1 Astyanax mexicanus surface and cavefish chromosome-scale assemblies for trait

2 variation discovery

**Formatted:** Numbering: Continuous

3 Wesley C. Warren[1,2], Edward S. Rice[1], Maggs X[1], Emma Roback[3], Alex Keene[4], Fergal

4 Smith[5], Denye Ogeh[5], Leanne Haggerty[5], Rachel A. Carroll[1], Suzanne McGaugh[3],

5 Nicolas Rohner[6,7]

6

7 [1]Department of Animal Sciences, Department of Surgery, University of Missouri, Bond

8 Life Sciences Center, Columbia, MO

9 [2]Institute for Data Science and Informatics, University of Missouri, Columbia, MO

10 [3]Department of Ecology, Evolution, and Behavior, University of Minnesota, Saint Paul,

11 MN

12 [4]Department of Biology, Texas AM University, College Station, TX

13 [5]European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome

14 Genome Campus, Hinxton, Cambridge CB10 1SD, UK

15 [6]Stowers Institute for Medical Research, Kansas City, MO

16 [7]Department of Molecular and Integrative Physiology, KU Medical Center, Kansas City,

17 KS

18

19 Authors for correspondence: Wesley C. Warren (warrenwc@missouri.edu); Nicolas

20 Rohner (nro@stowers.org)

21 Keywords: *Astyanax mexicanus*, cavefish, chromosome assembly

22

23    **Abstract**

24    The ability of organisms to adapt to sudden extreme environmental changes produces

25    some of the most drastic examples of rapid phenotypic evolution. The Mexican Tetra,

26    *Astyanax mexicanus*, is abundant in the surface waters of northeastern Mexico, but

27    repeated colonizations of cave environments have resulted in the independent evolution

28    of troglomorphic phenotypes in several populations. Here, we present three chromosome-

29    scale assemblies of this species, for one surface and two cave populations, enabling the

30    first whole-genome comparisons between independently evolved cave populations to

31    evaluate the genetic basis for the evolution of adaptation to the cave environment.  Our

32    assemblies represent the highest quality of sequence completeness with predicted protein-

33    coding and non-coding gene metrics far surpassing prior resources and, to our

34    knowledge, all long-read assembled teleost genomes, including zebrafish. Whole genome

35    synteny alignments show highly conserved gene order among cave forms in contrast to a

36    higher number of chromosomal rearrangements when compared to other phylogenetically

37    close or distant teleost species. By phylogenetically assessing gene orthology across

38    distant branches of amniotes, we discover gene orthogroups unique to *A. mexicanus.*

39    When compared to a representative surface fish genome, we find a rich amount of

40    structural sequence diversity, defined here as the number and size of insertions and

41    deletions as well as expanding and contracting repeats across cave forms. These new

42    more complete genomic resources ensure higher trait resolution for comparative,

43    functional, developmental, and genetic studies of drastic trait differences within a species.

44

45    **Introduction**

46    Natural trait alterations in the Mexican tetra cavefish *Astyanax mexicanus*, some with

47    extreme phenotypic consequences, have offered insight into how genetic adaptation

48    works (ROHNER 2018) (MCGAUGH *et al.* 2020). This species exists in two distinct forms:

49    a traditional river-dweller and a blind, depigmented cave-dweller. Many studies have

50    highlighted the use of the *A. mexicanus* surface fish to cave-life as an adaptation model to

51    investigate the unique molecular mechanisms underlying various disease traits, such as

52    obesity, sleep disorders, diabetes, heart regeneration, and many others, which has

53    elevated its importance as an evolutionary model and led to its rapidly expanding use

54  (DUBOUE et al. 2011; ASPIRAS et al. 2015) (OJHA AND WATVE 2018) (STOCKDALE et al.

55  2018) (BILANDŽIJA 2019). Furthermore, we envision its wider use with the expansion of

56  *A. mexicanus* genetic resources by taking advantage of a relatively short evolutionary

57  time of transition, ~150,000 years from surface to cave and clearly demarcated

58  phenotypic differences (HERMAN et al. 2018).

59  High-quality and nearly complete chromosomes of many species are becoming

60  more readily available, but challenges remain for genomes with difficult structural

61  features, for example, the highly repetitive zebrafish genome (CHERNYAVSKAYA et al.

62  2022). Further research has refined best practices across broad phyla to resolve

63  difficulties in building multiple genome assemblies with near gap-free representation

64  (JARVIS et al. 2022). To date, few aquatic species have reached these higher levels of

65  contiguity and representation, but recently developed long-read hybrid methods improve

66  de novo assembly (RAUTIAINEN et al. 2023). Many recently published aquatic genomes,

67  although significantly more complete in sequence representation, have relied on error-

68  prone long-read reads and their correction with short-reads (MOORE et al. 2023)

69  (ROBERTS et al. 2023). This new generation of genome assemblies has significantly

70  advanced our knowledge of historically underrepresented sequences, in particular sex

71  chromosomes (IMARAZENE et al. 2021) (DU et al. 2022).

72  In cavefish, a recent long-read based assembly improved identification of

73  candidate genes underlying QTLs, discovered fixed or variable deletions in each cave

74  form, and guided gene editing experiments (WARREN et al. 2021). However, this

75  assembly of a surface-dwelling individual suffered from a high level of gene

76  fragmentation due to the use of high-error long reads. An individual from the Pachón

77  cave-dwelling population was recently assembled with high levels of contiguity and

78  facilitated the discovery of novel sex chromosome origins (IMARAZENE et al. 2021).

79  However, given the high phenotypic divergence of the various populations, references of

80  multiple populations are necessary to best perform comparative genomics between

81  surface and cave morphs. To this end, we present here three highly continuous

82  chromosome-scale assemblies of individuals from the Río Choy surface population, as

83  well as the Molino and Tinaja cave populations to complement Pachón and for the first

84  time present more complete references that span demographically unique and

85   independent cave populations.  The generation of these resources in *A. mexicanus* will

86   support the community's use of this aquatic model of human disease, particularly those

87   wishing to compare trait diversity in the well-studied zebrafish model.

88

89   **Methods**

90   *Sequencing and assembly*. The *A. mexicanus* DNA samples were obtained from female

91   fishes of surface, Molino, and Tinaja origins (Fig. 1A) reared at the Stowers Institute

92   aquatic facility under IACUC approved protocol (Protocol ID: 2021-126). High

93   molecular weight DNA was isolated using the salting out method described in the 10X

94   genomics demonstrated protocol (10X Genomics; Pleasanton, CA) from muscle tissue to

95   generate single molecular real time (SMRT) sequences using HiFi sequence base calling

96   output mode from the Sequel II instrument (Pacific Biosciences) according to the

97   manufacturer's protocols. From all SMRT sequences, total HiFi coverage ranged from

98   27-46x using a genome size estimate of 1.4Gb with average read lengths across all

99   samples of 21kb (Suppl. Table 1).

