# GigaScience

# A 3-way hybrid approach to generate a new high quality chimpanzee reference genome (Pan_tro_3.0)

## --Manuscript Draft--

| | |
|---|---|
| Manuscript Number: | GIGA-D-16-00140R1 |
| Full Title: | A 3-way hybrid approach to generate a new high quality chimpanzee reference genome (Pan_tro_3.0) |
| Article Type: | Data Note |
| Funding Information: | |

| | |
|---|---|
| Abstract: | **Background:**<br>The chimpanzee is arguably the most important species for the study of human origins. A key resource for these studies is a high quality reference genome assembly, however, as most mammalian genomes, the current iteration of the chimpanzee reference genome assembly it is highly fragmented. In the current iteration of the chimpanzees reference genome assembly (Pan_tro_2.1.4), the sequence is scattered across more then 183,000 contigs and incorporating over 159,000 gaps, with a genome wide contig N50 of 51 Kbp.<br><br>**Findings:**<br>In this work we produce an extensive and diverse array of sequencing datasets to rapidly assemble a new chimpanzee reference that surpasses previous iterations in bases represented and organized in large scaffolds. To this end, we show substantial improvements over the current release of the chimpanzee genome (Pan_tro_2.1.4) by several metrics, such as: increased contiguity by >750% and 300% on contigs and scaffolds, respectively; closure of 77% of gaps in the Pan_tro_2.1.4 assembly gaps spanning >850 Kbp of novel coding sequence based on RNASeq data. We furthermore report over 2,700 genes that had putatively erroneous frame-shift predictions to human in Pan_tro_2.1.4 and show a substantial increase in the annotation of repetitive elements.<br><br>**Conclusions:**<br>We apply a simple 3-way hybrid approach to considerably improve the reference genome assembly for the chimpanzee, providing a valuable resource to study human origins. We furthermore produced extensive sequencing datasets that are all derived from the same cell line, generating a broad non-human benchmark dataset. |

| Corresponding Author: | Lukas F. K. Kuderna |
| --- | --- |
| | SPAIN |
| Corresponding Author Secondary Information: | |
| Corresponding Author's Institution: | |
| Corresponding Author's Secondary Institution: | |
| First Author: | Lukas F. K. Kuderna |
| First Author Secondary Information: | |
| Order of Authors: | Lukas F. K. Kuderna |
| | Chad Tomlinson |
| | LaDeana W. Hillier |
| | Annabel Tran |
| | Ian Fiddes |
| | Joel Armstrong |
| | Hafid Laayouni |
| | David Gordon |
| | John Huddleston |
| | Raquel Garcia Perez |
| | Inna Povolotskaya |
| | Aitor Serres Armero |
| | Jèssica Gómez Garrido |
| | Daniel Ho |
| | Paolo Ribeca |
| | Tyler Alioto |
| | Richard E. Green |
| | Benedict Paten |
| | Arcadi Navarro |
| | Jaume Betranpetit |
| | Javier Herrero |
| | Evan E. Eichler |
| | Andrew J. Sharp |
| | Lars Feuk |
| | Wesley C. Warren |
| | Tomas Marques-Bonet |
| Order of Authors Secondary Information: | |

| Response to Reviewers: | Dear Dr. Hans Zauner, |
| --- | --- |
| | Please find attached the resubmission of our revised manuscript GIGA-D-16-00140; A 3-way hybrid approach to generate a new high quality chimpanzee reference genome (Pan_tro_3.0), together with a point by point response of the comments of all |

reviewers. We believe that we satisfactorily address all the comments in there.

We also wanted to report an issue that has been brought to us after a revision on the assembly. We have detected several small inversions affecting either full contigs within a scaffold, or the end of a contig. Because these events are predominantly rather small, they have previously escaped our notice when assessing (large scale) structural errors using clones, as they are only detectable by fine scale comparisons to the previous chimpanzee assembly, as well as the human genome assembly. Nevertheless, in the spirit of full disclosure, we did not want to resubmit our revised manuscript without addressing the issue first, now included in the final version of the assembly. Altogether, this affects 2,990 fully inverted contigs, amounting for 20.5 Mb of sequence. These were flipped and left in place in a new version of the assembly. Furthermore, there are 1,505 pieces of contigs, where the breakpoint of the inversion lays within the contigs, amounting for 11.1 Mb of sequence. In these cases, we manually went through the flips and decided to move some off to the unplaced portion of the assembly, and inverted some in place, depending on how clearly we could assign them to belonging to the chromosome in question. Cumulatively, there are now 31.6Mb of sequence in the assembly were we have changed the orientation. This affects 1,938 of coding exons amounting for 276,641 bp of coding sequence. Please note that the overall sequence content of the assembly has remained the same.

We are currently in the process of resubmitting our revised assembly to NCBI. For the paper, we have chosen to explicitly state which version of the NCBIs accessions system is being used.

We are looking forward to your reply.

Kind Regards,

Lukas Kuderna & Tomas Marques-Bonet


########

#Please note we are attaching a formatted document for the reviews with all graphics in place. Below you will find the copied text only from that document.


Point by point rebuttal:

Reviewer 1:
Reviewer #1: The authors present several large new datasets of chimpanzee genome sequencing data, and they combine these datasets into a novel, high-quality genome assembly of Pan troglodytes.  As the authors state, this is a valuable addition to the set of available genome sequence resources and a vast improvement in genome quality over the existing Pan troglodytes assembly.  The manuscript needs some editing and cleaning up, but overall I believe it represents a significant contribution to the field and should eventually be published.

In the "Data description" section, the paper gives an overview of the datasets the authors used.  This section would benefit from a clear introduction and description of the sequencing strategies they employed to process these datasets.  I suggest that a new figure, in the form of a simple flowchart, could be a helpful visual aid: it would describe the assembly methods that were used to combine the various types of sequencing libraries, and it would illustrate the process of creating the 3-way hybrid intermediary assembly as well as the final (3.0) assembly.  Additionally, the "Data description" section is mostly devoid of citations.  More citations should be added in order to give proper attribution to the developers of the assembly methods, and to enable the reader to seek more information.

