

A high-quality reference genome for the common creek chub, *Semotilus atromaculatus*

Amanda V. Meuser^{1,2}, Amy R. Pitura^{1,2}, Elizabeth G. Mandeville¹

¹ Department of Integrative Biology, University of Guelph, Guelph, Ontario, Canada

² Co-first authors

Corresponding author: Amy Pitura

University of Guelph

50 Stone Road East

Guelph, Ontario N1G 2W1, Canada

apitura@uoguelph.ca

1-519-824-4120 ext. 52843

Running title: *Semotilus atromaculatus* reference genome

Keywords: *Semotilus atromaculatus*, reference genome, creek chub, synteny, cyprinid, leuciscid

1 Abstract

2 Creek chub (*Semotilus atromaculatus*) are a leuciscid minnow species commonly found in an-
3 thropogenically disturbed environments, making them an excellent model organism to study
4 human impacts on aquatic systems. Genomic resources for creek chub and other leuciscid
5 species are currently limited. However, advancements in DNA sequencing now allow us to
6 create genomic resources at a historically low cost. Here, we present a high quality 239 con-
7 tig reference genome for the common creek chub, created with PacBio HiFi sequencing. We
8 compared the assembly quality of two pipelines: Pacific Biosciences' Improved Phase Assem-
9 bly (IPA; 873 contigs) and Hifiasm (239 contigs). Quality and completeness of this genome
10 is comparable to the zebrafish (Danioninae) and fathead minnow (Leuciscidae) genomes.
11 The creek chub genome is highly syntenic to the zebrafish and fathead minnow genomes,
12 and while our assembly does not resolve into the expected 25 chromosomes, synteny with
13 zebrafish suggests that each creek chub chromosome is likely represented by 1-4 large contigs
14 in our assembly. This reference genome is a valuable resource that will enhance genomic bio-
15 diversity studies of creek chub and other non-model leuciscid species common to disturbed
16 environments.

17 Introduction

18 Genomic studies of non-model species have become increasingly feasible in the past two
19 decades (Narum *et al.* 2013, Lou *et al.* 2021). Efforts are now under way to sequence
20 the tree of life (Fan *et al.* 2020), including organisms of no known economic importance,
21 whose genomes are likely to be used primarily for conservation or evolutionary research.
22 While genomic studies of non-model organisms are proceeding at a rapid pace, progress is
23 still limited by the lack of suitable reference genomes for many species and clades. Using
24 a reference genome from a closely related species is sometimes possible when there is no
25 available reference for focal taxa (e.g. Mandeville *et al.* 2019). However, this approach can
26 produce misleading results under some circumstances, including when the goal of a study is
27 to examine within-species differentiation or similar and species-specific variation may be lost
28 (as when wolf vs. dog reference genomes were used for a study of wolves; Gopalakrishnan *et al.*
29 2017). These concerns are especially relevant when considering clades where a substantial
30 amount of structural genomic variation exists.

31 Fish genomes are quite diverse, and vary in size from 0.5–2 pg (excluding polyploids;
32 Smith & Gregory 2009). On a locus-specific functional level, essential processes like sex de-
33 termination can have an incredibly diverse genetic basis in fish (Bachtrog *et al.* 2014, Pennell
34 *et al.* 2018). Within North American teleost fish, a large proportion of fish biodiversity is
35 encompassed by the family Leuciscidae within the order Cypriniformes (Holm *et al.* 2022,
36 Stout *et al.* 2016), but previous genomic work on leuciscid minnow species has been lim-
37 ited. However, being extremely numerous and geographically widespread, these species have
38 great potential for use as model species to study the effects of anthropogenic disturbance
39 and overall population genetic structure of stream fishes. Some species are quite tolerant,
40 and persist or even thrive in disturbed environments (Stammler *et al.* 2008). Additionally,
41 these species are known to hybridize from morphological data, but hybridization patterns
42 have not been described in detail using genetic or genomic data (Corush *et al.* 2021).

43 One major limitation for future work is that no leuciscid reference genomes have been
44 available until recently, and the highest quality available reference genome would be a ze-
45 brafish (*Danio rerio* genome), which is not very closely related to many wild leuciscid taxa of
46 interest (Fig. 1; Schönhuth *et al.* (2018)). A recently published fathead minnow (*Pimephales*
47 *promelas*) genome provides one resource in the leuciscid family (Martinson *et al.* 2022), but
48 there is still a dearth of genomic resources for this hyper-diverse and geographically ubiqui-
49 tous clade. In consequence, previous genomic studies of creek chub have relied on artificial
50 reference genomes or reference genomes from distantly related taxa (e.g. Meuser *et al.* 2022).

51 We sequenced the genome of the common creek chub, *Semotilus atromaculatus*, to pro-
52 vide a genomic resource for future studies of creek chub and other leuciscid minnows in North
53 America. We chose to sequence creek chub because of the broad range, abundant popula-
54 tion sizes, and generally tolerant life history of this species. Creek chub are expected to
55 have $2n=50-52$ chromosomes and a genome size of 1.25pg, similar to many leuciscid species
56 (Gold & Amemiya 1987, Legendre & Steven 1969). We also compared the outcomes of two
57 assembly pipelines: Pacific Biosciences' Improved Phase Assembly (IPA) HiFi Genome As-
58 ssembler pipeline (github.com/PacificBiosciences/pbipa) and the software HiFiasm (Cheng

59 *et al.* 2021). Finally, we assessed synteny between the creek chub reference genome and the
60 zebrafish and fathead minnow genomes (Martinson *et al.* 2022).

61 Methods

62 Sampling was accomplished under animal utilization protocol #4237 approved by the Univer-
63 sity of Guelph Animal Care Committee, with permits from the Ontario Ministry of Natural
64 Resources and Forestry (Licence No. 1100698), and with private landowner permission.
65 One wild-caught, 10.6 cm (total length) creek chub was sampled to generate this reference
66 genome (see Fig S1 for a photo). We sampled this fish using a beach seine from Swan Creek
67 in southern Ontario, Canada, in August 2022. The sampled individual was identified mor-
68 phologically by expert field personnel and later identified genetically as a creek chub with
69 DNA barcoding. Following capture, we euthanized the target individual with an overdose
70 of MS-222, then sampled and flash froze muscle tissue in liquid nitrogen within 5 minutes of
71 euthanasia to ensure preservation of high molecular weight DNA. We stored the flash-frozen
72 muscle and remainder of the specimen in a -80°C freezer, except for 2 fin clips preserved
73 in 95% ethanol. We extracted DNA from the fin clips using a DNeasy Blood & Tissue
74 kit (Qiagen) and quantified the concentration using a NanoDrop 8000 Spectrophotometer
75 (Thermo Scientific). We used this DNA to verify our phenotypic identification with DNA
76 barcoding at the University of Guelph’s Advanced Analysis Center. The COI-3 region of
77 the mitochondrial genome was amplified and sequenced using thermalcycler conditions and
78 primers from Ivanova *et al.* (2007). The forward fasta sequence was input on BOLD’s Identi-
79 fication Engine (www.barcodinglife.org/index.php/IDS_OpenIdEngine; (Ratnasingham
80 & Hebert 2007)) and confirmed to belong to creek chub.

81 We sent flash frozen muscle tissue to the University of Delaware’s DNA Sequencing
82 & Genotyping Center, in Newark, Delaware, USA. High molecular weight (HMW) DNA
83 extraction was completed using the MagAttract HMW DNA kit (Qiagen), then the extracted
84 DNA was quantified using a Qubit Fluorimeter and DNA fragment sizes were assessed by
85 Femto Pulse system instrument (Agilent). Next, a Megaruptor 2 (Diagenode) was used to
86 shear 3 μg of DNA to 15kb fragments. Then, a SMRTbell DNA library was constructed
87 according to the Pacbio HiFi SMRTbell protocol using SMRTbell Express Template Prep
88 Kit 3.0 (Pacbio, 102-182-700). After BluePippin size selection (Sage Science, PAC20KB)
89 removed fragments smaller than 8 kb, the average size in the library was 18 kb based on
90 Femto Pulse System (Agilent) analysis. Finally, sequencing was performed on 2 SMRT 8m
91 cells on Sequel IIe instrument with 30 hours movie, using both the Sequel II Binding kit 2.2
92 and Sequel II Sequencing kit 2.0.