100      We assembled these reads into contigs using hifiasm v0.13-r308 with default

101   options (CHENG *et al.* 2021). To look for haplotigs, we aligned CLR (Circular Long

102   Reads) reads to the contigs using minimap v2.20 with arguments "-ax map-pb". We ran

103   purge_haplotigs v1.1.1 (GUAN *et al.* 2020) on all three assemblies; the Molino and Tinaja

104   assemblies had detectable signals of haplotigs but surface did not, so we removed

105   haplotigs from Molino and Tinaja using coverage cutoffs (100, 197, 250) and (30, 193,

106   250), respectively. Haplotigs represent assembled sequences that are essentially a copy of

107   the other haplotype and thus are removed to avoid redundancy in genome representation.

108   To perform scaffolding a sibling of the reference fish was used to generate and sequence

109   a HiC library developed with the Proximo Hi-C kit (Phase Genomics: Seattle, WA)

110   according to the manufacturers protocol. We next assembled the contigs into scaffolds

111   using a custom pipeline (see Code Availability) that aligns the Hi-C reads to the contigs

112   with bwa mem v0.7.17 (LI AND DURBIN 2009), postprocesses the alignments, and finally

113   scaffolds with SALSA v2.2 (GHURYE *et al.* 2019) and Juicebox (version 1.11.08)

114   (DURAND *et al.* 2016), as previously described (GHURYE AND POP 2019). We curated the

115   scaffolds using a combination of synteny and manual examination of the Hi-C heatmaps

116     (ROBINSON *et al.* 2018; RHIE *et al.* 2021; WARRENLAB. 2022). Agptools was used to

117     finalize chromosome assignments to be consistent with chromosome nomenclature of the

118     surface assembled genome reported in Warren et al (WARREN *et al.* 2021).

119     Chromosomes were numbered by aligned synteny to assembled chromosomes of the

120     previously assembled surface form of *A. mexicanus* (WARREN *et al.* 2021). The

121     completeness of gene representation was assessed using BUSCO v4.1.2

122     (Actinopterygii_odb10) (SIMAO *et al.* 2015).

123

124     *Whole genome interspecies and intraspecies analysis*. To better understand shared gene

125     organization with other teleost experimental models and across cave morphs, we first

126     built and visualized syntenic orthology of surface fish chromosomes (AstMex3) to widely

127     used fish experimental models including zebrafish (GRCz11), medaka (ASM223467v1),

128     and platyfish (X_maculatus-5.0) chromosomes using the Genespace R package (LOVELL

129     *et al.* 2022). Using the parse annotation function of Genespace we generated the gene bed

130     file from the gff files of each species. After files were parsed, formatted correctly and all

131     headers matched in the protein and gene bed file, we initiated Genespace as described in

132     Lovell (LOVELL *et al.* 2022). We repeated this same Genespace workflow per morph

133     compared to our surface fish assembly to evaluate possible assembly errors or natural

134     structural variation (SV).

135

136     *Gene annotation*. Gene annotation of all of our assemblies, in addition to the previously

137     assembled Pachón (IMARAZENE *et al.* 2021), was carried out by the use of the

138     standardized Ensembl workflows (AKEN *et al.* 2016). However, only the surface genome

139     described in this study was also annotated with the NCBI pipeline (PRUITT *et al.* 2014).

140     Numerous RNAseq data sets that exist in the NCBI sequence archive for surface and

141     cavefish from different tissue sources were used to improve the accuracy of protein-

142     coding and non-coding gene model builds. A full accounting of all gene annotations was

143     generated using AGAT v1.0.0 (J. 2020). We used standard parsing tools to compare lists

144     of gene symbols in all four annotations to one another. Here, we assumed homology of

145     matching gene symbols and ignored duplicates.

146

147     *Gene orthology.* To compare gene orthology with other vertebrates, we ran OrthoFinder

148     version 2.5.4 (EMMS AND KELLY 2019) on fifteen species total (

149     Table 2). To mitigate the effect of multiple transcripts per gene, we used primary

150     transcripts (the longest version) only (Suppl. Table 2), per OrthoFinder recommendations

151     (EMMS AND KELLY 2019). MAFFT was used for multiple sequence alignment and

152     Fasttree for gene tree inference in the OrthoFinder analysis. We ran a secondary analysis

153     for comparison using IQtree for gene tree inference. We generated orthogroups across all

154     species which included protein sequences from primary transcripts of *A. mexicanus*

155     surface genome, zebrafish (*Danio rerio*), Japanese medaka (*Oryzias latipes*), platyfish

156     (*Xiphophorus maculatus*), spotted gar (*Lepisosteus oculatus*), eastern brown snake

157     (*Pseudonaja textilis*), common wall lizard (*Podarcis muralis*), chicken (*Gallus gallus*),

158     human (*Homo sapiens*), chimpanzee (*Pan troglodytes*), mouse (*Mus musculus*), rat

159     (*Rattus norvegicus*), gray short-tailed opossum (*Monodelphis domestica*), western clawed

160     frog (*Xenopus tropicalis*), and a tunicate, *Ciona intestinalis,* as the outgroup. These taxa

161     were chosen 1) because they all have available Ensembl annotations, thus easily

162     compatible with Orthofinder algorithms, and 2) establish an even phylogenetic sampling

163     across amniotes and anamniotes. In addition to the four teleosts, we include the spotted

164     gar because this species originated prior to the teleost whole genome duplication, making

165     it an important bridge for ortholog predictions between teleosts and other vertebrates

166     (BRAASCH *et al.* 2016). We collated orthogroups that contain *A. mexicanus* sequences

167     lacking gene symbols, so that gene symbol annotation can be inferred from that of other

168     members in the orthogroup.

169

170     *Structural variation analysis.* To initially estimate the distribution and number of SVs

171     present in these independently derived cave morph genomes we used Assemblytics 1.2.1

172     (NATTESTAD AND SCHATZ 2016), that estimates tandem or repeat sequence contractions

173     and expansions, as well as deletions or insertions when compared to the surface genome

174     (NATTESTAD AND SCHATZ 2016). This approach can reflect uniqueness to a cave morph

175     or similarity across morphs when compared to the surface genome. The nucmer

176     algorithm of MUMmer release 4.x (KURTZ *et al.* 2004) was run to align assemblies at

177     their contig level to minimize false positives (NATTESTAD AND SCHATZ 2016). The

178 resulting delta file was input into the Assemblytics browser and run using default

179 parameters: a unique contig sequence anchoring size of 10,000 bp, and the variants

180 classified by type and size ranging from 50-500 bp and 500-10,000 bp for plot

181 visualization.

