

1 **Chromosomal assembly of the nuclear genome of the endosymbiont-bearing**

2 **trypanosomatid *Angomonas deanei***

3 John W. Davey, Carolina M. C. Catta-Preta, Sally James, Sarah Forrester, Maria

4 Cristina M. Motta, Peter D. Ashton, Jeremy C. Mottram

5 **Supplemental File Legends**

6 **Supplemental Methods**

7 1 Genome Assembly Edits

8 1.1 Symbiont

9 1.2 Kinetoplast DNA minicircle

10 1.3 Kinetoplast DNA maxicircle

11 1.4 Translocation

12 1.5 Inversion

13 1.6 Palindromic misassembly

14 1.7 Incomplete chromosome tig00306615

15 1.8 Incomplete chromosome tig00003599

16 1.9 Telomere edits

17 2 Validation of translocation and inversion

18 2.1 Validation with read alignments

19 2.2 Validation with PCR

20 3 Genome annotation

21 3.1 Transfer original annotations with flo

22 3.2 Filter duplicate annotations and fix sequence errors

23 3.3 Companion run

24 **SUPPLEMENTAL FILES**

25 **File S1** Supplemental Information

26 This file, describing the Supplemental Files S1 to S13 and containing Supplemental
27 Methods for the genome editing and annotation.

28 **File S2** Raw genome assembly

29 The raw genome assembly generated by Canu, in FASTA format, containing 212
30 contigs, 27 914 394 bp long.

31 **File S3** Tapestry report for the raw genome assembly

32 Summary statistics for all contigs in the raw genome assembly including read and
33 contig alignments, generated by Tapestry (Davey et al. 2020). This is a HTML file
34 which should open in any modern web browser.

35 **File S4** Tapestry contig order file

36 A table of raw genome contigs with annotations, generated manually using File S3,
37 in Comma-Separated Values (CSV) format. This file can be viewed by loading File
38 S3 in a web browser, clicking the 'Choose File...' button at the top, and choosing File
39 S4 to load. Alternatively it can be loaded into any program that parses CSV files (R,
40 Excel etc). This file matches the contig order listed in Table S1.

41 **File S5** Polished genome assembly

42 The polished genome assembly after editing and polishing with Oxford Nanopore
43 and Illumina reads, in FASTA format, containing 31 contigs, 21 826 979 bp long (29
44 chromosomes, symbiont and maxicircle).

45 **File S6** fix_annotation_errors.py

46 Python script to select annotations and correct errors in the nanopore assembly
47 sequence following transfer of the gene annotation with flo. See Supplemental
48 Methods below for full details. Arguments:

49 **-r, --reference_fasta**: reference genome in FASTA format

50 **-a, --assembly_fasta**: new genome assembly in FASTA format

51 **-g, --reference_gff**: reference annotation in GFF format

52 **-t, --transferred_gff**: transferred annotation in GFF format

53 **-o, --outputstub**: prefix for output files, default 'fixed'

54 **-w, --overwriteds**: overwrite gffutils databases if they already exist, default False

55 **-p, --protein_attribute**: name of the GFF attribute containing protein name, default
56 'product'

57 **File S7** Annotation transfer details

58 Output of File S6 in Tab-Separated Values (TSV) format describing the original CDS
59 features from the GCA_000442575.2 gene annotation and their destinations in the
60 new nuclear genome. Fields:

61 **ID**: ID attribute from CDS feature from reference GFF

62 **RefContig**: contig name from reference GFF

63 **RefStart, RefEnd:** start and end positions of CDS feature from reference GFF

64 **RefLength:** length of CDS feature in reference GFF

65 **RefProteinName:** name of protein from GFF attribute given as --protein_attribute

66 argument to fix_annotation_errors.py ('product' for GCA_000442575.2)

67 **RefNs:** number of Ns in reference DNA sequence

68 **RefStatus:** assessment of quality of reference protein (one or more of OK, Ns,

