

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-16-00140	
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:	Ministerio de Economía y Competitividad (BFU2014-55090-P)	Mr Tomas Marques-Bonet
	Stiftelsen för Strategisk Forskning (F06-0045)	Mr Lars Feuk
	National Institutes of Health (DA033660)	Mr Andrew J. Sharp
	National Institutes of Health (HG006696)	Mr Andrew J. Sharp
	National Institutes of Health (HD073731)	Mr Andrew J. Sharp
	National Institutes of Health (MH097018)	Mr Andrew J. Sharp
	March of Dimes Foundation (6-FY13-92)	Mr Andrew J. Sharp
	Ministerio de Economía y Competitividad (BFU2015-7116-ERC)	Mr Tomas Marques-Bonet
	Ministerio de Economía y Competitividad (BFU2015-6215-ERC)	Mr Tomas Marques-Bonet
	National Institutes of Health (HG002385)	Mr Evan E. Eichler
	National Institutes of Health (HG007990)	Mr Benedict Paten
	National Institutes of Health (HG007234)	Mr Benedict Paten
Abstract:	<p>Background: 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 than 183,000 contigs and incorporating over 159,000 gaps, with a genome wide contig N50 of 51 Kbp.</p> <p>Findings: 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.</p> <p>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.</p>	

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:	
Opposed Reviewers:	
Additional Information:	
Question	Response
Are you submitting this manuscript to a	No

special series or article collection?	
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum Standards Reporting Checklist?</p>	Yes

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

3 **Lukas F.K. Kuderna¹, Chad Tomlinson², LaDeana W. Hillier², Annabel**
4 **Tran³, Ian Fiddes⁴, Joel Armstrong⁴, Hafid Laayouni¹, David Gordon⁵, John**
5 **Huddleston⁵, Raquel Garcia Perez¹, Inna Povolotskaya¹, Aitor Serres Armero¹,**
6 **Jèssica Gómez Garrido⁶, Daniel Ho⁷, Paolo Ribeca⁸, Tyler Alioto⁶, Richard E.**
7 **Green^{9,12}, Benedict Paten⁴, Arcadi Navarro^{1,6,10}, Jaume Betranpetit¹, Javier**
8 **Herrero³, Evan E. Eichler⁵, Andrew J. Sharp⁷, Lars Feuk^{11,*}, Wesley C.**
9 **Warren^{2,*}, Tomas Marques-Bonet^{1,6,10*}**

10
11 (1) Institut de Biologia Evolutiva, (CSIC-Universitat Pompeu Fabra), PRBB, Doctor Aiguader 88,
12 Barcelona, Catalonia 08003, Spain.

13 (2) McDonnell Genome Institute, Department of Medicine, Department of Genetics, Washington
14 University School of Medicine, St. Louis, MO 63108, USA.

15 (3) Bill Lyons Informatics Centre, UCL Cancer Institute, University College London, London, UK.

16 (4) Genomics Institute, University of California Santa Cruz and Howard Hughes Medical Institute,
17 Santa Cruz, CA 95064, USA.

18 (5) Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA
19 98195, USA.

20 (6) CNAG-CRG, Centre for Genomic Regulation (CRG), Baldori i Reixac 4, 08028, Barcelona, Spain.

21 (7) Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New
22 York, NY 10029, USA

23 (8) The Pirbright Institute, Ash Road, Pirbright, Woking, GU24 0NF, United Kingdom

24 (9) Department of Biomolecular Engineering, University of California Santa Cruz, 1156 High Street,
25 Santa Cruz, CA 95060, USA.

26 (10) Institutio Catalana de Recerca i Estudis Avançats (ICREA), Barcelona, Catalonia 08010, Spain

27 (11) Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala
28 University, Uppsala, Sweden

29 (12) Dovetail Genomics, Santa Cruz, CA 95060, USA

30 **Abstract**

31 **Background**

32 The chimpanzee is arguably the most important species for the study of human
33 origins. A key resource for these studies is a high quality reference genome assembly,
34 however, as most mammalian genomes, the current iteration of the chimpanzee
35 reference genome assembly it is highly fragmented. In the current iteration of the
36 chimpanzees reference genome assembly (Pan_tro_2.1.4), the sequence is scattered
37 across more than 183,000 contigs and incorporating over 159,000 gaps, with a
38 genome wide contig N50 of 51 Kbp.