182

183 **Results and Discussion**

184 *Chromosome-scale assembly of the surface and cavefish forms*. We generated reference

185 genomes from single lab-reared *A. mexicanus* surface or cave morphs that are descended

186 from varying Mexico localities, representing independent origins of the cave form from

187 Tinaja and Molino populations as well as a surface Río Choy population. We sequenced

188 and assembled each genome using SMRT CCS (circular consensus sequencing) with a

189 genome coverage depth of 27.4 to 46.2 (RUAN AND LI 2020) (Suppl. Table 1). The final

190 *de novo* assemblies resulted in ungapped assembly length ranges of 1.36-1.41 Gb, total

191 contigs of 123-292, and N50 contig lengths of 12-47Mb (Table 1). Using on average 200

192 million 150bp Hi-C reads for chromosome-scale scaffolding resulted in the generation of

193 25 total chromosomes for the surface and each cave morph. Due to the exceptional

194 contiguity of the surface assembly (47 Mb N50 contig length), we established the

195 chromosome structure of this assembly first, then investigated discrepancies of the other

196 cavefish assemblies using multiple pairwise alignments, including the recently published

197 long-read assembly of the Pachón morph (IMARAZENE *et al.* 2021). Few order or

198 orientation errors were found and corrected; however, three chromosomes with complete

199 orientation differences in all cave forms must be further investigated for possible curation

200 fixes (Suppl. Fig. 1). The surface genome contiguity surpassed all our cave morphs and,

201 to our knowledge, all long-read assembled teleost genomes to date found in the NCBI

202 assembly archive, including the zebrafish (fDanRer4.1). Across all assemblies, only 3.7

203 to 7.4% of sequences could not be properly assigned to chromosomes. These new

204 assemblies contribute 90 Mb in new sequence on average and show a substantial

205 reduction (25-fold) in assembly gaps, when evaluated against the prior version of the

206 surface fish genome (A. mexicanus-2.0; Table 1). Our surface assembly exhibits a high

207 level of contiguity and completeness, despite that the total interspersed repeats for the

208 surface genome was 45.8%, which is only slightly lower than zebrafish's 48.4% with a

209     similar estimated genome size of 1.4 Gb.

210

211     *Structural differences*. One question we wished to address was: is gene order highly

212     conserved among surface and cave morph chromosomes despite their demographic origin

213     differences (Fig. 1A) and up to 190,000 generations of cave morph divergence from their

214     surface ancestor (MORAN R.L. 2022). We find no major chromosomal discrepancies in

215     chromosome order when performing pairwise comparisons of the Pachon, Tinaja, or

216     Molino cave morphs to the representative surface fish ancestor (Suppl. Fig. 1). Also, of

217     interest was: are there examples of unexpected teleost interspecies conservation in

218     chromosome gene order that can aid future studies aimed at understanding cave morph

219     standing genetic variation. The alignments of all 25 *A. mexicanus* chromosomes to three

220     distantly related teleost genomes (zebrafish, platyfish, and medaka) using aligned genes

221     are expected to display patterns of organizational divergence dependent on their

222     phylogenetic relationships. Overall, 56% of *A. mexicanus* chromosomes display complex

223     genomic rearrangements relative to all three distantly related teleost species (Fig. 1B). In

224     one example, the largest *A. mexicanus* chr1 (134 Mb), aligns with four zebrafish

225     chromosomes: 5, 16, 20, and 24 (Fig. 1C), suggesting multiple fissions or fusions

226     occurred throughout teleost genome evolution. In contrast, there were some

227     chromosomes that are nearly syntenic for gene order, such as *A. mexicanus* chr14 versus

228     zebrafish chr17 (Fig. 1C). A separate pairwise alignment of the surface fish and zebrafish

229     genomes (Danio rerio GRCz11) support these findings of variable synteny (Suppl. Fig.

230     2). These teleost interspecies alignments confirm earlier studies, and highlight the

231     complex trajectory in the evolution of the teleost genome (VOLFF 2005) (BRAASCH *et al.*

232     2016).

233

234     *Gene annotation*. We first assessed the sequence completeness of the assemblies using

235     benchmarking universal single-copy ortholog (BUSCO) (SIMAO *et al.* 2015) scores and

236     found on average 96.8% of BUSCOs present in their complete and unfragmented form,

237     2.4% missing, and 1.3% duplicated (Suppl. Table 3). Protein-coding genes from all four

238     genomes were predicted using Ensembl (AKEN *et al.* 2016) with the average being

239     26,974. The surface genome was also annotated using the NCBI workflow, with small

240    differences in the total protein-coding genes count (Suppl. Table 4). In contrast, large

241    differences were seen in the surface genome pseudogene annotation, 1,376 versus 183,

242    when comparing the NCBI and Ensembl output (Suppl. Table 4). Interestingly, a small

243    increase in Ensembl predicted pseudogenes in cave morphs compared to the surface is

244    observed (Suppl. Table 4) apart from Pachon. The correct annotation of pseudogenes

245    "non-functional" genes across species is a persistent problem, and especially deserves

246    further attention in cave species genomes given its importance in understanding

247    troglomorphic genetic adaptation (HARRISON 2021). While all these newly annotated *A.*

248    *mexicanus* genomes are improved over prior assemblies (WARREN *et al.* 2021), estimates

249    of non-coding genes were substantially improved (a 36-fold increase; Suppl. Table 4),

250    owing mostly to improved assembly contiguity and accuracy. In total, all these surface

251    and cave morph gene sets show Molino, Tinaja, and Surface have similar numbers of

252    mRNAs, exons, CDSs, and total CDS lengths, but in most comparisons, the increased

253    surface genome contiguity resulted in improved gene annotation, supporting the use of

254    AstMex3_surface as the standard reference for future *A. mexicanus* computational

255    experiments (Suppl. Table 5).

256        In pairwise comparisons, the surface genome had 1,041, 978, and 1,054 gene

257    symbols not found in Pachon, Molino, and Tinaja annotations, respectively (Suppl. Table

258    6). Similarly, each cave population had over 800 gene symbols not found in the surface

259    genome.

