

<b>Manuscript Number:</b>	GIGA-D-17-00324R1	
<b>Full Title:</b>	Cost-effective assembly of the African wild dog genome using linked reads.	
<b>Article Type:</b>	Research	
<b>Funding Information:</b>	John Stuelpnagel	Not applicable
<b>Abstract:</b>	<p>A high-quality reference genome assembly is a valuable tool for the study of non-model organisms. Genomic techniques can provide important insights about past population sizes, local adaptation, and aid in the development of breeding management plans. This information is 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 (&lt; \$3000), thereby facilitating future studies of this endangered canid. We generated assemblies for three individuals using the linked-read 10x Genomics Chromium system. The most continuous assembly had a scaffold and contig N50 of 21 Mb and 83 Kb, respectively, and completely reconstructed 95% of a set of conserved mammalian genes. Additionally, we estimate the heterozygosity and demographic history of African wild dogs, revealing that although they have historically low effective population sizes, heterozygosity remains high. We show that 10x Genomics Chromium data can be used to effectively generate high-quality genomes from Illumina short-read data of intermediate coverage (~25-50x). Interestingly, the wild dog shows higher heterozygosity than other species of conservation concern, possibly due to 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 conservationists to better understand the ecology and adaptability of those species in a changing environment.</p>	
<b>Corresponding Author:</b>	Ellie Armstrong Stanford University UNITED STATES	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Stanford University	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	Ellie Armstrong	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	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
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Response to Editor Comments</p> <p>We have included additional commentary regarding the sample preparation and processing. We have discussed further reasons for possible differences between the assemblies, as well as noted which parameters we are unable to investigate as a result of this study (e.g. the relationship between estimated molecule input length and percent genome phased). We have also changed the title, as requested.</p> <p>Reviewer 1</p> <p>Discretionary Revision: Perhaps it would be useful to run a PSMC-type analysis using multiple wild dog genomes to assess trends in historical population sizes in recent times for African wild dogs. This might produce useful results with conservation applications. There are several methods that have come out recently that can do a decent job with estimating population size in recent times.</p> <p>We have added a PSMC analyses of our three genomes. The results show comparative historical population sizes to those estimated in Campana et al. (2016) (Figure 1). The most notable differences are in the recent population size estimates and the timing of the beginning of the population decline, but are overall consistent.</p> <p>Edit: Line 444. The word "Heterozygosity" at the end of the paragraph seems out of place.</p> <p>This sentence has been revised.</p> <p>Reviewer 2</p> <p>Line 84 - 'The lineage is the only surviving member of a lineage of wolf-like canids' is I guess true to some degree, but that could be said of other wolf-like canids like the dhole, Ethiopian wolf, African Golden Wolf etc. Perhaps consider rewriting.</p> <p>This sentence and others have been revised as suggested from this comment, as well as comments from Reviewer 3 to reflect more accurate predictions of the divergence of the African wild dog lineage from other canids. We have included more up to date estimates for this timing.</p> <p>Line 171 and elsewhere, term 'high quality' is used. I agree that the scaffold size is excellent, but high quality also can refer to long contig sizes (in particular if one wants to study repeats, duplication etc). It would be useful if the authors could undertake a comparison of the contig sizes recovered here to those other genomes of similar SCAFFOLD quality (in particular genomes generated with different methods) so that readers can get a feel for how the contig size varies when using this approach as opposed to much more expensive methods (e.g. deep PacBio sequencing, or mate pair Illumina). Of the top of my head, one comparison in this regard could be to look at the recently published purely Illumina (mate pair) based wolf de novo genome (Gopalakrishnan et al. 2017 BMC Genomics). Unfortunately that genome is not annotated so other comparisons cannot be made (e.g. gene completeness) but simply what I suggest would be interesting.</p> <p>We have added an analyses comparing contig and scaffold sizes of our genomes with the wolf genome. We ran analyses on all genomes using the Assemblathon scripts (Table S2) and BUSCO v2 (Table 2). We also annotated the wolf genome for comparison of gene completeness with the same methods as we annotated the African wild dog genomes.</p> <p>Line 360-361 - perhaps give sequencing price per GB or per 100GB instead of per lane? As many readers may not know the lane output.</p> <p>We have noted the output of the sequencer and hope this provides a reference to the</p>

reader.

Reviewer 3

We especially thank Reviewer 3 for their extensive time and comments to our manuscript. Below we have outlined responses to these comments, as well as clarification on certain aspects of the manuscript.

1. Lines 1-2: The title should be revised because we've already been in the 'era of conservation genomics' for several years now, so this idea is out of date. How about just shortening the title to: "Cost-effective assembly of the African wild dog (*Lycaon pictus*) genome using linked reads"

Revised as suggested.

2. Line 80: Add a comma after "Taken together" so that the sentence reads: "Taken together, genomic tools are poised..."

Revised as suggested.

3. Line 82: "The African wild dog..." The species is also known by two other common names that are commonly applied to *Lycaon pictus* - African painted dog and Cape hunting dog. The former is especially used by many researchers and canid conservationists. Therefore, the authors should include these alternative names: "The African wild dog, also known as the African painted dog or Cape hunting dog (*Lycaon pictus*) is a medium-sized (18-34kg)..."

Revised as suggested.

4. Line 83: "sub Saharan should be hyphenated.

Revised as suggested.

5. Lines 123-125: "The groups containing the African wild dog and the domestic dog..." The authors cite the Nyakatura and Bininda-Emonds (2012) paper on the updated supertree analyses of the Carnivora to support the phylogenetic grouping and divergence time of the African wild dog in relation to the domestic dog. However, the supertree results are inconsistent with more direct assessments of canid phylogenetic history based on analyses of DNA sequences from multiple nuclear and mitochondrial loci. Supertree analyses have been empirically shown to produce inaccurate results regarding relationships. Direct assessment of DNA sequences indicate that the African wild dog and domestic dog, its wild counterpart, the gray wolf, and other wolf-like canids, are grouped together in the same clade (Tribe Canini, the wolf-like-canids). Furthermore, recent estimates of divergence times suggest that the African wild dog lineage and domestic dog lineage split only about 2.5 - 4 Mya (less than have the age suggested by Nyakatura and Bininda-Emonds, 2012). The authors should instead cite the following references: Lindblad-Toh et al. 2005 Nature 438: 803; Perini et al. 2010 Journal of Evolutionary Biology 23: 311; . The authors should then revise this sentence accordingly.

Associated sentences revised and inferences revised accordingly.

6. Lines 138-139: "...it has been impossible to assemble highly-contiguous genomes from only these short sequences." This statement is incorrect, in particular, the use of the word "impossible." Many mammalian genome assemblies with high continuity (e.g., human, dog, cow, Tasmanian devil, cheetah) have been generated using Illumina short read data. Short read data per se is not the problem. Given that enough paired-end shotgun and mate pair libraries are constructed and sequenced, the resulting short read data can be assembled to produce draft assemblies with high continuity despite the high content of repetitive sequences (comparable to or greater than those generated by the 10X Genomics Chromium System). Therefore, the comparison is a relative one and mostly depends on input. I suggest the authors revise the sentence as follows: "Because large proportions of typical mammal genomes consist of repetitive sequences, it has been challenging to obtain complete or highly

continuous genome assemblies using only these short sequences."

Revised as suggested.

7. Lines 173-175: "Thus, in order for it to be useful for conservation purposes the technology needs to be (a) cost-effective and (b) user-friendly." This sentence doesn't make sense and doesn't accord with the facts. Genomes of multiple endangered species (e.g., tiger - Cho et al. 2013 Nat Comm; crested ibis - Li et al 2014 Genome Biol.; cheetah - Dobrynin et al. 2015 Genome Biol.' Iberian lynx - Abascal et al. 2016 Genome Biol.) have been generated and directly useful for conservation purposes regardless of their cost-effectiveness or user-friendliness. The authors' statement precludes other potential sequencing technologies that may not be as cost-effective (e.g. PacBio long reads) but yet still may be used to obtain high quality genome assemblies for conservation genomic applications. And most surprisingly, why should user-friendliness with regards to analysis of next generation sequencing data (i.e., bioinformatics) ever be a criterion on whether it is useful or not for conservation? Please delete this sentence.

We have revised this sentence with an emphasis on the practicality of using genomics as a wide-spread tool in the conservation world. We would defend that it still remains elusive or out of reach for many conservation biologists to assemble a genome de novo, despite desiring to use what a reference assembly provides downstream for everyday conservation practice. We direct the reviewers to a recent study (Taylor et al. (2017) Bridging the conservation genetics gap by identifying barriers to implementation for conservation practitioners. Global Ecology and Conservation), which describes a common disconnect between managers desiring to use genetic and genomic resources, but lacking the funds and expertise to use such technologies.

8. Line 184: "and are presumed to be sisters..." The authors should indicate that the details behind this presumption are included in the supporting information and cite Appendix S1.

Revised accordingly.

9. Lines 202 - 204: The authors need to cite Hoepfner et al. 2014 here; e.g., "...from the most recent dog genome (267kb and 45.9Mb, respectively [48]),"

Revised accordingly.

10. Line 216: Same comment as point 9; need to cite the Hoepfner et al. 2014 paper.

Revised accordingly.

11. Lines 240-241: "Furthermore, repeat content of all wild dog assemblies was qualitatively similar to canFam3.1." Given that African wild dog and domestic dog share a relatively close evolutionary ancestry (see point #5 above), it's not surprising that their repeat contents would be similar. The authors should qualify their findings in these terms.

Revised accordingly.

12. Lines 242-245: "...the similarity in repeat content between the African wild dog compared to that of the domestic dog, highlights the value of using 10x Genomics Chromium technology to produce accurate and continuous assemblies." This seems like a specious conclusion. The canFam3.1 assembly was not generated using 10x Genomics data, yet it has a repeat content similar to the African wild dogs. This is likely due to the recent common ancestry (point #11) and not because of the technology used to sequence/assemble the genome. The repeat content of the two species would be similar regardless of the continuity of the assembly or how that was achieved. I recommend the authors delete the last sentence in this paragraph.

Revised as suggested.

13. Line 254: "...multi copy..." should be hyphenated (multi-copy).

Revised as suggested.

14. Line 255: "...and 37 not present in one individual." Specify which individual was missing these multi-copy genes (paralogues). Any reason why these 37 multi-copy genes were missing? Lower coverage? Assembly problem?

We re-phrased this sentence to more accurately reflect the results. Thirty-seven total singletons were missing across the three individuals, with the lowest coverage genome missing the most and the highest coverage genome missing the least.

