**The RAG transposon is active through the deuterostome evolution and domesticated in jawed vertebrates**.

**Jose Ricardo Morales Poole [1], Sheng Feng Huang [2], Anlong Xu [2, 3], Justine Bayet [1], Pierre Pontarotti [1]**

Jose Ricardo Morales Poole and Sheng Feng Huang contributed equally to the work.

1 Aix Marseille Université, CNRS, Centrale Marseille, I2M UMR 7373, équipe évolution biologique modélisation, 13453, Marseille, France.

2 State Key Laboratory of Biocontrol, Guangdong Key Laboratory of Pharmaceutical Functional Genes, School of Life Sciences, Sun Yat-sen University, Guangzhou, 510275, People's Republic of China.

3 Beijing University of Chinese Medicine, Dong San Huan Road, Chao-yang District, Beijing, 100029, People's Republic of China.

**Abstract**

RAG1 and RAG2 are essential subunits of the V(D)J recombinase required for the generation of the variability of antibodies and T-cell receptors in jawed vertebrates. It was demonstrated that the amphioxus homologue of RAG1-RAG2 is encoded in an active transposon, belonging to the transposase DDE superfamily. The data provided supports to the possibility that the RAG transposon has been active through the deuterostome evolution and is still active in several lineages. The RAG transposon corresponds to several families present in deuterostomes. RAG1-RAG2 V(D)J recombinase evolved from one of them, partially due to the new ability of the transposon to interact with the cellular reparation machinery. Considering the fact that the RAG transposon survived millions of years in many different lineages, in multiple copies, and that DDE transposases evolved

27  their association with proteins involved in repair mechanisms, we propose that the apparition of

28  V(D)J recombination machinery could be a predictable genetic event.

29

30  **Introduction**

31  The recombination-activating gene products known as RAG1 and RAG2 proteins constitute the

32  enzymatic core of the V(D)J recombination machinery of jawed vertebrates. The RAG1-RAG2

33  complex catalyzes random assembly of variable, diverse and joining gene segments that are present

34  in the jawed vertebrates genomes in numerous copies and together, with hyper-mutation, generate

35  the great diversity of the assembled antibodies and T-cell receptors. Therefore, the RAG1-RAG2

36  role in the V(D)J rearrangement of antigen receptors is crucial for the jawed vertebrates adaptive

37  immunity (Teng and Schatz 2015). Concerning the origins of RAG1-RAG2, it remains elusive for

38  more than 30 years as the genes were only found in jawed vertebrates (Danchin E. *et al.* 2004). On

39  the other hand, striking similarities between RAG1 and DDE transposase has been noted: common

40  reaction chemistry for DNA cleavage, similar organization of protein domain structure and

41  similarities between recombination signal sequences (RSSs) and terminal inverted repeat (TIRs)

42  targeted by transposases (Kapitonov and Jurka 2005; Fugmann 2010). The hypothetical transposon

43  ancestry of RAG was further supported upon the demonstration of RAG1-RAG2 mediated

44  transposition *in vitro* (Agrawal *et al.* 1998; Hiom *et al.* 1998) and *in vivo* (Chatterji *et al.* 2006;

45  Curry *et al.* 2007; Ramsden *et al.* 2010; Vanura *et al.*, 2007), thought the efficiency of such

46  reactions in vivo is highly disfavored comparing to recombination.

47  A next step in the understanding of RAG1-RAG2 recombinase evolution was the discovery of a

48  RAG1-RAG2-like locus in purple sea urchin genome, where genes for both proteins are oriented in

49  close proximity in a head-to-head manner as RAG1-RAG2 locus in vertebrates. However this locus

50  lacks TIR and thus does not show the typical features of a transposon (Fugmann *et al*. 2006).Due to

51  the similarity between RAG1 and *Transib* transposon (a family from the DDE transposon

52  superfamily) and the fact that RAG2 lacks similarity to any known transposon protein, even though

53   it harbors Kelch-like repeats and PHD domains as other eukaryotic proteins, led several authors to

54   propose that a *Transib*-like transposon joined the deuterostomian ancestor genome followed by

55   exons shuffling events bringing *Transib* and the ancestor of RAG2 together (Fugmann 2010). As a

56   result, the RAG1-RAG2 locus was then recruited for an unknown function. A second much more

57   recent recruitment as RAG1-RAG1 V(D)J recombinase most likely occurred at the base of the

58   jawed vertebrates evolution. Kapitonov and Koonin (2015) went a step further and provided *in*

59   *silico* evidence that RAG1 and RAG2 subunits of the V(D)J recombinase evolved from two proteins

60   encoded in a single transposon as they found three sequences that could correspond to fossilized

61   RAG1-RAG2 transposon (including TIRs) in one starfish genome. A major step in the

62   understanding of the RAG1-RAG2 evolution was reported by our group (Huang *et al*. 2016)

63   showing for the first time the presence of an active RAG transposon in the cephalochordate

64   *Branchiostoma belcheri* named ProtoRAG. The full length ProtoRAG transposon is bound by 5 bp

65   target site duplications (TSDs) and a pair of terminal inverted repeats (TIRs) resembling V(D)J

66   recombination signal sequences (RSSs). Between the TIRs reside tail-to-tail oriented, intron-

67   containing and co-transcribed, RAG1-like and RAG2-like genes. The RAG transposon has been

68   recently active in amphioxus as shown by indel polymorphisms. Furthermore the amphioxus

69   RAG1-RAG2-like proteins could mediate TIR-dependent transposon excision, host DNA

70   recombination, transposition and even signal joint formation at low frequency, using reaction

71   mechanisms similar to those used by vertebrates RAGs (Huang *et al*. 2016).

72   Here we bring more information about the evolution of RAG transposons. We show that beside *B.*

73   *belchieri,* an active RAG transposon is found in the hemichordate *Ptychodera flava,* that several

74   fossilized transposons are found in several deuterostomes species suggesting that RAG transposon

75   has been active through the history of the deuterostome lineage.

76

77   **Results**

78   **Description of an active RAG transposon in *P. flava* and many fossilized transposons in**

79    **deuterostomes**

80    Due to the discovery of an active RAG transposon in amphioxus *B. belchieri*, we screened all the

81    available genome and EST projects using the query sea urchin RAG1L and RAG2L sequences.

82    Many hits in several deuterostomians species were found, hits are found in protosomians but they

83    show low similarity and correspond to the transib transposons (Panchin and Moroz 2008) and the

84    chapaev transposon family (Kapitonov and Jurka 2007). The family reported by Panchin and Moroz

85    (2008) as well as many other families were found during our survey. However the connection

86    between these families and the RAG1-RAG2 is not clear even if they are related.

87    Among the hits found in deuterostomes, one of them corresponds to a complete transposon and

88    other several fossilized transposons (see Figure 1 and Supplementary Table 1) in the hemichordate

89    *P. flava*. In other deuterostome species, we found evidence for RAG1L-RAG2L structures without

90    TIRs but with many fragment copies of the RAG1L-RAG2L locus. some species with an

91    incomplete transposon with TIR and RAGL sequences and many other copies of RAG1L-RAG2L

92    fragments. The presence of TIR on many of these copies might indicate that they correspond to

93    fossilized transposons. Transcribed sequences database are available for several deuterostomes and

94    in most of the case RAG1L and RAG2L transcripts are found, complete or incomplete, thus

95    revealing the domestication of the transposon or their activity.