We have sought to clarify on the different sequencing strategies, and have added references where adequate in this section (Goodwin S et al. ,2016; Kuleshov et al, 2014; Putnam et al, 2016). We have also added the flowchart of the assembly process in supplementary figure 1 (see below).  Nevertheless, we have only slightly modified

the introductory section, as it is our understanding that this is adequate for the data note format.

The authors discuss the sequence content they have added to the chimpanzee genome.  It's interesting to see the length distribution of the gaps they have filled (Figure 1C), and I would be curious to see comparative length distributions for gaps they failed to fill, or for gaps they added.

We have added plots of length distributions for both of these cases. Supplementary figure S19 (first figure below) shows the length distribution of gaps we can identify as corresponding between Pantro_2.1.4 and Pan_tro_3, but fail to fill. Supplementary Figure S20 (second figure below) shows the length distribution of gaps present within Pan_tro_3.0. We note, that the overall shape of the distribution is similar, with peaks at small gap sizes.

The detail on the repeat resolution is also fascinating.  I think the authors sell themselves short by noting that the repeat fraction of the assembly increases from 50.9% to 52.2%: given that they only increase the assembly sequence length by ~8%, this actually shows that most of the sequence they've added is repeat sequence, which is a useful indicator of the new assembly's added value.  Similarly, Figure 1D, which shows the quantities of added repeat sequence for various repeat types, would be stronger if it also showed the quantities of already-existing repeat sequence for each type.

We have added a plot showing the full comparative repeat content of both Pan_tro2.1.4 and Pan_tro3 at supplementary Figure S21 (first figure below) , as well as a scaled version for repeat families with fewer annotations at Figure S22 (second figure below).

The authors compare the new (3.0) genome assembly to the existing (2.1.4) assembly. They observe a 99.9% overall sequence similarity and note that the 0.1% differences could be explained by SNPs; it would be interesting to see a deeper analysis of these SNPs, although this may be outside the scope of the manuscript.

While we agree that this would be an interesting exploration, we believe, in accordance with the reviewer, that an analysis of these SNPs is out of the scope of this manuscripts' format, especially considering its submission in the form of a GigaScience Data Note.

 Also, in the section "Gene annotation", they note a large number of genes with frameshift mutations between the 2.1.4 assembly and the human genome assembly. This is striking, but a fully fair comparison would also mention the number of genes that also contain frameshift mutations (perhaps newly added frameshift mutations) in the 3.0 assembly.

We have now also included the count of frame-shift mutations specific for Pantro_3 for a fully fair comparison in this section (674). We have furthermore clarified, that these frame shifts are not necessarily due to sequence errors in Pan_tro_2.1.4, but might

also constitute allelic variation, as the assembly only randomly captures one of two alleles at a given locus.

The conclusion is strong, but it would be stronger with some additional context that describes the achievement in this manuscript. The genome assembly is higher quality. But is it also more efficient, or more economic? Have the authors innovated any new genome assembly methods? Have they demonstrated a technique that could be easily applied to other genome assemblies?

We clarify that our approach should easily be applicable to genomes of similar complexity. Nevertheless, we would like to refrain to comment on efficiency or economical value of the assembly for two reasons: First, the price for sequencing on several platforms has been shown to be extremely dependent on the time of sequence production, and has even dramatically decreased within the timeframe of this manuscript. Second, efficiency in the context of genome assembly is a rather subjective issue, as it not straightforward to decide what to measure efficiency against.

Minor errors:

Section "Assembly generation": "These reads are derived from a 400 bps library, resulting in pairs that overlap over a ~50 bps region". If a 400-bp fragment is sequenced to 250 bp from both ends, wouldn't that result in an overlap of ~100 bp rather than ~50 bp?

We thank the reviewer for catching this error, it is now corrected. The library was size selected to around 450bp, resulting in an overlap of around 50 bp when sequencing 250bp on each side.

Section "Assembly generation": "we observed superior connectivity". The word "connectivity" is unclear in this context; it might be better to simply repeat "contiguity".

We have clarified this sentence by rephrasing it.

Section "Repeat resolution": "We furthermore added 38.2 Mbp of LINE to the assembly, corresponding to over 44,791 additional copies of L1 elements." First of all, this should say "LINEs" rather than "LINE". Secondly, these numbers do not add up. A typical L1 element is 6 Kbp in length; thus, 44,791 copies of L1 elements should necessarily occupy over 260 Mbp of sequence.

We have corrected this flaw by leaving out that number. We erroneously counting L1 annotations, that do not necessarily correspond to fully resolved copies of L1 elements, as repeatmasker annotates partial matches of repetitive elements.

Section "Resolution of segmental duplications": This section should contain more citations, especially for the WGAC and WSSD methods, which are named but not described at all.

We have added references to both these methods in the corresponding section (Bailey et al. 2001; Bailey et al 2002).

Reviewer #2:

The genome of a chimpanzee is an important asset for the study of human evolution, and an high quality reference assembly is long overdue. The authors were able to significantly improve on the previous assemblies. I have no problem with the 3-way hybrid approach they have taken. The paper is written very clearly, and is fully appropriate for the data note format. My recommendation is that this manuscript is accepted without reservation.

We would like to thank the reviewer for the positive feedback.

Reviewer 3:

1. Lines 191 to 212. Is this an exploration of segmental duplications in the genome or over-representation and under-representation in the assembly? What is the support for the conclusion that "segmental duplications are well resolved." The statement of the conclusion is very confusing: "we are likely to be overestimating ... by including an elevated rate of false positive paralogous regions…"

In this section, we explore the representation of segmental duplications within our assembly. We sought to clarify the confusion by rephrasing the concluding sentences of the paragraph to the following:

We then compared Pant_tro_3.0 to the human reference genome assembly GRCh38, an assembly that is based on a BAC hierarchical shotgun assembly strategy and may therefore be considered of gold standard with respect to representation of segmental duplications. We note similar proportions of bases in segmental duplications on chromosomal scaffolds (4,46% in Pan_tro_3.0 vs. 5,56% in GRCh38), however, we note an elevated genome wide rate of bases in duplications when including unplaced and unlocalized scaffolds, suggesting that our assembly includes false-positive paralogous regions within them (see supplementary Table 1).

