

1 **Improved contiguity of the threespine stickleback genome using long-read sequencing**

2

3 Shivangi Nath^{*}, Daniel E. Shaw^{*}, and Michael A. White^{*}

4

5 ^{*} Department of Genetics, University of Georgia, Athens, GA, 30602

6

7

8 **Running title:** Closing gaps in the stickleback genome

9 **Key words:** threespine stickleback fish, long-read sequencing, genome assembly, genome
10 browser, telomere sequence, centromere sequence

11

12 **Corresponding author:**

13 Michael A. White

14 Department of Genetics, University of Georgia

15 120 Green St.

16 Athens, GA 30602

17 whitem@uga.edu

18

19

20

21

22 **Abstract**

23 While the cost and time for assembling a genome have drastically reduced, it still
24 remains a challenge to assemble a highly contiguous genome. These challenges are rapidly
25 being overcome by the integration of long-read sequencing technologies. Here, we use long
26 sequencing reads to improve the contiguity of the threespine stickleback fish (*Gasterosteus*
27 *aculeatus*) genome, a prominent genetic model species. Using Pacific Biosciences sequencing,
28 we were able to fill over 76% of the gaps in the genome, improving contiguity over five-fold.
29 Our approach was highly accurate, validated by 10X Genomics long-distance linked-reads. In
30 addition to closing a majority of gaps, we were able to assemble segments of telomeres and
31 centromeres throughout the genome. This highlights the power of using long sequencing reads
32 to assemble highly repetitive and difficult to assemble regions of genomes. This latest genome
33 build has been released through a newly designed community genome browser that aims to
34 consolidate the growing number of genomics datasets available for the threespine stickleback
35 fish.

36

37 **Introduction**

38 Reference genome assemblies have been invaluable in the discovery of genes, the
39 annotation of regulatory regions, and for providing a scaffold for understanding genetic
40 variation within a species. With the advent of new sequencing technologies and the reduction
41 of cost, there has been a rapid increase in the total number of reference genomes available
42 across taxa. Although it has become much simpler to produce a draft genome assembly, the
43 completion of a high-quality, contiguous assembly remains a great challenge. There are many

44 regions within individual genomes that are unassembled. These regions are enriched for highly
45 repetitive sequence that cannot be assembled using sequencing technologies that produce
46 short fragments (Nagarajan and Pop 2013; Gnerre et al. 2011). Even the most highly refined
47 genomes, like the human genome till have many gaps, which often are composed of long
48 segmental duplications (Schneider et al. 2017).

49 Long reads from third-generation sequencing technologies have shown promise in
50 spanning highly repetitive regions of genomes, bridging previously intractable gaps in
51 assemblies to improve overall contiguity. Within the human genome, many highly repetitive
52 regions have been resolved, such as pericentromeres (Vollger et al. 2020), complete
53 centromeres (Jain et al. 2018b), telomeres (Jain et al. 2018a) and the entire major
54 histocompatibility complex (Jain et al. 2018a). *De novo* assemblies of highly repetitive Y
55 chromosomes have also become feasible using long-read sequencing (Mahajan et al. 2018;
56 Peichel et al. 2019). Overall, long-read sequencing has enabled telomere-to-telomere
57 chromosome assemblies in multiple species (Liu et al. 2020; Miga et al. 2019). It is clear that
58 hybrid assembly approaches incorporating long-read sequencing have greatly improved
59 contiguity of genomes.

60 Here we use long-read sequencing to improve the contiguity of the threespine
61 stickleback fish (*Gasterosteus aculeatus*) genome. The threespine stickleback fish has been an
62 important model system to understand evolution, ecology, physiology and toxicology (Wootton
63 1976; Bell and Foster 1994). The identification of the genetic mechanisms underlying many
64 adaptative traits was facilitated by the release of a high-quality reference genome assembly
65 (Jones et al. 2012). This genome assembly was constructed using paired-end Sanger sequencing

66 reads from multiple genomic libraries. Contigs were scaffolded to genetic linkage maps, which
67 resulted in 21 chromosome-level scaffolds (400.4 Mb), with 60.7 Mb of unplaced scaffolds. The
68 assembly has undergone several revisions, using high-density genetic linkage maps (Roesti et al.
69 2013; Glazer et al. 2015), and a Hi-C proximity-guided assembly (Peichel et al. 2017). Despite
70 multiple revisions, the latest version of the assembly (v. 4) still contains 13,538 gaps and 20.6
71 Mb of unplaced scaffolds (Peichel et al. 2017). The gaps between contigs in the chromosome
72 scaffolds likely represent repetitive regions or GC-rich regions, which have been shown to be
73 recalcitrant to traditional assembly methods (Benjamini and Speed 2012; Ross et al. 2013). In
74 order to improve the assembly, we used long-read sequencing to fill gaps in the threespine
75 stickleback genome assembly. We were able to close 76.7% of the gaps, incorporating 13.5% of
76 the previously unplaced scaffolds. Closed gaps were highly accurate, verified through long-
77 distance linked-read information. In addition, we were able to extend sequence of many of
78 chromosomes into telomeres. This assembly represents a noteworthy improvement, allowing
79 researchers to interrogate many previously inaccessible repetitive regions, and highlights the
80 power of long-read sequencing to substantially improve genome contiguity.

