



HAL
open science

Insights into RAG evolution from the identification of “missing link” family A RAGL transposons

Eliza C Martin, Lorlane Le Targa, Louis Tsakou-Ngouafo, Tzu-Pei Fan,
Che-Yi Lin, Jianxiong Xiao, Ziwen Huang, Shaochun Yuan, Anlong Xu,
Yi-Hsien Su, et al.

► **To cite this version:**

Eliza C Martin, Lorlane Le Targa, Louis Tsakou-Ngouafo, Tzu-Pei Fan, Che-Yi Lin, et al.. Insights into RAG evolution from the identification of “missing link” family A RAGL transposons. *Molecular Biology and Evolution*, 2023, 10.1093/molbev/msad232 . hal-04243785v2

HAL Id: hal-04243785

<https://hal.science/hal-04243785v2>



Submitted on 4 Nov 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

Insights into RAG Evolution from the Identification of “Missing Link” Family A RAGL Transposons

Eliza C. Martin ^{†1}, Lorlane Le Targa,^{†2} Louis Tsakou-Ngouafo,^{†2} Tzu-Pei Fan,³ Che-Yi Lin,³ Jianxiong Xiao,¹ Ziwen Huang,⁶ Shaochun Yuan,^{6,7} Anlong Xu ^{6,8}, Yi-Hsien Su,³ Andrei-Jose Petrescu,^{*,4} Pierre Pontarotti,^{*,2,5} and David G. Schatz^{*,1}

¹Department of Immunobiology, Yale School of Medicine, New Haven, CT 06520-8011, USA

²Aix-Marseille Université, IRD, APHM, MEPHI, IHU Méditerranée Infection, Marseille 13005, France

³Institute of Cellular and Organismic Biology, Academia Sinica, Taipei 11529, Taiwan

⁴Department of Bioinformatics and Structural Biochemistry, Institute of Biochemistry of the Romanian Academy, 060031 Bucharest, Romania

⁵CNRS SNC 5039, 13005 Marseille, France

⁶State Key Laboratory of Biocontrol, Guangdong Key Laboratory of Pharmaceutical Functional Genes, Southern Marine Science and Engineering Guangdong Laboratory (Zhuhai), School of Life Sciences, Sun Yat-sen University, Guangzhou 510275, China

⁷Laboratory for Marine Biology and Biotechnology, Qingdao National Laboratory for Marine Science and Technology, Qingdao 266237, China

⁸School of Life Sciences, Beijing University of Chinese Medicine, Beijing 100029, China

[†]These authors contributed equally.

*Corresponding authors: E-mails: andrei.petrescu@biochim.ro; pierre.pontarotti@univ-amu.fr; david.schatz@yale.edu.

Associate editor: Dr. Harmit Malik

Abstract

A series of “molecular domestication” events are thought to have converted an invertebrate RAG-like (RAGL) transposase into the RAG1–RAG2 (RAG) recombinase, a critical enzyme for adaptive immunity in jawed vertebrates. The timing and order of these events are not well understood, in part because of a dearth of information regarding the invertebrate RAGL-A transposon family. In contrast to the abundant and divergent RAGL-B transposon family, RAGL-A most closely resembles RAG and is represented by a single orphan RAG1-like (RAG1L) gene in the genome of the hemichordate *Ptychodera flava* (*PfRAG1L-A*). Here, we provide evidence for the existence of complete RAGL-A transposons in the genomes of *P. flava* and several echinoderms. The predicted RAG1L-A and RAG2L-A proteins encoded by these transposons intermingle sequence features of jawed vertebrate RAG and RAGL-B transposases, leading to a prediction of DNA binding, catalytic, and transposition activities that are a hybrid of RAG and RAGL-B. Similarly, the terminal inverted repeats (TIRs) of the RAGL-A transposons combine features of both RAGL-B transposon TIRs and RAG recombination signal sequences. Unlike all previously described RAG2L proteins, RAG2L-A proteins contain an acidic hinge region, which we demonstrate is capable of efficiently inhibiting RAG-mediated transposition. Our findings provide evidence for a critical intermediate in RAG evolution and argue that certain adaptations thought to be specific to jawed vertebrates (e.g. the RAG2 acidic hinge) actually arose in invertebrates, thereby focusing attention on other adaptations as the pivotal steps in the completion of RAG domestication in jawed vertebrates.

Key words: recombination activating gene (RAG), V(D)J recombination, evolution, transposition, DDE transposase, transposon molecular domestication.

Introduction

V(D)J recombination is essential for adaptive immunity in jawed vertebrates and in many species is responsible for generating the vast repertoire of antigen receptors expressed by developing lymphocytes (Gellert 2002; Flajnik 2014). A heterotetramer composed of RAG1 and RAG2 (hereafter, RAG) initiates V(D)J recombination by cleaving DNA at specific recombination signal sequences (RSSs) that flank each V, D, and J gene segment that participates

in the reaction (Fig. 1A) (Schatz and Swanson 2011; Kim et al. 2015). RSSs consist of conserved heptamer and nonamer components separated by either a 12 or 23 bp spacer and cleavage occurs efficiently only in a synaptic complex containing a 12RSS/23RSS pair (the 12/23 rule) (Fig. 1A) (Gellert 2002). RAG1 is composed of a core region essential for DNA cleavage (aa 384–1008; mouse RAG aa numbers are used unless otherwise specified) flanked by a long N-terminal region (NTR; aa 1–383) and a short C-terminal

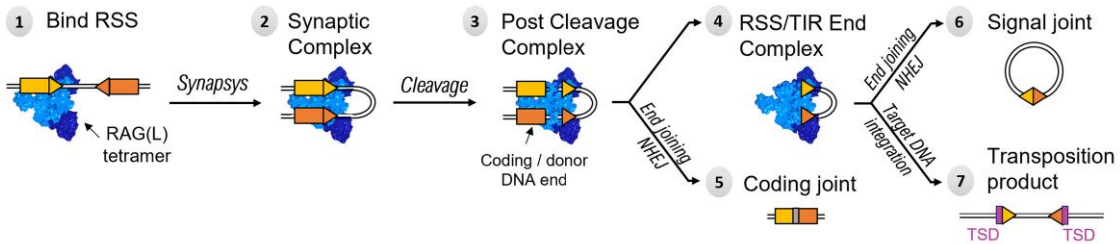
Received: August 09, 2023. **Revised:** September 28, 2023. **Accepted:** October 10, 2023

© The Author(s) 2023. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution.

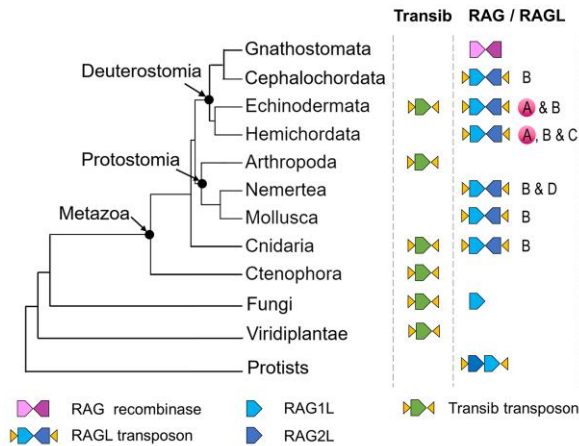
This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<https://creativecommons.org/licenses/by-nc/4.0/>), which permits non-commercial re-use, distribution, and reproduction in any medium, provided the original work is properly cited. For commercial re-use, please contact journals.permissions@oup.com

Open Access

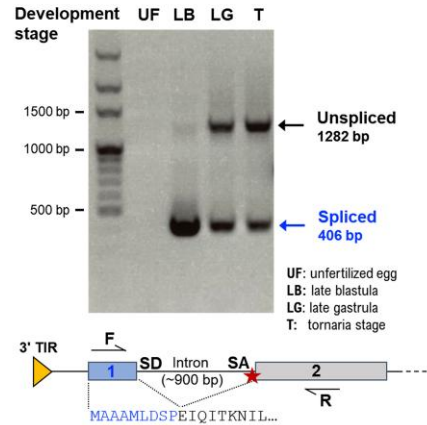
A V(D)J recombination versus transposition



B Taxonomic distribution



D PflRAG2L-A RT-PCR



C RAGL-A genomic loci diagrams

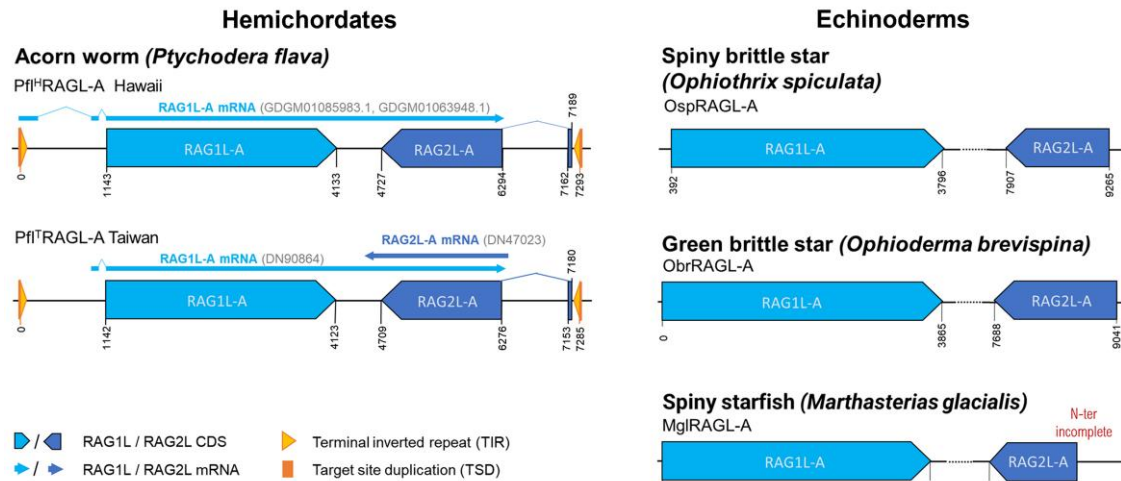


Fig. 1. Overview of RAGL-A transposons in hemichordates and echinoderms. (A) Schematic of V(D)J recombination and transposition. The RAG tetramer binds 2 RSSs (triangles) flanking the gene segments (rectangles) to form a synaptic complex (complex 2), within which cleavage takes place (complex 3). Subsequently, the RSS flanking regions (coding ends) are processed by nonhomologous end joining (NHEJ) to yield a coding joint (product 5). RAG remains bound to the RSS ends (complex 4) with 2 possible outcomes: in V(D)J recombination, processing by NHEJ enzymes to yield a signal joint (product 6), or in transposition, the TIR end complex inserts the mobile element into a new locus (product 7), generating a 5 bp target site duplication (TSD). (B) Taxonomic distribution of RAG/RAGL and *Transib* in eukaryotes. The presence of *Transib* and/or RAG/RAGL families (A to D) in clades of eukaryotes is indicated (clades lacking RAG/RAGL/*Transib* elements are omitted). RAGL-A elements identified in this study in hemichordates and echinoderms are highlighted in magenta. (C) Genomic loci diagrams of the most conserved copies of RAGL-A identified in hemichordates (left) and echinoderms (right). Transcriptomic support, whenever present, is mapped above gene diagrams with arrows. (D) *PflRAG2L-A* RT-PCR illustrating PCR products of the size expected from spliced and unspliced mRNA, with mRNA samples from Taiwan *P. flava* of 4 developmental stages: unfertilized egg (UF), late blastula (LB), late gastrula (LG), and tornaria (T). Replicates and controls are shown in [supplementary Fig. S2](#). Diagram below illustrates the location of the PCR primers, the 9 aa contributed by exon 1, and an in-frame stop codon upstream of exon 2 (star).

tail (CTT; aa 1009-1040). RAG2 consists of a core region required for cleavage activity (aa 1-350), an acidic “hinge” (AH) region of approx. 60 aa, and a C-terminal plant homeodomain (PHD) (Matthews et al. 2007; Schatz and Swanson 2011; Kim et al. 2015).

Structures of the RAG1 core/RAG2 core tetramer bound to DNA or in apo form have provided numerous insights into the mechanism of DNA binding and cleavage (Kim et al. 2015, 2018; Ru et al. 2015; Chen, Cui, Best 2020). DNA cleavage is performed by an RNaseH-fold DDE catalytic domain in the RAG1 core that shares structural similarity with the catalytic domains of cut and paste transposases and retroviral integrases (Montano and Rice 2011; Kim et al. 2015). The RAG1 core also contains the nonamer binding domain (NBD), which forms a dimer that binds the nonamers of the 12RSS and 23RSS in the synaptic complex. The NBD dimer is able to accommodate the length asymmetry of the 12/23RSS pair by pivoting on a flexible linker and is responsible for enforcing the 12/23 rule (Kim et al. 2015, 2018; Lapkouski et al. 2015; Ru et al. 2015). The RAG2 core is a 6-bladed kelch domain (Kim et al. 2015).

Many findings support the model that RAG evolved from a cut and paste transposon, including biochemical and structural similarities with transposases (Carmona and Schatz 2017; Liu et al. 2022) and the ability of RAG to perform transposition in vitro (Agrawal et al. 1998; Hiom et al. 1998). Transib mobile elements were the first to be suggested to share a common ancestor with RAG due to sequence similarities between Transib transposases and the catalytic core of RAG1 and between Transib terminal inverted repeats (TIRs) and RSSs (Kapitonov and Jurka 2005). Subsequent analyses revealed functional and structural similarities between RAG1 and the Transib transposase from the insect *Helicoverpa zea* (HzTransib) (Hencken et al. 2012; Carmona and Schatz 2017; Liu et al. 2019). Transib transposons do not, however, encode a protein resembling RAG2. The discovery of a transposon encoding RAG1-like (RAG1L) and RAG2-like (RAG2L) proteins in *Branchiostoma belcheri* (*Bbe*) provided a closer intermediate in evolution to jawed vertebrate RAG (Huang et al. 2016). The *BbeRAG1L/2L* gene pair preserves the convergent transcriptional orientation of jawed vertebrate RAG1/2 and is flanked by TIRs with heptamer sequences that resemble those of RSSs and Transib TIRs. The *BbeRAG1L/2L* transposase has numerous parallels with RAG and HzTransib including structural similarities, cleavage of DNA by a nick/hairpin mechanism, and the generation of a 5 bp target site duplication (TSD) during integration (Hencken et al. 2012; Huang et al. 2016; Liu et al. 2019; Zhang et al. 2019). RAGL transposons have now been found in numerous invertebrate clades—deuterostomes, protostomes, cnidarians, protists—some of which possess the full complement of expected transposon components: TSD-TIR5'-RAG1L-RAG2L-TIR3'-TSD (Morales Poole et al. 2017; Martin et al. 2020; Tao et al. 2022).