39 **Findings**

40 In this work we produce an extensive and diverse array of sequencing datasets to
41 rapidly assemble a new chimpanzee reference that surpasses previous iterations in
42 bases represented and organized in large scaffolds. To this end, we show substantial
43 improvements over the current release of the chimpanzee genome (Pan_tro_2.1.4) by
44 several metrics, such as: increased contiguity by >750% and 300% on contigs and
45 scaffolds, respectively; closure of 77% of gaps in the Pan_tro_2.1.4 assembly gaps
46 spanning >850 Kbp of novel coding sequence based on RNASeq data. We
47 furthermore report over 2,700 genes that had putatively erroneous frame-shift
48 predictions to human in Pan_tro_2.1.4 and show a substantial increase in the
49 annotation of repetitive elements.

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
66
67
68
69
70
71
72
73

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, all 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 (~400 bps fragment) on the Illumina HiSeq 2500 platform; ~9-fold sequence coverage from 43 Pacific Biosciences SMRT-Cells with P5-C3 chemistry on the RSII instrument; Illumina TruSeq Synthetic long reads at around 2-fold coverage; 1 lane of *in vitro* proximity ligation read pairs (prepared as a Chicago library by Dovetail Genomics) sequenced on the Illumina HiSeq 2000 platform.

These diverse datasets complement the resources that were already available for the same cell line, namely 6-fold coverage of ABI Sanger capillary reads used for the initial chimpanzee genome assembly, a 100 bps paired Illumina HiSeq data, a fosmid library at 6-fold physical coverage with available end sequences, a BAC library at 3-

1 74 fold physical coverage with available end sequences and around 700 finished BACs
2 75 [1]. Altogether, these data constitute an extensive and, to our knowledge,
3
4 76 unprecedented non-human, non-model organism benchmarking dataset for different
5
6
7 77 sequencing strategies.
8
9

10 78

11 12 13 79 **Assembly generation** 14

15 80 We generated a complete *de novo* assembly for the chimpanzee with a combination of
16
17 81 the datasets. At each step of our assembly we measured increase in contiguity by
18
19
20 82 means of the N50 statistic, which is defined as the length of a contig or scaffold such
21
22 83 that 50% of the assembly bases are contained in contigs or scaffolds of at least that
23
24
25 84 length. The starting point of our assembly scaffolding efforts are contigs generated
26
27 85 with DISCOVAR *de novo* [2] from 250 bps paired end reads. These reads are derived
28
29 86 from a 400 bps library, resulting in pairs that overlap over a ~50 bps region, a feature
30
31
32 87 that is exploited by the assembler. While based on Illumina sequencing, these libraries
33
34
35 88 have recently been shown to produce assemblies superior in contiguity when
36
37 89 compared to assemblies derived from conventional Illumina libraries [3]. Our base
38
39 90 assembly had a contig N50 of 87 Kbp, and was then scaffolded using proximity
40
41
42 91 ligation read-pairs generated by the Chicago method [4] and sequenced on the
43
44
45 92 Illumina platform. These data increased the N50 to 26 Mbp. Notably, individual
46
47 93 scaffolds exceed lengths of 75 Mbp and therefore already reach the order of
48
49
50 94 magnitude of full chromosomal arms. Despite a substantial gain in scaffold
51
52 95 continuity, remaining gap structure required us to attempt closure with long-read
53
54 96 single molecule sequences by PacBio using PBJelly [5]. By this means, we filled over
55
56
57 97 38,000 gaps (or 55%) among all scaffolds and in so doing increased the contig N50
58
59 98 by over 320% to 283 Kbp when compared to the base assembly (see Table 1). While
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

