

**Supporting Information Appendix for**

**Expansions, diversification and inter-individual copy number variations of  
AID/APOBEC family cytidine deaminase genes in lampreys**

Stephen J. Holland, Lesley M. Berghuis, Justin J. King, Lakshminarayan M. Iyer,  
Katarzyna Sikora, Heather Fifield, Sarah Peter, Emma M. Quinlan, Fumiaki Sugahara,  
Prashant Shingate, Inês Trancoso, Norimasa Iwanami, Elena Temereva, Christine  
Strohmeier, Shigeru Kuratani, Byrappa Venkatesh, Guillaume Evanno, L. Aravind,  
Michael Schorpp, Mani Larijani, and Thomas Boehm

*Lampetra planeri/Lampetra fluviatilis*

*L. planeri*

Lp#236

CDA1	
CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	
CDA1L2_1	
CDA1L2_2	

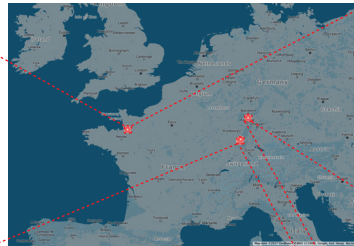
*L. fluviatilis*

CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	

Lf#29 (PCR)

CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	

Lf#33 (PCR)



*Lampetra planeri*

*Petromyzon marinus*

Lp#173

Lp#196 (PCR)

CDA1	
CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	
CDA1L2_1	
CDA1L2_2	

Lp#242

CDA1	
CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	
CDA1L2_1	
CDA1L2_2	

Lp#175

CDA1	
CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	
CDA1L2_1	
CDA1L2_2	

Pm#144

CDA1	
CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	
CDA1L2_1	
CDA1L2_2	

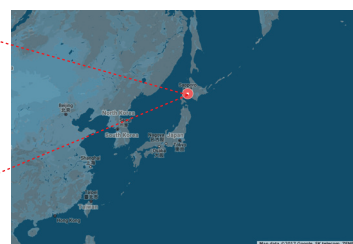
Pm#1

CDA1	
CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	
CDA1L2_1	
CDA1L2_2	

*Lethenteron japonicum*

Lj#1

CDA1	
CDA1L1_1	
CDA1L1_2	
CDA1L1_3	
CDA1L1_4	
CDA1L2_1	
CDA1L2_2	



**Fig. S1.** *CDAI*-like gene copy number in lamprey varies based on geographic location. Geographic locations of capture site for *Lampetra planeri* and *Lampetra fluviatilis* specimens analyzed here. The presence of individual *CDAILL1* genes is indicated by green color, the presence of *CDAILL2* genes by blue color; absence (or presence of pseudogenized versions) of a gene is indicated by grey color. For comparison, the *CDA* gene complements of *Petromyzon marinus* and *Lethenteron japonicum* specimens are also shown. Further details are given in Table S1.

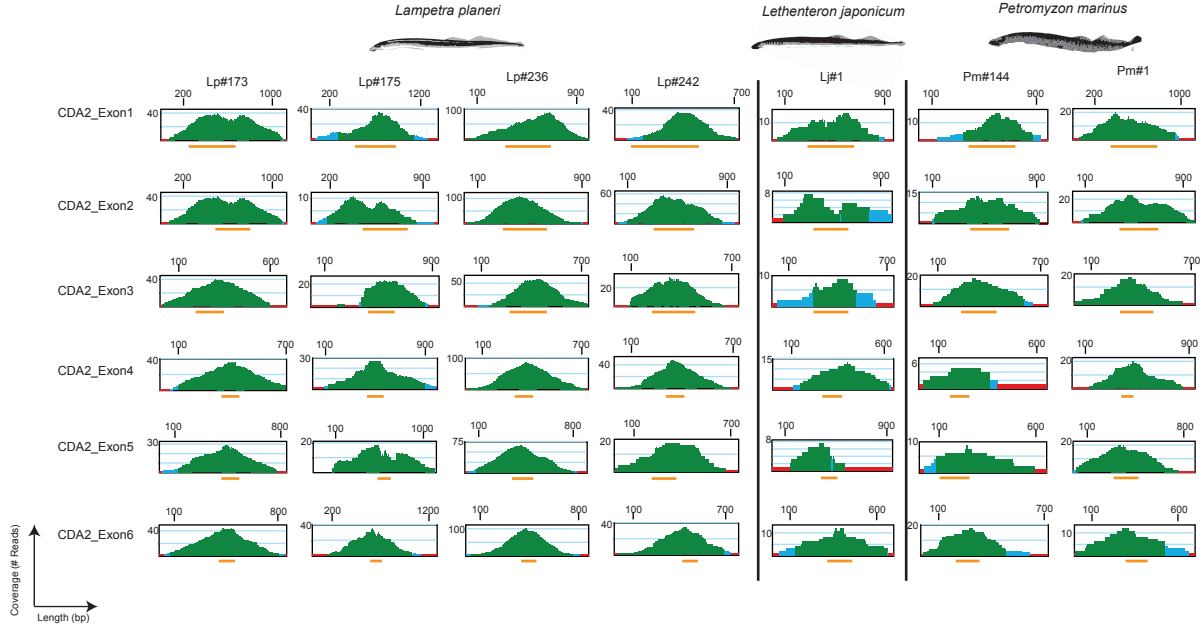
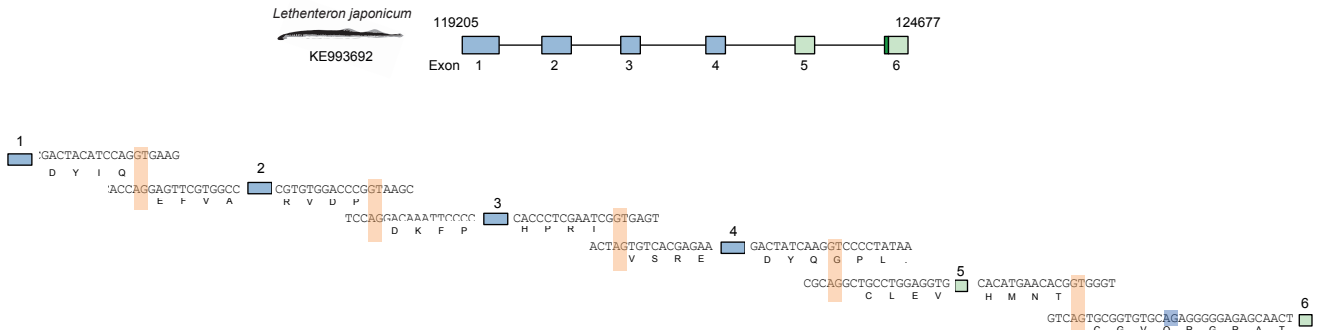
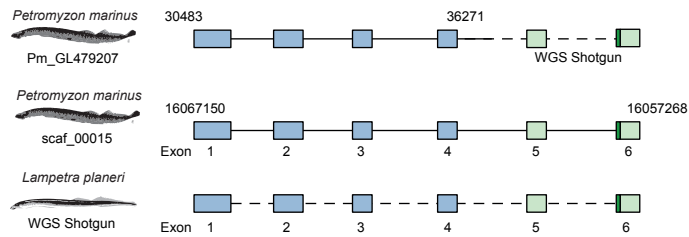
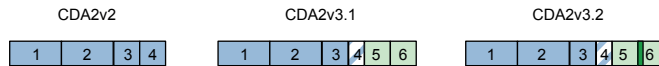
**A**