93 The initial assembly of the reference genome was performed by the University of Delaware’s
94 DNA Sequencing & Genotyping Center. They used Pacific Biosciences’ Improved Phase As-
95 sembly (IPA) HiFi Genome Assembler pipeline (github.com/PacificBiosciences/pbipa). In
96 addition, we used HiFiasm (Cheng *et al.* 2021) to create a genome assembly. This, and all
97 subsequent computation, was performed on Digital Research Alliance of Canada’s Cedar
98 high performance computing cluster. We ran HiFiasm (v0.16.1) with 32 CPUs to create a

99 HiFi-only assembly, as we did not have parental short reads or Hi-C reads to create either
100 the trio-binning or Hi-C integrated assemblies (Cheng *et al.* 2021).

101 We assessed genome assembly quality using custom R and shell scripts to quantify distri-
102 bution of assembled contig and scaffold lengths, and the number of unique assembled scaffolds
103 (Fig. 2). From these data, we calculated N50, N90, L50, and L90, as well as maximum,
104 mean, and median contig length (Table 1). As this assembly is comprised completely of long-
105 read PacBio data, there are no gaps in our assembled contigs, and we hereafter refer to these
106 fragments of the genome simply as contigs. We also ran this analysis on the most recent ver-
107 sions of the zebrafish (GCF_000002035.4_GRCz11) and fathead minnow (GCF_016745375,
108 Martinson *et al.* 2022). We assessed completeness of the creek chub reference genome using
109 BUSCO v5.2.2 with actinopterygii_odb10 as the database (Simão *et al.* 2015), as well as the
110 zebrafish and fathead minnow reference genomes for the sole purpose of comparison with the
111 same database. We used kraken 2 v2.1.2 (Wood *et al.* 2019) to assess contamination using a
112 custom database containing virus, plasmid, protozoa, archaea, bacteria, human, plant, and
113 fungi sequences.

114 We examined synteny between creek chub and zebrafish, using SynMap, from the plat-
115 form CoGe (Comparative Genomics, genomeevolution.org, Lyons & Freeling 2008). CoGe
116 DAGChainer outputs were used to create circular plots with *circos* (Krzywinski *et al.*
117 2009). We used a hard masked version of the genome, uploaded to CoGe (NCBI Window-
118 Masker (Hard) (v1.0,id65989;genomic). The exact zebrafish organism used was *Danio rerio*
119 (*zebrafish*;id43752) and the genome was *unmasked* (v11, id66058; CDS). This is the
120 most recent zebrafish genome (GRCz11) created by the Reference Genome Consortium, re-
121 leased May 9, 2017 (ncbi.nlm.nih.gov/assembly/GCA_000002035.4). All default analysis and
122 display options were used in the Legacy Version, with the exception of Syntenic Path Assem-
123 bly (SPA) being selected, contigs without synteny hidden, diagonals coloured by syntenic
124 block, contigs sorted by name, and minimum chromosome size set to 2,830,400bp, which
125 is the length of the 50th largest contig in the creek chub assembly. This SynMap analysis
126 can be generated at any time at this link: genomeevolution.org/r/1oxpo. We additionally
127 created a SynMap between only the 25 largest creek chub contigs and the zebrafish genome,
128 by setting the minimum chromosome length to that of the 25 largest contig in the assembly
129 (20,130,130bp, Fig S2). It can be viewed at this link: genomeevolution.org/r/1oxpw

130 We also used SynMap to assess synteny between creek chub and fathead minnow. The
131 zebrafish genome is more complete than the fathead minnow genome (Martinson *et al.* 2022).
132 However, creek chub are more closely related to fathead minnow than to zebrafish (Fig.
133 1). The version of the genome used was *unmasked* (v2,id66042;CDS) of GCA_016745375,
134 recently published by Martinson *et al.* (2022). We used the same analysis and display options
135 as mentioned above for the zebrafish SynMap. This SynMap analysis can be regenerated
136 by following this link: genomeevolution.org/r/1oxpx. We also created a SynMap with the
137 25 largest creek chub contigs and the fathead minnow genome (Fig S3). It can be viewed
138 at this link: genomeevolution.org/r/1oxq3. SynMap labels are based on the creek chub
139 and fathead minnow genomes' contig/scaffold codes from the FASTA headers and many of
140 these codes do not intuitively match the contig/scaffold's corresponding number. See Tables
141 S1 and S2 for a breakdown of the creek chub and fathead minnow genome's contig/scaffold
142 codes and corresponding contig/scaffold numbers.

143 We created a simple phylogeny to display the relationship between the creek chub, ze-
144 brafish, fathead minnow, and several other model teleost fish (Fig. 1). We created this
145 phylogeny using the R package *fishtree* (Chang *et al.* 2019), which pulls phylogenetic data
146 from its pre-assembled online database. We added the fish photo with Adobe Photoshop
147 (v22.0.0).

148 Results

149 PacBio HiFi sequencing on two SMRT cells produced 133GB of raw data in FASTQ for-
150 mat. This corresponded to 4,313,794 raw reads with a mean length of 16,406 base pairs,
151 corresponding to a coverage of 64x. An initial genome assembly was constructed by the
152 University of Delaware sequencing facility's bioinformatics team using Pacific Biosciences's
153 IPA pipeline, and resulted in an assembly that consisted of 873 contigs, with mean contig
154 length of 1,257,076, and an N50 of 5,722,762 (Table 1). We then improved upon this initial
155 assembly using the HiFiasm pipeline (Cheng *et al.* 2021), resulting in an assembly with 239
156 contigs, with a mean contig length of 4,599,676, and an N50 of 30,568,897, which is halfway
157 between the N50 of the zebrafish and fathead minnow genomes (Table 1). BUSCO analy-
158 sis indicates that in addition to being highly contiguous, this genome is largely complete,
159 with a score of 98.0% and 97.9% respectively for the HiFiasm and IPA assemblies (Table
160 2). BUSCO values were similar between the two assemblies with the exception of a higher
161 proportion of genes designated as complete and duplicated in the HiFiasm assembly relative
162 to the IPA assembly (2.5% versus 1.6%, respectively), however both values are similar to the
163 fathead minnow and lower than the zebrafish (Table. 2). Kraken2 analysis did not find any
164 contigs to be entirely contaminated with non-creek chub DNA. As the HiFiasm assembly
165 was comparable to or improved over the IPA assembly in all respects - high completeness
166 and low contamination, but fewer contigs, higher N50, and larger mean contig size - we used
167 the HiFiasm assembly for all subsequent analyses.

168 Comparative Genomics (CoGe)'s SynMap (Lyons & Freeling 2008) analysis produced
169 2274 syntenic blocks and 24532 syntenic matches with zebrafish. We see few major chromo-
170 somal rearrangements in creek chub relative to zebrafish (Fig. 3). The haploid chromosome
171 number is expected to be the same ($n=25$) for zebrafish, fathead minnow, and creek chub
172 (Gold & Amemiya 1987), and most contigs in our assembly corresponded in part or whole
173 to zebrafish chromosomes (Fig. 4). While our assembly is less contiguous than the zebrafish
174 genome, in many cases zebrafish chromosomes map to 1–4 larger assembled contigs of the
175 creek chub genome (Fig. 4) and the 50 largest contigs in the creek chub assembly contain
176 just over 95% of the total assembly content (Table 1).

177 Although the large scale pattern is synteny with zebrafish, there are a few regions of
178 the creek chub genome that appear more sharply divergent. In particular, the creek chub
179 contig 9 (the largest complete contig) showed synteny with both chromosomes 4 and 7 in
180 the zebrafish genome (Fig. 4), suggesting there may have been an interchromosomal fusion
181 event in creek chub. Contig 9 also shows two small inversions on the zebrafish chromosome
182 7 (Fig. 3), indicating intrachromosomal rearrangements. However, as Fig. 3 is coloured by
183 syntenic block, it is obvious that there have been many minor chromosomal rearrangements

184 or mutations. There are no stretches of synteny along any chromosome or contig that are
185 greater than a few dots – with each dot representing a window of 20 genes in which at least
186 5 genes are syntenic between species. Quantified in a different way, no stretches of synteny
187 along any creek chub contig are greater than 12000 nucleotides, with the average being 407
188 nucleotides.