81

82 **Methods**

83 **Ethics statement**

84 All procedures using threespine stickleback fish were approved by the University of
85 Georgia Animal Care and Use Committee (protocol A2018 10-003-Y2-A5).

86

87 **Closing gaps in the reference assembly**

88 Version four of the threespine stickleback reference genome assembly contains 1263
89 unplaced contigs (chr. Un) that were narrowed to chromosomes but were not placed into
90 specific gaps (there was a total of 3378 chr. Un contigs: 1263 contigs were narrowed broadly to
91 chromosomes and 2115 could not be localized to any chromosome). We used recently available
92 high-coverage long-read sequencing in combination with the 1263 chr. Un contigs that were
93 previously narrowed to chromosomes to fill the remaining gaps in the reference assembly. A
94 male Paxton Lake benthic threespine stickleback fish (Texada Island, British Columbia) was
95 sequenced to approximately 75x coverage (NCBI BioProject database accession PRJNA591630)
96 and assembled into contigs using Canu (Koren et al. 2017). This assembly had a total of 3593
97 contigs (N50: 683 kb) from across the genome. We closed gaps in the v. 4 threespine
98 stickleback reference assembly using these contigs and LR_Gapcloser with the parameter -a 1
99 (Xu et al. 2019). We increased the allowed deviation between gap length and the inserted
100 sequence length in order to provide additional flexibility for gap size that was not inferred
101 accurately in the v. 4 genome assembly. LR_Gapcloser fills existing gaps in the genome
102 assembly by identifying contigs which span a gap completely or partially from either end.
103 Three Canu contigs caused a reduction in total chromosome size after placement into gaps.
104 These contigs were likely misassembled by Canu, causing complications in the gap closing. We
105 omitted these three contigs from further analysis. We used BLAT (v. 3.5; Kent 2002) to identify
106 which of the 1263 previously narrowed chr. Un contigs were placed within a gap. We filtered
107 for stringent alignments by only retaining matches where at least 90% of the query length
108 aligned to the assembly and the total aligned region had 2% or less mismatches.

109 Many chr. Un contigs that were not placed in the v. 4 genome assembly may be
110 represented in the v. 5 assembly if they were contained completely within a PacBio Canu contig
111 (Peichel et al. 2019). To test this, we used BLAT to align the 3378 chr. Un contigs to the new v. 5
112 assembly. We filtered for stringent alignments by only retaining matches where at least 90% of
113 the query length aligned to the assembly and the total aligned region had 2% or less
114 mismatches. Chr. Un contigs that did not align to the assembly were retained as unassembled
115 and concatenated into a single fasta sequence, with each contig separated by 100 base pairs.

116

117 **Validation of the genome assembly using long-distance linked-reads**

118 We verified the revised genome assembly using 10X Genomics long-distance linked-read
119 sequencing from a single female freshwater threespine stickleback fish (Lake Washington,
120 Washington, USA). We extracted high molecular weight DNA from blood using alkaline lysis.
121 Blood was collected from euthanized fish into 0.85x SSC buffer. The cells were collected by
122 centrifuging for two minutes at 2000 xg. Pelleted cells were resuspended in five ml of 0.85x SSC
123 and 27 μ l of 20 μ g/ml Proteinase K solution. To lyse the cells, five ml of 2x SDS buffer (80mM
124 EDTA, 100mM Tris pH 8.0, and 1% SDS) was added to the suspension and the solution was
125 incubated at 55°C for two minutes. After incubation, 10 ml of buffered
126 phenol/chloroform/isoamyl-alcohol was added to the suspension. The suspension was
127 incubated at room temperature under slow rotation for 30 minutes. The suspension was
128 centrifuged for one minute at 2000 xg at 4°C to separate phases. The aqueous phase was
129 extracted, 10 ml of chloroform was added, and the suspension was rotated for one hour. The
130 chloroform extraction step was repeated twice. After all extractions, the aqueous phase was

131 separated and mixed with ice cold 100% ethanol and one ml of 3M sodium-acetate (pH 5.5).
132 The overall alignment rate of linked-reads to the assembly was 84.4%, resulting in a genome-
133 wide mean read depth of 26.1X.

