



A family of unusual immunoglobulin superfamily genes in an invertebrate histocompatibility complex

Aidan L. Huene^{a,b}, Steven M. Sanders^{a,b}, Zhiwei Ma^{a,b}, Anh-Dao Nguyen^c, Sergey Koren^c, Manuel H. Michaca^{a,b}, James C. Mullikin^{c,d}, Adam M. Phillippy^c, Christine E. Schnitzler^{e,f}, Andreas D. Baxevasis^c, and Matthew L. Nicotra^{a,b,g,1}

Edited by Ruslan Medzhitov, Yale University, New Haven, CT; received May 13, 2022; accepted August 11, 2022

Most colonial marine invertebrates are capable of allorecognition, the ability to distinguish between themselves and conspecifics. One long-standing question is whether invertebrate allorecognition genes are homologous to vertebrate histocompatibility genes. In the cnidarian *Hydractinia symbiolongicarpus*, allorecognition is controlled by at least two genes, *Allorecognition 1* (*Alr1*) and *Allorecognition 2* (*Alr2*), which encode highly polymorphic cell-surface proteins that serve as markers of self. Here, we show that *Alr1* and *Alr2* are part of a family of 41 *Alr* genes, all of which reside in a single genomic interval called the Allorecognition Complex (ARC). Using sensitive homology searches and highly accurate structural predictions, we demonstrate that the *Alr* proteins are members of the immunoglobulin superfamily (IgSF) with V-set and I-set Ig domains unlike any previously identified in animals. Specifically, their primary amino acid sequences lack many of the motifs considered diagnostic for V-set and I-set domains, yet they adopt secondary and tertiary structures nearly identical to canonical Ig domains. Thus, the V-set domain, which played a central role in the evolution of vertebrate adaptive immunity, was present in the last common ancestor of cnidarians and bilaterians. Unexpectedly, several *Alr* proteins also have immunoreceptor tyrosine-based activation motifs and immunoreceptor tyrosine-based inhibitory motifs in their cytoplasmic tails, suggesting they could participate in pathways homologous to those that regulate immunity in humans and flies. This work expands our definition of the IgSF with the addition of a family of unusual members, several of which play a role in invertebrate histocompatibility.

allorecognition | *Hydractinia* | AlphaFold | gene complex | nonself recognition

Allorecognition is the ability to distinguish self from nonself within the same species. Most encrusting colonial marine invertebrates, including sponges, corals, hydroids, bryozoans, and ascidians, are capable of allorecognition (1). This enables colonies to compete with conspecifics for space and prevents them from competing with themselves as they grow on three-dimensional surfaces (2). Allorecognition also reduces the risk of stem cell parasitism, which can occur if unrelated colonies fuse and one colony's germline contributes disproportionately to the gametic output of the chimera (3).

Allorecognition has long attracted the attention of marine ecologists interested in spatial competition (2), population geneticists interested in the generation and maintenance of allelic diversity (4), and evolutionary biologists interested in units of selection and the origins of multicellularity (5, 6). In addition, ever since immunologists learned that corals and sea squirts exhibit allorecognition, they have wondered whether the genes that underlie this ability might be homologous to vertebrate histocompatibility genes (7). If so, studying invertebrate allorecognition could help resolve the evolutionary history of immunity and perhaps lead to novel therapies in immunity and transplantation. Together, these interests have motivated the study of allorecognition genes in several species, including the poriferan *Amphimedon queenslandica* (8), the protochordate *Botryllus schlosseri* (9–11), and the cnidarian *Hydractinia symbiolongicarpus* (12–15).

In *Hydractinia*, allorecognition is controlled by the allorecognition complex (ARC), which encodes two linked genes called *Allorecognition 1* (*Alr1*) and *Allorecognition 2* (*Alr2*) (12, 13). In laboratory strains, each gene has two alleles, and together they control allorecognition using a “missing-self” strategy. Colonies that share at least one allele at both genes recognize each other as self and fuse to create a larger colony (Fig. 1*A*). Colonies that do not share alleles at either locus recognize each other as nonself and fight by discharging harpoon-like organelles called nematocysts into their opponents (Fig. 1*B*). Colonies that share alleles only at one locus—either *Alr1* or *Alr2*—fuse but later separate.

Alr1 and *Alr2* encode transmembrane proteins with highly polymorphic extracellular domains (14, 15). In nature, there are tens to hundreds of alleles for each gene (15, 16). In vitro studies have shown that the *Alr1* protein is capable of *trans* (cell-to-cell)

Significance

The immunoglobulin superfamily (IgSF) is one of the largest and most functionally versatile domain families in animal genomes. Although their amino acid sequences can vary considerably, IgSF domains have been traditionally defined by conserved residues at several key positions in their fold. Here, we sequenced an invertebrate histocompatibility complex and discovered a family of IgSF genes with amino acid sequences that lack most of these residues yet are predicted to adopt folds virtually identical to canonical V-set and I-set IgSF domains. This work broadens the definition of the IgSF and shows that the V-set domain was present earlier in animal evolution than previously appreciated.

Author contributions: A.L.H., S.K., J.C.M., A.M.P., C.E.S., A.D.B., and M.L.N. designed research; A.L.H., Z.M., A.-D.N., S.K., M.H.M., J.C.M., A.M.P., C.E.S., and M.L.N. performed research; A.L.H., S.M.S., A.-D.N., S.K., M.H.M., J.C.M., A.M.P., C.E.S., and M.L.N. analyzed data; and A.L.H., S.M.S., M.H.M., C.E.S., A.D.B., and M.L.N. wrote the paper.

The authors declare no competing interest.

This article is a PNAS Direct Submission.

Copyright © 2022 the Author(s). Published by PNAS. This article is distributed under Creative Commons Attribution-NonCommercial-NoDerivatives License 4.0 (CC BY-NC-ND).

¹To whom correspondence may be addressed. Email: matthew.nicotra@pitt.edu.

This article contains supporting information online at <http://www.pnas.org/lookup/suppl/doi:10.1073/pnas.2207374119/-/DCSupplemental>.

Published September 26, 2022.

homophilic binding, which only occurs between allelic variants with similar extracellular sequences (17). The same is true for Alr2 (17, 18). This variant-specific homophilic binding is hypothesized to be the mechanism of self/nonself discrimination *in vivo*.

The homology of *Alr1* and *Alr2* to other genes is unresolved. When they were originally identified, it was not possible to identify orthologs for either gene outside of *Hydractinia*. However, BLAST searches did return statistically significant alignments between domains in the Alr1 and Alr2 extracellular regions and immunoglobulin (Ig) domains (Fig. 1C). These hits had low sequence identity, making it difficult to determine whether the Alr domains belonged to the immunoglobulin superfamily (IgSF). These Ig-like domains were followed by a region referred to as the extracellular spacer (ECS), which had no detectable similarity to other domains.

It is also unclear how many *Alr* genes exist in *Hydractinia*. *Alr1* is flanked by several *Alr1*-like sequences (15), and the genomic region immediately upstream of *Alr2* contains two *Alr2* pseudogenes (14, 19). However, the full extent of this gene family is unknown because only a fraction of the ARC has been sequenced. Identifying additional *Alr* genes is of particular interest because the ARC probably contains at least one additional allodeterminant. Evidence for this comes from the fact that *Alr1* and *Alr2* can fail to predict allorecognition responses in field-collected colonies (14, 15, 20). The unidentified allodeterminant(s) likely reside in the ARC because genetic studies have shown that all *Hydractinia* allodeterminants, including dominant and codominant modifiers, are linked to *Alr1* and *Alr2* (21).

Here, we report the discovery of a family of 41 *Alr* loci, all encoded in the ARC. These genes show evidence of ancient and recent duplications. There is also evidence of alternative splicing that could give rise to functionally distinct isoforms. A majority of these genes encode single-pass transmembrane proteins with V-set and I-set Ig domains with highly unusual amino acid sequences. This indicates that the V-set domain was present in the last common ancestor of cnidarians and bilaterians, which is significant because V-set domains play a central role in the vertebrate adaptive immune system but have not been previously described outside of bilaterians. Unexpectedly, we find that the ECS of Alr proteins encodes a fibronectin III (Fn3)-like fold. Several Alr proteins also have immunoreceptor tyrosine-based activation motifs (ITAMs) or immunoreceptor tyrosine-based inhibitory motifs (ITIMs) in their cytoplasmic tails, suggesting a role for a conserved ITAM/ITIM-mediated signaling pathway in *Hydractinia*. Together, our results outline the full extent of a family of IgSF proteins, several of which are candidates for additional allodeterminants in *Hydractinia*.