96    Based on the phylogeny of the RAG1L and RAG2L protein sequences (see Figure 2 and

97    supplementary table 1 for the phylogenetic analysis and the families description), we can find

98    several RAG families in *P. flava*. Among them, B and C families have unambiguous TIR and TSD

99    structure. Two copies of B family show a TSD-5TIR-RAG1L-RAG2L-3TIR-TSD structure. While

100    one of these two copies encodes a complete RAG1L and RAG2L protein, the other one corresponds

101    to RAG1L and RAG2L pseudogenes. However its presence confirms that the authentic RAG

102    transposon appears in this family. The C family has one copy with TSD-5TIR-(RAG1L-RAG2L)-

103    3TIR-TSD structure, this copy seems to be inactivated (several in frame stop codons,

104    Supplementary Table 1). In addition,  three 5TIR-3TIR copies with no recognizable RAG1/2 genes

105 and one of those copies has both TSD. We also found 12 structures having the 5' or 3'TIR. We

106 failed to find TSD-TIR structure for other RAG-like families (A and unclassified families) in *P.*

107 *flava*, this could be due to the poor genome assembly or to the fact that some families have become

108 inactive. Anyway, these findings are sufficient to show that multiple families of RAG transposon

109 have been and are thriving in *P. flava*. Moreover, we found several fossilized transposons in the case

110 of *Patiria minata* as partially described in 2015 (Kapitonov and Koonin 2015), a 5TIR-

111 RAG1L_fragment-3TIR structure containing TSD and no RAG2L protein, a 5TIR adjacent to

112 RAG1L structure (TSD-5TIR-RAG1L) and other several structures having the 5' and 3'TIR but

113 without internal RAG coding sequences. These structures indicate that RAG was an active

114 transposon during the echinoderms evolution. Afterwards a comparative sequence analysis was

115 made in *B. belcheri*, *Branchiostoma floridae*, *P. flav*a (Pfl) and *P. minata* (Pmi) TIR sequences

116 (Figure 3) showing no identity between different *Transib*, vertebrate RSS and amphioxus, Pmi and

117 Pfl species except the first CAC nucleotides. Nonetheless both sequences analyzed in amphioxus,

118 share TIR similarity, suggesting a possible common origin of RAG transposon between these two

119 species of amphioxus. However, there is no identity between B and C RAG transposon families in

120 *P. flav*a, suggesting, despite the similarity between RAG-like proteins of both families, no TIR

121 similarity between each other, as they may be not reactive or functionally compatible. Previously,

122 an equivalent of RSS nonamer, a stretch of nine highly conserved nucleotides has been found in the

123 amphioxus ProtoRAG TIR, though this ProtoRAG nonamer has no similarity with the nonamer

124 found in RSS (Huang *et al*. 2016). However, there is no such nonamer or equivalently conserved

125 oligomer found in *P. minata* and *P. flava* B and C ProtoRAG family. All this suggests that the

126 nonamer structure is not important in echinoderms and hemichordates phyla, but became important

127 in amphioxus and vertebrates.

128 The species tree in Figure 1 (see also Supplementary Table 1) shows a summary of RAG1L-RAG2L

129 sequences distribution in deuterostomes according to the available data. When genomic and

130 transcription data are available the species names appear in red, whereas when only genomic data

131  are available the species names are shown in blue, and when only available expressed sequence data

132  corresponds to the species name are black. It is likely that the transposon is active if bona fide

133  sequences are present in the genome in several copies and fragments and if the putative transposons

134  are transcribed as in the case for *P. flava* RAGL-B and *B. belcheri*. On the other hand *P. miniata*

135  seems not to be transcribed since only fossilized transposons are found in the genome. In two

136  species of sea urchin, *Eucidaris tribuloides* and *Lytechinus variegatus,* no transcribed sequences are

137  found, but many copies of RAG1L-RAG2L are present on the genome without TIRs indicating that

138  might be fossilized transposons that became inactivated by the loss of the TIR sequences.

139  The case of *S. purpuratus* is more difficult to understand: the published RAG1L-RAG2L locus

140  (Fugmann *et al*. 2006) renamed here RAG1L-RAG2L B1, was believed to be domesticated, as the

141  RAG1L and RAG2L coding sequences are not interrupted by stop codons, RAG1L and RAG2L are

142  transcribed. And could be functional, but because no TIR sequences has been identified they cannot

143  be a transposon (Fugmann *et al*. 2006). However, we found many fragments which were highly

144  similar to this sequence in the *S. purpuratus* genome, which could reveal a recent transposition

145  event followed by the domestication of one of its copies (see supplementary data and Figure 3). We

146  found another RAGL copy which arose from a duplication event which occurred at the origin of the

147  echinoderms, named RAG1L-B2. The RAG1L-B2 copy is only found fragmented with multiple

148  recent copies in the genome whereas it is complete as RAG1L transcript. A possible explanation for

149  this second locus could be the existence of an active transposon with the genome sequence not well

150  assembled or otherwise a domesticated or recent fossilized transposon. For most of the species we

151  do not have information at the genomic level, but if we find RAGL transcript, this sequence could

152  correspond to an active transposon, domesticated transposon or recent pseudogene. This shows that

153  the transposon has been present in their ancestors.

154

155  **Features of the proteins encoded by the RAG-like proteins**

156  In ambulacraria (echinoderm and hemichordate) the deuterostome RAG1-like, 816-1136 aa-long

157  shares around 26.52% sequence identity between RAG1L-B family and vertebrate RAG1, around

158  33.21% between the orthologous RAGL-A family and the vertebrate RAG1 and only 27.79%

159  between RAG1L-A and RAG1L-B, while inside RAG1L-B family are sharing 48.75% of sequence

160  identity and only 20.13% respect to *Transib* transposase in terms of core region. As regards to

161  RAG1 lancelet, 30.47% and 37.62% sequence identity are shared with A and B families

162  respectively and only 27.45% with RAG1 vertebrate (see Supplementary Figure 2A). Clusters of

163  high identity are found between RAG1L and vertebrate RAG1 along much of their length,

164  suggesting conservation of multiple functional elements. Vertebrate RAG1 uses four acidic residues

165  to coordinate critical divalent cations at the active site (Ru *et al.* 2015) and all four are conserved in

166  RAG1L (Supplementary Figure 1A, red highlight). In addition, many cysteine and histidine residues

167  that coordinate zinc ions and play a critical role in proper folding of RAG1 (Kim *et al.* 2015), are

168  conserved between RAG1L and vertebrate RAG1 (Supplementary Figure 1A, * and # symbols).

169  However, RAG1L does not share much identity with vertebrate RAG1 in the region corresponding

170  to the nonamer binding domain, consistent with the fact that RAG transposons TIRs have no clear

171  similarity to the RSS nonamer. In fact, different families of RAG1-like have little similarity to each

172  other in the putative nonamer binding domain, consistent with the fact that different ProtoRAG

173  families have very different TIR sequences and no obvious nonamer regions, excluding the

174  amphioxus TIR. Finally, there are also some RAG1-like specific conserved regions (see

175  Supplementary Figure 1, underlined by *). It should be noted that PflRAG1L-A and jawed

176  vertebrate RAG1 show conserved position in the alignment not shared with other RAG1L families,

177   RAG2L 366-535aa long, shares low sequence identity between B family and vertebrate RAG2

178  (18.69%) and between B family and lancelet RAG2L (25.02%). On the other hand the RAG2L-B

179  family shares around 45.90% while lancelet RAG2L shares only 20.24% identity with RAG2

180  vertebrate (supplementary Figure 2B). However, the N-terminal six-bladed β-propeller domain (six

181  Kelch-like repeats), which is conserved in both vertebrate RAG2 and ProtoRAG RAG2L, can be

182  discerned in RAG2L. Strikingly, amphioxus RAG2L lacks the entire RAG2 C-terminal region,

183    including the PHD domain as shown previously (Huang *et al*. 2016). However, this PHD domain is

184    present in all other echinoderm and hemichordate RAG2 proteins (Supplementary Figure 1B). Thus,

185    the absence of this region in amphioxus RAG transposon might be a secondary loss.