By this means, we hope to clarify the questions of the reviewer: Our previous statements about segmental duplications being 'well resolved' referred to the comparable number of bases in segmental duplications between the Chimp and the Human assembly.

2. Lines 229 to 233. The finding is called "most striking" but it is accompanied by weak interpretation ("majority of … putative"). A little more investigation would probably support a strong claim of improvement. To estimate the veracity of the old frameshifts, clarify in what sense both assemblies are "mainly" based on data from the same individual, and measure if any frameshifted genes relied on reads from another individual. To estimate veracity of the new assemblies of these genes, rule out allelism in spanning reads, count framehshifts relative to human present in both assemblies, and count frameshifts exclusively in the new assembly. Presumably these counts are low.

We have now toned down this claim. We clarify to what extend the assemblies are based on the genomes of the same individual, and what proportions are derived from a different in Pan_tro_2.1.4 . We furthermore clarify, that these frameshifts don't necessarily constitute fixed changes, but might also be due to allelic variation. Furthermore, we have included the number of genes with predicted frameshifts only in Pan_tro_3, but not in Pan_tro_2.1.4

3. Line 125 "17 scaffold errors". Extrapolate the overall structural error rate using the number of bases spanned by fosmids. Extrapolate the likely number of remaining structural errors. Of errors fixed, where they attributable to the contig or the scaffold process specifically? Were the errors near particularly repetitive sequence?

Out of 671,716 fosmids with available end sequences for both ends from the CHORI-1251 library, 545,788 (81%) mapped with both ends and high quality (mismapping rate < 0.00001, MQ<=60) onto the same scaffold. Out of these, we find 539,315 (99%) to map with both ends in concordant orientation, and 6,473 (1%) with both ends in the discordant orientation. Cumulatively, these concordant mappings cover 2.7 Gbps (85%) of the whole assembly length. As PBJelly changes the naming scheme of the sequences after filling gaps, we could not deduce at which stage the structural errors we fixed where introduced.

4. Line 129 and 232 "500,000 SNPs … frameshift". Explanation, investigation, or speculation would help. What caused SNPs in contigs relative to the Illumina reads used to create the base assembly? What caused frameshifts in the prior assembly?

We now try to clarify by speculating that most of these corrected errors are due to regions where PacBio data was incorporated into the assembly (line 129). Given the relatively low coverage of PacBio data, we did not apply self correction on the PacBio reads, but rather corrected the assembly after filling in gaps. This leads to an elevated rate of residual errors in regions derived from PacBio data. We hypothesize that comparatively few of these corrected errors lay within regions derived from Illumina data. Given the constant change of coordinate system between the assemblies (with each incorporated platform) it is not straightforward to know which region in the assembly is finally derived from what platform.

We clarify that frameshifts with respect to Pantro2.1.4 are either because of allelic variation or sequencing errors in the Sanger data used to assemble Pantro2.1.4.

5. Line 65 "SMRTcells … synthetic long reads". The manuscript does not address issues of integrating these technologies. Did either require error correction? Was there ever any overlap or disagreement between these two types of long reads? Was either more helpful than the other?

We had included an analysis of gap-filling performance and repeat resolution for PacBio and TruSeq SLR in the supplementary section S2, where we compare how well gaps in Pantro-2.1.4 are resolved using only either technology, as well as a combination of both. We did not incorporate the Truseq SLR data into the assembly based on the observed high rate of repeat collapse (see supplementary sections S2, supplementary Figures S2 and S3). Indeed, we see that many common high identity repeats are under-represented in sequencing data derived from this platform. We did not pre-correct the PacBio data, but rather run a post-assembly error correction, as described in the manuscript.

6. Line 223 "paralogous coding duplications are better represented". What was the read coverage of these regions? Are these duplications specific to the chimp lineage or ancestral to primates? How is a paralogous coding duplication different from other kinds?

We now refrain to claim that, paralogous coding duplications are better represented' as although we observe a shift to higher read depths in these regions, some of the newly added paralogs do no validate by excessive read depth.

7. Line 145 "bringing its continuity to the range". The X chromosome N50 (422K) is actually larger than the average (385K). If the old assembly had smaller contigs at X due to half coverage of the X chromosome, then why isn't that factor at play in the new assembly?

We believe this to be the case mainly due to two reasons: First, following a back of the envelope calculation with the Lander-Waterman equation, there are about 5% of bases without a single read when sampling them at a 3X coverage. This relationship is non-linear with respect to coverage, and the number of unsampled bases drops to essentially 0 at a 30X coverage. Second, because of the initially poor assembly quality on the X, many BACs have been finished and integrated into the assembly. These BACs were also used for the final AGP creation of our assembly, boosting contiguity on this chromosome.

8. Line 126. Were SNPs concentrated in the gap fill regions? In the gaps filled with low-coverage long reads?

We speculate that this is most likely the case. However, given the constant change in coordinate system during the assembly, and also during the correction process itself, it is not straigthforward for us to keep track of the origin of each genomic region within the assembly.

9. Line 74 "finished BAC". The BACs never get mentioned again. Do the old and new assemblies agree with the BACs?

We now clarify that these BACs were integrated into Pan_tro_2.1.4 as well as the finished version of Pan_tro_3. Thus, by definition, the final assemblies agree with the BACs

| | 10. Line 76 "unprecedented". Is there any need to make this controversial claim? |
| --- | --- |
| | We have rephrased the sentence to tone down the claim. |
| | 11. Lines 89, 98, 128 "base assembly". The DISCOVAR assembly is referenced by several names, some of which I confused with 2.1.4. Assign it a name or number? |
| | The DISCOVAR base assembly is now consistently refered to as ‚DISCOVAR base assembly' |
| | 12. Line 92 should clarify this is a scaffold N50. |
| | We clarified this regards the scaffold N50. |
| | 13. Line 95 "remaining gap structure required us to". What is a gap structure? In what way was a response required? |
| | We have rephrased this sentence to make it clearer. |
| **Additional Information:** | |
| **Question** | **Response** |
| Are you submitting this manuscript to a special series or article collection? | No |
| **Experimental design and statistics**<br><br>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.<br><br>Have you included all the information requested in your manuscript? | Yes |
| **Resources**<br><br>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.<br><br>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist? | Yes |
| **Availability of data and materials** | Yes |

All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the "Availability of Data and Materials" section of your manuscript.

Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?

1  # A 3-way hybrid approach to generate a new high quality chimpanzee

2  # reference genome (Pan_tro_3.0)

3  Lukas F.K. Kuderna[1,2], Chad Tomlinson[3], LaDeana W. Hillier[3], Annabel

4  Tran[4], Ian Fiddes[5], Joel Armstrong[5], Hafid Laayouni[1,6], David Gordon[7], John

5  Huddleston[7], Raquel Garcia Perez[1], Inna Povolotskaya[1], Aitor Serres Armero[1],

6  Jèssica Gómez Garrido[2], Daniel Ho[8], Paolo Ribeca[9], Tyler Alioto[2], Richard E.

7  Green[10,11], Benedict Paten[5], Arcadi Navarro[1,2,12], Jaume Betranpetit[1], Javier

8  Herrero[4], Evan E. Eichler[7], Andrew J. Sharp[8], Lars Feuk[13,*], Wesley C.

9  Warren[3,*], Tomas Marques-Bonet[1,2,12*]

10

11  (1) Institut de Biologia Evolutiva, (CSIC-Universitat Pompeu Fabra), PRBB, Doctor Aiguader 88,

12  Barcelona, Catalonia 08003, Spain.

13  (2) CNAG-CRG, Centre for Genomic Regulation (CRG), Baldiri i Reixac 4, 08028, Barcelona, Spain.

14  (3) McDonnell Genome Institute, Department of Medicine, Department of Genetics, Washington

15  University School of Medicine, St. Louis, MO 63108, USA.

16  (4) Bill Lyons Informatics Centre, UCL Cancer Institute, University College London, London, UK.

17  (5) Genomics Institute, University of California Santa Cruz and Howard Hughes Medical Institute,

18  Santa Cruz, CA 95064, USA.

19  (6) Bioinformatics Studies, ESCI-UPF, Pg. Pujades 1, 08003, Barcelona, Spain.

20  (7) Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA

21  98195, USA.

22  (8) Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New

23  York, NY 10029, USA

24  (9) The Pirbright Institute, Ash Road, Pirbright, Woking, GU24 0NF, United Kingdom

25  (10) Department of Biomolecular Engineering, University of California Santa Cruz, 1156 High Street,

26  Santa Cruz, CA 95060, USA.

27  (11) Dovetail Genomics, Santa Cruz, CA 95060, USA

28   (12) Institucio Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Catalonia 08010, Spain

29 (13) Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala

30 University, Uppsala, Sweden

31 **Abstract**

32 **Background**

33 The chimpanzee is arguably the most important species for the study of human

34 origins. A key resource for these studies is a high quality reference genome assembly,

35 however, as most mammalian genomes, the current iteration of the chimpanzee

36 reference genome assembly it is highly fragmented. In the current iteration of the

37 chimpanzees reference genome assembly (Pan_tro_2.1.4), the sequence is scattered

38 across more then 183,000 contigs and incorporating over 159,000 gaps, with a

39 genome wide contig N50 of 51 Kbp.

40 **Findings**

41 In this work we produce an extensive and diverse array of sequencing datasets to

42 rapidly assemble a new chimpanzee reference that surpasses previous iterations in

43 bases represented and organized in large scaffolds. To this end, we show substantial

44 improvements over the current release of the chimpanzee genome (Pan_tro_2.1.4) by

45 several metrics, such as: increased contiguity by >750% and 300% on contigs and

46 scaffolds, respectively; closure of 77% of gaps in the Pan_tro_2.1.4 assembly gaps

47 spanning >850 Kbp of novel coding sequence based on RNASeq data. We

48 furthermore report over 2,700 genes that had putatively erroneous frame-shift

49 predictions to human in Pan_tro_2.1.4 and show a substantial increase in the

50 annotation of repetitive elements.

2

## Conclusions

We apply a simple 3-way hybrid approach to considerably improve the reference genome assembly for the chimpanzee, providing a valuable resource to study human origins. We furthermore produced extensive sequencing datasets that are all derived from the same cell line, generating a broad non-human benchmark dataset.

## Keywords

Chimpanzee reference genome, Assembly, Genomics

## Data description

### Creating a non-human sequencing benchmark dataset

To test the potentially combinatorial power of varied sequencing and mapping strategies, we created several different datasets on different platforms, to try to leverage the advantages of each, as the shortcomings of one sequencing strategy might be compensated by another one [1]. All datasets are derived from a single male western chimpanzee ('Clint', Coriell identifier S006007), the same individual used to generate the current Chimpanzee genome assembly. We produced ~120-fold sequence coverage of overlapping 250 bps reads (~450 bps fragment) on the Illumina HiSeq 2500 platform, offering high accuracy and throughput, but comparatively short reads; ~9-fold sequence coverage from 43 Pacific Biosciences SMRT-Cells with P5-C3 chemistry on the RSII instrument, offering long reads at lower accuracy; Illumina TruSeq Synthetic long reads at around 2-fold coverage, offering long range information derived from local assemblies of ~10Kb fragments [2]; 1 lane of *in vitro* proximity ligation read pairs (prepared as a Chicago library by Dovetail Genomics)

74 [3] sequenced on the Illumina HiSeq 2000 platform, offering spatial contact

75 information of the chromatin, that can be exploited for scaffolding.

76 These diverse datasets complement the resources that were already available for the

77 same cell line, namely 6-fold coverage of ABI Sanger capillary reads used for the

78 initial chimpanzee genome assembly, a 100 bps paired Illumina HiSeq data, a fosmid

79 library at 6-fold physical coverage with available end sequences, a BAC library at 3-

80 fold physical coverage with available end sequences and around 700 finished BACs

81 [4]. Altogether, these data constitute an extensive non-human, and non-model

82 organism benchmarking dataset for different sequencing strategies.