Results

The ARC Spans at Least 11.8 Mb. In previous work, we generated a BAC library from a colony homozygous for the ARC-F haplotype (colony 833-8 in *SI Appendix*, Fig. S1) and used it to perform chromosome walks starting from markers in the ARC linkage map (Fig. 1D) (13–15). The minimum tiling path of each walk was sequenced, resulting in six BAC contigs with a total length of 2.9 Mb (Fig. 1D and *SI Appendix*, Fig. S2A). Here, we sequenced and assembled the genome of colony 236-21,

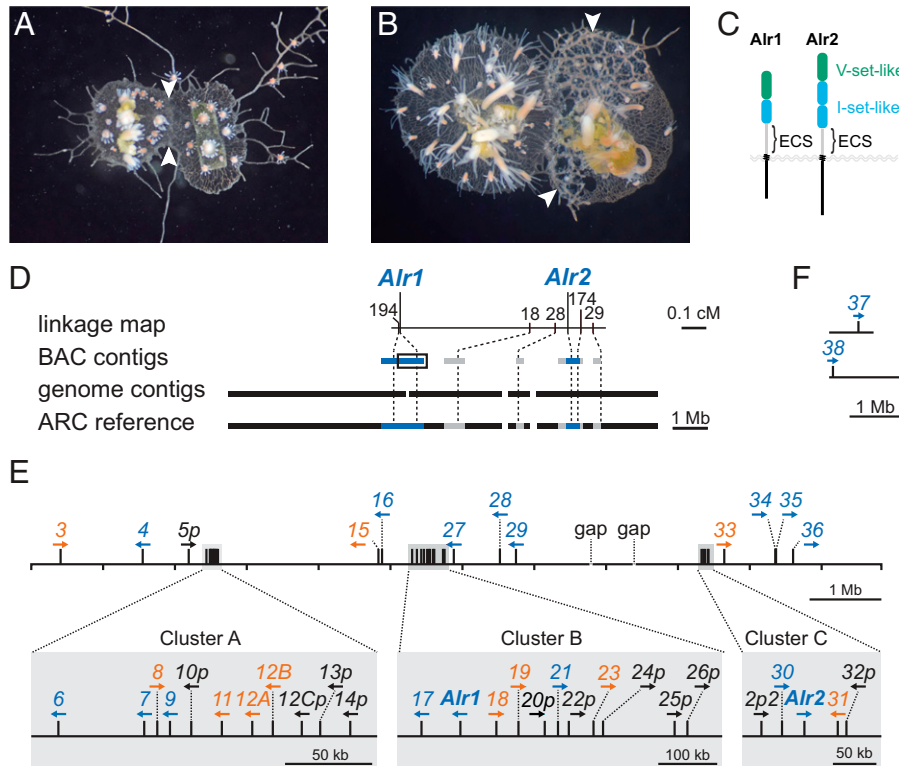


Fig. 1. Allorecognition in *Hydractinia* and the assembly of the *Alr* gene complex. (A) Fusion between two compatible colonies. Arrows point to the region of fusion. (B) Rejection between two incompatible colonies. The colony on the left has grown over the colony on the right. Arrows point to specialized structures called hyperplastic stolons, which are destroying the underlying tissue. (C) Domain architecture of Alr1 and Alr2. (D) Chromosome walks from the ARC linkage map (*Top*) generated six BAC contigs (below; blue = previously published; gray = this work). These were aligned to contigs from the assembled genome of an animal homozygous across the ARC (black). The resulting 11.83-Mb reference sequence was constructed by concatenating the BAC and genome assemblies (*Bottom*). (E) Identity, location, and orientation of *Alr* family members within the ARC reference (blue, bona fide gene; orange, putative gene; black, pseudogene). (F) Two *Alr* genes located in genome contigs that could not be physically linked to the ARC reference sequence.

an ARC-F homozygote and a descendant of colony 833-8 (*SI Appendix*, Fig. S1). Genomic DNA was sequenced via PacBio long-read sequencing and polished with Illumina data to create a genome assembly that was 431 Mb long, with 5,697 contigs and an N50 of 224 kb. It is available for download via Zenodo (22). We then aligned the original BAC contigs to this new assembly using Nucmer (23). We identified five genome contigs that overlapped the BAC contigs with >99% sequence identity (Fig. 1D and *SI Appendix*, Fig. S3 and Table S1). The only major discrepancies between the BAC contigs and the genome contigs were in repeat regions. Therefore, we merged these sequences by filling the gaps between the BAC contigs with sequences from the genome assembly. The resulting ARC-F reference sequence spans 11.83 Mb and contains two gaps of unknown size (Fig. 1D and *Dataset S1*).

The ARC Contains a Large Family of *Alr* Genes and Pseudogenes.

We next annotated all *Alr*-like genes, guided by ab initio gene predictions and sequence similarity to *Alr1* and *Alr2*. We also generated strand-specific RNA-sequencing (RNAseq) data from colony 236-21 feeding and reproductive polyps. These reads were mapped to the entire genome assembly then visualized to aid our identification of expressed sequences. Sequences without similarity to *Alr1* or *Alr2* were not annotated. As new gene models were created, we used them in iterative TBLASTX searches to identify *Alr* genes that might not have been detected in earlier similarity searches. Finally, to identify *Alr* genes that exist outside the ARC-F reference sequence, we used TBLASTX to query the full genome assembly with the amino acid translation of each *Alr* gene model. All gene models were then numbered sequentially, with pseudogenes receiving a lowercase “p” at

the end of their name (e.g., *Alr5p*). Alternative splice variants were indicated with a decimal number (e.g., *Alr1.1* and *Alr1.2*). Gene models whose full-length predicted amino acid sequences had >80% sequence identity were given the same number followed by a letter (e.g., *Alr12A* and *Alr12B*).

In total, we created 41 gene models (Fig. 1E). All but two were in the ARC reference sequence (*Dataset S2*). More than half (27/41) were encoded in one of three *Alr* clusters, which we named A, B, and C (Fig. 1E). The remaining genes, *Alr37* and *Alr38*, were on contigs not contiguous with the reference sequence (Fig. 1E and *Datasets S3–S6*). The expression level of each gene model was estimated from our RNAseq data (Fig. 2A). Gene models with less than one fragment per kilobase mapped (FPKM) were deemed unexpressed. To aid our analysis, we then classified each gene model as a bona fide gene, putative gene, or pseudogene.

A gene model was classified as a bona fide gene if it had one open reading frame (ORF) and each exon was expressed and properly spliced. Eighteen models fit this definition, including *Alr1* and *Alr2* (*Datasets S7 and S8*). As shown in Fig. 2B, *Alr* genes generally encoded single-pass transmembrane proteins with one to three domains similar to the *Alr1* and *Alr2* Ig-like domains, an ECS, a transmembrane helix, and a cytoplasmic tail. Without exception, individual Ig-like domains, ECS sequences, or transmembrane helices were encoded by single exons.

A gene model was classified as a putative gene if it had features that made us hesitant to call it a bona fide gene but was not obviously a pseudogene. Eleven gene models fit this definition (*Datasets S9 and S10*). Six had clear sequence similarity to bona fide *Alr* genes but were unexpressed (Fig. 2A and *SI Appendix*, Fig. S3A). We did not call them pseudogenes because

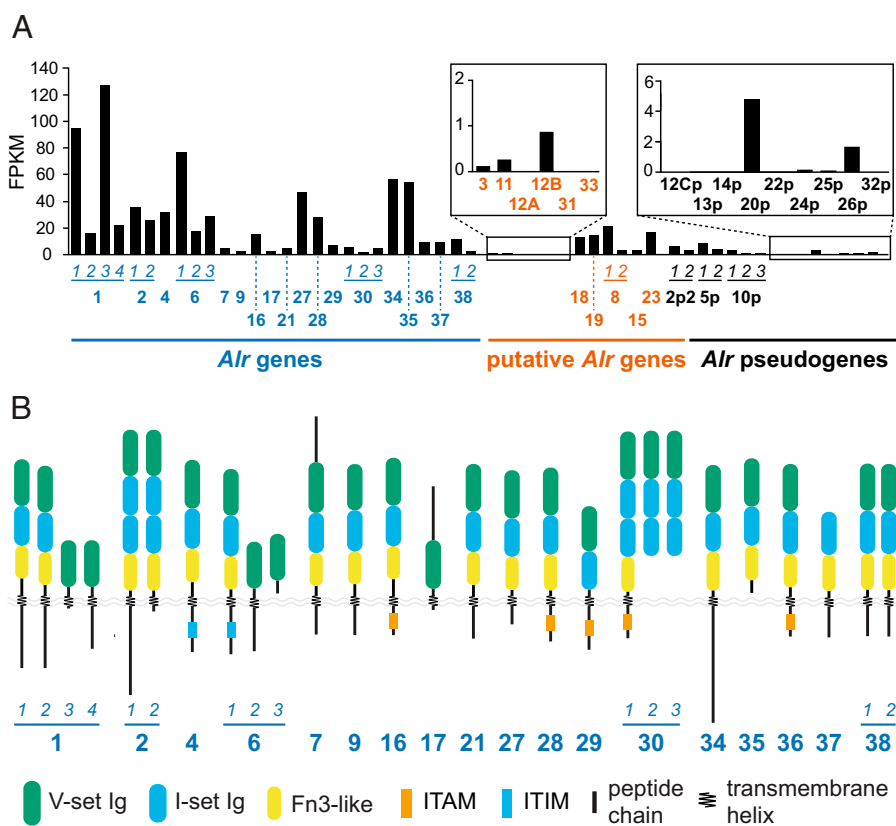


Fig. 2. Expression, domain architecture, and alternative splicing of *Alr* genes. (A) Estimated expression level of each bona fide gene, putative gene, and pseudogene. Genes are identified by bold numbers. Splice variants are indicated by horizontal lines and numbers in italics. The order of the putative genes follows their appearance in *SI Appendix*, Fig. S3. (B) Domain architecture of *Alr* proteins. The final domain predictions and the presence of ITAM and ITIM motifs are indicated here and described later in the text.

they could be expressed at developmental time points or in tissues not represented in our RNAseq dataset. Two gene models had ORFs that would encode a full Alr protein, but there was no evidence of splicing between exons 1, 2, and 3 (*SI Appendix, Fig. S3B*). Three gene models did not encode a signal peptide (*SI Appendix, Fig. S3C*).

A gene model was classified as a pseudogene if it was similar to a bona fide or putative *Alr* gene but was truncated by nonsense or frameshift mutations. Twelve gene models fit this definition (*SI Appendix, Table S2*). Several pseudogenes were expressed at modest levels relative to other *Alr* genes (Fig. 2A).

We paid particular attention to the region directly upstream of *Alr2* because it contained two pseudogenes reported in previous publications (14, 19). The first, immediately upstream of *Alr2*, was named CDS6P by Nicotra et al. (14) and *alr2P1* by Rosengarten et al. (19). It was assumed to be a nonfunctional partial duplication of *Alr2*. Here, we identified additional exons encoding a transmembrane domain, cytoplasmic tail, and 3' untranslated region. We therefore classified this locus as a bona fide gene and named it *Alr30*. The second pseudogene was called CDS5P by Nicotra et al. (14) and *alr2P2* by Rosengarten et al. (19). Here, we also concluded the locus was a pseudogene. For consistency with the previous work, we have named this pseudogene *Alr2p2*.

