# The repertoire of vertebrate STAT transcription factors: Origin and variations in fish

Pierre Boudinot, Steve Bird, Louis Du Pasquier, Bertrand Collet

1        # The repertoire of vertebrate STAT transcription factors:

2        ## origin and variations in fish

3

4        Pierre Boudinot[1], Steve Bird[2], Louis Du Pasquier[3] and Bertrand Collet[1]

5

6        [1] Université Paris-Saclay, INRAE, UVSQ, VIM, 78350, Jouy-en-Josas, France.

7        [2] Biomedical Unit, School of Science, University of Waikato, Hamilton 3240, New Zealand

8        [3] Zoology and Evolutionary Biology, University of Basel, 4051 Basel, Switzerland

9
10
11
12
13
14
15
16
17
18
19
20

21    **Correspondence:**
22
23    Pierre Boudinot phone: 00331 34652585; email: Pierre.Boudinot@inrae.fr
24    Bertrand Collet phone: 00331 34652637; email: Bertrand.Collet@inrae.fr
25    Université Paris-Saclay, INRAE, UVSQ, VIM, 78350, Jouy-en-Josas, France
26
27

28    **Abbreviations:**
29    CCD: coiled coil domain; CBP: CREB-binding protein**;** GAS: Gamma interferon activation site;
30    ISRE: interferon-sensitive responsive element; SH2: Src homology 2 domain; STAT: signal
31    transducer and activator of transcription; TAD: C terminal transactivation domain; TAZ:
32    Transcription Adaptor putative Zinc finger; WGD: whole genome duplication;
33
34

37
38

39    **Abstract**

40    The *stat* gene family diversified during early vertebrate evolution thanks to two rounds of

41    whole genome duplication (WGD) to produce a typical repertoire composed of 6 STAT

42    factors (named 1-6). In contrast, only one or two *stat* genes have been reported in *C.*

43    *elegans* and in *D. melanogaster*. The main types of STAT found from bony fish to mammals

44    are present in Agnathan genomes, but a typical STAT1-6 repertoire is only observed in jawed

45    vertebrates. Comparative syntenies showed that STAT6 was the closest to the ancestor of

46    the family. An extensive survey of *stat* genes across fish including polyploid species showed

47    that whole genome duplications did not lead to a uniform expansion of *stat* genes. While 2

48    to 5 *stat*1 are present in salmonids, whose genome duplicated about 35My ago, only one

49    copy of *stat*2 and *stat*6 is retained. In contrast, common carp, with a recent whole genome

50    duplication (5-10My), possesses a doubled *stat* repertoire indicating that the elimination of

51    *stat*2 and *stat*6 additional copies is not immediate. Altogether our data shed light on the

52    multiplicity of evolutionary pathways followed by key components of the canonical cytokine

53    receptor signalling pathway, and point to differential selective constraints exerted on these

54    factors.

55

56    **1.  Introduction**

57         Animals have evolved a number of efficient strategies to combat a large diversity of

58    pathogens. In mammals, complex immune mechanisms are orchestrated and regulated by a

59    network of cytokines acting through cognate ligand/receptor on multiple specialised

60    immunocytes [1]. In mammals, a large number of these cytokines signal through the

61    JAK/STAT signalling factors [2] composed of a particular combination of four Janus kinases

62    (JAK1-3, Tyrosine Kinase TYK2) and one of 7 Signal Transducer and Activator of Transcription

63    (STAT1-4, 5A, 5B, 6) [3]. The pathway involves a cascade of phosphorylation reactions [4],

64    multimeric complex formation and nuclear translocation [5] resulting in the induction of a

65    particular set of genes responsible for a specific cellular response [6]. Gamma interferon

66    activation site (GAS) is the core genomic motif targeted by STAT1 homodimers [7-9].  STATs

67    heterodimers, associated with additional transcription factors, can bind variants of GAS

68    motifs such as interferon-sensitive responsive element (ISRE) resulting in transcriptional

69    regulation of large gene sets. Such gene sets leading to particular immune responses were

70    associated to different STAT-dependant signalling. In addition, variations in the epigenetic

71    status of genomic elements and in the type of immune cells involved explain, at least in

72    part, the "specificity paradox" of the JAK/STAT signalling pathway, namely, how a 7-member

73    protein family can ensure the specificity of response of dozens of cytokines [10]. The human

74    STAT repertoire is composed of 7 transcription factors encoded by genes located on 3

75    chromosomes: *STAT1* and *STAT4* closely linked on chromosome 2, *STAT2* and *STAT6* on

76    chromosome 12, and *STAT3*, *STAT5A* and *STAT5B* closely linked on chromosome 17.

77        All these proteins share four domains: a N-terminal Protein interaction domain

78    ("STAT-i"), a coiled coil domain ("STAT a", CCD), a DNA-binding domain ("STAT b", DBD) and

79    a Src homology 2 domain ("STAT-SH2"). Additionally, STAT1 and STAT2 comprise a C

80    terminal transactivation (TAD) domain: the STAT1 transactivation domain (IPR022752) binds

81    selectively to the Transcription Adaptor putative Zinc finger (TAZ)2 domain of C CREB-

82    binding protein (CBP)/p300, while the STAT2 transactivation domain (IPR022756) binds to

83    the TAZ1 domain of this protein [16-18] (Figure 1A). This domain confers to STAT1 and

84    STAT2 an additional capacity to regulate gene expression since CBP and P300 are histone

85    acetyltransferases that control acetylation of histones in nucleosomes, thus regulating

86    chromatin remodelling and gene transcription.

87

88       The vertebrate *stat* repertoire emerged from an ancestral sequence present in the

89 common ancestor of protostomes and deuterostomes with all STAT typical domains [11],

90 through WGD, tandem duplication and dispersion [12]. Liongue et al. [13], proposed that

91 vertebrate *stat* genes originated from a set of two paralogs produced by local duplication,

92 subsequently duplicated "en bloc" by the two rounds of WGD that occurred during early

93 vertebrate evolution, leading to four copies of this cluster. Three of these copies (*STAT3-*

94 *STAT5, STAT2-STAT6,* and *STAT1-STAT4*) have been retained in human and most vertebrates.

95 In zebrafish, additional copies of *stat1* and *stat5* were found, likely due to the additional,

96 teleost-specific WGD [13].

97       In this work, we revisited the origins and the evolutive dynamic of the vertebrate

98 *stat* gene repertoire. To find out whether duplicated *stat* copies were retained or lost, we

99 focused on groups and species in which additional WGD occurred. We thus focused on ray

100 finned fish because their genomes were subjected to several WGD events including a

101 teleost-specific WGD event ("3R") that occurred at the root of this lineage about 350 million

102 years ago (Myr) and more recent events for example in salmonids 50-60 Myr ago [14-15]

103 and carps 5-10 Myr ago. In addition, salmonid fish such as Atlantic salmon and Rainbow

104 trout are the most relevant species for the fish farming industry in Europe and worldwide,

105 and their genomes are among the best characterized in teleost fish. We also characterized

106 *stat* genes from Chondrichthyans, Agnathans and non-vertebrate deuterostomians to clarify

107 how these transcription factors evolved during the emergence of vertebrates.

108

109     **2. Results and discussion**

110

111    *The repertoire of stat genes is well conserved across tetrapods*

112        A fundamental repertoire of six *stat* genes is well-conserved across all tetrapod

113    classes and in the coelacanth, as illustrated in Figure 1B (see also Table S1). One-to-one

114    orthology relationships between tetrapod and coelacanth genes are also supported by

115    conserved synteny groups comprising several markers flanking all *stat* gene clusters (as

116    shown for *stat2* in Figure 1C).

117

118    *Loss and retention of stat genes after WGD during fish evolution reveal contrasted*

119    *constraints on different stat subtypes.*

120        In ray-finned fishes, the *stat* repertoire comprises the same types as in the

121    coelacanth and tetrapods, with *stat1, stat2, stat3, stat4, stat5* and *stat6* present in all

122    species across teleosts. After a WGD occurred some 350 Myr during the early evolution of

123    this group, two copies of each *stat* gene should have been generated [19-21]. This *stat*

124    repertoire has been reshaped by further duplication and gene loss.

125        In fish groups that did not undergo additional WGD, such as herring (*Clupea*

126    *harengus*), pike (*Esox lucius*), zebrafish (*Danio rerio*), stickleback (*Gasterosteus aculeatus*)

127    and the marine species fugu (*Fugu rubripes*) and sea bream (*Sparus aurata)* (Figure 2 and

128    Table S1), *stat*2-6 could be found as single copy in contrast with two or more, *stat1*

129    paralogs. One *stat1* paralog (named "b") is always linked to *stat4* as observed across

130    tetrapods, while the other copy (named "a") is located on another chromosome. This was

131    also the case of Atlantic cod (*Gadus morhua)*, a gadiform species with a particular immune

132    system lacking CD4 and a functional MHC class II pathway. In some cases, an additional

133    *stat1* can be found like in herring on a third chromosome (Figure 2). In zebrafish, a *stat1*

134    pseudogene has been described close to *stat1b*, [13] but is not present in the last genome

135    assembly. Only one copy of *stat3, 4* and *5* was generally present in these species with a few

136    exceptions as a double *stat5* in zebrafish, produced by a local duplication. In contrast, the

137    retention of two functional *stat1* genes across multiple families of ray-finned fish suggests

138    that different types of selection pressures may affect this gene, compared to other stat

139    family members.