69 BadStart, BadStop, BadLength; see Supplemental Methods below for definitions)

70 **NewContig, NewStart, NewEnd:** contig, start and end positions in new genome to

71 which feature has been transferred

72 **NewLength:** length of feature in the new genome

73 **NewStrand:** orientation of feature in the new genome (+/- for forward/reverse)

74 **DNADiff:** difference in length in basepairs between the transferred and reference

75 feature DNA sequences (positive means longer in new genome)

76 **DNAScore:** pairwise alignment score of reference and transferred feature DNA

77 sequences, divided by the length of the transferred DNA sequence (NewLength)

78 **DNAProp:** the proportion of the difference in DNA sequence length (DNADiff)

79 compared to length of the new feature (NewLength)

80 **ProteinDiff:** difference in length in amino acids between the transferred and

81 reference feature protein sequences (positive means longer in new genome). 0 if

82 either protein is not well-formed.

83 **ProteinScore:** pairwise alignment score of reference and transferred feature protein

84 sequences, divided by the length of the transferred protein sequence. '-' if either

85 protein is not well-formed.

86 **NewStatus**: assessment of quality of transferred protein (one or more of OK,
87 Changed, NewLength, BadLength, BadStart, BadStop, ExtraStop; see Supplemental
88 Methods below for definitions).

89 **GroupBegin**: features are grouped together if they overlap. This is the left-most
90 position of the features in the current feature's group.

91 **GroupEnd**: The right-most position of the features in the current feature's group.

92 **GroupFeatures**: number of features in the current feature's group.

93 **GroupFeatureNames**: list of feature names in the current feature's group.

94 **FeatureStatus**: decision on the current feature based on the other features in the
95 group. A feature can be Chosen or Reject. Chosen features can be Accept
96 (sequence is fine, accept as is) or Replace (new sequence is bad, replace with
97 reference sequence). Reject features may be ignored because another feature is
98 higher quality (Prefer), because both the transferred and reference features are bad
99 and so cannot be fixed (BadRef), or because the reference and transferred features
100 differ in length by more than 10% (LenDiff).

101 **File S8** Companion weight function

102 Lua script passed to Companion as the WEIGHT_FILE option and based on the
103 default Companion weight__kinetoplastid.lua function. See Supplemental Methods.

104 **File S9** Annotated genome assembly

105 The final nuclear genome assembly with 29 chromosomes, after genome editing,
106 polishing and fixing of gene sequences during annotation transfer, in FASTA format,
107 containing 20 976 081 bp.

108 **File S10** Companion GFF3 annotation

109 Full gene annotation output by Companion in GFF3 format.

110 **File S11** transfer_gff3_info_to_embl.py

111 Python script transfer additional attributes from original reference genome to

112 Companion's EMBL file. Arguments:

113 **-e, emblgz**: gzipped EMBL file containing full genome and annotation, from

114 Companion output

115 **-g, gff**: Companion output GFF file

116 **-r, refgff**: reference GFF file

117 **-o, output**: name for output gzipped EMBL file

118 **File S12** Final assembly and annotation

119 Assembly and annotation submitted to the European Nucleotide Archive in EMBL

120 format.

121 **File S13** mosdepth_genome_redundancy.py

122 Python script to assess redundancy of genome assemblies, which takes a single

123 argument, **-m (--mosdepth)**, a gzipped BED file output by mosdepth of per-base

124 contig depths.

125 SUPPLEMENTAL METHODS

126 **1 Genome assembly edits**

127 The raw Canu assembly of 212 contigs (File S2) was manually filtered and edited to
128 produce an unpolished, close-to-complete genome assembly, based on the Tapestry
129 report for the raw assembly (File S3, File S4, Table S1) and minimap2 alignments of
130 the nanopore reads to the raw assembly and the raw assembly to itself.