Lp\_CDA2\_Exon1 MELREVVDALGSCVRHEPLGRAAFILRCFAAPSRKPRGTVILFDVDGAGRGLSGGHVVNNKQGTSIHAEVLLLSAVRAALPQR--CEGDAEEAPRGCTVHCYSTYSPCRDCVDYIQ  
 Lj\_CDA2\_Exon1 MELREVVDALGSCVRHEPLGRAAFILRCFAAPSRKPRGTVILFDVDGAGRGLSGGHVVNNKQGTSIHAEVLLLSAVRAALPQR--CEGDAEEAPRGCTVHCYSTYSPCRDCVDYIQ  
 Pm\_CDA2\_Exon1 MELREVVDALGSCVRHEPLGRAAFILRCFAAPSRKPRGTVILFDVDGAGRGLSGGHVVNNKQGTSIHAEVLLLSAVRAALLRRRRC--DGEATRGTCTVHCYSTYSPCRDCVEYIQ

Lp\_CDA2\_Exon2 EFVASTGVRVAMRCCRLYELDVTRRRPEAEGLRSLSLGRDFRLMGRDAIALLGGRLA---DGEASGSG-----GDAEPLVEMAGFGDEQLHAQVQRNRQIVEAYARYAGAVSLVLGELRVPD  
 Lj\_CDA2\_Exon2 EFVASTGVRVAMRCCRLYELDVTRRRPEAEGLRSLSLGRDFRLMGRDAIALLGGRLA---DGEASGSG-----GNAEPLVEMAGFGDEQLHAQVQRNRQIVEAYARYAGAVSLVLGELRVPD  
 Pm\_CDA2\_Exon2 EFGASTGVRVVIHCCRLYELDVNRRSEAEGLRSLSLGRDFRLMGRDAIALLGGRLANTADGESGASGNAWTETNVVEPLVDMTGFDEDLHAQVQRNKQIREAYANYASAVSLMLGELHVPD

Lp\_CDA2\_Exon3 DKPFFLADFLAQTSEVPSGTPRGARGPRGASSRGGPGIGRQRPADFERALGAYGLFLHPR  
 Lj\_CDA2\_Exon3 DKPFFLADFLAQTSEVPSGTPRGARGPRGASSRGGPGIGRQRPADFERALGAYGLFLHPR  
 Pm\_CDA2\_Exon3 DKPFFLAEFLAQTSEVPSGTPRETGRGPRGASSRGGPEIGRQRPADFERALGAYGLFLHPR

Lp\_CDA2\_Exon4 VSREADREEIKRDLIVAMRKHNYQGPL.  
 Lj\_CDA2\_Exon4 VSREADREEIKRDLIVAMRKHNYQGPL.  
 Pm\_CDA2\_Exon4 VSREADREEIKRDLIVAMRKHNYQGPL.

**B****C****D****E**

PmCDA2 MELREVVDALGSCVRHEPLSRVAFILRCFAAPSRKPRGTVILFDVDGAGRGLSGGHVVNNKQGTSIHAEVLLLSAVRAALLRRRRC--DGEATRGTCTVHCYSTYSPCRDCVEYIQEFVASTGVRVVIHCCRLYELDVT  
 LpCDA2\_v2 MELREVVDALGSCVRHEPLGRAAFILRCFAAPSRKPRGTVILFDVDGAGRGLSGGHVVNNKQGTSIHAEVLLLSAVRAALPQR--CEGDAEEAPRGCTVHCYSTYSPCRDCVDYIQEFVASTGVRVAMRCCRLYELDVT  
 LpCDA2\_v3.1 MELREVVDALGSCVRHEPLGRAAFILRCFAAPSRKPRGTVILFDVDGAGRGLSGGHVVNNKQGTSIHAEVLLLSAVRAALPQR--CEGDAEEAPRGCTVHCYSTYSPCRDCVDYIQEFVASTGVRVAMRCCRLYELDVT  
 LpCDA2\_v3.2 MELREVVDALGSCVRHEPLGRAAFILRCFAAPSRKPRGTVILFDVDGAGRGLSGGHVVNNKQGTSIHAEVLLLSAVRAALPQR--CEGDAEEAPRGCTVHCYSTYSPCRDCVDYIQEFVASTGVRVAMRCCRLYELDVT

PmCDA2 RRRSEAEGLRSLSLGRDFRLMGRDAIALLGGRLANTADGESGASGNAWTETNVVEPLVDMTGFDEDLHAQVQRNKQIREAYANYASAVSLMLGELHVPD  
 LpCDA2\_v2 RRRPEAEGLRSLSLGRDFRLMGRDAIALLGGRLA---DGEASGSG-----GDAEPLVEMAGFGDEQLHAQVQRNRQIVEAYARYAGAVSLVLGELRVPD  
 LpCDA2\_v3.1 RRRPEAEGLRSLSLGRDFRLMGRDAIALLGGRLA---DGEASGSG-----GDAEPLVEMAGFGDEQLHAQVQRNRQIVEAYARYAGAVSLVLGELRVPD  
 LpCDA2\_v3.2 RRRPEAEGLRSLSLGRDFRLMGRDAIALLGGRLA---DGEASGSG-----GDAEPLVEMAGFGDEQLHAQVQRNRQIVEAYARYAGAVSLVLGELRVPD

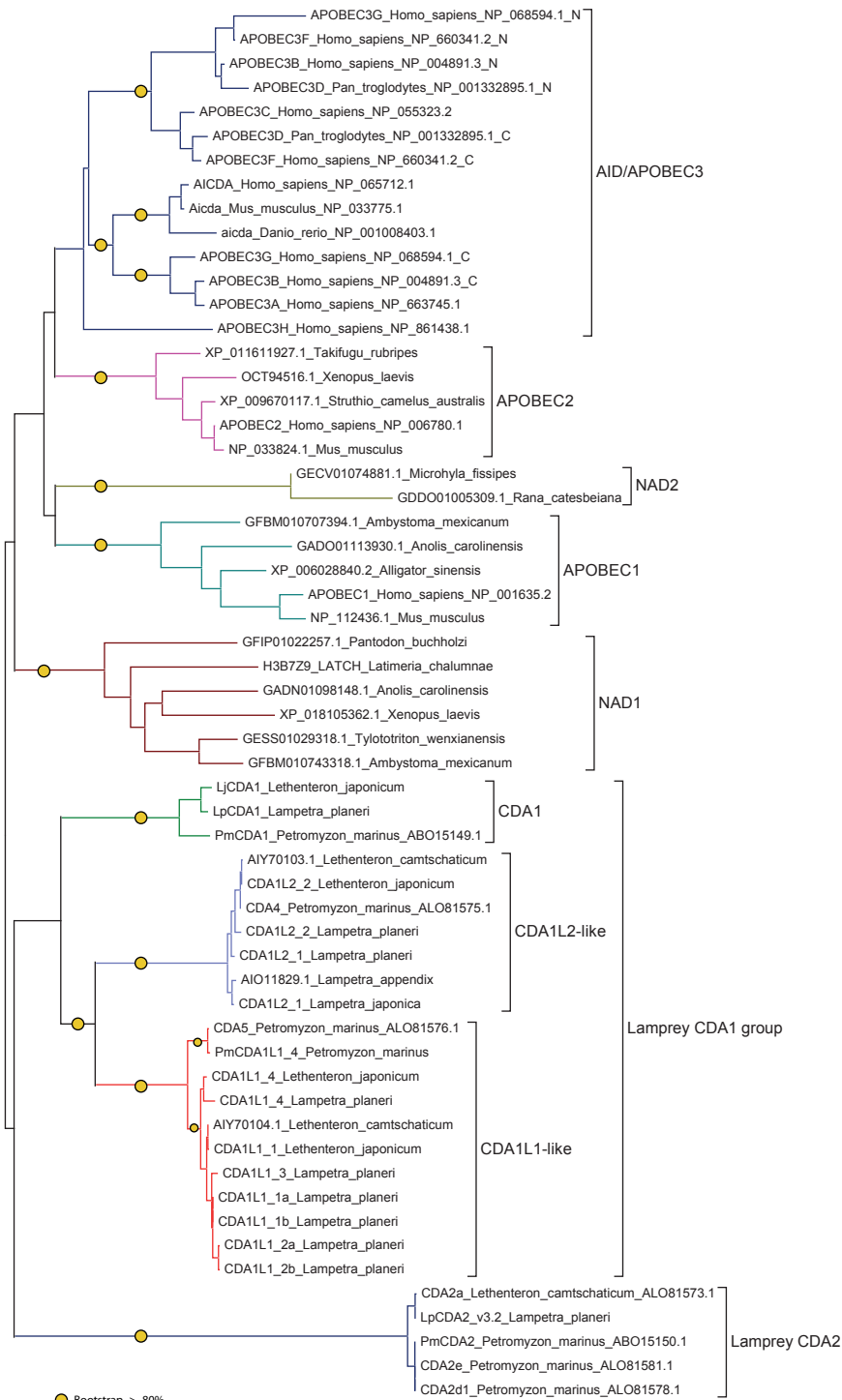
PmCDA2 PEIGRQRPADFERALGAYGLFLHPRIVSREADREEIKRDLIVAMRKHNYQGPL.  
 LpCDA2\_v2 PEIGRQRPADFERALGAYGLFLHPRIVSREADREEIKRDLIVAMRKHNYQGPL.  
 LpCDA2\_v3.1 PEIGRQRPADFERALGAYGLFLHPRIVSREADREEIKRDLIVAMRKHNYQGPL.  
 LpCDA2\_v3.2 PEIGRQRPADFERALGAYGLFLHPRIVSREADREEIKRDLIVAMRKHNYQGPL.