260

261    *Gene Orthology.* We targeted a specific span of model species across the vertebrate

262    phylogeny to classify potential orthologs (Fig. 2A). By analyzing all detected gene

263    ortholog relationships, we were able to better understand the utility of *A. mexicanus* for

264    comparative inference across aquatic models and other vertebrates (Fig. 2B). In the

265    OrthoFinder analysis, 95% of all genes were assigned to 20,612 orthogroups (Suppl.

266    Table 7; Fig. 2C). Fifty percent of genes assigned to orthogroups are in orthogroups of 19

267    genes or more, and over 25% of orthogroups contained all species (n = 5,285) (Fig. 2D),

268    of which 20.8% are single copy orthogroups (n = 1,097) (Suppl. Table 7; Fig. 2E). An

269    additional 14.4% of orthogroups (n = 2,988) contained at least one copy from every

270    vertebrate species and zero copies from the outgroup (Fig. 2B). All vertebrates had >96%

271 of their genes assigned to orthogroups (Fig. 2C; Suppl. Table 8). The outgroup, *C.*

272 *intestinalis*, had 66.8% of genes assigned to orthogroups (Suppl. Table 8). For *A.*

273 *mexicanus*, 96.3% of genes were assigned to orthogroups (Suppl. Table 8). We found

274 more one to one orthologs of *A. mexicanus* with humans than for zebrafish with humans

275 (n = 6,341 and n = 5,118, respectively; Fig. 2E). We identified 378 orthogroups between

276 *A. mexicanus* and humans that lack zebrafish orthologs (Fig. 2E). Overall, more

277 orthogroups contain genes from *A. mexicanus* and humans (n = 11,444) than zebrafish

278 and humans (n = 11,272). Addressing specifically orthologous relationships, 16,582

279 human genes (70.5%) have an ortholog with *A. mexicanus*, whereas 16,221 human genes

280 (68.9% ) have an ortholog with *Danio rerio*. These findings increase the breadth of

281 orthologs available in an alternative fish model for comparative trait dissection. Zebrafish

282 have more many-to-one and many-to-many orthologs with humans and more species-

283 specific orthogroups than *A. mexicanus* (n = 245 and n = 157, respectively; Fig. 2E). In

284 general, the availability of higher-quality genome assemblies is revealing teleost whole

285 genome duplication to be a pervasive source of genetic variation across these taxa

286 (ALBALAT AND CANESTRO 2016), and understanding the main evolutionary forces

287 impacting *A. mexicanus* gene losses or gains is beyond the scope of this study (ADRIAN-

288 KALCHHAUSER *et al.* 2020). In identifying orthologs, we also estimated the number of

289 duplications across the species tree, both at the tips and internal nodes which could

290 represent differential genome fractionation across lineages stemming from the teleost

291 whole genome duplication or lineage-specific duplications. At the base of teleosts, our

292 analysis identified 2,451 duplications. In the tip branches of teleosts, we identified 2,647,

293 2,825, 4,812, and 8,509 duplicates for *X. maculatus*, *O. latipes*, *A. mexicanus*, and *D.*

294 *rerio*, respectively. We suspect that teleost whole genome duplication, annotation bias,

295 phylogenetic sampling, and multiple evolutionary events or forces underpin these results.

296 Five orthogroups (OG0000022, OG0000023, OG0000045, OG0000050, and

297 OG0000086) each have over 100 genes from *A. mexicanus*. These uniquely large copy

298 numbers suggest that these genes experienced unique evolutionary pressures in *A.*

299 *mexicanus*. The orthogroup with the largest number of *A. mexicanus* gene copies is

300 OG0000022 (n = 215). This group and the group with the second largest number of *A.*

301 *mexicanus* genes, OG0000023 (n = 205), are dominated by *A. mexicanus* zinc finger

302   proteins. Additional studies using independent tests of gene family expansions and

303   contractions, such as CAFE (Mendes et al. 2021), will be needed to truly resolve

304   duplication events prior to explorations of their functional roles in the *A. mexicanus* cave

305   adaptation. We provide gene symbols for human and *A. mexicanus* orthologs as a

306   resource for comparative studies (Suppl. Table 9). These *A. mexicanus* protein-coding

307   gene resources further efforts to improve genome annotation for *A. mexicanus*, offering

308   another aquatic model beyond zebrafish for human comparative studies. Future research

309   to identify gene duplicates, pseudogenization events, and nonsynonymous protein coding

310   changes unique to *A. mexicanus* within a framework of surface and independent cave

311   morphs will address hypotheses about how this species was adaptable to extreme

312   environmental change.

313

314   *Structural variation.* Cave morphs display many phenotypic differences from their

315   surface ancestor that motivated us to evaluate genome evolution on a finer scale beyond

316   chromosomal gene order evaluated herein and prior studies that have focused on protein-

317   coding gene changes (WARREN *et al.* 2021) (MORAN R.L. 2022). Discovery of SVs

318   among cave morphs compared to surface has only been estimated using short reads to a

319   less contiguous reference genome (WARREN *et al.* 2021).  To evaluate moderately-sized

320   SVs (50 to 10,000 bp) among these cave morphs compared to the surface genome, we

321   estimated the presence of deletions, insertions, and contractions and expansions of

322   repeats using Assemblytics (NATTESTAD AND SCHATZ 2016). The number of insertions

323   and deletions shows comparable total bases affected across these cave morphs relative to

324   surface regardless of their sequence length distribution (500-10,000 bp Fig. 3;50-500 bp

325   Suppl. Fig. 3; Suppl. Table 10). The average total size of detected cave morph deletions

326   and insertions was 19.4 MB (1.42% of the genome size) when aligned to the surface

327   genome (Suppl. Table 10). SV divergences unique to each cave morph were also evident.

328   For example, the Tinaja morph had the highest total sequence size for all insertions (12.9

329   MB), followed by Pachón (12.8 Mb) and Molino (10.8 Mb) (Suppl. Table 10). Similarly

330   for all deletions, Tinaja was highest (7.54 MB), then Pachón (7.10 Mb), and Molino

331   (7.08Mb) (Suppl. Table 10). Interestingly, total inserted sequence by size distribution

332   changes cave morph order with Pachón (9.9 Mb) having the largest amount for the 500 to

333   10,000 bp range followed by Tinaja (9.7 Mb) (Suppl. Table 10). When evaluating the

334   cave morph repeat landscape, we find an even larger percentage of the genome impacted

335   relative to insertions and deletions, with the average repeat expansion or contraction

336   being 3.05 and 2.91% of the total genome, respectively (Suppl. Table 10). These broadly

337   classified SV and repeat results highlight their differences vary by category and cave

338   morph origin, which may be the result of numerous factors, including assembly

339   completeness and accuracy, mixed haplotype assembly architecture, and the diversified

340   origins of each cave morph reference (HERMAN *et al.* 2018).