Alternative Splicing Alters the Domain Architecture of Several *Alr* Gene Products. Several *Alr* genes were alternatively spliced in ways that would change their gene product's domain architecture. Evidence for alternative splicing came from both the assembled transcriptome and the observation of individual RNAseq reads spanning alternative introns. *Alr1*, for example, had four splice variants. *Alr1.1* and *Alr1.2* were previously described (15), but in *Alr1.3* and *Alr1.4*, exon 2 was spliced to new exons encoding alternative transmembrane domains and cytoplasmic tails (*SI Appendix, Fig. S4A*). This resulted in two isoforms lacking domain 2 and the ECS (Fig. 2B). *Alr6* was also alternatively spliced in a similar manner (*SI Appendix, Fig. S4B*). Notably, in *Alr6.3*, exon 2 was spliced to exon 11 and lacked a transmembrane helix, raising the possibility that its gene product is secreted (Fig. 2B). A similar splicing pattern, potentially leading to secreted gene products, was observed in *Alr30* and *Alr35* (Fig. 2B and *SI Appendix, Fig. S4C*). At *Alr2*, transcripts lacking the 22-bp exon 7 were also detected, which

would introduce a frameshift that truncated the cytoplasmic tail (Fig. 2B and *SI Appendix, Fig. S4D*).

Sequences of *Alr* Family Members Are Highly Diverse. The shared domain architecture of *Alr* proteins suggested a history of gene duplication. To investigate the evolutionary relationships between *Alr* family members, we attempted to create a single multiple sequence alignment for the gene products of all *Alr* genes and putative genes. However, their sequences were so divergent that it was impossible to obtain a high-quality alignment even after restricting ourselves to sequences of similar length. This led us to assess overall sequence similarity within the family by performing all possible pairwise alignments. We found the average percent identity between any two amino acid sequences (excluding splice isoforms of the same gene) was $24.3\% \pm 8.6\%$ (*SI Appendix, Fig. S5*). Only 2% of pairwise alignments had more than 50% identity. Thus, a substantial amount of sequence evolution has occurred since the origin of the *Alr* family.

Next, we subdivided each sequence into its constitutive extracellular domains and produced multiple sequence alignments (*Datasets S11–S13*) and neighbor-joining trees (Fig. 3) for each domain type. This revealed a pattern in which domains were more similar if they were encoded close to each other in the genome (Fig. 3). While this analysis has limited power to elucidate the history of the *Alr* gene family, it does suggest that the duplications within Cluster C occurred after it split from Clusters A/B.

The Cytoplasmic Tails of Many *Alr* Proteins Contain ITAMs or ITIMs. Unlike their extracellular domains, the *Alr* cytoplasmic tails were too diverse to be included in a single alignment. Therefore, we used CD-HIT to cluster them at 20% sequence identity. This placed 16/32 into three groups, for which we created separate alignments (*SI Appendix, Fig. S6*). Of the remaining 16 tails, three were <14 amino acids and the rest could not be grouped with other sequences. Thus, the cytoplasmic tails of *Alr* proteins are more divergent than their extracellular domains.

The domain architecture of most *Alr* proteins suggested they might be receptors with intracellular signaling functions. To investigate this, we searched their cytoplasmic tails for signaling motifs. We found ITAMs in the tails of six bona fide and eight putative *Alr* proteins (Fig. 4A and *SI Appendix, Fig. S7*). ITAMs, which have a consensus sequence of $Yxx[I/L]_x(6-9)Yxx[L/I]$ (24),

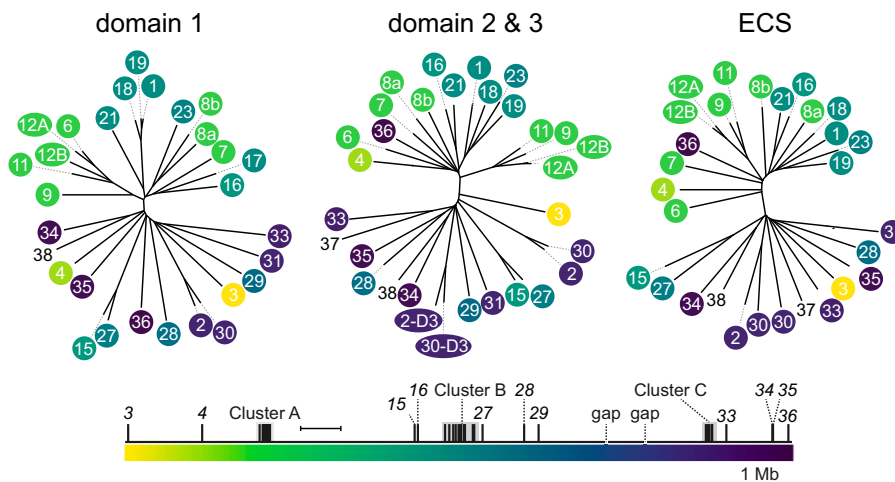


Fig. 3. Sequence similarity between *Alr* extracellular domains. Neighbor-joining trees of *Alr* extracellular domains. Leaves of each tree are color-coded according to their genomic position. Domains from *Alr37* and *Alr38* are not color-coded. Branch lengths calculated according to the BLOSUM26 matrix. (Scale bar, 100 units.)



Fig. 4. Alr ITAM and ITIM motifs and presence of select signaling molecules. (A) Cytoplasmic tails of Alr proteins with ITAMs (orange background). Overlapping ITAMs are indicated with heavier shading. (B) Truncated alignment of cytoplasmic tails with ITIMs (blue background). All tyrosines have a black background.

are found in receptors that activate immune responses in vertebrates (25) and phagocytosis of damaged cells in *Drosophila* (26).

Phosphorylated ITAMs are bound by the dual SH2 domains of a kinase called *Syk* in vertebrates and *Shark* in insects (27). *Syk* and *Shark* are related proteins and differ in that *Shark* has a set of ankyrin repeats between its two SH2 domains. To determine whether the *Hydractinia* genome encodes *Syk* or *Shark*-like kinases that might bind to these ITAMs, we performed a TBLASTN search of the complete genome assembly with the amino acid sequences of human *Syk* and *Drosophila Shark*, identifying *Hydractinia* homologs of each (SI Appendix, Figs. S7D and S8). This is consistent with previous work that has identified *Syk*-like and *Shark*-like kinases in *Hydra* (28, 29).

A second motif in vertebrates called the ITIM is found in receptors that counteract ITAM-mediated signaling and down-regulate immune responses (30). We found ITIMs, defined as [I/L/V/S]xYxx[I/V/L] (31), in two Alr tails, both from group 3 (Fig. 4B). In mammals, phosphorylated ITIMs are bound by the SH2 domains of SHP-1 and SHP-2, two phosphatases that dephosphorylate ITAM-bearing receptors, *Syk*-like proteins, and other components of activating pathways (32, 33). ITIMs are also bound by the phosphoinositide phosphatases SHIP1 and SHIP2, which dampen immune cell activation via the PI3K pathway (34, 35). We searched the *Hydractinia* genome and identified four SHP homologs and one SHIP homolog (SI Appendix, Figs. S7E, S9, and S10).

Together, these data show that many Alr genes have ITAMs and ITIMs, motifs that regulate the recognition of self and nonself in other animals. Moreover, the *Hydractinia* genome includes homologs of the enzymes that bind phosphorylated ITAMs or ITIMs and act as effectors of cellular activation or inhibition.

Domain 1 Is a V-Set Ig Domain. Domain 1 of Alr1 and Alr2 was originally described as V-set-like (14, 15). To determine whether domain 1 was similar to V-set Ig domains in Alr3–Alr38, and to explore the possibility that they were, in fact, homologous, we first used HMMER to compare each sequence to Pfam, a database of hidden Markov models for protein families (36). At an E-value cutoff of <0.01, only 6/29 sequences were similar to V-set domains (SI Appendix, Table S3). We also used HHpred, which is able to detect remote homologies (37), to search the

Structural Classification of Proteins extended (SCOPe) database, which classifies protein domains according to structural and evolutionary relationships (38). Using this approach, we found that 19/29 sequences had a >95% probability of homology to V-set domains (SI Appendix, Table S3).

Homologous proteins can evolve such that their primary sequences become highly divergent but their structures remain relatively unchanged (39). Therefore, we predicted the tertiary structure of each domain 1 with Colabfold (40) to further investigate its homology [deposited in Zenodo (22)]. Colabfold is a Google Colaboratory implementation of AlphaFold2 (41), which is capable of producing structural predictions with sub-angstrom root mean square deviation from experimental structures (42). Each residue in a model produced by Colabfold is assigned a predicted local distance difference test (pLDDT) score, which estimates how well the prediction would agree with an experimental structure. Residues with pLDDT >90 are considered highly accurate and have their side chains oriented correctly 80% of the time (41, 42). Residues with pLDDT >70 generally have their backbones predicted correctly. For the Alr domain 1 sequences, Colabfold produced the structural predictions with average (model-wide) pLDDT scores ranging from 80.6 to 97.4 (Fig. 5A and SI Appendix, Table S4).

Next, we performed structural alignments against the Protein Data Bank (PDB) [https://www.rcsb.org/; (43)] with Dali (44) and PDBeFOLD (45). The top hits from both methods were to V-set Ig domains (e.g., Fig. 5B) (SI Appendix, Table S4). Structural alignments produced by Dali are assigned a Z-score, which is used to estimate the likelihood that the two proteins are homologous. Z-scores between 8 and 20 indicate likely homologs (44). The domain 1 alignments had Z-scores ranging from 12.6 to 16.7, indicating probable homology with V-set domains (Fig. 5A), even though their overall sequence identities were 9 to 20% (SI Appendix, Table S4). Thus, our analysis of the primary, secondary, and tertiary structures of domain 1 all indicated they belong to the V-set family of Ig domains.