Four families of RAG/RAGL protein sequences have been identified (Morales Poole et al. 2017; Martin et al.

2020) (Fig. 1B). Family A is represented by jawed vertebrate RAG and orphan RAG1L-A open reading frames in the hemichordate *Ptychodera flava* (acorn worm) while family B encompasses virtually all RAGL elements identified in invertebrates. Families C and D are minor variants restricted to 1 lineage or species. The finding of a divergent RAGL element with an atypical organization of its RAG1L and RAG2L genes in the protist *Aureococcus anophagefferens* (*Aan*) (Tao et al. 2022) suggests that additional families might be identified.

RAG has undergone multiple adaptations during its evolution from transposase to recombinase, resulting in a tightly regulated complex whose enzymatic activities can be tolerated by the host. These adaptations, which include “coupled” cleavage of its 2 substrates, a requirement for asymmetric substrates (12/23 rule), and strong suppression of transposition activity in vivo, involved numerous, seemingly unrelated changes to the RAG1 and RAG2 proteins, arguing that RAG domestication occurred in multiple steps (Liu et al. 2022). The order and timing of these steps and whether they occurred before or after jawed vertebrate speciation are unknown.

A significant impediment to our understanding of RAG evolutionary history has been the large gap that exists between known RAGL transposases and RAG. Specifically, intact RAGL-A transposons have not been described and RAGL-B proteins exhibit multiple differences with RAG (Morales Poole et al. 2017; Martin et al. 2020). This latter point is illustrated by 3 adaptations that suppress transposition, arginine 848 in RAG1 and the AH and an inhibitory loop in RAG2, all of which have thus far been identified only in jawed vertebrates, leading to the model that they arose specifically in jawed vertebrates to protect against deleterious transposition events (Zhang et al. 2019, 2020; Liu et al. 2022).

Using iterative search algorithms developed in our previous study (Martin et al. 2020), we have identified previously overlooked RAG2L open reading frames, revealing the first complete RAGL transposons of the A family in *P. flava* (*PffRAGL-A*) and nearly complete A family elements in several species of echinoderms: *Ophiothrix spiculata* (spiny brittle star) (*OspRAGL-A*), *Ophioderma brevispina* (green brittle star) (*ObrRAGL-A*), and *Marthasterias glacialis* (spiny starfish) (*MglRAGL-A*). The encoded RAGL-A proteins intermingle domains and sequence features of jawed vertebrate RAG and RAGL-B transposases while the *PffRAGL-A* TIRs combine features of RSSs and RAGL-B TIRs. Furthermore, unlike all previously described RAG2L proteins, both hemichordate and echinoderm RAG2L-A contain AHs, which we demonstrate are capable of suppressing RAG-mediated transposition. These findings demonstrate that the AH did not arise uniquely in jawed vertebrates and suggest that inhibition of transposition activity began to arise prior to jawed vertebrate speciation. The identified invertebrate RAGL-A elements help bridge the gap between RAGL-B and jawed vertebrate RAG and provide insight into the order and timing of events during RAG domestication.

Results

Multiple Complete RAGL-A Transposon Copies in 2 Populations of *P. flava*

Using a sensitive iterative search strategy described previously (Martin et al. 2020), we performed searches for RAG2L-A sequences in 2 publicly available sequence scaffolds previously reported to contain *PffRAG1L-A* (Morales Poole et al. 2017) (see Materials and Methods). This led to the identification of 2 RAG2L-A sequences, 97% identical at the nucleotide level, encoding a 6-bladed kelch-like domain, an AH, and a PHD-like domain approximately 600 bp downstream of and in convergent orientation with RAG1L-A. Using these RAG2L-A sequences for further searches of the *P. flava* genome revealed 2 additional RAG2L-A loci (97% to 98% identity) for which no linked RAG1L-A counterpart could be found, although their location near scaffold boundaries precludes firm conclusions in this regard.

While the RAG1L/2L-A genomic loci thus identified all appeared to be pseudogenes containing frameshifts, searches of *P. flava* transcriptomic (TSA) data revealed mRNA sequences containing intact RAG1L-A open reading frames, several of which contained partial but intact RAG2L-A sequences on the reverse noncoding strand. This suggested the existence of at least 1 additional transposon copy, not covered by the public whole genome sequence (WGS) data, in the genome of *P. flava* (Hawaiian population; *P. flava*^H), from which both the WGS and TSA data were derived. Indeed, targeted sequencing of *P. flava*^H bacterial artificial chromosome (BAC) clones resulted in the identification of a RAGL-A transposon with the expected complete TSD-TIR5'-RAG1L-RAG2L-TIR3'-TSD configuration (Fig. 1C, supplementary Files S1 and S2). We then searched for RAGL-A sequences in WGS data generated from *P. flava* from Taiwan (*P. flava*^T), which identified 4 RAGL-A loci, one of which had a TSD-TIR5'-RAG1L-RAG2L-TIR3'-TSD configuration with intact RAG1L-A and RAG2L-A genes (Fig. 1C).

Notably, in both *P. flava*^H and *P. flava*^T, the various RAGL-A cassettes reside in different genomic locations, are flanked by different TSDs (supplementary Fig. S1A and B), and contain numerous unique mutations (supplementary file S2). Furthermore, where sufficient genome assembly data were available, we found that integration sites in *P. flava*^H were “empty” in the *P. flava*^T genome, and vice versa, with the empty site displaying a single copy of the TSD (Fig. S1C). To further investigate the dynamics of RAGL-A in *P. flava*, a second worm from the *P. flava*^T population (hereafter, isolate 2) was analyzed and found to contain a RAGL-A transposon in its genome inserted in a position in *P. flava*^T genome sequence scaffold 14 (transposon hereafter referred to as *PffRAGL-A.14*). Notably, this position is empty in the isolate (isolate 1) used for generating the *P. flava*^T genome assembly. The empty locus in isolate 1 displays intact flanking regions and a single copy of the TSD found at the site of integration in isolate 2 (Supplementary Fig. S1C), suggesting

that insertion of RAGL-A at this site is polymorphic in the *P. flava*^T population. Together, our findings argue that *PffRAGL-A* transposition events continued to occur after the Hawaii and Taiwan populations split.

A 5' Coding Exon Provides the *PffRAG2L-A* Start Codon

The *PffRAG2L-A* genomic loci and mRNA sequences from Hawaii and Taiwan populations encode kelch-AH-PHD configurations but all lacked an ATG start codon at the beginning of the kelch domain. The first in-frame ATG codon was within the second blade of the kelch domain and initiating protein synthesis at this site would almost certainly undermine the structural stability and functionality of the domain and would omit numerous upstream in-frame codons. Analysis of the region between the large RAG2L-A exon and the 3' TIR revealed the presence of several potential mRNA splice donor/acceptor motifs, with 1 pair displaying good agreement with the canonical motifs. The possibility of mRNA splicing in this region was investigated by reverse transcription combined with PCR (RT-PCR) using *P. flava*^T RNA samples purified from 4 development stages: unfertilized eggs, late blastula, late gastrula, and tornaria (a larval stage). A PCR product consistent with the predicted spliced mRNA was consistently detected in late blastula, late gastrula, and tornaria stages and was undetectable in unfertilized eggs (Fig. 1D; all 3 biological replicates shown in supplementary Fig. S2A). Control reactions showed that detection of the splice product was dependent on reverse transcription and that genomic DNA contamination was present in some samples, explaining strong signals seen for the unspliced product in some reactions (supplementary Fig. S2B). mRNA quality was verified by amplification of *Vasa* using intron-spanning primers, confirming the presence of intact mRNA in unfertilized egg samples (supplementary Fig. S2C) and supporting the conclusions that maternal *PffRAG2L-A* transcripts are absent in the unfertilized egg and that expression is induced during early development. Sequencing of the *PffRAG2L-A* spliced PCR product revealed the predicted splice junction and confirmed the presence of an upstream exon capable of adding 9 aa to the N-terminus of the protein (supplementary Fig. S2D). This upstream exon and the splice donor and acceptor sites are conserved in the intact RAGL-A transposon from *P. flava*^H (supplementary File S1). We conclude that *PffRAG2L-A* is induced and undergoes mRNA splicing during early development and that an upstream exon encodes the N-terminal residues of the *PffRAG2L-A* protein.

The RAGL-A Family Is Also Present in 2 Echinoderm Classes

Public WGS and TSA database screening led to the identification of additional RAG1L/2L-A gene pairs in another invertebrate phylum—echinoderms—in both the Ophiuroidea class (spiny brittle star *O. spiculata* [*Osp*] and green brittle star *O. brevispina* [*Obr*]) and the Asteroidea class

(spiny starfish *M. glacialis* [Mgl]) (Fig. 1C). The 3 echinoderm *RAG1L-A* genes encode complete and potentially functional endonucleases while a complete *RAG2L-A* counterpart was found only in *O. spiculata* and *O. brevispina* (the identified *MglRAG2L-A* locus lacks coding information for a portion of the *RAG2L* N-terminus). Multiple putative *RAGL-A* loci were identified in another Asteroidea member (spiny starfish *Zoroaster* sp.), but all appear to be degraded, nonfunctional loci and were not analyzed further (supplementary File S1).

All identified echinoderm *RAG2L-A* genes are found downstream of *RAG1L-A* in the expected reverse orientation (Fig. 1C) and encode *RAG2L-A* proteins containing an AH between kelch and PHD domains, as in *PflRAG2L-A*. *OspRAG1L-A* (but not *OspRAG2L-A*) was also recently reported by Tao et al. (2022). The WGS database was also found to contain 2 additional *OspRAG2L-A* loci (*OspRAG2L-A.2*, *OspRAG2L-A.3*) that are apparently unlinked to a *RAG1L-A* counterpart and that encode proteins with a complete kelch-AH-PHD configuration (supplementary File S1).

Phylogenetic Analysis: *RAGL-A* Clusters with Jawed Vertebrate *RAG*

Phylogenetic analyses based on *RAG1(L)* catalytic core sequences found that *RAG1L-A* clusters with jawed vertebrate *RAG1* rather than with *RAG1L-B* (Fig. 2A, supplementary Fig. S3A and B), consistent with previous studies (Morales Poole et al. 2017; Martin et al. 2020; Tao et al. 2022). Bootstrap support of the *RAG1L* phylogenetic tree indicates a statistically significant separation between the *RAG1L-A* and *-B* families and *RAG1L-A* catalytic core sequences show greater sequence identity with *RAG1* (35% to 39%) than with *RAG1L-B* (26% to 33%) (supplementary Fig. S4A). Notably, echinoderm *RAG1L-A*

sequences consistently form an outgroup to hemichordate and jawed vertebrate *RAG1L-A* (Fig. 2A, supplementary Fig. S3A and B), with potential implications for the origin of jawed vertebrate *RAG* (see Discussion).

RAG2(L) phylogenetic trees (computed using blades 2 to 5 of the *RAG2(L)* kelch domain) are less reliable than those for *RAG1L* with lower bootstrap support at many branches (supplementary Fig. S3C and D), as expected given the weaker conservation of *RAG2(L)* compared to *RAG1(L)* (supplementary Fig. S4) (Morales Poole et al. 2017; Martin et al. 2020). *RAG2L* trees maintain the same overall structure as for *RAG1(L)* (supplementary Fig. S3C and D). The orphan *OspRAG2L-A.2* and *.3* sequences group well with *OspRAG2L-A.1* (which is paired with *OspRAG1L-A*) (supplementary Fig. S3) but the 3 sequences share only 42% to 52% identity in their most conserved region (kelch domain blades 2 to 5). This high divergence might be due to either the existence of different *A* subfamilies or an increased rate of change in the loci hosting the isolated *RAG2L-A* genes.

Notably, both the *P. flava* and *O. spiculata* genomes also contain *RAGL-B* transposons (supplementary File S1), which encode *RAG1L-B* proteins with 42% to 45% identity in the catalytic core region with *RAGL-B* in other species, but lower identity (30% to 31%) with their intraspecies *RAGL-A* counterparts (supplementary Fig. S4). This argues that the *RAGL-A* and *RAGL-B* transposon lineages evolved independently and in parallel in both the *P. flava* and *O. spiculata* genomes.