15. Lines 270-272: "As expected, we see a higher number of singletons in these two individuals..." Here the authors should be more explicit about the discrepancy in the number of singleton SNPs in the two African wild dogs sequenced by Campana et al. 2016 and the three individuals sequenced by the authors. Please provide numbers or percentages about the differences and then cite the Appendix S1 for the detailed methods used for variant calling. Coverage in and of itself may not be the sole reason for the higher number of singletons in the two African wild dogs sequenced by Campana et al. More stringent filtering methods applied to these two individuals would likely have resulted in a comparable number of SNPs to the three individuals sequenced by the authors. The authors should discuss these alternatives. Also, the Nielsen et al. 2011 and 2012 references are not included in the references (main text or Appendix S1). Also, the authors should consider the following papers: Bryc et al. 2013 *Genetics* 195: 553 and Kousathanas et al. 2017 *Genetics* 205: 317.

We agree with the reviewer that there is much to be said for the different ways to estimate heterozygosity, but would add that this is difficult to do without introducing additional biases. Indeed, data-preprocessing, the choice of a reference genome (this particular issue is documented in Gopalakrishnan et al. 2017 using the wolf data), mapping tools, and filtering, may all introduce unknown biases in heterozygosity estimates. Our intention in this paper was not to estimate heterozygosity using multiple different methods, but rather use a single method and estimate differences. However, this would be a pertinent follow-up study in the future using a more controlled data set and we will certainly consider this. We have adjusted the language here to acknowledge the limitations of our analyses.

16. Lines 280-281: "Our estimates show that, while being heavily threatened, African Wild dogs seem to still retain a relatively high within individual heterozygosity." First "Wild" in this sentence should be revised as "wild." Second, the conclusion of "relatively high within individual heterozygosity" is impossible to judge without context to some reference/metric or other species. Relative to what exactly? The per site heterozygosities measured by the authors should be compared to those obtained from other species listed as endangered or critically endangered on the IUCN Red List. The paper by Robinson et al. 2016 *Current Biol.* 26: 1183 would be of use for this. Furthermore, it would be useful to compare the per site heterozygosities obtained for the three African wild dogs with those of gray wolves reported by Gopalakrishnan et al. 2017 *BMC Genomics* 18: 495 (see their Table S1).

We have included comparisons to those reported for several endangered species in Dobrynin et al. 2016, Gopalakrishnan et al. 2017, and Robinson et al. 2016.

17. Lines 299-301: "This may indicate that input molecule length is a key factor for scaffolding, while coverage is a key factor for contig assembly." Input molecule length is indeed likely to have a strong effect on assembly quality for the 10X Genomics platform. In fact, this is directly stated by 10X Genomics: "DNA quality. By far the most common cause of subpar assembly results is poor input DNA quality" (<https://support.10xgenomics.com/de-novo-assembly/software/pipelines/latest/troubleshooting>). In fact, the Chromium library preparation process may nick the DNA and thus cause fragmentation (smaller molecule lengths). The authors should include and cite the weblink above. It is somewhat surprising that the assemblies of the three African wild dogs were so different in terms of their assembly metrics (e.g., contig and scaffold N50s). Given that the 10X Genomics linked-read technology is still relatively new, it's difficult to judge

whether these results are common or not. The authors' findings do not accord with my own experience using 10X, where assembly metrics from multiple individuals of the same species were more consistent (mostly identical). The authors should discuss in one or two additional sentences other factors that may have influenced their results: 1) sample handling, storage, and/or preparation; 2) library preparation - were the three libraries prepared by the same lab or technician? The authors state in Appendix S1 that the three individuals were sequenced at two different sequencing facilities/vendors; 3) sequencing platforms, chemistries used (HiSeq X for two individuals vs. HiSeq4000 for the third).

We have included the link as part of our revisions and added this as a commentary. We do emphasize that the three assemblies were sequenced at different depths, which may also result in some of the stochasticity among our assemblies. We hope that what comes across is not that the assemblies are wildly different, but rather that as an assembly service which is cost-effective, that the results across individuals are more or less consistent.

18. Lines 357-375: Cost effectiveness: The authors should list the US sequencing facilities examined and their corresponding prices for Chromium library preparation and sequencing in the Supporting Information- Appendix 1 in a table. This will provide readers with the explicit information to gauge different costs associated with these services. This information is also usually provided on the websites of sequencing facilities and vendors. Also, the authors should indicate the pricings for the library preparation and sequencing at the two sequencing facilities they used to generate the data of the three African wild dogs. Also, how much would the cost be for if the authors had used generated and sequenced Illumina shotgun and mate pair libraries to obtain genome assemblies comparable in quality to those generated using the 10X Chromium platform?

We have included details on the prices we paid for each assembly. We are reluctant to include a survey of current costs because the cost for sequence services changes rapidly, and the prices posted on websites are not always representative of negotiated prices. We believe the prices we paid are within 15% of prices currently offered by most sequencing providers.

We have more explicitly listed the cost of each of our genomes by their components (the price of a lane and the price of the library prep) in comparison with the approximate cost to prepare the libraries and sequencing of the wolf genome, a comparable Illumina library based genome.

19. Lines 408-411: See my previous comments with respect to this issue in point # 15 above. It would be useful to cite Nielsen et al. 2011 and 2012 here.

We have incorporated the Nielsen et al. 2011 & 2012 citations where appropriate. We thank the reviewer for bringing this oversight to our attention.

20. Line 414: "other threatened large bodied carnivores..." - Neither the Iberian lynx nor dwarf Channel island fox would be considered large-bodied. I suggest the authors revise this just as: "other threatened carnivores..."

Revised as suggested.

21. Line 421: a comma should be added after "dogs" in this sentence.

Revised as suggested.

22. Line 433: "...as part of the assembly process, however, when the fasta consensus sequence..." This is a run-on sentence and should be broken into two sentences: "...as part of the assembly process. However, when the fasta consensus sequence..."

Revised as suggested.

23. Line 473: "DNA was extracted 9 days after the sample was taken." The authors should provide details about how this sample was stored prior to DNA extraction. Also,

what type of blood tubes (e.g., Vacutainer) were the samples collected into? These details are important to document given the importance of the HMW input DNA to the success of the 10X Genomics Chromium technology (and in the interests of reproducibility).

We had described the storage and processing of the samples in detail, but failed to reference appendix S1. We have corrected this error.

24. Line 486 (and in Appendix S1): In the interests of reproducibility, the default assembly parameters should be listed or described.

There are no assembly parameters for Supernova and it is simply 'supernova run' in the same directory as the fastq files.

25. Line 492: "lineage specific" should include a hyphen.

Revised as suggested.

26. Line 496: "BAC end" should include a hyphen.

Revised as suggested.

27. Lines 524-527: The 10X Genomics Supernova assembler outputs four FASTA data files (raw, megabubbles, pseudohap and pseudohap2); see: <https://support.10xgenomics.com/de-novo-assembly/software/pipelines/latest/output/generating>. Given that there are only two outputs that provide the phased information (pseudohap and pseudohap2), how could this choice for estimating heterozygosity possibly be described a random? In the interests of reproducibility, the authors should indicate which pseudo-haplotype file was used for which individual African wild dog. Also, the authors should at least take one individual (Sister 2, the one with the most continuous assembly) and estimate the heterozygosity from the other pseudo-haplotype file to check that there is no difference in the inferred number of heterozygous sites (this acts as a control).

We have included an analysis of the two distinct pseudohaplotypes from the --style=pseudohap2 output for Sister 2 and have included a more thorough description of which files were used for each. We do note, however, that the software does a randomized pseudohaplotype when the option --style=pseudohap is chosen and is noted here in the Supernova manual: "For pseudohap...Megabubble arms are chosen arbitrarily so many records will mix maternal and paternal alleles." However, for --style=pseudohap2, the maternal and paternal arms are separated. We have made efforts to make this more clear in the text.

28. Line 529: The Samtools and Picard programs should be capitalized.

Revised as suggested.

29. Literature cited: The authors should carefully check the formatting of their references so that they consistently conform to the journal standards (e.g., journal titles are often not properly capitalized).

Revised as suggested.

30. Methods (main text and Appendix S1): Samples. Given the requirement of input DNA with long molecule lengths and its importance to the 10X Genomics technology, no details or information is provided on how the HMW genomic DNA was assayed following extraction. This is absolutely crucial and related to the issue of experimental reproducibility. Such HMW DNA is usually assessed using pulse-field electrophoresis techniques or variations thereof. Since the authors used two different sequencing facilities to generate the libraries and sequencing data, different methods may have been used for the assays. In any case, the authors should provide the details about how the HMW DNA was assessed and evaluated prior to Chromium library preparation.

	<p>We have added additional information on the assays performed following extraction in the supplement.</p> <p>31. Phased assemblies: Even though the percentage of the assemblies that were phased is presented in Table S1, this feature is never discussed in detail in the main text. However, this is one of the most noteworthy (and marketed) features of the 10X Genomics platform. Phased assemblies also have a dramatic impact on the downstream population genetic analyses and provide additional information for these analyses compared to technologies that do not yield phased assemblies. The authors should include a description of the phasing results of the three African wild dog assemblies in the Data Description &amp; Analyses section as well as discuss this important feature of the 10X Genomics platform.</p> <p>We considered this point extensively during analyses, but unfortunately are not able to address this point with the data in hand. Although we can produce phased vcf files, the genomes produced from the Sister 1 and Sister 2 individuals by independent Supernova runs are still too fragmented for us to consider the phasing of any certain haplotype or position, nor to investigate whether the sisters share the expected amount of variation. We are continuing this project with population-level sequencing of individuals from Zimbabwe and hope to address this point further when we have additional information on the expected allele frequencies.</p>
<b>Additional Information:</b>	
<b>Question</b>	<b>Response</b>
Are you submitting this manuscript to a special series or article collection?	No
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p><b>Resources</b></p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p>	Yes



<p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>

[Click here to view linked References](#)

1 **Cost-effective assembly of the African wild dog (*Lycaon pictus*) genome using**  
2 **linked reads.**

3  
4  
5  
6 4 Ellie E. Armstrong<sup>1\*</sup>, Ryan W. Taylor<sup>1\*</sup>, Stefan Probst<sup>1,2</sup>, Peter Blinston<sup>3</sup>, Esther van der  
7  
8 5 Meer<sup>3</sup>, Hillary Madzikanda<sup>3</sup>, Olivia Mufute<sup>4</sup>, Roseline Mandisodza-Chikerema<sup>4</sup>, John  
9 6 Stuelpnagel<sup>5</sup>, Claudio Sillero-Zubiri<sup>6</sup>, Dmitri Petrov<sup>1</sup>

10  
11  
12  
13  
14  
15 8 <sup>1</sup>Program for Conservation Genomics, Department of Biology, Stanford University,  
16 9 Stanford, CA, USA

20 10 <sup>2</sup>Department of Integrative Biology, University of California, Berkeley, CA, USA

22 11 <sup>3</sup>Painted Dog Conservation, Dete, Zimbabwe

24 12 <sup>4</sup>The Zimbabwe Parks & Wildlife Management Authority, Zimbabwe

26 13 <sup>5</sup>10x Genomics, Inc., Pleasanton, CA

28 14 <sup>6</sup>Wildlife Conservation Research Unit, Zoology, University of Oxford, The Recanati-  
30 15 Kaplan Centre, Tubney, UK014

33 16

36 17 \* These authors contributed equally to this work.