V-set domains have nine β -strands named A, B, C, C', C'', D, E, F, and G according to their position in the primary amino acid sequence. Strand C'' is only found in V-set domains. Strand A is split into A and A'. These β -strands are arranged as a Greek key to form a β -sandwich, with one β -sheet consisting of strands A, B, E, and D and the other

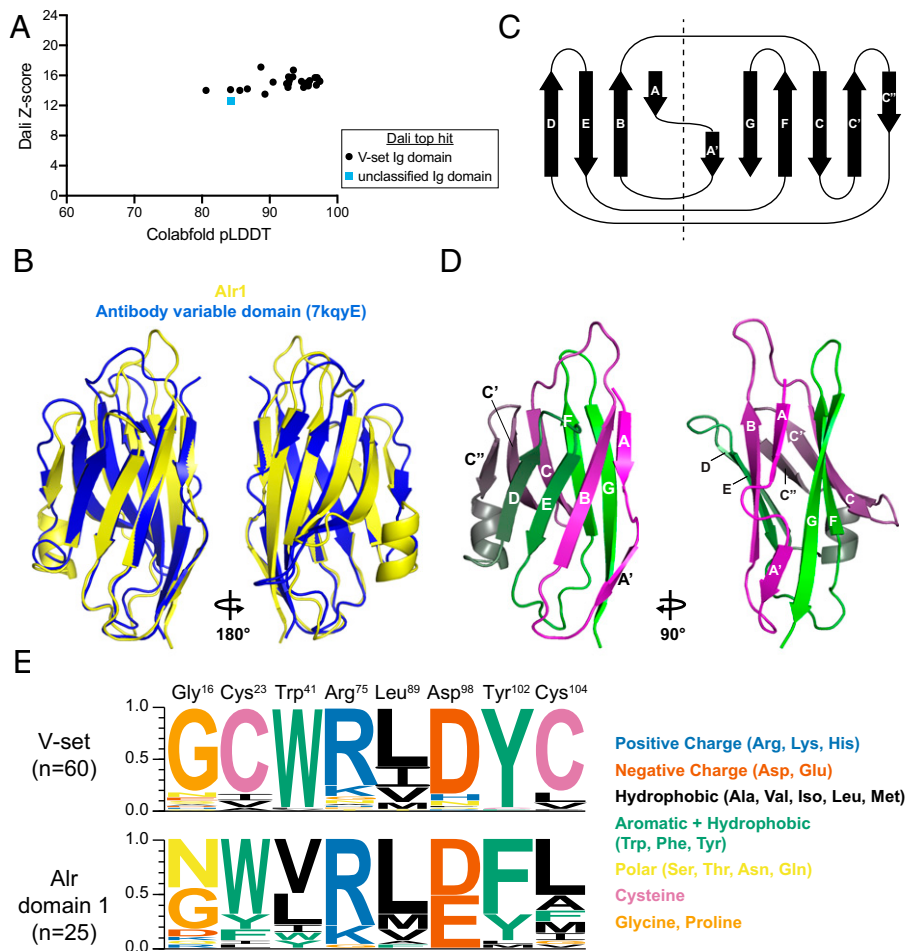


Fig. 5. Sequence analysis and structural predictions of domain 1. (A) Plot of each domain 1 model's average pLDDT score versus its alignment Z-score to the top PDB model identified by DALI. (B) Structural alignment of Alr1 domain 1 to a V-set Ig domain from human heavy-chain antibody (PDB ID 7KQY). (C) Topology of β -strands in V-set folds. Dotted line separates the strands within the same β -sheet. (D) V-set β -strands labeled on the predicted structure of Alr1 domain 1. (E) Sequence logo comparing frequency of amino acids at eight conserved positions in V-set Ig domains (*Top*) and Alr domain 1 sequences (*Bottom*).

consisting of strands A', G, F, C, C', and C'' (Fig. 5C). This is often referred to as the V-frame (46).

To determine whether domain 1 has the V-frame, we used STRIDE (47) to determine the secondary structure of our models, then assigned letters to the β -strands (*SI Appendix, Fig. S11*). Twenty-five models had all nine β -strands (e.g., Fig. 5D). Four models were missing a strand in either the A or A' position (*SI Appendix, Fig. S11*). Notably, all models had the V-set-specific C'' strand (*SI Appendix, Fig. S11*). In addition, 17/29 models included a seven-amino-acid α -helix between the C'' and D strands, which is not typically found in V-set domains.

These results led us to question why HMMER did not identify many Alr domains as V-set domains. Previous studies of V-set domains have identified a set of eight residues that are highly conserved, even across domains with as little as 20% sequence identity (46, 48). According to the nomenclature of Cannon et al. (48), they are Gly¹⁶, Cys²³, Trp⁴¹, Arg⁷⁵, Leu⁸⁹ (or other hydrophobic), Asp⁹⁸, Tyr¹⁰², and Cys¹⁰⁴. To determine whether these residues are conserved in domain 1, we generated a multiple sequence alignment between them and 60 canonical V-set sequences from the Pfam V-set sequence profile (pf07686). We then identified the residues that corresponded to the eight V-set residues (*SI Appendix, Fig. S12*). Our findings are summarized in Fig. 5E.

In V-set domains, Cys²³, Trp⁴¹, and Cys¹⁰⁴ form a nearly invariant structural motif called the “pin” (49). The cysteines form a disulfide bridge between β -strands B and F, while the

tryptophan packs against the bond to stabilize the hydrophobic core of the β -sandwich. All Alr domain 1 sequences, however, lacked these Cys residues, and only two had the Trp. Instead, Cys²⁵ was replaced by bulky, aromatic amino acids (Trp, Phe, or Tyr). Cys¹⁰⁴ and Trp⁴¹ were replaced by hydrophobic amino acids. Thus, in domain 1, the “pin” is replaced by a set of bulky hydrophobic residues that might serve a similar function by stabilizing the core of the β -sandwich.

The fourth and fifth V-set residues, Arg⁷⁵ and Asp⁹⁸, form a salt bridge between the CD and EF loops. The salt bridge is thought to stabilize the “bottom” of the domain and is found only in the V-set and I-set immunoglobulin domains (46, 48). We found the salt bridge in all but three (26/29), although the negatively charged Asp was often replaced with similarly charged Glu. Thus, the salt bridge, a hallmark of V-set and I-set domains, is also present in domain 1.

The sixth canonical residue, Tyr¹⁰², forms the “tyrosine corner,” a structural motif located at the start of the F strand and found only in Greek key proteins (50). While Tyr¹⁰² is highly conserved in V-set Ig-like domains (97% in our seed alignment), it was found in only 7/29 Alr domains (Fig. 5E). Instead, 20/29 had Phe, with its aromatic ring occupying the same location as that of Tyr¹⁰². Mutational studies have shown that a Tyr→Phe mutation has no effect on the ability of V-set Ig domains to fold properly (51). Thus, the residues at position 102 are consistent with domain 1 folding like a V-set domain.

The seventh canonical residue is Gly¹⁶, which is part of a β -turn between strands A' and B. A β -turn is a series of four residues that reverses 180° on itself such that the distance between C α (*i*) and C α (*i*+3) is less than 7 Å (52). β -turns often feature a hydrogen bond between CO(*i*) and NH(*i*+3)—the carboxyl and amine groups on the first and fourth residues, respectively—but this is not a requirement (53). We found this β -turn in all Alr domain 1 models (*SI Appendix*, Fig. S13). Twenty-eight featured a hydrogen bond, defined by the criterion that O(*i*) is <3.5 Å from N(*i*+3) (53). However, position *i* + 2 was Gly in only eleven sequences. Thus, like V-set domains, domain 1 is predicted to have a β -turn between strands A' and B, but in most cases it does not involve a glycine.

The eighth canonical V-set residue is a hydrophobic amino acid, typically leucine, at position 89. This residue resides at the center of the hydrophobic core. In Alr domain 1, 21/29 sequences had a leucine residue at this position. The remaining six had other hydrophobic residues. Thus, this canonical residue is shared between V-set and most domain 1 sequences.

In summary, the Alr V-set domains share some, but not all, sequence motifs commonly considered diagnostic of the V-set family. In the case of the pin motif, tyrosine corner, and β -turn,

the Alr sequences differ in a way that likely preserves the structural motif. Thus domain 1 appears to be a V-set Ig domain with a novel sequence profile.

Domains 2 and 3 Are I-Set Ig Domains. We next explored the relationship between Ig domains and the 29 domain 2 and two domain 3 sequences encoded by bona fide and putative *Alr* genes. At an E-value cutoff of <0.01, HMMER identified 14 as I-set Ig-like domains (pf07679) and another 4 as Ig-like domains (pf13927) (*SI Appendix*, Table S5). HHpred indicated 23/31 had a >95% probability of being homologous to the I-set family (*SI Appendix*, Table S5).

