Center, Eureka, Missouri from a captive born individual (here
473 named Eureka). DNA was extracted 9 days after the sample was taken. Though the
474 Chromium library preparation does not require large amounts of DNA, the DNA
475 should have a mean molecule length > 200kb (high-molecular weight, or HMW). DNA

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

476 from all individuals was extracted from blood samples using the QIAGEN MagAttract
477 HMW DNA kit following the provided instructions.

478

479 *Genome Assembly*

480 We constructed one sequencing library per individual using the 10x
481 Genomics Chromium System with 1.2ng of HMW input DNA. All libraries were then
482 sequenced on the Illumina HiSeqX (Sister 2, Eureka) or HiSeq 4000 (Sister 1)
483 platform. We subsequently assembled the three genomes using the 10x Genomics
484 genome assembler Supernova 1.1.1 Weisenfeld, Kumar [30];
485 <http://support.10xgenomics.com/de-novo-assembly/software/overview/welcome>)
486 using default assembly parameters.

487

488 *Assembly Quality Assessment*

489 We used the Supernova assembler as well as QUAStv4.3 to determine
490 continuity statistics, such as the scaffold N50 and the total number of scaffolds [46].
491 We further applied the program BUSCO v2 [47] to assess the presence of nearly
492 universal lineage specific single-copy orthologous genes in our assemblies using the
493 mammalian gene set from OrthoDB v9 (4104 genes; available at
494 <http://busco.ezlab.org>). We compare these results to the high-quality canFam3.1
495 assembly of the domestic dog (Hoeppner, Lundquist [48]; *Canis familiaris*). The
496 canFam3.1 assembly was built on 7x coverage of Sanger reads and BAC end
497 sequencing and has a scaffold N50 of 46Mb. We also inferred the number of
498 BUSCO's in the recently published Hawaiian monk seal genome (which was
499 assembled using a combination of 10x Genomics Chromium and Bionano Genomics
500 lrys data) and the two previously published African wild dog genomes (sequenced
501 with basic short read Illumina technology at low coverage and assembled using the
502 domestic dog; [21]).

503

504 *Repeat Identification and Masking*

1
2 505 We next identified repetitive regions in the genomes as another comparative
3
4 506 measure of assembly quality and to prepare the genome for annotation. Repeat
5
6 507 annotation was carried out using both homology-based and *ab-initio* prediction
7
8 508 approaches. We used the canid RepBase (<http://www.girinst.org/repbase/>; [49])
9
10 509 repeat database for the homology-based annotation within RepeatMasker
11
12 510 (<http://www.repeatmasker.org>). We then carried out *ab-initio* repeat finding using
13
14 511 RepeatModeler (<http://repeatmasker.org/RepeatModeler.html>).
15
16 512

17
18
19 513 *Gene Annotation*

20
21
22 514 Gene annotation for the three assemblies was performed with the genome
23
24 515 annotation pipeline Maker3 [50], which implements both *ab-initio* prediction and
25
26 516 homology-based gene annotation by leveraging previously published protein
27
28 517 sequences from dog, mouse, and human.

29
30
31 518 Orthologous genes between the three African wild dog assemblies, as well as
32
33 519 paralogous genes within each individual, were inferred using proteinortho [51].
34
35 520 Proteinortho applies highly parallelized reciprocal blast searches to establish
36
37 521 orthology and paralogy for genes within and between gene annotation files.
38
39 522

40
41
42 523 *Variant rates*

43
44 524 In order to estimate within individual heterozygosity, we selected a single
45
46 525 pseudo-haplotype (in cases where genomic regions were phased into haplotypes,
47
48 526 one of the two was chosen randomly) from Sister 2 to represent the reference
49
50 527 sequence. Next we mapped the raw reads from all three individuals to the reference
51
52 528 using bwa mem [52]. We then converted the resulting sam files to bam format using
53
54 529 samtools [53], and sorted and indexed them using picard
55
56 530 (<http://broadinstitute.github.io/picard/>). Realignment around insertion/deletion (indel)
57
58 531 regions was performed using GATK, and finally, we called heterozygous sites using
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