PflRAGL-A TIRs Are Chimeras of RSSs and *RAGL-B* TIRs

The TIRs elements of *Pfl^HRAGL-A* and *Pfl^TRAGL-A* are identical, indicative of the strong conservation of *RAGL-A* in the 2 populations, and exhibit a mixture of features of

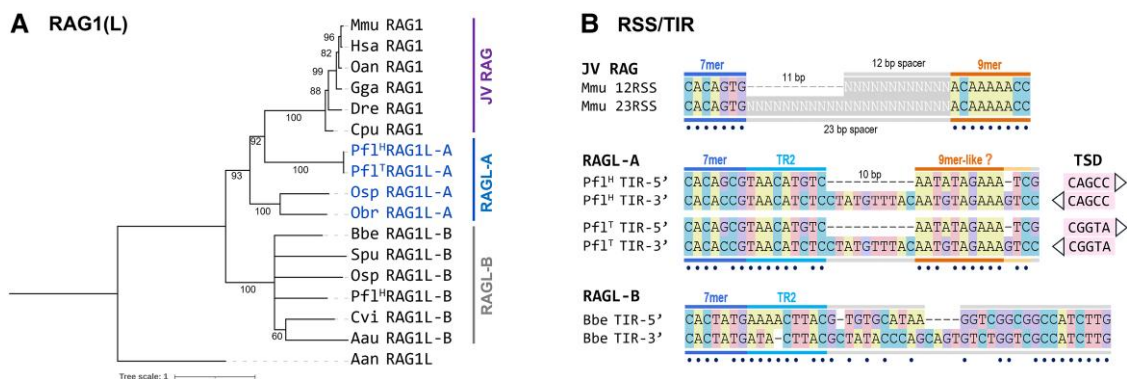


Fig. 2. Phylogeny and TIRs of *RAGL-A* elements. (A) Phylogeny trees of *RAG1/RAG1L* computed using Maximum Likelihood method using IQtree as described in Materials and Methods. Branch support (1,000 UFBoot replicates) is indicated next to each branch. (B) Comparison of selected RSSs and TIRs. Heptamer, TR2, spacer, and nonamer regions are indicated. Matching nucleotide pairs are depicted with dots. Target site duplication (TSD) pairs are shown next to the *PflRAGL-A* TIR sequences. Species abbreviations: Hsa, *Homo sapiens* (human); Mmu, *Mus musculus* (mouse); Oan, *Ornithorhynchus anatinus* (duckbill platypus); Gga, *Gallus gallus* (chicken); Dre, *Danio rerio* (zebrafish); Cpu, *Chiloscyllium punctatum* (shark); Pfl, *Ptychodera flava* (acorn worm); Osp, *Ophiothrix spiculata* (spiny brittle star); Obr, *Ophioderma brevispina* (green brittle star); Bbe, *Branchiostoma belcheri* (amphioxus); Spu, *Strongylocentrotus purpuratus* (purple sea urchin); Cvi, *Crassostrea virginica* (eastern oyster); Nge, *Notospermus geniculatus* (ribbon worms); Aau, *Aurelia aurita* (moon jellyfish); Aan, *Aureococcus anophagefferens* (alga). Analyses were performed on the most conserved core regions of *RAG1* (mouse: aa 388-1008).

jawed vertebrate RSSs and invertebrate RAGL-B TIRs (Fig. 2B and supplementary File S1). Like 12/23RSSs, the 5' and 3' TIRs contain a conserved heptamer closely resembling that of the consensus RSS separated from an AT-rich nonamer-like sequence by a “spacer” region with a 10 bp length asymmetry. Such asymmetry is not observed in BbeRAGL-B TIRs, which also lack an AT-rich nonamer-like sequence (Fig. 2B). Unlike 12/23RSSs, however, the “spacer” regions of the PflRAGL-A 5' and 3' TIRs are highly conserved for their first 9 bp, and in addition, the conserved region begins with a sequence rich in adenines. This conserved 9 to 10 bp region immediately adjacent to the heptamer, referred to as TIR region 2 (TR2), is observed in many deuterostome and protostome RAGL-B TIRs (Martin et al. 2020) and is important for DNA cleavage by BbeRAGL-B (Zhang et al. 2019). PflRAGL-A TIRs therefore display a hybrid heptamer-TR2-asymmetric spacer-nonamer organization.

The TIRs of several PflRAGL-A cassettes are flanked by 5 bp TSDs and these TSDs differ in sequence between *P. flava*^H and *P. flava*^T (Fig. 2B, supplementary Fig. S1), further supporting ongoing RAGL-A transposition activity in the *P. flava* genome after the divergence of the Hawaiian and Taiwanese populations. PflRAGL-A TSDs are GC-rich (82.5% in supplementary Fig. S1A) as is the case for TSDs generated by RAG, BbeRAGL-B, and Transib (Tsai et al. 2003; Kapitonov and Jurka 2005; Hencken et al. 2012; Huang et al. 2016), perhaps reflecting similar requirements for target site distortion during integration (Zhang et al. 2020).

In *O. spiculata*, the single RAG1L/2L-A gene pair identified was located close (within 400 bp) to 1 end of the sequence scaffold, preventing identification and proper validation of TIR elements, in contrast to *P. flava* where multiple transposon copies could be used to validate the transposon cassette margins. Similarly, no TIR elements were identified in the regions flanking RAGL-A gene pairs in the other echinoderm species despite the fact that for several loci, the sequence of substantial amounts of flanking genomic DNA is available.

RAG1L-A Proteins Display a Mixture of Invertebrate RAG1L-B and Jawed Vertebrate RAG1 Traits

In *P. flava*, *O. spiculata*, *O. brevispina*, and *M. glacialis*, RAG1L-A protein sequences display all known essential catalytic core domain components, including the catalytic DEDE tetrad and 4 Zn-coordinating residues that orchestrate folding of a zinc-binding domain (ZnB) that makes up much of the C-terminal portion of the catalytic core (Fig. 3A and B, and supplementary File S2). The proteins also contain 3 conserved cysteine-rich motifs (C1, C2, and C3) found in the RAG1 N-terminal region and all except OspRAG1L-A possess a RING-ZnF domain. Loss of the RING-ZnF domain has previously been reported in *Strongylocentrotus purpuratus* RAG1L-B and in several protostome RAGL-B proteins (Fugmann et al. 2006; Martin et al. 2020).

In contrast to RAG1L-B proteins, *P. flava* RAG1L-A exhibits 2 features characteristic of jawed vertebrate RAG1. The first is an NBD containing a GRPR/K motif (hereafter,

NBD_{GRPR/K}) (Fig. 3A and B). In RAG1, this motif makes direct contact with the A/T-rich portion of the nonamer and is required for RAG cleavage activity (Difilippantonio et al. 1996; Spanopoulou et al. 1996; Yin et al. 2009; Schatz and Swanson 2011). The presence of NBD_{GRPR/K} in PflRAG1L-A together with asymmetric TIRs containing an A/T-rich nonamer-like sequence is consistent with the possibility that PflRAGL-A mediates nonamer recognition by a mechanism similar to that of RAG. In echinoderm RAG1L-A, the corresponding sequence is GRPP₂ or GRRP₂, whose effect on DNA binding activity is difficult to predict and, in the case of GRPP₂, might compromise binding due to loss of electrostatic interactions between the R/K residue and the DNA backbone (Yin et al. 2009). BbeRAGL TIR elements are almost symmetrical in length and lack an adenine-rich nonamer region, which is mirrored by the fact that the BbeRAG1L NBD-equivalent domain lacks the GRPR/K motif and makes only a modest contribution to cleavage activity (Zhang et al. 2019).

The second feature shared uniquely between PflRAG1L-A and RAG1 is glutamate at the position corresponding to mouse RAG1 E649 (Fig. 3A and B). DNA cleavage by RAG occurs in a synchronous, or “coupled” fashion and only when both of its substrates are bound in the synaptic complex. The enforcement of coupled cleavage is dependent on residue E649, which is thought to exert its influence in part through hydrogen bond formation with RAG1 S963 (Kriatchko et al. 2006; Zhang et al. 2019). The E649/S963 pair, highly conserved in jawed vertebrate RAG1, is E/K in PflRAG1L-A and Q/H in Osp, Obr, and Mgl RAG1L-A proteins (Fig. 3B). These residues possess bulky side chains that could engage in hydrogen bonds and interact electrostatically, consistent with the possibility that PflRAGL-A and echinoderm RAGL-A possess some degree of coupled cleavage.

RAG1 E649 also suppresses RAG-mediated transposition modestly (approx. 2-fold) in a manner that is independent of S963 (Zhang et al. 2019), and transposition by PflRAGL-A and echinoderm RAGL-A might similarly be modestly downregulated by the E and Q residues, respectively, they possess at this position (Q shares many physicochemical properties with E). A second highly conserved residue in RAG1, arginine 848, also contributes to transposition suppression, but in this case, almost completely eliminates transposition activity (Zhang et al. 2019). This residue is a hydrophobic aa, most often methionine, in invertebrate RAG1L proteins, and all identified RAG1L-A proteins retain the transposition-permissive M at this position (Fig. 3A and B).

Interestingly, the C-terminal tails (CTTs) of the hemichordate and echinoderm PflRAG1L-A proteins possess a conserved CX₂CX₃GHX₄C motif (CCGHC motif hereafter) (Fig. 3A and B), found in nucleic acid-binding “zinc-knuckle” domains (De Guzman et al. 1998; Klein et al. 2000). The CCGHC motif is present in the CTTs of virtually all invertebrate RAG1L proteins but not jawed vertebrate RAG1 (where CTT is acidic). In BbeRAG1L-B, CTT_{CCGHC} is a DNA binding “clamp” that interacts with

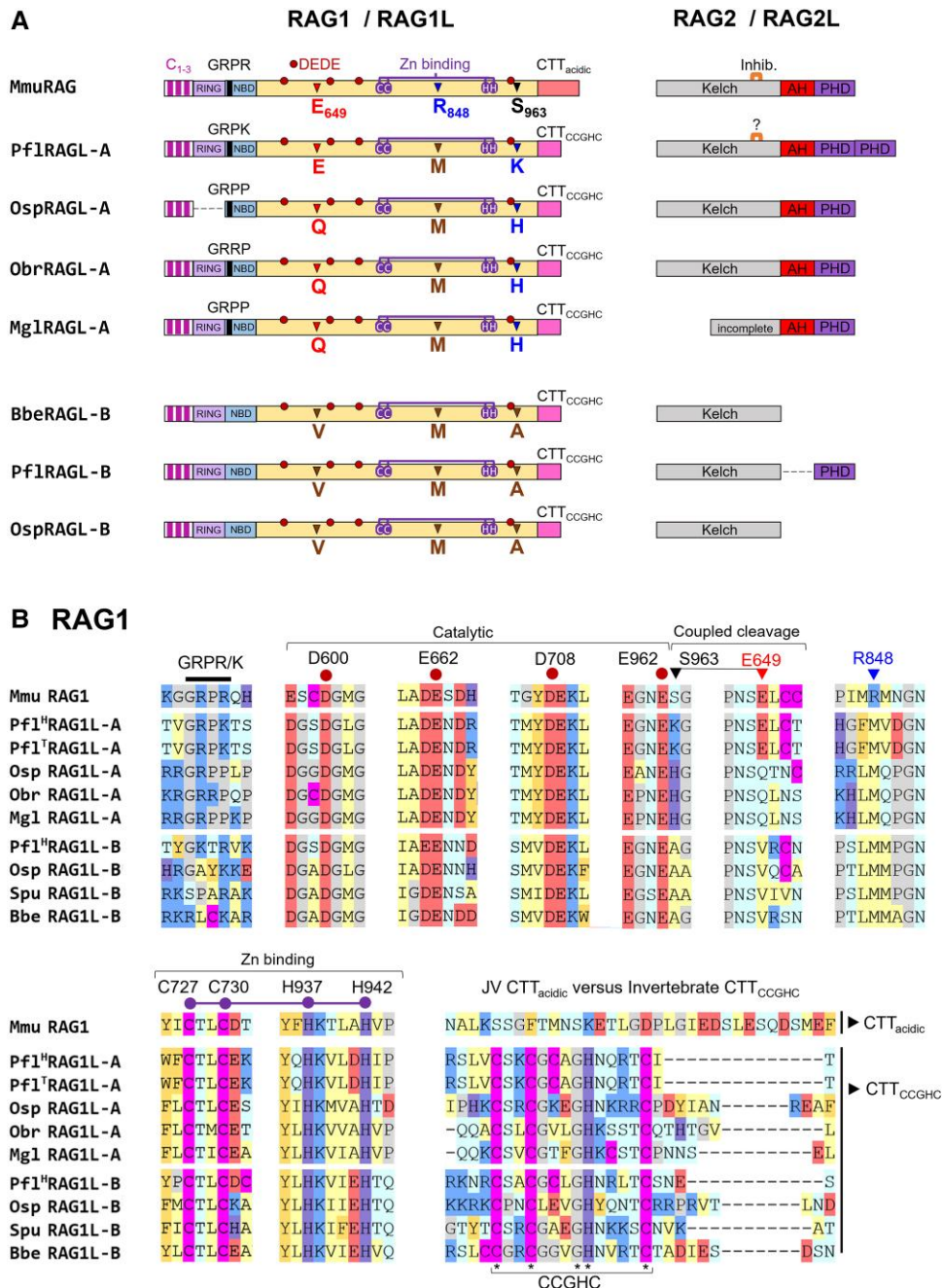


Fig. 3. RAG1L-A and RAG2L-A proteins. (A) Domain organization of RAG1(L)–RAG2(L) pairs from the most preserved copies of RAGL-A and RAGL-B from *P. flava*, *O. spiculata*, *O. brevispina*, and *M. glacialis*, mouse (Mmu) RAG, *B. belcheri*. The RAG1(L) conserved catalytic DEDE tetrad and Zinc-binding residues are depicted with circles, while regulatory adaptations E649, R848, and S963 identified in mouse RAG1 are indicated with triangles. Kelch, RAG2 kelch domain encoding blades 1 to 6; AH, acidic hinge; PHD, plant homeodomain; Inhib, B6–C6 transposition inhibitory loop in RAG2. Species abbreviations are as in the legends of Fig. 2 and supplementary Fig. S3. (B) Sequence conservation and variation at key positions and motifs in RAG1(L) proteins.

TR2 and is critical for cleavage activity (Zhang et al. 2019). The presence of a GRPR/K-containing NBD, asymmetric TIRs with an AT-rich nonamer-like region, and CTT_{CCGHC} in Pf1RAG1L-A suggests competing, and potentially redundant, modes of DNA binding (see Discussion).

We used artificial intelligence-based methods (AlphaFold, OmegaFold) and homology 3D modeling pipelines to compute predictive structural models of the core regions of Pf1RAG1L-A and OspRAG1L-A, revealing that

they can adopt structures similar to those of RAG1 and BbeRAG1L-B and that their predicted DNA binding surfaces exhibit strikingly similar charge distributions to those of RAG1/BbeRAG1L-B (supplementary Fig. S5). This observation suggests substantial parallels between the mechanisms of DNA engagement by RAG1L-A proteins and that by RAG1 and BbeRAG1L-B.

In summary, Pf1RAG1L-A and Osp/Obr/Mg1RAG1L-A are potentially catalytically active proteins possessing

features distinctive of jawed vertebrate RAG1 (GRPR/K or related motif, E/Q649) as well as features found previously only in invertebrate RAG1L-B (M848, CTT_{CCGHC}).

RAG2L-A Proteins Contain an Acidic Hinge Capable of Suppressing Transposition

A unique feature of RAG2L-A sequences, not observed in any other invertebrate RAG2L protein reported thus far, is an acidic region between the kelch and PHD domains. The RAG2L-A AH is somewhat longer (66 to 78 aa) than that of jawed vertebrate RAG2 (52 to 62 aa) and has nearly as high a density of acidic residues (34% to 40%) as in jawed vertebrate RAG2 AHs (39% to 44%) (Fig. 4A). The distribution of D/E residues shows some similarities between the different AHs, including a region of high D/E density near the AH C-terminus, but other sequence similarities are not apparent.

