1 **Expansion, retention and loss in the Acyl-CoA Synthetase *"Bubblegum"* (*Acsbg*)**

2 **gene family in vertebrate history**

3

4 Mónica Lopes-Marques[1*], André M. Machado[1], Raquel Ruivo[1], Elza Fonseca[1,2], Estela

5 Carvalho[1], and L. Filipe C. Castro[1,2*]

6

7 [1]Interdisciplinary Centre of Marine and Environmental Research (CIIMAR/CIMAR),

8 University of Porto (U.Porto), Matosinhos, Portugal

9 [2]Faculty of Sciences (FCUP), Department of Biology, University of Porto (U.Porto), Porto,

10 Portugal

11 * Address correspondence to:

12 Mónica Lopes-Marques & Luís Filipe Costa Castro

13 Interdisciplinary Centre for Marine and Environmental Research (CIIMAR)

14 Av. General Norton de Matos s/n, 4450-208 Matosinhos, Portugal

15 Tel.: +351223401800

16 Fax: +351223390608

17 Email: monicaslm@hotmail.com; filipe.castro@ciimar.up.pt

18 **Keywords**: Acyl-CoA synthetases; fatty acid activation; ACSBG; gene duplication;

19 vertebrates

1

20    **ABSTRACT**

21    Fatty acids (FAs) constitute a considerable fraction of all lipid molecules with a

22    fundamental role in numerous physiological processes. In animals, the majority of

23    complex lipid molecules are derived from the transformation of FAs through several

24    biochemical pathways. Yet, for FAs to enroll in these pathways they require an

25    activation step. FA activation is catalyzed by the rate limiting action of Acyl-CoA

26    synthases. Several Acyl-CoA enzyme families have been previously described and

27    classified according to the chain length of FA they process. Here, we address the

28    evolutionary history of the ACSBG gene family which activates, FA with more than 16

29    carbons. Currently, two different ACSBG gene families, *ACSBG1* and *ACSBG2*, are

30    recognized in vertebrates. We provide evidence that a wider and unequal *ACSBG* gene

31    repertoire is present in vertebrate lineages. We identify a novel *ACSBG-like* gene

32    lineage which occurs specifically in amphibians, ray finned fish, coelacanths and

33    chondrichthyes named *ACSBG3*. Also, we show that the *ACSBG2* gene lineage

34    duplicated in the Theria ancestor. Our findings, thus offer a far richer understanding on

35    FA activation in vertebrates and provide key insights into the relevance of comparative

36    and functional analysis to perceive physiological differences, namely those related with

37    lipid metabolic pathways.

## 1. INTRODUCTION

Lipids represent a complex group of biomolecules present in all living organisms, playing a key role in numerous biological processes, such as inflammatory response, reproduction, biological membranes and energy sourcing and storage. Additionally, they participate in the overall the homoeostasis as signal molecules, cofactors, and endogenous ligands for nuclear receptors (Wall et al., 2010; Robinson and Mazurak, 2013; Grygiel-Górniak, 2014). Aside from sterol lipids, all remaining lipids are obtained from the endogenous elaboration of fatty acid (FA) molecules (Watkins et al., 2007). Yet, for FAs to enroll in any anabolic and catabolic process they require an activation step. Thus, FA activation is a critical rate limiting step of FA metabolism. FA activation was first recognized in 1948 and referred to as "sparking" or "priming" (Grafflin and Green, 1948; Knox et al., 1948). This enzymatic step consists of a two-step thioesterification reaction catalyzed by Acyl-CoA synthetase (ACS), resulting in a thioester with coenzyme A (CoA) (Watkins et al., 2007).

Several ACS involved in FA activation have been previously identified and organized according to the degree of unsaturation and chain length of the FAs favored as substrate: the short-chain ACS-Family (*ACSS*), medium-chain ACS-Family (*ACSM*), long-chain ACS-Family (*ACSL*), very long-chain ACS-Family (*ACSVL*), Bubblegum ACS-Family (*ACSBG*) and ACSFamily (*ACSF*) (Watkins et al., 2007; Soupene and Kuypers, 2008). Although some substrate preference overlap is observed, these enzymes also differ in tissue distribution and subcellular location, an indication of their highly specific role in FA metabolism (Watkins et al., 2007). ACS enzymes have been found to have a wide taxonomic distribution, with homologues ranging from Eubacteria to Plants and Metazoa, a clear indication of their pivotal role in lipid metabolism (Hisanaga et al., 2004; Soupene and Kuypers, 2008).

Despite their wide taxonomic occurrence, the genetic repertoire of ACS has been found to vary, namely in vertebrates (Castro et al., 2012; Lopes-Marques et al., 2013). For example, some studies have disclosed that both ACSL and ACSS gene family composition and function were shaped by events of gene/genome duplication in combination with differential loss. Moreover, multi-genome comparisons across a wide range of vertebrate species revealed novel and previously uncharacterized ACS

3

69   enzymes (Castro et al., 2012; Lopes-Marques et al., 2013). The present work seeks to

70   build on previous findings and further extend, the knowledge regarding the genetic

71   repertoire and distribution ACS in vertebrates namely the ACS *Bubblegum* (ACSBG)

72   gene family.

73   ACSBG enzymes, also known as lipidosin, activate FA with C16 to C24 (Moriya-Sato et

74   al., 2000; Steinberg et al., 2000; Pei et al., 2003). Presently, 2 members of the ACSBG

75   gene family have been identified and characterized in mammals, ACSBG1 and ACSBG2

76   (Pei et al., 2003; Watkins et al., 2007). Similarly, to the previously described ACS

77   enzymes, both ACSBG members display conserved sequence motifs, such as the

78   putative ATP-AMP signature motif for ATP binding (Motif I) and a motif for FA binding,

79   characteristic of the ACS gene family (Motif II) (Moriya-Sato et al., 2000; Watkins et al.,

80   2007). Notably, all known ACS, with the exception of human ACSBG2, contain a highly

81   conserved arginine (Arg-R) in Motif II (Pei et al., 2006). The replacement of this Arg by

82   histidine (His-H) in Human ACSBG2 was found to confer a biphasic pH optimum (pH 6.5

83   and pH 7.5) to the enzyme, in contrast to the monophasic activity at pH 7-7.5 of the

84   mouse orthologue (Pei et al., 2006). Yet, due to the degree of conservation of both

85   Motifs I and II, these have previously been used to seek and identify potential ACS

86   enzymes (Steinberg et al., 2000; Watkins et al., 2007).