341

342   *Study conclusions.* The surface-to-cave genomic transitions that occurred in *A. mexicanus*

343   establish a unique model for the study of natural polygenic trait adaptation. Here, we

344   provide initial evidence that these more complete genomes will substantially advance our

345   capability to resolve these signatures of genetic adaption. The availability of nearly

346   complete genome copies of a surface and the independently evolved cave morphs will

347   drive future reevaluations of all types of segregating single nucleotide variants and SVs

348   in a pangenome-dependent manner (SIREN *et al.* 2021).

349

350

351

361

362

363    **Data availability**. All raw and processed data for this study are available by querying

364    NCBI BioProject accession numbers PRJNA807270, PRJNA819394, and

365    PRJNA819399. In addition, each assembly is available under GenBank numbers

366    GCA_023375975.1, GCA_023375835.1, and GCA_023375845.1. The full availability of

367    orthofinder results are attainable on figshare.

368

369    **Code availability.** Scripts used for this study are available at the following GitHub
370    repositories:
371    https://github.com/esrice/hic-pipeline
372    https://github.com/WarrenLab/purge-haplotigs-nf
373    https://github.com/WarrenLab/agptools

374

375

## Literature cited

Adrian-Kalchhauser, I., A. Blomberg, T. Larsson, Z. Musilova, C. R. Peart *et al.*, 2020 The round goby genome provides insights into mechanisms that may facilitate biological invasions. BMC Biol 18: 11.

Aken, B. L., S. Ayling, D. Barrell, L. Clarke, V. Curwen *et al.*, 2016 The Ensembl gene annotation system. Database 2016.

Albalat, R., and C. Canestro, 2016 Evolution by gene loss. Nat Rev Genet 17: 379-391.

Aspiras, A. C., N. Rohner, B. Martineau, R. L. Borowsky and C. J. Tabin, 2015 Melanocortin 4 receptor mutations contribute to the adaptation of cavefish to nutrient-poor conditions. Proc Natl Acad Sci U S A 112: 9668-9673.

Bilandžija, H., Hollifield, B., Steck, M., Meng, G., Ng, M., Koch, A.D., Gračan, R.,Ćetković, H., Porter, M.L., Renner, K.J., Jeffery, W.R., 2019 Phenotypic plasticity as an important mechanism of cave colonization and adaptation in Astyanax cavefish. bioRxiv.

Braasch, I., A. R. Gehrke, J. J. Smith, K. Kawasaki, T. Manousaki *et al.*, 2016 The spotted gar genome illuminates vertebrate evolution and facilitates human-teleost comparisons. Nat Genet 48: 427-437.

Cheng, H., G. T. Concepcion, X. Feng, H. Zhang and H. Li, 2021 Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. Nat Methods 18: 170-175.

Chernyavskaya, Y., X. Zhang, J. Liu and J. Blackburn, 2022 Long-read sequencing of the zebrafish genome reorganizes genomic architecture. BMC Genomics 23: 116.

Du, K., M. Pippel, S. Kneitz, R. Feron, I. da Cruz *et al.*, 2022 Genome biology of the darkedged splitfin, Girardinichthys multiradiatus, and the evolution of sex chromosomes and placentation. Genome Res 32: 583-594.

Duboue, E. R., A. C. Keene and R. L. Borowsky, 2011 Evolutionary convergence on sleep loss in cavefish populations. Curr Biol 21: 671-676.

Durand, N. C., J. T. Robinson, M. S. Shamim, I. Machol, J. P. Mesirov *et al.*, 2016 Juicebox Provides a Visualization System for Hi-C Contact Maps with Unlimited Zoom. Cell Syst 3: 99-101.

Emms, D. M., and S. Kelly, 2019 OrthoFinder: phylogenetic orthology inference for comparative genomics. Genome Biol 20: 238.

Ghurye, J., and M. Pop, 2019 Modern technologies and algorithms for scaffolding assembled genomes. PLoS Comput Biol 15: e1006994.

Ghurye, J., A. Rhie, B. P. Walenz, A. Schmitt, S. Selvaraj *et al.*, 2019 Integrating Hi-C links with assembly graphs for chromosome-scale assembly. PLoS Comput Biol 15: e1007273.

Guan, D., S. A. McCarthy, J. Wood, K. Howe, Y. Wang *et al.*, 2020 Identifying and removing haplotypic duplication in primary genome assemblies. Bioinformatics 36: 2896-2898.

Harrison, P. M., 2021 Computational Methods for Pseudogene Annotation Based on Sequence Homology. Methods Mol Biol 2324: 35-48.

420    Herman, A., Y. Brandvain, J. Weagley, W. R. Jeffery, A. C. Keene *et al.*, 2018 The role of
421            gene flow in rapid and repeated evolution of cave-related traits in Mexican
422            tetra, Astyanax mexicanus. Mol Ecol 27**:** 4397-4416.
423    Imarazene, B., K. Du, S. Beille, E. Jouanno, R. Feron *et al.*, 2021 A supernumerary "B-
424            sex" chromosome drives male sex determination in the Pachon cavefish,
425            Astyanax mexicanus. Curr Biol 31**:** 4800-4809 e4809.
426    J., D., 2020 AGAT: Another Gff Analysis Toolkit to handle annotations in any
427            GTF/GFF format., pp.  in *Zenodo*.
428    Jarvis, E. D., G. Formenti, A. Rhie, A. Guarracino, C. Yang *et al.*, 2022 Semi-automated
429            assembly of high-quality diploid human reference genomes. Nature 611**:** 519-
430            531.
431    Kurtz, S., A. Phillippy, A. L. Delcher, M. Smoot, M. Shumway *et al.*, 2004 Versatile and
432            open software for comparing large genomes. Genome Biol 5**:** R12.
433    Li, H., and R. Durbin, 2009 Fast and accurate short read alignment with Burrows-
434            Wheeler transform. Bioinformatics 25**:** 1754-1760.
435    Lovell, J. T., A. Sreedasyam, M. E. Schranz, M. Wilson, J. W. Carlson *et al.*, 2022
436            GENESPACE tracks regions of interest and gene copy number variation
437            across multiple genomes. Elife 11.
438    McGaugh, S. E., J. E. Kowalko, E. Duboue, P. Lewis, T. A. Franz-Odendaal *et al.*, 2020
439            Dark world rises: The emergence of cavefish as a model for the study of
440            evolution, development, behavior, and disease. J Exp Zool B Mol Dev Evol
441            334**:** 397-404.
442    Moore, B., M. Herrera, E. Gairin, C. Li, S. Miura *et al.*, 2023 The chromosome-scale
443            genome assembly of the yellowtail clownfish Amphiprion clarkii provides
444            insights into the melanic pigmentation of anemonefish. G3 (Bethesda) 13.
445    Moran R.L., R., E.J., Ornelas-García, C.P., Gross, J.B., Donny, A., Wiese, J., Keene, A.C.,
446            Kowalko, J.E., Rohner, N., McGaugh, S.E., 2022 Selection-driven trait loss in
447            independently evolved cavefish populations. bioRxiv.
448    Nattestad, M., and M. C. Schatz, 2016 Assemblytics: a web analytics tool for the
449            detection of variants from an assembly. Bioinformatics 32**:** 3021-3023.
450    Ojha, A., and M. Watve, 2018 Blind fish: An eye opener. Evol Med Public Health 2018**:**
451            186-189.
452    Pruitt, K. D., G. R. Brown, S. M. Hiatt, F. Thibaud-Nissen, A. Astashyn *et al.*, 2014
453            RefSeq: an update on mammalian reference sequences. Nucleic Acids Res 42**:**
454            D756-763.
455    Rautiainen, M., S. Nurk, B. P. Walenz, G. A. Logsdon, D. Porubsky *et al.*, 2023
456            Telomere-to-telomere assembly of diploid chromosomes with Verkko. Nat
457            Biotechnol.
458    Rhie, A., S. A. McCarthy, O. Fedrigo, J. Damas, G. Formenti *et al.*, 2021 Towards
459            complete and error-free genome assemblies of all vertebrate species. Nature
460            592**:** 737-746.
461    Roberts, M. B., D. T. Schultz, R. Gatins, M. Escalona and G. Bernardi, 2023
462            Chromosome-level genome of the three-spot damselfish, Dascyllus
463            trimaculatus. G3 (Bethesda) 13.