The mouse and human RAG2 AHs suppress transposition potently (>50-fold) (Zhang et al. 2019). To test whether invertebrate RAG2L-A AH sequences also possess transposition inhibitory activity, we appended the PflRAG2L-A or OspRAG2L-A.1 AH to the mouse RAG2 core region (Fig. 4B) and assayed the resulting chimeric

proteins for transposition activity together with mouse RAG1 bearing transposition activating mutations R848M and E649V. The results demonstrate that the *P. flava* and *O. spiculata* AH regions potently suppress transposition activity, to an extent equivalent or nearly equivalent to that observed with a partial (aa 351-387) or full (aa 351-418) mouse AH (Fig. 4C). An assay for recombination confirmed that all of the RAG2 fusion proteins tested support DNA cleavage, though as expected (Lu et al. 2015; Zhang et al. 2019), appending any AH region reduced recombination somewhat (Fig. 4D). We conclude that AH domains with transposition suppressive potential exist in RAG2L-A proteins in hemichordates and echinoderms and that the AH is not a jawed vertebrate-specific adaptation.

The RAG2 Transposition Inhibitory Loop Is Partially Restored in PflRAG2L-A

Jawed vertebrate RAG2 contains a 10 aa loop between β -strands B and C of blade 6 of the kelch domain (Fig. 5A). This B6–C6 loop and the B6 GG motif (a glycine doublet that is characteristic of beta strand B) are shifted out of the plane defined by the other 5 blades; this is not

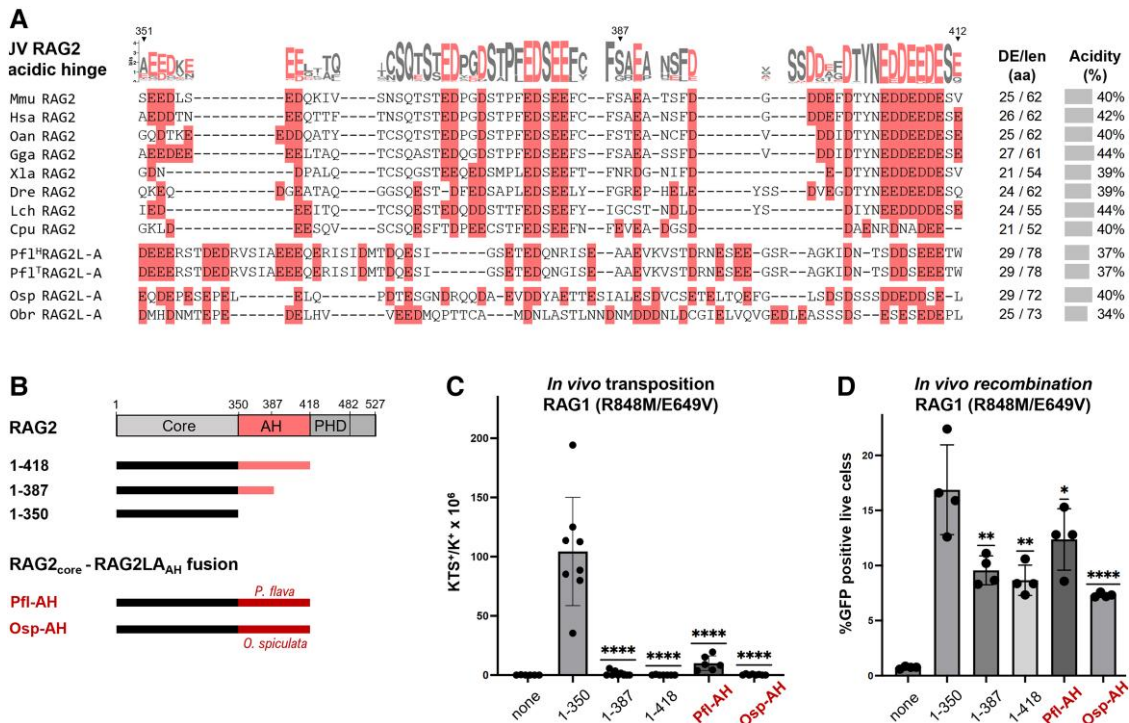


Fig. 4. RAG2L-A acidic hinge sequence and transposition inhibitory activity. (A) Comparison of the acidic hinge (AH) of jawed vertebrate (JV) RAG2 and RAG2L-A sequences in hemichordates (*P. flava*) and echinoderms (*O. spiculata* and *O. brevispina*), with acidic residues Asp and Glu highlighted. Overall length, number of acidic residues, and percent acidic residues shown at right. Conservation profile of jawed vertebrate RAG2 was computed as KL divergence as described in Materials and Methods (letter height is proportional to conservation). (B) Schematic diagrams of RAG2 proteins tested for activity. Mouse RAG2 is shown at top. Mouse RAG2 core alone or fused to either its own AH or that of PflRAG2L-A or OspRAG2L-A were tested. (C, D) In vivo transposition (C) and recombination (D) assays performed in human 293 cells upon expression of full length mouse RAG1 containing transposition activating mutations R848M and E649V and the indicated RAG2 fusion proteins (diagramed in panel B). Data points are biological replicates derived from independent experiments. Statistical significance calculated compared to RAG2 1-350 using 2-tailed T test ($P < 0.05$ (*), < 0.01 (**), < 0.0001 (****)). Species abbreviations are as in the legends of Fig. 2 and supplementary Fig. S3.

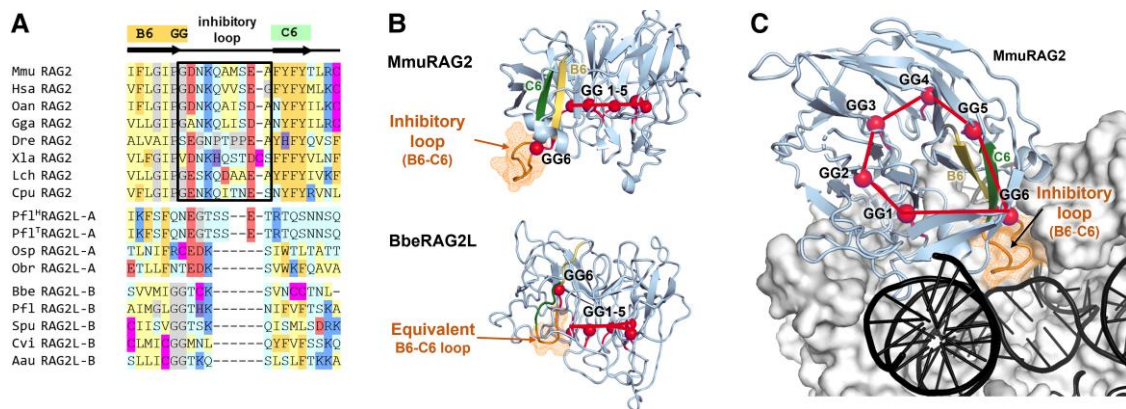


Fig. 5. RAG2 kelch blade 6 transposition inhibitory loop. (A) Sequence of jawed vertebrate RAG2 loop between β -strands B6 and C6 compared to the equivalent region in RAG2L-A and RAG2L-B proteins. GG, a double glycine motif (sometimes PG) frequently present at the end of β -strand B in kelch domain blades. (B) Structural differences in blade 6 as observed in the cryo-EM structures of mouse RAG2 and BbeRAG2L. The GG motif from blade 6 (GG6) shifts in opposite directions in the 2 proteins with respect to the plane generated by the rest of the GG motifs (GG1-5) (PDB: 6XNY, 6B40). (C) Top view of mouse RAG strand transfer complex illustrating mouse RAG2 structure (ribbon model) and the downward projection of the sixth blade and the B6–C6 loop (mesh representation) to make contact with target DNA (black). RAG1 shown as gray surface. (PDB: 6XNY). Species abbreviations are as in the legends of Fig. 2 and supplementary Fig. S3.

the case in BbeRAG2L-B (Fig. 5B). The RAG2 B6–C6 loop makes contact with target DNA in the RAG target capture and strand transfer complexes (Fig. 5C) (Chen, Cui, Wang 2020; Zhang et al. 2020) and deletion of 4 amino acids at the loop tip increases RAG-mediated transposition 2- to 3-fold (Zhang et al. 2020). While the corresponding loop is only 5 aa long in RAG2L-B proteins and OspRAG2L-B and ObrRAG2L-B, it is 8 aa long in PflRAG2L-A (Fig. 5A).

PflRAG2L-A Encodes a Double PHD

The jawed vertebrate RAG2 PHD contains 2 zinc fingers that create a pocket that binds the N-terminal tail of histone 3 when lysine 4 is trimethylated (H3K4me3) (Liu et al. 2007; Matthews et al. 2007). Pfl, Osp, and ObrRAG2L-A sequences contain a PHD abutting the AH (hereafter PHD₁) in which the zinc-coordinating cysteine and histidine residues are readily identifiable, with the exception of OspRAG2L-A.3 which lacks the final 2 cysteine residues. The proteins also contain a highly conserved tryptophan residue (W453 in mouse RAG2) required for methylated lysine binding (Matthews et al. 2007) (Fig. 6A, supplementary File S2B). Despite this similarity with jawed vertebrate PHD₁, RAG2L-A PHD₁ more closely resembles PHD₁ of RAG2L-B in possessing an additional conserved cysteine residue (* in Fig. 6A) and aa changes at conserved methyl-lysine-binding residues of jawed vertebrate PHD₁ (Y415→W or L and M443→W (Ramon-Maiques et al. 2007) (Fig. 6A). Hence, the histone recognition profile of RAG2L-A PHD₁ might differ from that of vertebrate RAG2 PHD₁, consistent with the finding that PHD₁ of SpuRAG2L-B preferentially binds H3K4me2 instead of H3K4me3 (Wilson et al. 2008).

Unlike any RAG2(L) protein described to date, PflRAG2L-A contains a second complete PHD (hereafter, PHD₂) immediately following PHD₁, that, like PHD₁, can encode CCHC and CCCC zinc fingers (fingers 3 and 4 in

Fig. 6B). The double PHD₁–PHD₂ domain of PflRAG2L-A has limited sequence identity with 2 groups of structurally related double PHDs: (i) the double PHDs of histone acetyltransferase MOZ, chromatin remodeling complex subunit DPF3, and histone-lysine N-methyltransferase KMT2C (Zeng et al. 2010; Xiong et al. 2016), and (ii) the double PHDs of the histone-lysine N-methyltransferases NSD1 and NSD3 (He et al. 2013; Berardi et al. 2016) (Fig. 6B and C). While sequence homology between the 2 double PHD groups is largely limited to their zinc-coordinating C/H residues, their 3D architectures bear striking similarities, including their 4 zinc-binding pockets with similar C/H patterns (Fig. 6B and C). The C/H pattern in the double PHD of PflRAG2L-A precisely matches that of the MOZ/DPF3/KMT2C group, and modeling suggests that it could plausibly adopt a similar structure (Fig. 6C). We note that SpuRAG2L-B might contain a degenerate PHD₂ (Fig. 6B).

Discussion

The absence of identified RAGL-A transposons had left a substantial gap in our understanding of the evolutionary history of jawed vertebrate RAG, in particular, numerous uncertainties as to the order in which various domains/residues were gained and lost and whether key functional adaptations occurred prior to or after jawed vertebrate speciation. Our identification of complete RAGL-A elements in hemichordates and echinoderms helps clarify these issues and allows for a more nuanced description of the evolution of regulated DNA binding, DNA cleavage, and transposition activities of the RAG recombinase.

It is plausible that the intact RAG1L/2L-A open reading frames in *P. flava*, *O. spiculata*, and *O. brevispina* encode active DNA endonucleases, and in the case of *P. flava*, where the open reading frames are flanked by TIRs, that they

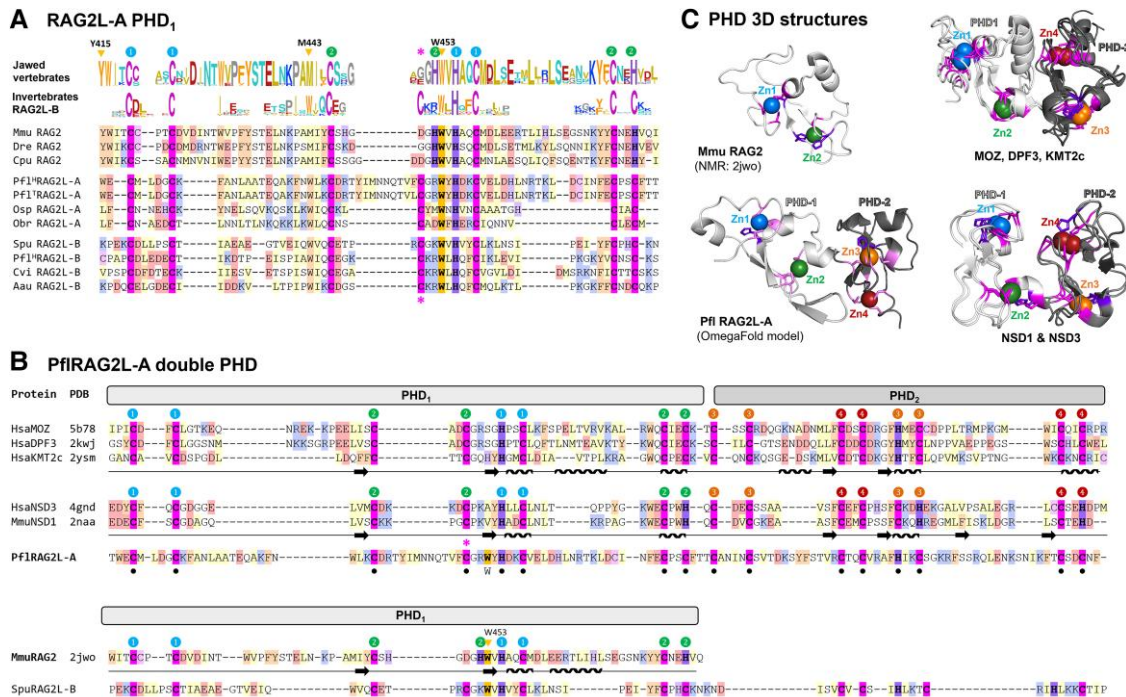


Fig. 6. RAG2L-A PHD region. (A) Comparison of PHD finger 1 (PHD₁) of RAG2L-A with other invertebrate RAG2L and jawed vertebrate RAG2 proteins. The Cys/His pattern is highlighted in magenta and purple, respectively with numbers indicating residues that coordinate zinc atoms 1 and 2. Triangles, conserved residues in jawed vertebrate RAG2 important for binding of H3K4me3; Asterisk, Cys residue conserved in RAG2L but not RAG2 proteins. Sequence conservation profile displayed above the alignment was computed as described in Fig. 4. (B) Sequence comparison of the double PHD region (PHD₁ and PHD₂) of PflRAG2L-A with 2 groups of double PHD domain proteins for which experimental structures are available and which display similar Cys/His patterns. Colored labels 1 to 4 indicate the Zn binding topology. The mouse single PHD and the incomplete extended PHD pattern of *S. purpuratus* RAG2L-B are shown below for comparison. Secondary structure elements indicated by the experimental 3D structures are shown below the amino acid sequence (arrow, β -strand, wavy line, α -helix). (C) Structural comparison of double PHD domains. 3D structures of mouse RAG2 single PHD (upper left, PDB: 2JWO), *P. flava* RAG2L-A double PHD (lower left, predicted model using OmegaFold), superimposition of double PHD domains of MOZ, DPF3, and KMT2c (upper right; PDB: 5B78, 2KWJ), 2YSM, respectively) and NSD1 and NSD3 (lower right; PDB: 4GDN, 2NAA, respectively). Zinc-binding Cys/His residues are displayed in magenta/purple, while Zn ions are colored as in labels in panels A and B.

remain active transposons. Recent transposition activity of *PflRAGL-A* is supported by the presence of multiple *RAGL-A* copies in both the Hawaiian and Taiwanese populations of *P. flava* and by the different TSDs and chromosome regions that flank them. *PflRAGL-A.14* appears to be present in some *P. flava*^T worms and not others; its allele frequency and propensity for mobilization remain to be explored.

