

1 **The germline-specific region of the sea lamprey genome plays a key role in spermatogenesis**

2 Tamanna Yasmin<sup>1</sup>, Phil Grayson<sup>1,2</sup>, Margaret F. Docker<sup>1</sup> and Sara V. Good<sup>1,3</sup>

3

4 <sup>1</sup>Department of Biological Sciences, University of Manitoba, Winnipeg, Manitoba R3T 2N2,  
5 Canada

6 <sup>2</sup>Present Address: Department of Biomedical Informatics, Harvard Medical School, Boston,  
7 Massachusetts, USA, 02115

8 <sup>3</sup>Department of Biology, University of Winnipeg, Winnipeg, Manitoba R3B 2E9, Canada

9

10 **Abstract**

11 The sea lamprey genome undergoes programmed genome rearrangement (PGR) in which ~20%  
12 is jettisoned from somatic cells soon after fertilization. Although the role of PGR in embryonic  
13 development has been studied, the role of the germline-specific region (GSR) in gonad  
14 development is unknown. We analysed RNA-sequence data from 28 sea lamprey gonads  
15 sampled across life-history stages, generated a genome-guided *de novo* superTranscriptome with  
16 annotations, and identified genes in the GSR. We found that the 638 genes in the GSR are  
17 enriched for reproductive processes, exhibit 36x greater odds of being expressed in testes than  
18 ovaries, show little evidence of conserved synteny with other chordates, and most have putative  
19 paralogues in the GSR and/or somatic genomes. Further, several of these genes play known roles  
20 in sex determination and differentiation in other vertebrates. We conclude that the GSR of sea  
21 lamprey plays an important role in testicular differentiation and potentially sex determination.

22

## 23 **Introduction**

24 The genetic structure and composition of germline and somatic cells typically remain constant  
25 throughout an organism's life span. However, under some conditions (e.g., cancer) or in some  
26 taxa, the genetic composition of cells varies by type and/or developmental stage<sup>1,2</sup>. Included in  
27 this is the unusual process of programmed genome rearrangement (PGR), in which either  
28 portions of chromosomes (chromosomal diminution) or entire chromosomes (chromosomal  
29 elimination) are removed during embryonic development, thereby reducing the genomic content  
30 of descendent cells by up to 90%<sup>3</sup>. Although the frequency of PGR across metazoans is  
31 unknown, it has been observed in more than 100 vertebrate and invertebrate species from nine  
32 major taxonomic groups<sup>3</sup>, including in lampreys<sup>4-7</sup>. In sea lamprey (*Petromyzon marinus*), flow  
33 cytometric measurements of DNA content in the germline (testes) vs. somatic (blood) cells  
34 indicate that ~20% (~500 Mb) of the germline genome is eliminated during PGR<sup>8</sup>. Further  
35 studies have shown that PGR in sea lamprey, which occurs ~3 days post-fertilization (dpf),  
36 shares conserved features with PGR in other agnathan lineages<sup>4-7</sup>. This event involves  
37 chromosomal elimination of repetitive and single-copy sequences and is enriched for genes  
38 involved in development or germline maintenance<sup>6,9</sup>. However, further research on the possible  
39 function of the germline-specific regions (GSR) in gonad development is needed.

40 Many hypotheses have been posited regarding the biological significance of PGR,  
41 including gene silencing, dosage compensation, position effects on gene expression, germline  
42 development, and sex determination<sup>1,10-13</sup>. In sea lamprey, it has been suggested that PGR  
43 permits the expression of genes beneficial to the germline during the early stages of embryonic  
44 development<sup>6,8</sup>, consistent with the high levels of gene silencing observed for genes in the  
45 GSR<sup>3,9</sup>. In the zebra finch (*Taeniopygia guttata*), chromosomal diminution of a germline-

46 restricted chromosome (GRC) occurs during early embryonic development; the genes in the  
47 GRC have higher expression in the ovary than the testis, and the GRC is later eliminated from  
48 mature sperm, being transmitted only through the oocytes<sup>14,15,16</sup>.

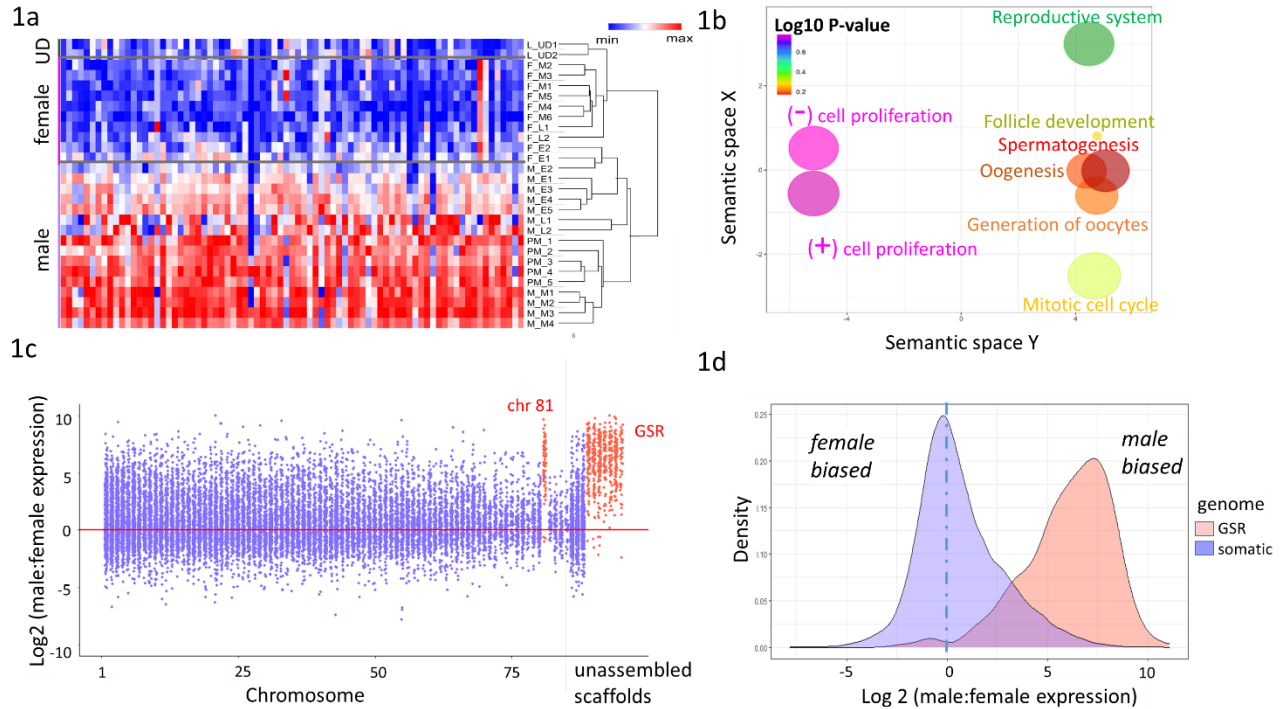
49 In sexually reproducing taxa without PGR, primordial germ cells (PGCs) are formed  
50 early during embryogenesis. PGCs typically develop through the coordination of three  
51 developmental cues: suppression of ongoing somatic differentiation, repression of DNA  
52 methylation, and inhibition of cell proliferation<sup>17</sup>. Once defined, PGCs subsequently exhibit  
53 tightly coordinated gene expression that leads to germ cell development and differentiation in  
54 both sexes. In lampreys, however, the germline cells are specified at fertilization, and somatic  
55 cell delineation occurs afterward, ~3 dpf, when PGR is initiated<sup>5</sup>. This intriguing reversal of  
56 events is heightened by the ongoing enigma of their sex determination. Lampreys do not have  
57 heteromorphic sex chromosomes and there is no evidence to date of genomic differences  
58 between males and females; sex may be determined by genetic factors in the germline genome,  
59 environmental factors, or a combination of the two (reviewed by<sup>18</sup>).

60 Here, we used RNA-sequence (RNA-seq) data from 28 sea lamprey gonads sampled at  
61 different life-history stages and in both sexes to generate a gonadal superTranscriptome and  
62 examined the function, expression, and evolutionary relationships of sex-biased genes,  
63 particularly in the GSR. We identified 638 germline-specific genes (GSGs), many of which were  
64 present in multiple germline-specific paralogues pertaining to 163 unique gene names that were,  
65 overall, very highly expressed during spermatogenesis, but lowly expressed during oogenesis  
66 and in undifferentiated larvae. The observation that the genes in the GSR appear to be present in  
67 undifferentiated larvae and females but are expressed at low levels suggests that the male-  
68 specific expression is due to regulatory changes, as opposed to there being a male-specific

69 germline sequence. Further, we found that ~55% of the GSGs also have paralogous copies in the  
70 somatic genome and ~19% have putative orthologues in other taxa including, most importantly,  
71 a core set of conserved genes involved in sex determination and spermatogenesis. Using publicly  
72 available RNA-seq data from 1–5 dpf embryos, we found that the genes expressed during  
73 gonadogenesis are either not expressed or lowly expressed during early embryo formation.  
74 Collectively, these results suggest that a major role of the GSR is in testicular differentiation and  
75 probably sex determination. PGR in sea lamprey may serve to reduce conflict of genes under  
76 sexual selection, a hypothesis further supported by the highly duplicated nature of genes in the  
77 GSR and their association in sexual differentiation and determination pathways in other taxa  
78 (see<sup>3</sup>).

## 79 **Results and discussion**

80 **GSGs show predominantly male-biased expression and have a key role in gametogenesis.** We  
81 used RNA-seq data from 28 sea lamprey gonads sampled across a range of developmental stages  
82 to generate a gonadal superTranscriptome using the Necklace pipeline<sup>19</sup>. Stages included  
83 undifferentiated larvae, female larvae following the onset of oogenesis and sexually mature (adult)  
84 females, prospective male larvae (i.e., those in which the gonad was still histologically  
85 undifferentiated but which were beyond the size at which ovarian differentiation is complete),  
86 males undergoing testicular differentiation following the onset of metamorphosis, and sexually  
87 mature (adult) males (see Supplementary Fig. 1 and Supplementary Table 1). This revealed a large  
88 number of genes that were highly expressed during male but not female gonad development; these  
89 genes were physically linked and mapped to chromosome 81 and many unplaced scaffolds based  
90 on the Vertebrate Genome Project (VGP) reference assembly. Thus, we sought to define which of  
91 the genes in our gonadal superTranscriptome mapped to the GSR.