87   ACSBG enzymes have been suggested to play a significant role in brain development

88   and reproduction (Moriya-Sato et al., 2000; Tang et al., 2001; Pei et al., 2006). Previous

89   reports with the *D. melanogaster* bubblegum mutant (termed bubblegum due to the

90   bubbly appearance of the lamina, a result of neurodegeneration and dilation of the

91   photoreceptor axons) and mouse, associated the disruption of *ACSBG1* to X-linked

92   adrenoleukodytrophy (X-ALD) (Min and Benzer, 1999; Moriya-Sato et al., 2000). X-ALD

93   is characterized by the accumulation of high levels of very long FA in plasma and

94   tissues, accompanied by neurodegeneration (Min and Benzer, 1999; Moriya-Sato et al.,

95   2000; Moser et al., 2002). On the other hand, *ACSBG2* plays an important role in

96   spermatogenesis and testicular development, being associated to male infertility

97   (Zheng et al., 2005; Fraisl et al., 2006). In agreement gene expression of *ACSBG1* is

98   found to be mainly restricted to brain, adrenal gland, gonads, spleen in mouse and

99   human (Moriya-Sato et al., 2000). In contrast, *ACSBG2* showed a more exclusive

4

100   expression pattern being highly expressed in the testis, followed by medulla and spinal

101   cord (Pei et al., 2006) .

102        Here, using a combination of phylogenetics, comparative genomics and gene

103   expression analysis we deduced the evolutionary history of the *ACSBG* gene family in

104   all major vertebrate lineages. Our findings illustrate the importance of comparative

105   analysis to address the role of adaptive evolution in the shaping of lipid metabolic

106   modifications between lineages.

## 2. MATERIALS & METHODS

### 2.1. DATABASE SEARCH AND PHYLOGENETIC ANALYSIS

NCBI GenBank release 220 June and release 221 August 2017 and Ensembl release 89 May and release 90 August 2017, databases were searched using tblastn and blastp to recover ACSBG-like sequences using human ACSBG1 (NP_055977) and ACSBG2 (NP_112186) amino acid sequences as query. All major vertebrate lineages such as mammals, birds, reptiles, amphibians, coelacanths, teleost fish, cartilaginous fish and cyclostomes were searched. Additionally, the following invertebrate lineages basal to chordates were also explored cephalochordates, hemichordates, Mollusca (Supplementary material 1).

Our search retrieved 121 ACSBG-like sequences; sharing a minimum 70% pairwise identity with corresponding query sequence for mammals; 60% for bird's reptiles and amphibians; 50% identity for teleost's and chondrichthyes and finally 40% identity for invertebrates. The collected sequences were aligned and inspected with partial sequences being removed, leaving a 119 full ORF or near full ORF sequences for phylogenetic analysis (Supplementary material 1). Amino acid sequences were aligned in MAFFT with the L-INS-I method (Katoh et al., 2005; Katoh and Toh, 2008). In the resulting alignment, all columns containing 90% gaps were stripped leaving a total of 787 positions for phylogenetic analysis. A second sequence alignment containing 121 ACSBG-like sequences including the truncated sequences of *Xenopus tropicalis* and *Xenopus leavis* was performed using the same method leaving a total of 788 positions for phylogenetic analysis. Both alignments were then individually submitted to PhyML3.0 server (Guindon et al., 2010), with evolutionary model determined automatically, resulting in the selection JTT+G+I in both cases. The branch support for phylogenetic trees was calculated using aBayes. The resulting trees were visualized and edited in Fig. Tree V1.3.1 available at http://tree.bio.ed.ac.uk/software/figtree/ and rooted with the invertebrate sequences.

### 2.2. SYNTENY AND PARALOGY ANALYSIS

Using as reference the human and teleost *ACSBG loci*, synteny maps of the genomic neighbourhoods of the *ACSBG1, ACSBG2* and *ACSBG3* gene were assembled in a set of

6

138  species representative of the major lineages analysed. The following genome

139  assemblies available in NCBI were accessed for *Homo sapiens* - GCF_000001405.33,

140  *Monodelphis domestica* - GCF_000002295.2, *Gallus gallus* - GCF_000002315.4,

141  *Pelodiscus sinensis* - GCF_000230535.1, *X. tropicalis* - GCF_000004195.3, *X. laevis* -

142  GCF_001663975.1, *Latimeria chalumnae* -  GCF_000225785.1, *Oryzias latipes* -

143  GCA_000313675.1, *Astyanax mexicanus* - GCA_000372685.2, *Danio rerio* -

144  GCA_000002035.4, *Lepisosteus oculatus* - GCA_000242695.1, *Callorhinchus milii* -

145  GCA_000165045.2, *Branchiostoma floridae* - GCA_000003815.1 and *S. kowalevskii* -

146  GCF_000003605.2. For *Tetraodon nigroviridis* and *Petromyzon marinus* genome

147  assemblies TETRAODON 8.0, Mar 2007 and Pmarinus_7.0, Jan 2011 available in

148  Ensembl were accessed. Paralogy analysis of the *ACSBG locus* was conducted using the

149  reconstructed ancestral chordate genome as described in (Putnam et al., 2008).

150

151  **2.3. RNA ISOLATION AND ASCBG TISSUE EXPRESSION PANEL IN *X. TROPICALIS***

152

153  *X. tropicalis* (African clawed frog) tissues (brain, skin, heart, liver, spleen, pancreas,

154  kidney, intestine, testis and ovary) were kindly provided by O. Brochain (CNRS, Orsay).

155  Total RNA was purified using the Illustra RNAspin Mini RNA Isolation Kit animal tissues

156  protocol (GE Healthcare) with on-column DNase I digestion. RNA quality was assessed

157  by electrophoresis and its concentration determined using a microplate

158  spectrophotometer (Take 3 and Synergy HT Multi-Mode Microplate Reader, Biotek).

159  First-strand cDNA was synthesized from 250ng RNA using the iScriptTM cDNA Synthesis

160  Kit (Bio-Rad), according to the manufacturer's instructions.

161  Forward and reverse primers sets were designed to flank an intron and to avoid

162  genomic DNA amplification. Primers sets were created for the following genes *ACSBG1*

163  - Forward-5'TTTGCCAGGATGTTGGAAGT3', Reverse-5'AAAGCTTCCACGTGCTCTGT 3',

164  annealing at 57°; *ACSBG2* - Forward-5' CTTTTCTGGGGACGTCATGT 3', Reverse-5'

165  TTGGAACCTGCTCTTTGAGG 3', annealing at 55° and *ASCGB3* - Forward-5'

166  TGCAGTCTTTGCTACGTTGG 3' reverse-5' ACAAACAGAGCTCCCCTGTG 3', annealing at

167  57°. To assess the quality of *X. tropicalis* cDNA two sets of primers targeting

168  housekeeping genes were included for β-actin – Forward-5' GGTCGCCCAAGACATCAG3',

169  Reverse-5'GCATACAGGGACAACACA annealing at 57º and for EEF1A1 – Forward-

170     5´TCGTTAAGGAAGTCAGCACA3´ and Reverse5´CATGGTGCATTTCAACAGAT3´ annealing

171     at 57º. PCR reactions were all performed using 2 µl of *X. tropicalis* cDNA and Phusion®

172     Flash high-fidelity Master Mix (FINNZYMES). PCR parameters were as follows: initial

173     denaturation at 98°C for 10 s, followed by 30 cycles of denaturation at 98°C for 1 s,

174     annealing for 5 s and elongation at 72°C for 10s and a final step of elongation at 72°C

175     for 1 min. PCR products were then loaded onto 2% agarose gel stained with GelRed and

176     run in TBE buffer at 80 V.