To further explore their potential homology to I-set Ig domains, we generated structural predictions for each domain [deposited in Zenodo (22)], with average pLDDT scores ranging from 81.1 to 95.1 (Fig. 6A and *SI Appendix*, Table S6). I-set domains are similar to V-set domains except that they lack a C' strand, and the C' strand is often shorter (Fig. 6B). We found that 21/31 domain 2 and 3 models had an I-set topology (e.g., Fig. 6C and *SI Appendix*, Fig. S11), and three more were only missing one of the two split A strands. Six others were missing the C' strand (*SI Appendix*, Figs. S11 and S15). These six models

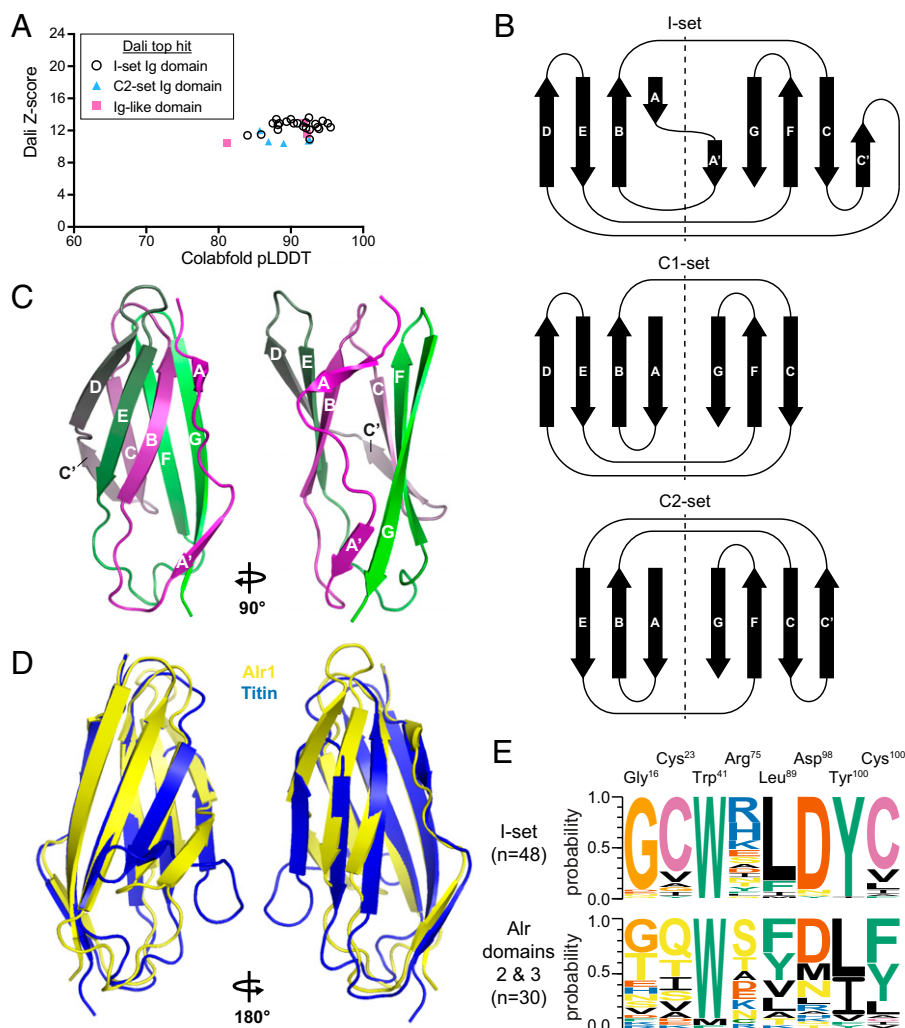


Fig. 6. Sequence analysis and structural predictions of domains 2 and 3. (A) Plot of each domain 2 and domain 3 model's average pLDDT score versus its alignment Z-score to the top PDB model identified by DALI. (B) Topology of β -strands in I-set, C1-set, and C2-set folds. Strand D is part of the EBA sheet in C1 Ig-domains but is part of the CFG sheet in Fn3 and C2-set domains, where it is often labeled C'. Dotted line separates the strands within the same β -sheet. (C) I-set β -strands labeled on the predicted structure of Alr1 domain 2. (D) Structural alignment of Alr1 domain 2 to an I-set Ig domain from titin (PDB ID 2RJM). (E) Sequence logo comparing frequency of amino acids at conserved positions in I-set Ig domains (Top) and Alr domains 2 and 3 (Bottom).

were I-set-like in having a split A strand but were similar to C1-set Ig domains in having a D strand and lacking a C' strand. The remaining domain, from Alr9, was predicted to have only six β -strands and lacked both C' and D strands. When we searched for structural homologs, the top hits from PDBE/FOLD and Dali were to I-set domains, although some hits were annotated as both I-set and C2-set domains, and one top hit for Alr15 was a filamin repeat (Fig. 6A and *SI Appendix*, Table S6). Filamin repeats are in the "Early" (E-set) superfamily of immunoglobulin-like β -folds and are possibly related to the IgSF or fibronectin type III superfamilies (38, 54).

We next investigated whether domains 2 and 3 had any of the conserved sequence motifs found in I-set domains. To do so, we aligned the Alr domains to 48 canonical I-set domains from the Pfam I-set sequence profile (pf07679) (*SI Appendix*, Fig. S16). We then searched for the sequence motifs common to V-frame Ig-like domains (46, 48). These results are summarized as a sequence logo in Fig. 6F. With respect to the pin motif (C–W–C), all Alr domain 2 sequences had the central tryptophan, but the two domain 3 sequences had methionine in its place. All domains lacked the paired cysteines. One cysteine was replaced by a hydrophobic residue, and the second was replaced by residues bearing no consistent physicochemical property (Fig. 6E and *SI Appendix*, Fig. S16). The Alr domains also lacked the salt bridge and tyrosine corner (Fig. 6E). The β -turn between β -strands A' and B was present in all but three structural models (*SI Appendix*, Fig. S17). The last of the eight conserved residues, a hydrophobic residue (typically leucine), was present, although in many it was a tyrosine (Fig. 6F).

More recently, Wang (55) defined the sequence signature of I-set domains via nine sequence motifs, denoted *i* through *ix*. The first four motifs include the C–W–C pin (in motifs *i*, *ii*, and *iv*), the tyrosine corner (part of motif *i*), and the tight turn in the A'B loop (motif *iii*) (*SI Appendix*, Fig. S16). We searched for the remaining five motifs and found that they were present in a majority of Alr domain 2 and 3 sequences (*SI Appendix*, Figs. S16 and S18).

Taken together, our data indicate most domain 2 and 3 sequences belong to the I-set family of Ig domains. Although these domains lack the disulfide bridge, salt bridge, and tyrosine corner, they have most other sequence motifs associated with I-set domains, including the conserved tryptophan. Most are predicted to have an I-set fold, although several appear to have lost the C' strand.

Part of the ECS Adopts a Fibronectin Type III-Like Fold. When Alr1 and Alr2 were described previously, no domains were identified in the ECS (14, 15). Here, we expanded our analysis to include all 29 ECS regions encoded by bona fide and putative *Alr* genes. Although HMMER searches against Pfam only returned two hits to domains of unknown function, HHpred indicated part of the ECS had a 60 to 88% probability of homology to fibronectin type III (Fn3) domains (*SI Appendix*, Table S7). To help us define this potential domain, we aligned the ECS sequences to the 98 Fn3 sequences in the seed alignment of the Pfam Fn3 profile (pf00041.23). We found that the N-terminal portion of the ECS aligned reasonably well to other Fn3 domains, but the C-terminal portion did not (*SI Appendix*, Fig. S19).

Because structure predictions are often better for single domains, we removed the C-terminal portion of the ECS sequences (*SI Appendix*, Fig. S19) then predicted their structures with Colabfold [deposited in Zenodo (22)]. Twenty-six models had average pLDDT >90, with the remaining three models >80 (Fig. 7A and *SI Appendix*, Table S8). Next, we investigated whether the Alr domains had a similar topology to Fn3 domains. Fn3 domains are an Ig-like fold with seven β -strands arranged in the same topology as C2-set immunoglobulin domains (56, 57) (Fig. 7B). All but one ECS model had seven β -strands (e.g., Fig. 7C), with the remaining model missing strand G (*SI Appendix*, Fig. S20). In all models, the β -strands adopted a C2-set/Fn3 topology. We next searched the PDB for proteins with similar structures. All hits from Dali and PDBE/FOLD were to Fn3 domains (Fig. 7A and *SI Appendix*, Table S8). Thus, the

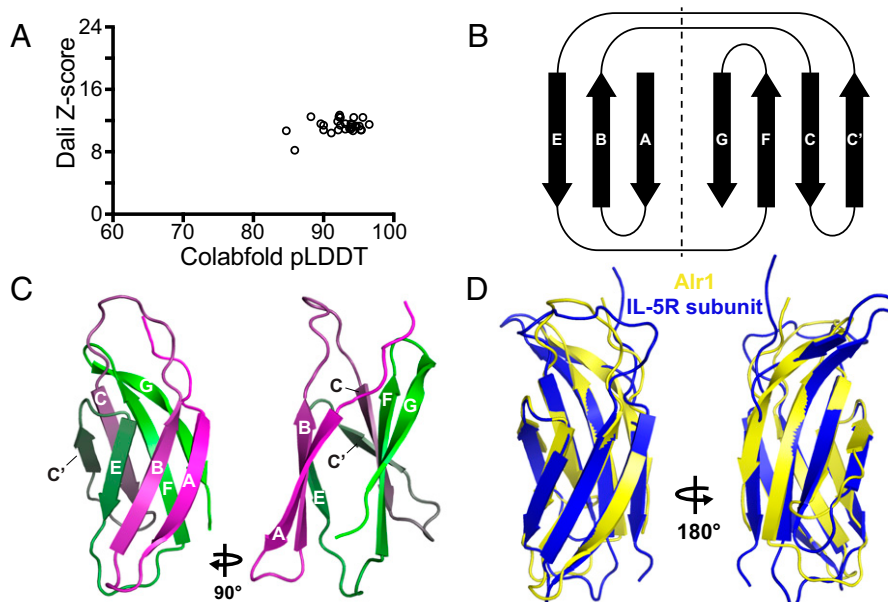


Fig. 7. Sequence analysis and structural predictions of the ECS. (A) Plot of each ECS model's average pLDDT score and corresponding alignment Z-score to the top PDB model identified by DALI. (B) Topology of β -strands in Fn3 folds. Dotted line separates the strands within the same β -sheet. (C) Predicted topology of β -strands in the Alr1 ECS folds. (D) Structural alignment of Alr1 ECS to a Fn3 domain from Interleukin-5 receptor subunit alpha (PDB ID 6H41).

secondary and tertiary structures of this Ig-like fold are predicted to be most similar to Fn3 domains.