PmCDA2  
 LpCDA2\_v2  
 LpCDA2\_v3.1 P.  
 LpCDA2\_v3.2 P.

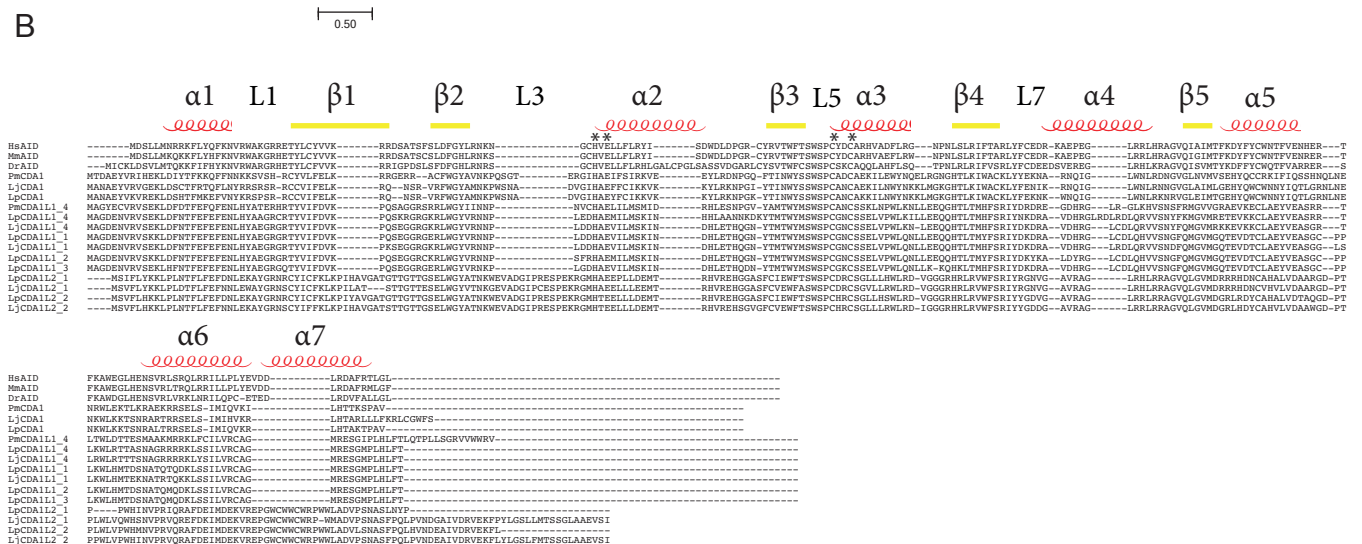
**Fig. S2.** Identification of novel splice variants of *CDA2* genes in *Lampetra planeri* and *Lethenteron japonicum*. (A) Comparison of translated *CDA2* exon sequences identified in the genomes of *Lampetra planeri* and *Lethenteron japonicum* with those previously characterized in *Petromyzon marinus*. Conserved APOBEC catalytic HxE and PCxxC motifs are highlighted. (B) Read coverage plots retrieved from whole genome sequences using the constituent *Petromyzon marinus* *CDA2* exon sequences as queries. Green color indicates a coverage by >5 reads; blue, 2-5 reads, and red, single read. Orange bars correspond to region of the contigs containing the open reading frames of the exons. Sequences corresponding to the 4 exons of *PmCDA2* were readily identifiable in the genomic sequence collections of all four individuals. Notably, the exon sequences of *CDA2* genes of Pm#1 and Pm#144 were identical to the previously described *Petromyzon marinus* sequence and clearly distinguishable from the corresponding *Lampetra planeri* and *Lethenteron japonicum* *CDA2* gene sequences (see panel A). Similarly, the exon sequences of *CDA2* genes from Lp#236 and Lp#242 were identical to those identified in Lp#173 and Lp#175. (C) Schematic of *CDA2* exon structure as established by the *Lethenteron japonicum* genome assembly (version LetJap 1.0) (top panel); the exon/intron junctions for each exon are indicated (lower panel). Splice donor/acceptor sites are highlighted in orange; the alternative splice acceptor site inside exon 6 is highlighted in blue. (D) Schematic of *CDA2* exon structure in *Petromyzon marinus*; in the genome assembly (version 7.0)(47) the first 4 exons of *CDA2* are present on a single scaffold (Pm\_GL479207, solid line), whereas sequences corresponding to exons 5 and 6 are not contained within this scaffold, but could be identified in our sequence collection obtained with a different *Petromyzon marinus* individual, Pm#1 (dashed line). A recent genome assembly from germline DNA of *P. marinus* (62) confirmed the presence of all 6 exons on a single scaffold (scaf\_00015). In the shotgun genome libraries of *Lampetra planeri* individuals, all exons could be identified but were not assembled into a contig (dashed line); the order of exons was established by cDNA sequences. (E) Alternatively spliced variants of *CDA2* genes are shown as schematic (top panel) and as amino acid alignments (lower panel). The striped part of exon 4 is lacking in *CDA2v3.1* and *CDA2v3.2* variants as a result of intra-exonic splicing between exons 4 and 5. The solid dark green box in *CDA2v3.2* refers to the presence of an additional segment encoding 4 amino

acids as a result of the use of a different splice acceptor signal (see panel C). Pm, *Petromyzon marinus*, Lp, *Lampetra planeri*.

A



B



**Fig. S3.** Sequence comparisons of AID/APOBEC deaminases. (A) Maximum-likelihood tree of lamprey CDAs and jawed vertebrate AID/APOBEC deaminases. Nodes with bootstrap support greater than 80% are marked. NAD1 and NAD2 are novel AID/APOBEC-like deaminases, which are discussed in the companion paper. For lamprey *CDAI*-like genes, related Genbank database entries were also used; small letters at the end of gene names correspond to presumptive allelic variants. (B) Sequence alignment of predicted protein sequences encoded by *CDAI*, *CDAIL1* and *CDAIL2* genes in *Lampetra planeri* (Lp), *Lethenteron japonicum* (Lj), and *Petromyzon marinus* (Pm). Conserved APOBEC catalytic glutamic acid and zinc-coordinating and cysteine and histidine residues are highlighted with an asterisk (\*) (c.f., Fig. 2A).