177

178

179

180     **2.4. ACSBG EXPRESSION ANALYSIS THROUGH RNA-SEQ**

181     The RNA-Seq analysis was performed using a collection of tissues datasets from seven

182     species Human (*H. sapiens*), mouse (*M. musculus*), chicken (*G. gallus*), western clawed

183     frog (*X. tropicalis*), zebrafish (*D. rerio*), spotted gar (*L. oculatus*) and elephant shark (*C.*

184     *milii*), available in National Centre for Biotechnology (NCBI) Sequence Read Archive

185     (SRA) (https://www.ncbi.nlm.nih.gov/sra/) (Supplementary material 2). To standardize

186     datasets from different sources, all files were converted to FASTQ file format and

187     sequence quality trimming was performed using Trimmomatic v 0.36 (Bolger et al.,

188     2014). Reads with 36bp in length and an average score of 20 phred were selected for

189     further analysis.

190     Reference sequences and respective annotation files of each specie were collected

191     from NCBI and Ensembl (Release 89) (Yates et al., 2016) (supplementary material 3).

192     For both elephant shark and western clawed frog the reference sequences for RNAseq

193     mapping were retrieved from NCBI (ftp://ftp.ncbi.nih.gov/genomes/refseq/), while for

194     the remaining species the reference sequences for mapping were retrieved from

195     Ensembl database (ftp://ftp.ensembl.org/pub/release-89/) (Supplementary material

196     3). Trimmed and groomed reads from each dataset were mapped to their respective

197     reference using Bowtie2 (Langmead and Salzberg, 2012), and the transcript

198     quantification was calculated in transcript per million (TPM), with RSEM v.1.2.31

199     software (Li and Dewey, 2011). TPM values for each gene were taken as evidence of

200     relative gene expression, low TPM values (< 0.5) were considered unreliable and

201 substituted with zero. To complete this exploratory gene expression analysis, the TPM

202 values were $\log_2$-transformed after adding a value of one.

203 ## 3. RESULTS AND DISCUSSION

204 ### 3.1. DATABASE MINING AND PHYLOGENETIC ANALYSIS REVEALS NOVEL MEMBERS ACSBG GENE

205 #### FAMILY

206 Initial blast searches identified ACSBG-like sequences and recovered a larger than

207 anticipated number of sequence hits. ACSBG1 and ACSBG2-like sequences were found

208 in species from the following vertebrate lineages: mammals, birds, reptiles,

209 amphibians, holostei, coelacantiforms and teleostei. In chondrichthyans and

210 cyclostomes no ACSBG1-like sequences were retrieved. Moreover, an additional

211 uncharacterized ACSBG2-like sequence was also retrieved in some mammalian species.

212 Database searches also recovered a novel set of ACSBG-like sequences in four

213 gnathostome lineages: amphibians (Western clawed frog and African clawed frog)

214 chondrichthyans (elephant shark), holostei (spotted gar), coelacantiforms (coelacanth)

215 and teleostei. However, amphibian sequences were considerably shorter than the

216 remaining ACSBG, thus being excluded from the main phylogenetic analysis.

217 To disclose the orthology of these various sequences a phylogenetic analysis was

218 performed. The resulting tree topology displays 3 well supported clades in vertebrates.

219 The first group contained all ACSBG1 sequences from mammals, reptiles, birds,

220 amphibians, coelacanths and teleost fish, with no representatives of chondrichtyes and

221 cyclostomes. Besides the ACSBG1 clade, we find a sister clade comprising ACSBG2

222 sequences. This contains ACSBG2 previously described in mammals (ACSBG2a) and an

223 uncharacterized ACSBG2-like (ACSBG2b) identified in the present work. The tree

224 topology suggests that the both mammalian ACSBG2 sequences are related by a

225 duplication event that took place in the ancestor of Theria. Out grouping the

226 mammalian ACSBG2 sequences we find the ACSBG2 sequences from birds, reptiles,

227 amphibians, coelacanths, chondrichthyes, teleost fish and cyclostomes. Thus, *bona fide*

228 orthologues of ACSBG2 are represented across all major vertebrate classes. Within the

229 third clade, we find a novel uncharacterized group of ACSBG sequences which have

230 sequence representatives in coelacanths, teleost fish, holostei and chondrichthyes. We

231 name this novel sequence ACSBG3. Finally, placed basally to all vertebrate sequences

232 we find invertebrate ACSBG sequences. The present tree topology provides robust

233 indications that the diversification of the ACSBG gene family occurred in the vertebrate

10

234   ancestor. A second phylogenetic analysis was run separately to include the

235   uncharacterized short *Xenopus sp.* ACSBG sequences (Supplementary material 4).

236   Although the overall tree topology is conserved with the main analysis (Fig. 1), we find

237   that the *Xenopus sp.* ACSBG sequences are placed basally to all vertebrate clades

238   hindering the identification of their orthology. This placing of *Xenopus sp*. ACSBG-like

239   sequences correlates to the highly divergent nature observed in the sequence

240   alignment, with these amphibian sequences being considerably shorter and displaying

241   a poor conservation of the AMP-binding motif (see section 3.3).

242   Additionally, we find that mammalian ACSBG2a and 2b sequences are placed in long

243   branches in both phylogenetic analysis (Fig.1 and Supplementary material 4) suggesting

244   an accelerated evolution and divergence of these sequences further analysed in section

245   3.3.

246

247   **3.2. A NOVEL ACSBG GENE, ACSBG3, IS AN OHNOLOG GONE MISSING IN AMNIOTES**

248   Phylogenetic analysis suggests that the ACSBG gene family expanded in the vertebrate

249   ancestor. This time frame coincides with the proposed timing of two round of whole

250   genome duplication (2R WGD) in the ancestral vertebrate approximately 500MYA

251   (Ohno, 1970; Dehal and Boore, 2005). Yet, while it is generally accepted that all

252   gnathostomes underwent the 2R-WGD, the extent of these genome duplications in

253   cyclostomes still remains a matter of debate (Smith and Keinath, 2015).