140

141    *Multiple stat1 paralogs are also retained in tetraploid Salmonids*

142        To further test this hypothesis, we then focused on tetraploid species in which larger

143    *stat* repertoires have been produced by an additional WGD, providing the opportunity to

144    test their evolutionary fate.

145        We first performed a comprehensive survey of *stat* genes in salmonids, a fish family

146    tetraploidized by an additional WGD that occurred about 50-60 Myr ago. In these species,

147    we typically found two blocks *stat3+5*, two blocks *stat4+1,* four or five copies of *stat1,* but

148    only one *stat2* and one *stat6* gene (as for rainbow trout in Figure 2). A comprehensive

149    characterization of *stat* genes across salmonids is presented in table 1. Among the two

150    genera *Oncorhynchus* and *Salmo,* we analysed six species for which high quality genomes

151    were available: Sockeye salmon *O. nerka*, rainbow trout *O. mykiss*, chinook salmon *O.*

152    *tshawytscha*, Coho salmon *O. kisutch*, brown trout *S. trutta* and Atlantic salmon *S. salar*. A

153    total of 16 *stat* loci were found in these six salmonid genomes (Table 1, Figure 2). They were

154    located on 9 chromosomes corresponding to 6 chromosomes in zebrafish (Figure 3), a

155    diploid cyprinid. Linkage analyses showed that *stat1a1-3, stat1b-4* and *stat3-5b* duplicated

156    blocks generated by the salmonid-specific WGD were retained (Figure 3), while there was no

157    evidence of multiple copies of *stat2* and *stat6* (not even pseudogenes).

158        Phylogenetic and synteny block analyses across species provided consistent insights

159        into the origin of these *stat* genes (Figure 3 and Figure S1) and allowed unambiguous

160        identification and annotation. For example, all *stat1a* were linked to *ccr4not* and *ftcd* – as

161        zebrafish *stat1a* – while *stat1b* genes were associated to *stat4* and *slc40* genes. A number of

162        sequences encoding ORFs with size lower than 50% of the average size of STATs proteins

163        were additionally found in the rainbow trout, brown trout and Atlantic salmon (Table 2).

164        These, which likely are assembly artefacts or pseudogenes, were not included in the

165        phylogenetic analysis.

166        In salmonids, an additional (fifth) *stat1* gene that we named *stat1a3* was found

167        immediately downstream of *stat1a2,* suggesting it was generated by local duplication.

168        Interestingly, in the rainbow trout, brown trout and Coho salmon, the STAT1A3 protein is

169        twice the size of the normal size of the STAT1. These long STAT1 proteins contain twice the

170        typical set of domains in tandem [STATi- STATa- STATb- STATSH2-CTD- STATi- STATa- STATb-

171        STATSH2-CTD] and seem to be due to a local duplication-fusion of two *stat1* ORFs. The

172        double *stat1a3* was confirmed in the rainbow trout by the EST CA361350 covering the

173        junction area between the end of the putative first *stat1* and the beginning of the second,

174        which excludes that *stat1a3* has been produced by an assembly error. Further functional

175        studies are required to determine the function of the encoded protein, its potential

176        intramolecular dimerization and GAS elements binding abilities.

177        While the ancient WGD that in early teleost fish has left two *stat1* but only one of

178        the other *stat* paralogs in diploid species, the more recent salmonid-specific WGD resulted

179        in five *stat1* being retained. In contrast, only one copy of *stat2* and *stat6* were kept, either

180        due to an early complete loss post-WGD or because of consistent selection pressures in

181        favour of a single copy.

182

*Up-regulation of salmonid stat genes during antiviral responses*

Figure 4 shows the expression profile of all chinook salmon *stat* genes from an RNAseq experiment carried out on the EC cell line [22]. We checked whether the salmonid *stat1* and *stat2* genes were induced by type I IFN in a manner consistent with zebrafish where, in zebrafish larva, recombinant IFNϕ1 induces a robust up-regulation of *stat1b* and *stat2,* but not of *stat1a* [23]. In the chinook salmon cell line EC [22] *stat1b1* and *stat2* were induced with a FC>1.5 following stimulation by salmonid recombinant type I IFN. *Stat1a1* was also induced to some extent (Figure 4). However, *stat1b2* was not up-regulated. At steady state, s*tat1a* paralogs were more expressed than *stat1b,* as in zebrafish, a pattern consistent with a functional constitutive expression of *stat1a* genes [22].

Thus, there is no strict conservation of the *stat1/2* genes inducibility between salmonids and zebrafish, although the most upregulated genes are *stat1b* and *stat2* in both species. Overall, the paralogs of a given genes may be expressed at low levels in healthy cells, but can reach much higher levels after stimulation, offering opportunities for complex regulations. Whether this profile is different in other cells or under different stimulation conditions remains to be clarified. Similar variations of steady state expression levels were also observed for *stat5: stat5.1* and *stat5.2* were detected at low levels, while *stat5.3* transcripts were at least 10 times more abundant (Figure 4).

Duplicated genes in polyploid species are expected to be eliminated by deletion/accumulation of mutations, if they do not acquire new functions (neo/sub-functionalization) or are not kept by selection for gene dosage [24]. Our data about zebrafish and salmonid multiple *stat1* paralogs strongly suggest that they were indeed

205  subjected to neofunctionalization. More functional work will be necessary to establish if this

206  is also true for salmonid *stat3, 4* and *5* paralogs.

207

208  *Classification and nomenclature of stat genes in tetraploid species based on the example of*

209  *salmonids.*

210       The survey of the *stat* gene cluster in salmonid fish highlighted a nomenclature issue

211  for *stat* genes in polyploid species. The current annotation of such complex duplicated

212  genomes is often misleading because of assembly errors. Some annotations inherited the

213  nomenclature used at the time of the first and often single gene discovery by homology

214  cloning and lack consistency with annotation in other fish species. Regarding salmonid *stat*

215  genes, the rainbow trout *stat1a1* and *stat1a2* were annotated *stat1-1* and *stat1*-2 with no

216  reference to the *stat1a* group defined previously in non-salmonid teleost such as zebrafish.

217  The *stat1a3* was left annotated as "uncharacterized protein" whereas phylogeny and blast

218  against the mammalian protein database allocated it to the *stat1* group. Based on our

219  results from phylogeny and synteny conservation, we therefore established a coherent

220  nomenclature (Figure 3, table S1). A similar approach may be followed in other groups of

221  tetraploid vertebrates for example in Amphibians.

222

223  *Other tetraploid genomes tell more about stat evolutionary dynamics.*

224       We also studied the *stat* genes from the common carp (*Cyprinus carpio,* Ensembl

225  100: German_Mirror_carp_1.0), an allotetraploid teleost due to a recent WGD that occurred

226  relatively recently 5-10 Myr ago. In this species, all duplicated loci were retained, with

227  exactly twice as many genes as in the diploid cyprinid zebrafish with 4 *stat1*, 2 *stat2*, 2 *stat3*,

228  2 *stat4*, 4 *stat5* and 2 *stat6* (Table S1). The phylogenetic tree and the distribution of these

229    genes in contigs indicate that they correspond to a duplication of the blocks typically found

230    in zebrafish and other diploid teleosts (Table S1, Figure 4A).

231        Polyploid species also originate by allopolyploidization, *i.e.* by genome association

232    due to hybridization among different species. The availability of the genome for the frog

233    *Xenopus laevis* (2n=36) offers an opportunity to estimate the effect of evolution of the two

234    subgenomes of an allotetraploid species that were combined about 17-18 Myr ago, on the

235    diversification of the *stat* gene family [25]. In parallel, we analysed the *stat* repertoire from

236    the genome of *Xenopus tropicalis* (2n=20), which is not made of obvious pairs of

237    homoeologous chromosomes [26]. Thirteen (13) and seven (7) *stat* genes were identified in

238    the genome of *X. laevis* and *X. tropicalis*, respectively. *X. laevis* shows an almost perfect

239    duplication, with the exception of the loss of *stat4.S* (Table S1; Figure 4B), while 8.3% and

240    31.5% of *X. laevis* genes with clear 1:1 or 2:1 orthologs in *X. tropicalis* were lost,

241    respectively, from L and S subgenomes [25].

242        Interestingly, these observations in the common carp *Cyprinus carpio* and the African

243    clawed frog *Xenopus laevis* show that additional *stat* genes in polyploid species are not

244    rapidly eliminated, maybe because different copies can get specialized functions easily and

245    quickly. The presence of two functional *stat2* and *sta6* genes is tolerated in both cases, and

246    the loss of one copy is not necessarily immediate after duplication. Furthermore, the

247    pattern of evolution of duplicated genomes in salmonids suggest some selection pressures

248    possibly associated to viral subversions strategies [27].

249

250    *Agnathan-specific stat genes shed light on the origin of vertebrate STAT transcription factors*

251        The repertoire of *stat* genes is generally more diverse in Vertebrates than in other

252    Metazoans [13], likely due to the two cycles of WGD that occurred in the early evolution of

253    this lineage. To get insight into the early steps of *stat* evolution in vertebrates, we analysed

254    genomes from cartilaginous fish (*i.e.*, Chondrichthyans) and Agnathans.