186

187    **Phylogenetic relation between the RAG families**

188    The phylogenetic analysis with the complete RAG sequences from the available deuterostome

189    species are shown in Figure 2A and 2B and synthesized in Table 1. At least two sub-families have

190    been present in the ancestral deuterostome, named RAGL-B and RAGL-A. Other families such as

191    RAGL-C have not been included in the phylogenetic history as they are found only in one species

192    (Table 1).

193    In the case of the orthologous relation found between RAG1L-A of *P. flava* (hemichordate) and

194    vertebrates RAG1 recombinase, we can observe that RAGL-A was lost in many lineages excluding

195    hemichordates and jawed vertebrates. RAGL-B conversely, is lost in tunicates and in vertebrates

196    lineage but conserved in several lineages as cephalochordates, hemichordates and echinoderms.

197    Furthermore the phylogenetic analysis shows that RAGL-B has been duplicated in the echinoderms

198    ancestor after the hemichordates/echinoderms split, and both copies have been kept (even if most of

199    them have been inactivated) in most of the echinoderm species (Table 1 and Figure 2C).

200

201    **RAG transposon has been active during the deuterostome evolution**

202    From the RAG transposon status: active, fossilized, domesticated, absent (Figure 1 and see the

203    description of an active transposon in *P. flava* and many fossilized transposons chapters), we can

204    proposed the following evolutionary history (Figure 4). The transposon has been active in the

205    deuterostomes ancestor and in the branch that leads to the common ancestor of chordate, still active

206    in cephalochordates and domesticated as a RAG1-RAG2 V(D)J recombinase in the common

207    ancestor of jawed vertebrates. The transposon has been lost in the Petromyzon lineage. The

208    transposon has been active in the branch originated from the node between deuterostomes and

209 ambulacraria antecesors. It remains active in hemichordates inside the subphylum of Enteropneusta

210 (at least on the *P. flava* lineage) but is lost in the other enteropneusts as *S. kowalevskii.*

211 Unfortunately we do not have genome information for the other hemichordates subphyla:

212 Pterobranchia. In the case of the echinoderms lineage, the transposon has been present in the

213 echinoderms common ancestor, in the branch leading to the common ancestor of crinoid, in the

214 clade formed by the sea urchin and holothuroids and in the clade formed by starfishes/ophiures. It

215 has been then lost in the crinoid lineage. The transposon has been active in the branch that goes

216 from the common ancestor of echinoderms to the common ancestor of sea urchin/Holothuroids and

217 starfishes/brittle stars. Concerning the Asteroidea/Ophiuroidea group, the transposon has been

218 active in their common ancestor and has been active in the Ophiure lineage in particular in *O.*

219 *spicalatus* where the transposon is likely to be active or has lost its activity recently. The transposon

220 seems to have been inactive in the starfish lineage but fragments showing similarities to RAG1-L

221 and/or RAG2-L transposons are found in this species. Furthermore, transposons are clearly found

222 fossilized in *P. miniata*. In the case of sea urchin/holothuroids group, it seems that the transposon

223 has been active in their common ancestor and inactive in the holothurian lineage. We should also

224 note that the transposon seems to be active in some sea urchin lineages as in *E. tribuloide* but much

225 less in others.

226

227 **Discussion**

228 In this report we show that a RAG transposon has been present in the deuterostome common

229 ancestors and was active since then in some lineages, fossilized later during evolution and

230 domesticated at least in the case of jawed vertebrates. The structural and regulatory features that

231 cause the jawed vertebrate RAG V(D)J recombinase to favor deletional/inversional recombination

232 over transposition as in the case of the RAG transposase (Huang *et al*. 2016) is not yet resolved. It

233 could be explained by how the cleaved ends and particularly the signal ends are processed. The

234 RAG V(D)J recombinase binds signal ends tightly as excepted for a transposase but it has acquired

235    the possibility to give up these ends efficiently to the non-homologous end joining machinery. This

236    allows recombination and prevents the propagation (Teng and Schatz 2015). Thus the jawed

237    vertebrate V(D)J recombinase differs from the current RAG transposon, as well as its transposon

238    precursor, in how it interfaces with the DNA repair apparatus. This new property occurred likely in

239    the jawed vertebrates common ancestor.

240    DDE transposases have been shown to interact with repair proteins. For example the Sleeping

241    Beauty transposase interacts directly with the Ku70 repair protein (Izsvák *et al*. 2004) and the pogo

242    transposase of *D. melanogaster* interacts with the proliferating cell nuclear antigen (PCNA), a key

243    protein for DNA replication and repair (Warbrick *et al*. 1998). Therefore, the associations of DDE

244    transposon with DNA repair and replication factors appear to evolve in a convergent manner

245    (Feschotte and Pritham 2007). This characteristic and the fact that the transposon survived during

246    millions of years in multiple copies in different lineages increased the probability of the co-option

247    of the RAG transposon as V(D)J recombinase. Therefore the apparition of V(D)J recombination

248    machinery in the jawed vertebrates phyla could be labeled as a predictable genetic events.

249    Our results could also explain better the origins of the T-cell receptor and B-cell receptor gene

250    organization. The earlier proposed scenario (Fugmann 2010; Koonin and Krupovic 2015; Hsu and

251    Lewis 2015) involved an insertion of the RAG transposon into the ancestral IG/TCR V-gene, prior

252    to the externalization of the RAG1-RAG2 complex while leaving the RSS-like TIR within the

253    IG/TCR V-gene. This was followed by duplication of this new genetic structure: **V**RSS-RSS**J**. The

254    RAG transposon was then co-opted as V-J recombinase and the system started to work. However,

255    this scenario explains the V-J structure IG light chain, TCR alpha and gamma chain but not the VDJ

256    organization of IG heavy chain or TCR beta and delta chains. Hsu and Lewis (2015) proposed the

257    following scenario for the origin of the D segment: the duplication of the VJ unit, followed by J- to

258    V- recombination and the insertion of non-templated N-region into the signal joint generates a proto

259    D segment. We proposed here an alternative hypothesis to explain the VDJ organization: while one

260    RAG was domesticated (likely RAGL-A orthologue), other RAG transposons (likely RAGL-B

261 orthologue) were still active as one of them split the **VRSS-RSSJ** copy and gave rise to **VRSS-**

262 **RSSDRSS-RSSJ**. RAGL-B transposase became then extinct and finally was lost during vertebrate

263 evolution.

264

265 **Material and methods**

266 **Identification of RAG1 and RAG2-like sequence in different data bases**

267 RAG1-RAG2-like locus identified in the echinoderm *Strongylocentrotus purpuratus* and in the

268 vertebrates genome were used as a protein sequence to perform a TBLASTN-based search against

269 the NCBI nr protein, transcriptome shotgun assembly (TSA) and the WGS database as of June 2016

270 (Altschul *et al.* 1990). These retrieved sequences were extracted and translated by ExPASy

271 Translate tool. Potential open reading frames of RAG1-RAG2 elements used in this study were

272 predicted using FGENESH (Solovyev *et al.* 2006) with the sea urchin organism specific gene-

273 finding parameters. The mRNA sequences were then assembled into contigs by CAP3 (Huang and

274 Madan 1999).

275

276 **Phylogenetic analysis**

277 The alignment and trees were constructed using MEGA6 (complete deletion, WAG with Freqs. (+F)

278 correction model, 1000 bootstrap replicates in Tamura *et al.* 2013). Thus, whether they are active,

279 fossilized or domesticated was classified into families. Short sequence copies were analyzed one by

280 one           with           the           reference           data           set           .