38 18 Corresponding Author: Ellie E. Armstrong ([elliea@stanford.edu](mailto:elliea@stanford.edu))

40 19

42 20

44 21 **Abstract**

46 22

48 23 **Background**

50  
51 24 A high-quality reference genome assembly is a valuable tool for the study of non-  
52  
53 25 model organisms. Genomic techniques can provide important insights about past  
54  
55 26 population sizes, local adaptation, and aid in the development of breeding  
56  
57 27 management plans. This information is important for fields like conservation genetics,  
58  
59 28 where endangered species require critical and immediate attention. However,  
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 funding for genomic-based methods can be sparse for conservation projects, as  
30 costs for general species management can consume budgets.

31

## 32 **Findings**

33 Here we report the generation of high-quality reference genomes for the African wild  
34 dog (*Lycaon pictus*) at a low cost (< \$3000), thereby facilitating future studies of this  
35 endangered canid. We generated assemblies for three individuals using the linked-  
36 read 10x Genomics Chromium system. The most continuous assembly had a  
37 scaffold and contig N50 of 21 Mb and 83 Kb, respectively, and completely  
38 reconstructed 95% of a set of conserved mammalian genes. Additionally, we  
39 estimate the heterozygosity and demographic history of African wild dogs, revealing  
40 that although they have historically low effective population sizes, heterozygosity  
41 remains high.

42

## 43 **Conclusions**

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

52

## 53 **Keywords**

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

56  
61  
62  
63  
64  
65

## 57 **Background**

58 Major population declines have been observed in vertebrate groups over the  
59 past several hundred years, primarily due to anthropogenic change [1]. This decline  
60 has resulted in extinction rates unprecedented in recent history [1, 2]. The  
61 conservation of extant species will require major efforts in restoring and preserving  
62 habitat, along with protection, management, and investment by local stakeholders.  
63 While, by definition, all species of conservation concern exist as small populations,  
64 populations generally still retain genetic variation that was generated and maintained  
65 when population sizes were much larger.

66 The historic genetic variation contains signals of demographic history, gene  
67 flow, and natural selection which can inform efforts towards the long-term survival of  
68 species. In addition to signals of a species history, genetic information can be used  
69 to uncover important contemporary or very recent events and processes. Genetic  
70 markers can be used to track individual movement across landscapes either  
71 indirectly by measuring relatedness, or directly by genotyping scat or hair left by an  
72 individual as it moves. Additionally, the identification and assignment of individuals  
73 through genotyping can be an important tool for law enforcement to assign  
74 contraband and confiscated materials to their geographic origin [4]. Conservationists  
75 can also use fine grained measurements of reproductive success along with  
76 genotypes and environmental variables to gather a detailed understanding of the  
77 factors contributing to or limiting population growth, such as inbreeding depression.  
78 Taken together, genomic tools are poised to have a major contribution to  
79 conservation [5, 6].

80 The African wild dog, also known as the African painted dog or Cape hunting  
81 dog (*Lycaon pictus*), is a medium-sized (18-34kg), endangered carnivore that lives in  
82 scattered populations in sub-Saharan Africa (Fig. 1A). The species is a surviving  
83 member of a lineage of wolf-like canids, including other species such as the  
84 Ethiopian wolf and the dhole [7]. Wild dogs have been subject to intense recovery

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

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  
103 Prior genetic studies on wild dogs using a combination of mitochondrial,  
104 microsatellite, and MHC markers have resulted in varying estimates of the start of the  
105 species decline on the African continent [19, 20]. Consistent with expectation, the  
106 data shows strong structuring among populations due to habitat fragmentation and  
107 isolation, as well as low genetic diversity within populations [20, 21]. For species that  
108 are experiencing such rapid and alarming declines, estimates that are particularly  
109 important for management decisions, such as effective population size, inbreeding  
110 and local adaptation, are greatly improved by the use of whole-genome methods.  
111 Recently, Campana and colleagues [22] sequenced low-coverage genomes of two  
112 African wild dog individuals from Kenya and South Africa, respectively, to investigate

113 demographic history and signatures of selection of these two separate populations.  
114 By mapping these data to the domestic dog genome, they discovered approximately  
115 780,000 single nucleotide polymorphisms (SNPs) between their two individuals  
116 which could be used to develop SNP typing for the two populations. However, given  
117 the low coverage of their genomes (5.7-5.8x average coverage) and the small  
118 number of individuals sequenced, additional sequencing will be needed to verify the  
119 authenticity of those SNPs. Further, important structural variation can be overlooked  
120 when mapping against a reference genome from a different genus, and mapping can  
121 be hindered if the divergence is high between the sample and the reference (see e.g.  
122 [23]). The groups containing the African wild dog and the domestic dog are estimated  
123 to have split approximately 2.5-4 Mya and furthermore, the domestic dog has  
124 undergone significant genomic selection in recent time [24, 25,26].

125         Despite the ever-declining cost to sequence DNA, the routine use of genomic  
126 approaches in conservation is still far from a reality. One of the major remaining  
127 barriers is the lack of reference genomes for species of conservation concern.  
128 Generating a *de novo* reference genome generally requires the sequencing and  
129 assembly of billions of base pairs that make up a genome. The first mammalian  
130 genome (human) required a massive collaboration among hundreds of scientists and  
131 nearly \$3 billion US dollars (1990-2001; [27, 28]). Fortunately, the cost to sequence  
132 DNA is now low enough that every base-pair in a typical mammalian genome can be  
133 sequenced to high-coverage for a few thousand US dollars. However, these low-cost  
134 sequencing methods produce very short sequences of 150-300 base-pairs in length  
135 (for a review on sequencing methods see [29]). Because large proportions of typical  
136 mammal genomes consist of repetitive sequences, it has been challenging to obtain  
137 complete or highly-contiguous genomes using only these short sequences. In order  
138 to achieve higher continuity, more elaborate and expensive library preparation or  
139 alternative sequencing technologies have to be used [29, 30]. Among others, these  
140 include mate-pair libraries, chromatin folding based libraries, such as cHiCago [31] or

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 HiC [32], and long-read sequencing technologies, such as Pacific Biosciences and  
142 Oxford Nanopore Technology. While the resulting genomes can show high  
143 continuity, those methods substantially increase the costs of sequencing projects and  
144 thus can hinder the generation of genomes for conservation biology purposes.

145 Here we report the use of the Chromium system developed by 10x Genomics  
146 [33], a genomic library preparation technique that facilitates cost-effective  
147 assemblies using short sequencing reads, to assemble three African wild dog  
148 genomes. In brief, the 10x Genomics Chromium system is based on dilution of high  
149 molecular weight (HMW) DNA. It uses as little as 1ng of input DNA, which is well-  
150 suited for a variety of applications. During library preparation, gel beads, so-called  
151 GEMs, are mixed with DNA and polymerase for whole-genome amplification. Each  
152 gel bead has primer oligos (44nt long) attached to its surface. These contain a  
153 priming site (22nt partial R1), a 16nt barcode region, and a 6nt N-mer region that  
154 binds to different places on the original DNA fragment. The low amount of input DNA  
155 ensures that each gel bead only binds a single (up to ~100kb) DNA fragment. In the  
156 next step, amplification of short reads along the original DNA fragment is performed  
157 within each gel bead. In most cases, this amplification results in spotted read  
158 coverage along the fragment. However, all reads from a respective GEM contain  
159 identical barcodes and can later be assigned to groups originating from the same  
160 DNA molecule. The information about which molecule of DNA the sequence  
161 originated from greatly increases the ability to identify the location of repetitive  
162 sequences. The library is then sequenced on an Illumina platform and the raw read  
163 data is assembled by the 10x Genomics Supernova assembler. The data produced  
164 also can be phased, presenting another potentially useful addition to genome  
165 assemblies.

166 We *de novo* assembled three African wild dog genomes using the 10x  
167 Genomics Chromium platform to investigate whether this technology is suitable for  
168 conservation genomic purposes. For any endangered species, a genome can enable

169 studies with the potential for large conservation impacts, but high-quality genomes  
170 have historically been costly or impossible due to the sampling requirements and  
171 analysis. Thus, for an assembly to be a practical component of many conservation  
172 projects, the technology needs to be (a) cost-effective and (b) user-friendly. We test  
173 the 10x Genomics Chromium based assemblies for reproducibility, continuity,  
174 conserved gene completeness, and repetitive content, as compared to the previously  
175 published domestic dog genome [34] and several other genomes built with various  
176 technologies. We further estimate heterozygosity of the individuals and within the  
177 phased data from the 10x technology and estimate historical effective population size  
178 from each genome.

179

## 180 **Data Description & Analyses**

181

### 182 *Assembly of the African wild dog genome*

183       Using 10x Genomics Chromium technology, we generated DNA libraries for  
184 three African wild dog individuals, two of which were collected from a wild pack in  
185 Hwange National Park, Zimbabwe and are sisters from the same litter born in June of  
186 2013 (identified as Sister 1 and Sister 2, additional information can be found in  
187 Appendix S1), and a third unrelated individual from the Endangered Wolf Center,  
188 Eureka, Missouri (identified as Eureka). A summary of the assembly statistics output  
189 by the Supernova assembler can be found in Table 1 (detailed statistics for each  
190 genome assembly can be found in Table S1). We generated ~1.2 billion paired-end  
191 reads for Sister 1, ~0.8 billion reads for Sister 2, and ~0.4 billion reads for Eureka.  
192 We then used the reads to assemble each genome using the 10x Genomics  
193 Supernova assembler (as explained in <https://support.10xgenomics.com/de-novo-assembly/software/overview/welcome>). The mean input DNA molecule length  
194 reported by the Supernova assembler was 19.91kb for Sister 1, 196 77.03kb for  
195 Sister 2, and 52.00kb for Eureka. All three assemblies corroborate a genome size of  
196



197 approximately 2.3Gb, which is similar to that of the domestic dog (2.4Gb; [34]).  
198 These three assemblies together constitute the first reported *de novo* assemblies for  
199 the African wild dog species.