92 **Fig. 1:** Identification of genomic location of the germline-specific genes (GSGs) in the VGP  
 93 genome and their expression in the sea lamprey gonadal samples used in this study. **a)** Heatmap  
 94 showing GSG expression pattern of all samples used in the study; UD stands for undifferentiated  
 95 larvae (see Supplementary Fig. 3 for full heatmap) **b)** Gene ontology term enrichment analysis of  
 96 the GSGs where colours indicate the  $\log_{10}$  of the false discovery rate-corrected  $P$ -value  
 97 (PANTHER overrepresentation test, with a Fisher exact test for significance and filtering using a  
 98 false discovery rate of 0.05); circle size denotes fold enrichment above expected values. **c)**  
 99 Scatterplot showing the  $\log_2(\text{male:female})$  normalized gene expression across all chromosomes  
 100 and concatenated scaffolds in the VGP assembly of the sea lamprey genome; regions identified as  
 101 belonging to the GSR are coloured in red, while those in the somatic genome are coloured in blue.  
 102 **d)** Density plot of the  $\log_2(\text{male:female})$  ratio of normalized gene expression. When  $x = 0$ , average  
 103 expression across all females = males. For genes in the somatic genome, the density peaks at  $\sim x =$   
 104 0 but is right skewed, while for genes in the GSR, the density peaks at  $\sim x = 7.5$ , showing that genes  
 105 in the GSR are male-biased.

106 The GSR in sea lamprey was identified for an earlier release of the sea lamprey germline  
 107 assembly ([www.stowers.org](http://www.stowers.org))<sup>20</sup>. Thus, we used a modified version of the DifCover pipeline used  
 108 for that analysis<sup>21</sup> to define the coordinates of GSR in the VGP assembly. Accordingly, GSRs  
 109 were designated as regions in which the read coverage of sperm DNA was  $> 2$ -fold more than  
 110 the read coverage of blood DNA. Based on the VGP reference genome, a total of 5253 genomic  
 111 intervals were mapped by DifCover, of which 919 segments had an enrichment score

112 (log<sub>2</sub>(standardized sperm coverage/blood coverage) greater than 2. The total span of the GSR-  
113 inferred regions consisted of more than 27 Mbps (Supplementary Table 2, Supplementary Fig.  
114 2).

115 We then used the segment enrichment scores to assign genes from our gonadal  
116 superTranscriptome to either the GSR or somatic genomes. Using an earlier scaffold-based  
117 assembly of the sea lamprey germline genome (available at SIMRbase), Smith et al., (2018)  
118 identified ~13Mbps including 356 protein-coding genes in the GSR<sup>20</sup>. On the other hand, using  
119 the VGP assembly which consists of 85 chromosomes and 1195 unassembled scaffolds, we  
120 assigned the entirety of chromosome 81 as well as 177 scaffolds (Supplementary Fig. 3) to the  
121 GSR, while the remaining 84 chromosomes and 1018 scaffolds were not germline-enriched,  
122 suggesting that they are found in the somatic genome. In total, 638 genes from our gonadal  
123 superTranscriptome mapped to the GSR; these 638 genes corresponded to only 163 unique gene  
124 names based on our combined Trinotate and reference genome annotation pipeline  
125 (Supplementary Table 3), with approximately half of the GSGs occurring in a single copy but the  
126 other half occurring in 2–77 duplicated copies (Supplementary Table 4, Supplementary Fig. 4).  
127 Importantly, however, none of the GSR enriched scaffolds or chromosomes showed overlap  
128 between the GSR and the somatic regions. This supports previous work that determined that  
129 PGR in lampreys is more likely to involve chromosome elimination than diminution<sup>9</sup>.

130 The expression analysis of the GSGs revealed that out of 638 GSGs, 409 genes (64% of  
131 the genes in the GSR) are moderately to highly expressed in one or more stages of the  
132 developing testis, but only a few GSGs are expressed in undifferentiated larval gonads (Fig. 1a,  
133 see Supplementary Fig. 5 for full heatmap). Functional enrichment analysis of the GSGs from  
134 the gonadal superTranscriptome indicate that they are involved in 26 pathways of which *wnt*

135 signaling and *E-Cadherin* signaling pathways each represented 16.4% of the total hits  
136 (Supplementary Fig. 6). Other critical pathways include the insulin/insulin growth factor (*igf*)  
137 pathway, gonadotropin releasing hormone receptor (*gnrhr*) pathway, transforming growth factor  
138 beta (*tgfb*) signaling pathway, and fibroblast growth factors (*fgf*) pathway, which contained  
139 2.7%, 2.7%, 2.7%, and 1.4% of all hits, respectively (Supplementary Fig. 6). Next, we analyzed  
140 the GO terms associated with genes in the GSR to obtain further insight into their molecular  
141 function (Supplementary Table 5). Using an overrepresentation test, we found that the highest  
142 FDR terms were associated with reproductive system development, positive and negative  
143 regulation of cell population proliferation, ovarian follicle development, oogenesis and  
144 spermatogenesis. Collectively, this demonstrates that the functional ontology of GSGs is  
145 enrichment for GO terms related to reproductive developmental processes (Fig. 1b,  
146 Supplementary Table 6).

147 PGR has been proposed as a mechanism to reduce conflict between the somatic and  
148 germline genomes during early embryogenesis. Bryant et al. (2016) identified that the genes  
149 eliminated during PGR are expressed throughout lamprey embryogenesis and found ontological  
150 overrepresentation of these genes in germline development and oncogenesis<sup>4</sup>. However, PGR is  
151 closely tied to PGC specification in sea lamprey since germ cells are those that do not undergo  
152 PGR. Thus, we hypothesized that genes in the GSR might play roles in both early embryogenesis  
153 as well as gonadal differentiation and/or development. To address this, we assessed global  
154 differences in expression of genes in the GSR vs. somatic genome during gonadal development  
155 across differentiated ovaries and testes sampled in early, mid, and late developmental stages as  
156 well as in undifferentiated larvae and prospective males prior to testicular differentiation (See  
157 Supplementary Fig. 1 and Supplementary Table 1). This revealed the surprising result that

158 almost all of the genes in the GSR exhibit male-biased expression during gonad development  
159 (Figs. 1c, 1d, Supplementary Fig. 7), while genes in the somatic genome were, overall, equally  
160 likely to be expressed in the female or male gonad (as expected): i.e., the density of the  
161 male:female gene expression ratio peaks at  $x = 0$  for genes in the somatic genome (Fig. 1d). A  
162 possible explanation for this observation could be that females do not have the same GSR as  
163 males, since the reference genome for sea lamprey was generated using sperm DNA. To examine  
164 this possibility, we aligned individual BAM files from both male and female gonad samples to  
165 the indexed superTranscriptome and annotation file using the Integrative Genome Viewer  
166 (IGV)<sup>22</sup>. This revealed 410 transcripts from female gonad samples that mapped to either known  
167 or novel exons in the GSR (Supplementary Fig. 8a–8b). This suggests that females harbour the  
168 GSR but that it exhibits very low gene expression in female gonads, perhaps due to  
169 hypermethylation.

170         The mechanism of sex determination in lampreys remains unknown, and may involve  
171 both genetic and environmental factors<sup>18,23–26</sup>. The single elongated gonad remains histologically  
172 undifferentiated for up to several years, and the differentiation process is asynchronous in  
173 females and males (see<sup>18</sup>). Ovarian differentiation occurs in the larval stage, following  
174 synchronized and extensive meiosis and oocyte growth. A few small oocytes may also appear in  
175 future males, but testicular differentiation does not occur until the onset of metamorphosis ~2–3  
176 years later, when resumption of mitosis in the remaining undifferentiated germ cells produces  
177 spermatogonia<sup>23</sup>. It also appears that some larvae may be capable of undergoing sex reversal to  
178 males following ovarian differentiation<sup>24</sup>. Thus, a suite of genes could be turned on to initiate  
179 testicular differentiation. Our data suggest that female sea lamprey gonads harbour the same  
180 GSR as males but, with the exception of some rRNA and ribosomal protein-coding genes,



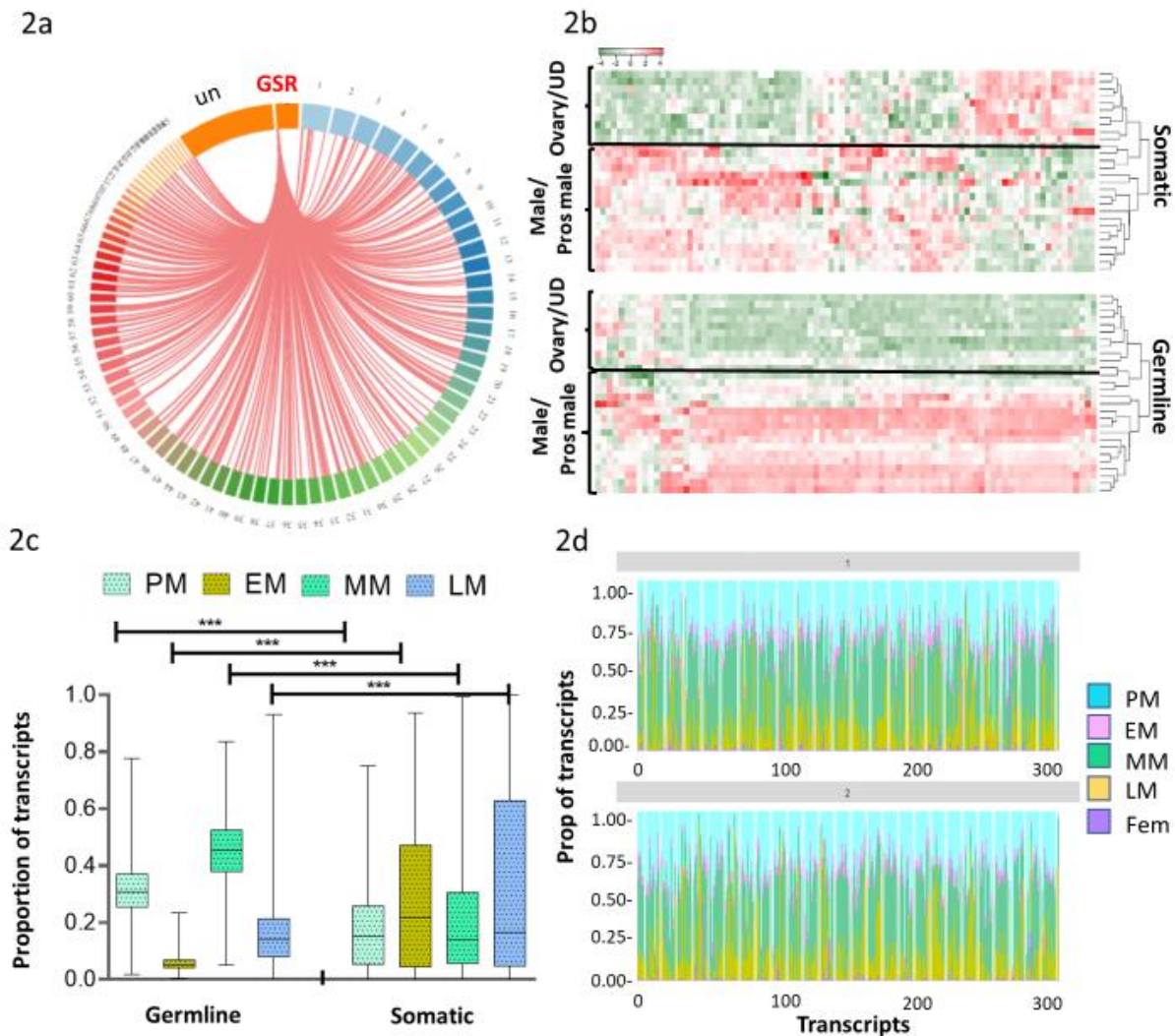
181 females exhibited very low expression of the GSGs (Supplementary Fig. 8a–8b). Male-biased  
182 sex ratios under conditions of high larval density or slow growth have led to suggestions that  
183 primary sex differentiation in lampreys is influenced by environmental factors<sup>25,26</sup>.  
184 Environmental factors that influence the activation or silencing of genes in the GSR could, at  
185 least partially, control sex determination. In this case, low expression of the GSGs would result  
186 in a phenotypically female lamprey, whereas high expression in late larval or early  
187 metamorphosing lamprey would produce a male.