254   To complement the phylogenetic analysis, validate events of gene duplication/loss,

255   resolve the orthology of the *Xenopus sp.* uncharacterized ACSBG and the origin of

256   ACSBG3 sequences, the genomic *locus* of each *Acsbg* gene was examined in a set of

257   representative species (Fig. 2). Comparative synteny analysis of the ACSBG1 *locus*

258   reveals a high degree of conservation of neighbouring gene families throughout all the

259   analysed lineages (Fig.2A). *ACSBG1* is localized in human chromosome 15, being

260   flanked by gene families such as the *IDH3A, CIB2, WDR61* and *CRABP1*. These flanking

261   gene families are also present in the vicinity of ACSBG1 *locus* in the all analysed

262   lineages (Fig. 2A). In the case of *C. milii* although no *ACSBG1* gene was found, synteny

263   analysis reveals that the *locus* organization is conserved, suggesting gene loss in this

264   lineage (Fig. 2A). Regarding the cyclostomes, extensive blast searches did not retrieve

265   an ACSBG1 sequence; synteny analysis uncovered a fragmented *locus* segregated into

11

266   at least two distinct scaffolds. Therefore, the absence of this gene in cyclostomes may

267   be attributed to poor genome coverage or to gene loss (Fig. 2A).

268   The human ACSBG2 gene resides in chromosome 19 and is flanked by the following

269   gene families: *RFX2, RANBP3, MLLT1* and *ACER1*. The *locus* architecture is conserved in

270   all species analysed. In cyclostomes, the *ACSBG2 locus* is disjointed and distributed

271   among several scaffolds thus the absence of ACSBG2 in cyclostomes remains similarly

272   to ACSBG1 unresolved (Fig. 2B). In the case of mammalian duplicates, we find that

273   *ACSBG2a* and *ACSBG2b* are located side by side in the *M. domestica.* Synteny analysis

274   of this *locus* in other mammals presenting both *ACSBG2a* and *2b* (data not shown) is

275   coincident with the observation for *M. domestica* supporting the hypothesis that

276   *ACSBG2a* and *2b* arose through tandem duplication in the ancestor of therian

277   mammals, with ACSBG2b being later lost in Haplorhini. Finally, we find that in the

278   *ACSBG3 locus*, despite the lesser conservation*,* some neighbouring gene families such

279   as *HINT2, SPAG8, RGP1* and *GBA2* are preserved in the majority of the analysed

280   lineages. Using these conserved neighbouring gene families, the corresponding *locus*

281   was mapped in birds and mammals to address the loss of *ACSBG3* in these lineages

282   (Fig. 2C). ACSBG3 is also absent in reptiles and the analysis of the *locus* revealed that it

283   is fragmented in various species examined (*Anolis carolinesis*, *Thamnophis sirtalis*,

284   *Alligator mississippiensis* and *Chrysemys picta*), hindering the validation of ACSBG3 loss

285   in this lineage.

286   We next investigated the synteny maps for the single ACSBG *locus* from two

287   invertebrate cephalochordates. Here we find that the *B. floridae locus* retains a

288   conserved gene family arrangement, namely with the presence of *HERC1-like* gene*,*

289   whose orthologue is found in the vicinity of vertebrate *ACSBG1* (Fig. 2 A and D

290   indicated in red). Similarly, the hemichordate *S. kowalevskii* also displays a conserved

291   neighbouring gene family, *CHRNA3* with the vertebrate orthologue placed in the

292   *ACSBG1 locus* (Fig. 2A and D indicated in red). Finally, to address the hypothesis that

293   ACSBG gene expansion took place with the 2R WGD the location of ACSBG and

294   neighbouring genes (with described paralogues underlined genes in Fig.2A, B, C and D)

295   were mapped to the predicted ancestral paralogons as described by Putnam *et al* 2008

296   (Putnam et al., 2008). Next, ancestral paralogons were mapped back to the same

297   ancestral linkage group, LG2 (Putnam et al., 2008) indicating that the *ACSBG loci* are

12

298 related by genome duplication, strongly suggesting that vertebrate ACBG diversity

299 arose with the 2R WGD (Fig. 2 F).

300

### 3.3. SEQUENCE ANALYSIS AND GENE EXPRESSION

302 To further characterize the novel *ACSBG2b* and *ACSBG3* a sequence alignment was

303 performed to inspect the typical ACS enzyme motifs (Watkins et al., 2007). The analysis

304 of this alignment revealed that the predicted AMP-binding domain (Motif I Fig. 3A and

305 3B), a highly conserved motif in all ACS enzymes from bacteria to humans (Black et al.,

306 1997; Steinberg et al., 2000; Weimar et al., 2002; Karan et al., 2003), is conserved in

307 the vast majority of the sequences collected with the exception of the novel ACSBG2b

308 (Fig.3B) and the ACSBG3 sequence in *Xenopus sp.*(see Supplementary material 5 for

309 alignment of the full 121 sequences). An indication that residues within Motif I play a

310 fundamental role in ACS catalytic activity was found in previous studies were the

311 mutation of residues within Motif I (positions 1, 2, 4, 5 and 10) in *E. coli* considerably

312 reduced the catalytic activity, while the replacement of residues 1 and 5 in *S. cerevisiae*

313 resulted in a minor reduction of enzymatic activity (Fig. 3A grey arrows) (Weimar et al.,

314 2002; Zou et al., 2002). Thus, the low conservation of this motif in the mammalian-

315 specific ACSBG2b strongly suggests that these enzymes may show an alternative

316 function or *modus operandi*. In the analysis of the *Xenopus sp.* ACSBG3, we find that

317 this motif differs from the remaining ACSBG3 identified, being disrupted with the

318 deletion of 3bp (Supplementary material 5). Again, this observation suggests an

319 alternative role for the enzyme given that AMP binding is essential for FA activation.

320 Regarding Motif II, also known as the ACS signature-motif and proposed to be involved

321 in acyl chain length specificity (Black et al., 1997), we find that again ACSBG2b displays

322 a divergent sequence when compared to the remaining ACSBG analysed here.

323 Interestingly, we observe that ACSBG2b presents an Asparagine residue (Asn-N) instead

324 of the highly conserved Arginine (Fig. 3A and B black arrow). Notably, human ACSBG2a

325 harbours a Histidine in this position, representing the single case described to date of

326 an ACSBG without an Arginine (Pei et al., 2006). Reverse mutation of the Histidine

327 within Motif II in human ACSBG2a showed that this residue assumes a critical role in

328 determining the optimal pH for this enzyme (Pei et al., 2006). Additionally, this

13

329    replacement (Asn) is only observed for placental mammals, with marsupials retaining

330    the conserved Arginine (Fig 3 B). Next, Motif V (KXX(R,K) is a conserved motif found in

331    several members of the ACS enzymes families and contains a conserved K- Lys

332    demonstrated to be essential for the catalytic function of ACS in *S. enterica* propionyl-

333    coA synthetase and ACS activity of murine ACSF2 (Horswill and Escalante-Semerena,

334    2002). Here we find that Motif V is conserved in all recovered ACSBG sequences with

335    the exception of *Xenopus sp* ACSBG3 due to the short nature of these sequences (see

336    Supplementary material 5). Finally, the Motifs III and IV, identified by Hisanaga *et al*

337    2008 (Hisanaga et al., 2004), are found to be conserved in the majority of analysed

338    sequences, with the exception of a conservative replacement in Motif IV in ACSBG2b.