99 we went on to further improve the assembly with additional data (see below), these
100 statistics give an approximation of the contiguity that can be expected for *de novo*
101 assemblies of previously unsequenced species using our three-way hybrid approach:
102 contigs derived from overlapping 250 bps paired end reads to scaffold with in vitro
103 HiC, and fill remaining gaps with PacBio data. When the contiguity metrics of this
104 intermediate assembly are compared to other representative non-human primate
105 genomes (as annotated by NCBI Refseq category, July 1, 2016; see supplementary
106 material), we observed superior connectivity in contig structure within our assembly
107 compared to all others. The only exception is the gorilla genome, recently assembled
108 from deep (~75-fold) long-read sequences [6]. However, our stepwise method offers
109 an approach that is considerably cheaper.

110

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

112 For the final release of the chimpanzee assembly, we created a reference assembly
113 that leveraged previous resources generated from the same individual [1]. First, we
114 merged in regions from Pan_tro_2.1.4 that were derived from Clint and gapped in our
115 assembly. It is known that Pan_tro_2.1.4 contains sequences from different
116 chimpanzees. To do so, we extracted flanking sequence regions of gaps in our
117 assembly and mapped all to Pan_tro_2.1.4, keeping only unique and concordant
118 mappings that do not span any gaps within Pan_tro_2.1.4, and merged the spanned
119 Pan_tro_2.1.4 sequence in.

120 To ensure accuracy was not sacrificed for continuity gains we utilized various
121 methods to measure error. Given that our assembly likely contained some erroneous
122 links between contigs or misassembled contigs as a result of *de novo* assembly,

123 conformational mapping or merging mistakes, we first used discordant mapping of
124 fosmid end sequences (~40 Kbp insert size) to identify any large misassemblies. We
125 identified 17 such scaffold errors and manually broke apart each. We also sought to
126 correct any remaining single base substitutions or small indels (<6 bps) with a series
127 of custom mapping and base integration programs (see supplementary material). With
128 the same Illumina data used to generate the first draft assembly (DISCOVAR *de*
129 *novo*) we corrected more than 500,000 single base or indel errors. As another measure
130 of quality we produced whole genome alignments to Pan_tro_2.1.4 and find our
131 assembly aligns with on average 99.9% identity, and the magnitude of remaining
132 differences can thus reasonably explained by the allelic diversity of western
133 chimpanzees [7].

134 Our final assembly, named Pan_tro_3.0, spans 2.95 Gbp in ordered and oriented
135 chromosomal sequences. An additional 140 Mbp of sequence is assigned to
136 chromosomes, but their order and orientation unknown, and 123 Mbp remain of
137 unknown chromosomal origin. Pan_tro_3.0 has a genome-wide contig and scaffold
138 N50 of 385 Kbp and 27 Mbp, respectively, constituting an improvement in contiguity
139 over Pan_tro_2.1.4 of 760% and 300%, respectively (see Figure 1a and Table1). We
140 observed this increase across all non-finished chromosomes, with the most
141 pronounced effect on the X chromosome (see Figure 1b). This chromosome shows the
142 highest degree of fragmentation in Pan_tro_2.1.4, likely due to the fact that the
143 effective sequence coverage on the sex chromosomes is only half that of the
144 autosomes, namely around 3-fold in the original assembly. We increased the contig
145 N50 on the X chromosome by 3,250% from 13 Kbp to 422 Kbp, thus bringing its
146 contiguity to the range observed on autosomes.

