

1 **Parallel evolution of amphioxus and vertebrate small-scale gene** 2 **duplications**

3 Brasó-Vives et al.

4 Supplementary note 1: *Proto-RAG*

5 The recombination-activating genes (RAG1 and RAG2) in jawed vertebrates encode the two subunits of a
6 protein complex essential for V(D)J recombination in immunoglobulin genes. They have their origin in a
7 transposon domestication in vertebrate evolution, a RAGB transposon. An active *Proto-RAG* transposon has
8 been described in the *B. belcheri* genome containing both RAG1 and RAG2 genes flanked by terminal inverted
9 repeats (TIR) (Huang *et al.* 2016). The TIR has a heptamer-spacer-nonamer structure similar to that of the
10 recombination signal sequence (RSS) essential for V(D)J recombination in jawed vertebrates.

11 We searched for RAGB (or *Proto-RAG*) transposons in *B. lanceolatum* genome, by performing a tBLASTn with
12 the sequences of *Proto-RAG1* (UniProt: A0A185KID9) and *Proto-RAG2* (UniProt: A0A185KIE0) from *B.*
13 *belcheri* against BraLan3. We encountered three regions containing both *Proto-RAG1* and *Proto-RAG2* blast hits
14 (referred to as sequences 12, 16 and 18 according to the chromosome of their respective location). TIR sequences
15 were found in the flanking sequences of all these regions by performing a BLASTn with known *Proto-RAG* TIR
16 sequences of *B. blecheri* (Tao *et al.* 2020) against these regions' sequences with 10kpb flanking sequences.

17 To predict completeness or fragmentation of *Proto-RAG* genes in the encountered sequences, we annotated the
18 transposon genes with Softberry FGENESH++ (Solovyev *et al.* 2006). This approach allowed us to see that
19 sequences 12 and 18 seemed complete and potentially active, while sequence 16 appeared incomplete and stopped
20 at the level of *Proto-RAG1*. Sequences 12 and 18 have 98.8% of sequence similarity and present the same structure
21 described for *Proto-RAG* transposon in *B. belcheri* with which they share 75% of sequence identity in the RAG1
22 and RAG2 genes sequences. They lack a PHD domain in RAG2 gene sequence, unlike what is seen in jawed
23 vertebrates' RAG2 sequences. The TIR sequences of the three *B. lanceolatum* transposons (including the
24 truncated sequence 16) are identical and show high conservation with other amphioxus *Proto-RAG* TIR
25 sequences (see alignments below).

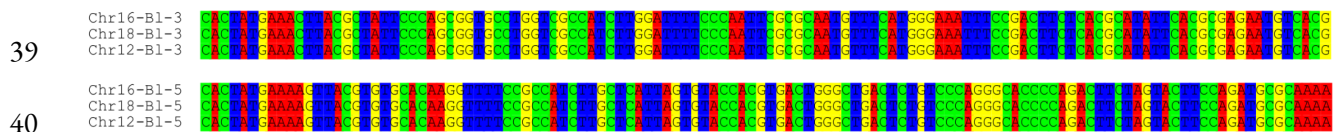
26 Finally, we searched BraLan3 for *Miniature Inverted-repeat Transposable Elements* (MITE) sequences containing
 27 5' and 3' TIRs separated by less than 1000 bp but containing neither RAG 1 and RAG 2. We found 13 MITE
 28 sequences in 10 different chromosomes. These results suggest that the *Proto-RAG* transposon is active in the
 29 amphioxus genome.

30 Huang S., X. Tao, S. Yuan, Y. Zhang, P. Li, *et al.*, 2016 Discovery of an Active RAG Transposon Illuminates the
 31 Origins of V(D)J Recombination. *Cell* 166: 102–114. <https://doi.org/10.1016/j.cell.2016.05.032>

32 Solovyev V., P. Kosarev, I. Seledsov, and D. Vorobyev, 2006 Automatic annotation of eukaryotic genes,
 33 pseudogenes and promoters. *Genome Biol.* 12.

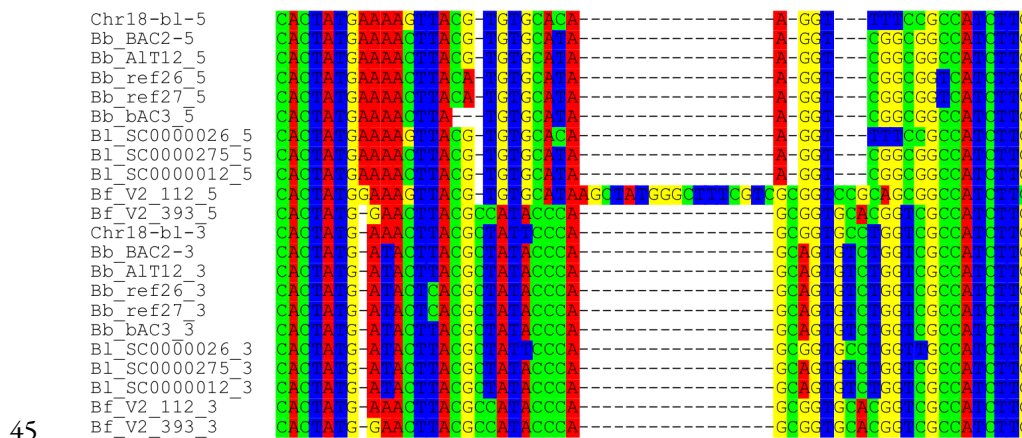
34 Tao X., S. Yuan, F. Chen, X. Gao, X. Wang, *et al.*, 2020 Functional requirement of terminal inverted repeats for
 35 efficient *ProtoRAG* activity reveals the early evolution of V(D)J recombination. *Natl. Sci. Rev.* 7: 403–
 36 417. <https://doi.org/10.1093/nsr/nwz179>

37
 38



41 Alignment of the TIR sequences of the tree *Proto-RAG* transposons found in BraLan3 (5' and 3' separately). Both
 42 the 3' and the 5' TIR sequences of the tree Proto-RAG transposons found in *B. lanceolatum* genome are identical
 43 in sequence.

44



46 Joined alignment of 5' and 3' TIR sequences of *Proto-RAG* transposons identified in Branchiostoma species.

47 Chr18-bl-5 and Chr18-bl-3 correspond to the 5' and 3' TIR sequences from the ProtoRAG transposon identified

48 here in BraLan3 chromosome 18 (as representatives of the BraLan3 TIR sequences). The other TIR sequences
49 are taken from (Tao *et al.* 2020). Bl, Bb and Bf correspond to *B. lanceolatum*, *B. belcheri* and *B. floridae*
50 respectively while the identification numbers correspond to their insertion scaffolds.