While additional analyses might reveal TIRs flanking echinoderm *RAGL-A* elements, their apparent absence parallels the structure of the first *RAG1L-RAG2L* gene pair to be identified in invertebrates (in the purple sea urchin *S. purpuratus*) (Fugmann et al. 2006). Indeed, to our knowledge, no intact *RAGL* element predicted to be capable of transposition has been identified in echinoderms despite the identification of *RAG1L-RAG2L* loci in multiple echinoderm species (this report and Fugmann 2010, Kapitonov and Koonin 2015, Morales Poole et al. 2017, Martin et al. 2020, Tao et al. 2022, and Yakovenko et al. 2022), and despite the identification of potentially active *RAGL* transposons in multiple lineages of deuterostomes, protostomes, and cnidarians, and even in a protist (Huang et al. 2016; Morales Poole et al. 2017; Martin et al. 2020; Tao et al.

2022). In contrast, *Transib* transposons with apparently intact TIRs are present in numerous species of echinoderms (our unpublished data and Kapitonov and Jurka 2005, Kapitonov and Koonin 2015, and Tao et al. 2022), suggestive of different selective pressures acting on *Transib* and *RAGL* transposons in this clade. Our findings indicate that *RAG1L/2L-A* loci might represent domesticated transposases performing novel functions for their echinoderm hosts, as has been proposed for *RAGL-B* loci in sea urchins (Fugmann et al. 2006; Yakovenko et al. 2022). Our findings also argue that *RAGL-A* and *RAGL-B* transposons evolved side-by-side in both hemichordates and echinoderms over extended evolutionary periods.

A particularly striking feature of the *RAGL-A* elements reported here is their chimeric nature, with the *RAGL-A* proteins and their TIRs exhibiting features distinctive of jawed vertebrate *RAG/RSSs* and of *RAGL-B* proteins/TIRs (summarized in Fig. 7A). Phylogenetic analyses demonstrate that invertebrate *RAGL-A* proteins are more closely related to jawed vertebrate *RAG* than are *RAGL-B* proteins (Fig. 2A and supplementary Fig. S4). Our findings have implications for the evolution of *RAG*'s RSS substrates and its DNA binding, DNA cleavage, and transposition activities.

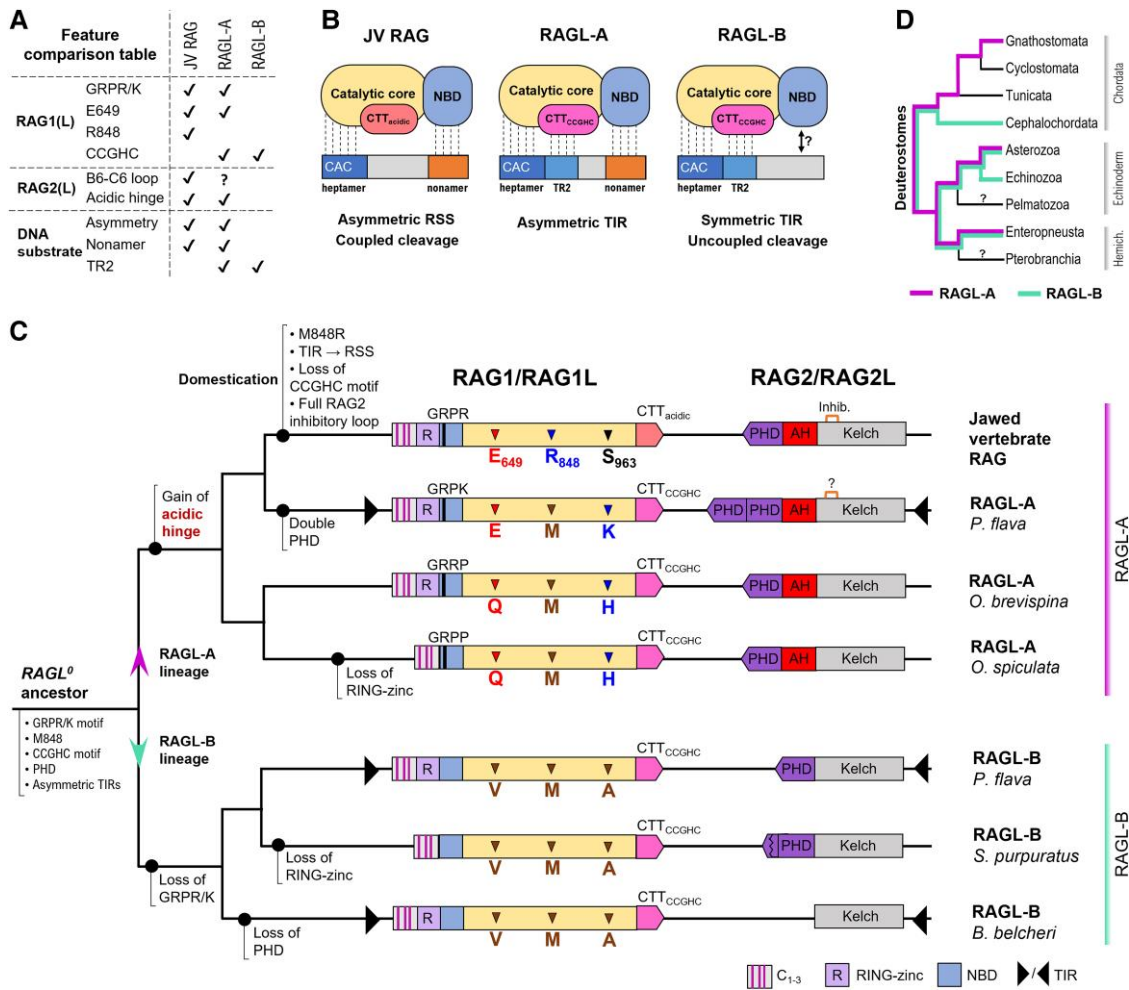


Fig. 7. Features, DNA binding modes, and model for evolution of RAG(L) proteins. (A) Feature comparison table of functionally important elements of the different RAG(L) lineage proteins and DNA substrates. (B) Differences in DNA binding modes between jawed vertebrate RAG and BbeRAGL-B and hypothesized DNA binding modes available to RAGL-A proteins. (C) Model for the evolution of RAGL-A, RAGL-B, and RAG. Beginning with $RAGL^0$, the presumed ancestral RAGL transposon, evolutionary events leading to the gain or loss of traits and changes in domain architecture are depicted. See text for additional details. RAG1(L) and RAG2(L) proteins for the indicated species are diagrammed in the tail-to-tail configuration observed for their respective genes. C_{1-3} , domain containing 3 cysteine pairs; R, RING-zinc finger domain; GRPR/K, motif involved in DNA binding at the N-terminus of the NBD; CTT_{acidic} , acidic C-terminal tail; CTT_{CCGHC} , C-terminal tail containing DNA binding domain with conserved CCGHC residues; Kelch, kelch domain constituting RAG2 core; AH, acidic hinge; PHD, plant homeodomain; black triangles, TIRs. For *ObrRAGL-A* and *OspRAGL-A*, TIRs have not been identified and hence are not shown, though for *OspRAGL-A*, the available DNA sequence assembly does not allow a definitive conclusion on this issue. Tree structure reflects the evolution of protein domains and features and is not meant to represent species phylogeny. (D) Schematic evolutionary tree of deuterostomes depicting the branches where RAGL-A and RAGL-B elements are found. Length of tree branches does not indicate degree of relatedness.

Evolution of RAG DNA Binding and RSS Substrates

RAG and BbeRAGL-B bind substrate DNA through distinct modes: while RAG is completely dependent on the RAG1 NBD for activity and CTT_{acidic} is dispensable, BbeRAGL relies heavily on BbeRAG1L CTT_{CCGHC} for activity and its NBD (which lacks a GRPR/K motif) makes a smaller contribution (Zhang et al. 2019) (Fig. 7B). To our knowledge, PflRAG1L-A is the first RAG1L protein to be described that contains both $NBD_{GRPR/K}$ and CTT_{CCGHC} . As a result, it would be predicted to be capable of both modes of DNA engagement (Fig. 7B), a conclusion supported by the presence of conserved TR2 and AT-rich nonamer-like sequences in its TIRs. How the PflRAGL-A enzyme might have exploited the availability of 2 DNA binding modes is

unknown. Biochemical experiments with chimeric RAG1–BbeRAG1L proteins indicate that at least in the context of RAG1 and BbeRAG1L, $NBD_{GRPR/K}$ and CTT_{CCGHC} are functionally redundant, with each rendering the other dispensable (Zhang et al. 2019). We hypothesize that such redundancy would cause a RAGL transposon harboring both domains to be prone to loss of either $NBD_{GRPR/K}$ and the nonamer (as in RAGL-B transposons) or CTT_{CCGHC} and TR2 (as in RAG and the RSS), and that such evolutionary instability might explain the dearth of RAGL enzymes, such as PflRAGL-A, that contain both $NBD_{GRPR/K}$ and CTT_{CCGHC} . Loss of CTT_{CCGHC} /TR2 by RAG/RSS likely facilitated the evolution of RAG's capacity for nuanced and flexible DNA recognition that enables it to cleave

heterogeneous RSS sequences at widely varying efficiencies (Ramsden et al. 1994; Yu et al. 2002; Feeney et al. 2004; Swanson 2004; Gopalakrishnan et al. 2013; Kim et al. 2018; Wu et al. 2020). Overall, our data strengthen the argument for co-evolution between RAG(L) DNA binding domains and their respective DNA recognition sites.

Our findings and those of other recent studies (Martin et al. 2020; Tao et al. 2022) demonstrate that RAGL transposons (of both the A and B families) can possess TIRs with a length asymmetry similar to that of 12/23 RSSs. In addition, our findings with PflRAGL-A indicate that jawed vertebrate asymmetric RSSs can readily be explained as arising from the TIRs of a RAGL-A transposon—an idea that contradicts the “Transib seed” hypothesis that proposed that RSSs arose directly from a *Transib* element (Yakovenko et al. 2021).

Evolution of Coupled Cleavage

Hydrogen bond formation between RAG1 E649 and S963 helps enforce coupled cleavage by mouse RAG, likely by regulating the structure of an α -helix containing active site residue E962, which in turn determines the integrity of the active site (Kriatchko et al. 2006; Zhang et al. 2019). While RAG1L-A proteins possess position 649/963 aa pairs (E/K or Q/H) that could serve a similar function, this aa pair is V/A in BbeRAG1L-B (Fig. 3B), which lacks hydrogen bond potential. Consistent with this, cleavage by BbeRAGL-B in biochemical assays is uncoupled (Zhang et al. 2019), though the behavior of the enzyme in its natural in vivo environment is not known. More generally, invertebrate RAG1L proteins of the B, C, and D families display a non-charged, hydrophobic residue (Val, Ile, Thr) at the position equivalent to E649 and a small amino acid (Ala, Gly, Ser, Cys) at the position equivalent to S963 (Martin et al. 2020). We hypothesize that coupled cleavage activity (and the necessary aa 649/963 pair) arose in a RAGL-A transposon prior to speciation of jawed vertebrates, thereby helping to ensure that TIR cleavage occurred in a coordinated fashion in a synaptic complex. Regulated cleavage within an organized synaptic complex is a common feature of bacterial and eukaryotic transposases and site-specific recombinases (Craig 2015) and might have provided a selective advantage to the RAGL-A transposon and/or its host by reducing the incidence of uncoupled DNA double strand breaks. An alternative scenario is that coupled cleavage was a property of the ancestral RAGL transposon and was subsequently lost in the RAGL-B lineage and retained in RAGL-A.

Evolution of a RAG Recombinase Lacking Transposase Activity

We previously proposed that suppression of RAG transposition activity began in the jawed vertebrate lineage after creation of the initial “split” antigen receptor gene (Liu et al. 2022). Based on the findings reported here, this hypothesis needs to be reconsidered. Hemichordate and echinoderm RAG2L-A proteins contain an AH of

approximately the same size and acidic amino acid content as the AH of jawed vertebrate RAG2 (Fig. 4A), and the Pfl and Osp RAG2L-A AH regions potently suppress RAG-mediated transposition when attached to the core region of mouse RAG2 (Fig. 4B to D). This is a striking finding given that there is little sequence similarity between the AHs of RAG2L-A and RAG2 and suggests that inhibition of transposition depends more on acidic amino acid content than on specific sequence motifs. A similar conclusion was reached regarding the sequence features of the RAG2 AH required to influence repair pathway choice in the post-cleavage phase of V(D)J recombination (Coussens et al. 2013).