200 The Sister 1 assembly resulted in a contig and scaffold N50 of 61.34 kb and  
201 7.91 Mb, respectively, the Sister 2 assembly achieved 83.47 kb contig and 21.34 Mb  
202 scaffold N50s, and the Eureka assembly had 50.15 kb contig and 15.31 Mb scaffold  
203 N50s (Table 1). While the scaffold N50's of these three 10x genomes are smaller  
204 than the ones from the most recent dog genome (267kb and 45.9Mb, respectively),  
205 they are still larger than most mammalian genomes assembled that used only short  
206 read data (see e.g. [36]). A recent *de novo* assembly of a wild wolf using Illumina  
207 mate-pair libraries of varying insert size resulted in a similar contig N50, but much  
208 lower scaffold N50 measurements than our results (Supporting Information Table S2;  
209 [35]). Interestingly, despite the molecule size being the highest for Sister 2, the  
210 highest percent phased data was obtained by Eureka (52.54% compared to 40.1%;  
211 Table S1).

212

### 213 *Conserved Genes*

214 The program BUSCO (Benchmarking Universal Single-Copy Orthologs) uses  
215 highly conserved single-copy orthologous genes from several different taxa and  
216 groups to test assemblies (both genomic and transcriptomic) for gene completeness,  
217 fragmentation, or absence as an indicator of assembly quality. Using BUSCO v2 on  
218 our assemblies, we found that the most continuous assembly, Sister 2, completely  
219 recovered 95.1% of conserved genes (Mammalia gene set; Table 2). Sister 1 and  
220 Eureka recovered 95.4% and 93.3% of complete conserved genes, respectively.  
221 Using the same analysis, we found 95.3% of complete conserved genes in the latest  
222 dog assembly (canFam3.1; [34]). This indicates that although the domestic dog  
223 assembly is more continuous overall, our assemblies recover nearly the same or  
224 even higher numbers of conserved genes. Surprisingly, Sister 1 had the fewest

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

225 missing genes out of all the assemblies assessed, despite lower continuity than  
226 Sister 2. We also ran BUSCO on the Hawaiian monk seal genome, generated  
227 through the combination of 10x Genomics Chromium and Bionano Genomics Irys  
228 data, and found it recovered 94.6% of conserved genes using BUSCO [37]. This  
229 suggests that using Bionano in addition to 10x does not greatly improve the  
230 reconstruction of the gene regions. However, the Hawaiian monk seal genome has a  
231 scaffold N50 of approximately 28Mb, so Bionano may improve the overall assembly  
232 continuity compared to 10x Genomics alone. The low-coverage genomes from  
233 Campana et al. 2016 achieved a BUSCO score of 92.8% for the individual from  
234 Kenya and 94.8% for the individual from South Africa [22]. The wolf genome also  
235 scored similarly (94.8%) [35].

236

### 237 *Repeat annotation*

238 We identified repetitive regions of the genome to discern how well these  
239 complex areas were assembled by the 10x Genomics Chromium technology. We  
240 found that for all three wild dog assemblies, total repeat content was evaluated to be  
241 within 3% of one another, which indicates consistency among assemblies from a  
242 single species (Supporting Information Table S3). No single repeat category was  
243 disproportionately affected during repeat annotation of the three genomes, which  
244 suggests that assembly quality was likely the most influential factor. Furthermore,  
245 repeat content of all wild dog assemblies was qualitatively similar to canFam3.1 [34]  
246 and the wolf genome [35], likely due to recent common ancestry between the two  
247 groups [24, 25, 26].

248

### 249 *Gene annotation*

250 Genome annotation resulted in very similar numbers of annotated genes  
251 between all three African wild dog individuals and the domestic dog [34]. Annotations  
252 ranged from 20,649 (Sister 2) to 20,946 (Sister 1) genes (Supporting Information

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

253 Table S4). Through detecting orthologous genes between individuals and paralogous  
254 genes within individuals, we found 12,617 one:one orthologs present in all three  
255 individuals and 6,462 one:one orthologs in two out of the three individuals. We found  
256 268 multi-copy genes present in all three individuals and 37 total not present in single  
257 individuals, likely due to their coverage differences (ten were missing in Sister 1,  
258 thirteen in Sister 2 and fourteen in Eureka). Overall, the number of annotated genes  
259 was comparable to those found in the domestic dog genome and the wolf genome  
260 (Supporting Information Table S4; [34,35]).

261

262 *Variant rates*

263 We found a high number of heterozygous sites to be shared between all three  
264 individuals (321k; here we report the heterozygous sites called using a posterior  
265 probability cutoff of 0.99; Supplementary Information Figure S2A). As expected,  
266 Sister 1 and Sister 2 share more heterozygous sites (344k) than either sister with  
267 Eureka (168k and 170k, for Sister 1 and Sister 2, respectively). Each individual  
268 shows a high number of singletons (heterozygous sites only found in one individual),  
269 with Sister 2 showing the highest number (1,100k), followed by Sister 1 (968k) and  
270 Eureka (825k). Even if we include the two low-coverage genomes from Campana et  
271 al. (2016) [21], we find a high number of shared heterozygous sites between all  
272 individuals (134k; Supporting Information Figure S2B). We see a higher number of  
273 singletons in these two individuals, most likely due to the lower reliability of the  
274 genotype calls caused by the low-coverage data (false positives caused by  
275 sequencing errors). We estimated a per site heterozygosity of 0.0008 to 0.0012 for  
276 Sister 1, 0.0009 to 0.0012 for Sister 2, and 0.0007 to 0.001 for Eureka using  
277 posterior cutoffs for genotype calls from 0.95 to 1 in ANGSD (Supporting Information,  
278 Fig. S1C). As can be seen in Supplementary Figure S2, except for a posterior  
279 probability cutoff of 1, where Sister 1 shows the highest heterozygosity, Sister 2  
280 always shows the highest, Sister 1 the second highest and Eureka the lowest

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

281 heterozygosity. Interestingly, Eureka shows a lower heterozygosity than the other  
282 two assemblies, even though its parents are thought to have originated from different  
283 localities (Supplement S1). With more stringent filtering, we likely could improve the  
284 heterozygosity estimates for the low-coverage individuals, but we did not investigate  
285 this further and maintained our methods across datasets for comparative purposes.

286 We did not see any major difference between heterozygosity estimates from  
287 repeat-masked and unmasked genomes [66]. The Supernova software estimated a  
288 heterozygous position every 2.6kb, 3.1kb, and 7.14kb for Sister 1, Sister2, and  
289 Eureka, respectively (Supporting Information Table S5). On the contrary, estimates  
290 based on genotype calls using ANGSD showed much more frequent heterozygous  
291 positions (850bp - 1.2kb, 814bp - 1.1kb and 999bp - 1.5kb depending on the  
292 posterior cutoff used; Supporting Information Table S5). Overall, our estimates show  
293 that, while being heavily threatened, African wild dogs seem to still retain a relatively  
294 high within-individual heterozygosity relative to other endangered species which have  
295 been estimated, such as those in the cheetah or the Amur tiger ( $> 0.0005$ ,  $0.0005$ ;  
296 [38]), or the island grey fox ( $>0.0005$ ; [39]). Additionally, the estimates here are  
297 comparable to those from several gray wolf individuals ( $0.0009$ - $0.0012$ ; [35]).

298 We also examined the phased data and its effect on heterozygosity estimates  
299 for one individual, Sister 2. We find that the estimates are relatively consistent  
300 between both the pseudohaplotypes, and the merged pseudohaplotype produced by  
301 the Supernova software (Supplementary Information Table S5) [66].

302

### 303 *Demographic history*

304

305 We estimated demographic history using the program PSMC [40]. Our results  
306 show similar demographic trends with those reported in Campana et al. (2016) [22],  
307 however, we observe declines beginning just over 1mya, as opposed to  
308 approximately 700,000 years ago (Figure 1C). From 1 million to 120,000 years ago

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

309 the population size steadily declines, resulting in a predicted  $N_e$  of approximately  
310 1,000-2,000 individuals. During the remainder of the African wild dog history, there  
311 are some small effective population size estimate fluctuations.

312 We also infer similar population histories from the genomes of the two sisters  
313 from Zimbabwe and furthermore, show very little difference between the inferred  
314 history of the third individual, Eureka (Figure 1C). This may be because the  
315 populations were formerly continuous and share their ancestral population history,  
316 but further analyses would be required to disentangle these hypotheses. We also do  
317 not detect additional large fluctuations as noted by Campana et al. (2016) [22], but  
318 more high coverage genomes from across populations would be needed to confirm  
319 that these do not exist, since our individuals are from distinct populations than those  
320 previously tested. Furthermore, population structure and short-term demographic  
321 incidents (e.g. populations bottlenecks) can affect PSMC estimations of historic  
322 population sizes [41]. In addition, the assumed mutation rate and generation times  
323 can have large effects on the resulting estimates. However, the data consistently  
324 reinforces that African wild dogs have existed at relatively low population sizes for a  
325 long time.

## 327 **Discussion**

### 329 *Assembly continuity and quality*

330 All three African wild dog assemblies produced with 10x Genomics Chromium  
331 data showed high continuity, high recovery rates of conserved genes, and expected  
332 proportions of repetitive sequence overall. The assembly for Sister 2, which has the  
333 highest mean molecule length, is also the most continuous (Contig N50: 83.47kb,  
334 Scaffold N50: 21.34Mb; Table 1). Interestingly, the Sister 1 genome has a higher  
335 contig N50 (61.34kb) than Eureka (50.15kb), but a lower scaffold N50 (7.91Mb and  
336 15.31Mb, respectively). This may indicate that input molecule length is a key factor

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

337 for scaffolding, while coverage is a key factor for contig assembly, and indeed, input  
338 DNA quality is noted as the most common cause of failed or substandard assemblies  
339 ([https://support.10xgenomics.com/de-novo-](https://support.10xgenomics.com/de-novo-assembly/software/pipelines/latest/troubleshooting)  
340 [assembly/software/pipelines/latest/troubleshooting](https://support.10xgenomics.com/de-novo-assembly/software/pipelines/latest/troubleshooting)). Furthermore, the percent of the  
341 genome able to be phased across genomes did not correspond to input molecule  
342 length (Table S1). More work will need to be done to determine the accuracy of the  
343 phased data and the wet lab methods and/or assembly parameters which influence  
344 these inferences.