134

135 **Assessing the completeness of the genome assembly**

136 In order to assess the completeness of the genome, we identified universal single copy
137 orthologs (BUSCO) throughout the new assembly, compared to the previous v. 4 assembly
138 (Peichel et al. 2017). BUSCO (v. 3.0.2) was run using default parameters, comparing against the
139 Actinopterygii lineage dataset (4584 total single copy orthologs; OrthoDB v. 9) (Simão et al.
140 2015). Actinopterygii was used because threespine stickleback fish are teleosts, which is the
141 largest infraclass of Actinopterygii.

142

143 **Identification of telomeric sequences**

144 PacBio long reads with highly repetitive regions are often not assembled into contigs. In
145 order to identify telomeric reads, we searched for the ancestral metazoan telomeric
146 motif 'TTAGGG' or 'CCCTAA' (Traut et al. 2007; Meyne et al. 1989; Moyzis et al. 1988) in the raw
147 PacBio reads. We searched for the motif and their respective counts in each read using the awk
148 command-line utility. Reads were considered for further analyses if they had more than 50
149 occurrences of the motif. These reads were aligned to the new genome assembly using
150 minimap2 (v. 2.17) (Li 2018) with default parameters to map to PacBio genomic reads (-ax map-
151 pb). Only the primary alignments were retained. Telomeric reads were assigned to a specific
152 chromosome if greater than 10 kb of unique sequence overlapped with one end of a

153 chromosome. Positive telomeric alignments were merged with the new genome assembly.
154 Repetitive sequence content within telomeres was visualized using the dotplot function in
155 Geneious Prime (v2019 1.1) (<https://www.geneious.com>).

156

157 **Identification of centromeric sequences**

158 BLAST+ (blastn; v. 2.7.1) (Camacho et al. 2009) was used to identify the 186 bp
159 threespine stickleback CENP-A monomer repeat (Cech and Peichel 2015) in the PacBio Canu
160 assembled contigs. Contigs containing CENP-A repeats were mapped to the new v. 5 repeat
161 masked assembly (see Genome annotation and repeat masking) using minimap2 (Li 2018) with
162 default parameters to map to PacBio genomic reads (-ax map-pb). Contigs were only retained if
163 greater than 10 kb of sequence mapped uniquely to one chromosome end. The number of
164 CENP-A repeats per chromosome were counted using blastn. Dotplots were generated using
165 Geneious Prime (v2019 1.1) (<https://www.geneious.com>).

166

167 **Genome annotation**

168 Genome features were lifted over from the previous assembly (v. 4) using a hybrid
169 approach. Genome features were first lifted over to the new assembly using the software
170 package flo (Pracana et al. 2017). Most of the features were lifted over successfully (98.1%). We
171 used BLAT to lift over the remaining fraction. The sequence for the features not lifted over with
172 flo were extracted from the version four assembly using samtools faidx. These sequences were
173 then aligned to the new assembly using BLAT. For each feature, the longest alignment was
174 chosen.

175 Many of the closed gaps were not represented in the previous genome assembly (v. 4)
176 and were therefore unannotated. We annotated these regions using the MAKER (v. 3.01.02)
177 genome annotation pipeline (Cantarel et al. 2008; Holt and Yandell 2011). These annotations
178 combined evidence from multiple RNA-seq transcriptomes, all predicted Ensembl protein
179 sequences (release 95), and *ab initio* gene predictions from SNAP (v. 2006-07-28) (Korf 2004)
180 and Augustus (v. 3.3.2) (Stanke et al. 2006). MAKER was run over three rounds using the RNA-
181 seq transcriptomes and methods previously described (Peichel et al. 2019).

182 Repeats were annotated across the genome using a combination of RepeatModeler (v.
183 1.0.11) and RepeatMasker (v. 4.0.5) (<http://www.repeatmasker.org>). Repeats were first
184 modeled using the default parameters of RepeatModeler. Repeats were then annotated and
185 masked using RepeatMasker with default parameters and the custom RepeatModeler
186 database.

187 We tested for enrichment of repeats and genes in closed gaps throughout the genome
188 by comparing to randomly drawn 10 Mb segments (we placed 9.9 Mb of sequence within gaps;
189 see Results). We also tested for enrichment of repeats and transposable elements in the
190 remainder of the unplaced chr. Un contigs by comparing to randomly drawn 20 Mb segments
191 throughout the assembled genome (19.6 Mb of chr. Un contigs remained unplaced; see
192 Results). We generated a null distributions by randomly drawing 10,000 segments throughout
193 the genome using bedtools (v. 2.29.2) shuffle (Quinlan and Hall 2010). We then used bedtools
194 intersect to count the number of repeats (with option -c for both 10Mb and 15Mb segments) as
195 well as the number of bases that overlapped genes (with option -wao for 10 Mb segments)
196 within each random segment.