281

282 **Sequence searches for TIR and TSD motifs**

283 We used three methods to search target site duplication (TSD) and terminal invert repeat (TIR)

284 sequences. In the first method, the upstream and downstream 20 Kb of sequence flanking the

285 RAG1-RAG2-like sequences were extracted and separated into a set of small fragments (using a

286 window size of 60 bp and a step size of 1 bp). In the first method, each upstream fragment was

287    compared with each downstream fragment for 4-6 bp TSDs and possible TIRs using a custom Perl

288    script. We required 40% identity for potential TIR pairs, and allowed only one mismatch for TSD

289    pairs. In the second method, all upstream fragments were compared against all downstream

290    fragments using BLAST. We required a minimum e-value of 100 and sequence identity of 40% in

291    the BLAST search. However, these two methods failed to work well and provided no reliable

292    results. Therefore, we turn to the third method. In this method, we posited that if there are multiple

293    copies of ProtoRAG transposons in the genome assembly, comparison between these copies could

294    help to determine their terminal sequences (TIR, etc.).

295

296       We focused on finding more complete elements that contain both TIR and RAG gene fragments,

297    such as "5TIR-RAGs-3TIR", "5TIR-RAGs" and "RAGs-3TIR".

298       Here is our procedure:

299    1. First we identified all genomic regions containing RAG1/2 fragments by using TBLASTN

300        and the amphioxus and vertebrate RAG1/2 proteins as queries;

301    2. The region containing RAG1/2 plus upstream 20kb and downstream 20kb was extracted,

302        which we called the RAG region;

303    3. Because there should have a clear border between the ProtoRAG and the host DNA, we

304        could determine the potential 5' and 3'-terminal of the ProtoRAG transposon by comparing

305        RAG regions with each other by using BLASTN (see the Figure below);

306    4. Finally, we examined the potential 5/3-terminal sequences of the RAG regions. Most of

307        them had been destroyed and therefore no detectable TIRs, but several of them shown clear

308        and intact TIR structure.

309    5. And the TSD if presents, should be right next to the TIR sequences.

310
311
312    Therefore, the sequences containing the RAG1/2-like fragments and the 20 Kb flanking regions

313    were compared to each other and also to the whole genome assembly using BLAST. The terminal

314    sequences were analyzed using a custom Perl script and then subjected to manual inspection.

315

316    **Summary of data availability**

317    In order to detect the absence or presence of a given structure in the genome or transcriptome, we

318    need to extract all the available taxonomic information from the NCBI database. It has to be noted

319    that even if the sequence for a given genome is not complete, when RAG1L-RAG2L seems to be an

320    active transposon, we should find an active or at least a fossilized transposons (in several copies).

321    Focusing on the genome database we can find species such as *Parastichopus parvimensis*,

322    *Acanthaster planci*, *Ophiothrix spiculata*, *Petromyzon marinus, Branchiostoma belcheri*,

323    *Oikopleura dioica, Botryllus schlosseri & Ciona savignyi.* Transcript sequences can be provided for

324    *Saccoglossus kowalevskii*, *Anneissia japonica*, *Psathyrometra fragilis*, *Abyssocucumis albatrossi*,

325    *Sclerodactyla briareus*, *Apostichopus japonicus*, *Parastichopus californicus*, *Echinarachnius*

326    *parma*, *Evechinus chloroticus*, *Paracentrotus lividus*, *Sphaerechinus granularis*, *Arbacia*

327    *punctulata*, *Henricia sp*. AR-2014, *Echinaster spinulosus*, *Peribolaster folliculatus*, *Leptasterias sp*.

328    AR-2014, *Pisaster ochraceus*, *Marthasterias glacialis*, *Asterias rubens*, *Asterias forbesi*, *Asterias*

329    *amurensis*, *Luidia clathrata*, *Patiria pectinifera & Ophiocoma echinata.* Finally, together with

330    genomic information and transcript expression we have *Ptychodera flava*, *Eucidaris tribuloides*,

331    *Strongylocentrotus purpuratus*, *Lytechinus variegatus*, *Patiria miniata*, *Homo sapiens*, *Mus*

332    *musculus*, *Gallus gallus*, *Xenopus tropicalis*, *Latimeria chalumnae*, *Danio rerio*, *Carcharhinus*

333    *leucas*, *Carcharhinus plumbeus*, *Branchiostoma floridae* and *Ciona intestinalis*.

334

335    **References**

336    *1.* Agrawal A. *et al.,* 1998. Transposition mediated by RAG1 and RAG2 and its implications

337      for the evolution of the immune system. Nature **394**, 744-751.

338    *2.* Altschul S.F. *et al*., 1990. Basic local alignment search tool. J. Mol. Biol. 215:403-410.

339    *3.* Chatterji M. *et al.,* 2006. Mobilization of RAG-generated signal ends by transposition and

340        insertion in vivo. *Mol. Cell. Biol*. **26**, 1558-1568.

341     *4.* Curry, J.D. *et al.*, 2007. Chromosomal reinsertion of broken RSS ends during T cell

342        development. J. Exp. Med., 204, 2293-2303.

343     *5.* Danchin E. *et al.,* 2004. The major histocompatibility complex origin. *Immunol. Rev.* **198**,

344        216-232.

345     *6.* Feschotte C., Pritham E.J, 2007. DNA transposons and the evolution of eukaryotic genomes.

346        *Annu. Rev. Genet.* **41**, 331-368.

347     *7.* Fugmann S.D. *et al*., 2006. An ancient evolutionary origin of the Rag1/2 gene locus. *Proc.*

348        *Natl. Acad. Sci. USA* **103**, 3728-3733.

349     *8.* Fugmann S.D., 2010. The origins of the Rag genes from transposition to V(D)J

350        recombination. *Semin. Immunol.* **22**, 10-16.

351     *9.* Hiom K. *et al.,* 1998. DNA transposition by the RAG1 and RAG2 proteins: a possible

352        source of oncogenic translocations. *Cell* **94**, 463-470.

353    *10.* Hsu, E., and Lewis, S.M., 2015. The Origin of V(D)J Diversification. In Molecular Biology

354        of B Cells, F.W. Alt, T. Honjo, A. Radbruch, and M. Reth, eds. Elsevier, Academic Press,

355        Amsterdam, pp. 133-148.

356    *11.* Huang S. *et al*., 2016. Discovery of an Active RAG Transposon Illuminates the Origins of

357        V(D)J Recombination. *Cell.* **166**, 102-114.

358    *12.* Huang X., Madan A., 1999. CAP3: A DNA sequence assembly program. *Genome Res.* **9**,

359        868-877.

360    *13.* Izsvák Z. *et al*., 2004. Healing the wounds inflicted by sleeping beauty transposition by

361        double-strand break repair in mammalian somatic cells. *Mol Cell.* **13**, 279-290.

362    *14.* Kapitonov V.V., Jurka J., 2005. RAG1 core and V(D)J recombination signal sequences were

363        derived from Transib transposons. *PLoS Biol.* **3**, 998-1011.

364    *15.* Kapitonov V.V., Jurka J., 2007. Chapaev-a novel superfamily of DNA transposons. Repbase

365         Reports. 7, 777–777.

366    *16.* Kapitonov V.V., Koonin E.V., 2015. Evolution of the RAG1-RAG2 locus: both proteins

367         came        from        the        same        transposon.        *Biol.        Direct*,

368         http://biologydirect.biomedcentral.com/articles/10.1186/s13062-015-0055-8.

369    *17.* Kim M.S. *et al.,* 2015. Crystal structure of the V(D)J recombinase RAG1-RAG2. *Nature*

370         **518**, 507-511.

371    *18.* Koonin E.V., Krupovic M., 2015. Evolution of adaptive immunity from transposable

372         elements combined with innate immune systems. *Nat. Rev. Genet.* **16**, 184-192.

373    *19.* Panchin Y., Moroz L.L., 2008. Molluscan mobile elements similar to the vertebrate

374         recombination-activating genes Biochem Biophys Res Commun. 369, 818–823.

375    *20.* Ramsden D.A. *et al*., 2010. Weed, B.D., Reddy, Y.V. V(D)J recombination: Born to be wild.

376         *Semin Cancer Biol.* **20**, 254-260.