189 In our SynMap analysis between creek chub and fathead minnow (Fig. 5), which produced
190 19230 syntenic matches in 2042 blocks, we can see that there have also been few major
191 chromosomal rearrangements since these species diverged. One obvious rearrangement is
192 potentially fission or fusion events, whereby scaffold 1 of the fathead minnow genome is
193 split between contigs 28 and 42 of the creek chub genome and scaffold 2 is split between
194 contigs 3, 36, and 44 (Fig. 6). The syntenic matches are less continuous and consistent
195 with some of the fathead minnow scaffolds when compared to zebrafish (Fig. 6 versus Fig.
196 4). However, this likely reflects the quality of the fathead minnow and creek chub genomes
197 compared to the zebrafish genome, rather than a closer phylogenetic relationship between
198 zebrafish and creek chub than fathead minnow and creek chub. Indeed, creek chub and
199 fathead minnow are more closely related to one another than to zebrafish (Fig. 1). The
200 three species are contained within the order Cypriniformes, with zebrafish in the family
201 Danioninae and fathead minnow and creek chub in the family Leuciscidae (Schönhuth *et al.*
202 2018, Stout *et al.* 2016). The increased number of syntenic matches over numerous different
203 contigs of each species (Fig. 6), as opposed to contained between a few as we see with
204 zebrafish and creek chub (Fig. 4), are most likely due to the lower continuity and quality of
205 annotation of the fathead minnow genome compared to the zebrafish genome. CoGe predicts
206 syntenic genes based off of sequence similarity, but with a lower quality annotation, is more
207 likely to identify transposable elements or repetitive regions as syntenic between genomes,
208 increasing background noise in the SynMap (Lyons & Freeling 2008).

209 Discussion

210 Our de novo sequencing approach relied entirely on PacBio data, which allowed us to suc-
211 cessfully assemble sequence data into a relatively small number of longer contigs (n=239 for
212 the HiFiasm assembly; Table 1). While this assembly is not quite chromosome scale, as the
213 expected haploid chromosome number is 25, there are larger scaffolds which likely approach
214 full chromosomes (Fig. 2), and synteny analyses with zebrafish suggest that each creek chub
215 chromosome is likely covered by 1–4 large contigs (Fig. 4). Analyses of completeness with
216 BUSCO confirm that a high proportion of expected genes are included (about 98% for both
217 assemblies), reinforcing that sequencing produced a high quality reference genome. The con-
218 tiguity and completeness of this assembly makes it a valuable resource for genomic studies
219 of non-model leuciscid fish. The high contiguity of sequence enabled by PacBio will allow
220 recovery of genetic architecture of traits where linkage of multiple loci might be extremely
221 relevant (for example, examining the genetic basis of sex determination; Meuser *et al.* 2022).
222 A high quality reference genome will also enable analyses that require whole-genome data,
223 such as identifying inversions between closely related species (Faria *et al.* 2019), or demo-
224 graphic inference (MSMC and PSMC; SFS; ABC; Li & Durbin 2011, Schiffels & Durbin

225 2014, Beichman *et al.* 2018).

226 In the interest of constructing the most complete and continuous assembly possible from
227 our data, we used two different assembly pipelines, IPA and HiFiasm (Cheng *et al.* 2021).
228 Using HiFiasm, we successfully reduced the number of contigs from 873 to 239, and increased
229 the N50 roughly six-fold (Table 1, Fig. 2). Much of the improvement in N50 and contig
230 number likely resulted from the linking of multiple long contigs to form contigs that approach
231 chromosome length, which enabled better understanding of synteny with other related species
232 (Fig. 4 and Fig. 6).

233 As expected, much of the creek chub genome is syntenic with previously published
234 genomes of model organisms, namely zebrafish. However, there are also some rearrange-
235 ments, including a number of inversions and regions that are not syntenic with the zebrafish
236 genome. We do not yet know what functions are encoded by those particular regions of
237 the genome, but structural genome changes, especially of the sex determining region, are
238 likely to play a major role in diversification of species-rich clades of fish like the Cyprini-
239 formes (Payseur *et al.* 2018, Huang *et al.* 2020). One particular region to note in the synteny
240 analysis was that approximately half of zebrafish chromosome 4 was not conserved between
241 species. This region is a sex determining region in zebrafish, which has been shown to exist
242 in wild – but not lab raised – strains of zebrafish (Wilson *et al.* 2014). Sex determination
243 systems vary widely across teleost fish species (Bachtrog *et al.* 2014, Pennell *et al.* 2018).
244 Creek chub are not known to have a large sex determining region (Meuser *et al.* 2022) or
245 heteromorphic sex chromosomes (Gold *et al.* 1979), which could be part of the reason why
246 there are no other large regions lacking synteny between the two genomes.

247 While our creek chub genome assembly does not quite have one large contig per chromo-
248 some, for each zebrafish chromosome there are 1–4 larger contigs in our genome assembly
249 that are highly syntenic and likely together comprise the creek chub chromosome (Fig. 4).
250 This tells us that our genome is nearly chromosome-resolution, less a few joins between
251 large contigs. This is especially apparent when comparing to the fathead minnow reference
252 genome; our creek chub reference genome has fewer and larger contigs than the fathead
253 minnow genome (Table 1, Fig. 6). While our creek chub genome is not yet annotated, it
254 is certainly nearly complete and of similar quality to other recently published fish genomes
255 (Martinson *et al.* 2022).

256 The high-quality creek chub reference genome presented in this paper will enable new
257 insights about the evolutionary history and genome function of leuciscid fish species. Initially,
258 we intend to use this reference genome to investigate the effects of anthropogenic disturbance
259 on a suite of leuciscid fish species. Creek chub and a number of closely related species are
260 widely distributed in North America, and are found in disturbed environments, which makes
261 them an ideal study species for assessing impacts of urbanization and agricultural land use
262 on fish species (similar to previous work in other taxa; Miles *et al.* 2019, Wei *et al.* 2021).
263 A future goal is to produce a genome annotation, which would allow analysis of functional
264 patterns of genomic variation and gene expression in a more meaningful way. More broadly,
265 we are now entering a new and exciting era for genomics of non-model organisms, when it is
266 possible to move beyond using genomes of model organisms as reference, and gain the more
267 fine-grain insights that can only be obtained with a conspecific or closely related reference
268 genome (Gopalakrishnan *et al.* 2017). Generating high quality reference genomes is essential

269 for quantifying genomic variation across the incredible biodiversity of fishes (Fan *et al.* 2020),
270 and will lead to new insights about the evolution of this species-rich group of vertebrates.

271 Acknowledgements

272 Computing was accomplished through an allocation from the Digital Research Alliance of
273 Canada to EGM. We would like to thank T. Frauley for assistance with fieldwork, B. Schultz
274 for productive discussions about analysis of synteny, and S.E. McFarlane for manuscript com-
275 ments and discussion of what should be in a genome paper. This manuscript was improved
276 by comments from the entire Mandeville lab at University of Guelph. We also thank B.
277 Kingham and O. Shevchenko of the University of Delaware DNA Sequencing & Genotyping
278 Center for coordinating the DNA extraction, library preparation, sequencing, and initial
279 assembly with the IPA pipeline. Finally, we would like to thank to E. Lyons and A. Nelson
280 from CoGe for assistance with creating the SynMap analyses, and the landowners of the
281 property bordering Swan Creek for allowing us access.

282 Author contributions

283 EGM, AVM, and ARP planned the project. ARP and AVM completed field sampling and
284 tissue dissections. AVM and ARP completed the analyses and made the figures, with assis-
285 tance from EGM. All authors contributed to writing and revising the manuscript.

286 Data Availability Statement

287 Supplemental files are available at FigShare. File S1 contains a photo of the creek chub used
288 to create the reference genome. File S2 contains a syntenic dot plot between zebrafish and
289 the creek chub assembly's largest 25 contigs. File S3 contains a syntenic dot plot between
290 fathead minnow and the creek chub assembly's largest 25 contigs. File S4 contains a table
291 of the creek chub assembly's contig headers, the associated contig number, and length of
292 the contig in base pairs. File S5 contains a table of the fathead minnow assembly's scaffold
293 code, the associated scaffold number, and length of the scaffold in base pairs. Upon the
294 acceptance of this manuscript, data and scripts used for analysis will be made publicly
295 available in Data Dryad. The genome is available on the NCBI genomes repository, under
296 accession number PRJNA994924. Custom scripts used in this work will be available on
297 Github: github.com/amanda-meuser/CreekChubGenome.

298 Conflict of Interest

299 The authors declare no conflict of interest.

300 **Funder Information**

301 This research was undertaken using a Resources for Research Groups (RRG) computing
302 allocation from the Digital Research Alliance of Canada. Sequencing for this project was
303 funded by the Canada First Research Excellence Fund, specifically, University of Guelph's
304 Food From Thought Research Support grant.