197

198 **Data Availability**

199 Long-distance linked-read sequencing and the updated genome assembly are available
200 on the NCBI BioProject database (<https://www.ncbi.nlm.nih.gov/bioproject/>) under accession
201 number PRJNA639125. The genome assembly is also available for download from the
202 threespine stickleback genome browser (<https://stickleback.genetics.uga.edu>). All
203 supplemental material has been uploaded to figshare. Reviewer link for accessing the SRA data:
204 <https://dataview.ncbi.nlm.nih.gov/object/PRJNA639125?reviewer=agjcut825ub04a0fvdjqj9vpa>

205 2

206

207 **Results and Discussion**

208 **A majority of gaps were closed across the threespine stickleback genome**

209 Version four of the threespine stickleback genome assembly (Peichel et al. 2017)
210 contained 13,538 gaps with an N50 of 91.7kb (total estimated length of gaps: 3,573,762 bp).
211 Using long-read PacBio contigs in conjunction with the 1263 chr. Un contigs that had been
212 narrowed to chromosomes from the previous assembly, we closed 10,394 of the gaps (76.8%),
213 leaving 3,144 gaps in the final assembly (File S1, File S2). In addition to the fully closed gaps, 146
214 gaps were partially closed. A total of 9,928,283 bases were added to gaps in the assembly. This
215 resulted in an overall greater contiguity of the genome, with a 5.57-fold greater N50 contig
216 length compared to the previous assembly (v. 5 N50: 510.8 kb; v. 4 N50: 91.7 kb) (Table 1).

217 Genome contiguity and annotation completeness is often assessed by BUSCO
218 (benchmarking universal single copy orthologs) statistics (Simão et al. 2015). We determined if

219 the additional sequence in the new assembly contained coding sequence that improved overall
220 BUSCO metrics. Of the 3640 genes within the database, we found a total of 3521 BUSCO genes
221 in the v.5 assembly. This represented an increase of 99 genes compared to the previous
222 assembly. In addition, the total number of fragmented BUSCO genes decreased to 14,
223 compared to 108 in the previous assembly (Table S1).

224 Of the 3378 chr. Un contigs from the previous assembly (v. 4), we determined how
225 many were represented in the closed gaps of the new assembly. Of the 3378 contigs, 457
226 contigs were placed within gaps (13.5%). The previous assembly used a Hi-C-based proximity-
227 guided assembly method that was able to narrow some of the chr. Un contigs (1263) to
228 chromosomes, but was not able to place these contigs into specific gaps (Peichel et al. 2017).
229 We used this information to verify whether our contig placement was corroborated by the Hi-C
230 sequencing. Of the 1263 previously narrowed chr. Un contigs, we placed 90 of into gaps. A
231 majority of these contigs (80.0%) fell within the same chromosome they were assigned to
232 previously by the Hi-C proximity-guided method. This high concordance further confirms our
233 methodology and closure of the gaps.

234 Across all closed gaps, we added 9.9 Mb of sequence to the genome. 1.0 Mb of this
235 newly added sequence was from chr. Un contigs previously sequenced, but not placed in
236 chromosomes. The remaining 8.9 Mb represented new regions from the long-read sequencing.
237 Many of the gaps in the genome likely represent highly repetitive regions that are challenging
238 to assemble. We compared the repetitive sequence content between the 9.9 Mb of newly
239 added sequence and the remainder of the genome. Indeed, we found newly closed gaps are
240 significantly enriched for repeat sequences (simple and interspersed repeats; 10,000

241 permutations; $P < 0.001$; Figure S1). Overall, 17.4% of newly added bases contained repetitive
242 DNA compared to 13.5% in the remainder of the genome. Across all newly added gap
243 sequence, we found a total of 1226 protein coding genes. The newly placed regions overall
244 exhibit lower density of coding sequence than the remainder of the genome, although this
245 result was not statistically significant (Figure S1; 10,000 permutations; $P < 0.083$). Only 8.3% of
246 the closed gap bases were contained within coding regions. Across the remainder of the
247 genome 28.3% of bases in the assembly were contained within coding regions. Combined, our
248 results suggest the highly repetitive nature of the sequence contained within these gaps may
249 have prevented assembly of these regions.

250 Although we closed a majority of gaps in the assembly, we were still unable to
251 determine where 2921 of the chr. Un contigs belonged in the assembly (total length:
252 19,587,834 bp). One possibility why we were unable to place these contigs is they are even
253 more difficult to assemble due to higher repetitive sequence content. Consistent with this, the
254 unplaced contigs were highly enriched for retrotransposons compared to the placed chr. Un
255 contigs ($P < 0.001$; Figure S2). It is also possible that these contigs represent segments of the
256 genome outside of gaps that are mis-assembled. Our method only focused on closing gaps
257 between contigs. Additional work will be necessary to determine whether these contigs
258 integrate elsewhere in the genome.