377    *21.* Ru H. *et al.,* 2015. Molecular Mechanism of V(D)J Recombination from Synaptic RAG1-

378         RAG2 Complex Structures. *Cell.* **163**, 1138-1152.

379    *22.* Solovyev V. *et al*., 2006. Automatic annotation of eukaryotic genes, pseudogenes and

380         promoters. *Genome Biol.* **7,** S10.1-S10.12.

381    *23.* Tamura K., *et al*., 2013. MEGA6: Molecular Evolutionary Genetics Analysis version 6.0.

382         *Mol. Biol. Evol.* **30**, 2725-2729.

383    *24.* Teng G., Schatz D.G., 2015. Regulation and Evolution of the RAG Recombinase. *Adv.*

384         *Immunol.* **128**, 1-39.

385    *25.*        Vanura K. *et al*., 2007. In vivo reinsertion of excised episomes by the V(D)J

386         recombinase:    a    potential    threat    to    genomic    stability.    PLoS    Biol.    5,

387         http://journals.plos.org/plosbiology/article?id=10.1371/journal.pbio.0050043

388    *26.*        Warbrick E. *et al*., 1998. PCNA binding proteins in Drosophila melanogaster: the

389    analysis of a conserved PCNA binding domain. *Nucleic Acids Res.* 26, 3925-3932.

390

391    We thank the EBM laboratory for advice and Olivier Loison for editing the manuscript.
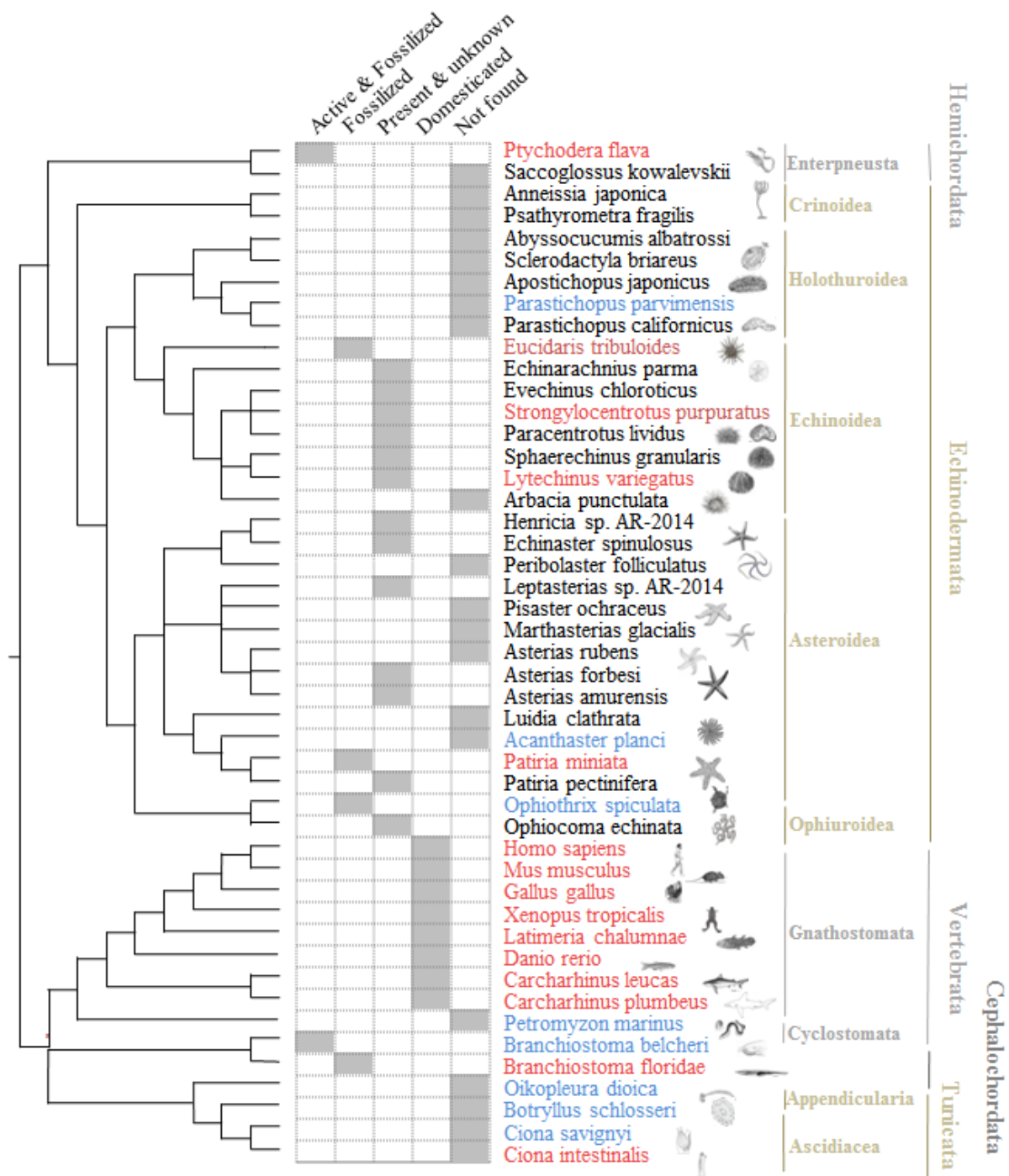
392

393    **Author's contribution: JRMP, PP and SFH conceived the project and design the study. JRMP,**

394    **PP and SFH analyzed the results. JRMP, PP, SFH and ALX wrote the manuscript.**

395

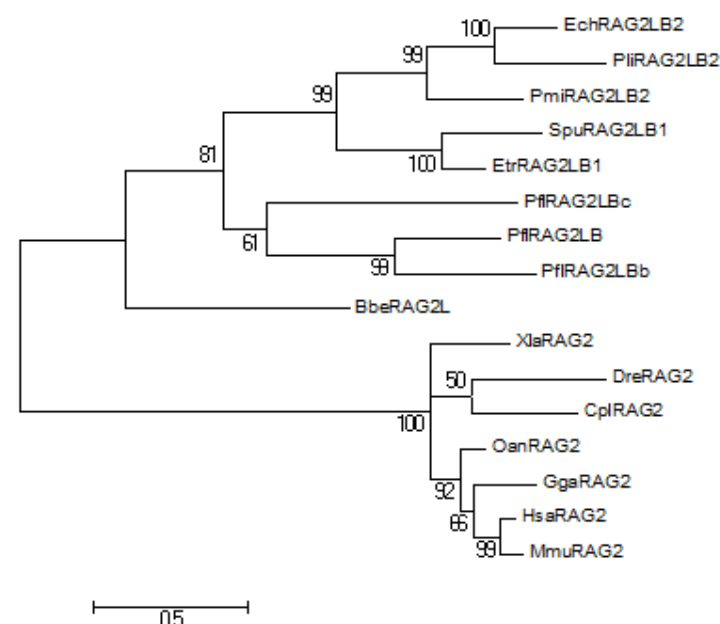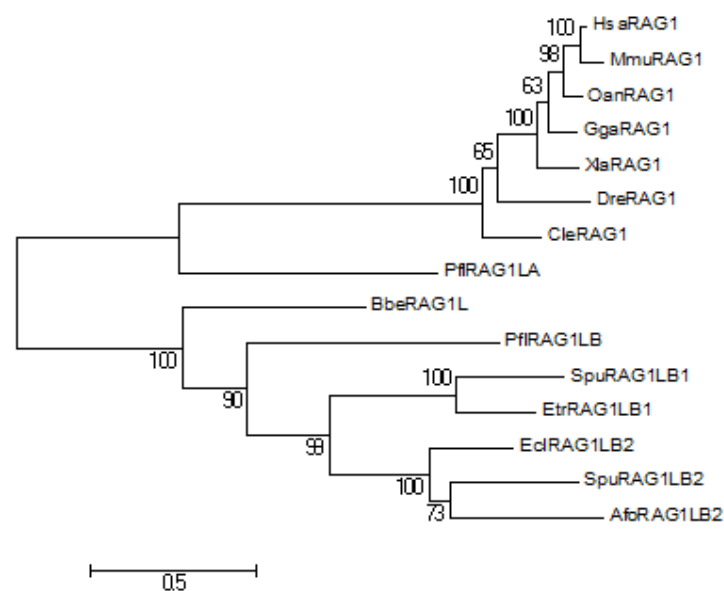396    The authors declare no competing financial interest.