131 Based on the Tapestry report, the 212 raw contigs were placed into 34 groups,
132 representing 29 chromosomes, the kinetoplast maxicircle, the symbiont, the
133 kinetoplast minicircle, a group of subtelomeric contigs, and a repeat contig (File S3,
134 File S4, Table S1). Further description of these groups and their special features
135 follow.

136 *1.1 Symbiont*

137 One contig, tig00000015, 915 769 bp long, had GC content 31.12% (as opposed to
138 the 47-52% GC contents of most long contigs in the assembly), no alignments to any
139 other contig, and one major self-alignment in forward orientation (1-96688 bp aligned
140 to 819069-915768 bp), indicating a circular contig. This contig was retained as the
141 symbiont genome, with the self-alignment removed to leave a raw 819 068 bp contig,
142 which was 821 860 bp long after polishing. This polished contig aligns in full to both
143 reference symbiont genomes (GenBank GCF_000319225.1 and GCF_000340825.1)
144 with >99.98% identity.

145 *1.2 Kinetoplast DNA minicircle*

146 127 contigs were short (between 1 094 bp and 45 325 bp), had GC content between
147 39.03% and 46.42%, and had many contig alignments between each other, but very
148 few alignments to other, longer contigs. These were assumed to be copies of the
149 kinetoplast DNA minicircle (Teixeiria et al. 2011), which typically occurs in thousands
150 of variable copies per kDNA network (Lukeš et al. 2002). Given the complexity of
151 these highly repetitive sequences, they were removed from the final assembly
152 without polishing; however, they are available in the raw assembly.

153 *1.3 Kinetoplast DNA maxicircle*

154 Three contigs, tig00000001 (57 346 bp, 31.86% GC), tig00000002 (30 918 bp,
155 31.90% GC), and tig00000005 (30 370 bp, 31.94% GC), had very similar GC
156 contents, clear alignments to each other, and no alignments to any other contigs.
157 The alignments (shown in Figure S1) show that both tig00000002 (red) and
158 tig00000005 (blue) had two full alignments to tig00000001 (black; arrows show
159 alignments wrapping from the end to the start of tig00000001), with both having
160 small overlaps which also align to themselves (for example tig00000002
161 29666-30918 aligns to tig00000002 1-1240). Based on these alignments, bases
162 1-29665 of tig00000002 were selected to represent one copy of the kinetoplast DNA
163 maxicircle genome. After polishing, this sequence was 29 845 bp long (File S5). The
164 reference kinetoplast maxicircle genome (GenBank KJ778684.1) has a full length
165 alignment to this polished sequence with >99.99% identity.

166 *1.4 Translocation*

167 Contigs tig00000126 (521 443 bp, red in Figure S2) and tig00000177 (219 070 bp,
168 blue in Figure S2) have telomeres at their ends but not at their starts. But their starts
169 align to the regions either side of 195233-195236 on tig00000104 (692 209 bp, black
170 in Figure S2), a contig with telomeres at both ends (Figure S2). This is consistent
171 with tig00000126 and tig00000177 being two chromosome arms of a chromosome
172 556 832 bp long that translocates with the two arms of tig00000104 (Haplotypes 1
173 and 2 in Figure S2).

174 If the chromosomes are translocating, there should be evidence of two further
175 haplotypes. Haplotype 3, consisting of the left arm of tig00000104 and the right arm
176 of tig00000126, is supported by tig00000126 containing most of the left arm of
177 tig00000104 (tig00000126:9748-109206 aligns to tig00000104:96043-195233).
178 Haplotype 4, consisting of the left arm of tig00000177 and the right arm of
179 tig00000104, is supported by tig00000417 (118 437 bp, orange in Figure S2), which
180 contains the regions of these arms that span the translocation breakpoint at
181 tig00000104:195233-195236 (tig00000417:58567-118415 aligns to
182 tig00000104:195236-258647, shown wrapping around the diagram by orange
183 arrows, Figure S2).