532 a probabilistic framework implemented in ANGSD [54]. We tested different posterior
533 probability cutoffs (1, 0.999, 0.99 and 0.95). To allow for comparison between all
534 individuals, we down-sampled all individuals to 20x mean nominal coverage (total
535 number of reads covering a position, independent of their barcode) for our analyses.
536 Heterozygosity was then simply calculated as the ratio of variable sites to the total
537 number of sites (variable and invariable). Furthermore, Supernova outputs the
538 distance between heterozygous sites as part of their assembly report. We further
539 downloaded the read data of Campana, Parker [21] and mapped them against our
540 Sister 2 assembly to compare heterozygosity estimates (using the approach outlined
541 above). Next, we estimated the number of shared heterozygous sites between a) our
542 individuals and b) our individuals and the two from Campana, Parker [21]. To do so,
543 we used the *gplots* library in R (<https://www.r-project.org>) to calculate the overlap
544 between the three sets and to display them in a Venn diagram.

545

546 **Availability of supporting data**

547 The data sets supporting the results of this article will be uploaded to the GigaDB
548 repository pending manuscript acceptance.

549

550 **Supporting Information**

551 Detailed information on methods, Supernova output, repeat annotation, gene
552 annotation, heterozygosity calculations, and different posterior probability cutoffs are
553 available online. The authors are solely responsible for the content and functionality
554 of these materials. Queries (other than absence of the material) should be directed to
555 the corresponding author.

556

557 **Competing Interests**

558 Author J. Stuelpnagel is a board member of 10x Genomics Inc. Author Ryan W.
559 Taylor is founder of End2End Genomics Inc.

560

1
2 **561 Authors' contributions**
3

4 562 Authors JS, CSZ, PB, SP, EA, and DP conceived the project. Authors EM, HM, OM,
5
6 563 and RMC contributed samples and insight to the project. RT assembled the genomes.
7
8 564 EA and SP performed the genome annotation and downstream analyses. EA, SP,
9
10 565 CST, DP, and RT wrote the paper. All authors read and approved the final
11
12 566 manuscript.
13

14
15 567
16

17 **568 Acknowledgements**
18

19
20 569 We thank M. Agnew, C. Asa, L. Padilla, and W. Warren for assistance in obtaining
21
22 570 the Eureka sample. T. Linderoth, T. Korneliussen, and K. Bi for help with the different
23
24 571 heterozygosity calculations. D. Church from 10x Genomics for discussion on how
25
26 572 SuperNova performs the heterozygous site calling. This work was funded by the a
27
28 573 donation to the Program for Conservation Genomics at Stanford University.
29

30
31 574
32

33 **575 Literature Cited**
34

35 576
36

- 37
38 577 1. Pimm, S.L., et al., *The biodiversity of species and their rates of extinction,*
39
40 578 *distribution, and protection.* Science, 2014. **344**(6187): p. 1246752.
41
42 579 2. Ceballos, G., et al., *Accelerated modern human-induced species losses:*
43
44 580 *Entering the sixth mass extinction.* Science Advances, 2015. **1**(5): p.
45
46 581 e1400253.
47
48 582 3. Epstein, B., et al., *Rapid evolutionary response to a transmissible cancer in*
49
50 583 *Tasmanian devils.* Nature communications, 2016. **7**: p. 12684.
51
52 584 4. Steiner, C.C., et al., *Conservation genomics of threatened animal species.*
53
54 585 *Annu. Rev. Anim. Biosci.*, 2013. **1**(1): p. 261-281.
55
56 586 5. Shafer, A.B., et al., *Genomics and the challenging translation into*
57
58 587 *conservation practice.* Trends in Ecology & Evolution, 2015. **30**(2): p. 78-87.
59
60
61
62
63
64
65