PflRAG2L-A possesses a B6–C6 connecting loop that is intermediate in size (8 aa) between that in RAG2L-B (5 aa) and that in jawed vertebrate RAG2 (10 aa). It is not known whether 8 aa is sufficient for transposition inhibition, but regardless, the PflRAG2L-A loop might be indicative of an intermediate in the evolution of the 10 aa loop of RAG2, which suppresses transposition 2- to 3-fold (Zhang et al. 2020).

Together, these findings indicate that at least 1 adaptation with the potential to suppress transposition arose in invertebrate RAGL-A transposases. Such adaptations, while reducing the mobility of the transposon, might have rendered them less damaging and more readily tolerated by their hosts—a common theme for transposons generally (Lohe and Hartl 1996; Davies et al. 2000; Levin and Moran 2011; Saha et al. 2015; Almeida et al. 2022). Only the change from methionine to arginine at RAG1 position 848 now appears to be a transposition-suppressive adaptation specific to jawed vertebrates. The suppressive effect of the M848R mutation is very strong, particularly in vivo, and appears to be due to the ability of methionine, but not arginine, to induce bends in target DNA needed for binding in a deep pocket in the RAG enzyme (Zhang et al. 2020). We refer to R848 as the “gatekeeper” residue for the regulation of RAG-mediated transposition to reflect both its potency and the likelihood that acquisition of arginine at this position was a pivotal event in RAG evolution that helped usher in the transition from RAGL transposase to RAG recombinase.

Evolution of RAG Chromatin Binding

Through its ability to bind H3K4me3, the RAG2 PHD finger plays an important role in specifying sites of chromatin binding by RAG (Teng et al. 2015; Maman et al. 2016) and in increasing RAG DNA cleavage activity, apparently by inducing allosteric changes in RAG1 (Shimazaki et al. 2009; Lu et al. 2015; Bettridge et al. 2017). Previous studies have established that the PHD was an early feature of RAG2L proteins, was lost in certain lineages (e.g. BbeRAG2L-B and perhaps amphioxus more generally (Braso-Vives et al. 2022)), and retains the ability to bind methylated lysine in the one case examined (H3K4me2 binding by SpuRAG2L-B) (Fugmann et al. 2006; Wilson et al. 2008; Huang et al. 2016; Martin et al. 2020; Tao et al. 2022). Our findings extend these observations by demonstrating that the PHD is a component of RAG2L-A

proteins and by revealing an unexpected “double PHD” in PflRAG2L-A. While certain residues important for methylated lysine binding are intact in RAG2L-A PHDs, other are not and it is difficult to predict whether RAG2L-A PHDs possess histone tail binding activity. The PHD₁–PHD₂ of PflRAG2L-A might adopt a structure similar to that of the double PHDs of several chromatin associated proteins (Fig. 6C) and it is plausible that it is capable of recognition of some form of acylated lysine. RAG1, very likely through the action of its NTR (aa 1–383), is also able to influence RAG binding in the genome (biasing binding to regions enriched in H3K27ac) (Maman et al. 2016). Some RAG1L proteins of both the A and B families contain all of the major elements found in the RAG1 NTR and hence might also contribute to chromatin binding.

An Updated Model of RAG/RAGL Evolution

Based on our findings, we propose the following refined model for RAG/RAGL evolution (Fig. 7C). It is thought that *Transib*, which was recently identified in bacteria (Tao et al. 2022), preceded RAGL in evolution and that the first RAGL transposon (which we designate RAGL⁰) was created early in eukaryotic evolution when a RAG2L gene became incorporated into a *Transib* transposon (Carmona and Schatz 2017; Liu et al. 2022; Tao et al. 2022). While the features of RAGL⁰ are unknown, we suggest that it contained asymmetric TIRs, CTT_{CCGHC}, and NBD_{GRPR/K}, all of which have been identified in *Transib* transposons (Kapitonov and Jurka 2005; Zhang et al. 2019; Tao et al. 2022). RAGL⁰ likely also contained a transposition-permissive aa (e.g. methionine) at the position corresponding to RAG1 848 and a PHD at the RAG2L C-terminus. In an early metazoan, RAGL⁰ gave rise to RAGL-B transposons, the defining feature of which was loss of the GRPR/K motif in the NBD and a strong reliance on CTT_{CCGHC} for DNA binding. RAGL-B was evolutionarily successful, being transmitted, primarily by vertical transmission, into numerous metazoan lineages including cnidarians, protostomes, and deuterostomes (Martin et al. 2020; Tao et al. 2022).

RAGL⁰ also gave rise to RAGL-A through acquisition of an AH in RAG2L and possibly E/Q at the RAG1L position equivalent to RAG1 649. A RAGL-A transposon, perhaps closely resembling *PflRAGL-A*, subsequently found its way into the jawed vertebrate lineage where it created the initial split antigen receptor gene and underwent several adaptations to facilitate its domestication: loss of CTT_{CCGHC} (creating strong reliance on NBD_{GRPR/K} and the RSS nonamer), the M848R gatekeeper mutation to terminate transposition, and acquisition of the full B6–C6 RAG2 inhibitory loop. As discussed previously (Zhang et al. 2019; Liu et al. 2022), suppression of RAG transposition activity involved multiple adaptations, perhaps because it was difficult to achieve the level of suppression needed to protect genome integrity in organisms in which millions of V(D)J recombination events occur each day. Acquisition of the AH was unlikely to have fully suppressed transposition activity given the evidence provided here for RAGL-A

activity in *P. flava* and the requirement for a transposition event in an early jawed vertebrate to create the first split antigen receptor gene.

Questions of particular interest raised by our data and this model are the timing of the emergence of RAGL-A and whether RAGL-A entered the jawed vertebrate lineage by vertical or horizontal transmission. Vertical inheritance would predict that RAGL-A arose in an early deuterostome prior to the divergence of hemichordates and echinoderms, about 560 million yrs ago (dos Reis et al. 2015). However, the “spotty” evolutionary distribution of RAGL-A elements—present in hemichordates and echinoderms but absent from tunicates, cephalochordates, and jawless vertebrates (Fig. 7D)—is suggestive of horizontal gene transfer (HGT). Phylogenetic clustering of jawed vertebrate RAG1 and hemichordate RAG1L-A (Fig. 2A) is consistent with this idea and with HGT of RAGL-A between a hemichordate and an early jawed vertebrate. Genome sequence data from additional species, particularly deuterostomes, might help address the possibility of HGT and provide an explanation for the absence of TIRs flanking RAGL elements in echinoderms.

Limitations of the Study

Our findings do not directly address the functionality of the identified RAGL-A proteins, preventing firm conclusions regarding their ability to perform DNA cleavage or other enzymatic functions. While evidence exists and is presented in our study for mRNA expression for some RAGL-A genes, endogenous RAGL-A protein expression has not been assessed. Aspects of our model for RAG(L) evolution might need to be revised as the genome sequences of more invertebrate organisms are reported.

Materials and Methods

Genomic Analysis

The genomic and transcriptomic public repositories (WGS/TSA) of all metazoan species were screened using tblastn (Gertz et al. 2006; Johnson et al. 2008) starting from the RAG1L-A sequence identified in *P. flava* (Morales Poole et al. 2017; Martin et al. 2020). The flanking regions of RAG1L-A loci were inspected for RAG2-like signatures and positive results were further added to the screening process. Putative RAG1L/2L-A loci identified in this way were then subjected to predictions of protein translation using Augustus (Stanke et al. 2008) and Softberry FGENESH+ (Solovyev et al. 2006) and classified as either complete or pseudogenized RAGL sequences depending on the completeness and compliance with canonical RAG1/2 domain organization.

The presence of TIRs, defining the margins of the transposon cassette, was investigated using a homology variation procedure described in Martin et al. (2020). Briefly, the regions flanking the RAGL loci were aligned to identify the transposon ends, based on an expectancy of higher sequence homology within the cassettes as compared to that of the insertion loci. The cassette end predictions

were further scrutinized based on compliance with the expected TIR consensus within the heptamer region and the presence of 5 bp TSDs.

BAC Screening, Selection, and Sequencing

Putative *RAGL1_A* containing *P. flava* clones were identified by screening a *P. flava* Bacterial Artificial Chromosome (BAC) library made from the Hawaiian population (Arshinoff et al. 2022). The presence of *RAG1L-A* was confirmed via colony-PCR followed by Sanger sequencing. The selected BACs were recovered and sequenced on an Illumina MiSeq using reagent kit Index Nextera XT kit V2-500 cycles and assembled with the CLC Genomics Workbench 20.0 v7.5 (QIAGEN). *P. flava* RAGL sequences of the Taiwanese population were retrieved from genome sequence deposited in GenBank (BioProject PRJNA747109).

PfRAGL^T-A.14 was identified using genomic DNA prepared from a *P. flava^T* isolate provided by Dr. Jr-Kai Yu (Institute of Cellular and Organismic Biology, Academia Sinica, Taiwan). Splinkerette PCR (Potter and Luo 2010; Tao et al. 2022) was used to attempt to identify novel element flanking sequences, with primers designed using the BCFJ01043787.1 sequence as reference (Morales Poole et al. 2017). This yielded genomic sequences flanking the 3' TIR which were used to identify a sequence in the NCBI database (BCFJ01045349.1) that lacked *PfRAGL-A* sequences but revealed a potential *PfRAGL-A* transposon insertion site. PCR primers were then designed that allowed identification of the 5' TIR with flanking sequences and subsequently the intact *PfRAGL-A* transposon. All primer sequences used in this and other experiments in our study can be found in [supplementary File S4](#).

N-terminal Splicing PCR and Sequencing

Adult *P. flava* were collected and embryo cultures were carried out as described previously (Lin et al. 2016). Total RNA was extracted from various developmental stages using the RNeasy Micro kit (Qiagen) and was reverse transcribed using the GoScript Reverse Transcription System (Promega) with oligo dT primers and PCR amplification of *RAG2L-A* was performed using the 2xKAPA LongRange HotStart ReadyMix (Kapa Biosystems) for 35 cycles with forward and reverse primers shown in [supplementary File S4](#). In some cases, products of the first PCR reaction were used as template for a second PCR. The amplicons were cloned and sequenced to confirm their identities. To examine whether the RNA samples were contaminated with genomic DNA, PCR was performed using equal amount of RNA as template (without reverse transcription). To examine the integrity of the RNA isolated from unfertilized eggs, RT-PCR was conducted to amplify a fragment of *vasa*, a known maternal transcript (Lin et al. 2021), using primers that match to sequences located on different exons ([supplementary File S4](#)).

The *OspRAG1L-A* locus on WGS scaffold JXSR01S0 03992.1 is interrupted by 2 contig merge areas, each flanked by duplicated segments of <100 bp. To recover the complete sequence, genomic DNA was prepared from a sample

of *O. spiculata* (generously provided by T. Arehart, Crystal Cove Conservancy) using SDS/proteinase K digestion and phenol/chloroform extraction followed by isopropanol precipitation. The *OspRAG1L-A* locus was amplified by PCR using 5'-TCGTTCTGTTTTAGGGACAAAGC and 5'-GTTGTGACCTCCTTGCCGCATCT as primers. The PCR reaction was carried out using GoTaq Long PCR 2× Master Mix (Promega) for 35 cycles. Amplicons were cloned and the plasmid inserts were initially sequenced using an Applied Biosystems 3730xL DNA Analyzer and then confirmed by the Whole Plasmid Sequencing service from Plasmidsaurus (<https://www.plasmidsaurus.com/>).

In Vivo Recombination and Plasmid-to-Plasmid Transposition Assays

The recombination assay was performed in Expi293 cells transfected with 1 µg of pTT5M-RAG1 R848M/E649V and pTT5M-RAG2 variants, and 2 µg of p290G using lipofectamine 2000 as described (Huang et al. 2016; Zhang et al. 2019). Cells were collected 72 h post-transfection and washed twice with PBS containing 2% FBS. The percentage of live cells expressing GFP was analyzed by flow cytometry as described (Zhang et al. 2019).

Transposition activity was measured using a plasmid-to-plasmid transposition assay as described previously (Zhang et al. 2019). Briefly, 293T cells were transfected with 4 µg each of pTT5M-RAG1 R848M/E649V and pTT5M-RAG2 variants, 6 µg of the donor plasmid (pTetRSS), and 10 µg of the target plasmid (pECFP-1) using polyethyleneimine. The cell medium was changed 24 h post-transfection and cells were collected after 48 h. Purified extrachromosomal DNA (300 ng) was used to transform electrocompetent MC1061 bacterial cells, which were plated onto kanamycin or kanamycin-tetracycline-streptomycin (KTS) plates. Transposition efficiency was calculated by dividing the number of colonies on KTS plates by the number of colonies on K plates, correcting for dilution factors. Plasmids from 30 colonies from KTS plates were sequenced to determine whether they contained a bona fide transposition event (3 to 7 bp TSD) and the results used to calculate a corrected transposition efficiency value by counting only the plasmids that contained a transposition event.

Phylogenetic Analyses

The phylogenetic analysis was performed using the Maximum Likelihood method implemented in IQtree (Minh et al. 2020) and PhyML (Guindon et al. 2010) on the most conserved regions of RAG1 (NBD plus catalytic core) and RAG2 (kelch blades 2 to 5). IQtree phylogeny was performed using automated substitution selection with the FreeRate model for heterogeneity (Soubrier et al. 2012), the ultrafast bootstrap approximation (Minh et al. 2013) and SH-aLRT branch test, both with 1,000 replicates, and the approximate Bayes test (Anisimova et al. 2011). In parallel, to assess robustness, phylogeny analyses

were performed using PhyML with the SMS model selection (Lefort et al. 2017) using either Akaike or Bayesian Information Criterion (AIC/BIC) and SH-aLRT branch support. Evolutionary relationships between species were retrieved from TimeTree (Kumar et al. 2017) and tree graphics were generated using iTOL (Letunic and Bork 2019).

