1   **Title:** Tissue-specific transcriptome for *Poeciliopsis prolifica* reveals evidence for genetic

2   adaptation related to the evolution of a placental fish.

3   **Authors and Affiliations:**

4   Nathaniel K. Jue*,[†], Robert J. Foley*, David N. Reznick[‡], Rachel J. O'Neill* and Michael J. O'Neill*

5

6   * Institute for Systems Genomics and Department of Molecular and Cell Biology, University of
7   Connecticut, Storrs, CT 06269
8   [†] current address: School of Natural Sciences, California State University, Monterey Bay,
9   Seaside, CA 93933
10  [‡] Department of Biology, University of California, Riverside, CA 92521
11

2

18  **Running Title:** Transcriptomics of Placental Fish

19  **Key words:** Transcriptome, Positive Selection, Gene Expression, Placenta, Fish

20  **Corresponding author:** Institute for Systems Genomics and Department of Molecular and Cell
21  Biology, University of Connecticut, Storrs, CT 06269, USA, Ph: (860)486-6856, Fax: (860)486-
22  1936, Email: michael.oneill@uconn.edu

## ABSTRACT

23

24        The evolution of the placenta is an excellent model to examine the evolutionary processes

25    underlying adaptive complexity due to the recent, independent derivation of placentation in

26    divergent animal lineages. In fishes, the family Poeciliidae offers the opportunity to study

27    placental evolution with respect to variation in degree of post-fertilization maternal provisioning

28    among closely related sister species. In this study, we present a detailed examination of a new

29    reference transcriptome sequence for the live-bearing, matrotrophic fish, *Poeciliopsis prolifica*,

30    from multiple-tissue RNA-seq data. We describe the genetic components active in liver, brain,

31    late-stage embryo, and the maternal placental/ovarian complex, as well as associated patterns of

32    positive selection in a suite of orthologous genes found in fishes. Results indicate the expression

33    of many signaling transcripts, "non-coding" sequences and repetitive elements in the maternal

34    placental/ovarian complex. Moreover, patterns of positive selection in protein sequence

35    evolution were found associated with live-bearing fishes, generally, and the placental *P.*

36    *prolifica*, specifically, that appear independent of the general live-bearer lifestyle. Much of the

37    observed patterns of gene expression and positive selection are congruent with the evolution of

38    placentation in fish functionally converging with mammalian placental evolution and with the

39    patterns of rapid evolution facilitated by the teleost-specific whole genome duplication event.

## INTRODUCTION

40

41        The study of the placenta provides insight into the evolutionary relationships of

42    biological phenomena such as complexity, live-birth and genetic conflict. A great deal of

43    research has focused on the function and development of mammalian placentas, uncovering the

44    unique regulatory, genetic, and evolutionary nature of this structure. Studies of gene regulation in

45    the mammalian placenta show a suite of unique features including genomic imprinting (Bressan

46    *et al.* 2009), non-coding RNAs (Koerner *et al.* 2009), and DNA methylation and histone-

47    modification mediated transcription (Maltepe *et al.* 2010). The placenta also has been shown to

48    be a tissue that utilizes genes derived from the co-option of retroelements for unique functional

49    purposes (Lavialle *et al.* 2013). Additionally, the placenta has been used as a model for

50    examining the evolution of tissue-specific novelties, such as newly derived cell-types (Lynch *et*

51    *al.* 2011), placental variation among eutherian mammals (Carter and Mess 2007), and genomic

52    imprinting related to viviparity (Renfree *et al.* 2013).

3

53      Placentation is typically studied in mammals, but fish present a compelling study system

54    for examining contributing factors to the evolution of this complex organ. The Neotropical fish

55    family Poeciliidae is comprised of approximately 200 species, all of which, with one exception,

56    give live birth. The majority of these poeciliids are lecithotrophic (i.e. yolk-feeding), wherein

57    eggs provide all necessary nutrients to support the embryo through development to birth.

58    However, placenta-like structures that permit post-fertilization maternal provisioning have

59    evolved independently in multiple poeciliid lineages, specifically within certain groups such as

60    species in the genus *Poeciliopsis*, within the last 750,000 years (Reznick *et al.* 2002). Unlike

61    comparisons between eutherian and marsupial mammals, who last shared an ancestor with their

62    non-placental monotreme counterparts (i.e. the egg-laying platypus and echidnas) ~200 million

63    years ago (Meredith *et al.* 2011), species within *Poeciliopsis* offer the opportunity to investigate

64    more "recent" changes leading to viviparity and placentation. The relatively recent adaptation of

65    placentation has resulted in wide variation among *Poeciliopsis* species with respect to the extent

66    of maternal provisioning. The extent of maternal investment across species ranges from highly

67    matrotrophic (i.e. placentotrophic) to lecithotrophic, including intermediate, or "partial",

68    placental species. These transitional states and independent evolutionary events make this system

69    particularly powerful for examining factors contributing to the evolution of placentation (see

70    Pollux *et al.* 2009 for review).

71       Although fish placentas exhibit functional convergence, they are diverse in structure,

72    with poeciliid placentas bearing features distinct from mammalian placentas. In poeciliids, the

73    maternal portion of the placenta is derived from the ovarian follicle. Fertilization occurs within

74    the ovarian follicle wherein the embryo will subsequently develop. Within placental *Poeciliopsis*

75    species, nutrient exchange occurs across an enlarged pericardial sac that contributes to a large,

76    highly vascularized belly sac (Turner 1940). In the closely related poeciliid species *Heterandria*

77    *formosa*, functionally convergent placental structures are notably divergent in structure; the

78    aforementioned sac structure covers regions more anterior on the developing embryo (Turner

79    1940). While specializations to the follicular epithelium, such as a thick, vascularized follicle

80    wall with dense microvilli and specialized cytoplasmic organelles are common features in the

81    maternal poeciliid placenta, much remains unknown about the ontogeny of the poeciliid

82    follicular placenta (Turner 1940; Grove and Wourms 1991).

4

83   To define the genetic components contributing to placental function and examine the

84 selective forces influencing the evolution of this unique poeciliid fish lineage, we constructed a

85 new reference transcriptome for the placental fish *Poeciliopsis prolifica*, the blackstripe

86 livebearer. A placental tissue-specific transcriptome profile was generated by comparison to non-

87 placental tissues from *P. prolifica*, while patterns of protein evolution were compared with other

88 closely and distantly related fish species. *P. prolifica* is a highly matrotrophic poeciliid fish that

89 shares a hypothesized lecithotrophic common ancestor with recently diverged lecithotrophic

90 sister taxa (Reznick *et al.* 2002), thus presenting a model system for examining evolutionary

91 genetic changes proximal to the emergence of the placenta. Notably, we find evidence indicating

92 genetic parallelism, both in function and evolution, of the fish placenta and the mammalian

93 placenta.

## METHODS AND MATERIALS

95 **Samples**

96   Tissue samples were harvested according to an IACUC approved protocol from captive

97 populations of *Poeciliopsis prolifica* raised at the University of Connecticut. Original stocks

98 were obtained from stock populations at the University of California-Riverside under care of Dr.

99 David Reznick and from Ron Davis, a live-bearer hobbyist in Florida. Both populations

100 originated from the same sample population from the Rio El Padillo in Mexico. Tissues were

101 isolated from fish dissected on ice, immediately snap frozen with liquid nitrogen, and stored at -

102 80° C. For this study, four sample types were isolated: female brain, liver, whole embryo, and

103 the maternal placental/ovarian tissue complex (MPC). Whole female brain was dissected from

104 the skull and is inclusive of the olfactory bulb, cerebrum, optic lobe, cerebellum and medulla

105 oblongata (to the tip of the spinal cord). Due to its delicate nature, maternal placental tissue was

106 isolated by dissecting whole ovary from pregnant females, excising any fertilized and observable

107 unfertilized eggs, tearing open ovarian follicles, removing developing embryos from those

108 follicles, and reserving the remaining maternal placental/ovarian tissue complex (MPC) that

109 included both ovarian follicles and some remaining ovarian tissue (Figure S1). Late-stage (i.e.

110 nearly full-term) whole embryos, identified by full pigmentation, large size, an ability to persist

111 after being excised from ovarian follicle, and being "late-eyed" (Stage 5 as described by

112 (Reznick 1981)) were sampled and stored with belly sacs intact.

113

**Sequencing**

115    Two types of sequencing platforms, Roche 454 and ABI SOLiD, were implemented in

116    this study. For 454 sequencing, RNA was isolated from 20 different individuals by

117    homogenizing and disrupting selected tissue samples with syringes in a Trizol solution. Due to

118    individual isolation yields, required template inputs for library construction, and to compensate

119    for among-individual variation, each RNA sample was then pooled by tissue type and mRNA

120    was isolated from 5-10 μg of total RNA using the Poly(A) Purist kit (Ambion). All RNA

121    samples were assessed for quality on a Bio-Rad Experion both pre- and post-Poly(A) extraction.

122    Sequencing libraries were made following standard RNA-Seq library construction protocol for

123    454 sequencing and sequenced on a Roche 454 Sequencer. To generate SOLiD sequencing data,

124    tissues for three individual MPCs and an embryo from one of these same females were first

125    stored in RNALater and then at -80° C. RNA was isolated by disruption and homogenization of

126    tissues using a Polytron and the RNAeasy mini kit (Qiagen). DNA was removed from each

127    sample by TurboDNAse (Ambion) and validated for sample integrity using an Agilent

128    Bioanalyzer. ERCC spike-in controls (Life Technologies) were then added to each sample and

129    ribosomal RNA (rRNA) was removed using the Ribozero kit (Epicenter). Final RNA-Seq

130    libraries were constructed from the resultant mRNA sample using standard SOLiD transcriptome

131    library construction protocols. Libraries were sequenced on an ABI SOLiD 5500xl.