References

- 305
- 306 Bachtrog D, Mank JE, Peichel CL, *et al.* (2014) Sex Determination: Why So Many Ways of
307 Doing It? *PLoS Biology*, **12**, e1001899.
- 308 Beichman AC, Huerta-Sanchez E, Lohmueller KE (2018) Using Genomic Data to Infer His-
309 toric Population Dynamics of Nonmodel Organisms. *Annual Review of Ecology, Evolution,*
310 *and Systematics*, **49**, 433–456.
- 311 Chang J, Rabosky DL, Smith SA, Alfaro ME (2019) An R package and online resource for
312 macroevolutionary studies using the ray-finned fish tree of life. *Methods in Ecology and*
313 *Evolution*, **10**, 1118–1124.
- 314 Cheng H, Concepcion GT, Feng X, Zhang H, Li H (2021) Haplotype-resolved de novo as-
315 sembly using phased assembly graphs with hifiasm. *Nature Methods*, **18**, 170–175.
- 316 Corush JB, Fitzpatrick BM, Wolfe EL, Keck BP (2021) Breeding behaviour predicts patterns
317 of natural hybridization in North American minnows (Cyprinidae). *Journal of Evolution-*
318 *ary Biology*, **34**, 486–500.
- 319 Fan G, Song Y, Yang L, *et al.* (2020) Initial data release and announcement of the 10,000
320 Fish Genomes Project (Fish10K). *GigaScience*, **9**, giaa080.
- 321 Faria R, Johannesson K, Butlin RK, Westram AM (2019) Evolving Inversions. *Trends in*
322 *Ecology and Evolution*, **34**, 239–248, publisher: Elsevier Ltd.
- 323 Gold JR, Amemiya CT (1987) Genome size variation in North American minnows
324 (Cyprinidae). II. Variation among 20 species. *Genome / National Research Council*
325 *Canada*, **29**, 481–489.
- 326 Gold JR, Whitlock CW, Karel WJ, Barlow JA (1979) Cytogenetic studies in North American
327 minnows (Cyprinidae). VI. Karyotypes of thirteen species in the genus *Notropis*.: VI.
328 Karyotypes of thirteen species in the genus *Notropis*. *CYTOLOGIA*, **44**, 457–466.
- 329 Gopalakrishnan S, Samaniego Castruita JA, Sinding MHS, *et al.* (2017) The wolf reference
330 genome sequence (*Canis lupus lupus*) and its implications for *Canis* spp. population ge-
331 nomics. *BMC Genomics*, **18**, 495.
- 332 Holm E, Mandrak NE, BurrIDGE ME (2022) *A Field Guide to Freshwater Fishes of Ontario*.
333 3rd edn., Royal Ontario Museum Press.
- 334 Huang K, Andrew RL, Owens GL, Ostevik KL, Rieseberg LH (2020) Multiple chromoso-
335 mal inversions contribute to adaptive divergence of a dune sunflower ecotype. *Molecular*
336 *Ecology*, pp. 1–15.
- 337 Ivanova NV, Zemlak TS, Hanner RH, Hebert PDN (2007) Universal primer cocktails for fish
338 DNA barcoding: BARCODING. *Molecular Ecology Notes*, **7**, 544–548.
- 339 Krzywinski M, Schein J, Birol *et al.* (2009) Circos: An information aesthetic for comparative
340 genomics. *Genome Research*, **19**, 1639–1645.

- 341 Legendre P, Steven M D (1969) Denombrement des chromosomes chez quelques cyprins. *Le*
342 *Naturaliste Canadien*, **96**, 913–918.
- 343 Li H, Durbin R (2011) Inference of human population history from individual whole-genome
344 sequences. *Nature*, **475**, 493–496.
- 345 Lou RN, Jacobs A, Wilder AP, Therkildsen NO (2021) A beginner’s guide to low-coverage
346 whole genome sequencing for population genomics. *Molecular Ecology*, **30**, 5966–5993.
- 347 Lyons E, Freeling M (2008) How to usefully compare homologous plant genes and chromo-
348 somes as DNA sequences: How to usefully compare plant genomes. *The Plant Journal*,
349 **53**, 661–673.
- 350 Mandeville EG, Walters AW, Nordberg BJ, Higgins KH, Burckhardt JC, Wagner CE (2019)
351 Variable hybridization outcomes in trout are predicted by historical fish stocking and
352 environmental context. *Molecular Ecology*, **28**, 3738–3755.
- 353 Martinson JW, Bencic DC, Toth GP, *et al.* (2022) De Novo Assembly of the Nearly Com-
354 plete Fathead Minnow Reference Genome Reveals a Repetitive but Compact Genome.
355 *Environmental Toxicology and Chemistry*, **41**, 448–461.
- 356 Meuser AV, Pyne CB, Mandeville EG (2022) Limited evidence of a genetic basis for sex de-
357 termination in the common creek chub, *Semotilus atromaculatus*. *Journal of Evolutionary*
358 *Biology*, **35**, 1635–1645.
- 359 Miles LS, Rivkin LR, Johnson MTJ, Munshi-South J, Verrelli BC (2019) Gene flow and
360 genetic drift in urban environments. *Molecular Ecology*, **28**, 4138–4151.
- 361 Narum SR, Buerkle CA, Davey JW, Miller MR, Hohenlohe PA (2013) Genotyping-by-
362 sequencing in ecological and conservation genomics. *Molecular Ecology*, **22**, 2841–2847,
363 iSBN: 1365-294X _eprint: NIHMS150003.
- 364 Payseur BA, Presgraves DC, Filatov DA (2018) Introduction: Sex chromosomes and speci-
365 ation. *Molecular Ecology*, **27**, 3745–3748.
- 366 Pennell MW, Mank JE, Peichel CL (2018) Transitions in sex determination and sex chro-
367 mosomes across vertebrate species. *Molecular Ecology*, **27**, 3950–3963.
- 368 Ratnasingham S, Hebert PDN (2007) BARCODING: bold: The Barcode of Life Data System
369 (<http://www.barcodinglife.org>): BARCODING. *Molecular Ecology Notes*, **7**, 355–364.
- 370 Schiffels S, Durbin R (2014) Inferring human population size and separation history from
371 multiple genome sequences. *Nature Genetics*, **46**, 919–925.
- 372 Schönhuth S, Vukić J, Šanda R, Yang L, Mayden RL (2018) Phylogenetic relationships and
373 classification of the Holarctic family Leuciscidae (Cypriniformes: Cyprinoidei). *Molecular*
374 *Phylogenetics and Evolution*, **127**, 781–799.

- 375 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM (2015) BUSCO: as-
376 sessing genome assembly and annotation completeness with single-copy orthologs. *Bioin-*
377 *formatics*, **31**, 3210–3212.
- 378 Smith EM, Gregory TR (2009) Patterns of genome size diversity in the ray-finned fishes.
379 *Hydrobiologia*, **625**, 1–25.
- 380 Stammler KL, McLaughlin RL, Mandrak NE (2008) Streams modified for drainage provide
381 fish habitat in agricultural areas. *Canadian Journal of Fisheries and Aquatic Sciences*,
382 **65**, 509–522.
- 383 Stout CC, Tan M, Lemmon AR, Lemmon EM, Armbruster JW (2016) Resolving Cyprini-
384 formes relationships using an anchored enrichment approach. *BMC Evolutionary Biology*,
385 **16**, 244.
- 386 Wei X, Huang M, Yue Q, *et al.* (2021) Long-term urbanization impacts the eastern golden
387 frog (*Pelophylax plancyi*) in Shanghai City: Demographic history, genetic structure, and
388 implications for amphibian conservation in intensively urbanizing environments. *Evolu-*
389 *tionary Applications*, **14**, 117–135.
- 390 Wilson CA, High SK, McCluskey BM, *et al.* (2014) Wild Sex in Zebrafish: Loss of the
391 Natural Sex Determinant in Domesticated Strains. *Genetics*, **198**, 1291–1308.
- 392 Wood DE, Lu J, Langmead B (2019) Improved metagenomic analysis with Kraken 2. *Genome*
393 *Biology*, **20**, 257.