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- 588 6. Girman, D., et al., *Molecular genetic and morphological analyses of the African wild dog (Lycaon pictus)*. Journal of heredity, 1993. **84**(6): p. 450-459.
 - 590 7. Woodroffe, R., J. Ginsberg, and D.W. Macdonald, *The African wild dog: status survey and conservation action plan*1997: IUCN.
 - 592 8. IUCN/SSC, *Regional conservation strategy for the cheetah and African wild dog in Southern Africa*, 2007, IUCN Species Survival Commission Gland.
 - 594 9. Woodroffe, R. and C. Sillero-Zubiri, *Lycaon pictus*. The IUCN Red List of Threatened Species 2012 (Downloaded August 2017), 2012. **e.T12436A167111116**.
 - 597 10. Woodroffe, R. and J.R. Ginsberg, *Edge effects and the extinction of populations inside protected areas*. Science, 1998. **280**(5372): p. 2126-2128.
 - 599 11. Courchamp, F., T. Clutton-Brock, and B. Grenfell, *Inverse density dependence and the Allee effect*. Trends in Ecology & Evolution, 1999. **14**(10): p. 405-410.
 - 602 12. Courchamp, F., T. Clutton-Brock, and B. Grenfell. *Multipack dynamics and the Allee effect in the African wild dog, Lycaon pictus*. in *Animal Conservation forum*. 2000. Cambridge University Press.
 - 605 13. McNutt, J.W. and J.B. Silk, *Pup production, sex ratios, and survivorship in African wild dogs, Lycaon pictus*. Behavioral Ecology and Sociobiology, 2008. **62**(7): p. 1061-1067.
 - 608 14. McNutt, J.W., *Sex-biased dispersal in African wild dogs, Lycaon pictus*. Animal behaviour, 1996. **52**(6): p. 1067-1077.
 - 610 15. Fanshawe, J.H. and C.D. Fitzgibbon, *Factors influencing the hunting success of an African wild dog pack*. Animal behaviour, 1993. **45**(3): p. 479-490.
 - 612 16. Creel, S. and N.M. Creel, *Six ecological factors that may limit African wild dogs, Lycaon pictus*. Animal Conservation, 1998. **1**(1): p. 1-9.
 - 614 17. Creel, S. and N.M. Creel, *Opposing effects of group size on reproduction and survival in African wild dogs*. Behavioral Ecology, 2015. **26**(5): p. 1414-1422.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- 616 18. Girman, D., et al., *Patterns of population subdivision, gene flow and genetic*
617 *variability in the African wild dog (Lycaon pictus)*. *Molecular Ecology*, 2001.
618 **10**(7): p. 1703-1723.
- 619 19. Marsden, C.D., et al., *Spatial and temporal patterns of neutral and adaptive*
620 *genetic variation in the endangered African wild dog (Lycaon pictus)*.
621 *Molecular Ecology*, 2012. **21**(6): p. 1379-1393.
- 622 20. Marsden, C.D., et al., *Highly endangered African wild dogs (Lycaon pictus)*
623 *lack variation at the major histocompatibility complex*. *Journal of heredity*,
624 2009. **100**(suppl_1): p. S54-S65.
- 625 21. Campana, M.G., et al., *Genome sequence, population history, and pelage*
626 *genetics of the endangered African wild dog (Lycaon pictus)*. *BMC genomics*,
627 2016. **17**(1): p. 1013.
- 628 22. Shapiro, B. and M. Hofreiter, *A paleogenomic perspective on evolution and*
629 *gene function: new insights from ancient DNA*. *Science*, 2014. **343**(6169): p.
630 1236573.
- 631 23. Nyakatura, K. and O.R. Bininda-Emonds, *Updating the evolutionary history of*
632 *Carnivora (Mammalia): a new species-level supertree complete with*
633 *divergence time estimates*. *BMC biology*, 2012. **10**(1): p. 12.
- 634 24. Lander, E.S., et al., *Initial sequencing and analysis of the human genome*.
635 2001.
- 636 25. Hayden, E.C., *The \$1,000 genome*. *Nature*, 2014. **507**(7492): p. 294.
- 637 26. Goodwin, S., J.D. McPherson, and W.R. McCombie, *Coming of age: ten*
638 *years of next-generation sequencing technologies*. *Nature Reviews Genetics*,
639 2016. **17**(6): p. 333-351.
- 640 27. Ekblom, R. and J.B. Wolf, *A field guide to whole - genome sequencing,*
641 *assembly and annotation*. *Evolutionary applications*, 2014. **7**(9): p. 1026-1042.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- 642 28. Putnam, N.H., et al., *Chromosome-scale shotgun assembly using an in vitro*
643 *method for long-range linkage*. Genome Research, 2016. **26**(3): p. 342-350.
- 644 29. Burton, J.N., et al., *Chromosome-scale scaffolding of de novo genome*
645 *assemblies based on chromatin interactions*. Nature biotechnology, 2013.
646 **31**(12): p. 1119-1125.
- 647 30. Weisenfeld, N.I., et al., *Direct determination of diploid genome sequences*.
648 Genome Research, 2017. **27**(5): p. 757-767.
- 649 31. Lok, S., et al., *De novo genome and transcriptome assembly of the Canadian*
650 *beaver (Castor canadensis)*. G3: Genes, Genomes, Genetics, 2017. **7**(2): p.
651 755-773.
- 652 32. Liu, S., et al., *Population Genomics Reveal Recent Speciation and Rapid*
653 *Evolutionary Adaptation in Polar Bears*. Cell, 2014. **157**(4): p. 785-794.
- 654 33. Huang, J., et al., *Analysis of horse genomes provides insight into the*
655 *diversification and adaptive evolution of karyotype*. Scientific reports, 2014. **4**.
- 656 34. Gnerre, S., et al., *High-quality draft assemblies of mammalian genomes from*
657 *massively parallel sequence data*. Proceedings of the National Academy of
658 Sciences, 2011. **108**(4): p. 1513-1518.
- 659 35. Mohr, D.W., et al., *Improved de novo Genome Assembly: Linked-Read*
660 *Sequencing Combined with Optical Mapping Produce a High Quality*
661 *Mammalian Genome at Relatively Low Cost*. bioRxiv, 2017: p. 128348.
- 662 36. Frantzen, M., et al., *Empirical evaluation of preservation methods for faecal*
663 *DNA*. Molecular Ecology, 1998. **7**(10): p. 1423-1428.
- 664 37. Taberlet, P. and G. Luikart, *Non-invasive genetic sampling and individual*
665 *identification*. Biological Journal of the Linnean Society, 1999. **68**(1-2): p. 41-
666 55.
- 667 38. Morin, P.A., G. Luikart, and R.K. Wayne, *SNPs in ecology, evolution and*
668 *conservation*. Trends in Ecology & Evolution, 2004. **19**(4): p. 208-216.