188 **Somatic paralogues of GSGs are expressed differently than germline paralogues.** We  
189 observed that many of the GSGs had duplicated copies: of the 163 GSGs, 92 were found to have  
190 one or more paralogous copies in the GSR (Supplementary Table 4) while 89 have putative  
191 paralogs in the somatic genome, suggesting that some of the GSGs may have been recruited to  
192 the GSR to play specific roles in gametogenesis. The somatic paralogues of the GSGs were  
193 found distributed throughout the entire somatic genome, on every chromosome except  
194 chromosome 49 (Fig. 2a, Supplementary Table 7). To assess whether the somatic paralogues of  
195 the GSR genes exhibit similar sex-biased expression, we selected one paralogous gene per  
196 genome (somatic and germline) and generated a heatmap to compare somatic vs. GSR  
197 expression of the paralogous genes (Fig. 2b). In keeping with the somatic-wide pattern (Fig. 1d),  
198 this demonstrated that the somatic paralogues of the GSR genes do not exhibit the same sex-  
199 biased expression (Fig. 2b).

200

201

202



203 **Fig. 2:** Somatic paralogues of GSGs are expressed differently than germline paralogues. **a)**  
 204 Circos plot indicating the link between genes in the GSR with putative somatic paralogues in the  
 205 sea lamprey genome. Chromosome 81 and enriched scaffolds are indicated as GSR and non-  
 206 enriched scaffolds are indicated as un (unplaced scaffolds in somatic genome). **b)** Heatmap  
 207 showing the relative expression of genes that have paralogues in both the GSR and somatic  
 208 genomes in males and prospective males (pros male), females (ovary), and undifferentiated  
 209 larvae (UD). **c)** Box plot showing the gene expression differences of somatic and GSR  
 210 paralogues of GSGs in prospective males (PM), and early, mid, and late males (EM, MM, and  
 211 LM, respectively). **d)** Comparison of the proportion of transcripts of GSGs in prospective,  
 212 early, mid, and late males (PM, EM, MM, LM) and females (Fem). X-axis represents the number of  
 213 transcripts present in GSR and Y-axis represents the proportion of transcripts in each stage.  
 214

215

216           Given the evidence of male-biased gene expression in the GSR, we next examined  
217 whether the GSGs had uniform expression across male gonadal developmental stages, using the  
218 number of male-biased genes per stage in the somatic genome as reference. In total, we  
219 identified 1270 male-biased genes (of 18,945 total genes), of which 409 (of 638 total genes)  
220 were found in the GSR and 861 (of 18,307 total genes) were found in the somatic genome,  
221 indicating that genes in the GSR have a 36x higher odds of exhibiting male-biased expression  
222 (OR = 36.5068 where  $P < 0.0001$ ). Using the normalized counts of transcripts exhibiting male-  
223 biased expression, we compared the proportion of total transcripts in early, mid, and late  
224 testicular development and in prospective males by examining the interaction between genome  
225 (somatic or GSR) and stage using a repeated measures mixed model design in which gene nested  
226 in genome was a random effect, and stage was a repeated measure (Supplementary Table 8,  
227 Supplementary Fig. 9). This showed that there was a higher proportion of genes expressed in  
228 males in mid-testicular development and in prospective males in the GSR compared to somatic  
229 genomes, and a significantly lower proportion of genes expressed in early and late testicular  
230 development in the GSR relative to the somatic genome (Fig. 2c, Supplementary Tables 8 & 9).

231           To visualize the stage-specific bias in gene expression of GSGs, we plotted the relative  
232 proportion of transcripts expressed in each of the three male gonadal stages (early, mid and late),  
233 as well as in prospective males and the pooled sum of transcripts expressed at any female stage  
234 (Fig. 2d, Supplementary Table 10). This underscores that there is a similar pattern of expression  
235 across all genes in the GSR: high gene expression in prospective and mid gonadal stage males,  
236 but zero to very low expression in females. These findings are similar but distinct from those in  
237 zebra finch: the chromosomes undergoing chromosomal diminution and the genes eliminated are  
238 not sex-biased; however, individual genes have showed expression in both testes and ovaries,

239 with overall greater enrichment for genes involved in ovarian development<sup>16</sup>. On the other hand,  
240 in a sciarid fly (*Sciara coprophila*), the elimination of one or two paternal X chromosomes in all  
241 somatic cells determines the sex of the embryo<sup>10</sup>. Here, we find evidence that the GSGs show  
242 comparatively higher expression in presumptive males when male sea lamprey are putatively  
243 undergoing sex determination and in a later stage of spermatogenesis when male gametes are  
244 generating spermatogonial Type B cells. This supports our hypothesis that gene expression in the  
245 sea lamprey GSR may function to control sex determination and/or differentiation, with high  
246 expression leading to testicular development in males and gene silencing resulting in ovarian  
247 differentiation in females.

248 Our IGV analysis showed that the GSR appears to be present in ovaries (Supplementary  
249 Fig. 8a-8b), suggesting that the genes in the GSR are turned off in females, while they are  
250 expressed in males throughout the sampled stages of spermatogenesis. One possibility is that  
251 differential DNA methylation is involved in sex determination/differentiation in sea lamprey.  
252 DNA methylation is a common process of epigenetic modification with known roles in gene  
253 regulation, embryogenesis and increasingly, sex determination<sup>27</sup> which has, interestingly,  
254 become more important throughout deuterostome evolution<sup>28</sup>. A recent study in zebrafish (*Danio*  
255 *rario*) found that DNA methylation plays important functions in germline development as well  
256 as in sexual plasticity<sup>29</sup>. Given the clear role for the GSR in male spermatogenesis, we wanted to  
257 probe the expression of the GSR during early development bracketing PGR itself. To this end,  
258 we analyzed publicly available RNA-seq data from sea lamprey embryos that span the PGR (1–5  
259 dpf). Of the 638 genes we identified in the GSR, only 186 were expressed during early  
260 embryogenesis. Of these 186 genes, 146 were expressed prior to PGR and 111 post-PGR, but  
261 only 20 had an average gene count >50 post-PGR and 18 pre-PGR (Supplementary Table 11),

262 while the five most abundantly expressed genes code for ribosomal proteins. We then compared  
263 the expression of the 186 GSGs expressed during pre- and post-PGR embryos with our male and  
264 female gonad samples, and find that they exhibit very low expression in females and embryos,  
265 but high expression in male gonads (Supplementary Fig. 10). This further supports the  
266 hypothesis that the role of the GSR in sea lamprey is predominantly to support male gonadal  
267 development.

### 268 **Evolutionary conservation of GSGs and their function in vertebrate spermatogenesis.**

269 Genes in the GSR are expected to be released from the dosage sensitivity constraints of genes in  
270 the somatic genome<sup>4</sup> and may show conservation of gene functions related to gonadal sex  
271 determination and differentiation in other vertebrates. We thus hypothesized that genes in the  
272 GSR 1) do not originate from a single linkage group in the pre-vertebrate ancestor, and do not  
273 map to a single linkage group in the post-2R vertebrate genome, 2) exhibit accelerated evolution  
274 either via high rates of duplication and/or amino acid change and 3) have known roles in sex  
275 determination or spermatogenesis in other vertebrates. To this end, we performed comparative  
276 mapping of genes in the sea lamprey GSR to an earlier chordate (*Branchiostoma belcheri*) and to  
277 nine post-2R taxa. Of the 163 unique gene names identified in the GSR, orthologues with  
278 variable levels of conservation across chordates were identified for 31 genes (Supplementary  
279 Table 12). Some of these genes are found predominantly as a single copy in most taxa, whereas  
280 in sea lamprey, we find a single copy in the GSR but multiple paralogues in the somatic genome  
281 (*rpab4*, *rlp37A*, *mid2bp*, and *hsop3*), multiple paralogues in both the GSR and somatic genomes  
282 (*scyp1*) or a single copy in the GSR and somatic genomes (*fgfr3*), or a single copy in the GSR  
283 but no copy in the somatic genome (e.g. *agrl3*, *cxbl*, *hsop3*, *rpab4*) (Supplementary Fig. 11,  
284 Supplementary Table 4).

285           On the other hand, some of the genes in the GSR are found in multiple paralogues in later  
286 vertebrate genomes, and in multiple copies in the GSR and/or somatic genomes in sea lamprey  
287 (*agrl3*, *cxbl*, *spop1*, *lpar1*, *cadh2*, *lrrn1*, *mlcl1*) (Supplementary Fig. 11). Of the 31 genes  
288 assigned to an orthogroup, 23 were also identified in *Branchiostoma* (Supplementary Table 12  
289 and Supplementary Fig. 11), and 9 are present only in the GSR (not the somatic genome)  
290 suggesting that some of the lamprey GSR genes are not novel. Lastly, a comparative syntenic  
291 analysis between all genes in the lamprey genome for which orthogroups were assigned to the  
292 pre-vertebrate ancestral genome (n = 9,850) or human genome (n = 19,701) found blocks of  
293 conserved synteny for somatic but not GSGs. For example, there is strongly conserved synteny  
294 between the sea lamprey somatic genome and the 17 linkage groups hypothesized to exist in the  
295 pre-2R vertebrate genome (see <sup>29</sup>), but the genes linked to the GSR are dispersed across all but  
296 two of the pre-2R linkage groups (Fig. 3a), and do not show conserved synteny in the human  
297 genome (Supplementary Fig. 12). This suggests that the genes involved in spermatogenesis in  
298 the GSR were independently duplicated into the GSR and were not part of an evolutionarily  
299 conserved paralogon.