394 **Tables and Figures**

Genome Statistics	Creek Chub (HiFiasm)	Creek Chub (IPA)	Zebrafish (GRCz11)	Fathead Minnow (GCF_016745375)
Number of contigs	239	873	NA	NA
Number of scaffolds	NA	NA	25 (1897 unplaced)	910 *
Longest contig or scaffold (bp)	58,351,558	23,528,990	78,093,715	59,790,976
Mean contig or scaffold length (bp)	4,599,676	1,257,076	873,221	1,170,614
Median contig/scaffold length (bp)	119,920	206,534	146,921	47,256
N50	30,568,897	5,722,762	52,186,027	11,952,773 *
N90	6,569,117	798,807	339,135	1,205,132
L50	15	49	14	23
L90	39	255	405	126
Percent of total genome assembly in 50 largest contigs	95.13	50.77	93.90	73.18
Percent of total genome assembly in 25 largest contigs	74.67	32.49	92.69	54.84

Table 1: Genome statistics for both the HiFiasm and IPA genome assemblies and the most recent versions of the zebrafish and fathead minnow genomes. Statistics for each assembly were generated using a custom script written in a combination of both shell and R. NA = not applicable. * Denotes fathead minnow statistics from Martinson *et al.* (2022).

BUSCO v5.2.2 (actinopterygii_odb10)	Creek Chub (HiFiasm)	Creek Chub (IPA)	Zebrafish (GRCz11)	Fathead Minnow (GCF_016745375)
Complete	3566 (98.0%)	3562 (97.9%)	3483 (95.6%)	3524 (96.9%)
Complete and single-copy	3476 (95.5%)	3505 (96.3%)	3434 (94.3%)	3431 (94.3%)
Complete and duplicated	90 (2.5%)	57 (1.6%)	49 (1.3%)	93 (2.6%)
Fragmented	25 (0.7%)	28 (0.8%)	57 (1.6%)	52 (1.4%)
Missing	49 (1.3%)	50 (1.3%)	100 (2.8%)	64 (1.7%)

Table 2: BUSCO (benchmarking universal single-copy orthologs) scores for both the HiFiasm and IPA genome assemblies and the most recent version of the zebrafish and fathead minnow genomes. Generated using BUSCO v5.2.2 (database: actinopterygii_odb10). Total number of BUSCO groups searched for each genome: 3640.

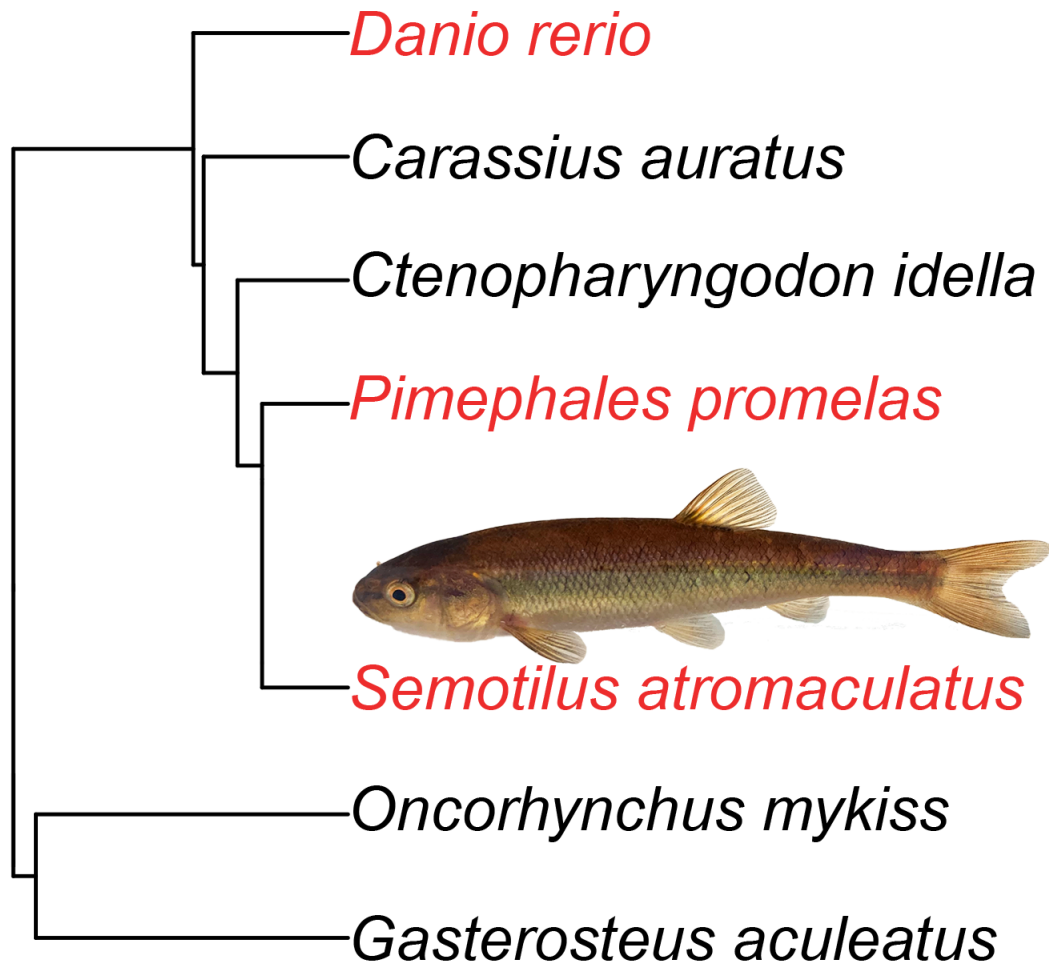


Figure 1: Phylogeny showing the relationship between zebrafish, fathead minnow, creek chub, and other fish commonly used as model species. The inset photo shows a creek chub individual. The phylogeny was created using data from the fishtree package in RStudio (Chang *et al.* 2019).

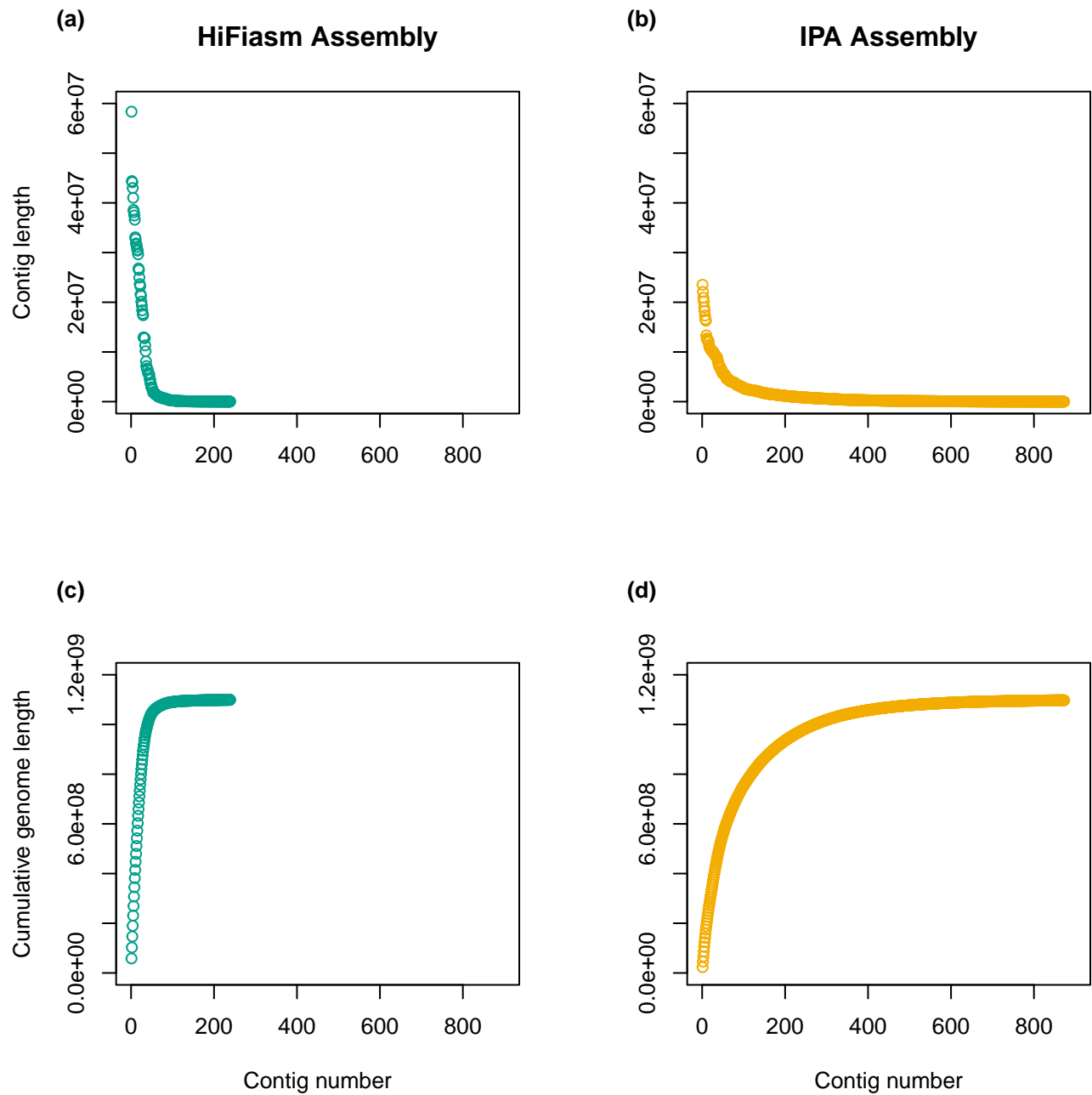


Figure 2: Visual comparison of the number of contigs, contig lengths, and cumulative genome lengths for both the HiFiasm and IPA genome assemblies. (a) Length of each contig in the HiFiasm assembly. (b) Length of each contig in the IPA assembly. (c) Cumulative genome length of HiFiasm assembly. (d) Cumulative genome length of IPA assembly.