259

260 **Long-distance linked-reads validated the closing of most gaps**

261 Long-distance linked-reads were used to validate placement of the new sequence.
262 Linked-read molecules that support closure of a gap would exhibit aligned short-reads

263 throughout the closed gap, whereas linked-read molecules that do not support closure of a gap
264 would have aligned short-reads outside of the gap, but a lack of alignment within the gap
265 (Figure 1). The gap closures were highly supported by the linked-read alignments. We only
266 observed 36 gaps (0.3%) that were not supported by linked-reads (i.e. a lack of short-read
267 alignments over the newly added sequence). The remainder of the 10,394 gaps in this analysis
268 that were closed (99.7%) were supported by the long-distance linked-read dataset (Figure 1).

269

270 **Telomere repeats and centromere repeats were identified within PacBio long reads**

271 The telomeres of threespine stickleback fish contain a tandemly repeated G-rich
272 hexanucleotide repeat that is conserved across metazoans (Meyne et al. 1989; Moyzis et al.
273 1988; Traut et al. 2007; Ocalewicz 2013). Although DNA probes targeting these repeats clearly
274 hybridize at the ends of all chromosomes in threespine stickleback fish, the underlying
275 sequence of these regions is missing from the genome assembly. We therefore searched for the
276 ancestral metazoan telomeric motif 'TTAGGG' or 'CCCTAA' in the raw PacBio reads to identify
277 putative telomere caps (Konrad et al. 2011; Ocalewicz et al. 2011). We identified 3525 PacBio
278 reads that contained telomere motifs. Seven of these reads contained enough unique sequence
279 to align to the end of individual chromosomes (chromosomes IV, VII, VIII, X, XIV, XV, XVII). These
280 reads showed an abundance of the ancestral metazoan telomeric motif at one end with little to
281 no higher order structure (Figure 2; Figure S3). The telomeric motif was repeated 114-492 times
282 throughout the sequence on different chromosomes.

283 We also searched for centromere repeats within the PacBio assembled contigs. We
284 identified the core 186 bp CENP-A repeat (Cech and Peichel 2015) within 91 contigs (the length

285 of repetitive DNA among contigs ranges from 12.61 – 125.17 kb). 48 of these contigs contained
286 enough unique sequence to align to all 21 chromosomes (Figure 3; File S3; File S4). 11
287 chromosomes had centromere contigs that map to both sides of the gap, 9 chromosomes had a
288 centromere contig only on one end of the centromere, and one contig contained a full
289 centromere sequence, spanning the entire gap (chromosome IX). Interestingly, on many of the
290 chromosomes, the repeat length was long enough to discern clear higher order repeat
291 structure (Figure 3; Figure S4). Our results are similar to the variability in higher order repeat
292 among the autosomes and X chromosome of humans (Hartley and O'Neill 2019; Willard 1985;
293 Willard et al. 1986; Alexandrov et al. 1993; Shepelev et al. 2015). We detected multiple contigs
294 mapping to either end of the centromere gap for all chromosomes (File S3, File S4), indicating
295 the male fish used for sequencing is likely heterozygous for centromeric arrays. This is
296 consistent with high polymorphism of centromere arrays observed within other species (Willard
297 et al. 1986; Greig et al. 1991; Mahtani and Willard 1990; Devilee et al. 1988; Wevrick and
298 Willard 1989).

299 Y chromosomes in mammals have also been documented to have highly variable
300 centromeric repeats that are divergent from their counterparts across the remainder of the
301 genome (Wolfe et al. 1985; Pertile et al. 2009; Miga et al. 2014). Assembly of segments of the
302 threespine stickleback Y chromosome centromere (Peichel et al. 2019) revealed an alpha
303 satellite monomer repeat that was divergent from the consensus monomeric repeat identified
304 from the remainder of the genome (Cech and Peichel 2015). With the assembly of larger tracks
305 of centromeric sequence from the autosomes and the X chromosome, we now show the Y
306 chromosome centromere is also divergent from the remainder of the genome at the level of

307 higher order repeats (Peichel et al. 2019), matching other rapidly evolving Y chromosomes.
308 Although our assembly has uncovered a large fraction of the centromeric sequence for each
309 chromosome, we were unable to assemble complete centromere sequences outside of the 46.5
310 kb centromere of chromosome IX. It therefore remains unknown how centromere length varies
311 throughout the threespine stickleback genome. Complete characterization of the centromeric
312 repetitive arrays will be aided by future sequencing of ultra-long reads (Jain et al. 2018b; Miga
313 et al. 2019).

314

315 **Conclusions**

316 By using long-read sequencing we were able to substantially improve the overall
317 contiguity of the threespine stickleback genome, increasing the N50 length of contigs over five-
318 fold. Our assembly also highlights the power of using long-read sequencing technologies to
319 assemble previously inaccessible regions of the genome, like centromeres and telomeres. We
320 have released this assembly through a new threespine stickleback fish community genome
321 browser (<https://stickleback.genetics.uga.edu>). This resource will be a useful addition to the
322 rapidly expanding functional genomics toolkit available in threespine stickleback fish.