300

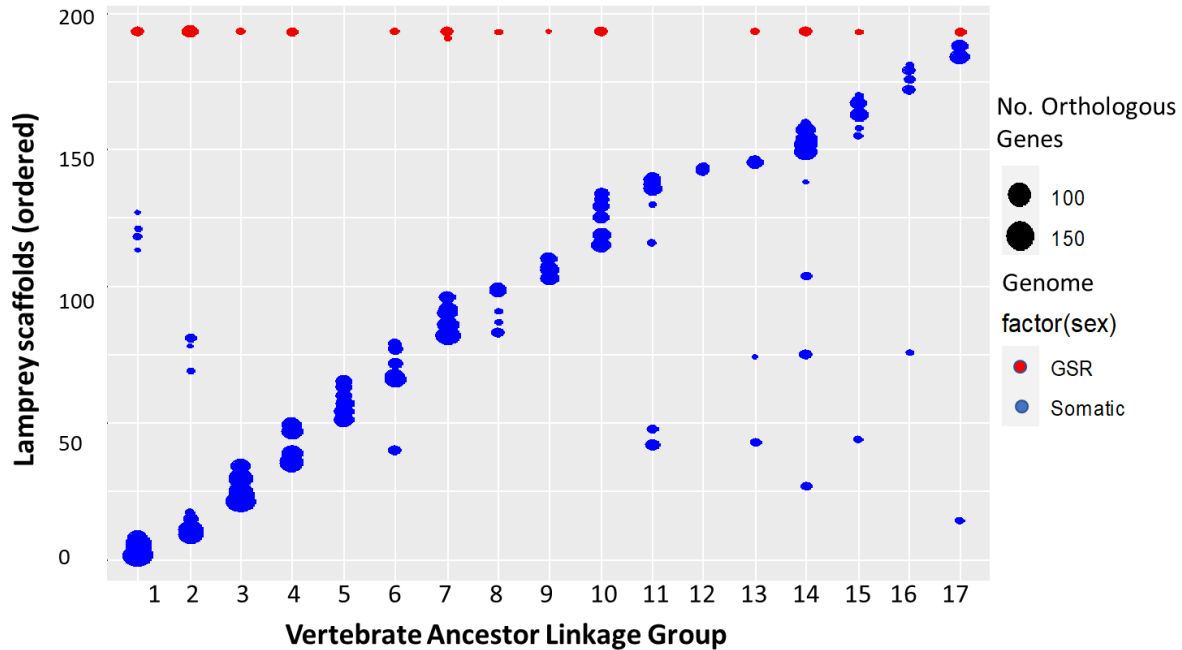
301

302

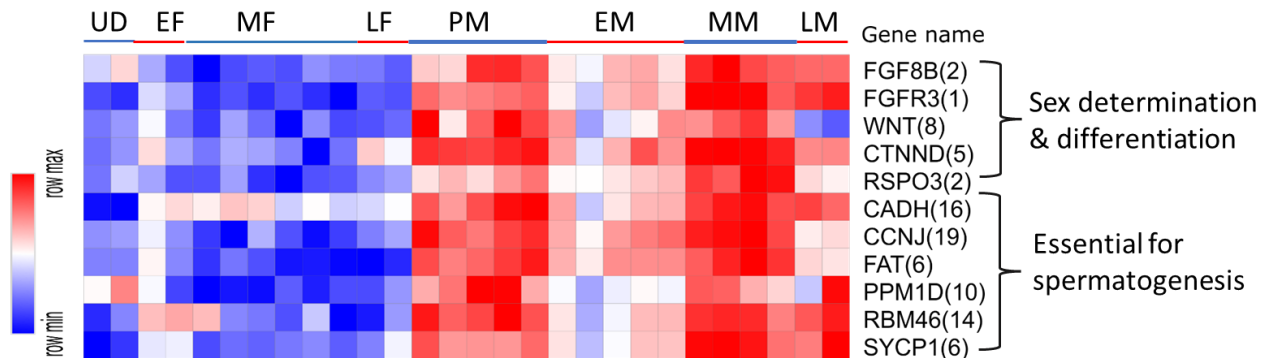
303

304

3a



3b



305 **Fig. 3:** The evolutionary relationship of GSGs with the pre-vertebrate ancestral genome  
 306 reconstructed by<sup>30</sup> and obtained from Genomicus webserver<sup>31</sup> and their functional  
 307 conservensess across vertebrate **a)** Chromosome plot showing comparative mapping in the sea  
 308 lamprey and ancestral genomes. **b)** Heatmap of the median expression of paralogs in gene  
 309 families present in GSR that have known roles in sex determination or spermatogenesis in other  
 310 vertebrates. The numbers in brackets are the number of paralogues of these genes present in  
 311 GSR. UD represents undifferentiated larvae, EF, MF and LF represent early, mid and late  
 312 females respectively and PM, EM, MM and LM represents prospective, early, mid and late males  
 313 respectively.

314

315 We searched the literature for evidence that any of the 163 unique gene names we

316 identified in the GSR are associated with sex determination and/or differentiation in other taxa

317 (Supplementary Table 13). Some of the GSGs have been found to exhibit female-biased

318 expression in later diverging vertebrates, and some are involved in ovarian development,  
319 suggesting that the tissue (gonad) of expression may be conserved, but the function (male vs.  
320 female gonadogenesis) is not. Importantly, however, we find orthologues or paralogues of most  
321 of the core genes involved in sex-determination across vertebrates e.g., fibroblast growth factor 8  
322 (*fgf8*), which is involved in sex determination in mice<sup>32-34</sup>, as well as fibroblast growth factor  
323 receptor 3 (*fgfr3*), which is involved in sex determination in sturgeon (*Acipenser dabryanus*)<sup>35</sup>.  
324 Other genes such as *scyp1*, which is important for early meiotic recombination during  
325 spermatogenesis<sup>36</sup>, R spondin (*rspo1*) and beta catenin 1 (*ctnbl1*) are important antagonists for  
326 Wnt pathway and initiating testicular differentiation<sup>37,38</sup>. Further, several of the gene families  
327 known to be essential for spermatogenesis are highly duplicated. For example, cadherins (*cadh*)  
328 are responsible for maintaining the integrity of testis structure<sup>39</sup>; cyclins (*ccnb*) are essential for  
329 cell progression during distinct phases of the male spermatogenesis pathway<sup>40</sup>; RNA binding  
330 proteins (*rbm*) play diverse and important roles in spermatogenesis including testis-specific  
331 splicing<sup>41</sup> and the absence of *rbm46* (present in 16 copies in the sea lamprey GSR) is associated  
332 with male infertility in mice<sup>42</sup>. Other important genes e.g., *sox9* and *cbx2* which play roles in  
333 stabilizing the male differentiation pathway, are present in the somatic genome of sea lamprey.  
334 We find that all of these genes are highly expressed in the gonads of prospective males and mid-  
335 males when gonadal germ cell specification and spermatogonial development are occurring  
336 respectively (Fig 3b). This suggests that the GSR is likely to play a role in gonadal sex  
337 determination and differentiation as well as spermatogenesis in sea lamprey.

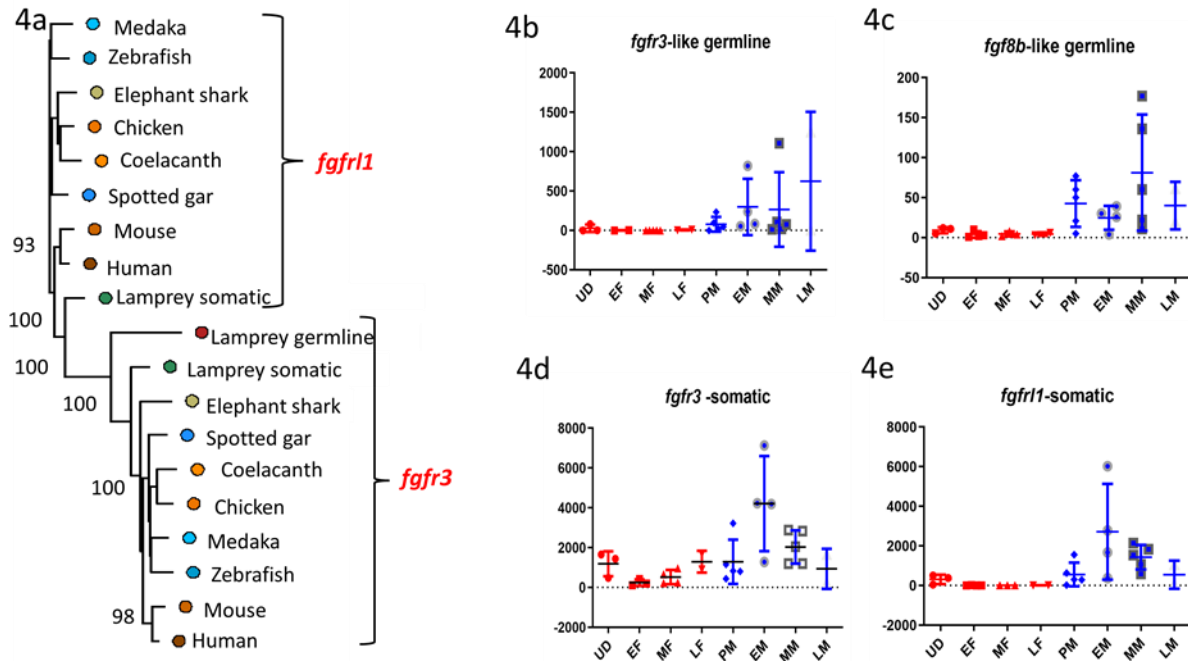
338 **Phylogenetic relationship of GSGs provides evidence of diversified genes involved in sex-**  
339 **determination pathway.** Lampreys diverged from the jawed vertebrate lineages more than 500  
340 million years ago<sup>43,44</sup>, either after the two rounds (2R) of whole genome duplication (WGD) that



341 occurred in early vertebrate evolution<sup>30,45</sup>, or more likely after 1R<sup>20,46,47</sup>. However, a recent study  
342 suggested that, after the 1R tetraploidization, lampreys underwent an additional  
343 hexaploidization<sup>48</sup>. Since lampreys have an unusual vertebrate ploidy state, it proved impossible  
344 to perform a reliable test of positive selection at the amino acid level (which requires essentially  
345 gapless alignments) for the germline genes in agnathans (lampreys and hagfishes) relative to  
346 other vertebrates. Thus, we selected a few genes which have important roles in gametogenesis in  
347 other species for phylogenetic analyses (see Supplementary Table 13).

348         Gene trees were reconstructed using the output from OrthoFinder and the orthologues of  
349 sea lamprey GSGs in 10 other chordates identified (see Supplementary Fig. 11) and combined  
350 with our data on gene annotations and genomic location (GSR vs. somatic) in sea lamprey. This  
351 revealed that the *cadh* gene family is highly duplicated in both the germline and somatic  
352 genomes of sea lamprey (16 vs. 15 duplicates, respectively) and hagfish (Supplementary Table  
353 12). In particular, *cadh2* has undergone a divergent expansion in the GSR in sea lamprey  
354 (Supplementary Fig. 13a); agnathans have witnessed an expansion of a somatic cluster of genes  
355 related to vertebrate *cadh1/cadh3/cadh13* as well as an expansion in both the somatic and  
356 germline genomes of a novel *cadh* paralogue (bottom of Supplementary Fig. 13a). Phylogenetic  
357 trees for *hykk* (Supplementary Fig. 13b), *scyp1* (Supplementary Fig. 13c), and *adgr1*  
358 (Supplementary Fig. 13d) depict similar patterns of one or more highly duplicated germline  
359 lineages that are sometimes interspersed with closely related somatic paralogues (*hykk* and  
360 *scyp1*), but overall, they show clades of highly diversified germline lineages marked by long

361 internal branch lengths, indicating that the GSGs exhibit independent evolution for variable  
 362 periods of time and may be subject to positive selection.