184 In the genome assembly, Haplotype 1 is represented by tig00000104, which is now
185 chr13 (polished length 698 360 bp, annotated length 698 408 bp). Haplotype 2 was
186 constructed by reversing tig00000177 and adding tig00000126:183682-521443,

187 making an additional chromosome 556 832 bp long, chr18, which was 561 060 bp
188 long after polishing and 561 137 bp after annotation.

189 If the Haplotype 2 edit has been made correctly, reads should align across the join.
190 Figure S3 shows reads aligned to the (unpolished) joined contig
191 tig00000177_tig00000126, with the join highlighted at 219070-219071 bp (red block
192 below base position axis). Although a number of SNPs and indels remain (as
193 expected in an unpolished genome), there are over 400 reads spanning this region,
194 supporting the accuracy of the edit (also, see 'Validation of translocation and
195 inversion with read alignments' and 'Validation of manual joins with PCR' below).

196 *1.5 Inversion*

197 Contig tig00000018 (1 076 494 bp) has a telomere at its end but not at its start. The
198 first 66.8 kb of this contig is also found, reversed, at 405404-472857 bp (Figure S4,
199 pink blocks). The region between these 67 kb sequences, between positions 67 kb
200 and 405 kb, has alignments to two other contigs; tig00003597 (222 300 bp) and
201 tig00000065 (104 319 bp).

202 The first 1-64742 bp of tig00003597 aligns to 340532-405402 bp of tig00000018 (the
203 tig000003597 region with red leftward arrows in Figure S4 aligning to the end of the
204 region with blue rightward arrows in tig00000018). But the remainder of tig00003597
205 (64743-222300 bp) is a different sequence ending with a telomere.

206 The first 1-50966 bp of tig00000065 aligns to tig00000018 66759-117661 bp (the
207 tig00000065 region with orange rightward arrows in Figure S4 aligning to the start of
208 the region with blue rightward arrows in tig00000018). But the remainder of
209 tig00000065 (50974-104319 bp) aligns to tig00003597 (65364-118716 bp).

210 These alignments suggest that the 67-405 kb region in tig00000018, ~338 kb long, is
211 an inversion, with both haplotypes present in the raw reads (Haplotype 1 and
212 Haplotype 2 in Figure S4). Labelling the four breakpoints from these haplotypes A, B,
213 C and D (see Figure S4), tig00000018 contains breakpoint D, but it also contains
214 breakpoint B; the assembler has confused the two haplotypes, assembled two
215 copies of the sequence at 405-472kb in tig00000018, and has then extended no
216 further into the unique material of tig00000018 upstream of 472 kb.

217 Breakpoint A is found in tig00003597, and breakpoint C in tig00000065, supporting
218 the presence of both haplotypes in the genome, as all four expected breakpoints are
219 present in the raw assembly (File S2). For further validation, see 'Translocation and
220 inversion validation' below.

221 Haplotype 2 has been included in the genome assembly, by taking
222 tig00003597:65364-222300 (reversed), then a short connecting region from
223 tig00000065 (50967-50973), also reversed, then tig00000018:66579-1076494. This
224 produced a chromosomal sequence 1 166 680 bp long, chr05, which was 1 174 890
225 bp long after polishing and 1 174 864 bp long after annotation.

226 1.6 Palindromic misassembly

227 Contig tig00000095 (569 734 bp) has a telomere at its end but not at its start. It has
228 a palindromic alignment to itself; the first 108 978 bases of the contig aligns to itself
229 in reverse orientation. There is a break in coverage at 54 226 bp, with no reads
230 spanning this position, and with telomeric sequence beginning from 54 226 onwards
231 (Figure S5). There are also few reads aligning to the first 54 kb of the contig
232 (Tapestry report, File S3). Therefore the contig has been cut at 54 226 bp, making a
233 chromosome with two telomeres 515 509 bp long, named chr22; this was 519 680
234 bp long after polishing and 519 842 bp long after annotation.