345         Despite having the highest continuity of all three assemblies, Sister 2 did not  
346 show the highest BUSCO completeness scores (see Table 2), although the  
347 differences were minor (with 95.1% complete BUSCOs compared to 95.4% for Sister  
348 1). Sister 1 achieved the highest BUSCO scores, even compared to the latest  
349 domestic dog genome assembly (CanFam3.1 [34]; 95.2%), which has three times  
350 higher contig N50 and an almost six times higher scaffold N50. The high scores are  
351 remarkable for the limited number of reads used for the assemblies (as low as 25x  
352 coverage). As expected, Sister 2, which showed the highest continuity also had the  
353 highest repeat content (see Supporting Information Table S3). All three assemblies  
354 resulted in similar repeat contents in terms of repeat composition as well as overall  
355 percentage (within 3% of each other), with the most continuous assembly (Sister 2)  
356 showing the highest number of repeats. Repeat composition in the African wild dog  
357 genomes was also similar to the domestic dog and the wolf [34, 35].

358         All assemblies yielded similar amounts of genes, with Sister 1 showing the  
359 highest number (see Supporting Information Table S4), which reflect its BUSCO  
360 scores. Closer investigations of one:one and one:many orthologs further showed a  
361 very good agreement between annotations obtained from all three individuals. The  
362 numbers of annotated genes for all three African wild dogs were similar to those  
363 calculated for the latest domestic dog assembly and wolf genome assembly [34, 35].  
364

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

365 *10x Genomics Chromium system: Feasibility and caveats*

366

367 Most mammal genomes published in the last several years use a mixture of  
368 paired-end (PE) and multiple mate pair (MP) Illumina libraries (e.g. [36] and [42]).  
369 While often resulting in good continuity (e.g. [42] or [43]), using different insert  
370 libraries considerably increases the cost per genome. On the contrary, 10x  
371 Genomics Chromium allows for assembly of a comparable or even more continuous  
372 genome using only a single library for a fraction of the cost (see below). Furthermore,  
373 as we show here, this library technology generates high-quality assemblies from as  
374 low as 25x coverage (see Eureka assembly), while the recommended coverage for  
375 PE plus MP assemblies is approximately 80x-100x [44]. We do note however, that  
376 the most recent wolf genome used a variety of PE and MP libraries to produce a  
377 highly continuous assembly with approximately 30x total coverage [35]. Recently,  
378 Mohr and colleagues [37] presented a highly continuous assembly of the endangered  
379 Hawaiian monk seal (~2.4Gb total genome assembly length) using a combination of  
380 10x Genomics Chromium and Bionano Genomics optical mapping. Interestingly, their  
381 10x Genomics Chromium (sans additional Bionano) assembly showed similar N50  
382 statistics to those reported here (scaffold N50 22.23Mb), showing that 10x Genomics  
383 Chromium technology alone consistently generates highly continuous mammalian  
384 genome assemblies.

385 A limitation of 10x Genomics Chromium technology is the requirement of  
386 fresh tissue samples for the isolation of HMW DNA. This can be difficult or  
387 impossible to obtain from some endangered species. Fortunately, small amounts of  
388 mammalian blood yield sufficient amounts of HMW DNA when properly stored.  
389 Additionally, DNA extraction kits such as the Qiagen MagAttract kit can extract  
390 sufficient amounts of HMW DNA from as little as 200µl (See Supplementary  
391 Information S1 and Supplementary Information Figure S1). For museum samples, or  
392 tissues stored for extended periods of time, reference-based mapping might be the

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

393 only option to extract long-range genomic information. However, for extant  
394 endangered species, especially those with individuals in captivity, 10x Genomics  
395 Chromium offers a cost-effective approach to sequence genomes. For species with  
396 genome sizes <1Gb and between ~3Gb and 5.8Gb special data processing will need  
397 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)  
398 [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  
399 for the 10x Chromium library preparation are designed for GC contents similar to  
400 human (~41%), implying that the method might not work as well for genomes that  
401 strongly divert from this GC content (e.g. for some invertebrates).

402

#### 403 *Cost-effectiveness*

404

405 Sequencing costs are steadily dropping. At the time the sequencing for this  
406 project was carried out a lane on the Illumina HiSeqX cost (output of ~120Gb)  
407 approximately \$1,500 - \$2,000 and a 10x Genomics library prep ranged from \$450 to  
408 \$1000, thus allowing the generation of high quality *de novo* genomes for less than  
409 \$3,000 total (2016-2017). As we have shown, the 10x method only requires a single  
410 library to be sequenced to an average coverage of 25x - 75x for comparable results.  
411 Furthermore, computational resources required to assemble the genome are very  
412 low. The current version of Supernova 1.2 only requires a minimum of 16 CPU cores  
413 and 244Gb of memory (for a human genome at 56x coverage;  
414 <https://www.10xgenomics.com/>), and the assembly can be carried out in only few  
415 days (depending on the number of available CPU cores). This is about a reduction of  
416 five times the memory requirement compared to the first version of Supernova.  
417 Additionally, Supernova does not require parameter input or tuning, thus allowing  
418 even novices to easily assemble 10x Genomics Chromium based genomes.

419 For a comparable Illumina assembly, such as the one produced in  
420 Gopalakrishnan et al. (2017), the cost would include two paired-end and two mate-



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

421 pair libraries plus the sequencing costs [35]. Although paired-end libraries are  
422 relatively cheap to produce (\$120-\$180 USD), mate-pair libraries can be much more  
423 expensive depending on their input size (\$2000-\$3000 for larger insert sizes, or  
424 \$700-\$1000 if non-size selected). In addition, mate-pair libraries require a much  
425 larger quantity of starting material compared to the 10x library prep.

426

#### 427 *Applications in conservation*

428

429           Traditionally, conservation biologists have obtained a great deal of genetic  
430 information from a few microsatellite markers and/or nuclear and mitochondrial loci.  
431 The analysis of microsatellite markers can provide a snapshot into contemporary  
432 population structure, but this method risks providing incomplete information on  
433 selection and migration and can be an unreliable way to identify individuals from  
434 degraded low-quality DNA samples (such as scat) due to the stochastic behavior of  
435 marker amplification (allelic dropout; [45]; [46]; [47]). Moreover, microsatellites can  
436 be difficult to successfully design and develop, which can quickly increase costs for  
437 species that have little to no genetic information available. The ability to rapidly and  
438 cost-effectively generate full genomes will allow conservation biologists to bridge this  
439 gap and harvest crucial fine-scale population information for population parameters  
440 such as inbreeding (e.g. [48]), load of deleterious mutations (e.g. [49]), gene flow  
441 (e.g. [50]) and population structure (e.g. [51]). Once a reference genome has been  
442 assembled, optional (low-coverage) re-sequencing data from several individuals  
443 allows for the typing of genome-wide information such as single-nucleotide  
444 polymorphisms (SNPs), potentially neutral microsatellite loci, and other genomic  
445 regions of interest. These data can then be used to investigate the aforementioned  
446 population parameters, but also further yield insights into adaptive genetic variation  
447 and perhaps the adaptive potential of different populations or species.

448

449 *Heterozygosity within African wild dog individuals*

1  
2 450

3  
4 451 A high number of heterozygous sites were shared between all three  
5  
6 452 individuals in this study, with Sister 1 and Sister 2 sharing more heterozygous sites  
7  
8 453 than either shared with Eureka. Each of the individuals further showed a high number  
9  
10 454 of singletons (heterozygous sites only found in one individual). Even when compared  
11  
12 455 to the two low-coverage genomes from Campana et al. (2016) we find a high number  
13  
14 456 of shared sites [22]. As expected, we see a much higher rate of singletons in these  
15  
16 457 two individuals. Due to the low-coverage (5.7 - 5.8x average coverage) we suspect a  
17  
18 458 higher proportion of the called heterozygous sites to be false positives due to  
19  
20 459 sequencing errors, which could potentially be removed with more stringent filtering.  
21  
22 460 Heterozygosity per site estimates indicate a high within individual diversity. Estimates  
23  
24 461 ranged from 0.0007 - 0.001 for Eureka to 0.0009 - 0.0012 for Sister 2, which are  
25  
26 462 similar to those obtained for lions (0.00074 – 0.00148) and tigers (0.00087 –  
27  
28 463 0.00104) [52]. Intriguingly, other threatened carnivores, such as the Iberian lynx  
29  
30 464 (*Lynx pardinus*), the cheetah (*Acinonyx jubatus*), and the island fox (*Urocyon*  
31  
32 465 *littoralis*) show nearly 10-fold lower heterozygosity (0.0001 [51], 0.0002 [38] and  
33  
34 466 0.000014 - 0.0004 [39], respectively). The high within-individual heterozygosity could  
35  
36 467 be a result of their social structure, as only unrelated individuals come together to  
37  
38 468 form new packs through dispersal. In addition, Hwange National Park is considered  
39  
40 469 to be a part of the most continuous population of African wild dogs, which may  
41  
42 470 explain the high heterozygosity of Sister 1 and Sister 2 [20]. Further sequencing of  
43  
44 471 other populations and additional unrelated individuals will be needed to assess  
45  
46 472 whether the high within-individual heterozygosity is a range-wide phenomenon in  
47  
48 473 African wild dogs.

49  
50  
51 474 The Supernova software reports distance between heterozygous site  
52  
53 475 estimates (see Supporting Information Table S1). Interestingly, those estimates were  
54  
55 476 much lower than the ones obtained based on the genotype calls produced with  
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

477 ANGSD. While Supernova estimated this distance to be 2.6kb in Sister 1, 3.1kb in  
478 Sister 2 and 7.1kb in Eureka, the ANGSD based estimates range from 850bp - 1.2kb  
479 for Sister 1, 814bp - 1.1kb for Sister 2 and 999bp - 1.5kb for Eureka, depending on  
480 the posterior cutoff used. Supernova calculates the distance between heterozygous  
481 sites as part of the assembly process. However, when the fasta consensus sequence  
482 is called part of the variation can get flattened (see e.g. [33]). This phenomenon is  
483 typically seen in regions between megabubbles, which are nominally homozygous,  
484 but could in fact have some variation that cannot be phased by Supernova. We also  
485 note that heterozygosity values obtained using genotype calls in ANGSD could also  
486 be biased, as they are based on the nominal and not the effective coverage. The  
487 nominal coverage is the total number of reads that cover a site in the assembly,  
488 whereas for the effective coverage only reads from different barcodes are included in  
489 the estimation. If individual barcoded regions amplified with different efficiency during  
490 the library preparation step, then heterozygosity estimates could be unreliable.  
491 However, this should not strongly affect genome-wide heterozygosity estimates, as  
492 we expect this issue to be rare.