323

324 **Acknowledgements**

325 This research was funded by the National Science Foundation IOS 1645170 (M.A.W.),
326 the Office of the Vice President of Research at the University of Georgia (M.A.W.), and the
327 University of Georgia Research Foundation (D.E.S.). We thank the Georgia Genomics and
328 Bioinformatics Core at the University of Georgia for help with the long-distance linked-read

329 sequencing. We also thank Brigitte Hofmeister and the Franklin College Office of Information
330 Technology at the University of Georgia for help building the threespine stickleback genome
331 browser.

332

333 **Figure Legends**

334 Figure 1. 10X Genomics linked-reads validate most of the closed gaps. 99.7% of closed gaps
335 exhibit linked-read alignments throughout the gap region, indicating a correctly closed gap (e.g.
336 Chr. XV: 8,942,107-8,942,560 bp). 0.03% of gaps were not validated by the linked-read
337 sequencing. In these regions, alignments of the linked-reads only occur outside of the gap (e.g.
338 Chr. XII: 20,435,502-20,437,016 bp). A representative schematic outlining how the linked-reads
339 should align is shown in black. The actual aligned linked-reads are shown by bolded color lines.
340 Thin lines indicate gaps between the linked-reads. Average read depth of linked-reads across
341 the genome was 26.1X. A subset of reads aligning is shown here for simplicity.

342

343 Figure 2. Telomeres exhibit a high density of the conserved metazoan telomere motif. Dots
344 represent 100% sequence identity between matching windows of 15 bp. The blue box
345 represents the end of chr. VII where the long read aligns uniquely to positions 30,722,876-
346 30,767,092. The green box denotes a ~10 kb segment rich with telomeric repeat sequence
347 (TRS). The remaining telomeres are shown in Figure S3.

348

349 Figure 3. Centromeres display higher order repeat structure. On chromosome I, the centromere
350 contig contains 673 186 bp monomer repeats. Sequence identity between repeats is depicted

351 by black dots matching windows of 300 bp with 100% sequence identity. The blue region
352 denotes the end of chromosome I (20,330,007-20,344,665 bp) with unique sequence. The
353 green region is the newly aligned centromere contig. The remaining centromeres are shown in
354 Figure S4.

355

356

357

358

359

360 **References**

- 361
- 362 Alexandrov, I.A., L.I. Medvedev, T.D. Mashkova, L.L. Kisselev, L.Y. Romanova *et al.*, 1993
- 363 Definition of a new alpha satellite suprachromosomal family characterized by
- 364 monomeric organization. *Nucleic Acids Res* 21 (9):2209-2215.
- 365 Bell, M., and S.A. Foster, 1994 *The evolutionary biology of the threespine sticklebacks*. Oxford
- 366 University Press.
- 367 Benjamini, Y., and T.P. Speed, 2012 Summarizing and correcting the GC content bias in high-
- 368 throughput sequencing. *Nucleic Acids Res* 40 (10):e72.
- 369 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos *et al.*, 2009 BLAST+: architecture
- 370 and applications. *BMC Bioinformatics* 10:421.
- 371 Cantarel, B.L., I. Korf, S.M. Robb, G. Parra, E. Ross *et al.*, 2008 MAKER: an easy-to-use
- 372 annotation pipeline designed for emerging model organism genomes. *Genome Res* 18
- 373 (1):188-196.
- 374 Cech, J.N., and C.L. Peichel, 2015 Identification of the centromeric repeat in the threespine
- 375 stickleback fish (*Gasterosteus aculeatus*). *Chromosome Res* 23 (4):767-779.
- 376 Devilee, P., T. Kievits, J.S. Wayne, P.L. Pearson, and H.F. Willard, 1988 Chromosome-specific
- 377 alpha satellite DNA: isolation and mapping of a polymorphic alphoid repeat from human
- 378 chromosome 10. *Genomics* 3 (1):1-7.
- 379 Glazer, A.M., E.E. Killingbeck, T. Mitros, D.S. Rokhsar, and C.T. Miller, 2015 Genome Assembly
- 380 Improvement and Mapping Convergent Evolutionary Traits in Sticklebacks with
- 381 Genotyping-by-Sequencing. *G3 (Bethesda)* 5 (7):1463-1472.
- 382 Gnerre, S., I. Maccallum, D. Przybylski, F.J. Ribeiro, J.N. Burton *et al.*, 2011 High-quality draft
- 383 assemblies of mammalian genomes from massively parallel sequence data. *Proc Natl*
- 384 *Acad Sci U S A* 108 (4):1513-1518.
- 385 Greig, G.M., S. Parikh, J. George, V.E. Powers, and H.F. Willard, 1991 Molecular cytogenetics of
- 386 alpha satellite DNA from chromosome 12: fluorescence in situ hybridization and
- 387 description of DNA and array length polymorphisms. *Cytogenet Cell Genet* 56 (3-4):144-
- 388 148.
- 389 Hartley, G., and R.J. O'Neill, 2019 Centromere Repeats: Hidden Gems of the Genome. *Genes*
- 390 *(Basel)* 10 (3).
- 391 Holt, C., and M. Yandell, 2011 MAKER2: an annotation pipeline and genome-database
- 392 management tool for second-generation genome projects. *BMC Bioinformatics* 12:491.
- 393 Jain, M., S. Koren, K.H. Miga, J. Quick, A.C. Rand *et al.*, 2018a Nanopore sequencing and
- 394 assembly of a human genome with ultra-long reads. *Nat Biotechnol* 36 (4):338-345.
- 395 Jain, M., H.E. Olsen, D.J. Turner, D. Stoddart, K.V. Bulazel *et al.*, 2018b Linear assembly of a
- 396 human centromere on the Y chromosome. *Nat Biotechnol* 36 (4):321-323.
- 397 Jones, F.C., M.G. Grabherr, Y.F. Chan, P. Russell, E. Mauceli *et al.*, 2012 The genomic basis of
- 398 adaptive evolution in threespine sticklebacks. *Nature* 484 (7392):55-61.
- 399 Kent, W.J., 2002 BLAT--the BLAST-like alignment tool. *Genome Res* 12 (4):656-664.
- 400 Konrad, O., W. Piotr, F.-S. Grazyna, and J. Malgorzata, 2011 Chromosomal location of Ag/CMA
- 401 3 -NORs, 5S rDNA and telomeric repeats in two stickleback species, pp. 12-19. Italian
- 402 Journal of Zoology.