235 1.7 Incomplete chromosome tig00306615

236 Contig tig00306615 (1 178 086 bp) has a telomere at its end but not at its start
237 (Figure S6). Bases 3-116067 of this contig align to bases 418613-534780 of contig
238 tig000003593 and to no other contig (Figure S6, File S3). Read alignments at this
239 region show only one read spanning the breakpoint at 116 067 bp with many
240 mismatches, despite good alignments to the surrounding areas (Figure S7). Also,
241 tig00000050 3-127103 (reversed) aligns just beyond this region, to
242 tig00306615:116132-243147 (Figure S6). As tig00000050 contains a telomere at its
243 end, and the alignment to tig000003593 appears to be an assembly error, the region
244 of tig00306615 aligning to tig000003593 was discarded, and a chromosome was
245 constructed from tig00000050:3-231117 (reversed) and
246 tig00306615:243148-1178086, making a sequence 1 166 054 bp long. Over 430
247 reads align cleanly across the join between tig00000050 and tig00306615 (Figure

248 S8), indicating that this join is accurate. This sequence was 1 174 919 bp long after
249 polishing, 1 175 096 bp long after annotation, and is now chr04.

250 *1.8 Incomplete chromosome tig00003599*

251 Contig tig000003599 (988 284 bp long) features a telomere at its start but not at its
252 end. It has two haplotypes that align full length to its end, tig00000047 (71 900 bp)
253 and tig00003600 (73 365 bp) (Figure S9); tig00000047 has a telomere. As these
254 contigs do not have major alignments anywhere else in the genome, they are likely
255 to reflect some structural variation at this chromosome end. In order to complete the
256 chromosome, tig00003599 was truncated up to and including 920 524 bp and
257 tig00000047:4-71900 was added, making a chromosome 992 421 bp long. Around
258 480 reads align cleanly across the join between tig00003599 and tig00000047
259 (Figure S10), indicating that the join is accurate. This chromosome is now chr07,
260 which was 999 268 bp long after polishing and 999 236 bp long after annotation.

261 *1.9 Telomere edits*

262 Five contig ends did not end with telomeric sequence. On inspection, three of these
263 contained telomeres upstream of a misassembled minicircle sequence, and the other
264 two had reads that aligned to the contig end and which contained telomere sequence
265 beyond the end of the contig.

266 The start and end of tig00000058 (767 463 bp) and the end of tig00003608 (422 011
267 bp) contain telomeres, but also have minicircle sequence beyond the telomere
268 sequence. The raw Tapestry report (File S3) shows these contigs aligning to

269 minicircle contigs (click on a contig name in the report diagram to show contig
270 alignments for that contig).

271 Minimap2 alignments of the first 1kb of tig00000058 showed 135 alignments to
272 minicircle contigs. Also, minimap2 alignments to the last 697 bp of tig00000058
273 showed 113 alignments to minicircle contigs. These minicircle sequences were
274 removed by editing the contig to bases 1220-766765, leaving a 765 546 bp contig
275 with a telomere at both ends of the sequence. This is now chr11, which was 770 936
276 bp long after polishing and 771 229 bp long after annotation.

277 Similarly, minimap2 alignments to the last 1kb of tig00003608 showed 17 alignments
278 to minicircle contigs. The 422 011 bp contig was edited to bases 1-420743, leaving a
279 420 743 bp contig now ending with a telomere. This is now chr25, which was 424
280 872 bp long after polishing and 424 834 bp long after annotation.