### 494 **Potential Implications**

495  
496 We find that the 10x Genomics Chromium system can be used to assemble  
497 highly continuous and accurate mammalian genome assemblies for less than \$3,000  
498 US dollars per genome (sequenced 2016 and 2017). The method can be easily  
499 applied to species of conservation concern for which genomic methods could greatly  
500 benefit their management and monitoring programs. For the African wild dog, these  
501 genomes will facilitate more reliable and cost-effective conservation efforts through  
502 the use of re-sequencing and SNP-typing methods. Compared to other species of  
503 conservation concern, the African wild dog has a relatively high heterozygosity.

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

504 Using demographic analyses, we also demonstrate that these wild dog populations  
505 appear to have been stable at lower effective population sizes for the past hundred  
506 thousand years. Additional studies should inquire whether this is consistent for  
507 populations across the African continent and evaluate current effective population  
508 sizes. More studies are also required to understand how both the social biology and  
509 recent precipitous population declines have impacted the population genomic  
510 structure of African wild dogs, and how management might use this information for  
511 the benefit and longevity of the species.

512

## 513 **Methods**

514

515 Detailed Methods can be found in Supporting Information (S1).

516

## 517 *Samples*

518 Blood samples from two individuals belonging to the same pack in Hwange  
519 National Park, Zimbabwe were provided by Painted Dog Conservation (CITES  
520 Export permit: ZW/0842/2015, ESA import permit: MA66259B-0, Research Council of  
521 Zimbabwe permit: 02553). These individuals were presumed to be sisters from direct  
522 observation of their litter at the den (here, named Sister 1 and Sister 2). DNA was  
523 extracted from samples two weeks after storage at -80°C. The third sample was  
524 provided by the Endangered Wolf Center, Eureka, Missouri from a captive born  
525 individual (here named Eureka). DNA was extracted 9 days after the sample was  
526 taken (additional information on sample storage can be found in appendix S1).  
527 Though the Chromium library preparation does not require large amounts of DNA,  
528 the DNA should have a mean molecule length > 200kb (high-molecular weight, or  
529 HMW). DNA from all individuals was extracted from blood samples using the  
530 QIAGEN MagAttract HMW DNA kit following the provided instructions.

531

532 *Genome Assembly*

1  
2 533 We constructed one sequencing library per individual using the 10x  
3  
4 534 Genomics Chromium System with 1.2ng of HMW input DNA. All libraries were then  
5  
6 535 sequenced on the Illumina HiSeqX (Sister 2, Eureka) or HiSeq 4000 (Sister 1)  
7  
8 536 platform. We subsequently assembled the three genomes using the 10x Genomics  
9  
10 537 genome assembler Supernova 1.1.1 [33]; [http://support.10xgenomics.com/de-novo-](http://support.10xgenomics.com/de-novo-assembly/software/overview/welcome)  
11  
12 538 [assembly/software/overview/welcome](http://support.10xgenomics.com/de-novo-assembly/software/overview/welcome)) using default assembly parameters.  
13  
14

15 539

17 540 *Assembly Quality Assessment*

19 541 We used the Supernova assembler as well as scripts from Assemblathon 2 to  
20  
21 542 determine continuity statistics, such as the scaffold N50 and the total number of  
22  
23 543 scaffolds [53]. We further applied the program BUSCO v2 (BUSCO,  
24  
25 544 RRID:SCR\_015008) [54] to assess the presence of nearly universal lineage-specific  
26  
27 545 single-copy orthologous genes in our assemblies using the mammalian gene set  
28  
29 546 from OrthoDB v9 (OrthoDB, RRID:SCR\_011980; 4104 genes; available at  
30  
31 547 <http://busco.ezlab.org>). We compare these results to the high-quality canFam3.1  
32  
33 548 assembly of the domestic dog ([34]; *Canis familiaris*). The canFam3.1 assembly was  
34  
35 549 built on 7x coverage of Sanger reads and BAC-end sequencing and has a scaffold  
36  
37 550 N50 of 46Mb. We also inferred the number of BUSCO's in the recently published  
38  
39 551 Hawaiian monk seal genome (which was assembled using a combination of 10x  
40  
41 552 Genomics Chromium and Bionano Genomics Irys data) and the two previously  
42  
43 553 published African wild dog genomes (sequenced with basic short read Illumina  
44  
45 554 technology at low coverage and assembled using the domestic dog for reference  
46  
47 555 mapping; [22]).  
48  
49  
50

51 556

53 557 *Repeat Identification and Masking*

55 558 We next identified repetitive regions in the genomes as another comparative  
56  
57 559 measure of assembly quality and to prepare the genome for annotation. Repeat  
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

560 annotation was carried out using both homology-based and *ab-initio* prediction  
561 approaches. We used the canid RepBase (<http://www.girinst.org/repbase/>; [56])  
562 repeat database for the homology-based annotation within RepeatMasker  
563 (RepeatMasker, RRID:SCR\_012954) [55]. We then carried out *ab-initio* repeat  
564 finding using RepeatModeler (RepeatModeler, RRID:SCR\_015027).

565

### 566 *Gene Annotation*

567 Gene annotation for the three assemblies was performed with the genome  
568 annotation pipeline Maker3 (MAKER, RRID:SCR\_005309) [57], which implements  
569 both *ab-initio* prediction and homology-based gene annotation by leveraging  
570 previously published protein sequences from dog, mouse, and human.

571 Orthologous genes between the three African wild dog assemblies, as well as  
572 paralogous genes within each individual, were inferred using Proteinortho [58].  
573 Proteinortho applies highly parallelized reciprocal blast searches to establish  
574 orthology and paralogy for genes within and between gene annotation files.

575

### 576 *Variant rates*

577 In order to estimate within-individual heterozygosity, we output a single  
578 pseudohaplotype using the 'style=pseudohap' parameter within Supernova from  
579 Sister 2 to represent the reference sequence. Next, we mapped the raw reads from  
580 all three individuals to the reference using BWA-MEM [52]. We then converted the  
581 resulting SAM files to BAM format using Samtools [53], and sorted and indexed them  
582 using Picard (Picard, RRID:SCR\_006525; <http://broadinstitute.github.io/picard/>).  
583 Realignment around insertion/deletion (indel) regions and duplicate marking was  
584 performed using GATK (GATK, RRID:SCR\_001876), and finally, we called  
585 heterozygous sites using a probabilistic framework implemented in ANGSD [54, 62,  
586 63]. We tested different posterior probability cutoffs (1, 0.999, 0.99, 0.98, and 0.95).  
587 To allow for comparison between all individuals, we down-sampled our three

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 assemblies to 20x mean nominal coverage (total number of reads covering a  
589 position, independent of their barcode) for our analyses. Heterozygosity was then  
590 simply calculated as the ratio of variable sites to the total number of sites (variable  
591 and invariable). Supernova also outputs the distance between heterozygous sites as  
592 part of their assembly report. We then used the read data of Campana et al. (2016)  
593 [21] and mapped them to our Sister 2 assembly to compare heterozygosity estimates  
594 (using the approach outlined above). Next, we estimated the number of shared  
595 heterozygous sites between a) our individuals and b) our individuals and the  
596 individuals from Campana et al. (2016) [21]. To do so, we used the *gplots* library in R  
597 (<https://www.r-project.org>) to calculate the overlap between the three sets and to  
598 display them in a Venn diagram.

599           Different pseudohaplotypes were obtained through the Supernova software  
600 by selecting either the '--style=pseudohap' or '--style=pseudohap2'. The two fasta  
601 files produced by 'pseudohap2' were then analyzed as described above.

602

### 603 *Demographic history*

604

605           We filtered each genome for putative X chromosome sequences by first  
606 aligning them to the domestic dog X scaffold [34]. Scaffolds showing significant  
607 alignment were then further filtered using the program BLAST [65]. The top hit for  
608 each alignment was chosen and all scaffolds which aligned with either the mouse,  
609 human, pig, domestic dog, or domestic cat X chromosome were removed. This was  
610 repeated for each assembly.

611

612           We then mapped the raw reads to the subset of scaffolds using BWA-MEM  
613 and called the consensus sequence using SAMtools and BCFtools  
614 (SAMtools/BCFtools, RRID:SCR\_005227) [59, 60]. Population history was  
615 reconstructed using PSMC and scaled using a mutations/site/generation rate of 6.0 x  
10<sup>-9</sup> and a generation time of 5 years [40]. This generation time a

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 mutation/site/generation rate was chosen because it was the average

617 mutation/site/generation rate inferred in Campana et al. (2016) [22].

618

#### 619 **Availability of supporting data**

620 Genomic and read data is available in the NCBI database under project accession

621 PRJNA488046. Further supporting data can be found in the *GigaScience* repository,

622 GigaDB [66].

623

#### 624 **Supporting Information**

625 Detailed information on methods, Supernova output, repeat annotation, gene

626 annotation, heterozygosity calculations, and different posterior probability cutoffs are

627 available online. The authors are solely responsible for the content and functionality

628 of these materials. Queries (other than absence of the material) should be directed to

629 the corresponding author.

630

#### 631 **Competing interests**

632 Author J. Stuelpnagel is a board member of 10x Genomics Inc. Author Ryan W.

633 Taylor is founder of End2End Genomics Inc.

634

#### 635 **Authors' contributions**

636 Authors JS, CSZ, PB, SP, EA, and DP conceived the project. Authors EM, HM, OM,

637 and RMC contributed samples and insight to the project. RT assembled the

638 genomes. EA and SP performed the genome annotation and downstream analyses.

639 EA, SP, CST, DP, and RT wrote the paper. All authors read and approved the final

640 manuscript.