255         Orthologs of all vertebrate *stat* were found in cartilaginous fish (Table 3, Figure 5A).

256    S*tat1, 2, 3, 5* and *6* have been annotated in most species of sharks and rays for which a

257    genome is available (Table S1). A typical *stat4* genomic sequence was not detected in shark

258    genomes except in the whale shark *Rhincodon typus* (Genbank ID XP_020376005). An EST

259    was also found in the dogfish shark *Squalus acanthias* (Genbank ID EE627912). Phylogenetic

260    analysis confirmed that *stat* genes from Chondrichthyans have human orthologs (Figure 5A).

261         In contrast, the list of *stat* genes was different in Agnathans: in two species, the sea

262    lamprey *Petromyzon marinus* and the hagfish *Eptatretus burgeri,* phylogenetic analysis

263    identified orthologs of human *STAT3* and *STAT5* (Figure 5A). Two other *stat* sequences

264    clustered with group1-4 (later referred as "*stat1-4*") and group5-6 (later referred as "*stat5-*

265    *6*") respectively but could not be assigned to a particular set, suggesting that the *stat*

266    repertoire of "modern" vertebrates was consolidated and standardized in Gnathostomes.

267    Additionally, the genomic neighborhood of agnathan *stat* did not fit the well-conserved

268    synteny blocks observed in jawed vertebrates (Figure 5B). These regions contain markers

269    located close to *stat* genes in vertebrates, such as in *ab1, gls, myo1b, cavin1, tmeff2,*

270    *slc39A10, dnah7.* However, these markers do not seem to be associated consistently with

271    *stat* sequences in agnathans and jawed vertebrates, suggesting that these regions were

272    produced by several duplications of an ancestral segment followed by extensive gene loss,

273    making the reconstitution of the history of this region difficult. Markers have been best

274    conserved in the regions encoding tetrapod *stat1, stat2, stat3 and stat4* and lamprey *stat1-*

275    *4* (Figure 5B). While two *stat* genes closely linked on lamprey scaffold 5 are most similar to

276    *stat3* and *stat5*/6 respectively, the markers found at close proximity do not match with

277     genes located close to human *stat3* and *5*: in human, *nab1*, *gls* and *myo1* homologs are

278     located on chromosome 2 close to *stat2* and *stat4*. Moreover, the lamprey *stat5* is linked to

279     *cavin*, a marker associated to human *stat*3/5, but also to *timeless* which is found close to

280     human *stat2* on chromosome 12.

281         Thus, all vertebrates seem to possess genes from both *stat1-4* and *stat5-6* groups,

282     encoded in genomic blocks inherited from an ancestral region containing *nab1, gls, myo1b,*

283     *cavin1, tmeff2, slc39A10,* and *dnah7* genes. However, the standardized *stat* repertoire

284     found in human was apparently established later in early gnathostomes. Further assemblies

285     of agnathan genomes will help to better understand the evolution of this region.

286

287     *Conserved linkages indicate that stat6 is a genomic environment closest to the ancestral stat*

288     *gene.*

289         We then analysed genomes from other deuterostomians. In these species, the

290     repertoire of *stat* genes was significantly smaller compared to vertebrates (Table 4, Figure

291     S2): three *stat* sequences were found in the cephalochordate lancelet *Branchiostoma*

292     *floridae* and in the tunicate *Ciona intestinalis,* and one in the appendicularia *Oikopleura*

293     *dioica,* in the hemichordate *Saccoglossus kowalevskii* and in the sea urchin

294     *Strongylocentrotus purpuratus.* Most sequences clustered in phylogenetic trees with STAT5

295     and STAT6 (data not shown), as reported previously for non-vertebrate STAT sequences

296     [13]. Only one sequence from the lancelet was more similar to the STAT1-4 group (Figure

297     S2). STAT5 and STAT6 have pleiotropic roles and are involved as transcription factors in the

298     biology of different cell types including epithelial and haematopoietic as well as immune

299     cells. Such critical functions in a wide range of contexts are consistent with a primordial

300    status of these genes within the family.   Overall, these results confirmed that a complete

301    STAT1-6 repertoire could not be found in these species.

302        As published previously [11], protostomians genomes also contain typical *stat* genes,

303    sometimes with multiple copies such as in the annelids *Helobdella* and *Capitella* (Table S1).

304    These sequences were most similar to vertebrate *stat5* and *stat6* as previously reported [6].

305    However, we were not able to find any *stat* synteny blocks shared between these species

306    and vertebrates. In contrast, linkage groups with *stat* genes from the placozoan *Trichoplax*

307    *adhaerens* stand out as an intriguing exception (Figure 6), which reminds of our previous

308    report about MHC [28]. In this species, *stat* genes were mainly located close to each other

309    on scaffold 2. Seven genes flanking this cluster were homologous to 7 markers located on

310    human chromosome 12, most of them in the close neighborhood of *STAT6* and *STAT2*, to 3

311    markers on human chromosome 2 close to *STAT1*, and to one marker on human

312    chromosome 17 close to *STAT5* and *STAT3*.  Interestingly, the best conserved set of linkages

313    involved the region of *stat6*, which appears to be most closely related to the ancestral *stat*

314    in phylogenetic analyses. These observations are consistent with the idea that both

315    vertebrate and *Trichoplax* genomes evolved relatively slowly while those of Protostomes

316    were subjected to extensive rearrangements. It also establishes a link between vertebrate

317    *stat* genes and basal bilaterians.

318    **Conclusions**

319        The canonical signalling pathway "cytokine receptor - JAK/STAT" contributes to many

320    functions in invertebrates, as illustrated by in depth studies in Drosophila. In this species,

321    this axis in involved in embryonic segmentation, in stem cell proliferation, in growth as well

322    as in immunity [29-36].. This repertoire of types of STAT transcription factors was

323    remarkably stable during tetrapod evolution. We found the same types of STAT in

324  Chondrichthyans but not in Agnathans, showing that this repertoire was likely standardized

325  with the emergence of Gnathostomes. In ray finned fish, successive WGD offered multiple

326  opportunities of further functional diversification and specialization. Our work shows that

327  only *stat1* paralogs were retained after the R3 WGD, with one being constitutive and the

328  other strongly induced by IFN. Focusing on Salmonids, we found several *stat-1*, *-3*, *-4* and *-5*

329  due to the most recent WGD, while one copy of the *stat2*/*6* block has been retained. With 5

330  paralogs and a remarkable long version with additional domains, *stat1* stands out as the

331  only member of the family prone to expansion and diversification. We have already

332  reported that chinook salmon cells in which *stat1a1* and *stat1a2* (with constitutive

333  expression) have been disrupted, completely lost type I IFN responsiveness [22]. Further

334  work is needed to dissect the specialized functions of these multiple *stat1* in various cell

335  types and infectious contexts. This evolutionary trend seems to be supported by the high

336  number of *stat1* genes in cyprinids which have been subjected to an independent WGD.

337  Overall, our work shows that the kinetics of *stat* loss is consistently variable across the

338  members of the family (Figure 7). Hence, the *stat* gene family is particularly suited to study

339  the fate of recently duplicated genes and in particular, loss-of function (or

340  pseudogeneization), dosage effect and neofunctionalization aspects [40]. Contrasted

341  inducibility of *stat* paralogs, which is a key mechanism of *stat* mediated immune responses,

342  provide a fast and efficient pathway towards neo/sub-functionalization for these critical

343  factors.

344

345  **Material and methods**

346  *Identification of stat sequences*

347     Genomes analyses were carried out using the Ensembl (Release 100) and NCBI web

348     interfaces. tBlastn and delta blast searches on the Refseq genomes, and genome

349     annotations searches were combined to pull out all the members of the *stat* gene family.

350     The NCBI genomes released version are as follows: Omyk_1.0 for *Oncorhynchus mykiss*,

351     Okis_V2 for *O. kisucht*, Otsh_v1.0 for *O. tshawytscha*, Oner_1.0 for *O. nerka*, fSalTru1.1 for

352     *Salmo trutta*, ICSASG_v2 for *S. salar*, fSpaAur1.1 for *Sparus aurata*, gadMor3.0 for *Gadus*

353     *morhua*, GRCz11 for *Danio rerio*, UCB_Xtro_10.0 for *Xenopus tropicalis*, Xenopus_laevis_v2

354     for *X. laevis*, GRCg6a for *Gallus gallus* and GRCh38.p13 for *Homo sapiens*. The domain

355     structure of the proteins encoded by *stat* genes was checked using SMART and pfam to look

356     for assembly problems and fragmentary sequences. MegaX software was used to carry out

357     phylogenetic analyses and confirm the homology relationships between sequences. The

358     evolutionary history was inferred by using the Maximum Likelihood method and JTT matrix-

359     based model. The bootstrap consensus tree inferred from 1000 replicates was taken to

360     represent the evolutionary history of the taxa analysed. Initial tree(s) for the heuristic

361     search were obtained automatically by applying Neighbor-Joining and BioNJ algorithms to a

362     matrix of pairwise distances estimated using a JTT model, and then selecting the topology

363     with superior log likelihood value.

364

365     *Microsynteny analysis*

366     Synteny were retrieved from Genomicus (version 100) and the orthology/paralogy

367     relationships available in Ensembl, and complemented by visual examination of the

368     graphical interface in both Ensembl and NCBI. A linkage was considered a conserved

369     microsynteny only when three or more such genes were linked in such a way in two species.