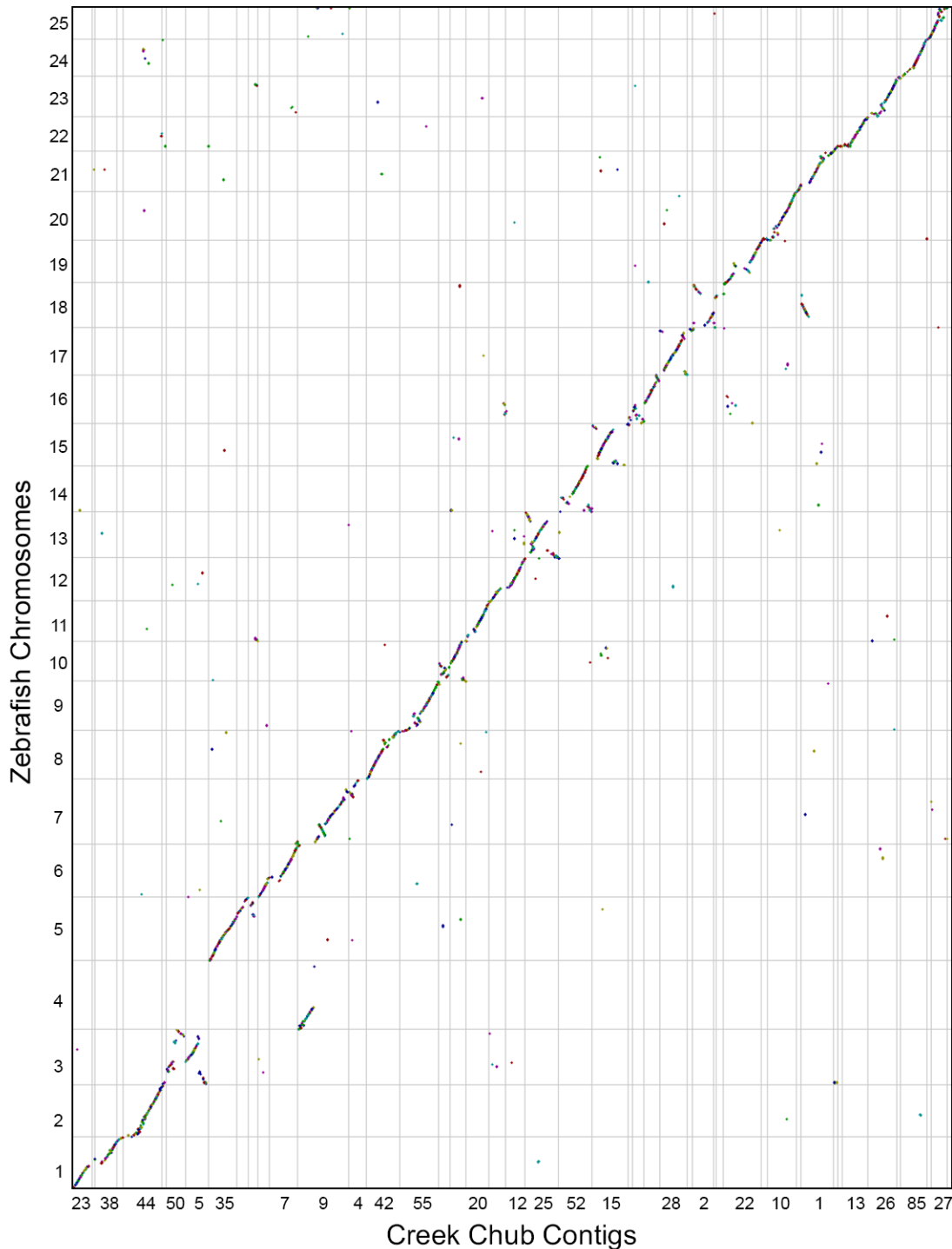


Figure 3: Dot plot showing synteny between creek chub (x-axis) and zebrafish (y-axis). All 25 zebrafish chromosomes from the GRCz11 version of the genome are present, while only the 50 largest contigs from the creek chub have been displayed, by setting the minimum contig length to 2,830,400 base pairs. The dot plot was made using CoGe's SynMap (Lyons & Freeling 2008). Each colour represents a different syntenic block. The figure can be regenerated at any time by following this link: genomevolution.org/r/1oxpo

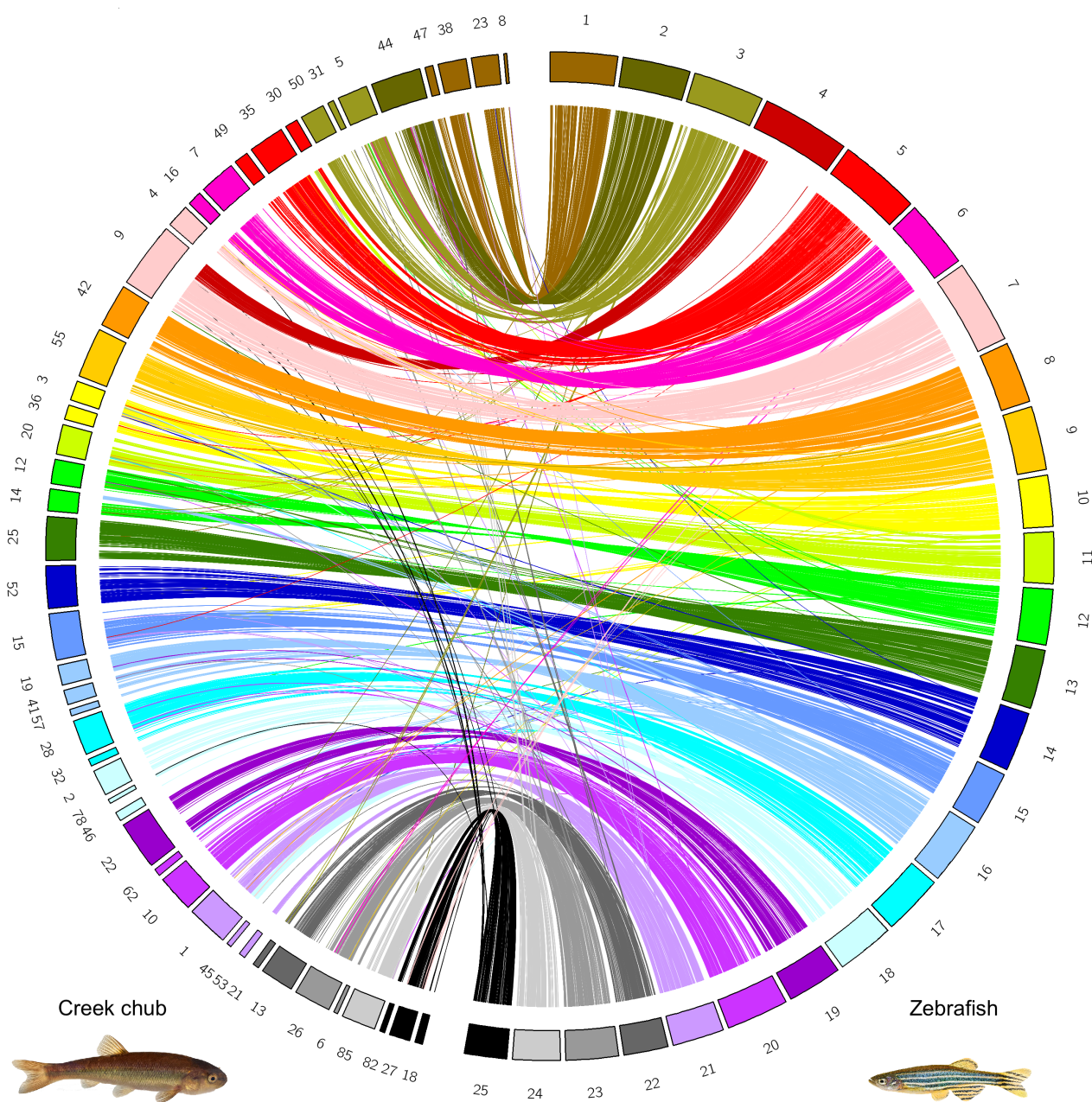


Figure 4: circos plot of syntenic matches between creek chub (left) and zebrafish (right). Creek chub contigs are coloured and mapped to reflect the zebrafish chromosome it has the majority of syntenic matches with. Zebrafish photo credit: Mirko_Rosenau.