363 **Fig. 4:** Phylogenetic tree of putative ligand and receptor of fibroblast growth factor (fgf) and  
 364 their genome-, sex- and stage-specific expression a) phylogenetic tree, b) *fgfr3*-like gene in the  
 365 GSR genome, c) *fgf8b* like gene in GSR, d) *fgfr3* in somatic genome, and e) *fgfr1* in the somatic  
 366 genome. In 4b–4e, the Y-axis is gene counts and the X-axis is stage, where UD is  
 367 undifferentiated larvae; EF, MF, and LF are early, mid, and late females; and PM, EM, MM, and  
 368 LM are prospective, early, mid, and late males, respectively.

369  
 370

371 We identified a novel *fgfr3*-like gene in the germline genome, and confirmed expression  
 372 of a possible ligand for it, *fgf8b*, also located in the germline genome. The *fgfr3*-like gene was  
 373 not identified by OrthoFinder as an orthologue of the somatic copy of *fgfr3*; thus, we  
 374 downloaded the canonical coding sequences for *fgfr3*, and a related gene also present in the  
 375 somatic genome, *fgfr11*, from eight post-2R taxa and reconstructed a ML tree with bootstrap  
 376 support (Fig. 4a). This revealed that the germline sequence of the *fgfr*-like coding sequence is  
 377 more closely related to *fgfr3* in the sea lamprey somatic genome and to the *fgfr3* in higher

378 vertebrates (bootstrap support 100%), while the somatic copy of *fgfr11* groups with the *fgfr11*  
379 sequences from later vertebrates and there is no paralogue in the GSR (100% bootstrap support).  
380 Examination of the expression of these three genes as well as the possible receptor for the  
381 germline gene, *fgf8b*, indicates that the germline copy of *fgfr3* and *fgf8b* have very low  
382 expression in female gonads, and somewhat higher expression in male gonads: notably, *fgf8b* is  
383 most highly expressed in prospective male and mid-stage male gonads (Fig. 3b, 4b–4e). Given  
384 their role in sex determination in other vertebrates, sea lamprey germline genes *fgf8b* and *fgfr3*  
385 warrant further investigation as possible loci involved in sex determination.

### 386 **Conclusion**

387 The study of PGR events and their effects on gonadal development and sex determination  
388 represent a burgeoning field in evolutionary biology. Our result suggests that the genes present in  
389 the GSR in sea lamprey are involved in the crucial processes of sex differentiation and testicular  
390 development, and could be involved in sex determination. We find GSGs are most highly  
391 expressed in prospective males and in males undergoing spermatogonial differentiation, but they  
392 have low overall expression in females. Given evidence that sex is partially determined by  
393 environmental factors in sea lamprey, the possible role of methylation in the GSR during early  
394 stages of gonad development in larval sea lamprey warrants further attention. We find low levels  
395 of syntenic or sequence conservation of genes in the GSR across chordates, but importantly,  
396 many of the genes identified in the GSR are known to play roles in gonad differentiation or sex  
397 determination in other vertebrates. Assuming females harbour the same GSR as males, our data  
398 suggests that the factors controlling epigenetic modification of the GSR are pivotal for sex  
399 determination and differentiation. Further work is needed to assess the presence and chromatin

400 accessibility of the GSR in females and to identify the function of the GSGs in sea lamprey sex  
401 determination and differentiation.

## 402 **Methods**

403 **Sample preparation and RNA extraction.** Sea lamprey from different life-history stages were  
404 collected by collaborators using these samples for other projects. An Abbreviated Protocol for  
405 Minimal Animal Involvement form completed at the University of Manitoba determined that an  
406 Animal Use Protocol (AUP) was not required because live sea lamprey were not handled by us  
407 for the purposes of this project, and no animals were sacrificed or manipulated solely to provide  
408 us with tissue.

409 Larval sea lamprey were collected by backpack, pulsed DC electrofishing in tributaries of  
410 the Richibucto River, New Brunswick, Canada, or in tributaries of Lake Huron and Lake  
411 Michigan in the Great Lakes basin (Supplementary Table 1). Larvae were transported or shipped  
412 live to Wilfrid Laurier University, Waterloo, ON, sorted according to size, and transferred to 110  
413 L holding tanks supplied with aerated well water at a flow rate of 1.0–2.0 L/min. The larvae were  
414 monitored for external signs of metamorphosis (e.g., changes in eye and oral disc morphology)  
415 and then euthanized at the desired stages. The brain and gills, required for other projects, were  
416 dissected and placed in RNAlater. With the remaining carcass (posterior to the last branchial  
417 pore), RNAlater was injected into the gut to perfuse the intestine, liver, gallbladder, kidneys, and  
418 gonad. The carcass with organs was then placed in a 10 mL Falcon tube and filled with RNAlater  
419 to saturate the tissues thoroughly. Dissections were completed as rapidly as possible to reduce  
420 any potential RNA degradation. Samples were kept at 4 °C for 24 h, stored at –80 °C, and then  
421 shipped to the University of Manitoba on dry ice, and stored at –80 °C upon arrival. The gonads  
422 were subsequently dissected out and placed in a 1.5 mL centrifuge tube with 1 mL RNAlater and

423 kept at  $-20^{\circ}\text{C}$ . Sex was identified during dissection based on physical inspection with the naked  
424 eye (i.e., the ovary is larger and has a different texture than the testis or undifferentiated gonad),  
425 and gonadal stage was identified by a combination of visual inspection and inferences based on  
426 larval size and stage of metamorphosis<sup>18</sup> (Supplementary Table 1).

427 Adult sea lamprey were captured in traps near the mouth of the Black Mallard River or  
428 Ocqueoc River, MI, during their upstream (spawning) migration (Supplementary Table 1).  
429 Lamprey were euthanized, length and weight measurements were taken, and  $\sim 35$  mg gonad was  
430 flash frozen in a 2.0 mL centrifuge tube and kept on dry ice (April 2018) or placed in a 1.5 mL  
431 centrifuge tube with 1 mL RNAlater and kept at  $-20^{\circ}\text{C}$  (June 2018). Samples were shipped to  
432 the University of Manitoba on dry ice, and stored at  $-80^{\circ}\text{C}$ .

433 Total RNA was isolated from  $\sim 30$  mg of gonadal tissue from each individual using the  
434 RNeasy Mini kit (Qiagen, USA) according to the manufacturer's protocol. The extracted RNA  
435 was treated with RNase-free DNase set (Qiagen, USA) to remove residual genomic DNA. RNA  
436 quantity and quality was assessed using a NanoVue Plus spectrophotometer. The RNA samples  
437 were preserved at  $-80^{\circ}\text{C}$ .

438 To obtain a comprehensive representation of gene expression, RNA from individuals at  
439 the same stage of development and same sex was pooled. Early males ( $n = 4$ ) were those  
440 identified by external morphological characteristics to be in the early to mid stages of  
441 metamorphosis and thus presumed to be in the early stages of spermatogonial differentiation, that  
442 is, in the process of producing Type A spermatogonia (Supplementary Fig. 1, Supplementary  
443 Table 1)<sup>18</sup>. Mid males (metamorphosing stage 7 and immediately post-metamorphosis;  $n = 6$ )  
444 were presumed to be undergoing spermatogonial proliferation and differentiation and producing  
445 Type A and Type B spermatogonia, while late males were sexually mature ( $n = 2$ ). In early

446 females (n = 2), ovarian differentiation had been initiated and/or completed (i.e., with a number  
447 of small growing oocytes in the gonad), mid-stage females (n = 6) had completed oocyte  
448 differentiation and were arrested in meiotic prophase with larger growing oocytes, and late  
449 females were sexually mature (n = 2). In addition to samples that were definitively male and  
450 female, larvae that had histologically undifferentiated gonads and were below the size at which  
451 ovarian differentiation occurs (n = 2) and presumptive male larvae with histologically  
452 undifferentiated gonads but beyond the size at which ovarian differentiation is complete (n = 4)  
453 were included.

454 **Library preparation, Illumina sequencing, and data filtering.** High-quality RNA from 28  
455 gonad samples was sent to Genome Quebec, McGill University, Montreal, to construct a cDNA  
456 library and perform RNA sequencing. Messenger RNA (mRNA) was isolated using poly-A  
457 isolation and non-normalized libraries prepared using the Illumina TruSeq DNA Kit and  
458 Epicentre Script Seq Kit. Sequencing was performed in both forward and reversed directions and  
459 100 base pair (bp) reads were generated on an Illumina Hi-Seq 4000 PE100. The resulting RNA-  
460 Seq paired-end (PE) reads were checked for quality control using FASTQC (v0.11.8)<sup>49</sup>, and low-  
461 quality sequences and adapters were trimmed with Trimmomatic (v0.36)<sup>50</sup>, using  
462 ILLUMINACLIP:TruSeq3-PE-2.fa:2:15:10 LEADING:5 TRAILING:5 SLIDINGWINDOW:4:5  
463 MINLEN:50 and a quality score threshold of Phred-33.

464

465 **Combining reference and *de novo* assemblies:**

466 ***Generating comprehensive gonadal superTranscriptome:*** The software pipeline Necklace<sup>19</sup>,  
467 was used to generate a merged superTranscriptome derived from three sources: 1) a genome-  
468 guided alignment using the sea lamprey reference genome, 2) a *de novo* assembly using Trinity,

469 and 3) a reference-based proteome from other chordate species. For the genome-guided  
470 assembly, the 28 gonadal transcriptomes were mapped to the Vertebrate Genome Project (VGP)  
471 sea lamprey reference germline genome  
472 ([https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/010/993/605/GCF\\_010993605.1\\_kPetMar1.pri/GCF\\_010993605.1\\_kPetMar1.pri\\_genomic.fna.gz](https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/010/993/605/GCF_010993605.1_kPetMar1.pri/GCF_010993605.1_kPetMar1.pri_genomic.fna.gz)) and associated gene annotation file  
473 ([https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/010/993/605/GCF\\_010993605.1\\_kPetMar1.pri/GCF\\_010993605.1\\_kPetMar1.pri\\_genomic.gff.gz](https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/010/993/605/GCF_010993605.1_kPetMar1.pri/GCF_010993605.1_kPetMar1.pri_genomic.gff.gz)) available at NCBI. Reads were aligned to the  
474 sea lamprey genome using HISAT2, and StringTie<sup>51</sup> was used to assemble transcripts, some of  
475 which map to known genes and some of which are novel (MSTRG IDs). For the third tier of the  
476 Necklace pipeline, reference proteomes from a non-teleost fish, spotted gar (*Lepisosteus*  
477 *oculatus*)  
478 ([https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/242/695/GCF\\_000242695.1\\_LepOcu1/GCF000242695.1\\_LepOcu1\\_protein.faa.gz](https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/242/695/GCF_000242695.1_LepOcu1/GCF000242695.1_LepOcu1_protein.faa.gz)), and a cartilaginous fish, elephant shark (*Callorhinchus*  
479 *milii*) ([https://www.ncbi.nlm.nih.gov/genome/689?genome\\_assembly\\_id=49056](https://www.ncbi.nlm.nih.gov/genome/689?genome_assembly_id=49056)) were used.