The primary amino acid sequences of Fn3 domains have six conserved amino acids (56, 57). To determine whether the Alr Fn3-like domain had these residues, we aligned them to sequences from the Pfam Fn3 profile (pf00041.23). The ECS sequences only aligned well to Fn3 domains at their N terminus (*SI Appendix, Fig. S21*). With respect to the six conserved amino acids, ECS and Fn3 sequences shared a proline at the beginning of strand A, a tryptophan at the end of strand B, and a tyrosine at the beginning of strand C (*SI Appendix, Figs. S21 and S22*). The fourth conserved residue in Fn3 domains, a tyrosine at the end of strand C, was present in 8/29 ECS sequences and replaced by phenylalanine in 14/29 ECS sequences. However, unlike Fn3 domains, the ECS sequences were missing the leucine in the EF loop and the tyrosine residue that forms the tyrosine corner in strand F. Fn3 domains also have six additional “topohydrophobic” positions (i.e., positions usually occupied by VILFMWY residues) (57), which were also present more than 50% of the time in the ECS sequences. (*SI Appendix, Figs. S21 and S22*).

Taken together, these data indicate that most Alr proteins have a Fn3-like fold between their I-set domain and transmembrane helix.

Alr Proteins Have Six Invariant Cysteines Likely to Form Disulfide Bridges. While investigating protein alignments of the Alr domains, we identified six cysteine residues that were conserved across all sequences. These residues are not typically found in immunoglobulin or Fn3 folds. Three were in the Ig domain immediately preceding the Fn3-like domain, and three were located within the Fn3-like domain itself (Fig. 8 *A* and *B*). Within the Ig domain, the first two cysteines were in β -strands A and B and were predicted to form a disulfide bridge in 28/31 models. Within the Fn3-like domain, the first and third cysteines, located in β -strands B and E, were predicted to form a disulfide bridge in 28/31 models. Thus, an intradomain disulfide bridge appears to stabilize the fold of most Alr I-set and Fn3-like folds.

The two remaining invariant cysteines were in the EF loop of the Ig domain and the BC loop of the Fn3-like domain. These domains appear in tandem in 28 of the Alr proteins, raising the possibility that they are stabilized by a disulfide bridge. To test this, we predicted structures for the tandem I-set Ig and Fn3-like domains from each Alr protein [deposited in Zenodo (22)]. All were predicted with high confidence (*SI Appendix, Table S9*), and in all models the two domains were predicted to be linked by a disulfide bridge (e.g., Fig. 8*C*).

Discussion

We have shown that the *Hydractinia* ARC contains a family of genes homologous to *Alr1* and *Alr2*, which we have named the *Alr* gene family. Most encode proteins with a receptor-like topology and domain architecture (Fig. 2*B* and *SI Appendix, Fig. S23*). Despite this similarity, individual Alr proteins have low sequence identity when compared to each other, suggesting the gene family is old, has experienced high rates of molecular evolution, or both.

Our data indicate that the N-terminal domains of Alr proteins are Ig domains. Specifically, domain 1 is either part of the V-set family or represents a new family within the IgSF that is most closely related to V-set domains. Likewise, domains 2 and 3 are either part of the I-set family or represent a new family most closely related to I-set domains. We propose including

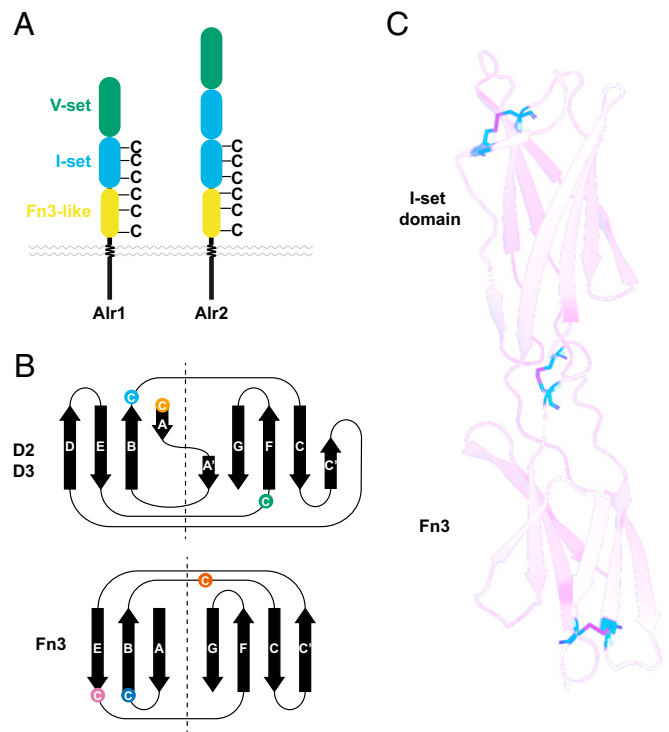


Fig. 8. Invariant cysteines in Alr proteins and final domain predictions. (*A*) Occurrence of invariant cysteine residues found in Alr proteins with an Alr1-like or Alr2-like domain architecture. (*B*) Position of invariant cysteines in domains 2, 3 (I-set), and the ECS fold (Fn3-like). (*C*) Structural prediction of tandem I-set and Fn3-like domains of Alr1 shows each invariant cysteine forms a disulfide bond.

the Alr domains within existing V-set and I-set families. One reason for doing so is that the sequence profiles often used to identify Ig domains may be biased against recognizing them in nonbilaterians. Indeed, profiles for V-set and I-set domains were originally defined using a handful of sequences, mostly from vertebrates and model organisms (46, 58). Today V-set and I-set profiles in Pfam contain thousands of sequences but are dominated by sequences from chordates and, to a much lesser extent, arthropods (36). Although there is considerable sequence diversity within these sequences, they still represent the descendants of only two evolutionary nodes within metazoans. A consequence, suggested by this study, is that many V-set and I-set Ig domains go undetected in nonbilaterians and other understudied taxa. As genomes continue to be sequenced, our ability to detect these domains has been enhanced by new methods for detecting remote homologs and producing highly accurate structural models (40–42). These new data will undoubtedly enhance our understanding of how Ig domains evolved and may force us to revisit how these families are defined.

Indeed, we know of no other V-set Ig domain identified outside of bilaterians. This is significant given the essential role that V-set domains play in the vertebrate adaptive immune system, where their sequences are rearranged somatically to generate sequence diversity in antibodies and T cell receptors (24). Our data indicate the V-set domain is likely older than previously suspected because the last common ancestor of *Hydractinia* and bilaterians appears to have had distinct V-set and I-set Ig domains. These domains appear to have followed different evolutionary trajectories to arrive at their current sequences. As in vertebrates, the *Hydractinia* V-set domain plays a critical role in self/nonself recognition (17, 18). However, it remains unclear whether this reflects a conserved, ancient function or is an example of convergent evolution.

We have also discovered that the region previously referred to as the ECS actually encodes a domain with an Fn3-like fold. However, the sequences of Alr and Fn3 domains differed substantially between strands D through G, and HHpred did not identify them as homologs with high confidence. It therefore remains unclear whether this Fn3-like domain belongs to the fibronectin III superfamily.

Many Alr Ig domains lack the disulfide bridge and tryptophan found at the core of most Ig domains. Although there are many examples of Ig domains lacking either the disulfide bridge or the tryptophan, we are unaware of any V-set or I-set Ig domains that lack both. The reason for this loss is unclear. Intriguingly, the Alr I-set domains, along with their neighboring Fn3-like domains, are predicted to have conserved disulfide bridges that link neighboring β -strands within one β -sheet. There is also a conserved disulfide bridge between neighboring I-set and Fn3-like domains, which may stabilize their orientation.

This study was made possible by the introduction of computational methods that were unavailable when Alr1 and Alr2 were originally described. First, our use of HHpred (37) enabled us to detect remote homologs from the SCOPe database. Second, recent and dramatic improvements in protein structure prediction (41, 42) made accessible through Colabfold (40) enabled us to generate models for most Alr domains that are predicted to accurately represent their true structure. Nonetheless, our conclusion that the Alr proteins contain V-set and I-set Ig domains awaits confirmation from experimentally derived structures when they become available.

One obvious function for the newly described *Alr* genes is allorecognition. It is therefore significant that many *Alr* genes are located outside the genomic region originally mapped by Cadavid et al. (12) and Powell et al. (13). If these *Alr* genes were rendered homozygous in the inbred lines used for mapping, any effect they might have as an allodeterminant would have been missed. In nature, however, these genes might be as polymorphic as *Alr1* and *Alr2*. Such undiscovered allodeterminants could explain the appearance of unexpected allorecognition responses between outbred colonies (14, 15, 20).

Alr6 is a good candidate for such an allodeterminant because it is very similar to *Alr1*. *Alr6* occupies a location in cluster A that is roughly syntenic with the location of *Alr1* in Cluster B. It has an *Alr1*-like intron/exon structure, and it encodes a cytoplasmic tail that is clearly homologous to that of *Alr1*. In addition, *Alr6* is alternatively spliced to generate isoforms that have a single Ig domain and an alternative cytoplasmic tail, a feature shared only with *Alr1*. We hypothesize this similarity could extend to *Alr6* functioning as a third allodeterminant in the ARC.

Alr30 could also be an allodeterminant. *Alr30* is immediately upstream of *Alr2* and was previously considered a pseudogene (15, 19). Here, we show that *Alr30* encodes a transmembrane protein with an extracellular region with recognizable sequence similarity to *Alr2* but a cytoplasmic tail that has no detectable similarity to other *Alr* genes. Since this gene resides within the genomic interval previously defined in Nicotra et al. (14), it is formally possible that it, too, contributes to allorecognition phenotypes.

The domain architecture of most Alr proteins also suggests alternative functions in extracellular protein–protein interactions. Tandem Ig domains are commonly found in cell adhesion molecules, proteins involved in cell-to-cell communication, and immune receptors (59). An adhesive function would also be consistent with that already described for Alr1 and Alr2 (17).

The ITAM motifs in some Alr cytoplasmic tails are also potentially significant. ITAM-mediated signaling activates inflammation and cellular immune responses in vertebrates (27). It also

activates phagocytosis of unwanted or damaged cells in *Drosophila* (26, 60) and appears to promote immune responses to bacteria in oysters (61). An immune function for some ITAM-bearing Alr proteins therefore seems plausible.