83

## Assembly generation

85 We generated a complete *de novo* assembly for the chimpanzee with a combination of

86 the datasets. At each step of our assembly we measured increase in contiguity by

87 means of the N50 statistic, which is defined as the length of a contig or scaffold such

88 that 50% of the assembly bases are contained in contigs or scaffolds of at least that

89 length. The starting point of our assembly scaffolding efforts are contigs generated

90 with DISCOVAR *de novo* [5] from 250 bps paired end reads. These reads are derived

91 from a 450 bps library, resulting in pairs that overlap over a ~50 bps region, a feature

92 that is exploited by the assembler. While based on Illumina sequencing, these libraries

93 have recently been shown to produce assemblies superior in contiguity when

94 compared to assemblies derived from conventional Illumina libraries [6]. The

95 DISCOVAR base assembly had a contig N50 of 87 Kbp, and was then scaffolded

96 using proximity ligation read-pairs generated by the Chicago method [3] and

97 sequenced on the Illumina platform. These data increased the scaffold N50 to 26

98 Mbp. Notably, individual scaffolds exceed lengths of 75 Mbp and therefore already

4

99   reach the order of magnitude of full chromosomal arms. We sought to take advantage

100  of this highly contiguous scaffolds and attempt closure of remaining gaps with long-

101  read single molecule sequences by PacBio using PBJelly (PBJelly,

102  RRID:SCR_012091) [7]. By this means, we filled over 38,000 gaps (or 55%) among

103  all scaffolds and in so doing increased the contig N50 by over 320% to 283 Kbp when

104  compared to the DISCOVAR base assembly (see Table 1). While we went on to

105  further improve the assembly with additional data (see below), these statistics give an

106  approximation of the contiguity that can be expected for *de novo* assemblies of

107  previously unsequenced species using our three-way hybrid approach: contigs derived

108  from overlapping 250 bps paired end reads to scaffold with in vitro HiC, and fill

109  remaining gaps with PacBio data. When the contiguity metrics of this intermediate

110  assembly are compared to other representative non-human primate genomes (as

111  annotated by NCBI Refseq category, July 1, 2016; see supplementary material), we

112  observed superior contiguity in contig structure within our assembly compared to all

113  others. The only exception is the gorilla genome, recently assembled from deep (~75-

114  fold) long-read sequences [8]. However, our stepwise method offers an approach that

115  is considerably cheaper.

116

117  **Assembly refinement and comparison to Pan_tro_2.1.4**

118  For the final release of the chimpanzee assembly, we created a reference assembly

119  that leveraged previous resources generated from the same individual  [4]. First, we

120  merged in regions from Pan_tro_2.1.4 that were derived from Clint and gapped in our

121  assembly. It is known that Pan_tro_2.1.4 contains sequences from different

122  chimpanzees. To do so, we extracted flanking sequence regions of gaps in our

123 assembly and mapped all to Pan_tro_2.1.4, keeping only unique and concordant

124 mappings that do not span any gaps within Pan_tro_2.1.4, and merged the spanned

125 Pan_tro_2.1.4 sequence in.

126 To ensure accuracy was not sacrificed for continuity gains we utilized various

127 methods to measure error. Given that our assembly likely contained some erroneous

128 links between contigs or misassembled contigs as a result of *de novo* assembly,

129 conformational mapping or merging mistakes, we first used discordant mapping of

130 fosmid end sequences (~40 Kbp insert size) to identify any large misassemblies. We

131 identified 17 such scaffold errors and manually broke apart each. We also sought to

132 correct any remaining single base substitutions or small indels (<6 bps) with a series

133 of custom mapping and base integration programs (see supplementary material). With

134 the same Illumina data used to generate the DISCOVAR base assembly, we corrected

135 more than 500,000 single base or indel errors. Most of these residual errors are

136 presumably derived from regions where PacBio data was incorporated into the

137 assembly, as this platform is known to have an elevated error rate. As another

138 measure of quality we produced whole genome alignments to Pan_tro_2.1.4 and find

139 our assembly aligns with on average 99.9% identity, and the magnitude of remaining

140 differences can thus reasonably explained by the allelic diversity of western

141 chimpanzees [9].

142 For our final assembly, named Pan_tro_3.0, we integrated previously available

143 finished clone sequences derived from Clint where possible. Pan_tro_3.0 spans 2.95

144 Gbp in ordered and oriented chromosomal sequences. An additional 140 Mbp of

145 sequence is assigned to chromosomes, but their order and orientation unknown, and

146 123 Mbp remain of unknown chromosomal origin. Pan_tro_3.0 has a genome-wide

147 contig and scaffold N50 of 385 Kbp and 27 Mbp, respectively, constituting an

6

148    improvement in contiguity over Pan_tro_2.1.4 of 760% and 300%, respectively (see

149    Figure 1a and Table1). We observed this increase across all non-finished

150    chromosomes, with the most pronounced effect on the X chromosome (see Figure

151    1b). This chromosome shows the highest degree of fragmentation in Pan_tro_2.1.4,

152    likely due to the fact that the effective sequence coverage on the sex chromosomes is

153    only half that of the autosomes, namely around 3-fold in the original assembly. We

154    increased the contig N50 on the X chromosome by 3,250% from 13 Kbp to 422 Kbp,

155    thus bringing its contiguity to the range observed on autosomes.

156    Overall, we decreased the number of contigs by more than 60% from 183,860 to

157    72,226 and the number of gaps by 83% from 156,857 to 26,715. As gap structures

158    between the assemblies may not correspond, we identified filled gaps from

159    Pan_tro_2.1.4 by extracting their flanking regions and mapping them onto

160    Pan_tro_3.0. By keeping only unique and concordant mappings that do not span any

161    gaps in Pan_tro_3.0, we estimate the sequences of 122,943 (77%) gaps to be filled,

162    amounting for 60.3 Mbp of sequence. The majority of these fill sequences are

163    comparably short (see Figure 1C) and significantly enriched in interspersed genomic

164    repeats with 58% of them  (p<0.0001, feature permutation test) into repeats. Of these,

165    around 16 Mbp are fully embedded within fill sequences corresponding to, amongst

166    others, over 29,650 novel short interspersed nuclear elements (SINE) annotations and

167    20,888 novel long interspersed nuclear elements (LINE) annotations.