**Assembly**

133    Post-sequencing, all 454 reads were trimmed using 454 Newbler software to remove bar

134    codes and the program CUTADAPT v1.2.1 (Martin 2011) to remove adapter sequences and trim

135    low quality regions of reads. Seqclean was then used to remove poly-A tails. CUTADAPT was

136    also used for trimming out all barcode and adapter sequences as well as quality trimming for

137    SOLiD libraries. All SOLiD libraries were then screened against an in-house database of rRNA

138    sequences to remove any rRNA sequences that may have not been removed in the rRNA-

139    depletion step. All remaining SOLiD reads were normalized using the Trinity-associated *in silico*

140    k-mer normalization protocols. All trimmed 454 reads and normalized SOLiD reads from all

141    tissues were then input into the Trinity transcriptome assembler (release 7/17/2014) (Grabherr *et*

6

142    *al.* 2011). Following the Trinotate pipeline (release 4/30/2015) for annotating predicted

143    transcripts (Haas *et al.* 2013), open-reading frames (ORFs) were predicted using Transdecoder

144    (release 1/27/2015). All transcripts and predicted proteins were then annotated via homology

145    against the SwissProt/Uniprot database and assigned any associated Gene Ontology (GO) terms

146    and eggNOG orthologs group membership. Predicted proteins were also searched for PFAM

147    protein domain and identification as a signaling protein using SignalP (v4.1) (Nielsen 2017),

148    transmembrane protein using TMHMM (v2.0) (Krogh *et al.* 2001), or ribosomal RNA using

149    RNAmmer (v1.2) (Lagesen *et al.* 2007). All transcripts were examined for any additional

150    homologies against the NCBI *nr* database using BLASTX and annotated using BLAST2GO

151    (v2.5.0) (Conesa *et al.* 2005). Any transcript without an *nr* BLASTX-hit was also searched

152    against the NCBI *nt* database with BLASTN. Finally, all transcripts were assessed with

153    BLASTN for homology with known non-coding RNAs (ncRNAs) identified in zebrafish (*Danio*

154    *rerio*) (Ulitsky *et al.* 2011). Databases versions for all homology searches were all updated on

155    7/1/15 before this analysis was completed.

156          Tissue-specific gene expression patterns were surveyed by mapping reads to the Trinity

157    assembled transcriptome sequence, quantifying read coverage among transcripts, and testing for

158    differences among comparison groups. Mapping was performed using BWA (v0.7.7)(SW

159    algorithm) (Li and Durbin 2010) for all 454 data, and Bowtie2 (v4.1.2) (Langmead and Salzberg

160    2012) for all SOLiD data. Gene expression and read counts were estimated for all transcripts

161    using the program eXpress v1.5.1 (Roberts and Pachter 2013). Count data from 454 mapping

162    was passed through R-based DESeq2 analysis (Love *et al.* 2014) to assess significant differences

163    in pairwise comparisons of gene expression patterns among tissue samples, while correcting p-

164    values for False Discovery Rates (FDR) due to multiple comparison tests. Since sequencing

165    libraries were generated from pooled samples, they were assumed to represent an "average"

166    perspective. Due to the lack of replicates of pooled samples, best practices outlined in the

167    DESeq2 manual were used to generate dispersion estimates by comparing counts among tissue

168    types as opposed to between replicates. This process should be conservative with respect to false

169    positives since it errs on the side of using larger than necessary dispersion values. FPKM

170    (fragments per kilobase per millions reads) values were then used in BioLayout Express3D

171    (v3.2) (Theocharidis *et al.* 2009), along with the MCL (v12-068) clustering algorithm (van

172    Dongen and Abreu-Goodger 2012), to generate a preliminary 3-D gene atlas of co-expressed

173  genes clusters. Due to modest read coverage of 454 sequencing libraries, only "highly"

174  expressed genes (an FPKM value > 50 in at least one tissue) were included in clustering

175  analyses.

**Evolutionary Rates**

177  Evidence of positive selection in the evolutionary rates of poeciliid genes was tested

178  using the branch-sites models implemented in the program PAML v4.7 (Yang 2007).  cDNA

179  resources for six other species of fish whose genome and gene models have already been

180  described were downloaded from ENSEMBL and compared to our sequences for *P. prolifica*.

181  These species included the following: *Danio rerio*, *Gadus morhua*, *Takifugu rubripes*,

182  *Oreochromis niloticus*. *Gasterosteus aculeatus*, and *Xiphophorus maculatus* (Figure S2). Of

183  these six species, *X. maculatus* is the most closely-related species to *P. prolifica*; both are in the

184  family Poeciliidae. However, *X. maculatus* differs significantly from *P. prolifica* in reproductive-

185  style since it is a lecithotrophic (yolk-feeding) live-bearer with no evidence of post-fertilization

186  maternal provisioning. *P. prolifica* is highly matrotrophic, with sufficient post-fertilization

187  maternal provisioning to sustain an eight fold increase in dry mass between the fertilization of

188  the egg and birth (Pires *et al.* 2007). Predicted coding sequence regions for *P. prolifica* were

189  compared to cDNA reference sequences for each species using reciprocal best BLAST hit

190  approaches (TBLASTX in this case) to identify orthologous genes between species. Once

191  orthologs were identified, all orthologous gene clusters that lacked a predicted ortholog for any

192  species (i.e. no reciprocal best BLAST hit found) or, when examining high-scoring segment pair

193  (HSP) alignment regions, that yielded a multiple sequence alignments less than <200 bp long

194  were discarded. Using in-house Python scripts, the remaining orthologs were passed through a

195  series of analysis steps. Groups of orthologs were first reconstructed in the same strand and

196  aligned using the codon-guided multiple sequence alignment (MSA) algorithm MACSE v 0.9b1

197  (Ranwez *et al.* 2011). MSAs were cleaned using trimAl (Capella-Gutiérrez *et al.* 2009) to

198  remove all gaps both from within, and at the ends of, the aligned sequences. MACSE includes

199  the convenient feature of assessing frameshift and stop codon issues associated with multiple

200  sequence alignment. Thus, in order to avoid confounding alignment problems related to poor

201  data quality, low scoring MSAs and true pseudogenized gene sequences, all of which would

202  contribute to false positives in subsequent PAML analyses, this feature was leveraged to identify

8

203 and remove any MSA with either a frameshift ambiguity or base ambiguity from further
204 analysis.

205     The remaining MSAs were then analyzed in PAML with three different phylogenetic
206 "foregrounds" to test for positive selection in rapid codon evolutionary rates: *P. prolifica* only, *X.*
207 *maculatus* only, and all poeciliids. These three levels of examination provided a proxy test of the
208 evolutionary changes possibly associated with three reproductive-styles, respectively:
209 matrotrophic viviparity, lecithotrophic vivparity, and vivparity (generally). Classification of sites
210 having significant evidence for being under positive selection required a significantly better fit of
211 the branch-sites alternative model of positive selection over the null model (implemented as
212 described in the PAML manual – Model 2A vs. Model 1A – with a $\chi^2$ test using p-value < 0.05 as
213 the threshold for identifying significant improvements in maximum likelihood model fit) and
214 identification using the Bayes empirical Bayes (BEB) method (p-value >0.95). All sites and
215 predicted proteins were compared among different "foreground" analyses to classify protein
216 evolution associated with the aforementioned reproductive-style that these species represent.

217     A distance-based gene family tree for the *RAB11 family-interacting protein* gene family
218 (*RAB11FIP*) was constructed using neighbor-joining tree methods to describe the general
219 patterns of gene duplication and evolution in fishes. Jukes-Cantor distances among protein
220 sequences were used to generate tree topology. All sequences included in this gene family tree
221 where gathered by identifying any *P. prolifica* predicted protein sequence with homology to
222 *RAB11FIP*s in *Danio rerio* using BLASTP (e-value < 1e-5) and using those predicted proteins to
223 identify any other existing protein sequences for *RAB11FIP* genes in fishes using BLASTP (e-
224 value <1e-5; taxonomically restricted search to "bony fishes" – taxid: 7898). MUSCLE v3.8.31
225 (Edgar 2004) was used to generate a multiple sequence alignment for all sequences and CLC
226 Genomics Workbench v7.5 was used to generate a tree with 100 bootstraps. To focus analysis on
227 *RAB11FIP* genes only, all clusters of genes identified as the protein *UNC-13* (a homologous
228 gene to *RAB11FIP*s) were trimmed from final tree.

229 **Data Availability:**

230 All read data was deposited in the NCBI SRA database under the following accession numbers:
231 SRR1639275, SRR1640127, SRR1640137, SRR1640160, SRR1640171, SRR1640200,

232    SRR1640209, SRR1640216, and SRR1640219 under the BioProject PRJNA266248. All custom

233    scripts are available here:

234    https://github.com/juefish/Jue_et_al_G3_P_prolifica_transcriptome.git.

235                                        **RESULTS**

236    **Assembly Statistics**

237         *De novo* assembly of 3,696,154 Roche 454 and 159,802,508 SOLiD reads (post-

238    trimming, see Table S1 for library details) yielded a transcriptome of 331,767,677 Mb (43.74%

239    GC) with 478,065 predicted transcripts (TSA Reference ID: GBYX00000000.1). Average contig

240    length was 639 bp and N50 was 885 bp. These contigs were grouped into 319,532 components,

241    which are analogous to estimated "genes" or groups of isoforms (Table 1). While some of these

242    predicted transcripts could be spurious or fragmented results from the assembler, 236,360

243    (49.4%) of these predicted transcripts were well-supported with read depth of coverage >10x,

244    representing a very diverse transcriptome (Table 1).

245         Within this assembled transcriptome, 113,240 transcripts (23.6% of total) were predicted

246    to have a protein open-reading frame (ORF) (Figure 1), with over 80% of these predicted

247    proteins (both total transcripts and genes) carrying homology with a protein in the

248    UniProtKB/Swiss-Prot database, and >75% of those showing associations with known Pfam

249    domains (Figure 1). Functional Gene Ontology (GO) annotations were identified for the majority

250    of these sequences with homology to *nr* database reference sequences, representing a multitude

251    of functional elements, spanning a range of categories in the Gene Ontology (Figure 2). Another

252    41,851 transcripts with no BLAST result at all (8.7%) showed similarity to REPBASE repetitive

253    element sequences, including 1,747 transcripts from 1,043 predicted genes that incorporated

254    repetitive element genes (Table S3). These transcripts span a wide-range of repetitive element

255    origins, including elements known to have specific placental function in mammals such as

256    *retrotransposon-derived protein PEG10-like*. Another 286 transcripts carry regions identified by

257    homology with non-coding RNAs from *D. rerio*. These transcripts represent a variety of non-

258    coding RNAs that may be involved in gene regulation (Table S4). For instance, one identified

259    transcript shows homology with *cyrano*, a lncRNA demonstrated to be necessary for proper

260    embryonic development and interacting with a known miRNA miR-7 (Ulitsky *et al.* 2011). A

10

261  small number (24) of these transcripts showed evidence for bidirectional transcription and, thus,

262  candidates for active functioning in gene regulation through complementary base-pairing with

263  coding transcripts.