Could ITAM-mediated signaling play a role in allorecognition responses? One possibility is that Alr proteins with ITAM motifs activate rejection responses when they bind nonpolymorphic, *Hydractinia*-specific ligands on opposing tissues. This rejection response would then be inhibited if polymorphic allodeterminants bind a compatible ligand. At present, this model is only supported by two seemingly disparate observations. First, *Hydractinia* mount the most vigorous and sustained rejection responses against other *Hydractinia* (62). This indicates colonies can identify the type of tissue they encounter. Second, the initial stages of rejection and fusion are morphologically indistinguishable (63). In both responses, nematocytes migrate to the point of contact and arrange their nematocysts as batteries pointed at their opponent. In rejection, the batteries fire, but in a fusion, the nematocytes migrate away as the tissues merge. This suggests rejection could be the default allorecognition response in *Hydractinia*. In this model, Alr proteins with ITAMs would activate rejection, which would be inhibited later by homophilic binding between compatible allodeterminants. This model would be analogous to the balance of ITAM and ITIM-mediated signaling that determines whether natural killer (NK) cells become activated in the vertebrate immune system. If true, it could indicate a deep evolutionary relationship between invertebrate and vertebrate self-recognition systems or, alternatively, the convergent co-opting of this signaling module for self/nonself recognition.

Our decision to classify some *Alr* sequences as putative genes relied heavily on whether the genes were expressed and correctly spliced. Two caveats are associated with our expression data. First, because the RNAseq experiment was primarily intended to guide our annotation, it did not include biological or technical replicates. The resulting expression levels should therefore be viewed as rough estimates. Second, the RNA used to generate these reads was extracted from a pool of feeding and reproductive polyps. Mat and stolon tissue—the normal sites of allorecognition responses—were not included because we and others have been unable to isolate high-quality RNA from these tissues. It is also possible that some *Alr* genes are expressed at developmental time points not represented in our dataset. Therefore, we expect to update our classification with additional data in the future.

The sequences of the Alr genes themselves do not appear to be orthologous to other invertebrate allorecognition genes. Nonetheless, this new ARC sequence reinforces similarities between *Hydractinia* and other species in which allorecognition is controlled by genomic clusters of related genes (5). This clustering may enhance the efficiency and coordination of allorecognition gene expression. It may also facilitate the generation of sequence diversity via gene conversion or unequal crossing over. Moreover, these invertebrate allorecognition complexes are remarkably similar to the complexes that control self/nonself recognition in vertebrates, namely the major histocompatibility complex (64), leukocyte receptor complex (65), and NK complex (66). Identifying evolutionary links between these systems may become possible in the future as we survey a broader swath of metazoan genomes and simultaneously deepen our molecular understanding of how invertebrate allorecognition works.

Materials and Methods

Colony 236-21 was maintained on glass microscope slides in 38-L aquaria filled with artificial seawater as previously described (67). DNA was extracted

as detailed in *SI Appendix*. PacBio and Illumina libraries were constructed and sequencing performed at the NIH Intramural Sequencing Center (NISC) via a whole-genome shotgun approach. All raw reads are available through BioProject PRJNA802249. The genome was assembled with the Celera Assembler version 8.3r2 (68) and was deposited in Zenodo (<https://zenodo.org/record/6546560>) (22). We then used NUCmer from the MUMmer package (23) to align the resulting assembly to the previously sequenced ARC BACs and merged these sequences to create a reference sequence for the ARC. Genome sequencing and assembly are detailed in *SI Appendix*.

RNA Extraction, Sequencing, and Mapping. RNA was harvested and extracted from a mixture of gastrozooids and gonozooids as detailed in *SI Appendix*. RNA-Seq libraries were constructed and sequenced at NISC as detailed in *SI Appendix*. Raw reads are available through BioProject PRJNA802249. To calculate expression levels of our annotated *Alr* genes, paired-end RNAseq reads were mapped to the entire genome assembly using HISAT2 (69). A reference-guided transcriptome was generated with Cufflinks (70). Transcript abundance was also estimated with Cufflinks with a correction for multiple read mappings as detailed in *SI Appendix*.

Annotation of *Alr* Genes and Sequence Comparisons and Analyses. *Alr* genes were annotated using Apollo (71). Methods for multiple sequence alignments, pairwise sequence alignments, sequence clustering, protein sequence annotation, and phylogenetic tree construction are provided in *SI Appendix*.

Structural Prediction and Alignment. For single-domain predictions, we generated a custom multiple sequence alignment, as detailed in *SI Appendix* which was submitted to Colabfold via the "AlphaFold2_mmseqs2" notebook,

version 1.1 (40) and were deposited in Zenodo (22). The secondary structure of each model was determined with STRIDE (47). Structure comparisons were performed with DALI (44) and PDBeFOLD (<https://www.ebi.ac.uk/msd-srv/ssm/>) (45). Models were visualized in Pymol 2.3 (72). Further details can be found in *SI Appendix*.

Data, Materials, and Software Availability. DNA sequences have been deposited in GenBank (PRJNA802249) (73). All other data are available as supplemental files or can be downloaded from Zenodo (<https://doi.org/10.5281/zenodo.6546559>) (22).

ACKNOWLEDGMENTS. We thank Leo Buss for sharing colony 236-21 and for helpful discussion. This work was funded by NSF grant 1557339 (M.L.N.), NSF grant 1923259 (C.E.S. and M.L.N.), NIH, Intramural Research Program of the National Human Genome Research Institute, grant ZIA HG000140 (A.D.B.), and NIH, Intramural Research Program of the National Human Genome Institute, grant ZIA HG200398 (A.M.P.). A.L.H. and S.M.S. were supported by NIH grant T32 AI074490.

Author affiliations: ^aThomas E. Starzl Transplantation Institute, University of Pittsburgh, Pittsburgh, PA 15261; ^bCenter for Evolutionary Biology and Medicine, University of Pittsburgh, Pittsburgh, PA 15261; ^cComputational and Statistical Genomics Branch, National Human Genome Research Institute, NIH, Bethesda, MD 20892; ^dNIH Intramural Sequencing Center, NIH, Rockville, MD 20892; ^eWhitney Laboratory for Marine Bioscience, University of Florida, St. Augustine, FL 32080; ^fDepartment of Biology, University of Florida, Gainesville, FL 32611; and ^gDepartment of Immunology, University of Pittsburgh, Pittsburgh, PA 15213