168

169

**Table 1 - Assembly statistics comparing the previous chimpanzee assembly, our intermediary assembly based on the 3-way hybrid and the finished assembly Pan_tro_3.0. In this context, we defined gaps at stretches of at least 10 consecutive "N" in the assembly. Contigs are defined as contiguous stretches of sequence without gaps.**

| | Pan_tro_2.1.4 | 3-way hybrid (intermediary) | Pan_tro_3.0 |
|---|---|---|---|
| **Scaffold N50 (bps)** | 8,925,874 | 26,681,610 | 26,972,556 |
| **Contig N50 (bps)** | 50,665 | 282,774 | 384,816 |
| **Contig N90 (bps)** | 7,231 | 41,655 | 53,112 |
| **Assembly length (bps)** | 3,309,577,923 | 2,992,696,208 | 3,231,154,112 |
| **Assembly length w\o N's (bps)** | 2,902,338,968 | 2,990,712,612 | 3,132,603,062 |
| **Scaffolds** | 24,129 | 45,000 | 44,448 |
| **Contigs** | 183,827 | 76,674 | 72,226 |
| **Gaps** | 159,698 | 31,674 | 26,715 |

174

### Repeat resolution

Large genomic repeats constitute a major confounding factor in genome assembly and are therefore one of the main reasons for their fragmentation and thus, the assembly repeat representation can be a proxy of its quality. To assess the repeat resolution of interspersed repeats, we masked Pan_tro_3.0 using RepeatMasker (RepeatMasker, RRID:SCR_012954) [10] selecting chimpanzee specific repeats, resulting in 1.64 Gbp (52.2%) being annotated as repeats. The proportion of repetitive elements is similar in Pan_tro_2.1.4 (50.9%), however, given the large amount of newly resolved sequences this translates into a substantial increase in annotated repeats. Specifically, we annotate 164 Mbp of novel repeats in Pan_tro_3.0, comprising around 10% of the

185 whole repeat annotation. We observe this increase consistently across all families of

186 interspersed repeats (see Figure 1D). The increases range as high as 300% for satellite

187 sequences, corresponding to an additional 68.2 Mbp of newly resolved sequence in

188 this category. We also increased the amount of annotated SINE by 27.9 Mbp,

189 including 83,637 additional resolved copies of *Alu* elements. We find the increase in

190 annotations to be negatively correlated with age for *Alu* elements, and thus find the

191 highest increase (8.8%) for the youngest and least divergent subfamily (*AluY*),

192 suggesting that common high identity repeats are now better resolved. We

193 furthermore added 38.2 Mbp of sequence annotated as LINEs to the assembly. We

194 also observed a noteworthy increase in annotated long terminal repeats (LTR), adding

195 15.9 Mbp to this repeat category, corresponding to 30,574 additional annotations of

196 endogenous retroviruses (ERV) in the genome. When comparing all types of

197 interspersed repeats between Pan_tro_2.1.4 and Pan_tro_3.0, we find a median

198 increase of 4.7% of sequence, highlighting that repeat resolution is much improved in

199 Pan_tro_3.0 (see supplementray table S4).

## Representation of segmental duplications

201 To analyze the representation of segmental duplications in Pan_tro_3.0, we applied

202 two alternative approaches: First, we performed a whole genome assembly

203 comparison (WGAC) to compare repeat-free sequences of the assembly to itself [11].

204 This method identifies duplicated sequence in blocks of at least 1 Kbp with 90%

205 identity or higher. Excluding unplaced contigs, we find 140 Mbp of non-redundant

206 duplicated sequence in Pan_tro_3.0 chromosomes, or 4.46% of the non-gap bases in

207 the assembly, results that are consistent with previous read-depth estimates for

208 chimpanzee [12] and analyses of high quality, finished human genome assemblies

209 (see supplementary material S3). Second, we identified duplications by whole-

210 genome shotgun sequence detection (WSSD) that identifies duplications at least 10

211 Kbp long with over 94% identity by detecting regions of increased read depth

212 compared to known unique regions [13]. We used 31,366,275 Sanger capillary reads

213 derived from Clint, and find 51 Mbp of duplicated sequence meeting these criteria on

214 placed chromosomes, compared to 68 Mbp detected by WGAC.

215 Genome wide, we discovered 178,245 redundant pairwise alignments corresponding

216 to 388 Mbp of non-redundant sequence above 1Kbp in length and 90% identity

217 (12.39% of the genome sequence excluding gaps) by WGAC, and 63 Mbp of

218 duplicated sequence by WSSD (compared to 284 Mbp WGAC ≥10 Kbp, >94%

219 identity). We then compared Pant_tro_3.0 to the human reference genome assembly

220 GRCh38, an assembly that is based on a BAC hierarchical shotgun assembly strategy

221 and may therefore be considered of gold standard with respect to representation of

222 segmental duplications. We note similar proportions of bases in segmental

223 duplications on chromosomal scaffolds (4,46% in Pan_tro_3.0 vs. 5,56% in

224 GRCh38), however, we note an elevated genome wide rate of bases in duplications

225 when including unplaced and unlocalized scaffolds. This suggests that our assembly

226 includes false-positive paralogous regions within them (see supplementary Table 1).

227 **Gene annotation**

228 We produced a new gene annotation based on projections from all human transcripts

229 in the GENCODE annotation V24 set combined with RNA-seq data derived from

230 brain, heart, liver and testis from three different individuals [14]. To quantify the

231 effect of the underlying sequence on the annotation, we annotated Pan_tro_2.1.4. with

232 the same data. We observe improvements in gene annotation in Pan_tro_3.0 in all

233 considered metrics: We increased the number of recovered consensus gene models for

234 protein coding transcripts by 2.7%, and are now able to project and annotate 89.5% of

10

235  the GENCODE human coding transcripts onto the new assembly. The average

236  coverage of these transcripts within the genome is 98.9%, a gain of 2%. We also

237  observe an increase of 6.6% in transcripts with multiple mappings. We checked for

238  newly resolved exonic sequences in filled gaps with respect to Pan_tro_2.1.4, and find

239  17,818 exons, amounting to 851 Kbp of non-overlapping sequence to be fully

240  embedded within them. Altogether, we retrieved models for 77,858 coding transcripts

241  corresponding to the isoforms of 20,373 coding genes.