Protein Sequence Analysis

The predicted RAG1L/2L protein sequences were investigated for structural compliance with the expected domain organization and fold characteristics derived from structures of RAG (e.g. Kim et al. 2015) and BbeRAGL (Zhang et al. 2019). Secondary structure, relative solvent accessibility, and intrinsic disorder predictions were generated using several methods (Cheng et al. 2005; Drozdetskiy et al. 2015; Jones and Cozzetto 2015; Wang et al. 2016; Buchan and Jones 2019) and further merged into a consensus structural profile. Predictive 3D models of the identified *P. flava* and *O. spiculata* sequences were generated using AlphaFold (Jumper et al. 2021), OmegaFold (Wu et al. 2022), and Modeller remote homology modeling (Webb and Sali 2016).

Identity/similarity matrices were computed on the most conserved regions of RAG1 (NBD plus catalytic core) and RAG2 (kelch blades 2 to 5) using Ugene (Okonechnikov et al. 2012) and in-house scripts. Similarity scores were derived from Blosum62 by considering as similar all amino acid pairs with Blosum62 scores above 0. Protein alignments were performed using MAFFT (Minh et al. 2020) and identity/similarity were computed excluding gap positions.

Sequence variability analysis of the acidic hinge domain and PHD was performed starting from a set of 421 RAG2 sequences retrieved from UniprotKB using Jackhmmer (Johnson et al. 2010) and mouse RAG2 as query. Sequence variability was expressed as KL divergence and computed using Weblogo (Crooks et al. 2004). Analysis of electrostatic 3D surface potential was performed using the Adaptive Poisson-Boltzmann Solver (APBS) and all displayed 3D graphics were generated using the PyMOL Molecular Graphics System, Version 2.2.3 Schrödinger, LLC.

Supplementary Material

Supplementary material is available at *Molecular Biology and Evolution* online.

Acknowledgments

The authors thank Tim Arehart for providing *O. spiculata* samples, Jr-Kai Yu for providing *P. flava* samples, Ellen Hsu for help in analysis of TIR sequences in the BAC sequence analysis, Katherine Buckley for sending *P. flava* BAC library clones, and Yi-Chih Chen for assistance with the RT-PCR analysis. This work was supported by a public grant overseen by the French National Research Agency (ANR) as part of the second “Investissements d’Avenir” program

(reference: ANR-17-RHUS-000X) (P.P.), Romanian Academy programs 1 and 3 of IBAR (A.-J.P.), Romanian Ministry of Education and Research, CNCS—UEFISCDI, project number PN-IV-P1-PCE-2023-0636 (A.-J.P.), grant 111-2326-B-001-018 from NSTC, Taiwan (Y.-H.S.), Project 31970852 to S.Y. from The National Natural Science Foundation of China, and NIH grant R01 AI137079 (D.G.S.).

Author Contributions

P.P., D.G.S., A.-J.P., Y.-H.S., S.Y., and A.X. provided overall direction for the experiments and analyses. E.C.M., L.L.T., and L.T.-N. performed database searches, sequence analyses, sequence alignments, and phylogenetic analyses. E.C.M. and L.T.-N. contributed to the development of hypotheses and experimental plans. L.T.-N., T.-P.F., and C.-Y.L. contributed to the screening of BAC libraries and sequencing and analysis of BAC clones. L.T.-N. and E.C.M. performed the analyses that led to the characterization of the complete *PffRAGL^H-A* transposon. E.C.M. performed structural predictions and variability analyses and E.C.M., L.L.T., and L.T.-N. created the figures. J.X. performed the transposition and V(D)J recombination assays. Z.H. performed the experiments that resulted in the identification of *PffRAGL^T-A.14*. All authors contributed to the data interpretation and analysis and to the writing of the paper.

Conflict of interest statement. None declared.

Data Availability

The data underlying this article are available in the article and in its online supplementary material.

References

- Agrawal A, Eastman QM, Schatz DG. Transposition mediated by RAG1 and RAG2 and its implications for the evolution of the immune system. *Nature*. 1998;**394**(6695):744–751. <https://doi.org/10.1038/29457>.
- Almeida MV, Vernaz G, Putman ALK, Miska EA. Taming transposable elements in vertebrates: from epigenetic silencing to domestication. *Trends Genet*. 2022;**38**(6):529–553. <https://doi.org/10.1016/j.tig.2022.02.009>.
- Anisimova M, Gil M, Dufayard JF, Dessimoz C, Gascuel O. Survey of branch support methods demonstrates accuracy, power, and robustness of fast likelihood-based approximation schemes. *Syst Biol*. 2011;**60**(5):685–699. <https://doi.org/10.1093/sysbio/syr041>.
- Arshinoff BI, Cary GA, Karimi K, Foley S, Agalakov S, Delgado F, Lotay VS, Ku CJ, Pells TJ, Beatman TR, et al. Echinobase: leveraging an extant model organism database to build a knowledgebase supporting research on the genomics and biology of echinoderms. *Nucl Acids Res*. 2022;**50**(D1):D970–D979. <https://doi.org/10.1093/nar/gkab1005>.
- Berardi A, Quilici G, Spiliotopoulos D, Corral-Rodriguez MA, Martin-Garcia F, Degano M, Tonon G, Ghitti M, Musco G. Structural basis for PHDVC5HCHNSD1-C2HRNizp1 interaction: implications for Sotos syndrome. *Nucl Acids Res*. 2016;**44**(7):3448–3463. <https://doi.org/10.1093/nar/gkw103>.
- Bettridge J, Na CH, Pandey A, Desiderio S. H3k4me3 induces allosteric conformational changes in the DNA-binding and catalytic regions

- of the V(D)J recombinase. *Proc Natl Acad Sci USA*. 2017;**114**(8): 1904–1909. <https://doi.org/10.1073/pnas.1615727114>.
- Braso-Vives M, Marletaz F, Echchiki A, Mantica F, Acemel RD, Gomez-Skarmeta JL, Hartasanchez DA, Le Targa L, Pontarotti P, Tena JJ, et al. Parallel evolution of amphioxus and vertebrate small-scale gene duplications. *Genome Biol*. 2022;**23**(1):243. <https://doi.org/10.1186/s13059-022-02808-6>.
- Buchan DWA, Jones DT. The PSIPRED protein analysis workbench: 20 years on. *Nucl Acids Res*. 2019;**47**(W1):W402–W407. <https://doi.org/10.1093/nar/gkz297>.
- Carmona LM, Schatz DG. New insights into the evolutionary origins of the recombination-activating gene proteins and V(D)J recombination. *FEBS J*. 2017;**284**(11):1590–1605. <https://doi.org/10.1111/febs.13990>.
- Chen X, Cui Y, Best RB, Wang H, Zhou ZH, Yang W, Gellert M. Cutting antiparallel DNA strands in a single active site. *Nat Struct Mol Biol*. 2020;**27**(2):119–126. <https://doi.org/10.1038/s41594-019-0363-2>.
- Chen X, Cui Y, Wang H, Zhou ZH, Gellert M, Yang W. How mouse RAG recombinase avoids DNA transposition. *Nat Struct Mol Biol*. 2020;**27**(2):127–133. <https://doi.org/10.1038/s41594-019-0366-z>.
- Cheng J, Randall AZ, Sweredoski MJ, Baldi P. SCRATCH: a protein structure and structural feature prediction server. *Nucl Acids Res*. 2005;**33**(Web Server):W72–W76. <https://doi.org/10.1093/nar/gki396>.
- Coussens MA, Wendland RL, Deriano L, Lindsay CR, Arnal SM, Roth DB. RAG2's acidic hinge restricts repair-pathway choice and promotes genomic stability. *Cell Rep*. 2013;**4**(5):870–878. <https://doi.org/10.1016/j.celrep.2013.07.041>.
- Craig NL. A moveable feast: an introduction to mobile DNA. In: Craig NL, Chandler M, Gellert M, Lambowitz AM, Rice PA, Sandmeyer SB, editors. *Mobile DNA III*. Washington D. C.: ASM Press; 2015. p. 3–39.
- Crooks GE, Hon G, Chandonia JM, Brenner SE. Weblogo: a sequence logo generator. *Genome Res*. 2004;**14**(6):1188–1190. <https://doi.org/10.1101/gr.849004>.
- Davies DR, Goryshin IY, Reznikoff WS, Rayment I. Three-dimensional structure of the Tn5 synaptic complex transposition intermediate. *Science*. 2000;**289**(5476):77–85. <https://doi.org/10.1126/science.289.5476.77>.
- De Guzman RN, Wu ZR, Stalling CC, Pappalardo L, Borer PN, Summers MF. Structure of the HIV-1 nucleocapsid protein bound to the SL3 psi-RNA recognition element. *Science*. 1998;**279**(5349):384–388. <https://doi.org/10.1126/science.279.5349.384>.
- Diflippantonio MJ, McMahan CJ, Eastman QM, Spanopoulou E, Schatz DG. RAG1 mediates signal sequence recognition and recruitment of RAG2 in V(D)J recombination. *Cell*. 1996;**87**(2): 253–262. [https://doi.org/10.1016/S0092-8674\(00\)81343-4](https://doi.org/10.1016/S0092-8674(00)81343-4).
- dos Reis M, Thawornwattana Y, Angelis K, Telford MJ, Donoghue PC, Yang Z. Uncertainty in the timing of origin of animals and the limits of precision in molecular timescales. *Curr Biol*. 2015;**25**(22): 2939–2950. <https://doi.org/10.1016/j.cub.2015.09.066>.
- Drozdetskiy A, Cole C, Procter J, Barton GJ. JPred4: a protein secondary structure prediction server. *Nucl Acids Res*. 2015;**43**(W1): W389–W394. <https://doi.org/10.1093/nar/gkv332>.
- Feeney AJ, Goebel P, Espinoza CR. Many levels of control of V gene rearrangement frequency. *Immunol Rev*. 2004;**200**(1):44–56. <https://doi.org/10.1111/j.0105-2896.2004.00163.x>.
- Flajnik MF. Re-evaluation of the immunological Big Bang. *Curr Biol*. 2014;**24**(21):R1060–R1065. <https://doi.org/10.1016/j.cub.2014.09.070>.
- Fugmann SD, Messier C, Novack LA, Cameron RA, Rast JP. An ancient evolutionary origin of the Rag1/2 gene locus. *Proc Natl Acad Sci USA*. 2006;**103**(10):3728–3733. <https://doi.org/10.1073/pnas.0509720103>.
- Fugmann SD. The origins of the Rag genes—from transposition to V(D)J recombination. *Semin Immunol*. 2010;**22**(1):10–16. <https://doi.org/10.1016/j.smim.2009.11.004>.
- Gellert M. V(D)J recombination: RAG proteins, repair factors, and regulation. *Annu Rev Biochem*. 2002;**71**(1):101–132. <https://doi.org/10.1146/annurev.biochem.71.090501.150203>.
- Gertz EM, Yu YK, Agarwala R, Schaffer AA, Altschul SF. Composition-based statistics and translated nucleotide searches: improving the TBLASTN module of BLAST. *BMC Biol*. 2006;**4**(1): 41. <https://doi.org/10.1186/1741-7007-4-41>.
- Gopalakrishnan S, Majumder K, Predeus A, Huang Y, Koues OI, Verma-Gaur J, Loguerio S, Su AI, Feeney AJ, Artyomov MN, et al. Unifying model for molecular determinants of the preselection Vbeta repertoire. *Proc Natl Acad Sci USA*. 2013;**110**(34): E3206–E3215. <https://doi.org/10.1073/pnas.1304048110>.
- Guindon S, Dufayard JF, Lefort V, Anisimova M, Hordijk W, Gascuel O. New algorithms and methods to estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst Biol*. 2010;**59**(3):307–321. <https://doi.org/10.1093/sysbio/syq010>.
- He C, Li F, Zhang J, Wu J, Shi Y. The methyltransferase NSD3 has chromatin-binding motifs, PHD5-C5HC, that are distinct from other NSD (nuclear receptor SET domain) family members in their histone H3 recognition. *J Biol Chem*. 2013;**288**(7): 4692–4703. <https://doi.org/10.1074/jbc.M112.426148>.
- Hencken CG, Li X, Craig NL. Functional characterization of an active Rag-like transposase. *Nat Struct Mol Biol*. 2012;**19**(8):834–836. <https://doi.org/10.1038/nsmb.2338>.
- Hiom K, Melek M, Gellert M. DNA transposition by the RAG1 and RAG2 proteins: a possible source of oncogenic translocations. *Cell*. 1998;**94**(4):463–470. [https://doi.org/10.1016/S0092-8674\(00\)81587-1](https://doi.org/10.1016/S0092-8674(00)81587-1).
- Huang S, Tao X, Yuan S, Zhang Y, Li P, Beilinson HA, Zhang Y, Yu W, Pontarotti P, Escrava H, et al. Discovery of an active RAG transposon illuminates the origins of V(D)J recombination. *Cell*. 2016;**166**(1):102–114. <https://doi.org/10.1016/j.cell.2016.05.032>.
- Johnson LS, Eddy SR, Portugaly E. Hidden Markov model speed heuristic and iterative HMM search procedure. *BMC Bioinformatics*. 2010;**11**(1):431. <https://doi.org/10.1186/1471-2105-11-431>.
- Johnson M, Zaretskaya I, Raytselis Y, Merezuk Y, McGinnis S, Madden TL. NCBI BLAST: a better web interface. *Nucl Acids Res*. 2008;**36**(Web Server):W5–W9. <https://doi.org/10.1093/nar/gkn201>.
- Jones DT, Cozzetto D. DISOPRED3: precise disordered region predictions with annotated protein-binding activity. *Bioinformatics*. 2015;**31**(6):857–863. <https://doi.org/10.1093/bioinformatics/btu744>.
- Jumper J, Evans R, Pritzel A, Green T, Figurnov M, Ronneberger O, Tunyasuvunakool K, Bates R, Zidek A, Potapenko A, et al. Highly accurate protein structure prediction with AlphaFold. *Nature*. 2021;**596**(7873):583–589. <https://doi.org/10.1038/s41586-021-03819-2>.
- Kapitonov VV, Jurka J. RAG1 Core and V(D)J recombination signal sequences were derived from Transib transposons. *PLoS Biol*. 2005;**3**(6):e181. <https://doi.org/10.1371/journal.pbio.0030181>.
- Kapitonov VV, Koonin EV. Evolution of the RAG1–RAG2 locus: both proteins came from the same transposon. *Biol Direct*. 2015;**10**(1): 20. <https://doi.org/10.1186/s13062-015-0055-8>.
- Kim MS, Chuenchor W, Chen X, Cui Y, Zhang X, Zhou ZH, Gellert M, Yang W. Cracking the DNA code for V(D)J recombination. *Mol Cell*. 2018;**70**(2):358–370.e4. <https://doi.org/10.1016/j.molcel.2018.03.008>.
- Kim MS, Lapkowski M, Yang W, Gellert M. Crystal structure of the V(D)J recombinase RAG1–RAG2. *Nature*. 2015;**518**(7540): 507–511. <https://doi.org/10.1038/nature14174>.
- Klein DJ, Johnson PE, Zollars ES, De Guzman RN, Summers MF. The NMR structure of the nucleocapsid protein from the mouse mammary tumor virus reveals unusual folding of the C-terminal zinc knuckle. *Biochemistry*. 2000;**39**(7):1604–1612. <https://doi.org/10.1021/bi9922493>.
- Kriatchko AN, Anderson DK, Swanson PC. Identification and characterization of a gain-of-function RAG-1 mutant. *Mol Cell Biol*. 2006;**26**(12):4712–4728. <https://doi.org/10.1128/MCB.02487-05>.