281 tig00000070 (852 128 bp) has two copies of the telomere sequence TTAGGG at its
282 end. However, inspection of soft-clipped regions of reads beyond the end of
283 tig00000070 shows many reads featuring long TTAGGG telomere sequences
284 (Figures S11 and S12). As there are some soft-clipped reads that appear to have
285 telomeric sequence immediately following the contig, and some that have
286 non-telomeric sequence, it may be that some sequence variation in this region has
287 prevented the assembler from completing the telomere. However, there is no
288 evidence for any other continuation of this contig, and so we can assume the contig

289 is almost a complete chromosome. This is now chr08, 859 818 bp long after
290 polishing and 859 978 bp long after annotation.

291 Similarly, tig00000134 (525 903 bp) has no telomeric sequence at its start, but
292 soft-clipped reads aligning to this region contain long telomeric sequence (Figure
293 S13). However, there again appears to be sequence variation in these reads which
294 perhaps has prevented the assembler from completing the telomere. As with
295 tig00000070, there is no evidence for this being anything other than a complete
296 chromosome. It is now chr20, 530 488 bp long after polishing and 530 564 bp long
297 after annotation.

298 **2 Validation of translocation and inversion**

299 *2.1 Validation with read alignments*

300 Only one inversion haplotype and two translocation haplotypes are included in the
301 genome assembly, as all unique material is contained in these sequences. However,
302 to demonstrate the existence of both inversion haplotypes and all four translocation
303 haplotypes, six contigs were constructed containing these haplotypes (Table S2), all
304 raw reads were aligned against them and the breakpoints and joins were examined
305 (Table S2, Figures S14-S21). Reads aligned across all breakpoints and joins and
306 throughout all contigs, confirming the presence and accuracy of each of these
307 haplotypes.

308 2.2 Validation of manual joins with PCR

309 After polishing of the genome, we validated the manual contig joins described in
310 Table S2 as features 'Translocation' (chr13, chr18), 'Inversion' (chr05), 'Incomplete 1'
311 (chr04) and 'Incomplete 2' (chr07) using PCR. We designed primers using Primer3
312 v2.3.7 via the Python package primer3-py v0.5.4 and tested primers for other
313 occurrences in the genome using BLAST 2.9.0 via Biopython 1.74. The primer
314 sequences and next best hits in the genome are listed in Table S3 and primer
315 products in Table S4; the primers were designed against the polished genome
316 assembly (File S5) and so the join locations in Table S4 do not match the raw edit
317 positions above. To validate the incomplete chromosomes, single PCR products from
318 one pair of primers spanning a single join location were required; the translocation
319 and inversion required more complex validation involving four different combinations
320 of each set of four primers, listed in Table S4 and visualised in Figure 2.

321 All primer pairs produced a single product with the expected length as listed in Table
322 S4, except for the product of Inversion I1+I3 (Figure 2A), which was ~800 bp long
323 (rather than the expected 158 bp) and shows some evidence of producing multiple
324 products (smear on gel). This indicates that inversion Haplotype 2 (chr05b in Table
325 S4) was not reconstructed accurately, but this is not surprising given the repetitive
326 content typically found at inversion breakpoints. Inversion Haplotype 1 is included in
327 the genome assembly and has been validated by these PCRs, as have the other
328 manual joins. These PCRs therefore provide further evidence for the existence of the
329 inversion and translocation and the structural accuracy of the genome assembly.

330 3. Genome annotation

331 3.1 Transfer original annotations with flo

332 We used flo (Pracana et al. 2017, <https://github.com/wurmlab/flo>, commit 41f5ae4)
333 to transfer the GCA_000442575 *A. deanei* annotations to our new genome
334 assembly, using BLAT options -fastMap and -oneOff=1 but default BLAT options
335 otherwise (for example, we used the BLAT default minIdentity=90 rather than the flo
336 suggestion of minIdentity=98, given the known errors in nanopore genome
337 assemblies). We included the polished nuclear, symbiont, maxicircle and raw
338 minicircle assemblies in our new assembly, as the reference annotation includes
339 non-nuclear genes and we wanted to avoid transferring these to the nuclear genome
340 by mistake.