370

371     *Expression analysis*

372         RNAseq transcriptome analysis on the chinook salmon STAT2-KO GS2 and CHSE-

373     EC cell lines was described by [22]. Briefly, STAT2 KO or control CHSE-EC cells were

374     stimulated (or not) in EMEM medium supplemented with 250 ng/ml of recombinant *O.*

375     *mykiss* IFNA2.  Three biological replicates (Flask 1-3) were used for library construction for

376     each group, and RNA-Seq libraries were prepared using TruSeq Stranded mRNA Sample

377     Preparation Kit (Illumina) according to the manufacturer's instructions. Libraries were

378     validated for quality on Agilent DNA1000 Kit, pooled in equimolar amounts and sequenced

379     in pair-ends 2x75 bp on Illumina NextSeq 500/550. For each library, a depth 20 M reads

380     were generated. Reads were then spliced-aligned to 47,898 genes (47,022 Gnomon, 876

381     RefSeq, GCF_002163495.1_Omyk_1.0_genomic.gff from the NCBI).

382

388

389     **Data availability**

390     All data generated or analysed during this study are included in this published article and its

391     supplementary information files

392     **Conflict of interests disclosure**

393     The authors declare no commercial or financial conflict of interest.

394

395

396 **References**

397 1. Aaronson, D.S., Horvath, C.M., 2002. A road map for those who don't know JAK-STAT.

398 Science (80-. ). https://doi.org/10.1126/science.1071545

399 2. Kiu, H., Nicholson, S.E., 2012. Biology and significance of the JAK/STAT signalling

400 pathways. Growth Factors 30, 88–106.

401 https://doi.org/10.3109/08977194.2012.660936

402 3. Villarino, A. V, Kanno, Y., O'Shea, J.J., 2017. Mechanisms and consequences of Jak–STAT

403 signaling in the immune system. Nat. Immunol. 18, 374–384.

404 https://doi.org/10.1038/ni.3691

405 4. Decker, T., Kovarik, P., 2000. Serine phosphorylation of STATs. Oncogene 19, 2628–2637.

406 https://doi.org/10.1038/sj.onc.1203481

407 5. Reich, N.C., 2013. STATs get their move on. JAK-STAT 2, e27080.

408 https://doi.org/10.4161/jkst.27080

409 6. Liongue, C., Sertori, R., Ward, A.C., 2016. Evolution of Cytokine Receptor Signaling. J.

410 Immunol. 197, 11–18. https://doi.org/10.4049/jimmunol.1600372

411 7. Nast, R., Staab, J., Meyer, T., 2019. Gene Activation by the Cytokine-Driven Transcription

412 Factor STAT1, in: Gene Regulation. IntechOpen.

413 https://doi.org/10.5772/intechopen.82699

414 8. Stark, G.R., Darnell, J.E., 2012. The JAK-STAT Pathway at Twenty. Immunity.

415 https://doi.org/10.1016/j.immuni.2012.03.013

416 9. Decker, T., Kovarik, P., Meinke, A., 1997. GAS elements: A few nucleotides with a major

417 impact on cytokine-induced gene expression. J. Interf. Cytokine Res. 17, 121–134.

418 https://doi.org/10.1089/jir.1997.17.121

419     10. Lin, J.-X., Leonard, W.J., 2019. Fine-Tuning Cytokine Signals. Annu. Rev. Immunol. 37,

420          295–324. https://doi.org/10.1146/annurev-immunol-042718-041447

421     11. Wang, Y., Levy, D.E., 2012. Comparative evolutionary genomics of the STAT family of

422          transcription factors. JAK-STAT 1, 23–36. https://doi.org/10.4161/jkst.19418

423     12. Copeland, N.G., Gilbert, D.J., Schindler, C., Zhong, Z., Wen, Z., Darnell, J.E., Mui, A.L.-F.,

424          Miyajima, A., Quelle, F.W., Ihle, J.N., Jenkins, N.A., 1995. Distribution of the

425          mammalian stat gene family in mouse chromosomes. Genomics 29, 225–228.

426          https://doi.org/10.1006/geno.1995.1235

427     13. Liongue, C., O'Sullivan, L.A., Trengove, M.C., Ward, A.C., 2012. Evolution of JAK-STAT

428          pathway components: Mechanisms and role in immune system development. PLoS

429          One 7, e32777. https://doi.org/10.1371/journal.pone.0032777

430     14. Lien, S., Koop, B.F., Sandve, S.R., Miller, J.R., Kent, M.P., Nome, T., Hvidsten, T.R., Leong,

431          J.S., Minkley, D.R., Zimin, A., Grammes, F., Grove, H., Gjuvsland, A., Walenz, B.,

432          Hermansen, R.A., Von Schalburg, K., Rondeau, E.B., Di Genova, A., Samy, J.K.A., Olav

433          Vik, J., Vigeland, M.D., Caler, L., Grimholt, U., Jentoft, S., Inge Våge, D., De Jong, P.,

434          Moen, T., Baranski, M., Palti, Y., Smith, D.R., Yorke, J.A., Nederbragt, A.J., Tooming-

435          Klunderud, A., Jakobsen, K.S., Jiang, X., Fan, D., Hu, Y., Liberles, D.A., Vidal, R., Iturra,

436          P., Jones, S.J.M., Jonassen, I., Maass, A., Omholt, S.W., Davidson, W.S., 2016. The

437          Atlantic salmon genome provides insights into rediploidization. Nature 533, 200–205.

438          https://doi.org/10.1038/nature17164

439     15. Pasquier, J., Cabau, C., Nguyen, T., Jouanno, E., Severac, D., Braasch, I., Journot, L.,

440          Pontarotti, P., Klopp, C., Postlethwait, J.H., Guiguen, Y., Bobe, J., 2016. Gene

441          evolution and gene expression after whole genome duplication in fish: the PhyloFish

442          database. BMC Genomics 17, 368. https://doi.org/10.1186/s12864-016-2709-z

443    16. Wojciak, J.M., Martinez-Yamout, M.A., Dyson, H.J., Wright, P.E., 2009. Structural basis

444          for recruitment of CBP/p300 coactivators by STAT1 and STAT2 transactivation

445          domains. EMBO J. 28, 948–958. https://doi.org/10.1038/emboj.2009.30

446    17. Bhattacharya, S., Eckner, R., Grossman, S., Oldread, E., Arany, Z., D'Andrea, A.,

447          Livingston, D.M., 1996. Cooperation of Stat2 and p300/CBP in signalling induced by

448          interferon- α. Nature 383, 344–347. https://doi.org/10.1038/383344a0

449    18. Zhang, J.J., Vinkemeier, U., Gu, W., Chakravarti, D., Horvath, C.M., Darnell, J.E., 1996.

450          Two contact regions between Stat1 and CBP/p300 in interferon γ signaling. Proc.

451          Natl. Acad. Sci. U. S. A. 93, 15092–15096. https://doi.org/10.1073/pnas.93.26.15092

452    19. Van de Peer, Y., 2004. Tetraodon genome confirms Takifugu findings: Most fish are

453          ancient polyploids. Genome Biol. https://doi.org/10.1186/gb-2004-5-12-250

454    20. Jaillon, O., Aury, J.M., Brunet, F., Petit, J.L., Stange-Thomann, N., Maucell, E., Bouneau,

455          L., Fischer, C., Ozouf-Costaz, C., Bernot, A., Nicaud, S., Jaffe, D., Fisher, S., Lutfalla, G.,

456          Dossat, C., Segurens, B., Dasilva, C., Salanoubat, M., Levy, M., Houdet, N., Castellano,

457          S., Anthouard, V., Jubin, C., Castelli, V., Katinka, M., Vacherie, B., Blémont, C., Skalli,

458          Z., Cattolico, L., Poulain, J., De Berardinis, V., Cruaud, C., Dupart, S., Brottler, P.,

459          Coutanceau, J.P., Gouzy, J., Parra, G., Lardier, G., Chapple, C., McKernan, K.J.,

460          McEwan, P., Bosak, S., Kellis, M., Volff, J.N., Gulgó, R., Zody, M.C., Mesirov, J.,

461          Lindblad-Toh, K., Birren, B., Nusbaum, C., Kahn, D., Robinson-Rechavi, M., Laudet, V.,

462          Schachter, V., Quétler, F., Saurin, W., Scarpeill, C., Wincker, P., Lander, E.S.,

463          Weissenbach, J., Roest Crollius, H., 2004. Genome duplication in the teleost fish

464          Tetraodon nigroviridis reveals the early vertebrate proto-karyotype. Nature 431,

465          946–957. https://doi.org/10.1038/nature03025

466    21. Christoffels, A., Brenner, S., Venkatesh, B., 2006. Tetraodon genome analysis provides

467        further evidence for whole-genome duplication in the ray-finned fish lineage. Comp.