- 1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
- 669 39. Vieira, F.G., et al., *Estimating inbreeding coefficients from NGS data: impact*
670 *on genotype calling and allele frequency estimation*. *Genome Research*, 2013.
671 **23**(11): p. 1852-1861.
- 672 40. Robinson, J.A., et al., *Genomic flatlining in the endangered island fox*.
673 *Current Biology*, 2016. **26**(9): p. 1183-1189.
- 674 41. Pazmiño, D.A., et al., *Genome-wide SNPs reveal low effective population*
675 *size within confined management units of the highly vagile Galapagos shark*
676 *(Carcharhinus galapagensis)*. *Conservation Genetics*, 2017: p. 1-13.
- 677 42. Hampton, J.O., et al., *Molecular techniques, wildlife management and the*
678 *importance of genetic population structure and dispersal: a case study with*
679 *feral pigs*. *Journal of Applied Ecology*, 2004. **41**(4): p. 735-743.
- 680 43. Abascal, F., et al., *Extreme genomic erosion after recurrent demographic*
681 *bottlenecks in the highly endangered Iberian lynx*. *Genome biology*, 2016.
682 **17**(1): p. 251.
- 683 44. Dobrynin, P., et al., *Genomic legacy of the African cheetah, Acinonyx jubatus*.
684 *Genome biology*, 2015. **16**(1): p. 277.
- 685 45. Kim, S., et al., *Comparison of carnivore, omnivore, and herbivore mammalian*
686 *genomes with a new leopard assembly*. *Genome biology*, 2016. **17**(1): p. 211.
- 687 46. Gurevich, A., et al., *QUAST: quality assessment tool for genome assemblies*.
688 *Bioinformatics*, 2013. **29**(8): p. 1072-1075.
- 689 47. Simão, F.A., et al., *BUSCO: assessing genome assembly and annotation*
690 *completeness with single-copy orthologs*. *Bioinformatics*, 2015. **31**(19): p.
691 3210-3212.
- 692 48. Hoepfner, M.P., et al., *An Improved Canine Genome and a Comprehensive*
693 *Catalogue of Coding Genes and Non-Coding Transcripts*. *PLoS one*, 2014.
694 **9**(3): p. e91172.
- 695 49. Jurka, J., et al., *Rebase Update, a database of eukaryotic repetitive*
696 *elements*. *Cytogenetic and genome research*, 2005. **110**(1-4): p. 462-467.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