483 In the second step, a *de novo* assembly of reads was generated with Trinity<sup>52</sup> for all 28  
484 samples. The assembled transcripts from genome-guided and *de novo* assembly were sorted into  
485 three groups: annotated transcripts that align to the reference genome (known genes), transcripts  
486 that align to the reference genome but are not found in the reference annotation (reference-based  
487 novel genes), and unmapped novel transcripts – those that align to the spotted gar/elephant shark  
488 proteome (*de novo*-specific genes). These three groups were merged into a single  
489 superTranscriptome and used for the second stage of the analysis: gene counting and differential  
490 expression analyses. The Necklace pipeline allows for the identification of novel transcripts yet  
491 generates a compact and comprehensive superTranscriptome, while preventing the introduction

492 of false chimeras generated during *de novo* assembly. The step-by-step workflow of Necklace  
493 pipeline is illustrated in Supplementary Fig. 14.

494 In total, we identified 42,479 genes in the sea lamprey germline genome, of which 20,630  
495 overlapped with those annotated by NCBI (representing ~94% of the total number of genes in  
496 the VGP annotation), 21,808 were identified *de novo* through StringTie, and 40 Trinity *de novo*  
497 assembled transcripts matched sequences in the spotted gar/elephant reference proteome by  
498 homology. However, since the genomic location of these 40 homology-based sequences could  
499 not be ascertained, they were discarded from further analyses. Of the remaining 42,439  
500 sequences, tRNA, rRNA, and lncRNAs (long non-coding RNAs), were removed, retaining  
501 18,945 protein-coding transcripts (16,328 from the VGP annotation and 2,617 novel transcripts,  
502 which is ~14% of the total gene list) (Supplementary Fig. 15). Those 18,945 genes pertain to  
503 12,583 unique gene names, which would be a lower limit on the actual number of genes  
504 identified, since paralogous genes may be assigned the same gene name.

505 **Gene-counts:** Reads from each of the 28 gonadal transcriptomes were subsequently aligned to  
506 the merged superTranscriptome, and gene counts extracted and filtered. These gene-counts are  
507 used for further downstream analysis, i.e., in differential gene expression analysis, identifying  
508 sex-biased and sex-specific transcripts and genes.

#### 509 **Functional annotation and identifying orthogroups:**

510 **Functional annotation:** All of the 18,945 putatively protein-coding genes generated from the  
511 Necklace pipeline were annotated using Trinotate pipeline (v3.2.0)<sup>53</sup> following the method  
512 described at (<http://trinotate.github.io/>). Initially, Transdecoder (v5.5.0) was used to obtain the  
513 expected start and stop sites of protein translation from the assembled superTranscriptome. Then



514 each transcript and protein sequence were searched against the SwissProt database using blastx  
515 and blastp. The HMMER algorithm was used to search PFAM (ran in hmmer (v3.2.1)) for  
516 protein domain identification, signalp (v4.1f) and tmHMM (v2.0c) were used to predict the  
517 signal peptide and transmembrane regions, respectively, and Rnammer (v1.2) was used to  
518 identify rRNA transcripts which were automatically removed in a later stage of the pipeline. In  
519 the final stage, the results from blast searches were combined with the other functional  
520 annotation data and loaded into the Trinotate.SQLite database: an e-value of 1e-5 was used as the  
521 threshold to generate the functional annotation report.

522 **Orthogroup identification:** The homology between the genes in our annotated sea lamprey  
523 gonadal superTranscriptome were compared to genes in 11 chordate species chosen to represent  
524 important time points in chordate evolution using the OrthoFinder pipeline<sup>54</sup>. Protein sequences  
525 were obtained from human ([ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-103/fasta/homo_sapiens/pep/Homo_sapiens.GRCh38.pep.all.fa.gz)  
526 [103/fasta/homo\\_sapiens/pep/Homo\\_sapiens.GRCh38.pep.all.fa.gz](103/fasta/homo_sapiens/pep/Homo_sapiens.GRCh38.pep.all.fa.gz)), mouse (*Mus musculus*)  
527 ([ftp://ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-102/fasta/mus_musculus/pep/Mus_musculus.GRCm38.pep.all.fa.gz)  
528 [102/fasta/mus\\_musculus/pep/Mus\\_musculus.GRCm38.pep.all.fa.gz](102/fasta/mus_musculus/pep/Mus_musculus.GRCm38.pep.all.fa.gz)), zebrafish  
529 ([ftp.ensembl.org/pub/release-103/fasta/danio\\_rerio/pep/Danio\\_rerio.GRCz11.pep.all.fa.gz](ftp://ftp.ensembl.org/pub/release-103/fasta/danio_rerio/pep/Danio_rerio.GRCz11.pep.all.fa.gz)),  
530 chicken ([ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-103/fasta/gallus_gallus/pep/Gallus_gallus.GRCg6a.pep.all.fa.gz)  
531 [103/fasta/gallus\\_gallus/pep/Gallus\\_gallus.GRCg6a.pep.all.fa.gz](103/fasta/gallus_gallus/pep/Gallus_gallus.GRCg6a.pep.all.fa.gz)), medaka (*Oryzias sinensis*)  
532 ([ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-103/fasta/oryzias_sinensis/pep/Oryzias_sinensis.ASM858656v1.pep.all.fa.gz)  
533 [103/fasta/oryzias\\_sinensis/pep/Oryzias\\_sinensis.ASM858656v1.pep.all.fa.gz](103/fasta/oryzias_sinensis/pep/Oryzias_sinensis.ASM858656v1.pep.all.fa.gz)), spotted gar  
534 ([ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-103/fasta/lepisosteus_oculatus/pep/Lepisosteus_oculatus.LepOcu1.pep.all.fa.gz)  
535 [103/fasta/lepisosteus\\_oculatus/pep/Lepisosteus\\_oculatus.LepOcu1.pep.all.fa.gz](103/fasta/lepisosteus_oculatus/pep/Lepisosteus_oculatus.LepOcu1.pep.all.fa.gz)), elephant shark  
536 ([25](ftp://ftp.ensembl.org/pub/release-</a></p></div><div data-bbox=)

537 [103/fasta/callorhinchus\\_milii/pep/Callorhinchus\\_milii.Callorhinchus\\_milii-6.1.3.pep.all.faa.gz](#)),  
538 coelacanths ([ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-)  
539 [103/fasta/latimeria\\_chalumnae/pep/Latimeria\\_chalumnae.LatCha1.pep.all.faa.gz/](#)), hagfish  
540 (*Eptatretus burgeri*) ([ftp.ensembl.org/pub/release-](ftp://ftp.ensembl.org/pub/release-)  
541 [103/fasta/eptatretus\\_burgeri/pep/Eptatretus\\_burgeri.Eburgeri\\_3.2.pep.all.faa.gz](#)), amphioxus  
542 (*Branchiostoma belcheri*)  
543 ([https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/001/625/305/GCF\\_001625305.1\\_Haploidv18h27/](https://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/001/625/305/GCF_001625305.1_Haploidv18h27/)  
544 [GCF\\_001625305.1\\_Haploidv18h27\\_protein.faa.gz](#)). OrthoFinder uses the complete list of  
545 known protein sequences from all included taxa to find putative orthologues, and then creates  
546 orthogroups with related sets of orthologues. OrthoFinder exports multiple sequence alignments  
547 and rooted gene trees for all orthogroups, which can be used to infer gene duplication events.  
548 Overall, in this study, 93.1% of the genes in the 12 chordate species were assigned to one of  
549 27,364 orthogroups, and 5606 orthogroups contained representatives of all 12 species.

#### 550 **Prediction of germline-specific region (GSR) and genes by enrichment analysis:**

551 **Identifying GSGs in the GSR:** Although the GSR in sea lamprey has been identified for a  
552 previous germline assembly (gPmar100)<sup>20</sup>, when the new VGP germline genome was deposited  
553 on NCBI, the corresponding positions were not available. Following the protocol from Smith et  
554 al. 2018, germline enrichment was calculated using the DifCover program by calculating  
555 differences in read depth between a single germline sample (sperm) and a single somatic sample  
556 (blood) from the same male. We downloaded and mapped the same sperm (SRR5535435) and  
557 blood (SRR5535434) samples they had used from their previous analysis to identify the GSR  
558 coordinates in the newly-deposited VGP genome in order to facilitate our downstream  
559 transcriptomic analyses. The DNACopy output file was generated by following step by step

560 workflow with default settings described in the DifCover pipeline  
561 (<https://github.com/timnat/DifCover>) (See Supplementary Fig. 14)<sup>21</sup>. This DNACopy output file  
562 was then used to identify GSR from the new chromosome level assembly, VGP germline  
563 genome and the DNACopy output file. Later, the GSGs were identified by extracting all genes  
564 that fell within regions having an enrichment score >2 using bedtools (v2.29.0) with the aid of  
565 the genome-based annotation file (generated in Necklace pipeline).

566         Initially, we identified 1845 GSGs by extracting the DNACopy output file from VGP  
567 genome; however, only 783 protein-coding GSGs were retained with gene counts after the initial  
568 filtration steps discussed in previous section. In the next step, we sorted genes based on their  
569 location: if two genes with the same name had overlapping start and end points, the canonical  
570 transcript was retained, which reduced the number of genes to 672. In the final step, we extracted  
571 the protein sequences associated with each of these genes from the transdecoder pep file and  
572 removed ambiguous sequences. The final list consisting of 638 GSGs was merged with the  
573 Trinotate annotation report to assign putative gene names for the novel genes, and with the  
574 reference annotation for genes identified by VGP. In total, 163 unique gene names were assigned  
575 to the 638 GSGs, 70 of them in a single copy, and the remaining 93 in 2–77 copies  
576 (Supplementary Fig. 2).