339    The highly conserved Histidine is replaced by biochemically similar residue, tyrosine,

340    having a minor or no predicted impact.

341    In an attempt to infer the function of the newly identified *ACSBG2b* and *ACSBG3*, and

342    address the retention of these genes in a restricted number of lineages, we next

343    performed an expression analysis using available RNA-Seq SRAs (Fig. 3C). Similarly, to

344    previous reports (Moriya-Sato et al., 2000; Tang et al., 2001; Pei et al., 2006), relative

345    expression profiles reveal that *ACSBG1* expression is mainly limited to brain and

346    gonads, with the exception of *D. rerio* for which the liver stands as the main expression

347    site. On the other hand, the expression profile of non-mammalian vertebrate *ACSBG2*

348    was found to be more extensive than in mammals (Pei et al., 2006) with expression

349    detected in all analysed tissues of *C. milii*, *D. rerio*, *L. oculatus*, *G. gallus*, and *X.*

350    *tropicalis*. Interestingly, the expression analysis of mammalian specific duplicates

351    *ACSBG2a* and *ACBG2b* shows a confined expression of the duplicates essentially in

352    testis, with a relatively low expression of *ACSBG2a* detected in human kidney.

353    Regarding the gene expression profile of *ACSBG3 in X. tropicalis, L. oculatus*, and in

354    *ACSBG3a* of *C. milli*, we find a localized and high relative expression in ovary and in

355    testis. Semi-quantitative PCR expression analysis of ACSBG3 from *X. tropicalis* is in

356    accordance with *in silico* RNAseq analysis (Fig. 3D), with ACSBG1 expression confined

357    to brain and testis, while ACSBG2 is detected all tissues except ovary and finally

358    ACSBG3 being restricted to testis and brain (Fig. 3D).

359  High expression of *ACSBG3* in gonads is indicative that this enzyme may play an

360  important role in reproduction similarly to the role of *ACSBG2* (Pei et al., 2006). Finally,

361  for *ACSBG3* in *C. milii* no expression was detected in any of the analysed tissues.

362

363  **3.4. Evolutionary history of ACSBG gene family**

364  Using a multi-comparative approach, including database searches, phylogenetic and

365  synteny analysis we have uncovered a larger than anticipated genetic repertoire of

366  *ACSBG* genes in vertebrates. We find that the initial expansion of the *ACSBG* gene

367  family from which arose *ACSBG1 ACSBG2* and *ACSBG3* is coincident with the 2R WGD,

368  with representative gene orthologues present in several gnathostome lineages (Fig. 4).

369  The detailed analysis of the *ACSBG* gene repertoire revealed a differential retention of

370  *ACSBG3,* with this paralogue being lost in birds, mammals and possibly in reptiles,

371  while being retained in teleosts, amphibians and chondrichthyes. The identification of

372  additional ACSBG enzymes in teleosts correlates with previous studies, where

373  differential paralogue retention led to the maintenance of extra ACS enzyme

374  paralogues, namely *ACSL2* and *ACSS1b* in teleosts (Fraisl et al., 2006; Wall et al., 2010).

375  The preservation of duplicated genes is often observed when the corresponding

376  transcript, in this case ACS, is in high demand (Zhang, 2003). Thus, one may

377  hypothesize that the preservation of additional ACS duplicates in teleosts is a means to

378  fulfil a high demand of FA activation given that FA oxidation is considered to be the

379  main energy source in this lineage (Tocher, 2003). Finally, further duplications were

380  observed in the ancestor of mammals with the tandem duplication of *ACSBG2* and also

381  in specific lineages such as the *ACSBG3* in *C. milii* and *T. nigroviridis* and *ACSBG1 L.*

382  *chalumnae* (Fig. 4).

383  **4. CONCLUSION**

384

385  Our findings suggest that FA activation metabolic modules, including the ACSBG gene

386  family, have significantly diversified upon vertebrate radiation as a consequence of

387  genome duplication, lineage specific duplication and losses.

388

389

398

399    **REFERENCES**

400

401

402    Black, P.N., Zhang, Q., Weimar, J.D. and DiRusso, C.C., 1997. Mutational Analysis of a Fatty Acyl-
403        Coenzyme A Synthetase Signature Motif Identifies Seven Amino Acid Residues That
404        Modulate Fatty Acid Substrate Specificity. Journal of Biological Chemistry 272, 4896-
405        4903.

406    Bolger, A.M., Lohse, M. and Usadel, B., 2014. Trimmomatic: a flexible trimmer for Illumina
407        sequence data. Bioinformatics 30, 2114-20.

408    Castro, L.F.C., Lopes-Marques, M., Wilson, J.M., Rocha, E., Reis-Henriques, M.A., Santos, M.M.
409        and Cunha, I., 2012. A novel Acetyl-CoA synthetase short-chain subfamily member 1
410        (Acss1) gene indicates a dynamic history of paralogue retention and loss in vertebrates.
411        Gene 497, 249-255.

412    Dehal, P. and Boore, J.L., 2005. Two Rounds of Whole Genome Duplication in the Ancestral
413        Vertebrate. PLOS Biology 3, e314.

414    Fraisl, P., Tanaka, H., Forss-Petter, S., Lassmann, H., Nishimune, Y. and Berger, J., 2006. A novel
415        mammalian bubblegum-related acyl-CoA synthetase restricted to testes and possibly
416        involved in spermatogenesis. Arch Biochem Biophys 451, 23-33.

417    Grafflin, A.L. and Green, D.E., 1948. Studies on the cyclophorase system; the complete
418        oxidation of fatty acids. J Biol Chem 176, 95-115.

419    Grygiel-Górniak, B., 2014. Peroxisome proliferator-activated receptors and their ligands:
420        nutritional and clinical implications - a review. Nutrition Journal 13, 17.

421    Guindon, S., Dufayard, J.-F., Lefort, V., Anisimova, M., Hordijk, W. and Gascuel, O., 2010. New
422        Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the
423        Performance of PhyML 3.0. Syst Biol 59, 307-321.

424    Hisanaga, Y., Ago, H., Nakagawa, N., Hamada, K., Ida, K., Yamamoto, M., Hori, T., Arii, Y.,
425        Sugahara, M., Kuramitsu, S., Yokoyama, S. and Miyano, M., 2004. Structural basis of the
426        substrate-specific two-step catalysis of long chain fatty acyl-CoA synthetase dimer. J
427        Biol Chem 279, 31717-26.