147 Overall, we decreased the number of contigs by more than 60% from 183,860 to
 148 72,226 and the number of gaps by 83% from 156,857 to 26,715. As gap structures
 149 between the assemblies may not correspond, we identified filled gaps from
 150 Pan_tro_2.1.4 by extracting their flanking regions and mapping them onto
 151 Pan_tro_3.0. By keeping only unique and concordant mappings that do not span any
 152 gaps in Pan_tro_3.0, we estimate the sequences of 122,943 (77%) gaps to be filled,
 153 amounting for 60.3 Mbp of sequence. The majority of these fill sequences are
 154 comparably short (see Figure 1C) and significantly enriched in interspersed genomic
 155 repeats with 58% of them ($p < 0.0001$, feature permutation test) into repeats. Of these,
 156 around 16 Mbp are fully embedded within fill sequences corresponding to, amongst
 157 others, over 29,650 novel short interspersed nuclear elements (SINE) and 20,888
 158 novel long interspersed nuclear elements (LINE) annotations.

160 **Table 1 - Assembly statistics comparing the previous chimpanzee assembly, our intermediary assembly**
 161 **based on the 3-way hybrid and the finished assembly Pan_tro_3.0. In this context, we defined gaps at**
 162 **stretches of at least 10 consecutive “N” in the assembly. Contigs are defined as contiguous stretches of**
 163 **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

164

165 **Repeat resolution**

166 Large genomic repeats constitute a major confounding factor in genome assembly and
167 are therefore one of the main reasons for their fragmentation and thus, the assembly
168 repeat representation can be a proxy of its quality. To assess the repeat resolution of
169 interspersed repeats, we masked Pan_tro_3.0 using RepeatMasker [8] selecting
170 chimpanzee specific repeats, resulting in 1.64 Gbp (52.2%) being annotated as
171 repeats. The proportion of repetitive elements is similar in Pan_tro_2.1.4 (50.9%),
172 however, given the large amount of newly resolved sequences this translates into a
173 substantial increase in annotated repeats. Specifically, we annotate 164 Mbp of novel
174 repeats in Pan_tro_3.0, comprising around 10% of the whole repeat annotation. We
175 observe this increase consistently across all families of interspersed repeats (see
176 Figure 1D). The increases range as high as 300% for satellite sequences,
177 corresponding to an additional 68.2 Mbp of newly resolved sequence in this category.
178 We also increased the amount of annotated SINE by 27.9 Mbp, including 83,637
179 additional resolved copies of *Alu* elements. We find the increase in annotations to be
180 negatively correlated with age for *Alu* elements, and thus find the highest increase
181 (8.8%) for the youngest and least divergent subfamily (*AluY*), suggesting that
182 common high identity repeats are now better resolved. We furthermore added 38.2
183 Mbp of LINE to the assembly, corresponding to over 44,791 additional copies of L1
184 elements. We also observed a noteworthy increase in annotated long terminal repeats
185 (LTR), adding 15.9 Mbp to this repeat category, corresponding to 30,574 additional
186 annotated copies of endogenous retroviruses (ERV) in the genome. When comparing

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

187 all types of interspersed repeats between Pan_tro_2.1.4 and Pan_tro_3.0, we find a
188 median increase of 4.7% of sequence, highlighting that repeat resolution is much
189 improved in Pan_tro_3.0 (see supplementray table S4).

190 **Representation of segmental duplications**

191 To analyze the representation of segmental duplications in Pan_tro_3.0, we applied
192 two alternative approaches: First, we performed a whole genome assembly
193 comparison (WGAC) to compare repeat-free sequences of the assembly to itself. This
194 method identifies duplicated sequence in blocks of at least 1 Kbp with 90% identity or
195 higher. Excluding unplaced contigs, we find 140 Mbp of non-redundant duplicated
196 sequence in Pan_tro_3.0 chromosomes, or 4.46% of the non-gap bases in the
197 assembly, results that are consistent with previous read-depth estimates for
198 chimpanzee [9] and analyses of high quality, finished human genome assemblies (see
199 supplementary material S3). Second, we identified duplications by whole-genome
200 shotgun sequence detection (WSSD) that identifies duplications at least 10 Kbp long
201 with over 94% identity by detecting regions of increased read depth compared to
202 known unique regions. We used 31,366,275 Sanger capillary reads derived from
203 Clint, and find 51 Mbp of duplicated sequence meeting these criteria on placed
204 chromosomes, compared to 68 Mbp detected by WGAC.