**Tissue Specific Gene Expression**

265       Using MCL clustering of gene expression estimates, we generated a preliminary gene

266  atlas for *P. prolifica* to identify clusters of co-expressed transcripts among four different sample

267  types: MPC, female brain, liver, and late-stage developing embryo, hereafter referred to as

268  "tissues". Before clustering, pairwise tests for significant differences (p-value <0.05 after

269  correction for FDR) in gene expression using DESeq2 were conducted across all transcripts in all

270  tissues and revealed 45, 108, 18, and 24 transcripts were specifically expressed in MPC, whole

271  embryo, brain and liver, respectively. For MCL clustering analysis and gene atlas construction, a

272  subsample of the 6,839 most highly expressed transcripts (FPKM values >50 in at least one of

273  the four tissues) were included in the analysis. This subset further reduced the number of

274  identifiable (via pairwise comparisons) tissue-specific transcripts included in the atlas that were

275  significant for tissue-specific expression to 24, 36, 4, and 5 for MPC, embryo, brain and liver,

276  respectively. Using the tissue-specific gene expression patterns of these transcripts (Figure S3)

277  and the MCL clustering algorithm, nine co-expressed gene clusters were identified (Figure 3).

278  Cluster 1 was the largest cluster and generally associated with transcripts that have high

279  expression in the brain, but showing some co-expression with other tissues, particularly MPC

280  and embryo. Cluster 2 was generally associated with transcripts highly expressed in embryo,

281  cluster 3 was associated with transcripts highly expressed in MPC, and cluster 4 was associated

282  with transcripts highly expressed in liver. Clusters 5 to 9 (which represented only 2.2% of the

283  transcripts in the atlas) were defined by expression across multiple tissue types, displaying gene

284  expression profiles indicative of "house-keeping"-like genes (Figure S3). Transcripts with

285  significant evidence for tissue-specific expression largely supported these cluster classifications

286  with 32 of the 36 aforementioned "embryo"-specific genes in cluster 2 and all 24 of the MPC

287  genes in cluster 3. Brain and liver clusters were less clearly supported with none of the four

288  "brain" genes in cluster 1 and only one of the five "liver" genes in cluster 4; however, the

289  number of transcripts in these clusters was so low that detectability may have been limited.

290  Transcripts involved in progesterone signaling pathways were observed as highly expressed in

291 placental tissues. Overall, 242 transcripts with ORFs were identified as having GO-associations

292 with progesterone regulatory pathways, including *Protein DEPP* (*decidual protein induced by*

293 *progesterone*), suggesting that similar developmental patterns in cell differentiation and

294 specialization maybe be occurring in fish as it does in mammals during pregnancy (Watanabe, et

295 al. 2005).

**Repetitive Element Transcripts**

297      Repetitive element gene expression was observed across various tissue samples and a

298 subset of the gene atlas clusters. Of the clustered 454 expression data, the MPC cluster (#3) had

299 the highest number of repetitive element transcripts, with a total of 9 transcripts; the "brain"

300 cluster (#1) had the second highest repetitive element transcript count at five transcripts. Cluster

301 2 (embryo), cluster 4 (liver), and cluster 5 (multiple tissues) had 2, 1, and 1 transcript(s),

302 respectively. Only one transcript of these 18 transcripts found in the gene atlas clusters

303 (identified as a *transposable element tc1 transposase*) showed no expression in placenta; all 17

304 other transcripts were expressed (>50 FPKM) in MPC (eight of these transcripts were also

305 identified as homologs to *transposable element tc1 transposases*). One transcript (a *reverse*

306 *transcriptase*) was also identified using the aforementioned pairwise significance testing

307 (DEseq2, p-value < 0.05) as more expressed in MPC as opposed to other tissues (FPKM $_{MPC}$ =

308 155.7 vs. FPKM $_{average\_other\_tissues}$ = 5.07). The three MPC SOLiD libraries also indicated high

309 levels of MPC gene expression of repetitive element-derived transcripts. From the SOLiD RNA-

310 Seq data, 98% of the 1,747 transcripts from the broader transcriptome reference sequence and

311 originating from repetitive elements were expressed in either MPC or embryonic tissues, with

312 227 predicted transcripts from 199 predicted genes expressed either only in the MPC or >5 fold

313 greater expression in MPC over embryonic tissues (Table S3). Approximately an equal number,

314 213 predicted transcripts and 199 predicted genes were found associated with embryonic tissues

315 using the same criteria (Table S3). Eight transcripts had an FPKM value of >50 across and were

316 identified as four gene families that included an envelope protein, a partial pol protein, a tc1

317 transposase and a tc3 element. In addition to the gene classes mentioned above, other repetitive-

318 element transcripts were identified as *retrotransposon-derived protein PEG10-like*, *120.7 kDa*

319 *protein in NOF-FB transposable element*, *retroviral polyprotein*, and *transposable element tcb1*

320 *transposases*. These transcripts appeared unique to the poeciliid lineage, showing between 50%

12

321    and 70% similarity to other repetitive element reference sequences from other species, with only

322    *retrotransposon-derived protein PEG10-like* showing high similarity (88%) with reference

323    sequences from the NCBI *nr* database.

**Transcripts with Unknown Function**

325    The majority of these clusters of highly expressed genes consisted of transcripts with no

326    known annotation. Of the highly expressed transcripts described in these clusters, 79.4%

327    (n=6260) were not identifiable via BLAST searches of SwissProt/UniProt, *nr* and *nt* databases

328    (e-value $< 1 \times 10^{-5}$). A large number (786, or 12.6%, of the total unknowns) of these predicted

329    transcripts had evidence for some type of repeat in their sequence, with 761 of the repeats

330    identified as either a simple repeat or low complexity sequence, indicating that the sequence may

331    be part of a non-coding region (Wren *et al.* 2000; Morgante *et al.* 2002; Liu *et al.* 2012). Many

332    of these sequences are likely either species-specific 5' or 3' UTRs or previously undescribed

333    non-coding RNAs. For example, another four of these transcripts in this cluster were associated

334    with known non-coding RNA sequence from *D. rerio* (3 with miRNAs and 1 with a lncRNA);

335    however, given that all of these sequences were much longer than miRNA size (312-982 bp) and

336    not readily identifiable as miRNA precursors (Liu *et al.* 2015), they are more likely to be binding

337    sites for such targets than host transcripts. Another 49 transcripts had predicted ORFs associated

338    with them, but no BLAST annotation and thus appear to be novel protein sequences. Of these 49

339    predicted proteins, two were identified as prospective signaling peptides, one of which was a

340    member of the MPC gene cluster. The other "signaling" peptide and two other predicted proteins

341    were identified as transmembrane proteins. The signaling/transmembrane protein was a member

342    of the "house-keeping gene" cluster (but most highly expressed in liver), while the other two

343    transmembrane proteins were associated with either the "brain" cluster or the "embryo" cluster.

344    Notably, the "embryo" cluster member was also highly expressed in MPC ($\text{FPKM}_{embryo}$=53.5;

345    $\text{FPKM}_{placenta}$=39.5). Given exhaustive attempts to annotate these sequences and the fact that they

346    are highly expressed transcripts, these sequences appear to be novel to this species.

**Protein Evolutionary Rates**

348    Reciprocal best BLAST hits of the cDNA coding sequence against the predicted and

349    known cDNAs for six fish species with sequenced genomes revealed predicted *P. prolifica*

13

350 transcripts to have 12,631 orthologs with *Danio rerio*, 14,761 orthologs with *Xiphophorus*

351 *maculatus*, 12,899 orthologs with *Takifugu rubripes*, 12,316 orthologs with *Gadus morhua*,

352 13,388 orthologs with *Gasterosteus aculeatus*, and 13,282 orthologs with *Oreochromis niloticus*.

353 Out of all of these orthologs, only 5,398 were shared orthologs for all seven species (including *P.*

354 *prolifica*). Within this shared ortholog set, 963 ortholog alignments showed evidence of open-

355 reading frame indels in at least one species' orthologous sequence, resulting in a frame-shift in

356 predicted codon sequences (Table S5). These frame-shifts could be the result of errors in a given

357 fish reference sequence or bona fide mutations in a specific species. While all species showed

358 evidence for frame-shifts, transcript sequences from *D. rerio*, *P. prolifica*, and *X. maculatus* had

359 a higher proportion of orthologs with an identified frame-shift than the remaining species (Table

360 S5). Additionally, 978 ortholog groups were discarded from the PAML analysis due to

361 ambiguous bases ("N") in the reference sequences; this was a disproportionately acute issue with

362 *G. morhua* sequences (912 orthologs).

363 Within the final set of 3,457 orthologs employed in our PAML analyses, 2,298 sites

364 across 404 predicted proteins were identified as undergoing positive selection. Of these sites,

365 917, 1104, and 247 were associated with *P. prolifica*, *X. maculatus*, and both poeciliids,

366 respectively (Figure 4, Table S6-S13). The predicted proteins carrying these sites covered a

367 wide-range of biological functions (Figure S4) with no overall significant enrichment for any

368 specific functional GO terms relative to the overall transcriptome annotation. Comparisons

369 between the matrotrophic *P. prolifica* and lecithotrophic *X. maculatus* orthologs with sites under

370 positive selection showed genes under positive selection in *P. prolifica* to be significantly

371 enriched for a variety of GO terms over those found in *X. maculatus* (Figure 5, FDR p-value <

372 0.05). The terms were generally associated with Biological Processes related to biosynthesis and

373 regulatory processes, Molecular Functions terms related to nucleic acid binding, and Cellular

374 Components terms related to the nucleus. Of these sites, 1,376 occurred in regions of these open-

375 reading frames that carried no discernable, previously known protein domain defined by Pfam

376 database searches. Thus, these sites indicate possible novel functional domains for these proteins

377 in *P. prolifica*.

378 While the majority of proteins undergoing positive selection (67%) had less than five

379 sites identified as under positive selection, many of the genes under positive selection exhibited

14

380    evidence for extensive rapid evolution (Table S7). For instance, the *GRAM domain-containing*

381    *protein 4*, *GRAMD4*, carries 94 sites identified as evolving rapidly in *X. maculatus*. These sites

382    account for 16% of the entire protein sequence for this gene. None of these sites overlap with the

383    known GRAM protein domain, indicating that this region may be an important novel functional

384    domain. *GRAMD4* is a membrane protein known to be a tumor suppressor in apoptotic pathways

385    associated with mitochondria (John *et al.* 2011). *Insulin-like growth factor 1a receptor* (*IGF1RA*)

386    is another gene that has a large number of sites under positive selection in *X. maculatus*. Overall,

387    96 sites within *IGF1RA* were shown to be under positive selection, with eight sites showing

388    changes in both poeciliids, while the remaining 88 were restricted to *X. maculatus* (Figure 6).