- R. D. Rosengarten, M. L. Nicotra, Model systems of invertebrate allorecognition. *Curr. Biol.* **21**, R82–R92 (2011).
- L. W. Buss, Competition within and between encrusting clonal invertebrates. *Trends Ecol. Evol.* **5**, 352–356 (1990).
- D. J. Laird, A. W. De Tomaso, I. L. Weissman, Stem cells are units of natural selection in a colonial ascidian. *Cell* **123**, 1351–1360 (2005).
- R. K. Grosberg, The evolution of allorecognition specificity in clonal invertebrates. *Q. Rev. Biol.* **63**, 377–412 (1988).
- L. F. Grice, B. M. Degan, "How to build an allorecognition system: A guide for prospective multicellular organisms" in *Evolutionary Transitions to Multicellular Life: Principles and Mechanisms*, I. Ruiz-Trillo, A. M. Nedelcu, Eds. (Springer, 2015), pp. 395–424.
- L. W. Buss, *The Evolution of Individuality* (Princeton University Press, 1987).
- F. M. Burnet, "Self-recognition" in colonial marine forms and flowering plants in relation to the evolution of immunity. *Nature* **232**, 230–235 (1971).
- L. F. Grice *et al.*, Origin and evolution of the sponge aggregation factor gene family. *Mol. Biol. Evol.* **34**, 1083–1099 (2017).
- A. W. De Tomaso *et al.*, Isolation and characterization of a protochordate histocompatibility locus. *Nature* **438**, 454–459 (2005).
- S. V. Nyholm *et al.*, *fester*, A candidate allorecognition receptor from a primitive chordate. *Immunology* **25**, 163–173 (2006).
- A. Voskoboinik *et al.*, Identification of a colonial chordate histocompatibility gene. *Science* **341**, 384–387 (2013).
- L. F. Cadavid, A. E. Powell, M. L. Nicotra, M. Moreno, L. W. Buss, An invertebrate histocompatibility complex. *Genetics* **167**, 357–365 (2004).
- A. E. Powell *et al.*, Differential effect of allorecognition loci on phenotype in *Hydractinia symbiolongicarpus* (Cnidaria: Hydrozoa). *Genetics* **177**, 2101–2107 (2007).
- M. L. Nicotra *et al.*, A hypervariable invertebrate allodeterminant. *Curr. Biol.* **19**, 583–589 (2009).
- S. F. Rosa *et al.*, *Hydractinia* allodeterminant *alr1* resides in an immunoglobulin superfamily-like gene complex. *Curr. Biol.* **20**, 1122–1127 (2010).
- A. Gloria-Soria *et al.*, Evolutionary genetics of the hydroid allodeterminant *alr2*. *Mol. Biol. Evol.* **29**, 3921–3932 (2012).
- U. B. Karadge, M. Gosto, M. L. Nicotra, Allorecognition proteins in an invertebrate exhibit homophilic interactions. *Curr. Biol.* **25**, 2845–2850 (2015).
- A. L. Huene, T. Chen, M. L. Nicotra, New binding specificities evolve via point mutation in an invertebrate allorecognition gene. *iScience* **24**, 102811 (2021).
- R. D. Rosengarten, M. A. Moreno, F. G. Lakkis, L. W. Buss, S. L. Dellaporta, Genetic diversity of the allodeterminant *alr2* in *Hydractinia symbiolongicarpus*. *Mol. Biol. Evol.* **28**, 933–947 (2011).
- H. Rodriguez-Valbuena, A. Gonzalez-Muñoz, L. F. Cadavid, Multiple *Alr* genes exhibit allorecognition-associated variation in the colonial cnidarian *Hydractinia*. *Immunogenetics*, 10.1007/s00251-022-01268-3 (2022). Correction in: *Immunogenetics*, 10.1007/s00251-022-01271-8 (2022).
- A. E. Powell *et al.*, Genetic background and allorecognition phenotype in *Hydractinia symbiolongicarpus*. *G3 (Bethesda)* **1**, 499–504 (2011).
- A. L. Huene *et al.*, *Hydractinia* strain 236-21 genome assembly and *Alr* domain predictions. Zenodo. <https://zenodo.org/record/6546560>. Deposited 13 May 2022.
- S. Kurtz *et al.*, Versatile and open software for comparing large genomes. *Genome Biol.* **5**, R12 (2004).
- K. Murphy, P. Travers, M. Walport, *Janeway's Immunobiology* (Garland Science, 2008), vol. 15, pp. 655–708.
- J. S. Bezbradica, R. Medzhitov, Role of ITAM signaling module in signal integration. *Curr. Opin. Immunol.* **24**, 58–66 (2012).
- J. S. Ziegenfuss *et al.*, Draper-dependent glial phagocytic activity is mediated by Src and Syk family kinase signalling. *Nature* **453**, 935–939 (2008).
- A. Mócsai, J. Ruland, V. L. J. Tybulewicz, The SYK tyrosine kinase: A crucial player in diverse biological functions. *Nat. Rev. Immunol.* **10**, 387–402 (2010).
- T. A. Chan *et al.*, Identification of a gene encoding a novel protein-tyrosine kinase containing SH2 domains and ankyrin-like repeats. *Oncogene* **9**, 1253–1259 (1994).
- R. E. Steele, N. A. Stover, M. Sakaguchi, Appearance and disappearance of Syk family protein-tyrosine kinase genes during metazoan evolution. *Gene* **239**, 91–97 (1999).
- L. L. Lanier, Evolutionary struggles between NK cells and viruses. *Nat. Rev. Immunol.* **8**, 259–268 (2008).
- A. D. Barrow, J. Trowsdale, You say ITAM and I say ITIM, let's call the whole thing off: The ambiguity of immunoreceptor signalling. *Eur. J. Immunol.* **36**, 1646–1653 (2006).
- U. Lorenz, SHP-1 and SHP-2 in T cells: Two phosphatases functioning at many levels. *Immunol. Rev.* **228**, 342–359 (2009).
- M. Garg, M. Wahid, F. Khan, Regulation of peripheral and central immunity: Understanding the role of Src homology 2 domain-containing tyrosine phosphatases, SHP-1 & SHP-2. *Immunobiology* **225**, 151847 (2020).
- S. D. Pauls, A. J. Marshall, Regulation of immune cell signaling by SHIP1: A phosphatase, scaffold protein, and potential therapeutic target. *Eur. J. Immunol.* **47**, 932–945 (2017).
- M. P. Thomas, C. Erneux, B. V. L. Potter, SHIP2: Structure, function and inhibition. *ChemBioChem* **18**, 233–247 (2017).
- S. El-Gebali *et al.*, The Pfam protein families database in 2019. *Nucleic Acids Res.* **47**, D427–D432 (2019).
- L. Zimmermann *et al.*, A completely reimplemented MPI bioinformatics toolkit with a new HHpred server at its core. *J. Mol. Biol.* **430**, 2237–2243 (2018).
- J.-M. Chandonia, N. K. Fox, S. E. Brenner, SCOPe: Classification of large macromolecular structures in the structural classification of proteins-extended database. *Nucleic Acids Res.* **47**, D475–D481 (2019).
- A. S. Yang, B. Honig, An integrated approach to the analysis and modeling of protein sequences and structures. II. On the relationship between sequence and structural similarity for proteins that are not obviously related in sequence. *J. Mol. Biol.* **301**, 679–689 (2000).
- M. Mirdita *et al.*, ColabFold: Making protein folding accessible to all. *Nat. Methods* **19**, 679–682 (2022).
- J. Jumper *et al.*, Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
- K. Tunyasuvunakool *et al.*, Highly accurate protein structure prediction for the human proteome. *Nature* **596**, 590–596 (2021).
- H. M. Berman *et al.*, The protein data bank. *Nucleic Acids Res.* **28**, 235–242 (2000).
- L. Holm, "Using Dali for protein structure comparison" in *Structural Bioinformatics: Methods and Protocols*, Z. Gáspári, Ed. (Springer, 2020), pp. 29–42.
- E. Krissinel, K. Henrick, Secondary-structure matching (SSM), a new tool for fast protein structure alignment in three dimensions. *Acta Crystallogr. D Biol. Crystallogr.* **60**, 2256–2268 (2004).
- Y. Harpaz, C. Chothia, Many of the immunoglobulin superfamily domains in cell adhesion molecules and surface receptors belong to a new structural set which is close to that containing variable domains. *J. Mol. Biol.* **238**, 528–539 (1994).
- M. Heining, D. Frishman, STRIDE: A web server for secondary structure assignment from known atomic coordinates of proteins. *Nucleic Acids Res.* **32**, W500–W502 (2004).
- J. P. Cannon, R. N. Haire, G. W. Litman, Identification of diversified genes that contain immunoglobulin-like variable regions in a protochordate. *Nat. Immunol.* **3**, 1200–1207 (2002).
- A. M. Lesk, C. Chothia, Evolution of proteins formed by beta-sheets. II. The core of the immunoglobulin domains. *J. Mol. Biol.* **160**, 325–342 (1982).
- J. M. Hemmingsen, K. M. Gernert, J. S. Richardson, D. C. Richardson, The tyrosine corner: A feature of most Greek key beta-barrel proteins. *Protein Sci.* **3**, 1927–1937 (1994).

51. S. J. Hamill, E. Cota, C. Chothia, J. Clarke, Conservation of folding and stability within a protein family: The tyrosine corner as an evolutionary cul-de-sac. *J. Mol. Biol.* **295**, 641–649 (2000).
52. K. C. Chou, Prediction of tight turns and their types in proteins. *Anal. Biochem.* **286**, 1–16 (2000).
53. J. S. Richardson, The anatomy and taxonomy of protein structure. *Adv. Protein Chem.* **34**, 167–339 (1981).
54. N. K. Fox, S. E. Brenner, J.-M. Chandonia, SCOPe: Structural classification of proteins—Extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res.* **42**, D304–D309 (2014).
55. J.-H. Wang, The sequence signature of an Ig-fold. *Protein Cell* **4**, 569–572 (2013).
56. D. J. Leahy, W. A. Hendrickson, I. Aukhil, H. P. Erickson, Structure of a fibronectin type III domain from tenascin phased by MAD analysis of the selenomethionyl protein. *Science* **258**, 987–991 (1992).
57. D. M. Halaby, A. Poupon, J. Mornon, The immunoglobulin fold family: Sequence analysis and 3D structure comparisons. *Protein Eng.* **12**, 563–571 (1999).
58. A. F. Williams, A. N. Barclay, The immunoglobulin superfamily—Domains for cell surface recognition. *Annu. Rev. Immunol.* **6**, 381–405 (1988).
59. A. N. Barclay, Membrane proteins with immunoglobulin-like domains—A master superfamily of interaction molecules. *Semin. Immunol.* **15**, 215–223 (2003).
60. I. R. Evans, F. S. L. M. Rodrigues, E. L. Armitage, W. Wood, Draper/CED-1 mediates an ancient damage response to control inflammatory blood cell migration in vivo. *Curr. Biol.* **25**, 1606–1612 (2015).
61. J. Sun, L. Wang, C. Yang, L. Song, An ancient BCR-like signaling promotes ICP production and hemocyte phagocytosis in oyster. *iScience* **23**, 100834 (2020).
62. R. G. Lange, M. H. Dick, W. A. Müller, Specificity and early ontogeny of historecognition in the hydroid *Hydractinia*. *J. Exp. Zool.* **262**, 307–316 (1992).
63. R. Lange, G. Plickert, W. A. Müller, Histoincompatibility in a low invertebrate, *Hydractinia echinata*: Analysis of the mechanism of rejection. *J. Exp. Zool.* **249**, 284–292 (1989).
64. J. Kaufman, Unfinished business: Evolution of the MHC and the adaptive immune system of jawed vertebrates. *Annu. Rev. Immunol.* **36**, 383–409 (2018).
65. J. Trowsdale, D. C. Jones, A. D. Barrow, J. A. Traherne, Surveillance of cell and tissue perturbation by receptors in the LRC. *Immunol. Rev.* **267**, 117–136 (2015).
66. J. Kelley, L. Walter, J. Trowsdale, Comparative genomics of natural killer cell receptor gene clusters. *PLoS Genet.* **1**, 129–139 (2005).
67. S. M. Sanders *et al.*, CRISPR/Cas9-mediated gene knockin in the hydroid *Hydractinia symbiolongicarpus*. *BMC Genomics* **19**, 649 (2018).
68. K. Berlin *et al.*, Assembling large genomes with single-molecule sequencing and locality-sensitive hashing. *Nat. Biotechnol.* **33**, 623–630 (2015).
69. D. Kim, J. M. Paggi, C. Park, C. Bennett, S. L. Salzberg, Graph-based genome alignment and genotyping with HISAT2 and HISAT-genotype. *Nat. Biotechnol.* **37**, 907–915 (2019).
70. C. Trapnell *et al.*, Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat. Biotechnol.* **28**, 511–515 (2010).
71. N. A. Dunn *et al.*, Apollo: Democratizing genome annotation. *PLoS Comput. Biol.* **15**, e1006790 (2019).
72. Schrödinger, LLC, The PyMOL Molecular Graphics System (Version 2.3, 2020).
73. A. L. Huene *et al.*, *Hydractinia symbiolongicarpus* strain 236-21 ARC sequencing and annotation. NCBI BioProject. <https://www.ncbi.nlm.nih.gov/bioproject?term=PRJNA802249>. Accessed 19 September 2022.