468        Biochem.    Physiol.    -    Part    D    Genomics    Proteomics    1,    13–19.

469        https://doi.org/10.1016/j.cbd.2005.06.001

470    22. Dehler, C.E., Lester, K., Della Pelle, G., Jouneau, L., Houel, A., Collins, C., Dovgan, T.,

471        Machat, R., Zou, J., Boudinot, P., Martin, S.A.M., Collet, B., 2019. Viral Resistance and

472        IFN    Signaling    in    STAT2    Knockout    Fish    Cells.    J.    Immunol.    203,    465–475.

473        https://doi.org/10.4049/jimmunol.1801376

474    23. Levraud, J.-P., Jouneau, L., Briolat, V., Laghi, V., Boudinot, P., 2019. IFN-Stimulated Genes

475        in Zebrafish and Humans Define an Ancient Arsenal of Antiviral Immunity. J.

476        Immunol. 203, 3361–3373. https://doi.org/10.4049/jimmunol.1900804

477    24. Levasseur, A., Pontarotti, P., 2011. The role of duplications in the evolution of genomes

478        highlights the need for evolutionary-based approaches in comparative genomics.

479        Biol. Direct. https://doi.org/10.1186/1745-6150-6-11

480    25. Session, A.M., Uno, Y., Kwon, T., Chapman, J.A., Toyoda, A., Takahashi, S., Fukui, A.,

481        Hikosaka, A., Suzuki, A., Kondo, M., Van Heeringen, S.J., Quigley, I., Heinz, S., Ogino,

482        H., Ochi, H., Hellsten, U., Lyons, J.B., Simakov, O., Putnam, N., Stites, J., Kuroki, Y.,

483        Tanaka, T., Michiue, T., Watanabe, M., Bogdanovic, O., Lister, R., Georgiou, G.,

484        Paranjpe, S.S., Van Kruijsbergen, I., Shu, S., Carlson, J., Kinoshita, T., Ohta, Y.,

485        Mawaribuchi, S., Jenkins, J., Grimwood, J., Schmutz, J., Mitros, T., Mozaffari, S. V.,

486        Suzuki, Y., Haramoto, Y., Yamamoto, T.S., Takagi, C., Heald, R., Miller, K.,

487        Haudenschild, C., Kitzman, J., Nakayama, T., Izutsu, Y., Robert, J., Fortriede, J., Burns,

488        K., Lotay, V., Karimi, K., Yasuoka, Y., Dichmann, D.S., Flajnik, M.F., Houston, D.W.,

489        Shendure, J., Dupasquier, L., Vize, P.D., Zorn, A.M., Ito, M., Marcotte, E.M.,

490    Wallingford, J.B., Ito, Y., Asashima, M., Ueno, N., Matsuda, Y., Veenstra, G.J.C.,

491    Fujiyama, A., Harland, R.M., Taira, M., Rokhsar, D.S., 2016. Genome evolution in the

492    allotetraploid frog Xenopus laevis. Nature 538, 336–343.

493    https://doi.org/10.1038/nature19840

494  26. Uno, Y., Nishida, C., Takagi, C., Ueno, N., Matsuda, Y., 2013. Homoeologous

495    chromosomes of Xenopus laevis are highly conserved after whole-genome

496    duplication. Heredity (Edinb). 111, 430–436. https://doi.org/10.1038/hdy.2013.65

497  27. Nan, Y., Wu, C., Zhang, Y.J., 2017. Interplay between Janus kinase/signal transducer and

498    activator of transcription signaling activated by type I interferons and viral

499    antagonism. Front. Immunol. https://doi.org/10.3389/fimmu.2017.01758

500  28. Suurväli, J., Jouneau, L., Thépot, D., Grusea, S., Pontarotti, P., Du Pasquier, L., Rüütel

501    Boudinot, S., Boudinot, P., 2014. The Proto-MHC of Placozoans, a Region Specialized

502    in Cellular Stress and Ubiquitination/Proteasome Pathways. J. Immunol. 193, 2891–

503    2901. https://doi.org/10.4049/jimmunol.1401177

504  29. Perrimon, N., Mahowald, A.P., 1986. l(1)hopscotch, a larval-pupal zygotic lethal with a

505    specific maternal effect on segmentation in Drosophila. Dev. Biol. 118, 28–41.

506    https://doi.org/10.1016/0012-1606(86)90070-9

507  30. Zeidler, M.P., Perrimon, N., Strutt, D.I., 1999. Polarity determination in the Drosophila

508    eye: A novel role for unpaired and JAK/STAT signaling. Genes Dev. 13, 1342–1353.

509    https://doi.org/10.1101/gad.13.10.1342

510  31. Gregory, L., Came, P.J., Brown, S., 2008. Stem cell regulation by JAK/STAT signaling in

511    Drosophila. Semin. Cell Dev. Biol. https://doi.org/10.1016/j.semcdb.2008.06.003

512     32. Morin-Poulard, I., Vincent, A., Crozatier, M., 2013. The Drosophila JAK-STAT pathway in

513         blood cell formation and immunity. JAK-STAT 2, e25700.

514         https://doi.org/10.4161/jkst.25700

515     33. Rajan, A., Perrimon, N., 2012. Drosophila cytokine unpaired 2 regulates physiological

516         homeostasis by remotely controlling insulin secretion. Cell 151, 123–137.

517         https://doi.org/10.1016/j.cell.2012.08.019

518     34. Vanha-aho, L.M., Valanne, S., Rämet, M., 2016. Cytokines in Drosophila immunity.

519         Immunol. Lett. https://doi.org/10.1016/j.imlet.2015.12.005

520     35. Stepkowski, S.M., Chen, W., Ross, J.A., Nagy, Z.S., Kirken, R.A., 2008. STAT3: An

521         important regulator of multiple cytokine functions. Transplantation.

522         https://doi.org/10.1097/TP.0b013e3181739d25

523     36. Stepkowski, S.M., Chen, W., Ross, J.A., Nagy, Z.S., Kirken, R.A., 2008. STAT3: An

524         important regulator of multiple cytokine functions. Transplantation.

525         https://doi.org/10.1097/TP.0b013e3181739d25

526     37. Whelan, S., Goldman, N., 2001. A general empirical model of protein evolution derived

527         from multiple protein families using a maximum-likelihood approach. Mol. Biol. Evol.

528         18, 691–699. https://doi.org/10.1093/oxfordjournals.molbev.a003851

529     38. Felsenstein, J., 1985. Confidence limits on phylogenies: An approach using the bootstrap.

530         Evolution (N. Y). 39, 783–791. https://doi.org/10.1111/j.1558-5646.1985.tb00420.x

531     39. Kumar, S., Stecher, G., Li, M., Knyaz, C., Tamura, K., 2018. MEGA X: Molecular

532         evolutionary genetics analysis across computing platforms. Mol. Biol. Evol. 35, 1547–

533         1549. https://doi.org/10.1093/molbev/msy096

534    40. Lan, X., Pritchard, J.K., 2016. Coregulation of tandem duplicate genes slows evolution of

535         subfunctionalization    in    mammals.    Science    (80-.    ).    352,    1009–1013.

536         https://doi.org/10.1126/science.aad8411

537    **Abbreviations**

538    **STAT** Signal Transducers and Activators of Transcription **GMCSF** Granulocyte-Macrophage

539    Colony–Stimulating Factor **R1-4 WGD** Round 1-4 of Whole Genome Duplication **IFN**

540    Interferon **GAS** IFN-γ Activation Sequence **Myr** Million Years **CCD** Coiled coil domain **DBD**

541    DNA-binding Domain **SH2** Src Homology 2 **CTD** C-Terminal Domain **CREB** C-AMP Response

542    Element-Binding **CCRNOT** Carbon Catabolite Repression—Negative On TATA-less **FTCD**

543    Formimidoyl Transferase CycloDeaminase **EST** Expressed Sequence Tag

544    **Correspondence:**

545    B. Collet and P. Boudinot

546

547     **Figure legends**

548

549     **Figure 1**. **Evolutionary history of STAT transcription factors across tetrapods. A**. Domain

550     structure of vertebrate STAT proteins. **B**. Maximum likelihood phylogenetic tree of STAT

551     amino-acid sequences from human *Homo sapiens* (hosa), chicken *Gallus gallus* (gaga), *Anolis*

552     *carolinensis* anole (anco), clawed frog *Xenopus tropicalis* (xetr) and *Latimeria chalumnae*

553     coelacanth (lach). Bootstrap values (in %) of key nodes are indicated. Bootstrap values lower

554     than 60% are not indicated. All sequences and sequence ID are provided in Table S1. **C**.

555     Conservation of genomic neighborhood of *stat2* genes in the same tetrapod species (based

556     on genome assemblies from Ensembl release 100).

557

558     **Figure 2**. **Repertoires of fish STAT amino-acid sequences and *stat* genes chromosomic**

559     **distribution**. Data from genome assembly Omyk_1.0 (Oncorhynchus mykiss, rainbow trout,

560     RefSeq GCF_002163495.1) and for other species from Ensembl release 100.

561

562     **Figure 3. *stat* genes from genomes of Salmonidae (4 *Oncorhynchus* species and 2 *Salmo***

563     ***species*): synteny block conservation analysis**. Data based on NCBI genome assemblies:

564     Okis_V2, *Oncorhynchus kisutch* (coho salmon): GCF_002021735.2 ; Oner_1.0, *Oncorhynchus*

565     *nerka* (sockeye salmon): GCF_006149115.1; Otsh_v1.0, *Oncorhynchus tshawytscha* (Chinook

566     salmon):    GCF_002872995.1;    Oket_V1,    *Oncorhynchus    keta*    (chum    salmon):

567     GCF_012931545.1 ; Omyk_1.0, *Oncorhynchus mykiss* (rainbow trout): GCF_002163495.1;

568     fSalTru1.1, *Salmo trutta* (river trout): GCF_901001165.1 ; ICSASG_v2, *Salmo salar* (Atlantic

569     salmon): GCF_000233375.1.