341 The GCA_000442575 annotation has 16 888 protein-coding genes, 45 tRNAs and 3
342 rRNAs. It has gene, mRNA, exon and CDS features for each protein-coding gene,
343 each with identical positions. flo transfers gene, exon and CDS features but not
344 mRNA features. Therefore, before running flo, we filtered these mRNA features,
345 other comments and region features from the annotation and updated the exon and
346 CDS Parent attributes using the following one-liner:

```
347 zcat GCA_000442575.2_Angomonas_deanei_Genome_genomic.gff.gz |  
348 grep -v "^##species" | awk '$3 !~ "region|RNA"' | sed -e  
349 's/Parent=rna/Parent=gene/g' > GCA_000442575.flo.gff3
```

350 flo produced a new GFF file containing 15 829 protein-coding genes transferred to
351 the new genome assembly. However, there were three problems with this transfer.
352 Firstly, many genes had duplicate annotations which needed to be collapsed to a
353 single annotation. Secondly, remaining errors in the new assembly meant some
354 transferred annotations did not produce valid protein sequences. Thirdly, flo only
355 transfers genes and not tRNAs and mRNAs.

356 *3.2 Filter duplicate annotations and fix sequence errors*

357 We wrote a Python script (File S6) to address duplicate annotations and errors in
358 gene sequences. This script takes the reference genome FASTA and annotation
359 GFF files, the new assembly FASTA and the flo-transferred GFF as input, and
360 produces updated FASTA and GFF files, as well as a TSV file describing how each
361 transferred GFF feature has been processed (File S7). The script does the following:

- 362 1. Build an interval tree for each chromosome using positions of CDS features to
363 identify sets of overlapping features.
- 364 2. For each set of overlapping features, choose one best feature to transfer (see
365 below for details).
- 366 3. If a chosen feature has a sequence in the new assembly that does not produce a
367 valid protein sequence, but the reference sequence does produce a valid protein,
368 replace the sequence with the sequence from the reference genome.
- 369 4. Output chosen features to a new GFF, updating coordinates to take replaced
370 sequences into account.
- 371 5. Output a new FASTA file containing fixed sequences.

372 A DNA sequence producing a valid protein is one that starts with a start codon, ends
373 with a stop codon, does not contain additional stop codons, does not contain Ns, and
374 whose length is divisible by 3.

375 Features were chosen from sets of overlapping features as follows:

376 1. Assign a status to each feature by examining the reference and transferred
377 sequences, including aligning and comparing the sequences. Sequences were
378 aligned with the Biopython pairwise2 module, using scores match=1, mismatch=-1,
379 open gap=-1, extend gap=-0.1. Possible statuses are:

380 - OK: reference and transferred protein sequences are valid and identical (although
381 the transferred DNA sequence may have synonymous substitutions)

382 - Changed: both sequences produce valid proteins of equal length but the
383 transferred protein sequence is different to the reference protein sequence

384 - NewLength: the transferred sequence produces a valid protein of a different
385 length to the reference sequence

386 - BadLength: transferred DNA sequence is not divisible by 3

387 - BadStart: first amino acid in transferred protein sequence is not M (methionine)

388 - BadStop: last amino acid in transferred protein sequence is not * (stop codon)

389 - ExtraStop: one or more stop codons (*) appear in transferred protein sequence

390 Valid sequences will be one of OK, Changed or NewLength, but invalid sequences

391 could have any combination of BadLength, BadStart, BadStop and ExtraStop

392 statuses.

393 2. Reject features where:

394 - the reference protein and the transferred protein are invalid (but consider features
395 with invalid reference proteins and valid transferred proteins, because in some cases
396 gaps have been filled in the new genome)

397 - the transferred sequence differs in length to the reference sequence by at least
398 10% (likely indicating a bad alignment)

399 3. Search for an acceptable feature within each group of overlapping features,
400 checking named features first, then hypothetical features; searching OK, Changed,
401 and NewLength features of each kind in that order; and ordering features of the
402 same type by largest alignment score first. Only consider NewLength features where
403 the reference DNA sequence contains Ns. Choose the first acceptable feature by
404 this ordering.