397    Correspondence and requests for materials should be addressed to **pierre.pontarotti@univ-amu.fr**

398    and **morales.poole@gmail.com.**
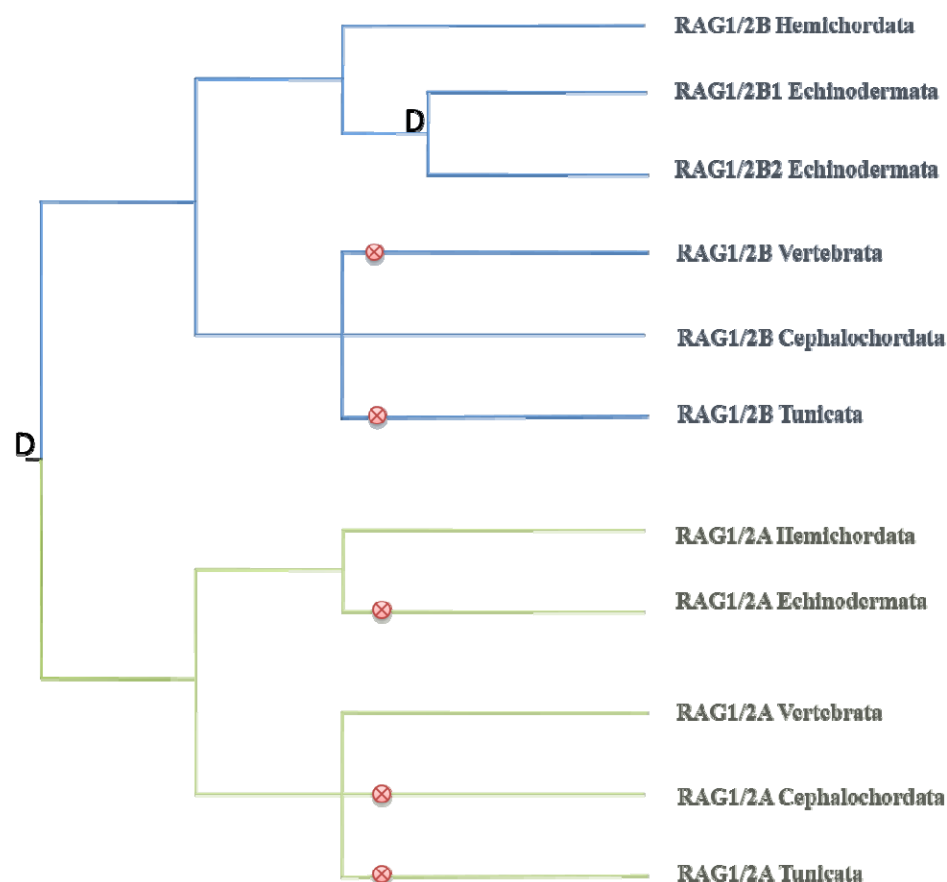
399

400    **Figures and tables**

401

**Figure 1 | Distribution of the RAG1-RAG2 sequences in deuterostomes.** Only species for which

the genomic and/or transcription data are available are represented in the phylogenetic tree. Species

are colored in red when genomic and transcription data are available, in blue when only genomic

data are available and species are colored in black only when expressed sequence data are available.

406

407

408



409

**Figures 2 |**

**Phylogenetic analysis with the complete RAG1 (2A) and RAG 2 sequences (2B).** See also table

1. Phylogeny of two families A and B and other families such as P. flava RAG C are only found in

one species (see also table 1). It is therefore difficult to decipher their story.

415

**Figure 2C | Schematic drawing of duplications and losses of the RAG families during the**

**deuterostome evolution: duplication (D) and lost (⊗).**

416

417

418

```
ccgCTGCtgc  : 11 : Pf1_B_BCFJ01017854_BCFJ01052781
gccCAATGtgc : 11 : Pf1_B_BCFJ01094280
cacTGTGG--- :  8 : Pf1_B_BCFJ01052780
cacCATCCgta : 11 : Pf1_C_BCFJ01036631
aacCCCGGctg : 11 : Pf1_C_BCFJ01107546
gtcTGGCA--- :  8 : Pf1_C_BCFJ01102604
ctcGGGTG--- :  8 : Pf1_C_BCFJ01084502
ttaCCTTC--- :  8 : Pf1_C_BCFJ01046932
tgcTGCCA--- :  8 : Pf1_C_BCFJ01016857
tcgCAGTG--- :  8 : Pf1_C_BCFJ01107167
gtgCATTG--- :  8 : Pf1_C_BCFJ01047137
tgcGGGCG--- :  8 : Pf1_C_BCFJ01031953
---CGCATtcg :  8 : Pf1_C_BCFJ01070588
---GGGTGcca :  8 : Pf1_C_BCFJ01150129
---GGCCGgtc :  8 : Pf1_C_BCFJ01103012
---CAGTGtgg :  8 : Pf1_C_BCFJ01107168
---CGCGCtcg :  8 : Pf1_C_BCFJ01287958
---CACGGgta :  8 : Pf1_C_BCFJ01048214
---CGCATtcg :  8 : Pf1_C_BCFJ01083599
----------- :  - : blank
cgtCCAGGgtc : 11 : Pmi_JH779599
aacCCATAccg : 11 : Pmi_JH774215
ctcTGTTAtat : 11 : Pmi_JH780459
gagTTTAG--- :  8 : Pmi_JH769343
gatTTTAG--- :  8 : Pmi_JH774292
tgtTGATG--- :  8 : Pmi_JH782081
gcgATGTG--- :  8 : Pmi_JH775549
aagCGGGA--- :  8 : Pmi_JH781149
---CCTGGcat :  8 : Pmi_AKZP01156453
---TTCAGcaa :  8 : Pmi_JH771625
---CATATgca :  8 : Pmi_AKZP01165822
---GTGCGgga :  8 : Pmi_AKZP01162400
```

419

```
ccgCTGCtgc  : 11 : Pf1_B_BCFJ01017854_BCFJ01052781
gccCAATGtgc : 11 : Pf1_B_BCFJ01094280
cacTGTGG--- :  8 : Pf1_B_BCFJ01052780
cacCATCCgta : 11 : Pf1_C_BCFJ01036631
aacCCCGGctg : 11 : Pf1_C_BCFJ01107546
gtcTGGCA--- :  8 : Pf1_C_BCFJ01102604
ctcGGGTG--- :  8 : Pf1_C_BCFJ01084502
ttaCCTTC--- :  8 : Pf1_C_BCFJ01046932
tgcTGCCA--- :  8 : Pf1_C_BCFJ01016857
tcgCAGTG--- :  8 : Pf1_C_BCFJ01107167
gtgCATTG--- :  8 : Pf1_C_BCFJ01047137
tgcGGGCG--- :  8 : Pf1_C_BCFJ01031953
---CGCATtcg :  8 : Pf1_C_BCFJ01070588
---GGGTGcca :  8 : Pf1_C_BCFJ01150129
---GGCCGgtc :  8 : Pf1_C_BCFJ01103012
---CAGTGtgg :  8 : Pf1_C_BCFJ01107168
---CGCGCtcg :  8 : Pf1_C_BCFJ01287958
---CACGGgta :  8 : Pf1_C_BCFJ01048214
---CGCATtcg :  8 : Pf1_C_BCFJ01083599
cgtCCAGGgtc : 11 : Pmi_JH779599
aacCCATAccg : 11 : Pmi_JH774215
ctcTGTTAtat : 11 : Pmi_JH780459
gagTTTAG--- :  8 : Pmi_JH769343
gatTTTAG--- :  8 : Pmi_JH774292
tgtTGATG--- :  8 : Pmi_JH782081
gcgATGTG--- :  8 : Pmi_JH775549
aagCGGGA--- :  8 : Pmi_JH781149
---CCTGGcat :  8 : Pmi_AKZP01156453
---TTCAGcaa :  8 : Pmi_JH771625
---CATATgca :  8 : Pmi_AKZP01165822
---GTGCGgga :  8 : Pmi_AKZP01162400
```