- 403 Koren, S., B.P. Walenz, K. Berlin, J.R. Miller, N.H. Bergman *et al.*, 2017 Canu: scalable and
404 accurate long-read assembly via adaptive. *Genome Res* 27 (5):722-736.
- 405 Korf, I., 2004 Gene finding in novel genomes. *BMC Bioinformatics* 5:59.
- 406 Li, H., 2018 Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34
407 (18):3094-3100.
- 408 Liu, J., A.S. Seetharam, K. Chougule, S. Ou, K.W. Swentowsky *et al.*, 2020 Gapless assembly of
409 maize chromosomes using long-read technologies. *Genome Biol* 21 (1):121.
- 410 Mahajan, S., K.H. Wei, M.J. Nalley, L. Gibilisco, and D. Bachtrog, 2018 De novo assembly of a
411 young *Drosophila* Y chromosome using single-molecule sequencing and chromatin
412 conformation capture. *PLoS Biol* 16 (7):e2006348.
- 413 Mahtani, M.M., and H.F. Willard, 1990 Pulsed-field gel analysis of alpha-satellite DNA at the
414 human X chromosome centromere: high-frequency polymorphisms and array size
415 estimate. *Genomics* 7 (4):607-613.
- 416 Meyne, J., R.L. Ratliff, and R.K. Moyzis, 1989 Conservation of the human telomere sequence
417 (TTAGGG)_n among vertebrates. *Proc Natl Acad Sci U S A* 86 (18):7049-7053.
- 418 Miga, K.H., S. Koren, A. Rhie, M.R. Vollger, A. Gershman *et al.*, 2019 Telomere-to-telomere
419 assembly of a complete human X chromosome. *bioRxiv*:735928.
- 420 Miga, K.H., Y. Newton, M. Jain, N. Altemose, H.F. Willard *et al.*, 2014 Centromere reference
421 models for human chromosomes X and Y satellite arrays. *Genome Res* 24 (4):697-707.
- 422 Moyzis, R.K., J.M. Buckingham, L.S. Cram, M. Dani, L.L. Deaven *et al.*, 1988 A highly conserved
423 repetitive DNA sequence, (TTAGGG)_n, present at the telomeres of human
424 chromosomes. *Proc Natl Acad Sci U S A* 85 (18):6622-6626.
- 425 Nagarajan, N., and M. Pop, 2013 Sequence assembly demystified. *Nat Rev Genet* 14 (3):157-
426 167.
- 427 Ocalewicz, K., 2013 Telomeres in fishes. *Cytogenet Genome Res* 141 (2-3):114-125.
- 428 Ocalewicz, K., P. Woznicki, G. Furgala-Selezniow, and M. Jankun, 2011 Chromosomal location of
429 Ag/CMA 3 -NORs, 5S rDNA and telomeric repeats in two stickleback species. *Italian*
430 *Journal of Zoology*:(12-19).
- 431 Peichel, C.L., S.R. McCann, J.A. Ross, A.F.S. Naftaly, J.R. Urton *et al.*, 2019 Assembly of a young
432 vertebrate Y chromosome reveals convergent signatures of sex chromosome evolution.
433 *bioRxiv*:2019.2012.2012.874701.
- 434 Peichel, C.L., S.T. Sullivan, I. Liachko, and M.A. White, 2017 Improvement of the Threespine
435 Stickleback Genome Using a Hi-C-Based Proximity-Guided Assembly. *J Hered* 108
436 (6):693-700.
- 437 Pertile, M.D., A.N. Graham, K.H. Choo, and P. Kalitsis, 2009 Rapid evolution of mouse Y
438 centromere repeat DNA belies recent sequence stability. *Genome Res* 19 (12):2202-
439 2213.
- 440 Pracana, R., A. Priyam, I. Levantis, R.A. Nichols, and Y. Wurm, 2017 The fire ant social
441 chromosome supergene variant Sb shows low diversity but high divergence from SB.
442 *Mol Ecol* 26 (11):2864-2879.
- 443 Quinlan, A.R., and I.M. Hall, 2010 BEDTools: a flexible suite of utilities for comparing genomic
444 features. *Bioinformatics* 26 (6):841-842.
- 445 Roesti, M., D. Moser, and D. Berner, 2013 Recombination in the threespine stickleback genome-
446 -patterns and consequences. *Mol Ecol* 22 (11):3014-3027.