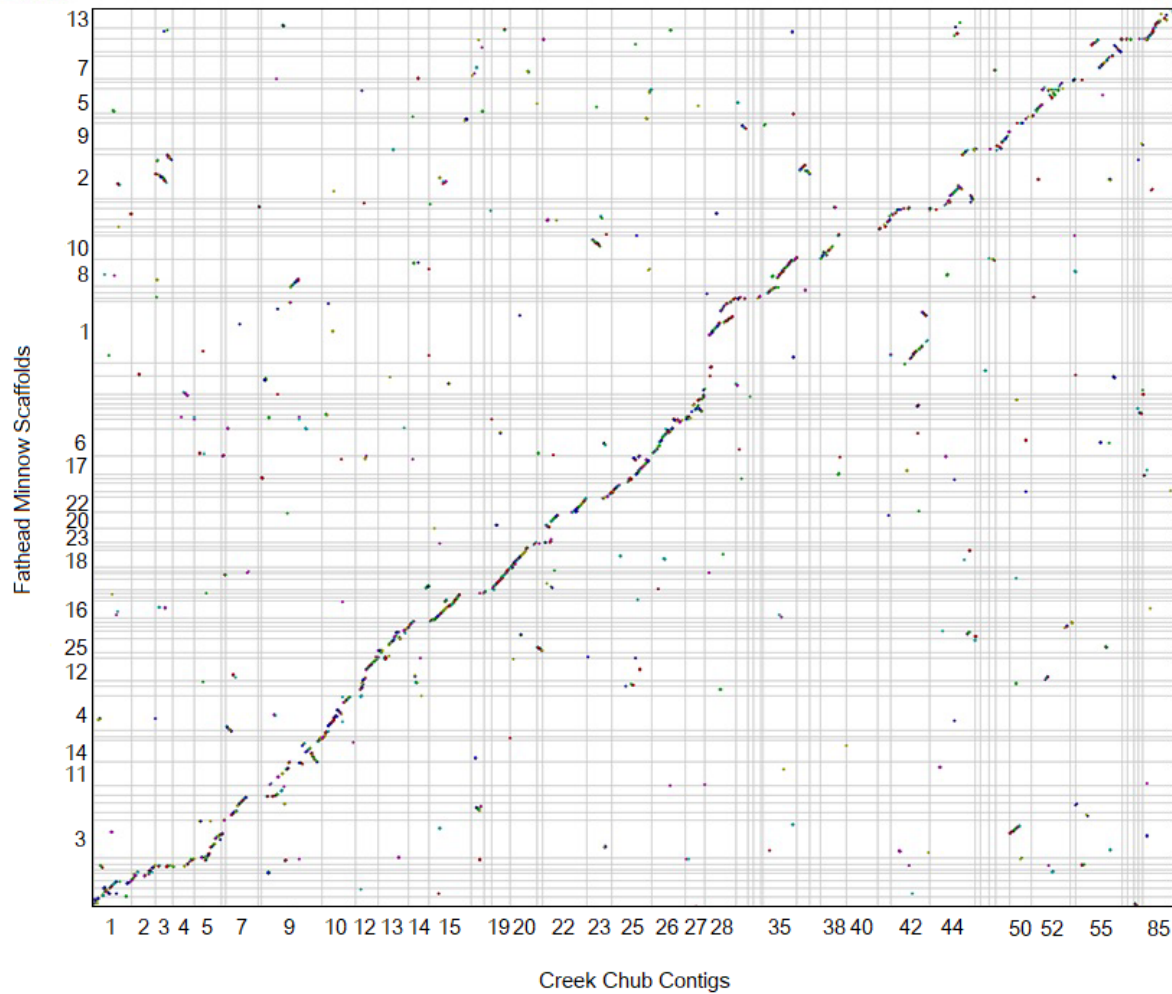


Figure 5: Dot plot made using CoGe's SynMap (Lyons & Freeling 2008) showing synteny between creek chub (x-axis) and fathead minnow (y-axis). Only the 50 largest contigs from the creek chub genome have been displayed, by setting the minimum chromosomes length to 2,830,400 base pairs. Each colour represents a different syntenic block. The figure can be regenerated at any time by following this link: genomevolution.org/r/1oxpx

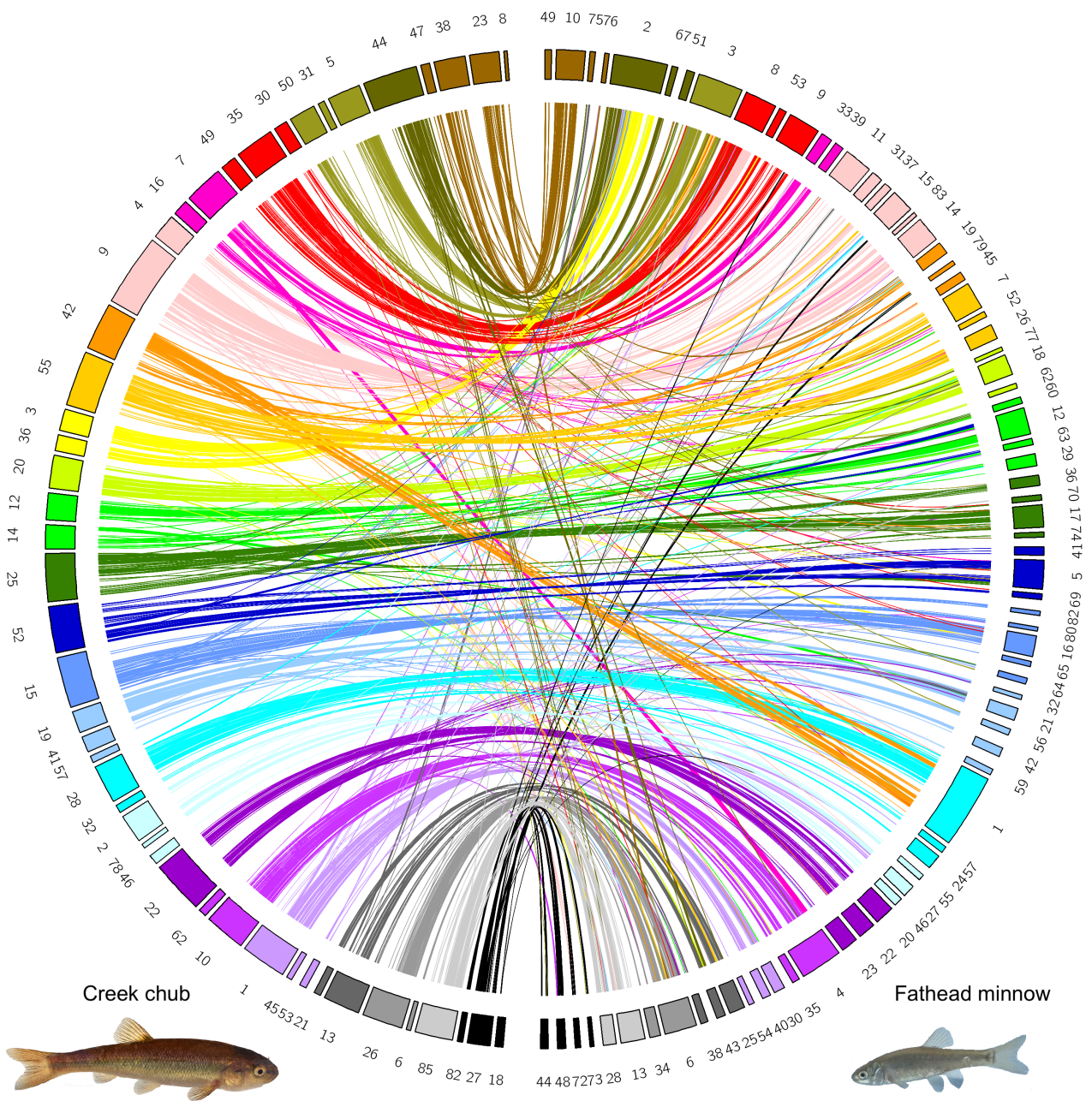


Figure 6: circos plot of syntenic matches between creek chub (left) and fathead minnow (right). Fathead minnow scaffolds are coloured and mapped to reflect the creek chub contig they have the majority of syntenic matches with.