697 50. Holt, C. and M. Yandell, *MAKER2: an annotation pipeline and genome-*
698 *database management tool for second-generation genome projects.* BMC
699 bioinformatics, 2011. **12**(1): p. 491.

700 51. Lechner, M., et al., *Orthology detection combining clustering and synteny for*
701 *very large datasets.* 2014.

702 52. Li, H. and R. Durbin, *Fast and accurate short read alignment with Burrows–*
703 *Wheeler transform.* Bioinformatics, 2009. **25**(14): p. 1754-1760.

704 53. Li, H., et al., *The sequence alignment/map format and SAMtools.*
705 Bioinformatics, 2009. **25**(16): p. 2078-2079.

706 54. Korneliussen, T.S., A. Albrechtsen, and R. Nielsen, *ANGSD: analysis of next*
707 *generation sequencing data.* BMC bioinformatics, 2014. **15**(1): p. 356.

708
709
710
711
712
713

714 **Tables**

715

716 **Table 1. Assembly Statistics.** Assembly statistics for the three African wild dog
 717 genomes reported by the Supernova assembler. Coverage was assessed using
 718 samtools depth.

		Sister 1	Sister 2	Eureka
Input	Reads (m)	1,200	801.56	427.6
	Average coverage	69	46	25
	Mean molecule size (kb)	19.91	77.03	52.00
Contig	N50 (kb)	61.34	83.47	50.15
	Longest (kb)	524.60	615.40	450.50
	Number (k)	78.62	68.64	108.00
Scaffold	N50 (mb)	7.91	21.34	15.31
	Longest (kb)	43.96	69.63	41.67
	Number (k)	11.78	17.64	25.78
Total size (gb)	Scaffolds >= 10kb	2.27	2.26	2.20
	Scaffolds >= 500bp	2.34	2.40	2.42

719

720

721 **Table 2. Conserved Gene Statistics.** Results of the BUSCO v2 gene annotation
 722 from three African wild dog genome assemblies, canFam3.1, low-coverage wild dog
 723 genomes [21], and the recently published Hawaiian Monk seal genome [35].

724

Assembly	Species	Complete	Single copy	Duplicated	Fragmented	Missing	Total searched
Sister 1	<i>L. pictus</i>	3914	3875	39	102	88	4104

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

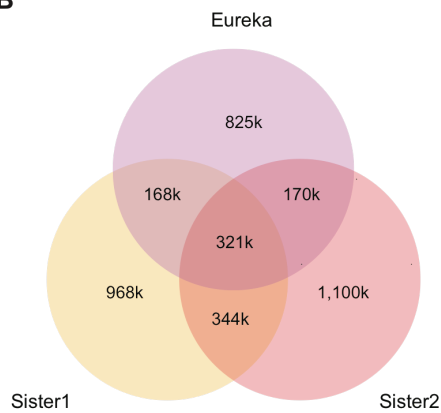
Sister 2	<i>L. pictus</i>	3903	3845	58	107	94	4104
Eureka	<i>L. pictus</i>	3829	3789	40	169	106	4104
canFam3.1	<i>C. familiaris</i>	3910	3857	53	98	96	4104
Kenya	<i>L. pictus</i>	3849	3823	26	136	119	4104
South Africa	<i>L. pictus</i>	3892	3867	25	104	108	4104
Hawaiian monk seal	<i>Neomonachus schauinslandi</i>	3881	3833	48	118	105	4104