205 Genome wide, we discovered 178,245 redundant pairwise alignments corresponding
206 to 388 Mbp of non-redundant sequence above 1Kbp in length and 90% identity
207 (12.39% of the genome sequence excluding gaps) by WGAC, and 63 Mbp of
208 duplicated sequence by WSSD (compared to 284 Mbp WGAC ≥ 10 Kbp, $>94\%$
209 identity). Altogether, these results suggest that segmental duplications are well
210 resolved in Pan_tro_3.0 on the chromosomal level, however, we are likely to be

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

211 overestimating the total amount of segmental duplications genome wide by including
212 an elevated rate of false-positive paralogous regions localized on unplaced scaffolds.

213 **Gene annotation**

214 We produced a new gene annotation based on projections from all human transcripts
215 in the GENCODE annotation V24 set combined with RNA-seq data derived from
216 brain, heart, liver and testis from three different individuals [10]. To quantify the
217 effect of the underlying sequence on the annotation, we annotated Pan_tro_2.1.4. with
218 the same data. We observe improvements in gene annotation in Pan_tro_3.0 in all
219 considered metrics: We increased the number of recovered consensus gene models for
220 protein coding transcripts by 2.7%, and are now able to project and annotate 89.5% of
221 the GENCODE human coding transcripts onto the new assembly. The average
222 coverage of these transcripts within the genome is 98.9%, a gain of 2%. We observe
223 an increase of 6.6% in transcripts with multiple mappings, suggesting that paralogous
224 coding duplications are better represented in this assembly. We checked for newly
225 resolved exonic sequences in filled gaps with respect to Pan_tro_2.1.4, and find
226 17,818 exons, amounting to 851 Kbp of non-overlapping sequence to be fully
227 embedded within them. Altogether, we retrieved models for 77,858 coding transcripts
228 corresponding to the isoforms of 20,373 coding genes.

229 Perhaps most strikingly, we found 5,039 human coding transcripts corresponding to
230 2,728 genes with predicted frameshift mutations in Pan_tro_2.1.4 to human, but not in
231 Pan_tro_3.0. Given that both assemblies are mainly based on data from the same
232 individual, the majority of these predictions constitute putative sequence errors in
233 Pan_tro_2.1.4.

234

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

235 In summary, we describe a hybrid assembly approach to obtain a more complete de
236 novo chimpanzee reference genome assembly, substantially increasing contiguity
237 metrics within it.

238

239

240 **Figure 1**

241 A: Genome wide distribution of contig lengths between Pan_tro_2.1.4 and
242 Pan_tro_3.0. The peak for Pan_tro_3.0 is shifted to higher values by an order of
243 magnitude.

244 B: Increase in contig N50 for all chromosomes that were not finished with clones in
245 Pan_tro_2.1.4 or Pan_tro_3.0.

246 C: Length distribution of filled gaps in Pan_tro_3. Negative values constitute wrongly
247 separated overlapping contig ends in Pan_tro_2.1.4.

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

249

250 **Declarations**

251 **Abbreviations**

252 bps: base pairs, Kbp: kilo base pairs, Mbp: mega base pairs, indel: insertion-deletion,
253 SINE: short interspersed nuclear element, LINE: long interspersed nuclear element,
254 LTR: long terminal repeat, ERV: endogenous retrovirus, WGAC: whole genome
255 assembly comparison, WSSD: whole-genome shotgun sequence detection.