A

```

Lp_CDA11_1      MAGDENVRVSKLDFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRGERLNGYVRRNPLDDHAEVIIMSKINDH--ETHQGNVTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTMFYSRIYDKDR
LpCDA11_1_PCR_#196  -----VRRNPLDDHAEVIIMSKINDH--ETHQGNVTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTMFYSRIYDKDR
Lp_CDA11_2      MAGDENVRVSKLDFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRCKRLMGYVRRNPSFRHAEMIIMSKINDH--ETHQGNVTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTMFYSRIYDKYK
LpCDA11_2_PCR_#196  -----SFRHAEMIIMSKINDH--ETHQGNVTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTMFYSRIYDKYK
Lp_CDA11_3      MAGDENVRVSEKLFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRGERLNGYVRRNPLDDHAEVIIMSKINDH--ETHQGNVTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTMFYSRIYDKDR
LpCDA11_3_PCR_#196  -----LKKQ--KLTMHFSRIYDKDR
Lp_CDA11_4      MAGDENVRVSEKLFNTFEFENLHYAAGRCRTYVIFDVKPSQKRGKRLMGYVRRNPLEDHAEMIIMSKINHHLAANNKDKYMTWYMSWSPCGNCSSSELPWLKILLEEQQHTLTMFYSRIYDKDR
LpCDA11_4_PCR_#196  -----AANNKDKYMTWYMSWSPCGNCSSSELPWLKILLEEQQHTLTMFYSRIYDKDR

Lp_CDA11_1      AVDHRGLCDL--QHVVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
LpCDA11_1_PCR_#196  -----QHVVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp_CDA11_2      ALDYHGLCDL--QHVVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
LpCDA11_2_PCR_#196  -----QHVVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp_CDA11_3      AVDHRGLCDL--QHVVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
LpCDA11_3_PCR_#196  -----QHVVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp_CDA11_4      AVDHRGLRDLRDLQRVSNYFKMGVMRETEVKKCLAEYVEASRR--LKWLRRTTASNAGRRLKLSLILVRCAGMRESGMPHLFT.
LpCDA11_4_PCR_#196  -----LKWLRRTTASNAGRRLKLSLILVRCAGMRESGMPHLFT.
  
```

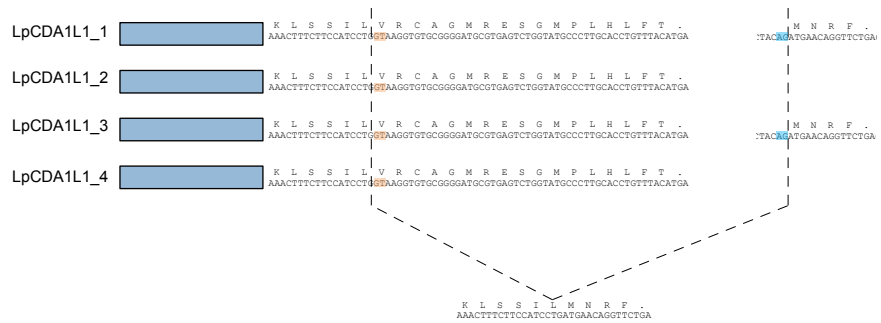
B

```

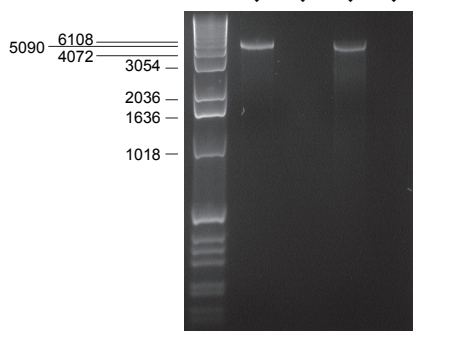
Lp#173_CDA11_1  MAGDENVRVSKLDFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRGERLNGYVRRNPLDDHAEVIIMSKINDHLETHQGN-YTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTM
Lj#1_CDA11_1    MAGDENVRVSKLDFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRGERLNGYVRRNPLDDHAEVIIMSKINDHLETHQGN-YTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTM
Lp#173_CDA11_2  MAGDENVRVSKLDFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRCKRLMGYVRRNPSFRHAEMIIMSKINDHLETHQGN-YTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTM
Lp#175_CDA11_2  MAGDENVRVSKLDFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRCKRLMGYVRRNPSFRHAEMIIMSKINDHLETHQGN-YTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTM
Lp#242_CDA11_2  MAGDENVRVSKLDFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRCKRLMGYVRRNPSFRHAEMIIMSKINDHLETHQGN-YTMTWYMSWSPCGNCSSSELPWLQNLLEEQQHTLTM
Lp#173_CDA11_3  MAGDENVRVSEKLFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRGERLNGYVRRNPLDDHAEVIIMSKINDHLETHQDN-YTMTWYMSWSPCGNCSSSELPWLQNLK-KQKILTM
Lp#242_CDA11_3  MAGDENVRVSEKLFNTFEFENLHYAEGRGRTYVIFDVKPSQEGGRGERLNGYVRRNPLDDHAEVIIMSKINDHLETHQDN-YTMTWYMSWSPCGNCSSSELPWLQNLK-KQKILTM
Lp#173_CDA11_4  MAGDENVRVSEKLFNTFEFENLHYAAGRCRTYVIFDVKPSQKRGKRLMGYVRRNPLEDHAEMIIMSKINHHLAANNKDKYMTWYMSWSPCGNCSSSELPWLKILLEEQQHTLTM
Lp#236_CDA11_4  MAGDENVRVSEKLFNTFEFENLHYAAGRCRTYVIFDVKPSQKRGKRLMGYVRRNPLEDHAEMIIMSKINHHLAANNKDKYMTWYMSWSPCGNCSSSELPWLKILLEEQQHTLTM
Lj#1_CDA11_4    MAGDENVRVSEKLFNTFEFENLHYAAGRCRTYVIFDVKPSQKRGKRLMGYVRRNPLEDHAEVIIMSKINDHLETHQGN-YTMTWYMSWSPCGNCSSSELPWLKILLEEQQHTLTM
Fm#14_CDA11_4  MAGYECVVRVSEKLFDTFEPQENLHYATERHRTYVIFDVKPSQAGGRSRLNGYIINNPNVCHAEILIMSMIDRHLESNFGV-YAMTWYMSWSPCGNCSSSELPWLKILLEEQQHTLTM

Lp#173_CDA11_1  YFSRIYDKDRAVDHRG---LCDLQHVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lj#1_CDA11_1    HFSRIYDKDRAVDHRG---LRLDLQRVSNDFQMGVMGQTEVDTCCLAEYVEASGGLSLKWLHMTDKNATRTQKLSLILVRCAGMRESGMPHLFT.
Lp#173_CDA11_2  YFSRIYDKYKALDYRG---LCDLQHVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp#175_CDA11_2  YFSRIYDKYKALDYRG---LCDLQHVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp#242_CDA11_2  YFSRIYDKYKALDYRG---LCDLQRVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp#173_CDA11_3  HFSRIYDKDRAVDHRG---LCDLQHVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp#242_CDA11_3  HFSRIYDKDRAVDHRG---LCDLQHVSNQFQMGVMGQTEVDTCCLAEYVEASGCPPLKWLHMTDSNATQTDKLSLILVRCAGMRESGMPHLFT.
Lp#173_CDA11_4  HFSRIYDKDRAVDHRGLRDLRDLQRVSNYFKMGVMRETEVKKCLAEYVEASRR--TLKWLRTTASNAGRRLKLSLILVRCAGMRESGMPHLFT.
Lp#175_CDA11_4  HFSRIYDKDRAVDHRGLRDLRDLQRVSNYFKMGVMRETEVKKCLAEYVEASRR--TLKWLRTTASNAGRRLKLSLILVRCAGMRESGMPHLFT.
Lp#236_CDA11_4  HFSRIYDKDRAVDHRGLRDLRDLQRVSNYFKMGVMRETEVKKCLAEYVEASRR--TLKWLRTTASNAGRRLKLSLILVRCAGMRESGMPHLFT.
Lj#1_CDA11_4    HFSRIYDKDRAVDHRG---LCDLQRVSNYFKMGVMRETEVKKCLAEYVEASRR--TLKWLRTTASNAGRRLKLSLILVRCAGMRESGMPHLFT.
Fm#14_CDA11_4  HFSRIYDRDREDDHRG---LRLGLKHSNDFRQMGVMRETEVKKCLAEYVEASRR--TLWLDRTTASNAGRRLKLSLILVRCAGMRESGMPHLFT.
  