577         Given our finding that genes in the GSR are highly expressed during gonad development  
578 (see Results and discussion), we wanted to assess whether all or a subset of the genes in the GSR  
579 are also expressed during early embryonic development. To this end, we downloaded paired-end  
580 RNA-seq read data for embryos sampled at 1 dpf (SRR3002837), 2 dpf (SRR3002840), 2.5 dpf  
581 (SRR3002843), 3 dpf (SRR3002846), 4 dpf (SRR3002849) and 5 dpf (SRR3002852) from the  
582 SRA database

583 ([https://www.ncbi.nlm.nih.gov/sra?linkname=bioproject\\_sra\\_all&from\\_uid=306044](https://www.ncbi.nlm.nih.gov/sra?linkname=bioproject_sra_all&from_uid=306044)). Reads  
584 were aligned to the VGP genome using HISAT2 (v 2.2.1)<sup>51</sup>, and assembled into transcripts, and  
585 gene and transcripts counts were obtained per sample using Stringtie (v 2.0)<sup>51</sup>. In the next step,  
586 we extracted the embryo-expressed GSGs using the merged annotation file generated in the  
587 previous step and the DNACopy output file generated from DifCover analysis (Supplementary  
588 Fig. 13). We considered only genes available in the reference annotation that mapped to our  
589 GSR, which resulted in 184 GSGs (after filtering ncRNA (non-coding RNAs), rRNA (ribosomal  
590 RNAs) and pseudogenes from the reference annotation) expressed in early embryonic  
591 development of which 149 genes overlapped with those expressed in our gonad samples. The  
592 gene count file was used to compare gene expression in the GSR pre-PGR (1dpf, 2 dpf and 2.5  
593 dpf) and post-PGR (3 dpf, 4 dpf and 5 dpf) and between gonads and embryos (both pre- PGR  
594 and post-PGR).

595 ***Identifying somatic copies of GSGs:*** To identify putative paralogues of GSGs in the somatic  
596 genome, the list of all genes was sorted based on the unique gene names obtained from either the  
597 Trinotate annotation report or from the reference annotation. Of the 163 unique genes identified  
598 in the GSR, 89 were found to have either a single or multiple putative paralogues in the somatic  
599 genome, of which 31 were matched to a unique OrthoGroup by OrthoFinder (Supplementary  
600 Fig. 16).

#### 601 **Identifying male-biased expression:**

602 To identify sex-biased genes, only the gonadal transcriptomes from definitive males were  
603 analyzed; the undifferentiated larvae and prospective males were removed from this comparison.  
604 Since traditional differential gene expression analyses using a threshold log-fold change between  
605 conditions would likely identify genes with low or no expression in one sex but also low

606 expression in the other sex, we aimed to target genes that were very highly expressed in one sex  
607 but had much lower or no expression in the other. The gene count data was obtained from the  
608 raw gene count data across samples with reference to the necklace superTranscriptome, and these  
609 were later filtered using normalized gene counts and log fold change (logFC) using DESeq2<sup>55</sup>  
610 and EdgeR<sup>56</sup>. That gave us an estimate of 6088 genes with higher logFC in males than females of  
611 which top 20% genes were considered to be male-biased as long as the total gene count is equal  
612 or more than 1000 (Supplementary Fig. 17).

### 613 **Comparison of GSG expression across sex and stage and functional enrichment analysis:**

614 *Comparison of GSG expression across sex and stage:* The genome-wide raw gene counts were  
615 converted to normalized counts using DESEQ2<sup>55</sup>, and the log<sub>2</sub> expression of all genes in the  
616 GSR compared in a sex- and stage-specific manner. To assess global differences in sex bias in  
617 gene expression, we compared the density of the relative log<sub>2</sub>(male:female normalised  
618 expression) of all genes in the somatic genome vs GSR. To compare the difference in expression  
619 between GSGs against their somatic paralogues across both sexes and stage, we calculated the  
620 mean normalized log<sub>2</sub> gene count and visualised the data with a heat map. To assess whether  
621 GSGs exhibit stage-specific sex-biased expression, we extracted the list of all genes exhibiting  
622 sex-biased expression (as defined above) in both the somatic genome and GSR, and assessed  
623 whether the proportional expression of genes differed between genomes and stage using a  
624 repeated measures mixed model in which proportion gene expression was the response and the  
625 model was gene(genome) + stage + genome\*stage, with gene as a random effect, and stage as  
626 the repeated measure.

627 *Functional enrichment:* The list of genes identified in the GSRs was submitted for pathway  
628 analysis using the human protein-coding genes as background using PANTHER (v14)

629 (<http://pantherdb.org>)<sup>57</sup>. The PANTHER GO-slim molecular process terms associated with each  
630 gene were used for an over-representation test<sup>57</sup> in which the Fisher exact test was performed to  
631 assess the significance of terms at an FDR of 0.05. Additionally, we used the gene ontology  
632 (GO) terms associated with the GSGs identified by Trinotate<sup>53</sup>, and visualized them in REVIGO  
633 (<http://revigo.irb.hr/>)<sup>58</sup> in a scatter plot that shows cluster representatives in a two-dimensional  
634 space derived by GO terms with semantic similarity measure and clustering set at 0.9 overall.  
635 Terms were plotted with size proportional to fold-enrichment above expected and color  
636 according to the log<sub>10</sub> of the FDR p-value (Fig. 1b; see Results and discussion).

### 637 **Comparative mapping and phylogenetic analysis:**

638 *Comparative mapping:* Lampreys, being an intermediate lineage between 1R and 2R WGD, are  
639 important model organisms for the study of evolution of genes as well as evolution of  
640 physiological process<sup>20,30,45–48</sup>. We compared the evolutionary origin and relationship of genes in  
641 the GSR to the pre-2R vertebrate genome and in later diverging taxa. To this end, we identified  
642 orthologues of the genes from all 85 assembled chromosomes as well as to scaffolds that we  
643 identified as enriched (GSR) or non-enriched (somatic) for germline DNA (see results) to those  
644 in human, chicken, and spotted gar and a reconstructed pre-2R vertebrate genome. Orthologous  
645 genes were identified using the output of OrthoFinder<sup>54</sup>, assigned to their chromosomal location  
646 using BioMart (ENSEMBL) and the number of co-orthologues per linkage group/chromosome  
647 calculated pairwise. To remove marginally supported synteny, we used the observed number of  
648 lamprey orthologues identified by chromosome for each species comparison as the maximum  
649 expected value and retained all syntenic chromosomal pairs; if there were >10 genes shared  
650 between species for lamprey chromosomes 1–69, or > 5 for chromosomes 70–85 (this criterion  
651 was set based on chromosome size), and retained all orthologous gene matches for the GSR. The

652 reconstructed pre-2R vertebrate genome was downloaded from the ftp site of the Genomicus  
653 webserver<sup>31</sup> (<ftp://ftp.biologie.ens.fr/pub/dyogen/genomicus/69.10> details of the reconstruction  
654 are described in<sup>30</sup>).

655 **Phylogenetic analysis:** Given that genes in the GSRs may have unique evolutionary histories,  
656 the phylogenetic relationships of a subset of the genes in the GSRs and their somatic orthologues  
657 were reconstructed along with orthologous/paralogous genes identified from the 11 taxa included  
658 in the OrthoFinder output. Phylogenetic trees were obtained from OrthoFinder which uses  
659 RaxML reconstruction<sup>54</sup>. Trees were not available for all GSGs, including the sea lamprey  
660 putative orthologue of *fgfr3*, which has been shown to be important for sex determination and  
661 differentiation in other taxa<sup>32,35,59,60</sup>. Thus, we obtained sequences for *fgfr3* for the same 11  
662 species employed in the OrthoFinder analyses, and then performed an alignment in Mafft  
663 (<https://mafft.cbrc.jp/alignment/server/>) followed by ML reconstruction with RAXML. We  
664 hypothesized that genes in the GSR are under relaxed evolutionary constraint and relaxed dosage  
665 sensitivity and thus may exhibit accelerated rates of sequence evolution. However, we were  
666 unable to employ tests of dN/dS due to the difficulty of obtaining sufficiently un-gapped  
667 alignments of the coding sequence of the sea lamprey genes relative to those from other jawed  
668 vertebrates. Nevertheless, phylogenetic trees were generated to understand the relationship of  
669 paralogous copies of the GSGs in the GSR to those in the somatic genome, as well as the  
670 relationship of the protein coding sequences in sea lamprey to those in other chordate taxa.

671 **Data availability:** The RNA-sequencing reads used for this study have been deposited in the  
672 NCBI repository under the BioProject accession number PRJNA749754 and will be available  
673 upon publication of manuscript.

674 **Acknowledgements:** We thank Dr. John B. Hume (Michigan State University), Dr. Nicholas  
675 Johnson (U.S. Geological Survey), Dr. Michael Wilkie (Wilfrid Laurier University), and Joshua  
676 Sutherby (University of Manitoba) for providing us with the samples used in this study. We also  
677 thank Arfa Khan (University of Manitoba) for her guidance during sample processing and RNA  
678 extraction.

679 **Funding:** This research was funded by the Natural Sciences and Engineering Research Council  
680 of Canada (NSERC) Discovery Grants program (MFD, SVG), the University of Manitoba  
681 Graduate Enhancement of Tri-Council Stipends program (MFD), and the Great Lakes Fishery  
682 Commission Sea Lamprey Research Program (MFD).

#### 683 **References:**

- 684 1. Kloc, M. & Zagrodzinska, B. Chromatin elimination - An oddity or a common mechanism  
685 in differentiation and development? *Differentiation* **68**, 84–91 (2001).
- 686 2. Zufall, R. A., Robinson, T. & Katz, L. A. Evolution of developmentally regulated genome  
687 rearrangements in eukaryotes. *J. Exp. Zool. B Mol. Dev. Evol. J EXP ZOOLOG PART B* **304**  
688 448–455 (2005).
- 689 3. Wang, J. & Davis, R. E. Programmed DNA elimination in multicellular organisms. *Curr.*  
690 *Opin. Genet.* **27**, 26–34 (2014).
- 691 4. Bryant, S. A., Herdy, J. R., Amemiya, C. T. & Smith, J. J. Characterization of  
692 somatically-eliminated genes during development of the sea lamprey (*Petromyzon*  
693 *marinus*). *Mol. Biol. Evol.* **33**, 2337–2344 (2016).
- 694 5. Smith, J. J., Antonacci, F., Eichler, E. E. & Amemiya, C. T. Programmed loss of millions  
695 of base pairs from a vertebrate genome. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 11212–11217  
696 (2009).
- 697 6. Smith, J. J., Baker, C., Eichler, E. E. & Amemiya, C. T. Genetic consequences of  
698 programmed genome rearrangement. *Curr. Biol.* **22**, 1524–1529 (2012).
- 699 7. Timoshevskiy, V. A., Lampman, R. T., Hess, J. E., Porter, L. L. & Smith, J. J. Deep  
700 ancestry of programmed genome rearrangement in lampreys. *Dev. Biol.* **429**, 31–34  
701 (2017).
- 702 8. Smith, J. J., Antonacci, F., Eichler, E. E. & Amemiya, C. T. Programmed loss of millions  
703 of base pairs from a vertebrate genome. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 11212–11217  
704 (2009).