641

#### 642 **Acknowledgements**



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

643 We thank M. Agnew, C. Asa, L. Padilla, and W. Warren for assistance in obtaining  
644 the Eureka sample. T. Linderoth, T. Korneliussen, and K. Bi for help with the different  
645 heterozygosity calculations. D. Church from 10x Genomics for discussion on how  
646 SuperNova performs the heterozygous site calling. We also thank the reviewers for  
647 their extremely helpful comments and suggestions in the improvement of this  
648 manuscript as well as the editor and GigaDB staff for the assistance in submitting the  
649 supporting data and refining the manuscript. This work was funded by a donation to  
650 the Program for Conservation Genomics at Stanford University.

651

## 652 **Literature Cited**

653

- 654 1. Pimm SL, Jenkins CN, Abell R, Brooks TM, Gittleman JL, Joppa LN,  
655 Raven PH, Roberts CM. and Sexton JO. The biodiversity of species and  
656 their rates of extinction, distribution, and protection. *Science*. 2014;  
657 344(6187):1246752.
- 658 2. Ceballos G, Ehrlich PR, Barnosky AD, García A, Pringle RM, Palmer TM.  
659 Accelerated modern human–induced species losses: Entering the sixth  
660 mass extinction. *Science advances*. 2015; 5:e1400253.
- 661 3. Epstein B, Jones M, Hamede R, Hendricks S, McCallum H, Murchison  
662 EP, Schönfeld B, Wiench C, Hohenlohe P, Storfer A. Rapid evolutionary  
663 response to a transmissible cancer in Tasmanian devils. *Nature*  
664 *communications*. 2016; 7:12684.
- 665 4. Harper C, Ludwig A, Clarke A, Makgopela K, Yurchenko A, Guthrie A,  
666 Dobrynin P, Tamazian G, Emslie R, van Heerden M, Hofmeyr M. Robust  
667 forensic matching of confiscated horns to individual poached African  
668 rhinoceros. *Current Biology*. 2018; 28(1):R13-4.
- 669 5. Steiner CC, Putnam AS, Hoeck PE, Ryder OA. Conservation genomics of  
670 threatened animal species. *Annu. Rev. Anim. Biosci*. 2013; 1(1):261-81.

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

671 6. Shafer AB, Wolf JB, Alves PC, Bergström L, Bruford MW, Brännström I,  
672 Colling G, Dalén L, De Meester L, Ekblom R, Fawcett KD. Genomics and  
673 the challenging translation into conservation practice. *Trends in Ecology &*  
674 *Evolution*. 2015; 30(2):78-87.6.

675 7. Girman DJ, Kat PW, Mills MG, Ginsberg JR, Borner M, Wilson V,  
676 Fanshawe JH, Fitzgibbon C, Lau LM, Wayne RK. Molecular genetic and  
677 morphological analyses of the African wild dog (*Lycaon pictus*). *Journal of*  
678 *Heredity*. 1993; 84(6):450-9.

679 8. Woodroffe R, Ginsberg J, and MacDonald DW. The African wild dog:  
680 status survey and conservation action plan. IUCN/SSC Canid Specialist  
681 Group. 1997; IUCN.

682 9. IUCN/SSC Regional conservation strategy for the cheetah and African  
683 wild dog in Southern Africa. IUCN. Species Survival Commission Gland.  
684 2007; IUCN.

685 10. Woodroffe R, Sillero-Zubiri C. *Lycaon pictus*. The IUCN Red List of  
686 Threatened Species. 2012;2012:e-T12436A167111116.

687 11. Woodroffe R and Ginsberg JR. Edge effects and the extinction of  
688 populations inside protected areas. *Science*. 1998; 280(5372):2126-2128.

689 12. Courchamp F, Clutton-Brock T, and Grenfell B. Inverse density  
690 dependence and the Allee effect. *Trends in Ecology & Evolution*. 1999;  
691 14(10):405-410.

692 13. Courchamp F, Clutton-Brock T, and Grenfell B. Multipack dynamics and  
693 the Allee effect in the African wild dog, *Lycaon pictus*. *Animal*  
694 *Conservation forum*. 2000; 3(4):277-285. Cambridge University Press.

695 14. McNutt JW and Silk JB. Pup production, sex ratios, and survivorship in  
696 African wild dogs, *Lycaon pictus*. *Behavioral Ecology and Sociobiology*.  
697 2008; 62(7):1061-1067.

- 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
- 698 15. McNutt JW. Sex-biased dispersal in African wild dogs, *Lycaon pictus*.  
699 Animal behaviour. 1996; 52(6):1067-1077.  
700 16. Fanshawe JH and Fitzgibbon CD. Factors influencing the hunting success  
701 of an African wild dog pack. Animal behaviour. 1993; 45(3):479-490.  
702 17. Creel S and Creel NM. Six ecological factors that may limit African wild  
703 dogs, *Lycaon pictus*. Animal Conservation. 1998; 1(1):1-9.  
704 18. Creel S and Creel NM. Opposing effects of group size on reproduction  
705 and survival in African wild dogs. Behavioral Ecology. 2015; 26(5):1414-  
706 1422.  
707 19. Girman DJ, Vila C, Geffen E, Creel S, Mills MG, McNutt JW, Ginsberg JK,  
708 Kat PW, Mamiya KH, Wayne RK. Patterns of population subdivision, gene  
709 flow and genetic variability in the African wild dog (*Lycaon pictus*).  
710 Molecular Ecology. 2001; 10(7):1703-23.  
711 20. Marsden CD, Woodroffe R, Mills MG, McNutt JW, Creel S, Groom R,  
712 Emmanuel M, Cleaveland S, Kat P, Rasmussen GS, Ginsberg J. Spatial  
713 and temporal patterns of neutral and adaptive genetic variation in the  
714 endangered African wild dog (*Lycaon pictus*). Molecular Ecology. 2012;  
715 21(6):1379-93.  
716 21. Marsden CD, Mable BK, Woodroffe R, Rasmussen GS, Cleaveland S,  
717 McNutt JW, Emmanuel M, Thomas R, Kennedy LJ. Highly endangered  
718 African wild dogs (*Lycaon pictus*) lack variation at the major  
719 histocompatibility complex. Journal of heredity. 2009; 100:S54-65.  
720 22. Campana MG, Parker LD, Hawkins MT, Young HS, Helgen KM, Gunther  
721 MS, Woodroffe R, Maldonado JE, Fleischer RC. Genome sequence,  
722 population history, and pelage genetics of the endangered African wild  
723 dog (*Lycaon pictus*). BMC genomics. 2016; 17(1):1013.

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

724 23. Shapiro B and Hofreiter M. A paleogenomic perspective on evolution and  
725 gene function: new insights from ancient DNA. *Science*. 2014;  
726 343(6169):1236573.

727 24. Lindblad-Toh K, Wade CM, Mikkelsen TS, Karlsson EK, Jaffe DB, Kamal  
728 M, Clamp M, Chang JL, Kulbokas III EJ, Zody MC, Mauceli E. Genome  
729 sequence, comparative analysis and haplotype structure of the domestic  
730 dog. *Nature*. 2005; 438(7069):803.

731 25. Perini FA, Russo CA, Schrago CG. The evolution of South American  
732 endemic canids: a history of rapid diversification and morphological  
733 parallelism. *Journal of evolutionary biology*. 2010; 23(2):311-22.

734 26. Koepfli KP, Pollinger J, Godinho R, Robinson J, Lea A, Hendricks S,  
735 Schweizer RM, Thalmann O, Silva P, Fan Z, Yurchenko AA. Genome-  
736 wide evidence reveals that African and Eurasian golden jackals are  
737 distinct species. *Current Biology*. 2015; 25(16):2158-65.

738 27. International Human Genome Sequencing Consortium. Initial sequencing  
739 and analysis of the human genome. *Nature*. 2001; 409(6822):860.

740 28. Hayden, E.C., The \$1,000 genome. *Nature*, 2014. 507(7492): p. 294.  
741 Hayden EC. Is the \$1,000 genome for real?. *Nature News*. 2014 Jan 15.

742 29. Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of  
743 next-generation sequencing technologies. *Nature Reviews Genetics*.  
744 2016; 17(6):333.

745 30. Eklom R, Wolf JB. A field guide to whole-genome sequencing, assembly  
746 and annotation. *Evolutionary applications*. 2014; 7(9):1026-42.

747 31. Putnam NH, O'Connell BL, Stites JC, Rice BJ, Blanchette M, Calef R,  
748 Troll CJ, Fields A, Hartley PD, Sugnet CW, Haussler D. Chromosome-  
749 scale shotgun assembly using an in vitro method for long-range linkage.  
750 *Genome research*. 2016; 26(3):342-50.

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

751 32. Burton JN, Adey A, Patwardhan RP, Qiu R, Kitzman JO, Shendure J.  
752 Chromosome-scale scaffolding of de novo genome assemblies based on  
753 chromatin interactions. *Nature biotechnology*. 2013; 31(12):1119.

754 33. Weisenfeld NI, Kumar V, Shah P, Church DM, Jaffe DB. Direct  
755 determination of diploid genome sequences. *Genome research*. 2017;  
756 27(5):757-67.

757 34. Hoepfner MP, Lundquist A, Pirun M, Meadows JR, Zamani N, Johnson J,  
758 Sundström G, Cook A, FitzGerald MG, Swofford R, Mauceli E. An  
759 improved canine genome and a comprehensive catalogue of coding  
760 genes and non-coding transcripts. *PLoS one*. 2014; 9(3):e91172.

761 35. Gopalakrishnan S, Castruita JA, Sinding MH, Kuderna LF, Räikkönen J,  
762 Petersen B, Sicheritz-Ponten T, Larson G, Orlando L, Marques-Bonet T,  
763 Hansen AJ. The wolf reference genome sequence (*Canis lupus lupus*)  
764 and its implications for *Canis* spp. population genomics. *BMC genomics*.  
765 2017; 18(1):495.

766 36. Lok S, Paton TA, Wang Z, Kaur G, Walker S, Yuen RK, Sung WW,  
767 Whitney J, Buchanan JA, Trost B, Singh N. De novo genome and  
768 transcriptome assembly of the Canadian beaver (*Castor canadensis*). *G3:  
769 Genes, Genomes, Genetics*. 2017; 7(2):755-73.

770 37. Mohr DW, Naguib A, Weisenfeld N, Kumar V, Shah P, Church DM, Jaffe  
771 D, Scott AF. Improved de novo Genome Assembly: Synthetic long read  
772 sequencing combined with optical mapping produce a high quality  
773 mammalian genome at relatively low cost. *bioRxiv*. 2017; 128348.

774 38. Dobrynin, P., et al., Genomic legacy of the African cheetah, *Acinonyx  
775 jubatus*. *Genome biology*, 2015. 16(1): p. 277.