```

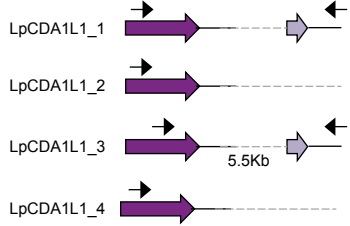
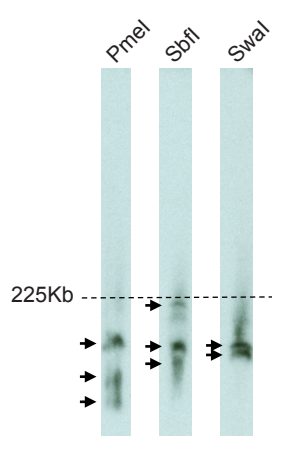
C



D



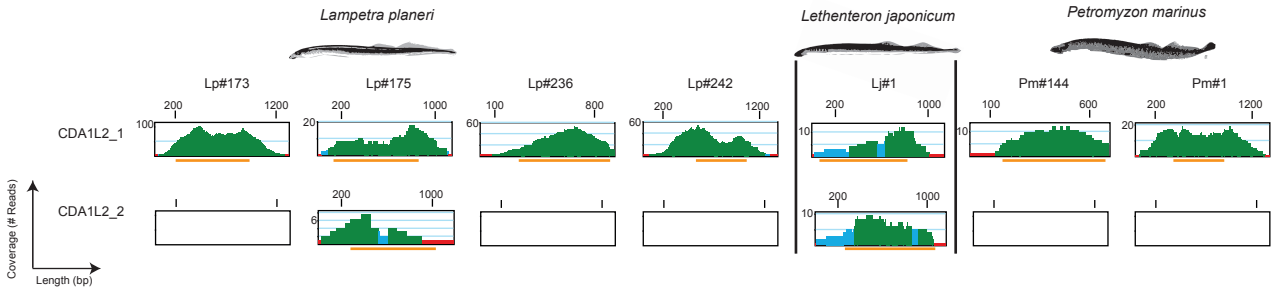
E



**Fig. S4.** Conservation of *CDAILL1\_4* gene and alternative splicing of *CDAILL1\_1* and 3.

(A) Validation of gene-specific PCR assays. Amino acid sequences deduced from whole genome sequences of *Lampetra planeri* #173 correspond perfectly to sequences derived from PCR products generated using version specific forward primers designed for each gene using genomic DNA of individual Lp#196. Primer binding sites are highlighted in red. (B) Comparison of deduced *CDAILL1* amino acid sequences derived from whole genome sequences of *Lampetra planeri* (Lp#173, 175, 236, 242), and *Lethenteron japonicum* (Lj#1) individuals. A *CDAILL1\_4* gene was also identified in a *Petromyzon marinus* individual (Pm#144) that lacked the canonical *PmCDAI* gene (c.f., Fig. 1). As is the case for many *Lampetra planeri* individuals, the Lp#236 shotgun genome sequences lacked detectable *PmCDAI* sequences, but contained a single *CDAILL1* gene, *CDAILL1\_4*; this result was confirmed by Southern filter hybridization and PCR analyses (c.f., Fig. 1). Indeed, additional variations in *CDAILL1* gene content were revealed upon further study; the genomes of *Lampetra planeri* individuals Lp#8a and Lp#8b contained three and two *CDAILL1* genes, respectively see (c.f., Fig. 1). One additional combination of three *CDAILL1* genes was observed in a specimen of another lamprey species, *Lampetra fluviatilis* (Lf#33) (Table S3); *Lampetra planeri* and *Lampetra fluviatilis* are considered to be highly related paired species (46). (C) Schematic depicting intra-exonic splicing of *CDAILL1\_1* and *CDAILL1\_3* genes. Splice donor and acceptor sequences are highlighted in pink and blue, respectively. (D) Long-range PCR of genomic DNA using *CDAILL1* gene-specific forward primers (see panel A) and a reverse primer designed to anneal to the 3'-UTR of the exon. The identity of the PCR products was confirmed by sequencing the 5'- and 3'-ends, resulting in the deduced genomic arrangement shown at the bottom of the panel. (E) Representative Southern filter hybridizations after separation of genomic restriction digests of genomic DNA of *Lampetra planeri* on pulsed-field gels, hybridized with the *CDAILL1\_4* probe (Accession MG495256). The position of the lowest size marker (225Kb) is indicated with a dashed line; fragments hybridizing to the probe are indicated by arrows. The autoradiographic images shown are taken from different parts of the same film.