464    Robinson, J. T., D. Turner, N. C. Durand, H. Thorvaldsdottir, J. P. Mesirov *et al.*, 2018
465        Juicebox.js Provides a Cloud-Based Visualization System for Hi-C Data. Cell
466        Syst 6**:** 256-258 e251.
467    Rohner, N., 2018 Cavefish as an evolutionary mutant model system for human
468        disease. Dev Biol 441**:** 355-357.
469    Ruan, J., and H. Li, 2020 Fast and accurate long-read assembly with wtdbg2. Nat
470        Methods 17**:** 155-158.
471    Simao, F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva and E. M. Zdobnov,
472        2015 BUSCO: assessing genome assembly and annotation completeness with
473        single-copy orthologs. Bioinformatics 31**:** 3210-3212.
474    Siren, J., J. Monlong, X. Chang, A. M. Novak, J. M. Eizenga *et al.*, 2021 Pangenomics
475        enables genotyping of known structural variants in 5202 diverse genomes.
476        Science 374**:** abg8871.
477    Stockdale, W. T., M. E. Lemieux, A. C. Killen, J. Zhao, Z. Hu *et al.*, 2018 Heart
478        Regeneration in the Mexican Cavefish. Cell Rep 25**:** 1997-2007 e1997.
479    Volff, J. N., 2005 Genome evolution and biodiversity in teleost fish. Heredity (Edinb)
480        94**:** 280-294.
481    Warren, W. C., T. E. Boggs, R. Borowsky, B. M. Carlson, E. Ferrufino *et al.*, 2021 A
482        chromosome-level genome of Astyanax mexicanus surface fish for comparing
483        population-specific genetic differences contributing to trait evolution. Nat
484        Commun 12**:** 1447.
485    Warrenlab., 2022 Nextflow workflow for purging haplotigs from a genome assembly.

486    **Author notes**

487    **Competing interests**. The authors declare no competing interests.
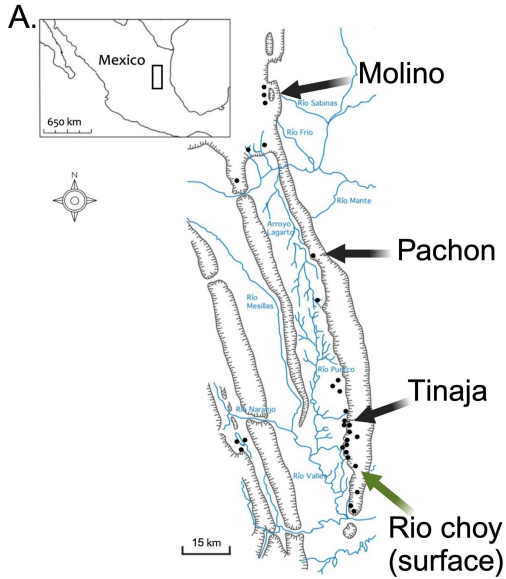
488

489    Figure legends

490    Figure 1. Summary of interspecies chromosomal synteny. A. Geographical representation
491    of the independent cave morph populations and their physical appearance used in this
492    study. Image courtesy of Alex Keene. B. Gene ortholog synteny of *A. mexicanus* against
493    three teleost species. C. Higher resolution image of *A. mexicanus* (surface genome)
494    aligned to the same three teleost species in B.

495

496    Figure 2. Summary of OrthoFinder analysis of vertebrate species *Astyanax mexicanus*,
497    *Danio rerio*, *Oryzias latipes*, *Xiphophorus maculatus*, *Lepisosteus oculatus*, *Pseudonaja*
498    *textilis*, *Podarcis muralis*, *Gallus gallus*, *Homo sapiens*, *Pan troglodytes*, *Mus musculus*,
499    *Rattus norvegicus*, *Monodelphis domestica*, *Xenopus tropicalis*, and *Ciona*
500    *intestinalis.* Bar charts describe data for each species, aligned to the matching species in
501    the tree. A. Phylogenetic tree built with all species using shared gene orthologs. B.
502    Number of orthogroups each classified by type per species. C. Percentage of genes by
503    orthogroups by species. D. Number of species-specific orthogroups per species. E.
504    Ortholog multiplicities for all species.

505

Figure 3. Genomic structural variation among *A. mexicanus* cave morph assemblies when compared to surface. Cave morph specific distribution by counts for A. Molino, B. Tinaja, and C. Pachón for the size distribution of 500 to 10,000 bp.

Table 1. Representative genome assembly metrics for sequenced *A. mexicanus* genomes.

| Common name | Assembled version | Total size Mb | Total contigs | Contig N50 length Mb | Unplaced Mb |
|---|---|---|---|---|---|
| Surface | Astyanax mexicanus-2.0 | 1,291 | 3,030 | 1.7 | 394 |
| Río Choy Surface | AstMex3_surface | 1,373 | 123 | 47 | 51 |
| Molino | AstMex3_Molino | 1,361 | 252 | 12 | 101 |
| Tinaja | AstMex3_Tinaja | 1,410 | 292 | 26 | 78 |
| Pachón | AMEX_1.1 | 1,378 | 529 | 14.8 | 29 |

A.