725

A



B



C

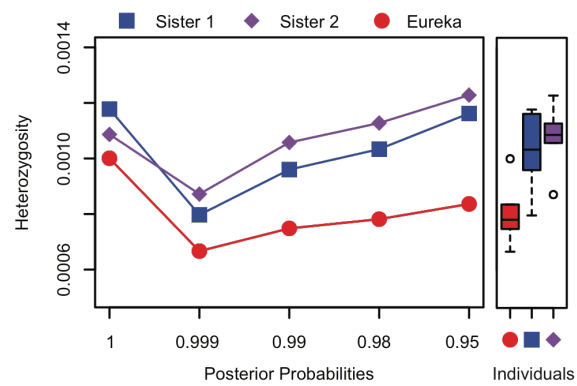
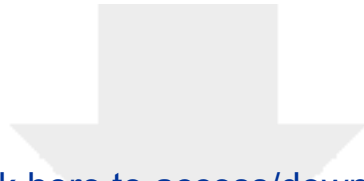


Figure 1. Shared heterozygous sites between the different African wild dog individuals. A) Pack of African wild dogs. B) Shared heterozygous sites between the three *de novo* assemblies (calculated using a posterior cutoff of 0.99). Many of the heterozygous sites are shared between all individuals and more heterozygous sites are shared between the two sisters than between each sister and Eureka. C) Comparison of heterozygosity estimates using different posterior probability cutoffs for all three assemblies. Boxplot of heterozygosity values (y-axis) calculated for different posterior probability cutoffs.

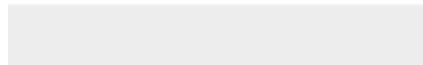
Supporting information for methods, gene and repeat annotations,
and heterozygosity calculations.



[Click here to access/download](#)

Supplementary Material

[Supporting_information_AWD_Gigascience.docx](#)



Stanford, 28th of November 2017

Dear Editor,

We would like to submit our manuscript titled “Entering the era of conservation genomics: Cost-effective assembly of the African wild dog genome using linked reads.” for consideration for publication in GigaScience.

The Anthropocene is currently impacting species across the globe at an unprecedented rate. Managers and scientists are searching for cost-effective methods to monitor population level changes and understand past distributions and adaptations to best plan for the future. Genomic tools can inform many of these questions and assist in developing SNP assays for monitoring, but historically genome sequencing has been problematic both because of its high cost and sample submission requirements. Thus, the demonstration of a system that can circumvent these issues is critical.

Here we present the results of the 10x Genomics Chromium platform for genome assembly of an endangered canid, the African wild dog (*Lycaon pictus*). Our work is the first *de novo* assembly for this species and demonstrates that this method of genome sequencing generates comparable or better assemblies than traditional Illumina based paired-end and mate-paired sequencing for a small fraction of the cost. We find that the results of the 10x Chromium technology are consistent and reproducible across the three individuals we sequenced. Additionally, we show that the wild dog has a higher heterozygosity than expected given its endangered status, which may be an outcome of its social biology.

The African wild dog, like many other endangered species, is being heavily targeted for recovery. Important measures to consider for such recovery programs will be statistics such as inbreeding and heterozygosity, and also a better understanding of the demography of the wild dogs. A well-assembled reference genome provides the basis for which these measures can be estimated using low-coverage population sequencing. As endangered species often have less genomic resources than model organisms, 10x Genomics Chromium assemblies will open up a new avenue of study without impacting conservation dollars that must be put towards on the ground monitoring.

In conclusion, we show that the 10x Genomics Chromium system produces a highly continuous and quality genome assembly for comparably less coverage and cost than other technologies. This work is particularly relevant for species of conservation concern that could benefit from genetic monitoring and in-depth studies of adaptive capabilities to mounting pressures such as climate change and shrinking population sizes.

We recently uploaded a pre-print to BioRxiv and the response has been overwhelmingly positive. We feel strongly that the manuscript is a good fit to the scope and mission of the GigaScience journal. We hope that it will reach a broad audience and facilitate the expansion of high-quality genomic datasets from many organisms.

We hope you find our manuscript of interest and look forward to hearing about your decision.

Sincerely,
Ellie Armstrong