570    **Figure 4**. **Expression levels of *stat* genes (basal and induced by recombinant type I**

571    **interferon) determined by RNAseq in CHSE-EC cell *O. tshawytscha* [22].** Data are on a log

572    scale and represent the average + Standard deviation (N = 3). When the induction is

573    statistically significant (*** $p<0.001$), the Fold Change is indicated

574

575    **Figure 5. STAT repertoires in other polyploid species**. **A.** Phylogenetic tree of STAT proteins

576    from common carp *Cyprinus carpio* (Cyca) and zebrafish *Danio rerio* (Dare). The evolutionary

577    history was inferred using the Maximum likelihood method (number of bootstrap tests

578    :1000 replicates). Bootstrap values (in %) of key nodes are indicated. Bootstrap values lower

579    than 60% are not shown. All ambiguous positions were removed for each sequence pair

580    (pairwise deletion option). The chromosome (for zebrafish) or the scaffold (for carp) are

581    indicated. **B.** Phylogenetic tree of STAT proteins from *Xenopus tropicalis* (Xetr), *Xenopus*

582    *laevis* (Xela) and *Gallus gallus* (gaga). The gene ID are indicated and refer to Table S1. The

583    evolutionary history was inferred using the Maximum likelihood method as for A.

584

585    **Figure 6**. **STAT amino-acid sequences from Chondrichthyans and Agnathans**. A. Maximum

586    likelihood phylogenetic tree of STAT amino-acid sequences from human, elephant shark (a

587    chimera), spiny dogfish and zebra bullhead sharks, the little skate (a ray) and sea lamprey

588    and inshore hagfish (Agnathans). Bootstrap values (in %) of key nodes are indicated.

589    Bootstrap values lower than 60% are not indicated. All sequences and sequence ID are

590    provided in Table S1. B. Genomic context of *stat* genes in human and sea lamprey based on

591    data from Ensembl release 100 (Human GRCh38.p13 and Sea Lamprey Pmarinus_7.0).

592

593　**Figure 7**. **Conserved genomic neighborhood between human STAT genes located on**

594　**chromosomes 2, 12 and 17, and *stat* genes found in *Trichoplax adhaerens***. The location of

595　markers is indicated besides gene names between brackets when relevant. Data from

596　genome assemblies in Ensembl release 100 (Human GRCh38.p13 and *Trichoplax adhaerens*

597　ASM15027v1).

598

599　**Figure 8. Evolutionary pathways of stat genes in Deuterostomians**. WGD are indicated by

600　red "X", and the date indicated for the most recent events. Salmonid *stat* genes for which

601　paralogs have not been retained are boxed in red.

602

603　**Figure S1. Phylogenetic tree of salmonids STAT amino-acid sequences.**

604　The human and zebrafish sequences were included in the analysis as reference and basis for

605　nomenclature. The evolutionary history was inferred by using the Maximum Likelihood

606　method and JTT matrix-based model. The bootstrap consensus tree inferred from 500

607　replicates is taken to represent the evolutionary history of the taxa analysed. This analysis

608　involved 95 amino acid sequences (hs: *Homo sapiens*; dr: *Danio rerio*, zebrafish; on:

609　*Oncorhynchus nerka*, Sockeye salmon; om: *Oncorhynchus mykiss*, rainbow trout, ot:

610　*Oncorhynchus tshawytscha*, chinook salmon; ok: *Oncorhynchus kisutch*, Coho salmon; st:

611　*Salmo trutta*, brown trout and ss: *Salmo salar*, Atlantic salmon). There were a total of 1742

612　positions in the final dataset. All sequences and sequence ID are provided in Table S1 and

613　are based on the NCBI accession number except for ss stat1a3, only annotated as a

614　translated transcript in ENSEMBL Evolutionary analyses were conducted in MEGA X [39].

615

616 **Figure S2. Phylogenetic tree of human and non-vertebrate deuterostomian STAT amino-**

617 **acid sequences.**

618 The evolutionary history was inferred by using the Maximum Likelihood method and JTT

619 matrix-based model. (Hosa: human; Brfl: lancelet, *Branchiostoma lanceolatum*; Ciin: *Ciona*

620 *intestinalis*; Sako: *Sacchoglossus kowalevsky*; Oidi: *Oikopleura dioica*; Stpu: sea urchin,

621 *Strongylocentrotus purpuratus*). All sequences and sequence ID are provided in Table S1.

622 Evolutionary analyses were conducted in MEGA X [39].
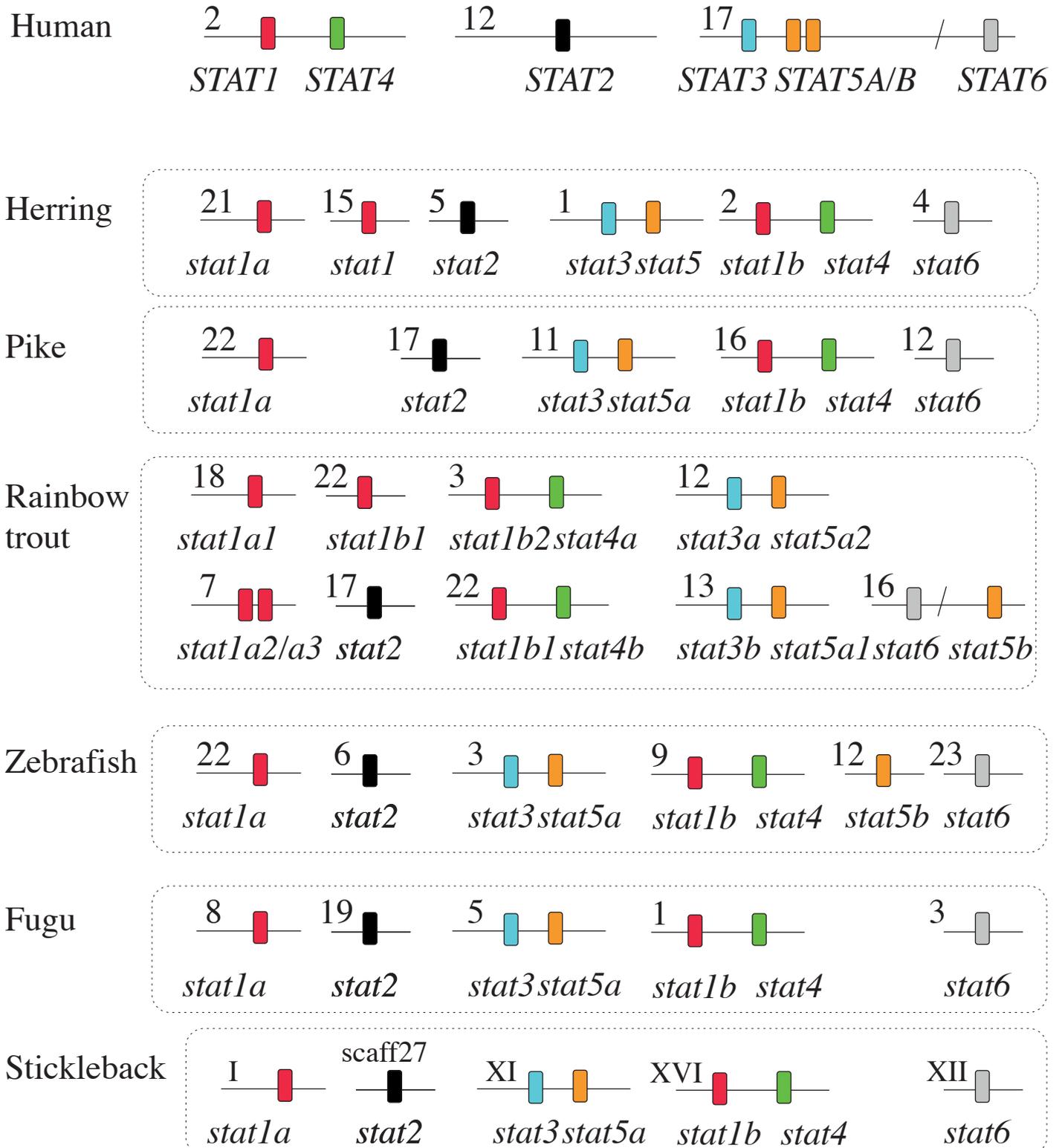
**Figure 1**

**Figure 2.**

**Figure 3**



Stat1a2 absent in ss, is a pseudogene in st
Stat1a3 has two copies in om, ok, sf
Stat5a1 and stat5a2 are pseudogenes in ot

**Figure 5**

# Figure 6

# Figure 7



*H. sapiens*

hChr 17 — *STAT5* *STAT3* *CAVIN1* *ATP6VOA1*

hChr 2 — *NEMP2* *NAB1* *STAT1* *MYO1B* *CAVIN2*

hChr 12 — *STAT2* *TIM1* *RBMS2* *RDH16* *MYO1A* *NEMP1* *STAT6* *NAB2* *SHMT2* (57.2) *ATP6VOA2* (123)
(56.3)