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

256

257 **Competing interests**

258 REG is co-founder of Dovetail Genomics

259

260 **Funding**

261 LFKK is supported by an FPI fellowship associated to BFU2014-55090-P (FEDER);

262 LF is supported by the Swedish Foundation for Strategic Research F06-0045 and the

263 Swedish Research Council; EEE is an investigator of the Howard Hughes Medical

264 Institute. AJS is supported by NIH grants DA033660, HG006696, HD073731 and

265 MH097018, and research grant 6-FY13-92 from the March of Dimes TMB is

266 supported by MINECO BFU2014-55090-P (FEDER), BFU2015-7116-ERC and

267 BFU2015-6215-ERC, Fundacio Zoo Barcelona and Secretaria d'Universitats i

268 Recerca del Departament d'Economia i Coneixement de la Generalitat de Catalunya.

269 This work was supported, in part, by grants from the U.S. National Institutes of

270 Health (NIH grants HG002385 to E.E.E., HG007990 and HG007234 to B.P.).

271

272 **Author's Contributions**

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

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

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

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

277 the manuscript with input from all authors.

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

278

279 **Acknowledgements**

280 We would like to acknowledge Bojan Obradovic and James Richardson for sequence
281 contribution. The authors acknowledge the use of the UCL Legion High Performance
282 Computing Facility ([Legion@UCL](#)), and associated support services, in the
283 completion of this work.

284

285 **Availability of supporting data**

286 Supporting data area available through the GigaDB database. This Whole Genome
287 Shotgun project has been deposited at DDBJ/ENA/GenBank under the accession
288 AACZ000000000. The version described in this paper is version AACZ040000000. The
289 assembly is available at https://www.ncbi.nlm.nih.gov/assembly/GCF_000001515.7
290 and at the UCSC genome browser under the identifier panTro5.

291

292 **References**

- 293 1. Mikkelsen TS, , Evan E. Eichler MCZ, Jaffe, David B., Yang S-P, , Wolfgang
294 Enard IH, Bork, Peer, Butler J, Fronick, et al. Initial sequence of the chimpanzee
295 genome and comparison with the human genome. *Nature*. 2005;437:69–87.
- 296 2. Weisenfeld NI, Yin S, Sharpe T, Lau B, Hegarty R, Holmes L, et al.
297 Comprehensive variation discovery in single human genomes. *Nat. Genet.*;
298 2014;46:1350–5.
- 299 3. Chaisson MJP, Wilson RK, Eichler EE. Genetic variation and the de novo
300 assembly of human genomes. *Nat. Rev. Genet.*; 2015;16:627–40.

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

301 4. Putnam NH, Connell BO, Stites JC, Rice BJ, Hartley PD, Sugnet CW, et al.
302 Chromosome-scale shotgun assembly using an in vitro method for long-range linkage.
303 Genome Res. 2016;1–25.

304 5. English AC, Richards S, Han Y, Wang M, Vee V, Qu J, et al. Mind the gap:
305 upgrading genomes with Pacific Biosciences RS long-read sequencing technology.
306 Liu Z, editor. PLoS One; 2012