420

421

**Figure 3 | Full length alignment of the TSD and TIR sequences.** (3A) Alignment of RAG transposon TSDs and flanking sequences from the *P. flava*, *P. miniata*, *B. belcheri*, *B. floridae* genome. The length of the TSDs is the same: 5bp for the Transib and the RAG-L transposon which indicates similar mechanisms of transposition (3B). Alignment of *ProtoRAG* TIR sequences with the consensus RSS and *Transib* TIR. IUPAC codes used in the alignment: N=A, C, G or T; K=G or T; W=A, T; V=A, C or G. Lower case indicates an undetermined nucleotide. Shading indicates sequence conservation, with darker gray indicating a higher degree of conservation. Bbe:
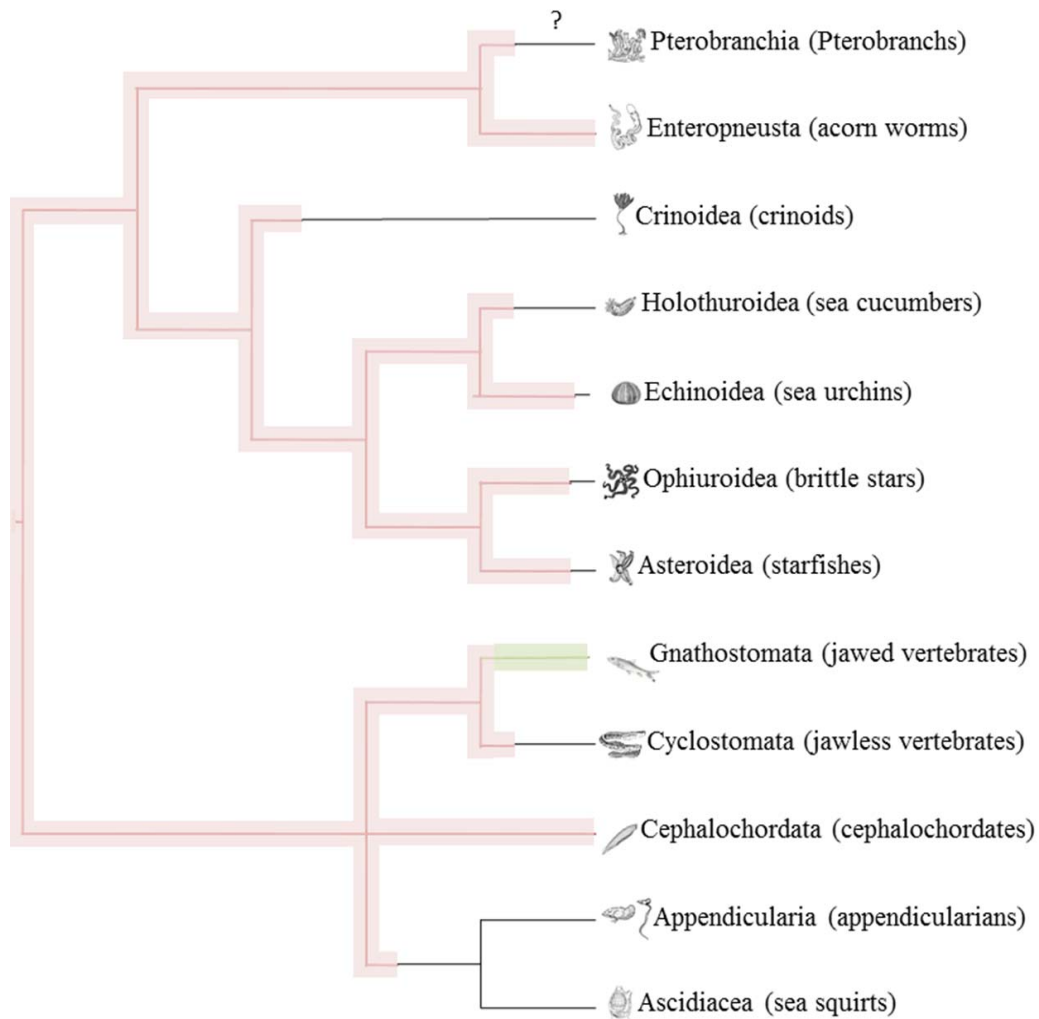
429    *B. belcheri*; Bfl: *B. floridae*. Pfl: *P. flava*, Pmi: *P. minata* RAG *transposon* copy identification

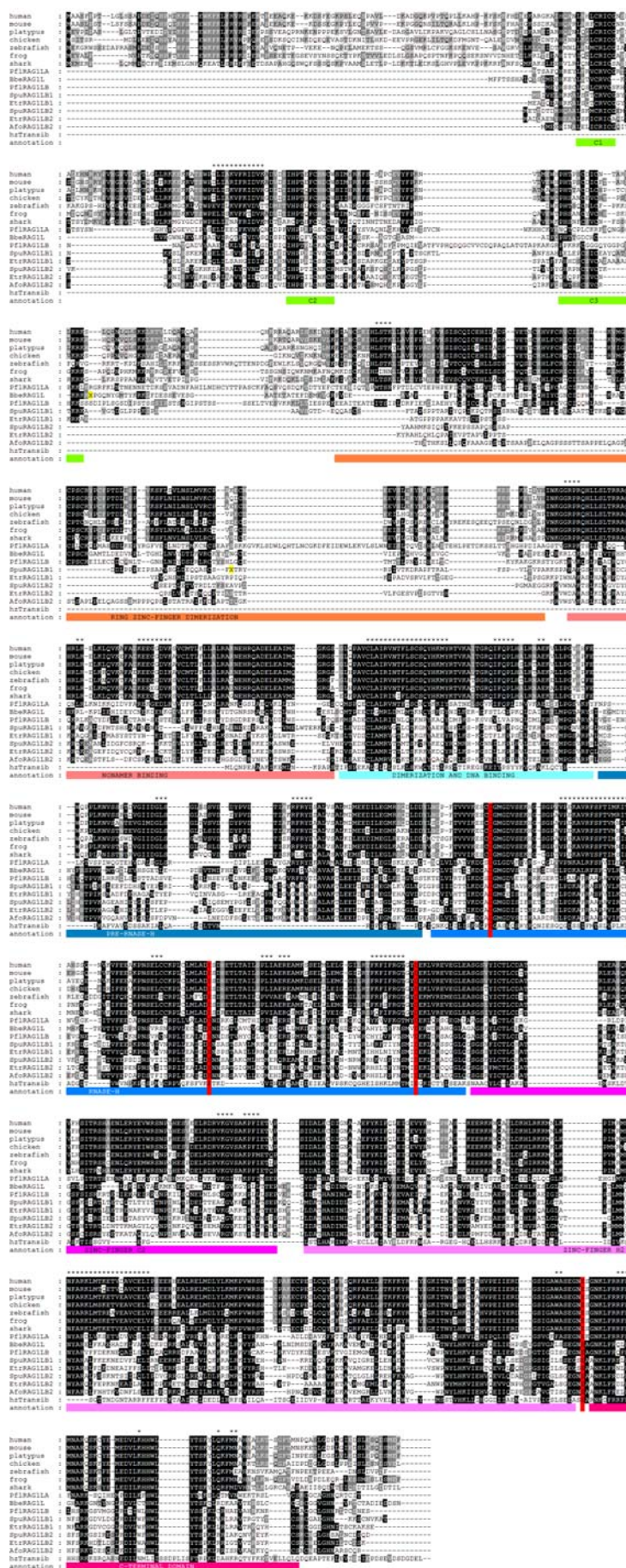430    numbers correspond to those listed in Table S1.