428    Horswill, A.R. and Escalante-Semerena, J.C., 2002. Characterization of the Propionyl-CoA
429        Synthetase (PrpE) Enzyme of Salmonella enterica:  Residue Lys592 Is Required for
430        Propionyl-AMP Synthesis. Biochemistry 41, 2379-2387.

431    Karan, D., Lesbats, M., David, J.R. and Capy, P., 2003. Evolution of the AMP-forming Acetyl-CoA
432        Synthetase Gene in the Drosophilidae Family. Journal of Molecular Evolution 57, S297-
433        S303.

17

434  Katoh, K., Kuma, K.-i., Toh, H. and Miyata, T., 2005. MAFFT version 5: improvement in accuracy
435        of multiple sequence alignment. Nucleic Acids Research 33, 511-518.

436  Katoh, K. and Toh, H., 2008. Recent developments in the MAFFT multiple sequence alignment
437        program. Brief Bioinform 9, 286-98.

438  Knox, W.E., Noyce, B.N. and Auerbach, V.H., 1948. Studies on the cyclophorase system: III.
439        Obligatory Sparking of the fatty acid oxidaton. Journal of Biological Chemistry 176, 117-
440        122.

441  Langmead, B. and Salzberg, S.L., 2012. Fast gapped-read alignment with Bowtie 2. Nat Methods
442        9, 357-9.

443  Li, B. and Dewey, C.N., 2011. RSEM: accurate transcript quantification from RNA-Seq data with
444        or without a reference genome. BMC Bioinformatics 12, 323.

445  Lopes-Marques, M., Cunha, I., Reis-Henriques, M.A., Santos, M.M. and Castro, L.F.C., 2013.
446        Diversity and history of the long-chain acyl-CoA synthetase (Acsl) gene family in
447        vertebrates. BMC Evolutionary Biology 13, 271.

448  Min, K.-T. and Benzer, S., 1999. Preventing Neurodegeneration in the *Drosophila* Mutant
449        *bubblegum*. Science 284, 1985-1988.

450  Moriya-Sato, A., Hida, A., Inagawa-Ogashiwa, M., Wada, M.R., Sugiyama, K., Shimizu, J., Yabuki,
451        T., Seyama, Y. and Hashimoto, N., 2000. Novel Acyl-CoA Synthetase in
452        Adrenoleukodystrophy Target Tissues. Biochemical and Biophysical Research
453        Communications 279, 62-68.

454  Moser, H.W., Smith, K.D., Watkins, P.A., Powers, J. and Moser, A.B., 2002. X-linked
455        adrenoleukodystrophy. In The metabolic and molecular bases of inherited disease.
456        McGraw Hill.

457  Ohno, S., 1970. Evolution by Gene Duplication, 1 ed. Springer Berlin Heidelberg, New York.

458  Pavlidis, P. and Noble, W.S., 2003. Matrix2png: a utility for visualizing matrix data.
459        Bioinformatics 19, 295-6.

460  Pei, Z., Jia, Z. and Watkins, P.A., 2006. The second member of the human and murine
461        bubblegum family is a testis- and brainstem-specific acyl-CoA synthetase. J Biol Chem
462        281, 6632-41.

463  Pei, Z., Oey, N.A., Zuidervaart, M.M., Jia, Z., Li, Y., Steinberg, S.J., Smith, K.D. and Watkins, P.A.,
464        2003. The acyl-CoA synthetase "bubblegum" (lipidosin): further characterization and
465        role in neuronal fatty acid beta-oxidation. J Biol Chem 278, 47070-8.

466  Putnam, N.H., Butts, T., Ferrier, D.E.K., Furlong, R.F., Hellsten, U., Kawashima, T., Robinson-
467        Rechavi, M., Shoguchi, E., Terry, A., Yu, J.-K., Benito-Gutierrez, E., Dubchak, I., Garcia-

468  Fernandez, J., Gibson-Brown, J.J., Grigoriev, I.V., Horton, A.C., de Jong, P.J., Jurka, J.,
469  Kapitonov, V.V., Kohara, Y., Kuroki, Y., Lindquist, E., Lucas, S., Osoegawa, K., Pennacchio,
470  L.A., Salamov, A.A., Satou, Y., Sauka-Spengler, T., Schmutz, J., Shin-I, T., Toyoda, A.,
471  Bronner-Fraser, M., Fujiyama, A., Holland, L.Z., Holland, P.W.H., Satoh, N. and Rokhsar,
472  D.S., 2008. The amphioxus genome and the evolution of the chordate karyotype.
473  Nature 453, 1064-1071.

474  Robinson, L. and Mazurak, V., 2013. N-3 Polyunsaturated Fatty Acids: Relationship to
475  Inflammation in Healthy Adults and Adults Exhibiting Features of Metabolic Syndrome.
476  Lipids 48, 319-332.

477  Smith, J.J. and Keinath, M.C., 2015. The sea lamprey meiotic map improves resolution of
478  ancient vertebrate genome duplications. Genome Research 25, 1081-1090.

479  Soupene, E. and Kuypers, F.A., 2008. Mammalian long-chain acyl-CoA synthetases. Exp Biol
480  Med (Maywood) 233, 507-21.

481  Steinberg, S.J., Morgenthaler, J., Heinzer, A.K., Smith, K.D. and Watkins, P.A., 2000. Very Long-
482  chain Acyl-CoA Synthetases: Human "Bubblegum" represents a new family of proteins
483  capable of activating Very Long chain Fatty Acids. Journal of Biological Chemistry 275,
484  35162-35169.

485  Tang, P.-Z., Tsai-Morris, C.-H. and Dufau, M.L., 2001. Cloning and characterization of a
486  hormonally regulated rat long chain acyl-CoA synthetase. Proceedings of the National
487  Academy of Sciences 98, 6581-6586.

488  Tocher, D.R., 2003. Metabolism and Functions of Lipids and Fatty Acids in Teleost Fish. Reviews
489  in Fisheries Science 11, 107-184.

490  Wall, R., Ross, R.P., Fitzgerald, G.F. and Stanton, C., 2010. Fatty acids from fish: the anti-
491  inflammatory potential of long-chain omega-3 fatty acids. Nutrition Reviews 68, 280-
492  289.

493  Watkins, P.A., Maiguel, D., Jia, Z. and Pevsner, J., 2007. Evidence for 26 distinct acyl-coenzyme A
494  synthetase genes in the human genome. Journal of Lipid Research 48, 2736-2750.

495  Weimar, J.D., DiRusso, C.C., Delio, R. and Black, P.N., 2002. Functional Role of Fatty Acyl-
496  Coenzyme A Synthetase in the Transmembrane Movement and Activation of
497  Exogenous Long-chain Fatty Acids: Amino Acid residues within the ATP/AMP signature
498  motif of Escherichia coli FadD are required for enzyme activity and Fatty acid transport.
499  Journal of Biological Chemistry 277, 29369-29376.