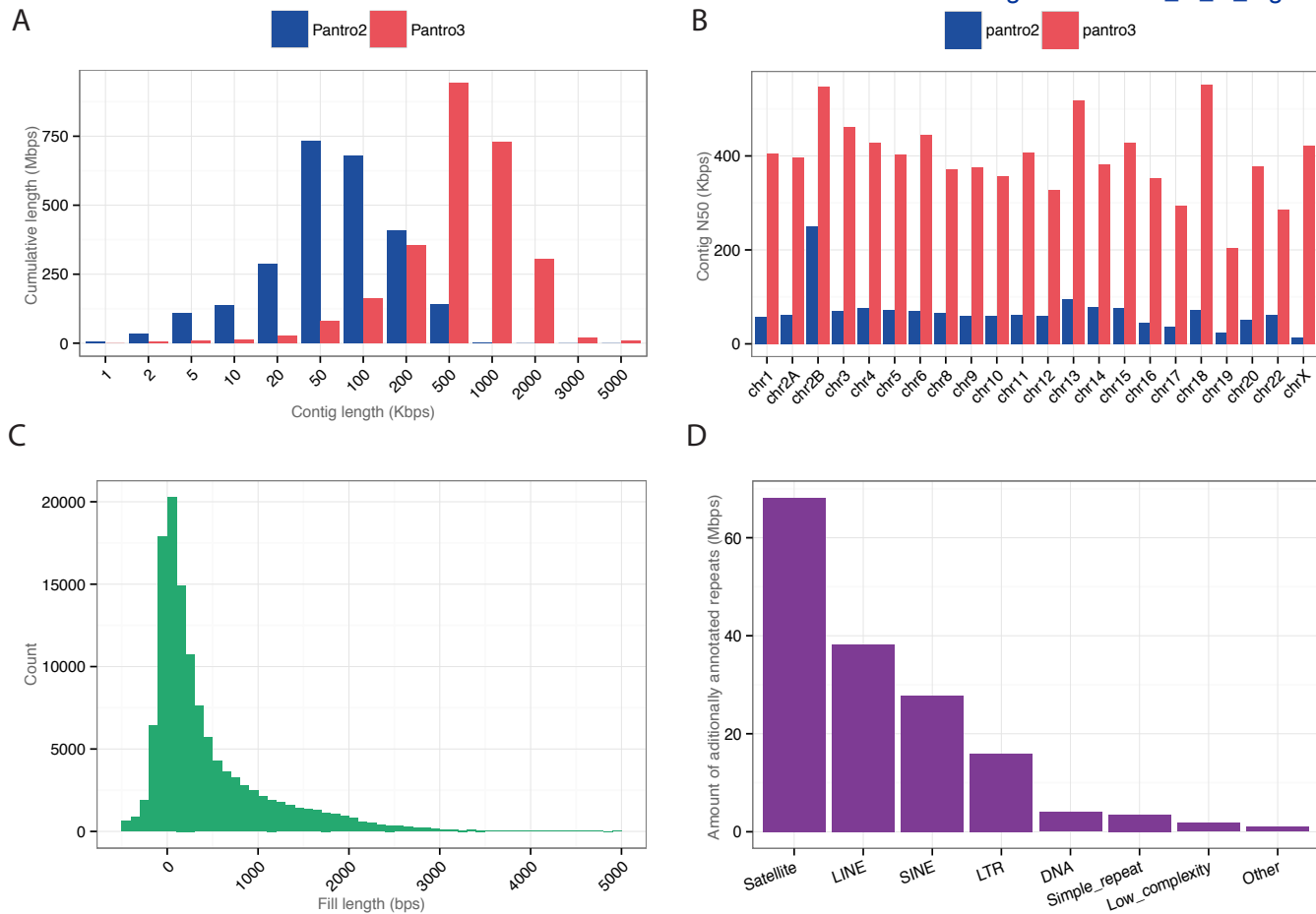
307 6. Gordon D, Huddleston J, Chaisson MJP, Hill CM, Kronenberg ZN, Munson KM,
308 et al. Long-read sequence assembly of the gorilla genome. Science;2016;0344.

309 7. Prado-Martinez J, Sudmant PH, Kidd JM, Li H, Kelley JL, Lorente-Galdos B, et al.
310 Great ape genetic diversity and population history. Nature. 2013;499:471–5.

311 8. Smit A, Hubley R, Green P. RepeatMasker Open-3.0. RepeatMasker. 1996. p.
312 www.repeatmasker.org.

313 9. Cheng Z, Ventura M, She X, Khaitovich P, Graves T, Osoegawa K, et al. A
314 genome-wide comparison of recent chimpanzee and human segmental duplications.
315 Nature. 2005;437:88–93.

316 10. Ruiz-Orera J, Hernandez-Rodriguez J, Chiva C, Sabidó E, Kondova I, Bontrop R,
317 et al. Origins of De Novo Genes in Human and Chimpanzee. Noonan J, editor. PLOS
318 Genet.;2015;11:e1005721.
319

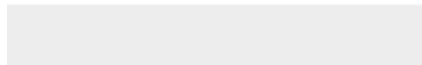




[Click here to access/download](#)

Supplementary Material

Kuderna_et_al.SUPPLEMENTARY.docx





Click here to access/download
Supplementary Material
DATA_SUBMISSION_FILE