431



432

433    **Figure 4 | Evolution of the RAG transposon.** Transposon activity is indicated in bold pink and

434    V(D)J recombinase activity is indicated in bold green.

435

| | P. flava | L. variegatus | E. tribuloides | P. lividus | S. purpuratus | E. parma | P. pectinifera | A. amurensis | Leptasterias sp. | A. forbesi | E. chloroticus | P. miniata | O. spiculata | Henricia sp. AR-2014 | S. granularis | E. spinulosus | O. echinata | A. japonica | P. fragilis | A. albatrossi | S. briareus | A. japonicus | P. parvimensis | P. californicus | A. punctulata | P. folliculatus | P. ochraceus | M. glacialis | A. rubens | L. clathrata | A. planci |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RAG1B1-like | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| RAG1B2-like | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| RAG1B-like | * | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Other families | ** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Not phylogenetically assigned | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Not found | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| RAG2B1-like | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| RAG2B2-like | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| RAG2B-like | *** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Other families | **** | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Not phylogenetically assigned | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| Not found | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

436

**Table 1 | Presence of RAG subfamilies in the different species.** Sequences were classified through phylogenetic analysis. Short sequence copies were analyzed one by one against the reference data described in Figure 2A and 2B  The classification as B family (or A family labeled with "**") is straightforward as it is based on orthologous relationships between different phyla (differences between echinoderms and hemichordates for example). Inside B family, two groups named B1 and B2 are found in several echinoderms. If an echinoderm sequence is classified as B family, but not as B1 or B2 we call it B-like (RAG1Bd-like is labeled with "*" while RAG2Bb-like and RAG2Bc-like are labeled with "***"). We have two specific cases, C family only found in *P. flava* (RAG1 labeled with "**" and RAG2 labeled with "****") and X family in *Ophiotrix spiculata*. The rest of species which do not belong to A or B family, are not phylogenetically assigned due to the fact that none enough phylogenetic signals are available.

448

449

450

451

452

**Supplementary figure and table**

454

455
456 **Figure S1 | Features of the proteins encoded by the RAG and RAG-like proteins.** (A) Protein

457 alignment of RAG1L with vertebrate RAG1. Repeat motifs in amphioxus and the purple sea urchin

458 RAG1L were removed and replaced with an "X" and highlighted in yellow. Three regions of

459 conserved cysteine and histidine residues that might bind zinc are underlined with green bars. The

460 N-terminal zinc binding dimerization domain is underlined with dark-red bars. The subdomains of

461 the RAG1 core region are indicated with colored bars. The conserved acidic catalytic residues are

462 highlighted with red shading (D600, E662, D708 and E962 on mouse RAG1). The PflRAG1LA is

463  more similar to vertebrate RAG1, and those regions were labeled with "*". GenBank accessions for

464  mouse RAG1, shark RAG1, lancelet RAG2L and sea urchin RAG1L are NP_033045,

465  XP_007886047, KJ748699 and NP_001028179, respectively.

466  (B) Protein alignment of RAG2L with vertebrate RAG2. Color shading shows the conservation of

467  physiochemical properties. The N-terminal amino acid sequences correspond to Kelch-like repeats.

468  The central conserved GG motifs of the six Kelch-like repeats are underlined in red. The plant

469  homeodomain (PHD) is also underlined below the alignment. GenBank accessions for mouse

470  RAG2, shark RAG2, lancelet RAG2L and sea urchin RAG2L are NP_033046, XP_007885835,

471  KJ748699 and NP_001028184, respectively.

472

473

```
Percent Identity  Matrix - created by Clustal2.1

 1: transib-1_HM 100.00   18.97   21.85   19.39   19.54   19.60   19.93   18.49   18.40   19.83   21.75
 2: BbeRAG1L       18.97  100.00   39.96   35.03   39.49   37.53   40.20   37.88   30.47   27.55   27.35
 3: PflRAG1LB      21.85   39.96  100.00   35.39   40.37   40.84   41.04   39.00   28.49   26.54   27.30
 4: SpuRAG1LB1     19.39   35.03   35.39  100.00   62.11   43.13   46.87   42.06   27.34   24.74   24.69
 5: EtrRAG1LB1     19.54   39.49   40.37   62.11  100.00   45.70   47.65   43.71   28.79   26.60   26.77
 6: AfoRAG1LB2     19.60   37.53   40.84   43.13   45.70  100.00   56.55   54.70   26.78   26.87   25.86
 7: EclRAG1LB2     19.93   40.20   41.04   46.87   47.65   56.55  100.00   58.77   27.55   28.35   26.75
 8: SpuRAG1LB2     18.49   37.88   39.00   42.06   43.71   54.70   58.77  100.00   26.80   27.03   25.79
 9: PflRAG1LA      18.40   30.47   28.49   27.34   28.79   26.78   27.55   26.80  100.00   33.26   33.15
10: HsaRAG1        19.83   27.55   26.54   24.74   26.60   26.87   28.35   27.03   33.26  100.00   64.41
11: CleRAG1        21.75   27.35   27.30   24.69   26.77   25.86   26.75   25.79   33.15   64.41  100.00
```

474

```
Percent Identity  Matrix - created by Clustal2.1

 1: HsaRAG2    100.00   55.34   19.93   17.10   19.21   16.34   19.74   18.62   16.75   18.54
 2: CplRAG2     55.34  100.00   20.54   17.05   20.21   16.75   20.16   19.05   19.90   21.04
 3: BbeRAG2L    19.93   20.54  100.00   29.50   28.61   21.05   23.91   27.46   20.53   25.07
 4: PflRAG2LBc  17.10   17.05   29.50  100.00   29.80   28.67   29.26   30.77   27.63   27.82
 5: PflRAG2LB   19.21   20.21   28.61   29.80  100.00   29.50   29.61   34.13   33.33   29.48
 6: StrRAG2LB1  16.34   16.75   21.05   28.67   29.50  100.00   64.16   41.29   37.20   37.01
 7: EtrRAG2LB1  19.74   20.16   23.91   29.26   29.61   64.16  100.00   44.74   42.64   39.91
 8: PmiRAG2LB2  18.62   19.05   27.46   30.77   34.13   41.29   44.74  100.00   53.79   49.33
 9: EchRAG2LB2  16.75   19.90   20.53   27.63   33.33   37.20   42.64   53.79  100.00   59.09
10: PliRAG2LB2  18.54   21.04   25.07   27.82   29.48   37.01   39.91   49.33   59.09  100.00
```

475  **Figure S2 | Percent Identity Matrix of RAG1 (S2A) and RAG2 (S2B).** In order to provide a

476  multiple alignment, Clustal-Omega requires a guide tree which defines the order in which

477  sequences/profiles are aligned. A guide tree in turn is constructed, based on a distance matrix.

478  Conventionally, this distance matrix is comprised of all the pairwise distances of the sequences. The

479  distance measure Clustal-Omega uses for pairwise distances of unaligned sequences is the k-tuple

480    measure. By default, the distance matrix is used internally to construct the guide tree and is then

481    discarded. By specifying, the internal distance matrix can be written to file.

482

483    **Table S1 | RAGL distribution in non-chordate genome and expressed sequence.** Distribution in

484    the cephalordate phyla: *B. belcheri* and *B. floridae* available in [13].

485
486