Scaff 2 —
*TRIADG53559 (5.8)*    TRIADG21854 (6.8)    *TRIADG20799(7.9)*
*TRIADG63669 (6.0)*    *TRIADG20417 (7.2)*
*TRIADG53910-12-14-18-20-47*
*TRIADG53906 (8.8)*
*TRIADG20682 (9.2)*

*T. adhaerens*

Scaff 1 —
*TRIADG52438 (9.8)*

# Figure 8

| | Stat1-4 | Stat5-6 | Stat1 | Stat2 | Stat3 | Stat4 | Stat5 | Stat6 |
|---|---|---|---|---|---|---|---|---|
| Sea urchin | | 1 | | | | | | |
| Accorn worm | | 2 | | | | | | |
| Lancelet | 1 | 1 | | | | | | |
| Sea squirt | | 3 | | | | | | |
| Oikopleura | | 1 | | | | | | |
| **Agnathans** | | | | | | | | |
| Lamprey | 1 | 1 | | | 1 | | 1 | |
| **Chondrichthyans** | | | | | | | | |
| Leucoraja erinacea | | | 1 | 1 | 1 | | 1 | 1 |
| Squalus acanthias | | | 1 | 1 | 1 | | 1 | 1 |
| Rhincodon typus | | | | | | 1 | | |
| **Bony fishes** | | | | | | | | |
| Zebrafish | | | 2 | 1 | 1 | 1 | 2 | 1 |
| Common carp | | | 4 | 2 | 2 | 2 | 4 | 2 |
| Atlantic Salmon | | | 5 | 1 | 2 | 2 | 3 | 1 |
| Rainbow trout | | | 5 | 1 | 2 | 2 | 3 | 1 |
| Chinook Salmon | | | 5 | 1 | 2 | 2 | 3 | 1 |
| **Amphibians** | | | | | | | | |
| Xenopus tropicalis | | | 1 | 1 | 2 | 1 | 1 | 1 |
| Xenopus laevis | | | 2 | 2 | 4 | 1 | 2 | 2 |
| **Other tetrapods** | | | 1 | 1 | 1 | 1 | 1 or 2 | 1 |

Tree annotations: X, X; X 5-10My; X; X 60My; X 17My

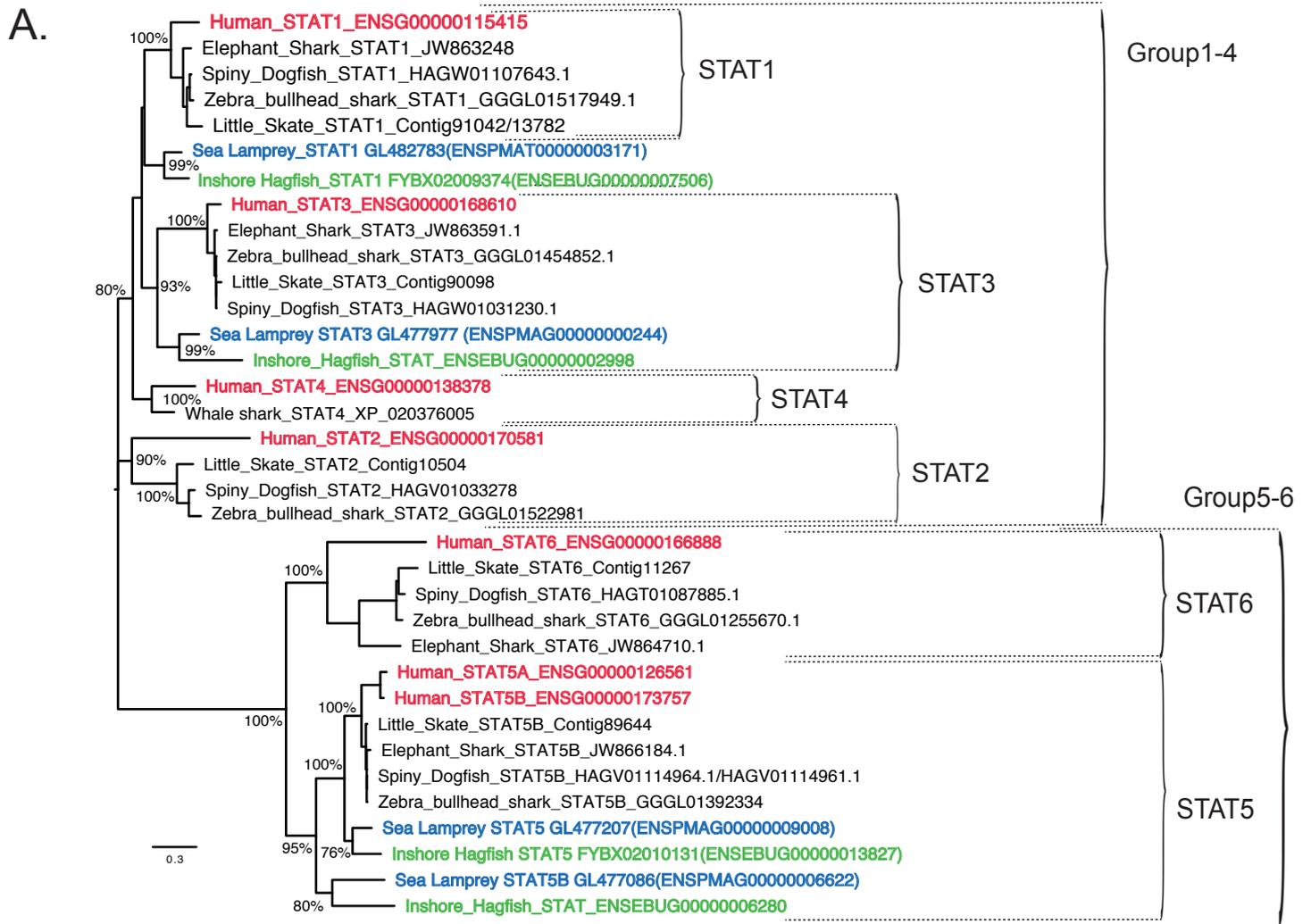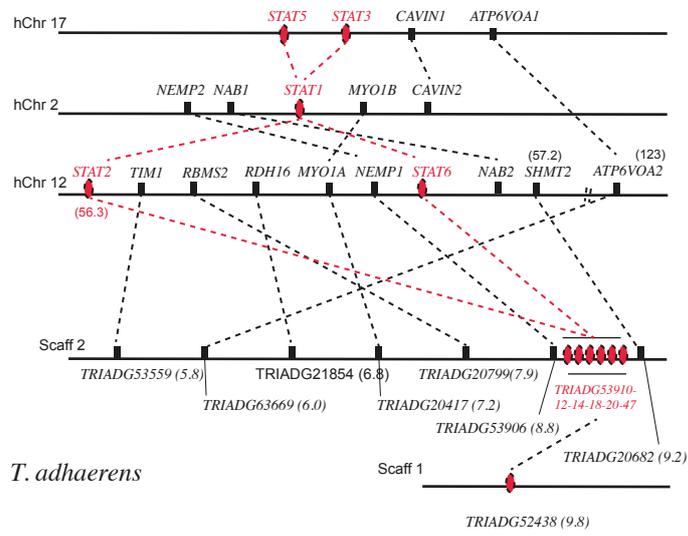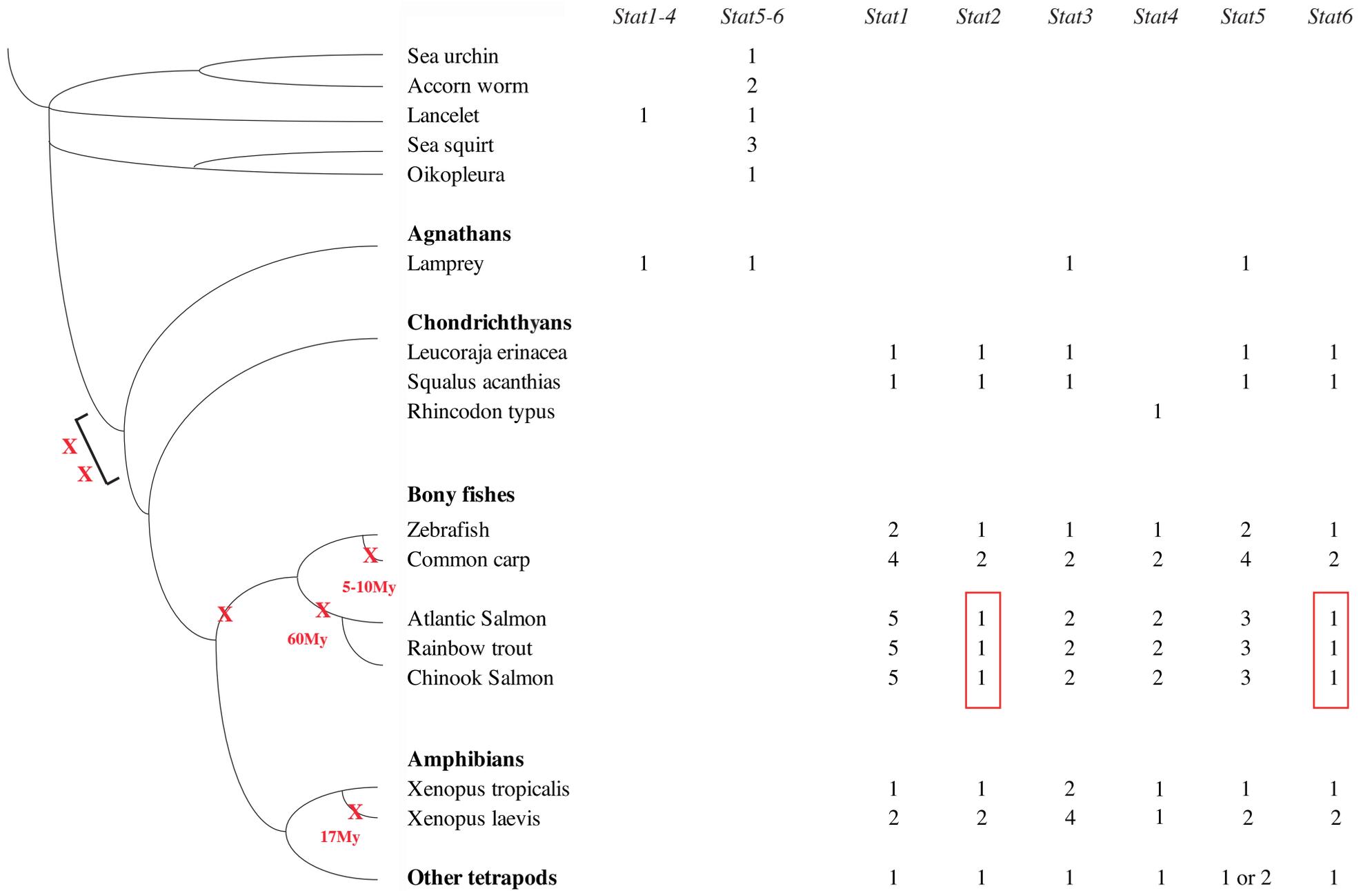**Table 1**: List of potentially functional stat genes in five species of salmonids fish. The nomenclature is extended from *Danio rerio*.