500  Yates, A., Akanni, W., Amode, M.R., Barrell, D., Billis, K., Carvalho-Silva, D., Cummins, C.,
501  Clapham, P., Fitzgerald, S., Gil, L., Girón, C.G., Gordon, L., Hourlier, T., Hunt, S.E.,
502  Janacek, S.H., Johnson, N., Juettemann, T., Keenan, S., Lavidas, I., Martin, F.J., Maurel,
503  T., McLaren, W., Murphy, D.N., Nag, R., Nuhn, M., Parker, A., Patricio, M., Pignatelli, M.,
504  Rahtz, M., Riat, H.S., Sheppard, D., Taylor, K., Thormann, A., Vullo, A., Wilder, S.P.,
505  Zadissa, A., Birney, E., Harrow, J., Muffato, M., Perry, E., Ruffier, M., Spudich, G.,

506  Trevanion, S.J., Cunningham, F., Aken, B.L., Zerbino, D.R. and Flicek, P., 2016. Ensembl
507      2016. Nucleic Acids Research 44, D710-D716.

508  Zhang, J., 2003. Evolution by gene duplication: an update. Trends in Ecology & Evolution 18,
509      292-298.

510  Zheng, Y., Zhou, Z.M., Min, X., Li, J.M. and Sha, J.H., 2005. Identification and characterization of
511      the BGR-like gene with a potential role in human testicular
512      development/spermatogenesis. Asian J Androl 7, 21-32.

513  Zou, Z., DiRusso, C.C., Ctrnacta, V. and Black, P.N., 2002. Fatty acid transport in Saccharomyces
514      cerevisiae. Directed mutagenesis of FAT1 distinguishes the biochemical activities
515      associated with Fat1p. J Biol Chem 277, 31062-71.

516

517

## Figure captions

519

**Figure 1:** Maximum likelihood phylogenetic analysis of ACSBG amino acid sequences rooted with the invertebrate clade. Numbers at nodes indicate posterior probabilities calculated using aBayes.

523

**Figure 2:** Comparative genomic maps of vertebrate *ACSBG1* (A) *ACSBG2* (B) and *ACSBG3* (C and D) gene *loci.* Paralogy analysis and invertebrate genomic maps of ACSBG (E and F)*.*

527

**Figure 3:** Sequence alignment and ACS Motif analysis. **A**- Sequence logo graphs of the consensus sequences of all ACSBG sequences recovered excluding mammal specific ACSBG2b acyl-coenzyme A sequences, totalizing 104 sequences. **B**- Sequence logo graphs of the consensus sequences of all mammalian specific ACSBG2b (17 sequences). Overall height of the stack reflects the degree of conservation, the height of each letter represents relative frequency of a given residue in a specific position. Black arrow highlights the highly conserved arginine residue (Pei et al., 2006) and corresponding position in ACSBG2b   **C**- Heatmap of the relative expression of *ACSBG1 ACSBG2* and *ACSBG3* obtained from RNA-seq analysis and visualized using  Matrix2png (Pavlidis and Noble, 2003). **D-** Tissue expression profile of *X. tropicalis* ACSBG genes.

538

**Figure 4:** Proposed evolutionary history of the *ACSBG* gene family in vertebrates. Yellow corresponds to *ACSBG1* blue *ACSBG2* and green *ACSBG3*. Grey full lined circles with question marks ind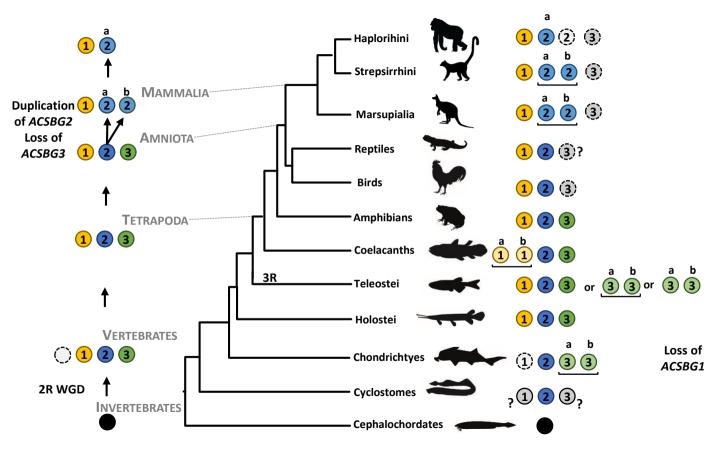icate unknown or unresolved if gene is present, grey circles dashed lined indicate gene loss. Black line under genes indicates tandem duplication.
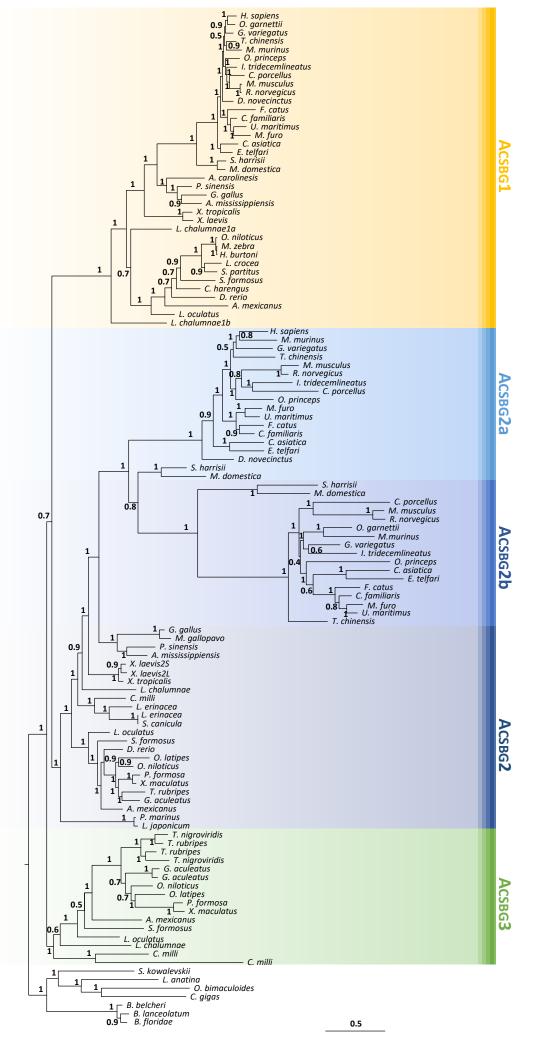
543