B.

*D. rerio*
(zebrafish)

*A. mexicanus*
(surface)

*X. maculatus*
(platyfish)

*O. latipes*
(medaka)

200 Mbp

C.

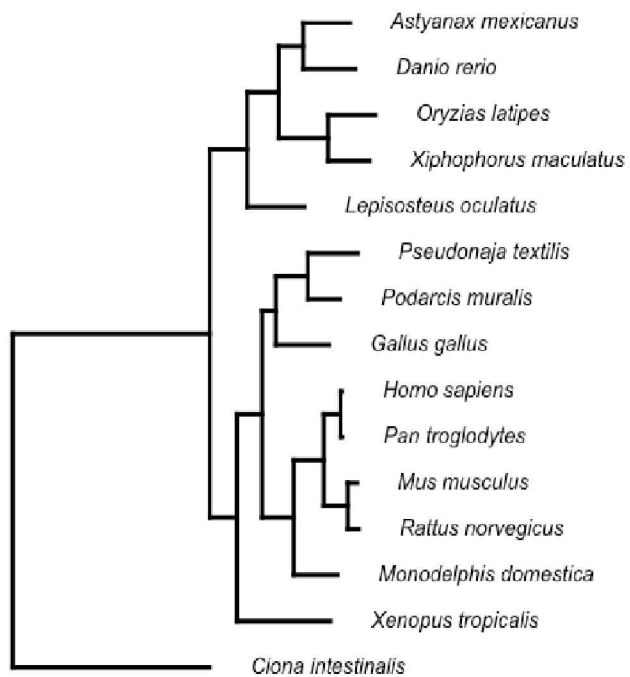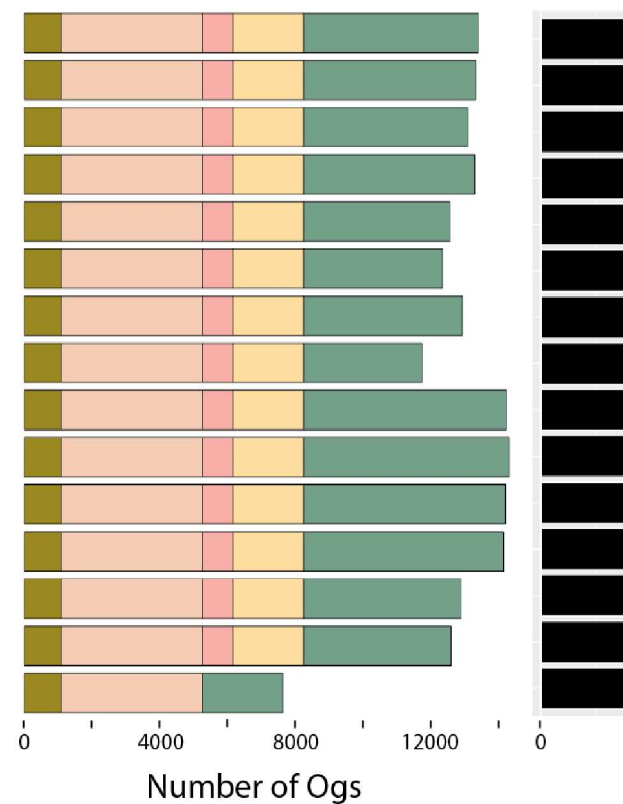Platyfish

Medaka

Zebrafish

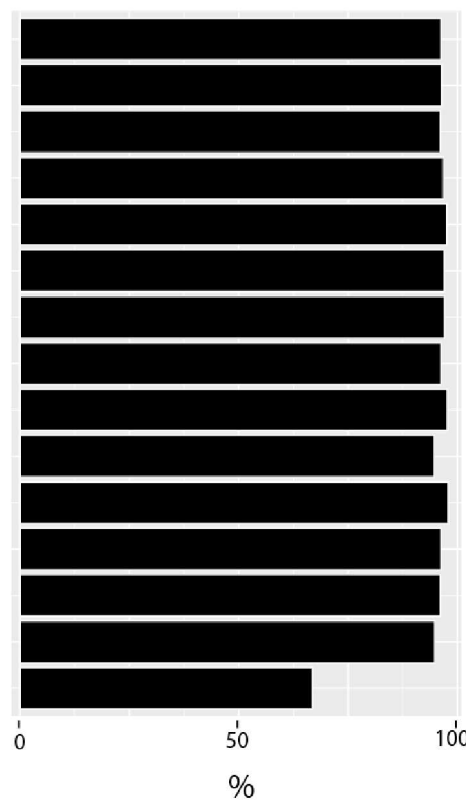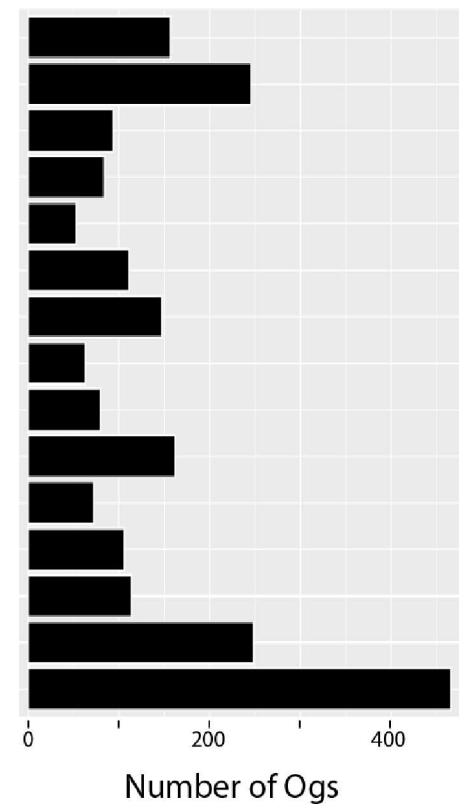Surface