| | *Oncorhynchus* | | | | *Salmo* | |
|---|---|---|---|---|---|---|
| | *nerka* | *mykiss* | *tshawytscha* | *Kisutch*** | *truttae* | *salar* |
| stat 1a1 | 115129673 (5/+) | 100136755 (18/-) | 112266551 (14/-)[a] | 109906681 (16/-) | 115195984 (6/+) | 100136558 (16/+) |
| stat 1a2 | 115133971 (2/-) | 100137016 (7/-) | 112253897 (7/-)[b] | 109890649 (5/-) | 115179269 (39/-)[##] | |
| stat1a3 | 115133986 (2/+) | 110527523 (7/+)[*,d,§] | 112253898 (7/+)[c] | 109890433 (5/+)[d] | 115179399 (39/+)[d] | 106575214[#] (17/-) |
| stat 1b1 | 115138513 (1/-) | 110501544 (22/+) | 112244575 (Un/+) | 109871062 (26/+) | 115156118 (20/+) | 100196256 (21/-)[e] |
| stat 1b2 | 115108041 (3/-) | 110520020 (3/-) | 112235369 (3/-) | 109865728 (2/+) | 115160666 (24/-) | 106586142 (25/-) |
| stat 2 | 115143174 (15/+) | 110494323 (17/+) | 112217577 (2/+) | 109895769 (1/+) | 115207827 (14/-) | 100270812 (12/+) |
| stat3a | 115110281 (26/-) | 110538194 (12/-) | 112225989 (27/-)[***] | 116374454 (6/-) | 115171116 (32/-) | 106601025 (3/-) |
| stat3b | 115104128 (21/-) | 100136756 (13/-) | 112258660 (9/-) | 109898422 (10/-) | 115197994 (1/-) | 106607297 (6/+) |
| stat4a | 115108086 (3/-) | 110520023 (3/-) | 112235400 (3/-) | 109865765 (2/-) | 115160669 (24/-) | 106586145(25/-) |
| stat4b | 115138591 (1/-) | 110501546 (22/+) | 112225151 (26/-)[***] | 109870786 (26/+) | 115156112 (20/+) | 100380385 (21/-) |
| stat5a1 | 115103495 (21/-) | 100135887 (13/-)[**] | 112258659 (9/-)[##] | 109897291 (10/-) | 115198021 (1/-) | 106607295 (6/+)[**] |
| stat5a2 | 115110283 (26/-) | 110538192 (12/-) | 112225727 (27/-)[##] | 109893324 (6/-) | 115171115 (32/-) | 100380532 (3/-) |
| stat5b | 115106643 (23/+) | 110491683 (16/+) | 112223777 (24/+) | 109865530 (20/+) | 115161213 (2/+) | 106579144 (19/+) |
| stat6 | 115102666 (20/+) | 110491929 (16/+) | 112221662 (22/+) | 109869625 (24/+) | 115167799 (30/+) | 106567004 (13/+) |

[*] annotated as "uncharacterised protein", [**] RefSeq status indicated as provisional, [***] annotated as "low quality protein", [#] annotated as pseudogene in NCBI but as coding for ENSSSAT00000091945 in Ensembl, [##] annotated as pseudogene, [a] duplicated due to genome assembly errors, identical to 112253955, [b] duplicated due to genome assembly errors, identical to 112253778, [c] duplicated due to genome assembly errors, identical to 112253779, [d] double size, [e] small size (not included in phylogenetic analysis).

[**] The chinook genome still contains some assembly errors resulting in artificially duplicated regions in particular between the chromosomes 7 and 14 and within the chromosome 7. In this species, only one stat5 was annotated as functional in contrast to 3 stat5 genes found in the five other species analysed, named stat5.1-3.

[§] stat1a3 is doubled in NCBI but separate in Ensembl as ENSOMYG00000034815 and ENSOMYG00000035706, both annotated as stat1a.

**Table 2.** Additional genes encoding short ORF with significant hit to signal transducer and activator of transcription. None could be identified in *O. tshawytscha* nor in *O. kisutch.*

| Species | GeneID (chromosome/orientation) | Size of the largest isoform (aa) |
|---|---|---|
| *O. mykiss* | 110503197 (24/+) | 243 |
| *S. truttae* | 115188323 (Un/+) | 184 |
| | 115189982 (Un/+) | 249 |
| | 115190002 (Un/-) | 136 |
| | 115156120 (20/+) | 249 |
| | 115156122 (20/+) | 186 |
| | 115156123 (20/+) | 179 |
| | 115156124 (20/+) | 147 |
| | 115156138 (20/+) | 186 |
| | 115156139 (20/+) | 184 |
| *S. salar* | 106583229 (22/-) | 314 |

**Table 3.** Presence and number of *stat* genes in genomes of Chondrichthyans and Agnathans

| Phylum | Species | Gene number | Gene name | Conserved domains |
|---|---|---|---|---|
| Chondrichthyans | Elephant shark *Callorhincus milli* | >=4 | ENSCMIG00000003696 (3) | STATi-STATa-STATb-SH2 |
| | | | ENSCMIG00000003757 (5b) | STATi-STATa-STATb-SH2 |
| | | | ENSCMIG00000010732 (1) | STATi-STATa-STATb-SH2 |
| | | | ENSCMIG00000015418 (1) | STATi-STATa-STATb-SH2** |
| | Whale shark *Rhyncodon typus* | | XP_020376005 | STATi-STATa-STATb-SH2 |
| Agnathans | Lamprey *Petromyzon marinus* | 4 | ENSPMAG00000000244 | STATi-STATa-STATb-SH2 |
| | | | ENSPMAG00000002770 | STATi-STATa-STATb-SH2 |
| | | | ENSPMAG00000006622 | … STATa-STATb… |
| | | | ENSPMAG00000009008 | STATi-STATa-STATb-SH2 |
| | Hagfish *Eptatretus burgeri* | 4/5 | ENSEBUG00000002998& | … STATb-SH2 |
| | | | ENSEBUG00000003439 | …-SH2 |
| | | | &ENSEBUG00000006280 | (separated by gls) |
| | | | ENSEBUG00000007506 | STATi-STATa-STATb-SH2 |
| | | | ENSEBUG00000013827 | STATi-STATa-STATb |

** there are 2 STAT genes in this entry !!

**Table 4.** Presence and number of *stat* genes in genomes of other deuterostomians.

| Phylum | Species | Gene number | Gene name | Conserved domains |
|---|---|---|---|---|
| Non vertebrate Deuterostomians | Sea urchin *Strongylocentrotus purpuratus* | 1 | SP-STAT | STATi-STATa-STATb-SH2 |
| | *Saccoglossus kowalevskii* (Hemichordates) | 1 | XP_006814941 | STATi-STATa-STATb-SH2 |
| | *Branchiostoma floridae* (Cephalochordates) | 2 | XP_019630041/BL18533 XP_002594129BL09530 | STATi-STATa-STATb-SH2 STATi-STATa-STATb-SH2 |
| | Oikopleura dioica (Appendicularia) | 1 | AAS21327 | STATi-STATa-STATb-SH2 |
| | *Ciona intestinalis* (Tunicates) | 3 | ENSCING00000004044 ENSCING00000010308 ENSCING00000024295 | STATi-STATa-STATb-SH2 …STATa-STATb-SH2 STATi-STATa |