A



B

```

#173_CDA1L2_1_Lp ATGTCGATCTTCCTTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
#175_CDA1L2_1_Lp ATGTCGATCTTCCTTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
#236_CDA1L2_1_Lp ATGTCGATCTTCCTTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
#175_CDA1L2_2_Lp ATGTCGGTCTTCCTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
#242_CDA1L2_1_Lp ATGTCGATCTTCCTTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
Lj#1_CDA1L2_1_Lj ATGTCGGTCTTCCTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
Lj#1_CDA1L2_2_Lj ATGTCGGTCTTCCTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
#144_CDA1L2_1_Pm ATGTCGGTCTTCCTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC
Pm#1_CDA1L2_1_Pm ATGTCGGTCTTCCTTACAAGAAGCTGCCCTCAACACGTTTCTCTGGAGTTGCAACCTCGAGAAGGCGTACGGAAGGAACAGATGCTACATTTGCTCAAGCTCAAACCCATCCAC

#173_CDA1L2_1_Lp GCGCTGGCGCCACCGGCACCCCGGCACACAGGATCCGAGCTCTGGGGTACGCCACCAACAAGTGGGAGTCCGCCAGGCATCCACCGGAGGCCCGGAGAAGCGCGGCATGCAC
#175_CDA1L2_1_Lp GCGCTGGCGCCACCGGCACCCCGGCACACAGGATCCGAGCTCTGGGGTACGCCACCAACAAGTGGGAGTCCGCCAGGCATCCACCGGAGGCCCGGAGAAGCGCGGCATGCAC
#236_CDA1L2_1_Lp GCGCTGGCGCCACCGGCACCCCGGCACACAGGATCCGAGCTCTGGGGTACGCCACCAACAAGTGGGAGTCCGCCAGGCATCCACCGGAGGCCCGGAGAAGCGCGGCATGCAC
#175_CDA1L2_2_Lp GCGCTGGCGCCACCGGCACCCCGGCACACAGGATCCGAGCTCTGGGGTACGCCACCAACAAGTGGGAGTCCGCCAGGCATCCACCGGAGGCCCGGAGAAGCGCGGCATGCAC
#242_CDA1L2_1_Lp GCGCTGGCGCCACCGGCACCCCGGCACACAGGATCCGAGCTCTGGGGTACGCCACCAACAAGTGGGAGTCCGCCAGGCATCCACCGGAGGCCCGGAGAAGCGCGGCATGCAC
Lj#1_CDA1L2_1_Lj ----TCG---CCACGACACCCCGGCACACAGGATCCGAGCTCTGGGGTACGCCACCAACAAGTGGGAGTCCGCCAGGCATCCACCGGAGGCCCGGAGAAGCGCGGCATGCAC
Lj#1_CDA1L2_2_Lj GCGCTGGCGCCACCGGCACCCCGGCACACAGGATCCGAGCTCTGGGGTACGCCACCAACAAGTGGGAGTCCGCCAGGCATCCACCGGAGGCCCGGAGAAGCGCGGCATGCAC
#144_CDA1L2_1_Pm GCGCTGGCGCCACCGGCACCCCGGC-----GATMCGAGCTCTGGGGTACGCCACCGGCAGCGAGGTCGCCGGCGCAGCCCGGAGGCCCGGAGAAGCGCGGCATGCAC
Pm#1_CDA1L2_1_Pm GCGCTGGCGCCACCGGCACCCCGGC-----GATMCGAGCTCTGGGGTACGCCACCGGCAGCGAGGTCGCCGGCGCAGCCCGGAGGCCCGGAGAAGCGCGGCATGCAC

#173_CDA1L2_1_Lp GCCGAGGAGCCGCTGCTGGATGAGATGACACGCCACGCTCCCGGAGCAGCGCG--CGCCAGCTTCTGATCGAGTGGT-TCACGTCGTGGAGCCCTCGCACCGCTGCTCGGGGCTGCTGC
#175_CDA1L2_1_Lp GCCGAGGAGCCGCTGCTGGATGAGATGACACGCCACGCTCCCGGAGCAGCGCG--CGCCAGCTTCTGATCGAGTGGT-TCACGTCGTGGAGCCCTCGCACCGCTGCTCGGGGCTGCTGC
#236_CDA1L2_1_Lp GCCGAGGAGCCGCTGCTGGATGAGATGACACGCCACGCTCCCGGAGCAGCGCG--CGCCAGCTTCTGATCGAGTGGT-TCACGTCGTGGAGCCCTCGCACCGCTGCTCGGGGCTGCTGC
#175_CDA1L2_2_Lp GCCGAGGAGCCGCTGCTGGATGAGATGACACGCCACGCTCCCGGAGCAGCGCG--CGCCAGCTTCTGATCGAGTGGT-TCACGTCGTGGAGCCCTCGCACCGCTGCTCGGGGCTGCTGC
#242_CDA1L2_1_Lp GCCGAGGAGCCGCTGCTGGATGAGATGACACGCCACGCTCCCGGAGCAGCGCG--CGCCAGCTTCTGATCGAGTGGT-TCACGTCGTGGAGCCCTCGCACCGCTGCTCGGGGCTGCTGC
Lj#1_CDA1L2_1_Lj GCCGAGGAGTCTGCTGGAGAGATGACACGCCACGCTCCCGGAGCAGCGCG--CGCCAGCTTCTGATCGAGTGGT-TCGCGTCGTGGAGCCCTCGCACCGCTGCTCGGGGCTGCTGC
Lj#1_CDA1L2_2_Lj GCCGAGGAGTCTGCTGGATGAGATGACACGCCACGCTCCCGGAGCAGCGCG--CGCTGGCTTCTGATCGAGTGGT-TCACGTCGTGGAGCCCTCGCACCGCTGCTCGGGGCTGCTGC
#144_CDA1L2_1_Pm GCCGGGAGCTCCGCTGGAGAGGTCGACGCCACGCTCCCGGAGCAGCGCGAGCCTTCTGCTGAGCGGTATCCCGTCTGGAGCCCGGAGCAGCGCCCTCGGGGCTGCGCC
Pm#1_CDA1L2_1_Pm GCCGGGAGCTCCGCTGGAGAGGTCGACGCCACGCTCCCGGAGCAGCGCGAGCCTTCTGCTGAGCGGTATCCCGTCTGGAGCCCGGAGCAGCGCCCTCGGGGCTGCGCC

#173_CDA1L2_1_Lp TCCACTGGCTGCGCGAGCTGCGCGCGCGGCACACCGGCTGCGGCTCTGGTCTCCCGAATCTACCGCGGAACCTCGGGGCGTGGCGCCGCGTGCCTGCCTCACTACCGCCGCGCGGGG
#175_CDA1L2_1_Lp TCCACTGGCTGCGCGAGCTGCGCGCGCGGCACACCGGCTGCGGCTCTGGTCTCCCGAATCTACCGCGGAACCTCGGGGCGTGGCGCCGCGTGCCTGCCTCACTACCGCCGCGCGGGG
#236_CDA1L2_1_Lp TCCACTGGCTGCGCGAGCTGCGCGCGCGGCACACCGGCTGCGGCTCTGGTCTCCCGAATCTACCGCGGAACCTCGGGGCGTGGCGCCGCGTGCCTGCCTCACTACCGCCGCGCGGGG
#175_CDA1L2_2_Lp TCCACTGGCTGCGCGAGCTGCGCGCGCGGCACACCGGCTGCGGCTCTGGTCTCCCGAATCTACTACGGAACCTCGGGGCGTGGCGCCGCGTGCCTGCCTCGCCGCGCGGGG
#242_CDA1L2_1_Lp TCCACTGGCTGCGCGAGCTGCGCGCGCGGCACACCGGCTGCGGCTCTGGTCTCCCGAATCTACCGCGGAACCTCGGGGCGTGGCGCCGCGTGCCTGCCTCACTACCGCCGCGGGG
Lj#1_CDA1L2_1_Lj TCCGTTGGCTGCGCGAGCTGCGCGCGCGGCACACCGGCTGCGGCTCTGGTCTCCCGAATCTACCGCGGAACCTCGGGGCGTGGCGCCGCGTGCCTGCCTCACTACCGCCGCGGGG
Lj#1_CDA1L2_2_Lj TCCGTTGGCTGCGCGAGCTGCGCGCGCGGCACACCGGCTGCGGCTCTGGTCTCCCGAATCTACTACGGAACCTCGGGGCGTGGCGCCGCGTGCCTGCCTCACTACCGCCGCGGGG
#144_CDA1L2_1_Pm GCGCGTGGCTGCGCGAGCTGCGCGCGCGGCACCGGCTGCGGCTCTGGTCTCCCGCATGACC-----GATMCGAGCTCTGGGGTACGCCACCGGCAGCGAGGTCGCCGGCGCAGCCCGGAGGCCCGGAGAAGCGCGGCATGCAC
Pm#1_CDA1L2_1_Pm GCGCGTGGCTGCGCGAGCTGCGCGCGCGGCACCGGCTGCGGCTCTGGTCTCCCGCATGACC-----GATMCGAGCTCTGGGGTACGCCACCGGCAGCGAGGTCGCCGGCGCAGCCCGGAGGCCCGGAGAAGCGCGGCATGCAC

#173_CDA1L2_1_Lp *****
#175_CDA1L2_1_Lp *****
#236_CDA1L2_1_Lp *****
#175_CDA1L2_2_Lp *****
#242_CDA1L2_1_Lp *****
Lj#1_CDA1L2_1_Lj *****
Lj#1_CDA1L2_2_Lj *****
#144_CDA1L2_1_Pm *****
Pm#1_CDA1L2_1_Pm *****

#173_CDA1L2_1_Lp TCGAGCTAGCGCTGATGGACAGAGGGCGGCACGATCTGCGCGCACCGCTTGGTGGAGCGCGCGCGGGCGATCCACGCGCC-----CGTGGCACATAAAGCTCCCCGCA
#175_CDA1L2_1_Lp TCGAGCTAGCGCTGATGGACAGAGGGCGGCACGATCTGCGCGCACCGCTTGGTGGAGCGCGCGCGGGCGATCCACGCGCC-----CGTGGCACATAAAGCTCCCCGCA
#236_CDA1L2_1_Lp TCGAGCTAGCGCTGATGGACAGAGGGCGGCACGATCTGCGCGCACCGCTTGGTGGAGCGCGCGCGGGCGATCCACGCGCC-----CGTGGCACATAAAGCTCCCCGCA
#175_CDA1L2_2_Lp TCGAGCTAGCGCTGATGGACAGAGGGCGGCACGATCTGCGCGCACCGCTTGGTGGAGCGCGCGCGGGCGATCCACGCGCC-----CGTGGCACATAAAGCTCCCCGCA
#242_CDA1L2_1_Lp TCGAGCTAGCGCTGATGGACAGAGGGCGGCACGATCTGCGCGCACCGCTTGGTGGAGCGCGCGCGGGCGATCCACGCGCC-----CGTGGCACATAAAGCTCCCCGCA
Lj#1_CDA1L2_1_Lj TCGAGCTAGCGCTGATGGACAGAGGGCGGCACGATCTGCGCGCACCGCTTGGTGGAGCGCGCGCGGGCGATCCACGCGCC-----CGTGGCACATAAAGCTCCCCGCA
Lj#1_CDA1L2_2_Lj TCGAGCTAGCGCTGATGGACAGAGGGCGGCACGATCTGCGCGCACCGCTTGGTGGAGCGCGCGCGGGCGATCCACGCGCC-----CGTGGCACATAAAGCTCCCCGCA
#144_CDA1L2_1_Pm TCCACGCGCCCTCCGACGAGATCTGGACGAGGCGA-GAGAGCAAGGAGCGTCTGCGCGGCTGCGCGCGCGTGGTGCAGCGTGGTGGCCAGCTGCGCGCGCGTGGTGCAGCGTGGTGGCCAGCTGCGCGCGG
Pm#1_CDA1L2_1_Pm TCCACGCGCCCTCCGACGAGATCTGGACGAGGCGA-GAGAGCAAGGAGCGTCTGCGCGGCTGCGCGCGCGTGGTGCAGCGTGGTGGCCAGCTGCGCGCGG

#173_CDA1L2_1_Lp TCAATTACCCGTGA
#175_CDA1L2_1_Lp TCAATTACCCGTGA
#236_CDA1L2_1_Lp TCAATTACCCGTGA
#175_CDA1L2_2_Lp TCAATTACACGTGAATGATGAGGCAATTTGGACAGAGTGGAGAAGTTCTCTGA
#242_CDA1L2_1_Lp TCAATTACCCGTGA
Lj#1_CDA1L2_1_Lj TCAATTACCCGTGAATGATGAGGCAATTTGGACAGAGTGGAGAAGTTCTCCGTACCTGGGCGAGTTGCTTATGACAGCTCTGGTTGGACGAGGAGTCTCAATCTGA
Lj#1_CDA1L2_2_Lj TCAATTACCCGTGAATGATGAGGCAATTTGGACAGAGTGGAGAAGTTCTCCGTACCTGGGCGAGTTGCTTATGACAGCTCTGGTTGGACGAGGAGTCTCAATCTGA
#144_CDA1L2_1_Pm TCAATTACCCGTGAATGATGAGGCAATTTGGACAGAGTGGAGAAGTTCTCCGTACTTGGGCGAGTTTCTACTGACAGCTCTGGTTGGCA
Pm#1_CDA1L2_1_Pm TCAATTACCCGTGAATGATGAGGCAATTTGGACAGAGTGGAGAAGTTCTCCGTACTTGGGCGAGTTTCTACTGACAGCTCTGGTTGGCA