- 705 9. Timoshevskiy, V. A., Herdy, J. R., Keinath, M. C. & Smith, J. J. Cellular and molecular  
706 features of developmentally programmed genome rearrangement in a vertebrate (Sea  
707 lamprey: *Petromyzon marinus*). *PLoS Genet.* **12**, e1006103 (2016).
- 708 10. Goday, C., Rosario, M. & Summary, E. Chromosome elimination in Sciarid flies. *Wiley*  
709 *Online Libr.* doi:10.1002/1521-1878(200103)23:3<242::AID-BIES1034>3.0.CO;2-P.
- 710 11. Müller, F. & Tobler, H. Chromatin diminution in the parasitic nematodes *Ascaris suum*  
711 and *Parascaris univalens*. *Int. J. Parasitol.* **30**, 391–399 (2000).
- 712 12. Tobler, H. The differentiation of germ and somatic cell lines in nematodes. *Results Probl.*  
713 *Cell Differ.* **13**, 1–69 (1986).
- 714 13. Tobler, H., Müller, F., Back, E. & Aeby, P. Germ line - soma differentiation in *Ascaris*: A  
715 molecular approach. *Experientia* **41**, 1311–1319 (1985).
- 716 14. Pigozzi, M. I. & Solari, A. J. The germ-line-restricted chromosome in the zebra finch:  
717 Recombination in females and elimination in males. *Chromosoma* **114**, 403–409 (2005).
- 718 15. Biederman, M. K. *et al.* Discovery of the first germline-restricted gene by subtractive  
719 transcriptomic analysis in the Zebra finch, *Taeniopygia guttata*. *Curr. Biol.* **28**, 1620-  
720 1627.e5 (2018).
- 721 16. Kinsella, C. M. *et al.* Programmed DNA elimination of germline development genes in  
722 songbirds. *Nat. Commun.* **10**, (2019).
- 723 17. Magnúsdóttir, E. *et al.* A tripartite transcription factor network regulates primordial germ  
724 cell specification in mice. *Nat. Cell Biol.* **15**, 905–915 (2013).
- 725 18. Docker, M. F., Beamish, F. W. H., Yasmin, T., Bryan, M. B. & Khan, A. The Lamprey  
726 Gonad. in *Lampreys: Biology, Conservation and Control Vol 2*, 1–186 (2019) edited by  
727 M. F. Docker. Springer.
- 728 19. Davidson, N. M. & Oshlack, A. Necklace: combining reference and assembled  
729 transcriptomes for more comprehensive RNA-Seq analysis. *Gigascience* **7**, 1–6 (2018).
- 730 20. Smith, J. J. *et al.* The sea lamprey germline genome provides insights into programmed  
731 genome rearrangement and vertebrate evolution. *Nat. Genet.* **50**, 270–277 (2018).
- 732 21. Timoshevskaya, N. Difcover. (2019) Available online: <https://github.com/timnat/DifCover>  
733 (accessed on 12 September 2019)
- 734 22. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).
- 735 23. Hardisty, M. Gonadogenesis, sex differentiation and gametogenesis. In: Hardisty MW,  
736 Potter IC (eds) *The biology of lampreys*. *Acad. Press. New York* 295–360 (1971).
- 737 24. Lowartz, S. M. & Beamish, F. W. H. Novel perspectives in sexual lability through  
738 gonadal biopsy in larval sea lampreys. *J. Fish Biol.* **56**, 743–757 (2000).
- 739 25. Docker, M. F., William, F. & Beamish, H. Age, growth, and sex ratio among populations  
740 of least brook lamprey, *Lampetra aepyptera*, larvae: an argument for environmental sex  
741 determination. *Environ. Biol. Fishes* **41**, 191–205 (1994).

- 742 26. Johnson, N. S., Swink, W. D. & Brenden, T. O. Field study suggests that sex  
743 determination in sea lamprey is directly influenced by larval growth rate. *Proc. R. Soc. B*  
744 *Biol. Sci.* **284**, (2017).
- 745 27. Capel, B. Vertebrate sex determination: evolutionary plasticity of a fundamental switch.  
746 *Nat. Rev. Genet.* 2017 1811 **18**, 675–689 (2017).
- 747 28. Xu, X. *et al.* Evolutionary transition between invertebrates and vertebrates via methylation  
748 reprogramming in embryogenesis. *Natl. Sci. Rev.* **6**, 993–1003 (2019).
- 749 29. Wang, X. *et al.* The Role of DNA Methylation Reprogramming during Sex Determination  
750 and Transition in Zebrafish. *Genomics. Proteomics Bioinformatics* (2021)  
751 doi:10.1016/j.gpb.2020.10.004.
- 752 30. Sacerdot, C., Louis, A., Bon, C., Berthelot, C. & Roest Crollius, H. Chromosome  
753 evolution at the origin of the ancestral vertebrate genome. *Genome Biol.* **19**, 166 (2018).
- 754 31. Nguyen, N. T. T., Vincens, P., Crollius, H. R. & Louis, A. Genomicus 2018: Karyotype  
755 evolutionary trees and on-the-fly synteny computing. *Nucleic Acids Res.* **46**, D816–D822  
756 (2018).
- 757 32. Windley, S. P. & Wilhelm, D. Signaling pathways involved in mammalian sex  
758 determination and gonad development. *Sex Dev.* **9**, 297–315 (2016).
- 759 33. Tao, W. *et al.* Characterization of gonadal transcriptomes from Nile Tilapia (*Oreochromis*  
760 *niloticus*) reveals differentially expressed genes. *PLoS One* **8**, e63604 (2013).
- 761 34. Jia, Y. *et al.* Transcriptome analysis of three critical periods of ovarian development in  
762 Yellow River carp (*Cyprinus carpio*). *Theriogenology* **105**, 15–26 (2018).
- 763 35. Chen, Y. *et al.* Gonadal transcriptome sequencing of the critically endangered *Acipenser*  
764 *dabryanus* to discover candidate sex-related genes. *PeerJ* **2018**, e5389 (2018).
- 765 36. Takemoto, K. *et al.* Sycp2 is essential for synaptonemal complex assembly, early meiotic  
766 recombination and homologous pairing in zebrafish spermatocytes. *PLOS Genet.* **16**,  
767 e1008640 (2020).
- 768 37. Lavery, R. *et al.* Testicular differentiation occurs in absence of R-spondin1 and Sox9 in  
769 mouse sex reversals. *PLOS Genet.* **8**, e1003170 (2012).
- 770 38. Barbara, N. & Humphrey H, Y. Gonadal identity in the absence of pro-testis factor SOX9  
771 and pro-ovary factor beta-catenin in mice. *Biol. Reprod.* **93**, (2015).
- 772 39. Piprek, R. P., Kloc, M., Mizia, P. & Kubiak, J. Z. The central role of cadherins in gonad  
773 development, reproduction, and fertility. *Int. J. Mol. Sci.* **21**, 1–21 (2020).
- 774 40. Wolgemuth, D. J., Manterola, M. & Vasileva, A. Role of cyclins in controlling  
775 progression of mammalian spermatogenesis. *Int. J. Dev. Biol.* **57**, 159 (2013).
- 776 41. Venables, J. P. *et al.* RBMY, a probable human spermatogenesis factor, and other hnRNP  
777 G proteins interact with Tra2 $\beta$  and affect splicing. *Hum. Mol. Genet.* **9**, 685–694 (2000).
- 778 42. Jungmin, L. *et al.* Developmental stage-specific expression of Rbm suggests its

- 779 involvement in early phases of spermatogenesis. *Mol. Hum. Reprod.* **10**, 259–264 (2004).
- 780 43. Dos Reis, M. *et al.* Using Phylogenomic Data to Explore the Effects of 3 Relaxed Clocks  
781 and Calibration Strategies on 4 Divergence Time Estimation: Primates as a Test Case.  
782 *Syst. Biol.* **67**, 594–615 (2018)
- 783 44. Kumar, S., Stecher, G., Suleski, M. & Hedges, S. B. TimeTree: A resource for timelines,  
784 timetrees, and divergence times. *Mol. Biol. Evol.* **34**, 1812–1819 (2017).
- 785 45. Kuraku, S., Meyer, A. & Kuratani, S. Timing of genome duplications relative to the origin  
786 of the vertebrates: Did cyclostomes diverge before or after? *Mol. Biol. Evol.* **26**, 47–59  
787 (2009).
- 788 46. Smith, J. J. & Keinath, M. C. The sea lamprey meiotic map improves resolution of ancient  
789 vertebrate genome duplications. *Genome Res.* **25**, 1081–1090 (2015).
- 790 47. Simakov, O. *et al.* Deeply conserved synteny resolves early events in vertebrate evolution.  
791 *Nat. Ecol. Evol.* **4**, 820–830 (2020).
- 792 48. Nakatani, Y. *et al.* Reconstruction of proto-vertebrate, proto-cyclostome and proto-  
793 gnathostome genomes provides new insights into early vertebrate evolution. *Nat.*  
794 *Commun.* **12**, 4489 (2021).
- 795 49. Andrews, S. FastQC: a quality control tool for high throughput sequence data. (2010).
- 796 50. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina  
797 sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
- 798 51. Pertea, M., Kim, D., Pertea, G. M., Leek, J. T. & Salzberg, S. L. Transcript-level  
799 expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat.*  
800 *Protoc.* **11**, 1650–1667 (2016).
- 801 52. Haas, B. J. *et al.* De novo transcript sequence reconstruction from RNA-seq using the  
802 Trinity platform for reference generation and analysis. *Nat. Protoc.* **8**, 1494–1512 (2013).
- 803 53. Bryant, D. M. *et al.* A tissue-mapped axolotl de novo transcriptome enables identification  
804 of limb regeneration factors. *Cell Rep.* **18**, 762–776 (2017).
- 805 54. Emms, D. M. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome  
806 comparisons dramatically improves orthogroup inference accuracy. *Genome Biol.* **16**, 157  
807 (2015).
- 808 55. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion  
809 for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 810 56. Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for  
811 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139  
812 (2010).
- 813 57. Mi, H., Muruganujan, A., Ebert, D., Huang, X. & Thomas, P. D. PANTHER version 14:  
814 More genomes, a new PANTHER GO-slim and improvements in enrichment analysis  
815 tools. *Nucleic Acids Res.* **47**, D419–D426 (2019).

- 816 58. Supek, F., Bošnjak, M., Škunca, N. & Šmuc, T. REVIGO summarizes and visualizes long  
817 lists of gene ontology terms. *PLoS One* **6**, e21800 (2011).
- 818 59. Makiyan, Z. Studies of gonadal sex differentiation. *Organogenesis* **12**, 42–51 (2016).
- 819 60. Quinn, A. & Koopman, P. The molecular genetics of sex determination and sex reversal in  
820 mammals. *Semin. Reprod. Med.* **30**, 351–363 (2012).
- 821