389    Protein lengths for *IGF1RA* vary among species. In our *P. prolifica* assembly, we have predicted

390    only 711 residues for this protein, but our sequence may be incomplete as it lacks a 3' UTR

391    region. Within *X. maculatus* where there is a complete predicted gene sequence (1,332 aa), these

392    96 sites account for ~7% of the gene sequence. Of the 96 sites, 48 are located within the Furin-

393    like domain of the protein, 36 are in one of the Receptor L-domains, one is in the Fibronectin

394    type III domain, and 11 are found outside of any known protein domain.

395        *P. prolifica* generally showed different genes under positive selection than *X. maculatus*

396    (Figure 4; only 17.1% of the 404 orthologs under positive selection showed positive selection in

397    both species). For example, *RAB11 family-interacting protein 4-like* (*RAB11FIP4*), one of the six

398    types of *RAB11 family-interacting proteins* found in fishes (Figure S5), has 16 sites under

399    positive selection in *P. prolifica*, but none in *X. maculatus*, while another member of that same

400    gene family, *RAB11 family-interacting protein 1-like* (*RAB11FIP1*), has 5 sites under positive

401    selection in *X. maculatus* and 1 in both *X. maculatus* and *P. prolifica* (the 2 species have

402    different residues at that site). These sites may be associated with novel functional domains

403    because each of these sites were identified as being extracellular for both *RAB11FIP1* and

404    *RAB11FIP4* using the transmembrane identification algorithm TMHMM; however, none of these

405    sites are located within any "known" functional domain (Figure 6). Patterns of gene evolution

406    across fish species show that the rapid gene evolution may be likely facilitated by multiple

407    incidences of gene duplication. Along with *IGF1RA*, *RAB11FIP* gene family members showed

408    family-wide evidence for gene duplication events and both *RAB11FIP* genes that were shown to

409    be under positive selection had expressed paralogs in the reference transcriptome sequence

410    (Figure S5). These duplications likely occurred after the whole genome duplication event

411  experienced by all fishes (Jaillon *et al.* 2004) since there is only one copy of each family member

412  found in the gar, *Lepisosteus oculatus*, (Figure S5) which has not undergone the teleost fish

413  whole genome duplication event.

414  **DISCUSSION**

415  We have developed the most thorough transcriptome reference for a placental fish to

416  date, providing a significant extension to earlier work in a sister taxa (Panhuis *et al.* 2011), in

417  order to better understand the genetics and evolution of placentation in fish. Our sequence

418  assembly has been extensively annotated for functional content and provides a solid foundation

419  for establishing genomic resources for this genus. Identified transcripts cover diverse functions

420  and, given the sampling of both poly-A selected and ribo-depleted RNAs across multiple

421  tissues, provide a comprehensive assessment of both protein-coding and non-coding RNA genes

422  organism-wide. In addition to its general descriptive characteristics, this transcriptome reference

423  has also provided us with important insights into the genetics of this placental species.

424  There appears to be parallels in placental evolution in eutherian mammals and *P.*

425  *prolifica*, highlighted by the extensive presence of expressed repetitive elements in fish MPC

426  tissues. Eutherian mammals often utilize repetitive element components as functional

427  contributions to placental and embryonic development, including endogenous retroviral envelope

428  proteins (Mi *et al.* 2000), DNA transposon regulatory machinery (Lynch *et al.* 2011), and/or *gag*

429  and *pol* domains of LTRs (Ono *et al.* 2001). A total of 98% of the transcripts associated with

430  retroelements exhibited high expression to either placental or embryonic tissues. These

431  transcripts included a variety of orthologous genes associated with placental function in

432  mammals, such as *PEG10*, an imprinted gene expressed in the placenta of mammals. The

433  extensive presence of progesterone signaling-related genes also parallels mammalian placental

434  function, particularly functions associated with the corpus luteum (Gemmell 1995) and decidual

435  cells (observed expression of *Protein DEPP* in fish MPC parallels that also described in

436  mammalian placental and embryonic tissues (Watanabe *et al.* 2005)). Alternatively, expressed

437  repetitive element transposases may be co-opted genes involved in more general gene regulation

438  as transcription factors or DNA-binding proteins with centromeric functional roles (Feschotte

439  2008). The identification of seemingly convergent gene expression of genetic elements of similar

440  type, but different lineage and an apparent implication in the function of the independently

16

441 derived placental tissues of fish and mammals leads to a hypothesis that similar molecular and
442 cellular adaptations are functioning in both systems.

443       There was also extensive evidence for placental tissue usage of novel genes and
444 transcripts as functional components specific to this family of fishes and, possibly, restricted to
445 this species. Tissue-specific patterns of high gene expression implicate many novel components
446 to be active in the MPC. Most of these novel, predicted transcripts lacked homology to genes in
447 any existing genetic resource, strengthening support for their designation as "novel". In total,
448 17.6% of the predicted protein sequences in the reference transcriptome could not be associated
449 with any existing reference sequence via exhaustive comparison to known protein and coding
450 sequence databases. Mis-assembly and/or chimeric reads could only explain a minority of these
451 "unknowns" as the depth of coverage was generally high for these genes and, as evidenced by
452 the clustering analysis, many of these transcripts are highly expressed. Many "unknowns" (~9%)
453 were found to contain repetitive elements or have sequence homology with non-coding RNAs,
454 implicating the co-option of rapidly evolving elements in the origins of this novel transcriptional
455 diversity. As our pairwise ortholog identification shows, the closer the phylogenetic species
456 comparison is, the greater the proportion of the transcriptome we could identify and annotate
457 (e.g. 14,761 orthologs were found in *X. maculatus* vs. 12,631 orthologs in *D. rerio* for a 16.9%
458 increase in the number of identified orthologs). Overall, novel transcripts would appear to be
459 significant contributors to placental function in *Poeciliopsis*. This prediction is also congruent
460 with mammalian placental systems, wherein many of the transcripts observed associated with
461 placental development and function are derived from lineage-specific co-option and
462 domestication of typically inactive retroelements (Emera and Wagner 2012).

463 **Role of Protein Evolution/Positive Selection**

464       Using our reference sequence, we identified genes under positive selection in both
465 matrotrophic (P. prolifica) and lecithotrophic (X. maculatus) species of livebearing poeciliid
466 fishes. Overall, genes identified as under positive selection did not disproportionally represent
467 any specific functional group, indicating that any genetic signal of adaptation identified in this
468 analysis covered a wide-array of functional components in the Poeciliidae. However, the
469 statistically significant differences in functional groups among the lecithotrophic *X. maculatus*,
470 and the matrotrophic *P. prolifica* undergoing positive selection indicate that there may be

471    selective bias in the types of genes contributing to the rapid evolution of placentation in this

472    group. The identification of genes related to biosynthesis and gene regulation, especially those

473    associated with DNA-binding in nuclear regions, are significantly over-represented in genes

474    under positive selection in our placental species. That these functional categories would be under

475    strong selective pressure is consistent with the inherent requirement for placental tissues to

476    develop quickly to support embryonic growth as well as the potential for parent-offspring

477    intragenomic conflict. The unexpectedly extensive protein-coding sequence evolution is highly

478    relevant to continued interest in the relative contribution of either changes at the protein-coding

479    level or those in gene regulation contributing to evolutionary patterns (Hoekstra and Coyne

480    2007; Lynch and Wagner 2008; Stern and Orgogozo 2008).

481          While the majority of genes under positive selection contain only a few sites that are

482    rapidly evolving, some genes exhibit evidence of surprisingly large regions of their coding

483    sequence under Darwinian selection. Evidence of gene duplication would appear to facilitate the

484    potential for positive selection. For example, while the insulin-like growth factor signaling axis

485    is a key regulator of embryogenesis and fetal growth in all vertebrates (Schlueter *et al.* 2007),

486    there is considerable redundancy in many of its components in fishes due to the ancient WGD

487    (Jaillon *et al.* 2004). Specifically, there are multiple copies of *insulin growth factor receptor 1*

488    (paralogs A and B). It has been established for the genus *Poeciliopsis* that *IGF2* has evolved

489    under positive selection that is hypothesized to be driven by parent-offspring conflict (O'Neill *et*

490    *al.* 2007). While *IGF2* is excluded in our analysis due to stringency filters (it has a large indel

491    region in *P. prolifica* and the HSP alignment region was too short for inclusion), *IGF1RA* (also

492    known to be expressed in fish gonadal tissues (Mei *et al.* 2014)) was shown to have extensive

493    evidence for rapid evolution in this group with eight sites evolving rapidly in all Poeciliidae and

494    84 sites under positive selection in just *X. maculatus*. These 92 sites cover both conserved

495    protein domains and unannotated regions of the protein. The signal for positive selection on both

496    *IGF2* and *IGF1RA* in poeciliids may reflect the opposing parent-specific expression (imprinting)

497    of *IGF2* and its antagonistic receptor *IGF2R* in mammals. However, if conflict is driving this

498    pattern, then it seems to be pushed to extremes in the non-placental *X. maculatus*, where *IGF2*

499    has also been shown to be under positive selection (Schartl *et al.* 2013). This would appear to

500    contradict assumptions of the hypothesis that parent-offspring conflict would be more extensive

501    in placental species (Zeh and Zeh 2008); alternatively, this observation may indicate that conflict

502     manifests itself differently in the absence of material exchange between mother and fetus. For

503     instance, selection may be acting on the duration of ovoviviparous development in *X. maculatus*,

504     where different paternal genomes compete for gestational space within the mother, while the

505     maternal genome dictates the length of her pregnancy and maximum occupancy of gestational

506     spaces.

507     While, speculatively, *IGF1RA* may exhibit evidence for selective pressure due to genetic

508     conflict in the lecithotrophic *X. maculatus*, it is possible that other biological processes specific

509     to viviparous reproduction are also under selection in *P. prolifica*. For example, the *RAB11FIP*

510     genes show lineage-specific patterns of protein evolution, indicating different selection pressures

511     in *P. prolifica* and *X. maculatus*. *RAB11FIP*-associated proteins are typically identified by the

512     presence of a C-terminal Rab-binding domain and are involved in vesicle transport and recycling

513     (Lindsay 2004), protein trafficking and sorting (Peden *et al.* 2004) and recycling of membranes

514     in cytokinesis (Wallace *et al.* 2002). It is unclear precisely why these genes are under positive

515     selection in these fish, but given their defined functions they may be responding to lineage-

516     specific selection pressure involving cellular transport related to the evolution of live-bearing.

517     Gene duplication is likely providing a considerable contribution to the potential for these genes

518     to undergo changes due to positive selection (e.g. Steinke *et al.* 2006). Just as in the *IGF1R*

519     genes, each of these two *RAB11FIP* genes has a closely related paralogous copy that showed no

520     evidence for positive selection (Figure S5).