242  We find 5,039 human coding transcripts corresponding to 2,660 genes with predicted

243  frameshift mutations in Pan_tro_2.1.4 to human, but not in Pan_tro_3.0. Conversely,

244  we find 674 genes with predicted frameshift mutations to human that are present in

245  Pan_tro_3, but not in Pan_tro_2.1.4. Given that both assemblies are mainly based on

246  data from the same individual (with the exception of chromosome 21 and around 28%

247  of chromosome 7 in Pan_tro_2.1.4, which where derived from a different individual),

248  the majority of these predictions constitute either allelic variation or putative sequence

249  errors in Pan_tro_2.1.4.

250

251  In summary, we describe a hybrid assembly approach to obtain a more complete de

252  novo chimpanzee reference genome assembly, substantially increasing contiguity

253  metrics within it. Our proposed assembly method should be easily applicable to

254  different organisms of similar genomic architecture.

255

256

257  **Figure 1**

258  A: Genome wide distribution of contig lengths between Pan_tro_2.1.4 and

259  Pan_tro_3.0. The peak for Pan_tro_3.0 is shifted to higher values by an order of

260　magnitude.

261　B: Increase in contig N50 for all chromosomes that were not finished with clones in

262　Pan_tro_2.1.4 or Pan_tro_3.0.

263　C: Length distribution of filled gaps in Pan_tro_3. Negative values constitute wrongly

264　separated overlapping contig ends in Pan_tro_2.1.4.

265　D: Increase in annotated interspersed repeats separated by repeat family.

266

## Declarations

### Abbreviations

269　bps: base pairs, Kbp: kilo base pairs, Mbp: mega base pairs, indel: insetion-deletion,

270　SINE: short interspersed nuclear element, LINE: long interspersed nuclear element,

271　LTR: long terminal repeat, ERV: endogenous retrovirus, WGAC: whole genome

272　assembly comparison, WSSD: whole-genome shotgun sequence detection.

273

### Competing interests

275　REG is co-founder of Dovetail Genomics

276

### Funding

288

## Author's Contributions

290  TMB, WCW and LF conceived the study; LFKK, CT, LWH and REG produced and

291  analyzed the assembly; IF, JA, JGG, TA, BP, AT, HL, JB, RGP, IP, ASA, JHe, PR,

292  DH, AN, and AJS produced, analyzed and interpreted the assembly and annotations;

293  DG, JHu and EEE analyzed segmental duplications; TMB, WCW and LFKK wrote

294  the manuscript with input from all authors.

295

## Acknowledgements

301

## Availability of supporting data

303 Supporting data are available through the GigaDB database (GigaDB,

304 RRID:SCR_004002) [15]. This Whole Genome Shotgun project has been deposited at

305 DDBJ/ENA/GenBank under the accession AACZ00000000. The version described in

306 this paper is version AACZ04000000. The assembly is available at

307 https://www.ncbi.nlm.nih.gov/assembly/GCF_000001515.7 and at the UCSC genome

308 browser under the identifier panTro5. The assembly denominated Pan_tro_2.1.4 in the

309 manuscript refers to Pan_troglodytes-2.1.4 with the RefSeq assembly accession

310 GCF_000001515.6

311

## References

313 1. Goodwin S, McPherson JD, McCombie WR. Coming of age: ten years of next-

314 generation sequencing technologies. Nat Rev Genet [Internet]. Nature Publishing

315 Group; 2016;17:333–51. Available from:

316 https://www.ncbi.nlm.nih.gov/pubmed/27184599

317 2. Kuleshov V, Xie D, Chen R, Pushkarev D, Ma Z, Blauwkamp T, et al. Whole-

318 genome haplotyping using long reads and statistical methods. Nat. Biotechnol.

319 [Internet]. 2014;32:261–6. Available from: http://dx.doi.org/10.1038/nbt.2833

320 3. Putnam NH, Connell BO, Stites JC, Rice BJ, Hartley PD, Sugnet CW, et al.

321 Chromosome-scale shotgun assembly using an in vitro method for long-range linkage.

322 Genome Res. 2016;1–25.

323 4. Mikkelsen TS, , Evan E. Eichler MCZ, Jaffe, David B., Yang S-P, , Wolfgang

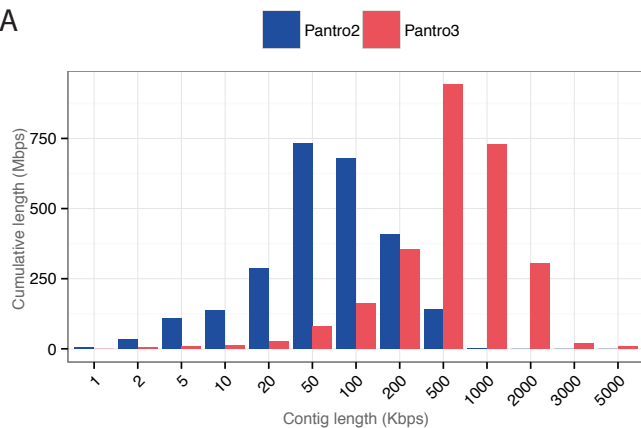324 Enard IH, Bork, Peer, Butler J, Fronick, et al. Initial sequence of the chimpanzee

325    genome and comparison with the human genome. Nature [Internet]. 2005

326    [cite;437:69–87. Available from: http://www.ncbi.nlm.nih.gov/pubmed/16136131

327    5. Weisenfeld NI, Yin S, Sharpe T, Lau B, Hegarty R, Holmes L, et al.

328    Comprehensive variation discovery in single human genomes. Nat. Genet. [Internet].

329    Nature Publishing Group; 2014;46:1350–5. Available from:

330    http://dx.doi.org/10.1038/ng.3121

331    6. Chaisson MJP, Wilson RK, Eichler EE. Genetic variation and the de novo

332    assembly of human genomes. Nat. Rev. Genet. [Internet]. Nature Publishing Group;

333    2015;16:627–40. Available from: http://dx.doi.org/10.1038/nrg3933

334    7. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the gap:

335    upgrading genomes with Pacific Biosciences RS long-read sequencing technology.