544

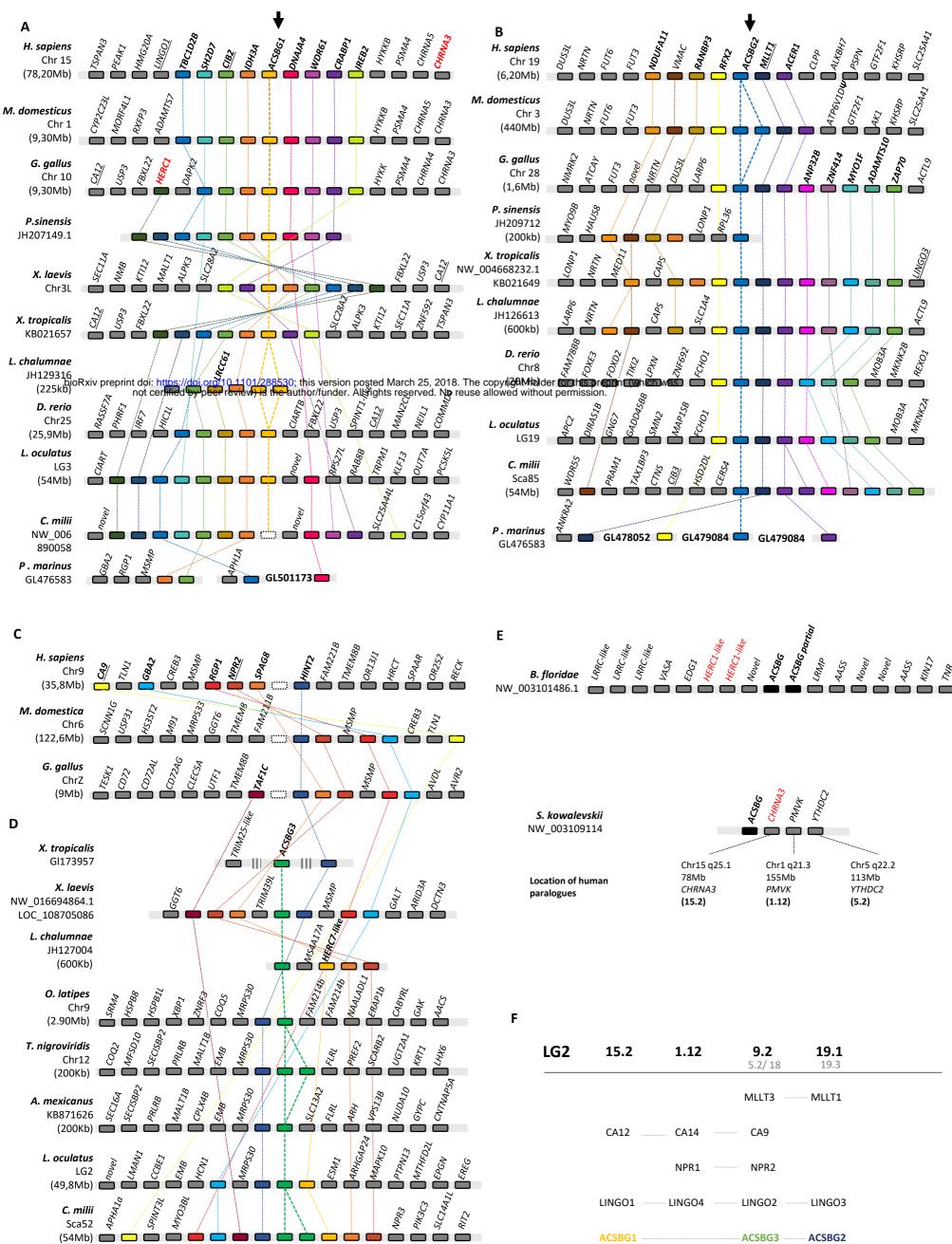**Supplementary material 1:** Table containing ACSBG sequences accession numbers.

20

546  **Supplementary material 2:** Accession numbers of the RNAseq files retrieved for
547  expression analysis.
548
549  **Supplementary material 3:** Genome and GTF files retrieved from Ensemble database
550  (Release 89) and Transcriptome files retrieved from NCBI used on this study and
551  accession numbers of reference genes.
552
553  **Supplementary material 4:** Complementary phylogenetic analysis of ACSBG sequences
554  including amphibian uncharacterized truncated ACSBG-like sequences indicated in red.
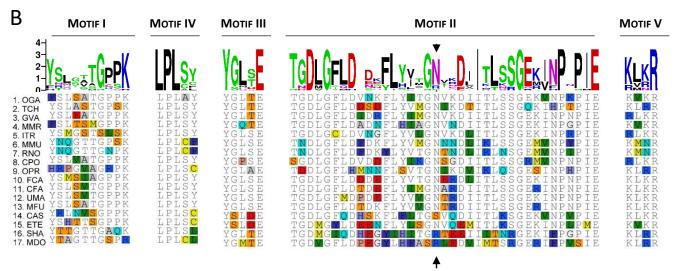555
556  **Supplementary material 5:** Motif sequence alignment of the full dataset used 121
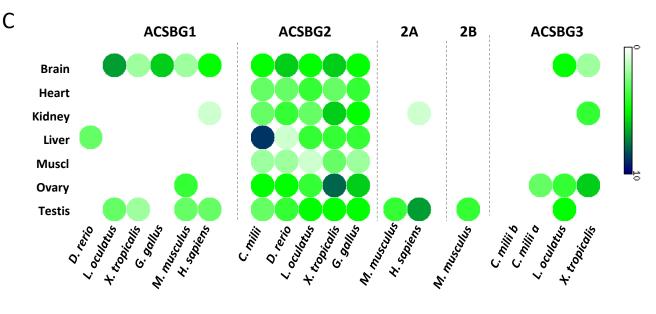557  ACSBG sequences. Red box highlights mammal specific ACSBG2b.
558

**A**

MOTIF I   MOTIF IV   MOTIF III   MOTIF II   MOTIF V

**B**

MOTIF I   MOTIF IV   MOTIF III   MOTIF II   MOTIF V

1. OGA
2. TCH
3. GVA
4. MMR
5. ITR
6. MMU
7. RNO
8. CPO
9. OPR
10. FCA
11. CFA
12. UMA
13. MFU
14. CAS
15. ETE
16. SHA
17. MDO

**C**

| | ACSBG1 | ACSBG2 | 2A | 2B | ACSBG3 |
|---|---|---|---|---|---|
| Brain | | | | | |
| Heart | | | | | |
| Kidney | | | | | |
| Liver | | | | | |
| Muscl | | | | | |
| Ovary | | | | | |
| Testis | | | | | |

*D. rerio   L. oculatus   X. tropicalis   G. gallus   M. musculus   H. sapiens*

*C. milii   D. rerio   L. oculatus   X. tropicalis   G. gallus*

*M. musculus   H. sapiens*

*M. musculus*

*C. milii b   C. milii a   L. oculatus   X. tropicalis*

**D**

MW  Brain  Skin  Heart  Liver  Spleen  Pancreas  Kidney  Intestine  Testis  Ovary  NTC

ACSBG1

ACSBG2

ACSBG3

β-ACTIN

EEF1A1