521     Overall, our study demonstrates patterns of both sequence and functional convergence of

522     the poeciliid placenta with the therian mammalian placenta. In contrast to predictions that genetic

523     components would be distinct to the poeciliid lineage given the relatively recent convergent

524     derivation of the fish placenta from the pericardial sac and it highly dissimilar structural form,

525     many of the genetic components that contribute to mammalian placental development and

526     function are also involved in the fish placenta. While it could be predicted that at least some of

527     the types of genes involved in placentation in both lineages would be similar with respect to

528     cellular function and functional requirements of any placenta in maternal-fetal exchange, it is

529     notable that we find parallel evolutionary mechanisms, beyond such cohorts of genes, evident in

530     the co-option of retroelements and gene duplication as key contributors to the evolution of this

531     complex organ.

19

532

REFERENCES

Bressan, F. F., T. H. C. De Bem, F. Perecin, F. L. Lopes, C. E. Ambrosio *et al.*, 2009 Unearthing the roles of imprinted genes in the placenta. Placenta 30: 823–834.

Capella-Gutiérrez, S., J. M. Silla-Martínez, and T. Gabaldón, 2009 trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. Bioinformatics 25: 1972–1973.

Carter, A. M., and A. Mess, 2007 Evolution of the placenta in eutherian mammals. Placenta 28: 259–262.

Conesa, A., S. Götz, J. M. García-Gómez, J. Terol, M. Talón *et al.*, 2005 Blast2GO: A universal annotation and visualization tool in functional genomics research. Application note. Bioinformatics 21: 3674–3676.

van Dongen, S., and C. Abreu-Goodger, 2012 Using MCL to extract clusters from networks. Bact. Mol. Networks Methods Protoc. 281–295.

Edgar, R. C., 2004 MUSCLE: multiple sequence alignment with high accuracy and high throughput. Nucleic Acids Res. 32: 1792–1797.

Emera, D., and G. P. Wagner, 2012 Transposable element recruitments in the mammalian placenta: Impacts and mechanisms. Brief. Funct. Genomics 11: 267–276.

Feschotte, C., 2008 Transposable elements and the evolution of regulatory networks. Nat. Rev. Genet. 9: 397–405.

Gemmell, R. T., 1995 A comparative study of the corpus luteum. Reprod. Fertil. Dev. 7: 303–12.

Grabherr, M. G., B. J. Haas, M. Yassour, J. Z. Levin, D. A. Thompson *et al.*, 2011 Full-length transcriptome assembly from RNA-Seq data without a reference genome. Nat. Biotechnol. 29: 644–652.

Grove, B. D., and J. P. Wourms, 1991 The follicular placenta of the viviparous fish, Heterandria formosa. I. Ultrastructure and development of the embryonic absorptive surface. J. Morphol. 209: 265–284.

Haas, B. J., A. Papanicolaou, M. Yassour, M. Grabherr, D. Philip *et al.*, 2013 De novo transcript sequence recostruction from RNA-Seq: reference generation and analysis with Trinity. Nat. Protoc. 8: 1–43.

Hoekstra, H. E., and J. A. Coyne, 2007 The locus of evolution: Evo devo and the genetics of adaptation. Evolution (N. Y). 61: 995–1016.

Jaillon, O., J.-M. Aury, F. Brunet, J.-L. Petit, N. Stange-Thomann *et al.*, 2004 Genome duplication in the teleost fish Tetraodon nigroviridis reveals the early vertebrate proto-karyotype. Nature 431: 946–957.

572  John, K., V. Alla, C. Meier, and B. M. Pützer, 2011 GRAMD4 mimics p53 and mediates the
573        apoptotic function of p73 at mitochondria. Cell Death Differ. 18: 874–86.
574  Koerner, M. V., F. M. Pauler, R. Huang, and D. P. Barlow, 2009 The function of non-coding
575        RNAs in genomic imprinting. Development 136: 1771–1783.
576  Krogh, A., B. Larsson, G. von Heijne, and E. L. Sonnhammer, 2001 Predicting transmembrane
577        protein topology with a hidden Markov model: application to complete genomes. J. Mol.
578        Biol. 305: 567-580.
579  Lagesen, K., P. Hallin, E.A. Rødland, H.-H. Stærfeldt, T. Rognes, *et al.*, 2007 RNAmmer:
580        consistent and rapid annotation of ribosomal RNA genes. Nucleic Acids Res. 35: 3100-
581        3108.
582  Langmead, B., and S. L. Salzberg, 2012 Fast gapped-read alignment with Bowtie 2. Nat.
583        Methods 9: 357–9.
584  Lavialle, C., G. Cornelis, A. Dupressoir, C. Esnault, O. Heidmann *et al.*, 2013 Paleovirology of
585        "syncytins", retroviral env genes exapted for a role in placentation. Philos. Trans. R. Soc. B
586        Biol. Sci. 368: 20120507–20120507.
587  Li, H., and R. Durbin, 2010 Fast and accurate long-read alignment with Burrows-Wheeler
588        transform. Bioinformatics 26: 589–595.
589  Lindsay, A. J., 2004 The C2 domains of the class I Rab11 family of interacting proteins target
590        recycling vesicles to the plasma membrane. J. Cell Sci. 117: 4365–4375.
591  Liu, B., L. Fang, F. Liu, X. Wang, J. Chen *et al.*, 2015 Identification of Real MicroRNA
592        Precursors with a Pseudo Structure Status Composition Approach (H. Budak, Ed.). PLoS
593        One 10: e0121501.
594  Liu, H., J. Yin, M. Xiao, C. Gao, A. S. Mason *et al.*, 2012 Characterization and evolution of 5'
595        and 3' untranslated regions in eukaryotes. Gene 507: 106–111.
596  Love, M. I., W. Huber, and S. Anders, 2014 Moderated estimation of fold change and dispersion
597        for RNA-seq data with DESeq2. Genome Biol. 15: 550.
598  Lynch, V. J., R. D. Leclerc, G. May, and G. P. Wagner, 2011 Transposon-mediated rewiring of
599        gene regulatory networks contributed to the evolution of pregnancy in mammals. Nat.
600        Genet. 43: 1154–1159.
601  Lynch, V. J., and G. P. Wagner, 2008 Resurrecting the role of transcription factor change in
602        developmental evolution. Evolution (N. Y). 62: 2131–2154.
603  Maltepe, E., A. I. Bakardjiev, and S. J. Fisher, 2010 The placenta: transcriptional, epigenetic, and
604        physiological integration during development. J. Clin. Invest. 120: 1016–1025.
605  Martin, M., 2011 Cutadapt removes adapter sequence from high-throughput sequencing reads.
606        EMBnet.journal 17: 10–12.
607  Mei, J., W. Yan, J. Fang, G. Yuan, N. Chen *et al.*, 2014 Identification of a gonad-expression
608        differential gene insulin-like growth factor-1 receptor (Igf1r) in the swamp eel (Monopterus
609        albus). Fish Physiol. Biochem. 40: 1181–1190.
610  Meredith, R. W., J. E. Janecka, J. Gatesy, O. A. Ryder, C. A. Fisher *et al.*, 2011 Impacts of the
611        Cretaceous Terrestrial Revolution and KPg extinction on mammal diversification. Science
612        334: 521–524.
613  Mi, S., X. Lee, X. Li, G. M. Veldman, H. Finnerty *et al.*, 2000 Syncytin is a captive retroviral
614        envelope protein involved in human placental morphogenesis. Nature 403: 785–789.
615  Morgante, M., M. Hanafey, and W. Powell, 2002 Microsatellites are preferentially associated
616        with nonrepetitive DNA in plant genomes. Nat. Genet. 30: 194–200.
617  Neilsen, H., 2017 Predicting Secretory Proteins with SignalP., pp. 59-73 in *Methods in*

618      *molecular biology (Clifton, N.J.)*, United States.

619    O'Neill, M. J., B. R. Lawton, M. Mateos, D. M. Carone, G. C. Ferreri *et al.*, 2007 Ancient and
620      continuing Darwinian selection on insulin-like growth factor II in placental fishes. Proc.
621      Natl. Acad. Sci. 104: 12404–12409.

622    Ono, R., S. Kobayashi, H. Wagatsuma, K. Aisaka, T. Kohda *et al.*, 2001 A Retrotransposon-
623      Derived Gene, PEG10, Is a Novel Imprinted Gene Located on Human Chromosome 7q21.
624      Genomics 73: 232–237.

625    Panhuis, T. M., G. Broitman-Maduro, J. Uhrig, M. Maduro, and D. N. Reznick, 2011 Analysis of
626      expressed sequence tags from the placenta of the live-bearing fish Poeciliopsis
627      (Poeciliidae). J. Hered. esr002.

628    Peden, A., E. Schonteich, J. Chun, J. Junutula, R. Scheller *et al.*, 2004 The RCP–Rab11 Complex
629      Regulates Endocytic Protein Sorting. Mol. Biol. Cell 15: 3751–3737.

630    Pires, M. N., K. E. McBride, and D. N. Reznick, 2007 Interpopulation variation in life-history
631      traits of Poeciliopsis prolifica: implications for the study of placental evolution. J. Exp.
632      Zool. A. Ecol. Genet. Physiol. 307: 113–125.

633    Pollux, B. J. A., M. N. Pires, A. I. Banet, and D. N. Reznick, 2009 Evolution of Placentas in the
634      Fish Family Poeciliidae: An Empirical Study of Macroevolution. Annu. Rev. Ecol. Evol.
635      Syst. 40: 271–289.

636    Ranwez, V., S. Harispe, F. Delsuc, and E. J. P. Douzery, 2011 MACSE: Multiple Alignment of
637      Coding SEquences accounting for frameshifts and stop codons. PLoS One 6: e22594.

638    Renfree, M. B., S. Suzuki, and T. Kaneko-Ishino, 2013 The origin and evolution of genomic
639      imprinting and viviparity in mammals. Phil. Trans. R. Soc. B 368: 20120151.

640    Reznick, D., 1981 "Grandfather Effects": The Genetics of Interpopulation Differences in
641      Offspring Size in the Mosquito Fish. Evolution (N. Y). 35: 941–953.

642    Reznick, D. N., M. Mateos, and M. S. Springer, 2002 Independent origins and rapid evolution of
643      the placenta in the fish genus Poeciliopsis. Science (80-. ). 298: 1018–1020.

644    Roberts, A., and L. Pachter, 2013 Streaming fragment assignment for real-time analysis of
645      sequencing experiments. Nat. Methods 10: 71–73.

646    Schartl, M., R. B. Walter, Y. Shen, T. Garcia, J. Catchen *et al.*, 2013 The genome of the
647      platyfish, Xiphophorus maculatus, provides insights into evolutionary adaptation and
648      several complex traits. Nat. Genet. 45: 567–572.