776 39. Robinson JA, Ortega-Del Vecchyo D, Fan Z, Kim BY, Marsden CD,  
777 Lohmueller KE, Wayne RK. Genomic flatlining in the endangered island  
778 fox. *Current Biology*. 2016; 26(9):1183.

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

779 40. Li H, Durbin R. Inference of human population history from individual  
780 whole-genome sequences. *Nature*. 2011; 475(7357):493.

781 41. Orozco-terWengel P. The devil is in the details: the effect of population  
782 structure on demographic inference. *Heredity*. 2016 Apr;116(4):349.

783 42. Liu S, Lorenzen ED, Fumagalli M, Li B, Harris K, Xiong Z, Zhou L,  
784 Korneliussen TS, Somel M, Babbitt C, Wray G. Population genomics  
785 reveal recent speciation and rapid evolutionary adaptation in polar bears.  
786 *Cell*. 2014; 157(4):785-94.

787 43. Huang J, Zhao Y, Shiraigol W, Li B, Bai D, Ye W, Daidiikhuu D, Yang L,  
788 Jin B, Zhao Q, Gao Y. Analysis of horse genomes provides insight into  
789 the diversification and adaptive evolution of karyotype. *Scientific reports*.  
790 2014; 4:4958.

791 44. Gnerre S, MacCallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ,  
792 Sharpe T, Hall G, Shea TP, Sykes S, Berlin AM. High-quality draft  
793 assemblies of mammalian genomes from massively parallel sequence  
794 data. *Proceedings of the National Academy of Sciences*. 2011;  
795 108(4):1513-8.

796 45. Frantzen MA, Silk JB, Ferguson JW, Wayne RK, Kohn MH. Empirical  
797 evaluation of preservation methods for faecal DNA. *Molecular Ecology*.  
798 1998; 7(10):1423-8.

799 46. Taberlet P, Luikart G. Non-invasive genetic sampling and individual  
800 identification. *Biological journal of the linnean society*. 1999; 68(1-2):41-  
801 55.

802 47. Morin PA, Luikart G, Wayne RK. SNPs in ecology, evolution and  
803 conservation. *Trends in Ecology & Evolution*. 2004; 19(4):208-16.

804 48. Vieira FG, Fumagalli M, Albrechtsen A, Nielsen R. Estimating inbreeding  
805 coefficients from NGS data: impact on genotype calling and allele  
806 frequency estimation. *Genome research*. 2013; 23(11):1852-61.

- 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
- 807 49. Pazmiño DA, Maes GE, Simpfendorfer CA, Salinas-de-León P, van  
808 Herwerden L. Genome-wide SNPs reveal low effective population size  
809 within confined management units of the highly vagile Galapagos shark  
810 (*Carcharhinus galapagensis*). *Conservation Genetics*. 2017; 18(5):1151-  
811 63.
- 812 50. Hampton JO, Spencer P, Alpers DL, Twigg LE, Woolnough AP, Doust J,  
813 Higgs T, Pluske J. Molecular techniques, wildlife management and the  
814 importance of genetic population structure and dispersal: a case study  
815 with feral pigs. *Journal of Applied Ecology*. 2004; 41(4):735-43.
- 816 51. Abascal F, Corvelo A, Cruz F, Villanueva-Cañas JL, Vlasova A, Marcet-  
817 Houben M, Martínez-Cruz B, Cheng JY, Prieto P, Quesada V, Quilez J.  
818 Extreme genomic erosion after recurrent demographic bottlenecks in the  
819 highly endangered Iberian lynx. *Genome biology*. 2016; 17(1):251.
- 820 52. Kim S, Cho YS, Kim HM, Chung O, Kim H, Jho S, Seomun H, Kim J,  
821 Bang WY, Kim C, An J. Comparison of carnivore, omnivore, and  
822 herbivore mammalian genomes with a new leopard assembly. *Genome*  
823 *biology*. 2016; 17(1):211.
- 824 53. Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I,  
825 Boisvert S, Chapman JA, Chapuis G, Chikhi R, Chitsaz H. Assemblathon  
826 2: evaluating de novo methods of genome assembly in three vertebrate  
827 species. *GigaScience*. 2013; 2(1):10.
- 828 54. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM.  
829 BUSCO: assessing genome assembly and annotation completeness with  
830 single-copy orthologs. *Bioinformatics*. 2015; 31(19):3210-2.
- 831 55. Smit AF, Hubley R, Green P. 1996–2010. RepeatMasker Open-3.0.
- 832 56. Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz  
833 J. Repbase Update, a database of eukaryotic repetitive elements.  
834 *Cytogenetic and genome research*. 2005; 110:462-7.

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

835 57. Holt C, Yandell M. MAKER2: an annotation pipeline and genome-  
836 database management tool for second-generation genome projects. BMC  
837 bioinformatics. 2011; 12(1):491.

838 58. Lechner M, Hernandez-Rosales M, Doerr D, Wieseke N, Thévenin A,  
839 Stoye J, Hartmann RK, Prohaska SJ, Stadler PF. Orthology detection  
840 combining clustering and synteny for very large datasets. PLoS One.  
841 2014; 9(8):e105015.

842 59. Li H, Durbin R. Fast and accurate short read alignment with Burrows-  
843 Wheeler transform. Bioinformatics. 2009;25(14):1754-60.

844 60. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G,  
845 Abecasis G, Durbin R. The sequence alignment/map format and  
846 SAMtools. Bioinformatics. 2009; 25(16):2078-9.

847 61. Korneliussen TS, Albrechtsen A, Nielsen R. ANGSD: analysis of next  
848 generation sequencing data. BMC bioinformatics. 2014; 15(1):356.

849 62. Nielsen R, Paul JS, Albrechtsen A, Song YS. Genotype and SNP calling  
850 from next-generation sequencing data. Nature Reviews Genetics. 2011;  
851 12(6):443.

852 63. Nielsen R, Korneliussen T, Albrechtsen A, Li Y, Wang J. SNP calling,  
853 genotype calling, and sample allele frequency estimation from new-  
854 generation sequencing data. PloS one. 2012; 7(7):e37558.

855 64. Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C,  
856 Salzberg SL. Versatile and open software for comparing large genomes.  
857 Genome biology. 2004; 5(2):R12.65.

858 65. Altschul SF, Madden TL, Schäffer AA, Zhang J, Zhang Z, Miller W,  
859 Lipman DJ. Gapped BLAST and PSI-BLAST: a new generation of protein  
860 database search programs. Nucleic acids research. 1997; 25(17):3389-  
861 402.



862 66. Armstrong E, Taylor RW, Prost S, Blinston P, van der Meer E,  
 863 Madzikanda H et al. Supporting data for "Entering the era of conservation  
 864 genomics: Cost-effective assembly of the African wild dog genome using  
 865 linked reads" GigaScience Database. 2018  
 866 <http://dx.doi.org/10.5524/100475>

867  
 868 **Tables**

869  
 870 **Table 1. Assembly Statistics.** Assembly statistics for the three African wild dog  
 871 genomes reported by the Supernova assembler. Coverage was assessed using  
 872 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 (mb)	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

873  
 874

875  
 1 876  
 2 877  
 3  
 4 878  
 5  
 6 879  
 7  
 8 880  
 9  
 10  
 11 881

**Table 2. Conserved Gene Statistics.** Results of the BUSCO v2 gene annotation from three African wild dog genome assemblies, canFam3.1, low-coverage wild dog genomes [22], the recently published wolf genome [35] and the Hawaiian monk seal genome [37].

Assembly	Species	Complete	Single copy	Duplicated	Fragmented	Missing	Total searched
Sister 1	<i>L. pictus</i>	3914	3875	39	102	88	4104
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
Wolf	<i>C. lupus</i>	3890	3849	41	110	104	4104
Hawaiian monk seal	<i>Neomonachus schauinslandi</i>	3881	3833	48	118	105	4104

882  
 29  
 30 883  
 31  
 32 884  
 33 885  
 34 886  
 35 887  
 36 888

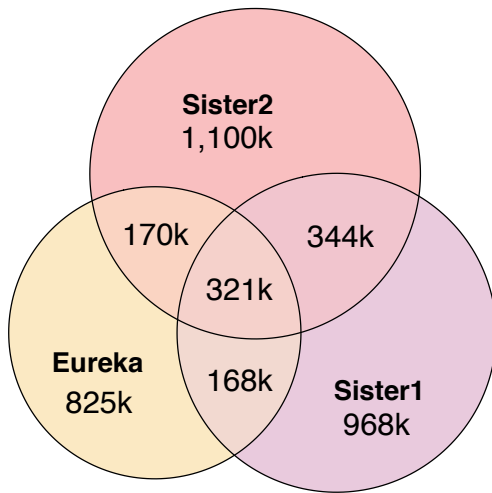
**Figure 1.** (A) Pack of African wild dogs. B) Shared heterozygous sites between the three *de novo* assemblies (calculated using a posterior cutoff of 0.99). More of the heterozygous sites are shared between the two sisters than between either sister and Eureka. C) PSMC reconstruction of the individuals' demographic history. Bootstrap replicates are plotted in lighter colors. Time is in years before present.

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

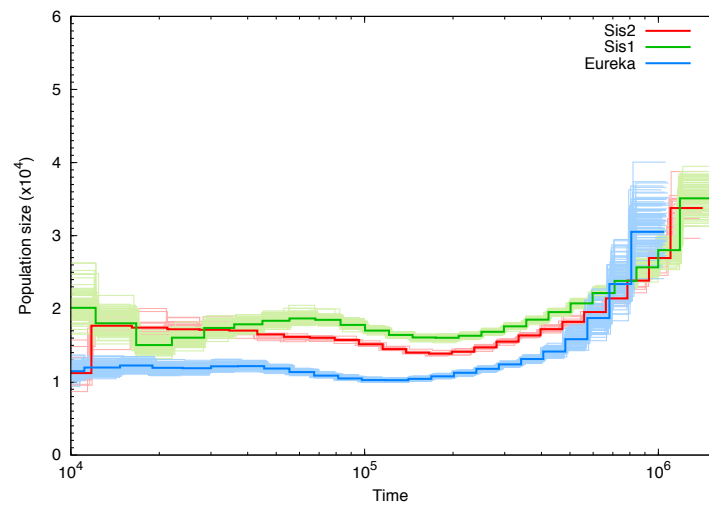
A



B



C





Click here to access/download

**Supplementary Material**

Supporting\_information\_AWD\_Gigascience\_final\_update  
.docx