405 4. If no acceptable feature is found, search through features with BadLength,
406 BadStart, BadStop and ExtraStop statuses, again checking named features first,
407 then hypothetical features, and sorting features of the same kind by smallest length
408 difference and then largest alignment score. Take the first feature by this ordering
409 and, if the reference protein is valid, choose this feature and mark the sequence for
410 replacement.

411 The 15 829 transferred CDS features were collapsed into 8 001 groups of
412 overlapping features, with between 1 and 15 features in each group, indicating the
413 highly redundant nature of the annotation; 3 878 groups contained more than one

414 overlapping feature. 748 features were rejected because the reference and
415 transferred proteins were invalid; 72 features were rejected because the transferred
416 length differed from the reference length by at least 10%. The remaining 15 009
417 features were considered for inclusion. Of these, 5 379 were accepted, 2 191 were
418 replaced with the reference sequence, and 7 439 were rejected in favour of another
419 feature; a feature was output for 7 570 of the 8 001 groups of overlapping features,
420 but good features could not be found for 431 groups. For the nuclear genome, 7 502
421 of 7 932 groups had features output, with 5 322 features accepted as is and 2 180
422 replaced with the reference sequence. The new nuclear genome assembly
423 increased in length from 20 975 274 bp to 20 976 081 bp long, an increase of 807
424 bp, with 2 917 803 bp of new sequence replaced by 2 918 610 bp of reference
425 sequence to ensure protein sequences were valid.

426 The fixed GFF was then updated to recover the mRNA features from the original
427 annotation, restoring the gene/mRNA/exon/CDS features for each of the 7 502
428 transferred protein coding genes.

429 *3.3 Companion run*

430 To search for novel genes, to annotate tRNAs and rRNAs (as flo had not transferred
431 the reference annotation's tRNAs and rRNAs), to annotate with Pfam and GO terms,
432 and to provide an EMBL format genome suitable for submission to public databases,
433 we ran Companion on the fixed genome assembly, using our transferred annotation
434 as a reference. This is an unusual Companion use case, as Companion usually
435 expects the reference to be from a different species, and this required some

436 modifications of the Companion process. Only the run_exonerate, make_embl,
437 use_reference and truncate_input_headers parts of the pipeline were run; the rest of
438 the pipeline was turned off. In particular, RATT was not run, because we did not
439 need to transfer annotations to the new assembly, and instead the pipeline was
440 edited to accept our new transfer as if it were RATT output. We also wrote a new
441 weight function, based on the default weight__kinetoplastid.lua function and passed
442 as the WEIGHT_FILE option (File S8). This function is intended to accept any
443 transferred gene over an Augustus prediction by giving it a 100-fold increase in
444 score, unless that gene is hypothetical, in which case it gets the standard 3-fold
445 increase in score relative to an Augustus prediction.

446 The Companion GFF3 output is in File S10. Companion also outputs EMBL files
447 suitable for submission to public databases; however, they do not include all
448 attributes from the reference and Companion GFF files, including old gene names.
449 We wrote a script (File S11) to add this information from the original and transferred
450 GFF files to the final annotation in EMBL format (File S12).

451 The final annotation contains 10 365 protein-coding genes (of which 7 199 were
452 transferred from the reference annotation and 3 166 were predicted by Augustus), 59
453 tRNAs, 26 rRNAs, 45 ncRNAs, 14 snoRNAs and 3 snRNAs. Although we did not
454 transfer the 45 tRNAs and 3 rRNAs from the reference annotation, alignments
455 showed that Companion has identified all of these RNA features and more. 303 of
456 the 7 502 previously transferred features were hypothetical proteins replaced by a
457 better Augustus prediction.