649    Schlueter, P. J., G. Peng, M. Westerfield, and C. Duan, 2007 Insulin-like growth factor signaling
650      regulates zebrafish embryonic growth and development by promoting cell survival and cell
651      cycle progression. Cell Death Differ. 14: 1095–1105.

652    Steinke, D., W. Salzburger, I. Braasch, and A. Meyer, 2006 Many genes in fish have species-
653      specific asymmetric rates of molecular evolution. BMC Genomics 7: 20.

654    Stern, D. L., and V. Orgogozo, 2008 The loci of evolution: how predictable is genetic evolution?
655      Evolution (N. Y). 62: 2155–2177.

656    Theocharidis, A., S. van Dongen, A. J. Enright, and T. C. Freeman, 2009 Network visualization
657      and analysis of gene expression data using BioLayout Express(3D). Nat. Protoc. 4: 1535–
658      1550.

659    Turner, C. L., 1940 Pseudoamnion, pseudochorion, and follicular pseudoplacenta in poeciliid
660      fishes. J. Morphol. 67: 59–89.

661    Ulitsky, I., A. Shkumatava, C. H. Jan, H. Sive, and D. P. Bartel, 2011 Conserved function of
662      lincRNAs in vertebrate embryonic development despite rapid sequence evolution. Cell 147:
663      1537–1550.

664  Wallace, D. M. E., A. J. Lindsay, A. G. Hendrick, and M. W. McCaffrey, 2002 Rab11-FIP4
665       interacts with Rab11 in a GTP-dependent manner and its overexpression condenses the
666       Rab11 positive compartment in HeLa cells. Biochem. Biophys. Res. Commun. 299: 770–
667       779.
668  Watanabe, H., K. Nonoguchi, T. Sakurai, T. Masuda, K. Itoh *et al.*, 2005 A novel protein Depp,
669       which is induced by progesterone in human endometrial stromal cells activates Elk-1
670       transcription factor. Mol. Hum. Reprod. 11: 471–476.
671  Wren, J. D., E. Forgacs, J. W. Fondon, A. Pertsemlidis, S. Y. Cheng *et al.*, 2000 Repeat
672       Polymorphisms within Gene Regions: Phenotypic and Evolutionary Implications. Am. J.
673       Hum. Genet. 67: 345–356.
674  Yang, Z., 2007 PAML 4: Phylogenetic analysis by maximum likelihood. Mol. Biol. Evol. 24:
675       1586–1591.
676  Zeh, J. A., and D. W. Zeh, 2008 Viviparity-driven Conflict. Ann. N. Y. Acad. Sci. 1133: 126–
677       148.
678
679
680

681

682    **Figure Legends:**

683    **Figure 1.** Distributions of various transcriptome annotations for *Poeciliopsis prolifica* reference

684    transcriptome predicted transcripts (blue) and alternatively-spliced variant groups, representing

685    "genes" (red). 140,709 transcripts (29.4% of total) exhibited identifiable homology (e-value < 1

686    x $10^{-5}$) with protein reference sequences in the NCBI *nr* database and another 29,199 (6.1%)

687    transcripts showed similarity (e-value < 1 x $10^{-5}$) with nucleotide reference sequence in the NCBI

688    *nt* database. 16,277 (11.6%) and 6,772 (4.8%) transcripts are associated with transmembrane

689    (TMHMM) and signaling (SignalP) proteins. 8,181 showed greater than 70% coverage of known

690    UniProtKB/Swiss-Prot orthologs; 3,785 transcripts were identified as containing the complete

691    ORFs of conserved orthologs in UniProtKB/Swiss-Prot database (Table S2).

692    **Figure 2.** Level 2 gene ontology term distributions for reference transcriptome of *Poeciliopsis*

693    *prolifica*.

694    **Figure 3.** Three-dimensional gene atlas derived from gene expression data for maternal

695    placental/ovarian complex (MPC), late-stage embryonic, brain, and liver tissue. Proximity in

696    space indicates similarity in gene expression profile across tissues. Clusters were defined using

697    MCL clustering algorithm on highly expressed genes (>50 FPKM in at least on tissue type) from

698    Roche 454 RNA-seq. Clusters 1-4 are mostly, though not exclusively, made up of transcripts that

699    are tissue-specifically expressed, while clusters 5-9 consist of transcripts that are highly

700    expressed across multiple tissues. Each of these clusters (1-9), had 2,940, 1,734, 1,638, 373, 65,

701    30, 28, 25, and 5 members, respectively.

702    **Figure 4.** Venn diagrams showing patterns of shared and unshared proteins and sites within

703    protein under positive selection among the three foreground taxon groupings tested with PAML.

704    **Figure 5.** Distribution of GO Terms that were differentially represented in genes identified to be

705    under positive selection in the matrotrophic/placental(PL) *Poeciliopsis prolifica* and

706    lecithotrophic(LC) *Xiphophorus maculatus*. GO terms include categories from all three main

707    ontologies (Biological Processes; Molecular Functions; Cellular Components).

24

708    **Figure 6.** Diagrams of *insulin growth factor-1 receptor-A* (*IGF1RA*) from *Xiphophorus*

709    *maculatus*, and *RAB11 family-interacting protein 1-lik*e and *4-like* (*RAB11FIP1* and

710    *RAB11FIP4*, respectively) from *Poeciliopsis prolifica* showing known protein domains, indel

711    regions among species (identified using regions of multiple sequence alignment), and sites

712    identified as being under positive selection from PAML analysis in live-bearing poeciliids, the