- Kumar S, Stecher G, Suleski M, Hedges SB. Timetree: a resource for timelines, timetrees, and divergence times. *Mol Biol Evol.* 2017;**34**(7):1812–1819. <https://doi.org/10.1093/molbev/msx116>.
- Lapkouski M, Chuenchor W, Kim MS, Gellert M, Yang W. Assembly pathway and characterization of the RAG1/2-DNA paired and signal-end complexes. *J Biol Chem.* 2015;**290**(23):14618–14625. <https://doi.org/10.1074/jbc.M115.641787>.
- Lefort V, Longueville JE, Gascuel O. SMS: smart model selection in PhyML. *Mol Biol Evol.* 2017;**34**(9):2422–2424. <https://doi.org/10.1093/molbev/msx149>.
- Letunic I, Bork P. Interactive Tree Of Life (iTOL) v4: recent updates and new developments. *Nucl Acids Res.* 2019;**47**(W1):W256–W259. <https://doi.org/10.1093/nar/gkz239>.
- Levin HL, Moran JV. Dynamic interactions between transposable elements and their hosts. *Nat Rev Genet.* 2011;**12**(9):615–627. <https://doi.org/10.1038/nrg3030>.
- Lin CY, Tung CH, Yu JK, Su YH. Reproductive periodicity, spawning induction, and larval metamorphosis of the hemichordate acorn worm *Ptychodera flava*. *J Exp Zool B Mol Dev Evol.* 2016;**326**(1):47–60. <https://doi.org/10.1002/jez.b.22665>.
- Lin CY, Yu JK, Su YH. Evidence for BMP-mediated specification of primordial germ cells in an indirect-developing hemichordate. *Evol Dev.* 2021;**23**(1):28–45. <https://doi.org/10.1111/ede.12361>.
- Liu C, Yang Y, Schatz DG. Structures of a RAG-like transposase during cut-and-paste transposition. *Nature.* 2019;**575**(7783):540–544. <https://doi.org/10.1038/s41586-019-1753-7>.
- Liu C, Zhang Y, Liu CC, Schatz DG. Structural insights into the evolution of the RAG recombinase. *Nat Rev Immunol.* 2022;**22**(6):353–370. <https://doi.org/10.1038/s41577-021-00628-6>.
- Liu Y, Subrahmanyam R, Chakraborty T, Sen R, Desiderio S. A plant homeodomain in RAG-2 that binds hypermethylated lysine 4 of histone H3 is necessary for efficient antigen-receptor-gene rearrangement. *Immunity.* 2007;**27**(4):561–571. <https://doi.org/10.1016/j.immuni.2007.09.005>.
- Lohe AR, Hartl DL. Autoregulation of mariner transposase activity by overproduction and dominant-negative complementation. *Mol Biol Evol.* 1996;**13**(4):549–555. <https://doi.org/10.1093/oxfordjournals.molbev.a025615>.
- Lu C, Ward A, Bettridge J, Liu Y, Desiderio S. An autoregulatory mechanism imposes allosteric control on the V(D)J recombinase by histone H3 methylation. *Cell Rep.* 2015;**10**(1):29–38. <https://doi.org/10.1016/j.celrep.2014.12.001>.
- Maman Y, Teng G, Seth R, Kleinstein SH, Schatz DG. RAG1 targeting in the genome is dominated by chromatin interactions mediated by the non-core regions of RAG1 and RAG2. *Nucl Acids Res.* 2016;**44**(20):9624–9637. <https://doi.org/10.1093/nar/gkw633>.
- Martin EC, Vicari C, Tsakou-Ngouafo L, Pontarotti P, Petrescu AJ, Schatz DG. Identification of RAG-like transposons in protozoans suggests their ancient bilaterian origin. *Mob DNA.* 2020;**11**(1):17. <https://doi.org/10.1186/s13100-020-00214-y>.
- Matthews AG, Kuo AJ, Ramon-Maiques S, Han S, Champagne KS, Ivanov D, Gallardo M, Carney D, Cheung P, Ciccone DN, et al. RAG2 PHD finger couples histone H3 lysine 4 trimethylation with V(D)J recombination. *Nature.* 2007;**450**(7172):1106–1110. <https://doi.org/10.1038/nature06431>.
- Minh BQ, Nguyen MA, von Haeseler A. Ultrafast approximation for phylogenetic bootstrap. *Mol Biol Evol.* 2013;**30**(5):1188–1195. <https://doi.org/10.1093/molbev/mst024>.
- Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R. IQ-TREE 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol Evol.* 2020;**37**(5):1530–1534. <https://doi.org/10.1093/molbev/msaa015>.
- Montano SP, Rice PA. Moving DNA around: DNA transposition and retroviral integration. *Curr Opin Struct Biol.* 2011;**21**(3):370–378. <https://doi.org/10.1016/j.sbi.2011.03.004>.
- Morales Poole JR, Huang SF, Xu A, Bayet J, Pontarotti P. The RAG transposon is active through the deuterostome evolution and domesticated in jawed vertebrates. *Immunogenetics.* 2017;**69**(6):391–400. <https://doi.org/10.1007/s00251-017-0979-5>.
- Okonechnikov K, Golosova O, Fursov M, team U. 2012. Unipro UGENE: a unified bioinformatics toolkit. *Bioinformatics.* 28:1166–1167.
- Potter CJ, Luo L. Splinkerette PCR for mapping transposable elements in *Drosophila*. *PLoS ONE.* 2010;**5**(4):e10168. <https://doi.org/10.1371/journal.pone.0010168>.
- Ramon-Maiques S, Kuo AJ, Carney D, Matthews AG, Oettinger MA, Gozani O, Yang W. The plant homeodomain finger of RAG2 recognizes histone H3 methylated at both lysine-4 and arginine-2. *Proc Natl Acad Sci USA.* 2007;**104**(48):18993–18998. <https://doi.org/10.1073/pnas.0709170104>.
- Ramsden DA, Baetz K, Wu GE. Conservation of sequence in recombination signal sequence spacers. *Nucl Acids Res.* 1994;**22**(10):1785–1796. <https://doi.org/10.1093/nar/22.10.1785>.
- Ru H, Chambers MG, Fu TM, Tong AB, Liao M, Wu H. Molecular mechanism of V(D)J recombination from synaptic RAG1–RAG2 complex structures. *Cell.* 2015;**163**(5):1138–1152. <https://doi.org/10.1016/j.cell.2015.10.055>.
- Saha A, Mitchell JA, Nishida Y, Hildreth JE, Ariberre JA, Gilbert WV, Garfinkel DJ. A trans-dominant form of Gag restricts Ty1 retrotransposition and mediates copy number control. *J Virol.* 2015;**89**(7):3922–3938. <https://doi.org/10.1128/JVI.03060-14>.
- Schatz DG, Swanson PC. V(D)J recombination: mechanisms of initiation. *Annu Rev Genet.* 2011;**45**(1):167–202. <https://doi.org/10.1146/annurev-genet-110410-132552>.
- Shimazaki N, Tsai AG, Lieber MR. H3k4me3 stimulates the V(D)J RAG complex for both nicking and hairpinning in trans in addition to tethering in cis: implications for translocations. *Mol Cell.* 2009;**34**(5):535–544. <https://doi.org/10.1016/j.molcel.2009.05.011>.
- Solovyev V, Kosarev P, Seledsov I, Vorobyev D. Automatic annotation of eukaryotic genes, pseudogenes and promoters. *Genome Biol.* 2006;**7**(Suppl 1):S10–S12. <https://doi.org/10.1186/gb-2006-7-s1-s10>.
- Soubrier J, Steel M, Lee MS, Der Sarkissian C, Guindon S, Ho SY, Cooper A. The influence of rate heterogeneity among sites on the time dependence of molecular rates. *Mol Biol Evol.* 2012;**29**(11):3345–3358. <https://doi.org/10.1093/molbev/mss140>.
- Spanopoulou E, Zaitseva F, Wang FH, Santagata S, Baltimore D, Panayotou G. The homeodomain region of Rag-1 reveals the parallel mechanisms of bacterial and V(D)J recombination. *Cell.* 1996;**87**(2):263–276. [https://doi.org/10.1016/S0092-8674\(00\)81344-6](https://doi.org/10.1016/S0092-8674(00)81344-6).
- Stanke M, Diekhans M, Baertsch R, Haussler D. Using native and syntactically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics.* 2008;**24**(5):637–644. <https://doi.org/10.1093/bioinformatics/btn013>.
- Swanson PC. The bounty of RAGs: recombination signal complexes and reaction outcomes. *Immunol Rev.* 2004;**200**(1):90–114. <https://doi.org/10.1111/j.0105-2896.2004.00159.x>.
- Tao X, Huang Z, Chen F, Wang X, Zheng T, Yuan S, Xu A. The RAG key to vertebrate adaptive immunity descended directly from a bacterial ancestor. *Natl Sci Rev.* 2022;**9**(8):nwac073. <https://doi.org/10.1093/nsr/nwac073>.
- Teng G, Maman Y, Resch W, Kim M, Yamane A, Qian J, Kieffer-Kwon KR, Mandal M, Ji Y, Meffre E, et al. RAG represents a widespread threat to the lymphocyte genome. *Cell.* 2015;**162**(4):751–765. <https://doi.org/10.1016/j.cell.2015.07.009>.
- Tsai CL, Chatterji M, Schatz DG. DNA mismatches and GC-rich motifs target transposition by the RAG1/RAG2 transposase. *Nucl Acids Res.* 2003;**31**(21):6180–6190. <https://doi.org/10.1093/nar/gkg819>.
- Wang S, Li W, Liu S, Xu J. RaptorX-Property: a web server for protein structure property prediction. *Nucl Acids Res.* 2016;**44**(W1):W430–W435. <https://doi.org/10.1093/nar/gkw306>.
- Webb B, Sali A. Comparative protein structure modeling using MODELLER. *Curr Protoc Bioinformatics.* 2016;**54**(1):5.6.1–5.6.37. <https://doi.org/10.1002/cpbi.3>.
- Wilson DR, Norton DD, Fugmann SD. The PHD domain of the sea urchin RAG2 homolog, SpRAG2L, recognizes dimethylated lysine

- 4 in histone H3 tails. *Dev Comp Immunol.* 2008;**32**(10): 1221–1230. <https://doi.org/10.1016/j.dci.2008.03.012>.
- Wu GS, Yang-lott KS, Klink MA, Hayer KE, Lee KD, Bassing CH. Poor quality Vbeta recombination signal sequences stochastically enforce TCRbeta allelic exclusion. *J Exp Med.* 2020;**217**(9): e20200412. <https://doi.org/10.1084/jem.20200412>.
- Wu R, Ding F, Wang R, Shen R, Zhang X, Luo S, Su C, Wu Z, Xie Q, Berger B, et al. High-resolution de novo structure prediction from primary sequence. *bioRxiv.* 2022. 10.1101/2022.07.21.500999.
- Xiong X, Panchenko T, Yang S, Zhao S, Yan P, Zhang W, Xie W, Li Y, Zhao Y, Allis CD, et al. Selective recognition of histone crotonylation by double PHD fingers of MOZ and DPF2. *Nat Chem Biol.* 2016;**12**(12):1111–1118. <https://doi.org/10.1038/nchembio.2218>.
- Yakovenko I, Agronin J, Smith LC, Oren M. Guardian of the genome: an alternative RAG/Transib co-evolution hypothesis for the origin of V(D)J recombination. *Front Immunol.* 2021;**12**:709165. <https://doi.org/10.3389/fimmu.2021.709165>.
- Yakovenko I, Tobi D, Ner-Gaon H, Oren M. Different sea urchin RAG-like genes were domesticated to carry out different functions. *Front Immunol.* 2022;**13**:1066510. <https://doi.org/10.3389/fimmu.2022.1066510>.
- Yin FF, Bailey S, Innis CA, Ciubotaru M, Kamtekar S, Steitz TA, Schatz DG. Structure of the RAG1 nonamer binding domain with DNA reveals a dimer that mediates DNA synapsis. *Nat Struct Mol Biol.* 2009;**16**(5):499–508. <https://doi.org/10.1038/nsmb.1593>.
- Yu K, Taghva A, Lieber MR. The cleavage efficiency of the human immunoglobulin heavy chain VH elements by the RAG complex: implications for the immune repertoire. *J Biol Chem.* 2002;**277**(7): 5040–5046. <https://doi.org/10.1074/jbc.M109772200>.
- Zeng L, Zhang Q, Li S, Plotnikov AN, Walsh MJ, Zhou MM. Mechanism and regulation of acetylated histone binding by the tandem PHD finger of DPF3b. *Nature.* 2010;**466**(7303): 258–262. <https://doi.org/10.1038/nature09139>.
- Zhang Y, Cheng TC, Huang G, Lu Q, Surleac MD, Mandell JD, Pontarotti P, Petrescu AJ, Xu A, Xiong Y, et al. Transposon molecular domestication and the evolution of the RAG recombinase. *Nature.* 2019;**569**(7754):79–84. <https://doi.org/10.1038/s41586-019-1093-7>.
- Zhang Y, Corbett E, Wu S, Schatz DG. Structural basis for the activation and suppression of transposition during evolution of the RAG recombinase. *EMBO J.* 2020;**39**(21):e105857. <https://doi.org/10.15252/emboj.2020105857>.