```

C

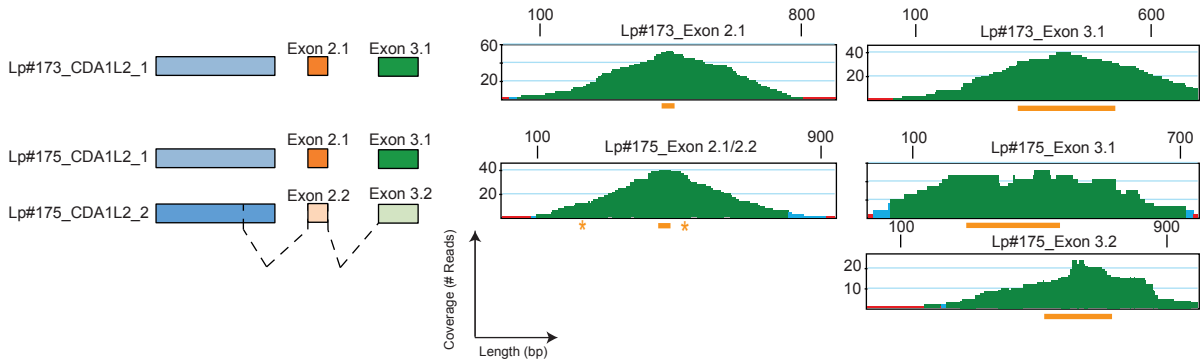
```

LpCDA1L2_1_Unspliced MSIFLYKPLPLNTFLPEFDNLEKAYGRNRCYICFKLPIHAVGATGTTGTTGSELNLYATNKWEVDGIPRESPEKRGMAHEPLLEDNTHRVHREHGGSFICIEWFTSWS
LpCDA1L2_1_Spliced -----FEFDNLEKAYGRNRCYICFKLPIHAVGATGTTGTTGSELNLYATNKWEVDGIPRESPEKRGMAHEPLLEDNTHRVHREHGGSFICIEWFTSWS
LpCDA1L2_2_Unspliced MSVFLHKKLPLNTFLPEFDNLEKAYGRNRCYICFKLPIYAVGATGTTGTTGSELNLYATNKWEVDGIPRESPEKRGMAHEPLLEDNTHRVHREHGGSFICIEWFTSWS
LpCDA1L2_2_Spliced MSVFLHKKLPLNTFLPEFDNLEKAYGRNRCYICFKLPIYAVGATGTTGTTGSELNLYATNKWEVDGIPRESPEKRGMAHEPLLEDNTHRVHREHGGSFICIEWFTSWS
LpCDA1L2_1_Unspliced PCDRC SGLLHSLWRDVGGRHRLRWFPSRIYRNVGAVRAGLRHLRRAVQLGVMDRRRHDCAHALVDAARGDPTP----PWHINVPRIQRAFEIMDEKVR-----E
LpCDA1L2_1_Spliced PCDRC SGLLHSLWRDVGGRHRLRWFPSRIYRNVGAVRAGLRHLRRAVQLGVMDRRRHDCAHALVDAARGDPTP----PWHINVPRIQRAFEIMDEKDDNGSDNSD
LpCDA1L2_2_Unspliced PCHRCSGLLHSLWRDVGGRHRLRWFPSRIYVDGAVRAGLRHLRRAVQLGVMDRLRDYCAHALVDTAQDPTPLWLVPWHMNVPRVQRAFEIMDEKVR-----E
LpCDA1L2_2_Spliced PCHRCSGLLHSLWRDVGGRHRLRWFPSRIYVDGAVRAGLRHLRRAVQLGVMDRLRDYCAHALVDTAQDPTPLWLVPWHMNVPRVQRAFEIMDEKDDNGSDNSD
LpCDA1L2_1_Unspliced <> Exon 2
PG-----NCWNCWRPFWL-----A-----DVPNSALNYP
LpCDA1L2_1_Spliced PGLSEISASGGHSHWDDDLHLPLEDLTVVVECTPSKQCPPEATAAPTLPRKRQEDPVDALTKARRALE
LpCDA1L2_2_Unspliced PG-----NCWNCWRPFWL-----A-----DVPNSALNYP
LpCDA1L2_2_Spliced PGLSEISASGGHSHWDDDLHLPLEDLTVVVECTPSKQCPPEATAAPTLPRKRQEDPVDALTKARRALE

```

**Fig. S5.** Alternative splicing of *CDAIL2* genes. (A) Read coverage plots from whole genome sequences (WGS) of several individuals, indicated above the plots. Green color indicates a coverage by >5 reads; blue, 2-5 reads, and red, single read. Orange bars correspond to region of contig containing open reading frame (ORF) of the gene. (B) Nucleotide alignment of ORFs of *CDAIL2* genes derived from whole genome sequences of *Lampetra planeri* (Lp#173, Lp#175, Lp#236, Lp#242), *Lethenteron japonicum* (Lj#1), and *Petromyzon marinus* (Pm#144 and Pm#1). *CDAIL2* genes in *Petromyzon marinus* are predicted to be pseudogenized based on numerous single nucleotide insertions (green asterisk) or deletions (red asterisk), and a large-scale deletion in the middle of the predicted ORF (blue asterisks). (C) Amino acid alignments of spliced and unspliced *CDAIL2* gene products from *Lampetra planeri*. Conserved APOBEC catalytic HxE and PCxxC motifs are highlighted in red. Exon 2 sequences are highlighted in orange, exon 3 sequences in green. Unspliced sequences are derived from whole genome shotgun sequences of lamprey Lp#175, the spliced sequences from RT-PCR or *de novo* transcriptome assembly.