713    lecithotrophic *Xiphophorus maculatus*, or the matrotrophic *Poeciliopsis prolifica*.
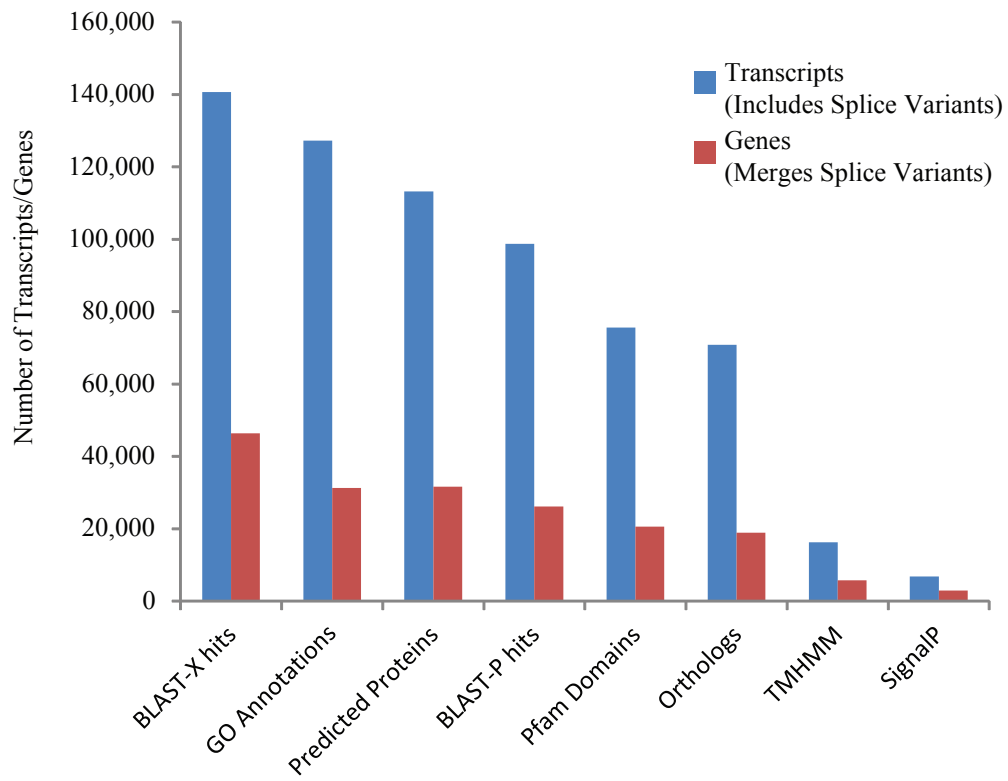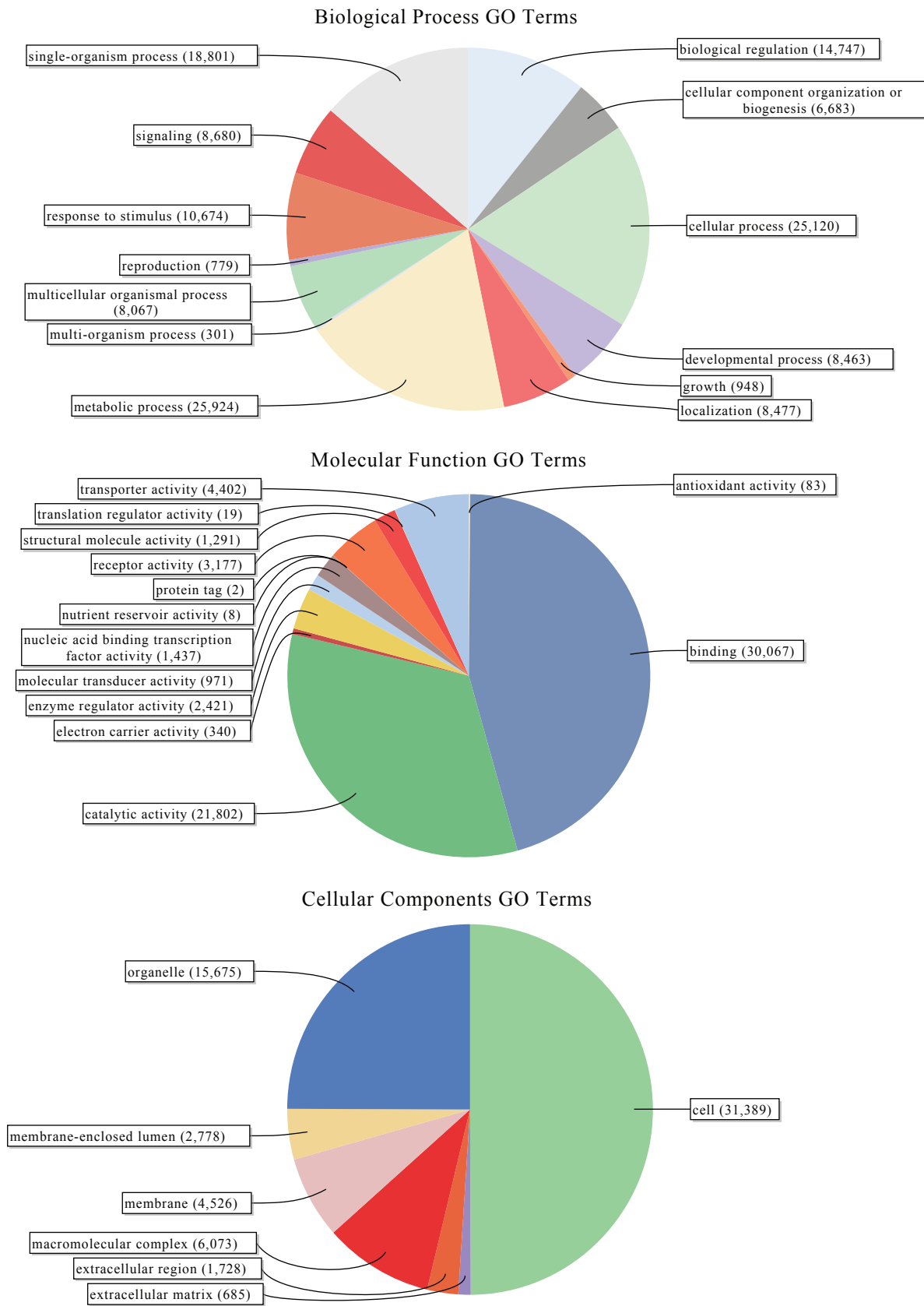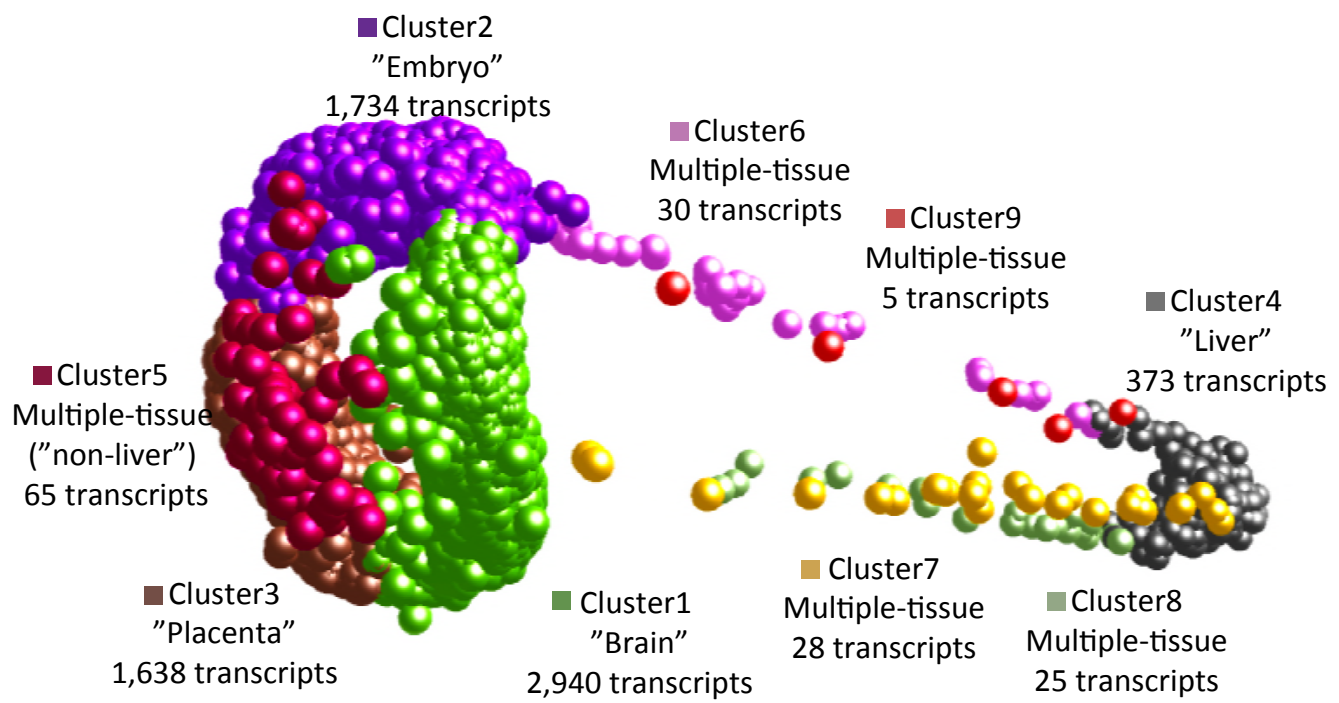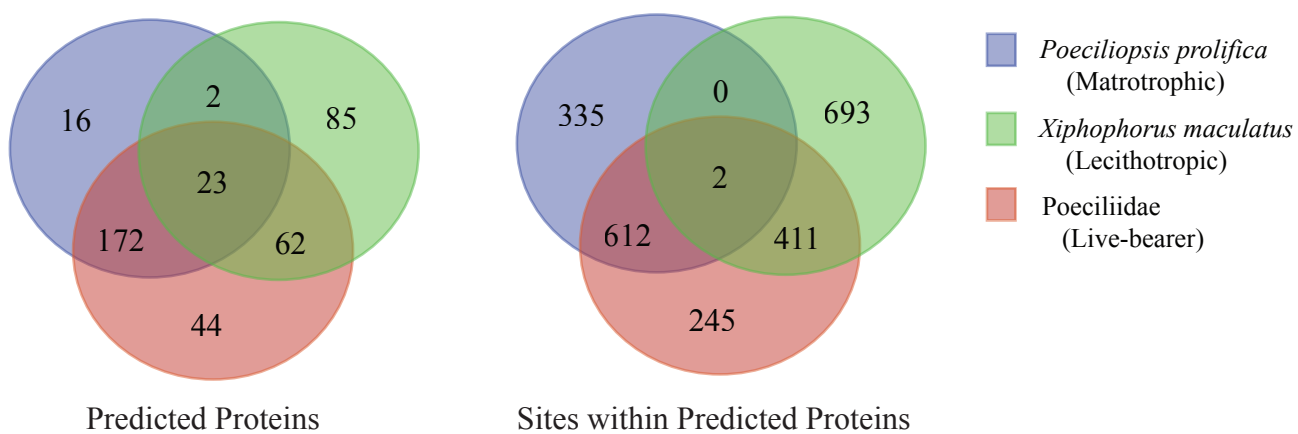
Figure 1.

## Biological Process GO Terms

single-organism process (18,801)
signaling (8,680)
response to stimulus (10,674)
reproduction (779)
multicellular organismal process (8,067)
multi-organism process (301)
metabolic process (25,924)
biological regulation (14,747)
cellular component organization or biogenesis (6,683)
cellular process (25,120)
developmental process (8,463)
growth (948)
localization (8,477)

## Molecular Function GO Terms

transporter activity (4,402)
translation regulator activity (19)
structural molecule activity (1,291)
receptor activity (3,177)
protein tag (2)
nutrient reservoir activity (8)
nucleic acid binding transcription factor activity (1,437)
molecular transducer activity (971)
enzyme regulator activity (2,421)
electron carrier activity (340)
antioxidant activity (83)
binding (30,067)
catalytic activity (21,802)

## Cellular Components GO Terms

organelle (15,675)
membrane-enclosed lumen (2,778)
membrane (4,526)
macromolecular complex (6,073)
extracellular region (1,728)
extracellular matrix (685)
cell (31,389)

Figure 2.

Figure 4.

Figure 5.

**Insulin Growth Factor-1 Receptor-A (IGF1RA)**

**RAB11 Family-Interacting Protein 1-like (RAB11FIP1)**

**RAB11 Family-Interacting Protein 4-like (RAB11FIP4)**

Legend:
- Multiple Sequence Alignment Region
- Among Species Indel
- Positive selection in *Xiphophorus maculatus*
- Positive selection in both Poeciliids
- Positive selection in *Poeciliopsis prolifica*
- Receptor L-domain
- Fibronectin type-III domain
- Protein tyrosine kinase
- C2 RAB11FIP1 Class 1 domain
- Transmembrane domain
- Signaling peptide
- Domain of unknown function (DUF4201)
- Coiled coil

Figure 6.