447 Ross, M.G., C. Russ, M. Costello, A. Hollinger, N.J. Lennon *et al.*, 2013 Characterizing and
448 measuring bias in sequence data. *Genome Biol* 14 (5):R51.

449 Schneider, V.A., T. Graves-Lindsay, K. Howe, N. Bouk, H.C. Chen *et al.*, 2017 Evaluation of
450 GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of
451 the reference assembly. *Genome Res* 27 (5):849-864.

452 Shepelev, V.A., L.I. Uralsky, A.A. Alexandrov, Y.B. Yurov, E.I. Rogaev *et al.*, 2015 Annotation of
453 suprachromosomal families reveals uncommon types of alpha satellite organization in
454 pericentromeric regions of hg38 human genome assembly. *Genom Data* 5:139-146.

455 Simão, F.A., R.M. Waterhouse, P. Ioannidis, E.V. Kriventseva, and E.M. Zdobnov, 2015 BUSCO:
456 assessing genome assembly and annotation completeness with single-copy orthologs.
457 *Bioinformatics* 31 (19):3210-3212.

458 Stanke, M., A. Tzvetkova, and B. Morgenstern, 2006 AUGUSTUS at EGASP: using EST, protein
459 and genomic alignments for improved gene prediction in the human genome. *Genome*
460 *Biol* 7 Suppl 1:S11.11-18.

461 Traut, W., M. Szczepanowski, M. Vítková, C. Opitz, F. Marec *et al.*, 2007 The telomere repeat
462 motif of basal Metazoa. *Chromosome Res* 15 (3):371-382.

463 Vollger, M.R., G.A. Logsdon, P.A. Audano, A. Sulovari, D. Porubsky *et al.*, 2020 Improved
464 assembly and variant detection of a haploid human genome using single-molecule, high-
465 fidelity long reads. *Ann Hum Genet* 84 (2):125-140.

466 Wevrick, R., and H.F. Willard, 1989 Long-range organization of tandem arrays of alpha satellite
467 DNA at the centromeres of human chromosomes: high-frequency array-length
468 polymorphism and meiotic stability. *Proc Natl Acad Sci U S A* 86 (23):9394-9398.

469 Willard, H.F., 1985 Chromosome-specific organization of human alpha satellite DNA. *Am J Hum*
470 *Genet* 37 (3):524-532.

471 Willard, H.F., J.S. Wayne, M.H. Skolnick, C.E. Schwartz, V.E. Powers *et al.*, 1986 Detection of
472 restriction fragment length polymorphisms at the centromeres of human chromosomes
473 by using chromosome-specific alpha satellite DNA probes: implications for development
474 of centromere-based genetic linkage maps. *Proc Natl Acad Sci U S A* 83 (15):5611-5615.

475 Wolfe, J., S.M. Darling, R.P. Erickson, I.W. Craig, V.J. Buckle *et al.*, 1985 Isolation and
476 characterization of an alphoid centromeric repeat family from the human Y
477 chromosome. *J Mol Biol* 182 (4):477-485.

478 Wootton, R., 1976 *The Biology of Sticklebacks*. Academic Press.

479 Xu, G.C., T.J. Xu, R. Zhu, Y. Zhang, S.Q. Li *et al.*, 2019 LR_Gapcloser: a tiling path-based gap
480 closer that uses long reads to complete genome assembly. *Gigascience* 8 (1).

481