336    Liu Z, editor. PLoS One [Internet]. Public Library of Science; 2012 ;7:e47768.

337    Available from: http://dx.plos.org/10.1371/journal.pone.0047768

338    8. Gordon D, Huddleston J, Chaisson MJP, Hill CM, Kronenberg ZN, Munson KM,

339    et al. Long-read sequence assembly of the gorilla genome. 2016;344.

340    9. Prado-Martinez J, Sudmant PH, Kidd JM, Li H, Kelley JL, Lorente-Galdos B, et al.

341    Great ape genetic diversity and population history. Nature . 2013;499:471–5.

342    Available from:

343    http://www.nature.com/nature/journal/v499/n7459/full/nature12228.html#/accessions

344    10. Smit A, Hubley R, Green P. RepeatMasker Open-3.0. RepeatMasker. 1996. p.

345    www.repeatmasker.org.

346    11. Bailey J a, Yavor AM, Massa HF, Trask BJ, Eichler EE. Segmental Duplications :

347    Organization and Impact Within the Current Human Genome Project Assembly

348    Segmental Duplications : Organization and Impact Within the Current Human

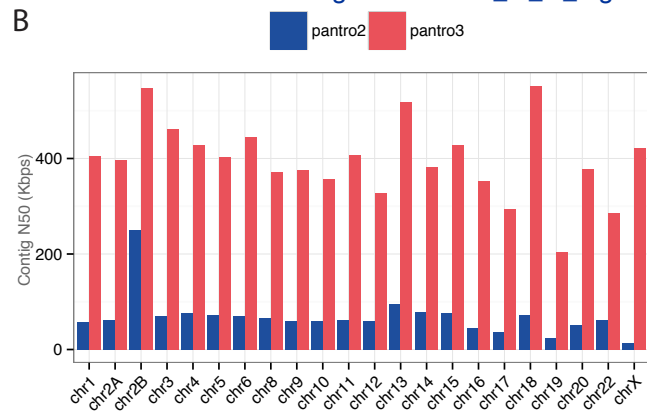349    Genome Project Assembly. Genome Res. 2001;11:1005–17.

15

350    12. Cheng Z, Ventura M, She X, Khaitovich P, Graves T, Osoegawa K, et al. A

351    genome-wide comparison of recent chimpanzee and human segmental duplications.

352    Nature. 2005;437:88–93.

353    13. Bailey J a, Gu Z, Clark R a, Reinert K, Samonte R V, Schwartz S, et al. Recent

354    segmental duplications in the human genome. Science. 2002;297:1003–7.

355    14. Ruiz-Orera J, Hernandez-Rodriguez J, Chiva C, Sabidó E, Kondova I, Bontrop R,

356    et al. Origins of De Novo Genes in Human and Chimpanzee. Noonan J, editor. PLOS

357    Genet. [Internet]. Public Library of Science; 2015 ;11:e1005721. Available from:

358    http://dx.plos.org/10.1371/journal.pgen.1005721

359    15. Kuderna LF, Tomlinson C, Hillier LW, Tran A, Fiddes I, Armstrong J et al. High

360    quality chimpanzee reference genome (Pan_tro_3.0) from hybrid assembly approach.
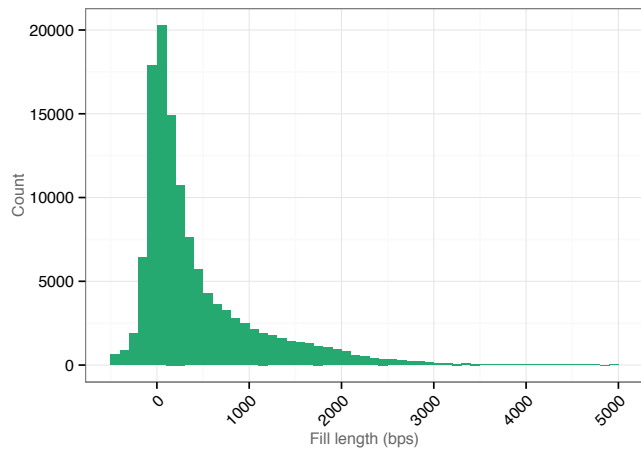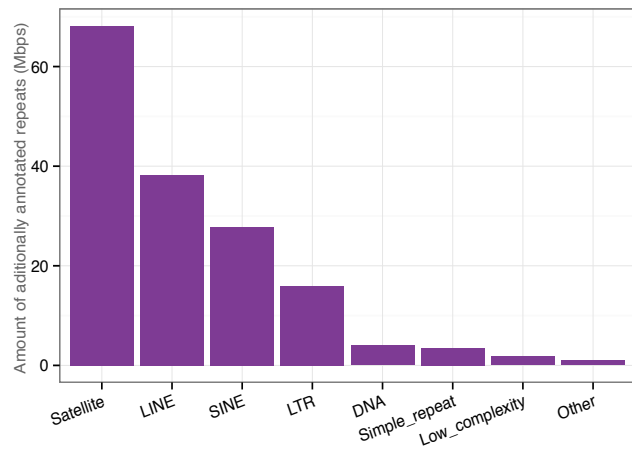
361    *GigaScience* Database. 2017**. http://dx.doi.org/10.5524/100327**

362

363

364
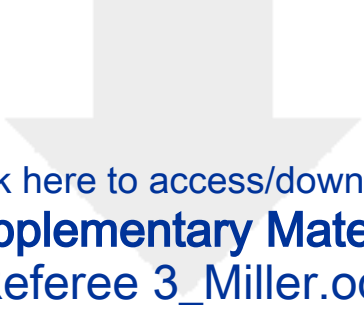
365

16

Figure 1

Point by point response to reviewers

Click here to access/download

**Supplementary Material**

Kuderna_et_al.SUPPLEMENTARY_resubmission.docx

Click here to access/download
**Supplementary Material**
Referee 3_Miller.odt