Figure S1. **A and B.** Intact ovary removed from gravid female. Scale bar = 0.5mm. **B.** Outlines of the different embryos within the ovary shown in **A. a.** maternal/placental ovarian tissue complex (MPC) with late stage embryo removed, **b.** Very Late-eyed stage embryo (i.e. nearly full term, Stage 6), **c.** Early-eyed stage embryo (Stage 3), **d.-f.** Late-eyed stage embryo (Stage 5).
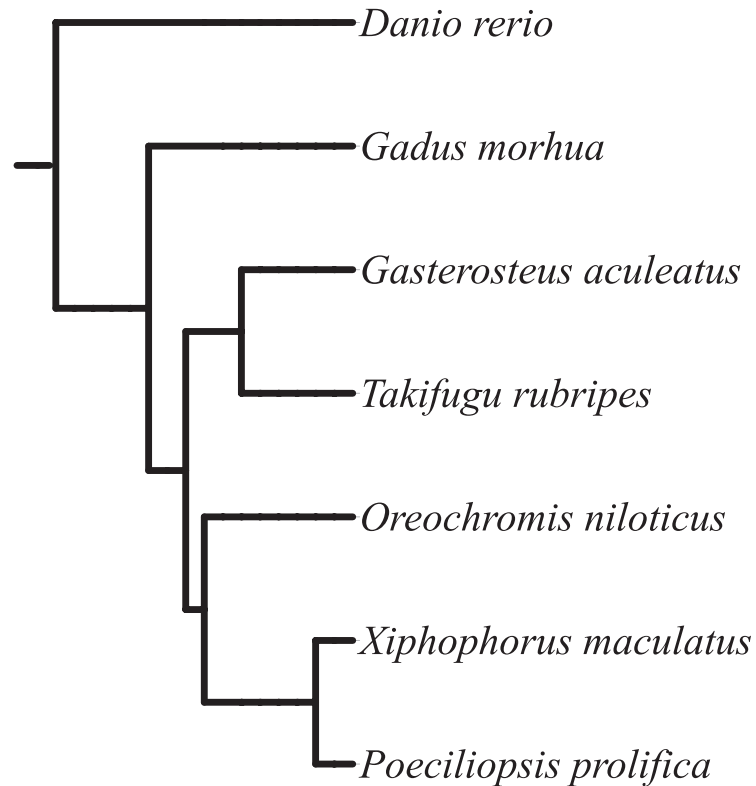
Figure S2. Species tree used in PAML. Tested foregrounds of live-bearing Poeciliids generally, lecithotrophic live-bearers, and matrotrophic live-bearers using the clade of Poeciliids, the *Xiphophorus maculatus* lineage, and the *Poeciliopsis prolifica* lineage, respectively for each case.
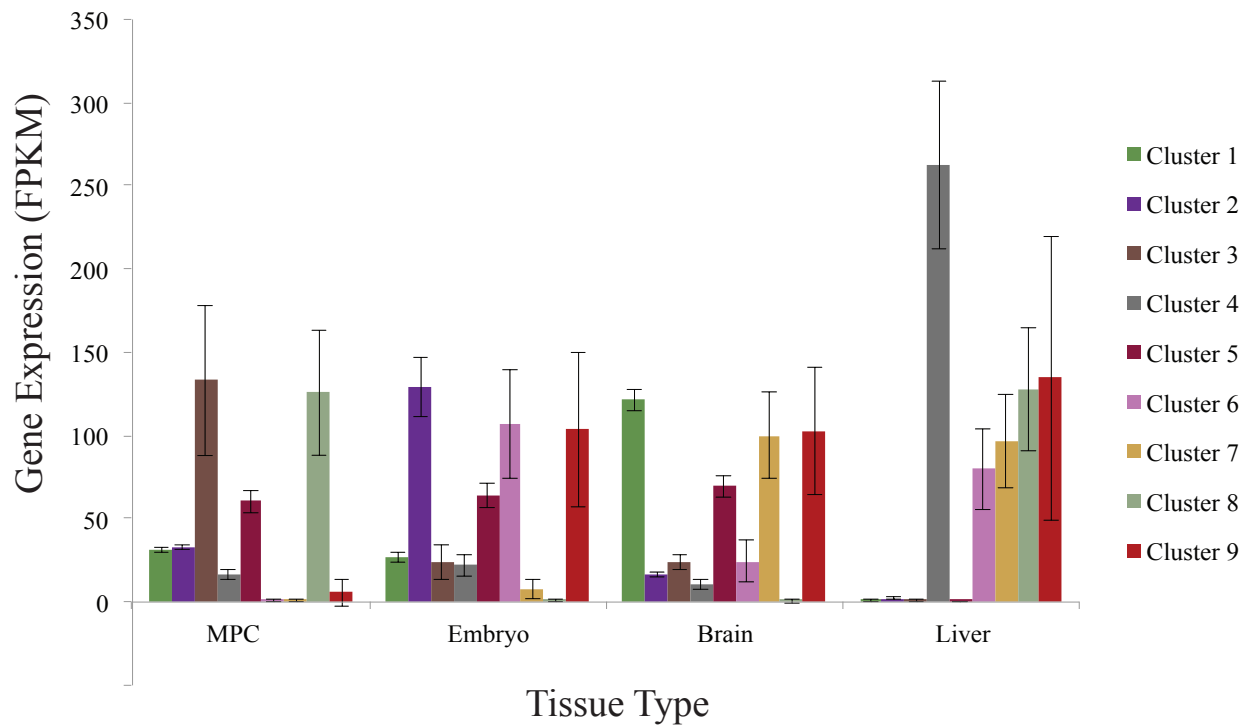
Figure S3. Average expression in fragments per kilobase per million base pairs (FPKM) in each tissue grouped by cluster genes indicating which clusters are associated with which tissues. Standard error bars shown. MPC indicates maternal placental/ovarian complex tissue sample.
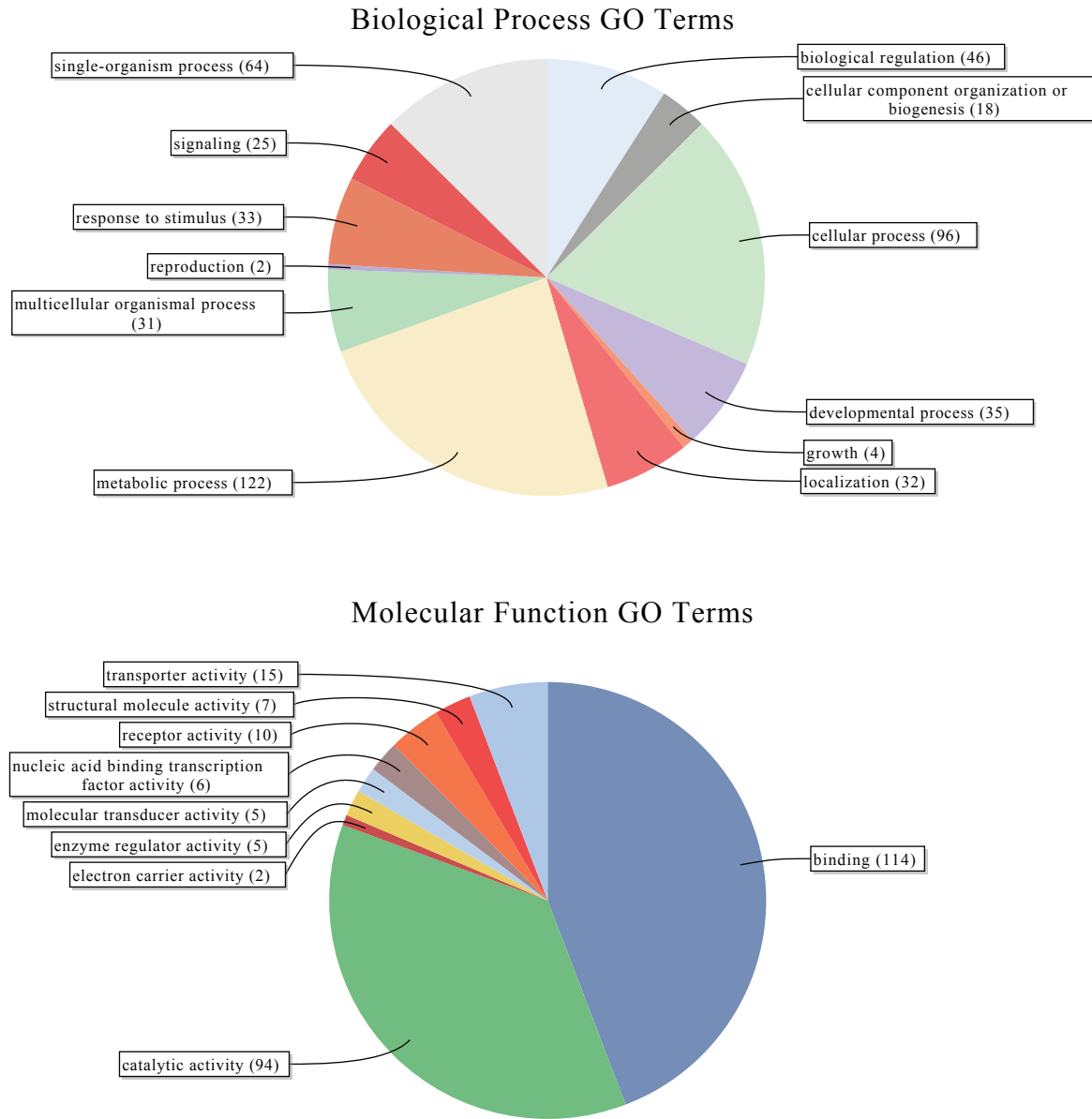
Figure S4. Level 2 Biological Process and Molecular Function Gene Ontology terms associated with gene identified as having sites under positive selection in poeciliids.
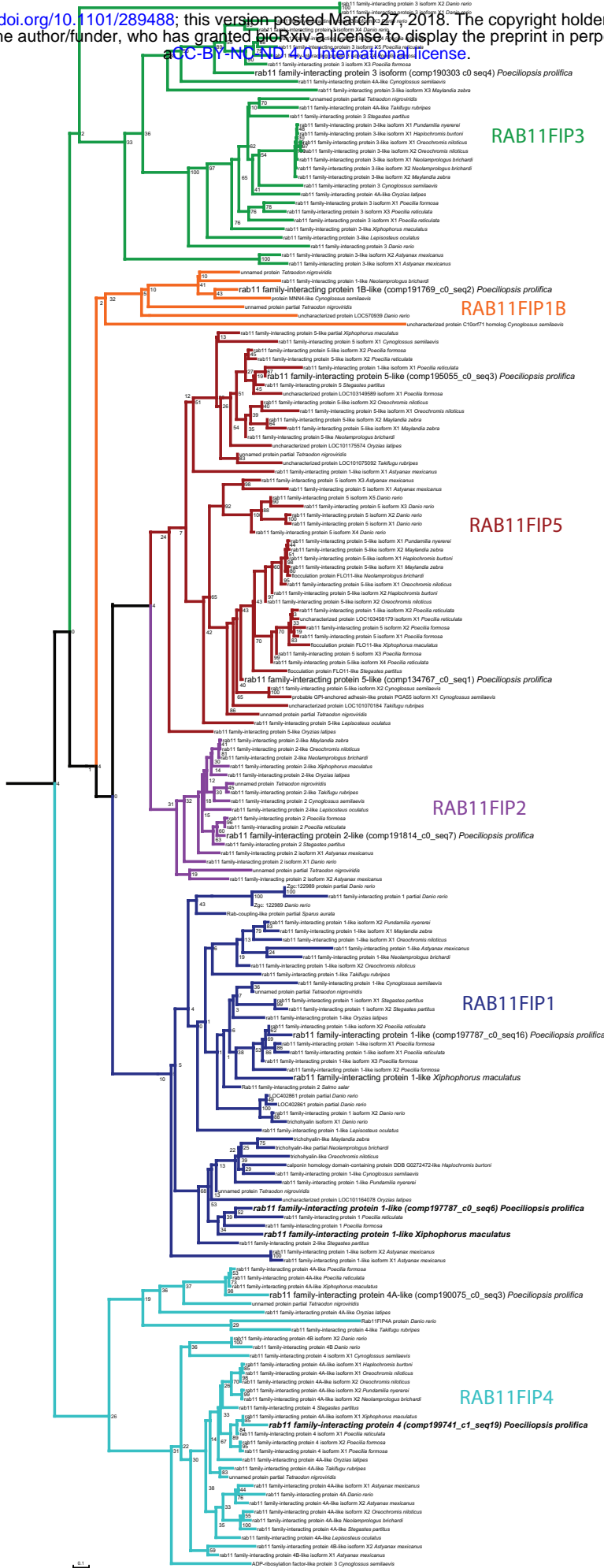
Figure S5. Phylogenetic distance-based gene-family tree for *RAB11 family-interacting proteins* (RAB11FIPs) in fishes. Each color represents different gene family member protein group. RAB11FIPs found in Poeciliopsis prolifica in enlarged fonts. Proteins with sites found to be under positive selection in PAML analysis in bold italics; specifically, RAB11FIP1 in *P. prolifica* and *X. maculatus* and RAB11FIP4 in *P. prolifica* only.