395 Supplemental Figures



Figure S1: The creek chub individual used to create the reference genome. This fish was sampled from Swan Creek, Ontario, Canada. Note the dot present at the base of the dorsal fin and intermediate scale size compared to similar species. Not visible in the photo are the small barbels in the groove of each side of the mouth and minimally visible is the large mouth and black “moustache”.

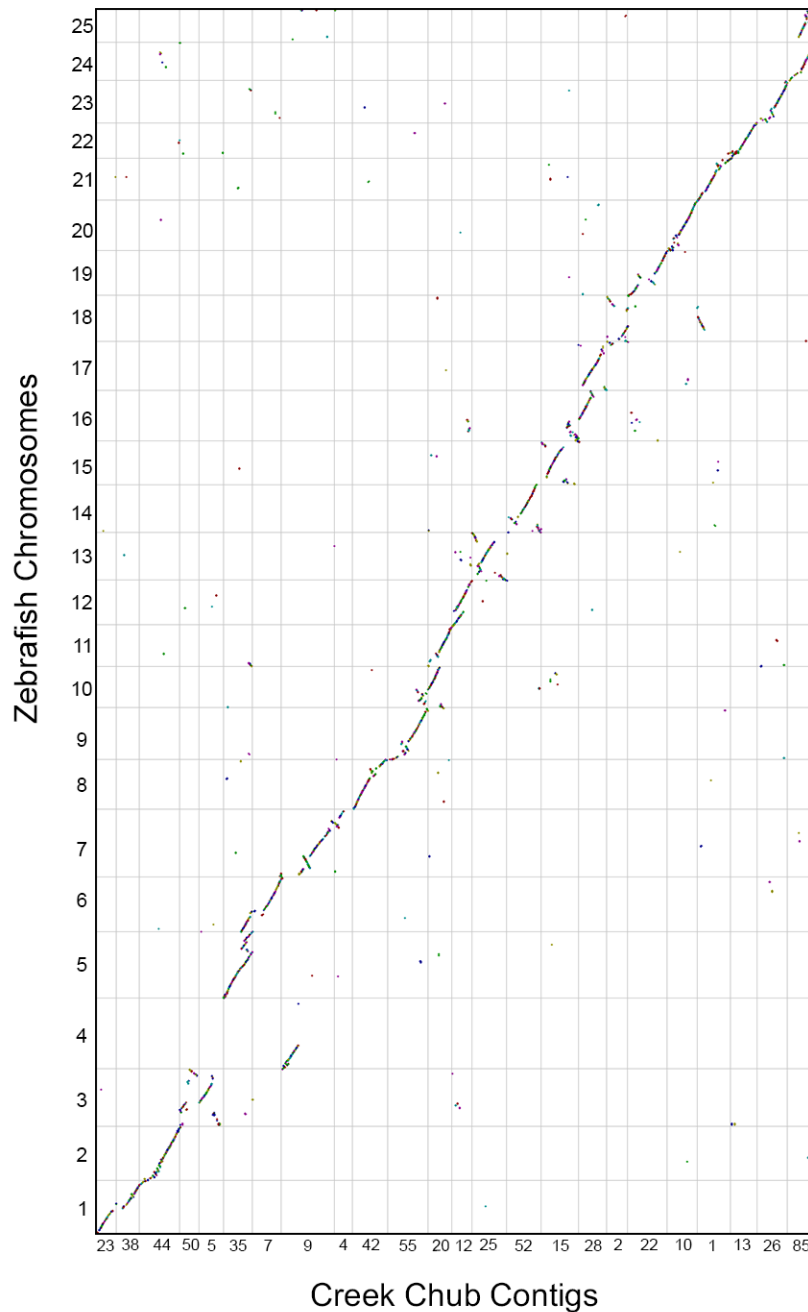


Figure S2: Dot plot made using CoGe's SynMap (Lyons & Freeling 2008) showing synteny between creek chub (x-axis) and zebrafish (y-axis). All 25 zebrafish chromosomes are present, while only the 25 largest contigs from the creek chub have been displayed, by setting the minimum chromosome length to 20,130,130 base pairs. Each colour represents a different syntenic block. The figure can be regenerated at any time by following this link: genomevolution.org/r/1oxpw

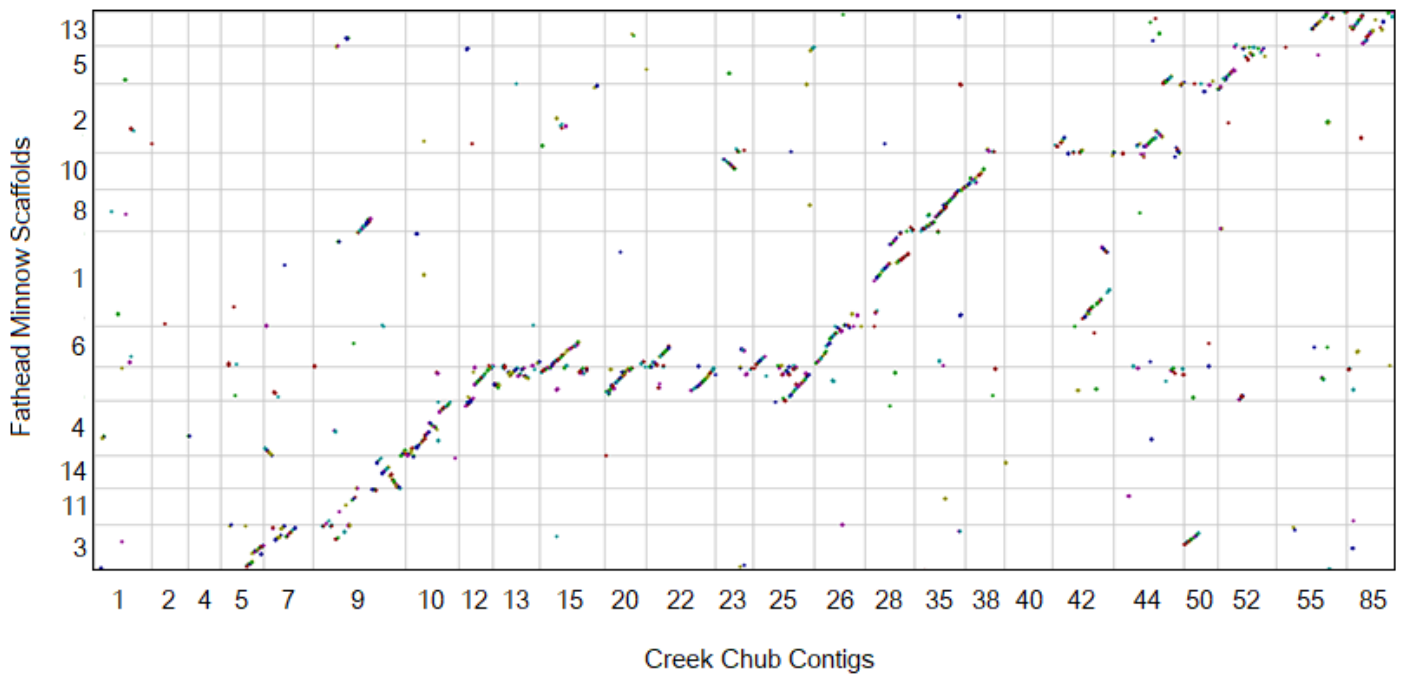


Figure S3: Dot plot made using CoGe's SynMap (Lyons & Freeling 2008) showing synteny between creek chub (x-axis) and fathead minnow (y-axis). Only the 25 largest contigs from the creek chub genome have been displayed, by setting the minimum contig length to 20,130,130 base pairs. Each colour represents a different syntenic block. The figure can be regenerated at any time by following this link: genomevolution.org/r/1oxq3