800 genes

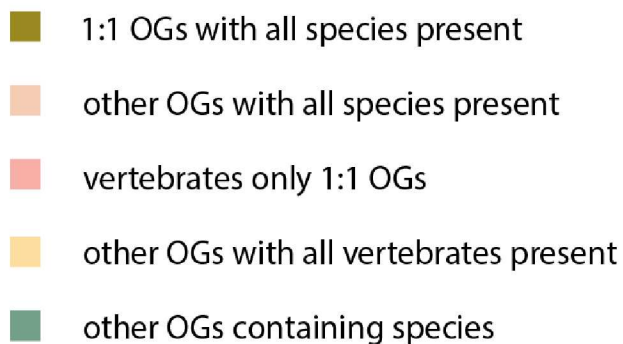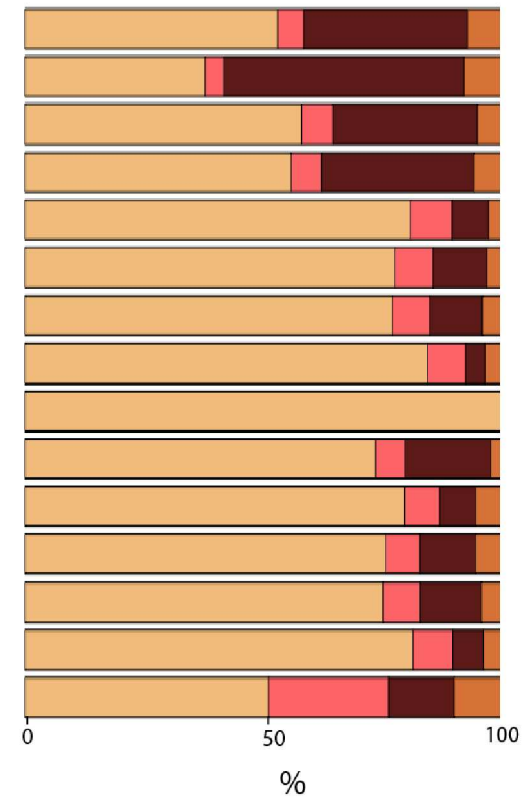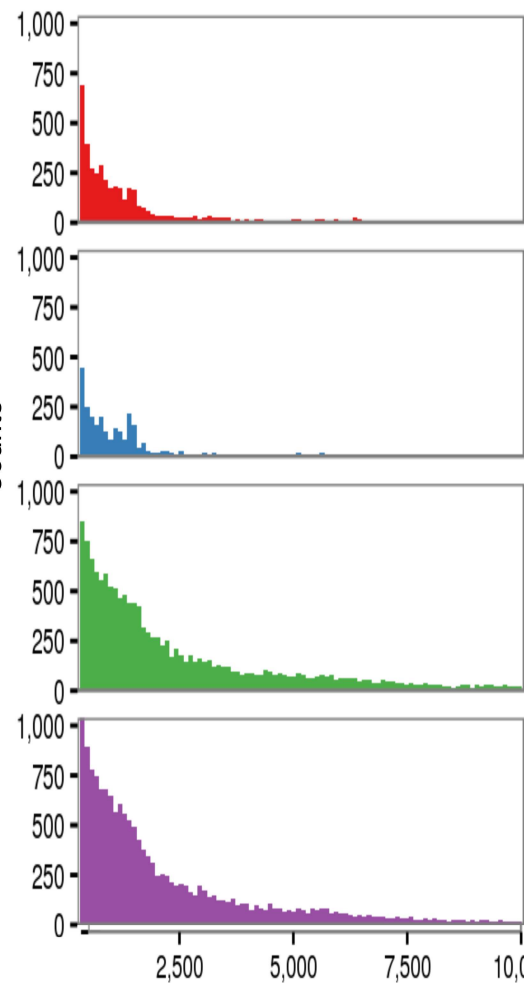| A. Species Tree | B. Orthogroup Types | C. % Genes in Orthogroups | D. Species-Specific Orthogroups | E. Multiplicity |
| --- | --- | --- | --- | --- |

*Homo sapiens*

Species Tree (A.):
- *Astyanax mexicanus*
- *Danio rerio*
- *Oryzias latipes*
- *Xiphophorus maculatus*
- *Lepisosteus oculatus*
- *Pseudonaja textilis*
- *Podarcis muralis*
- *Gallus gallus*
- *Homo sapiens*
- *Pan troglodytes*
- *Mus musculus*
- *Rattus norvegicus*
- *Monodelphis domestica*
- *Xenopus tropicalis*
- *Ciona intestinalis*

B. Number of Ogs (x-axis: 0, 4000, 8000, 12000)

C. % (x-axis: 0, 50, 100)

D. Number of Ogs (x-axis: 0, 200, 400)

E. % (x-axis: 0, 50, 100)

Legend (B.):
- 1:1 OGs with all species present
- other OGs with all species present
- vertebrates only 1:1 OGs
- other OGs with all vertebrates present
- other OGs containing species
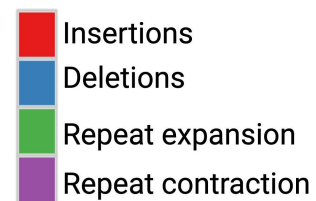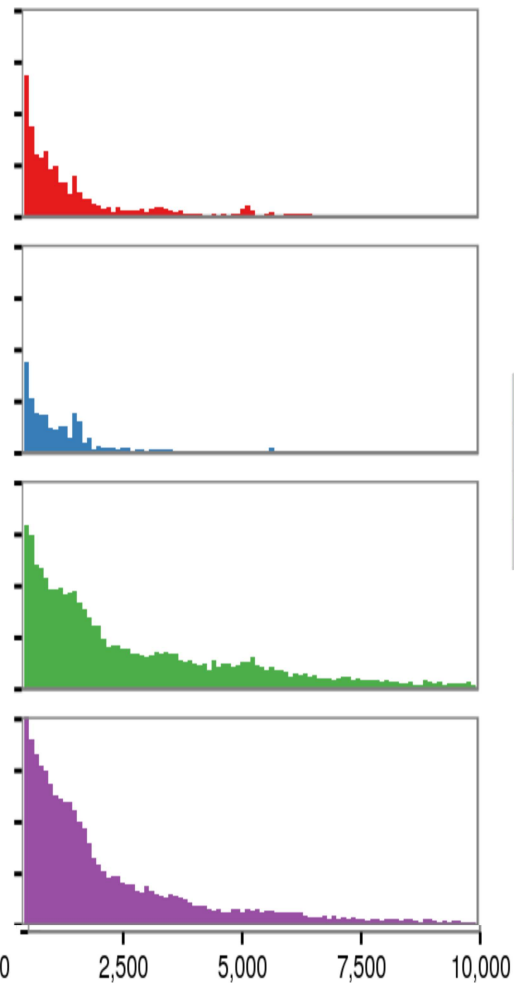
Legend (E.):
- 1 : 1
- 1 : many
- many : 1
- many : many

A. Molino   B. Tinaja   C. Pachon

Counts

Variant size (bp)

Insertions
Deletions
Repeat expansion
Repeat contraction