A



B

```

Lp173_CDA1L2_Exon2.1 -----TGGTTTGGTTTGGAGGACTCACAGATCCTTGAAAGATCTGGAACAGTACATCAGTGAACCTCCAGCAGGCGAGGCTATGTTGGCTGGGTCATGTAGCCTGCAT
Lp175_CDA1L2_Exon2.1 ACCACACTGGTTTGGTTGAGGAGCTCACAGATCCTTGAAAGATCTGGAACAGTACATCAGTGAACCTCCAGCAGGCGAGGCTATGTTGGCTGGGTCATGTAGCCTGCAT
Lp175_CDA1L2_Exon2.2 ACCACACTGGTTTGGTTGAGGAGCTCACAGATCCTTGAAAGATCTGGAACAGTACATCAGTGAACCTCCAGCAGGCGAGGCTATGTTGGCTGGGTCATGTAGCCTGCAT

Lp173_CDA1L2_Exon2.1 GCCCAGCCACCGCATGCCATGCAGCTTCTATTTGGCTGGATAAAGGAGGGCTAAACAGGCACGGCATGGGCTTCAGAAGAGAATGGGCTGATGGTATAGGAGAATTT
Lp175_CDA1L2_Exon2.1 GCCCAGCCACCGCATGCCATGCAGCTTCTATTTGGCTGGATAAAGGAGGGCTAAACAGGCACGGCATGGGCTTCAGAAGAGAATGGGCTGATGGTATAGGAGAATTT
Lp175_CDA1L2_Exon2.2 GCCCAGCCACCGCATGCCATGCAGCTTCTATTTGGCTGGATAAAGGAGGGCTAAACAGGCACGGCATGGGCTTCAGAAGAGAATGGGCTGATGGTATAGGAGAATTT

Lp173_CDA1L2_Exon2.1 GGGAGTGTGTGGACTTGCGAAGGATAAAGTATAGCGGTGTCAATATTGGACTCTGTGAAGGAGGCTCGTGAAGAGCGCTACCAGGCAGTTGGAGGCAACCCGCCCTCCAT
Lp175_CDA1L2_Exon2.1 GGGAGTGTGTGGACTTGCGAAGGATAAAGTATAGCGGTGTCAATATTGGACTCTGTGAAGGAGGCTCGTGAAGAGCGCTACCAGGCAGTTGGAGGCAACCCGCCCTCCAT
Lp175_CDA1L2_Exon2.2 GGGAGTGTGTGGACTTGCGAAGGATAAAGTATAGCGGTGTCAATATTGGACTCTGTGAAGGAGGCTCGTGAAGAGCGCTACCAGGCAGTTGGAGGCAACCCGCCCTCCAT

Lp173_CDA1L2_Exon2.1 CACCAGCGGACCATCGCCGCTGCATGCGCAGTGTCTCCCCGTGTTTCACTGCAATGCCACTACTGCCAAGTACATTGTCAACAAAATCACATCTCTTTTGTCTTTGCA
Lp175_CDA1L2_Exon2.1 CACCAGCGGACCATCGCCGCTGCATGCGCAGTGTCTCCCCGTGTTTCACTGCAATGCCACTACTGCCAAGTACATTGTCAACAAAATCACATCTCTTTTGTCTTTGCA
Lp175_CDA1L2_Exon2.2 CACCAGCGGACCATCGCCGCTGCATGCGCAGTGTCTCCCCGTGTTTCACTGCAATGCCACTACTGCCAAGTACATTGTCAACAAAATCACATCTCTTTTGTCTTTGCA

Lp173_CDA1L2_Exon2.1 GGACGACAACGGGAGTGATAACTCAGCTAAATGGATGGCAATGTATACTTTAAATATGCAGTATTTACATGAGAGGGATAATCCAGTCAATGGAATGTGAACCCAGTA
Lp175_CDA1L2_Exon2.1 GGACGACAACGGGAGTGATAACTCAGCTAAATGGATGGCAATGTATACTTTAAATATGCAGTATTTACATGAGAGGGATAATCCAGTCAATGGAATGTGAACCCAGTA
Lp175_CDA1L2_Exon2.2 GGACGACAACGGGAGTGATAACTCAGCTAAATGGATGGCAATGTATACTTTAAATATGCAGTATTTACATGAGAGGGATAATCCAGTCAATGGAATGTGAACCCAGTA

Lp173_CDA1L2_Exon2.1 AACATGCATTTGTCGGGCAAAATTTGGTCTGTGTTGGAATGTTTGCACCAATGTTGGTTCATATTTGGTCCAAATTTCCCTGTTGACTAGGAGAGTCTTTAGTTGGCT
Lp175_CDA1L2_Exon2.1 AACATGCATTTGTCGGGCAAAATTTGGTCTGTGTTGGAATGTTTGCACCAATGTTGGTTCATATTTGGTCCAAATTTCCCTGTTGACTAGGAGAGTCTTTAGTTGGCT
Lp175_CDA1L2_Exon2.2 AACATGCATTTGTCGGGCAAAATTTGGTCTGTGTTGGAATGTTTGCACCAATGTTGGTTCATATTTGGTCCAAATTTCCCTGTTGACTAGGAGAGTCTTTAGTTGGCT

Lp173_CDA1L2_Exon2.1 TGGAGCAACATATTTCAAAGGAAAAATTTGAATAGTAAAAATGCCCTTTATTTACCAAAAATATATCTCACTGATCTCAAGTACTCTTTTTCGGGAGGCGCGTCATTTG
Lp175_CDA1L2_Exon2.1 TGGAGCAACATATTTCAAAGGAAAAATTTGAATAGTAAAAATGCCCTTTATTTACCAAAAATATATCTCACTGATCTCAAGTACTCTTTTTCGGGAGGCGCGTCATTTG
Lp175_CDA1L2_Exon2.2 TGGAGCAACATATTTCAAAGGAAAAATTTGAATAGTAAAAATGCCCTTTATTTACCAAAAATATATCTCACTGATCTCAAGTACTCTTTTTCGGGAGGCGCGTCATTTG

Lp173_CDA1L2_Exon2.1 GGACGTTTGACCACTCCCATGAATGAGTGTAAATGATGCCACTGTCCGCTAGCCACTCCGCCATGAGCGGAGAAGAAGTGCAGTACAGTATGTTATCGTTCCACTTT
Lp175_CDA1L2_Exon2.1 GGACGTTTGACCACTCCCATGAATGAGTGTAAATGATGCCACTGTCCGCTAGCCACTCCGCCATGAGCGGAGAAGAAGTGCAGTACAGTATGTTATCGTTCCACTTT
Lp175_CDA1L2_Exon2.2 GGACGTTTGACCACTCCCATGAATGAGTGTAAATGATGCCACTGTCCGCTAGCCACTCCGCCATGAGCGGAGAAGAAGTGCAGTACAGTATGTTATCGTTCCACTTT

Lp173_CDA1L2_Exon2.1 TCGAAGAACGAGTA
Lp175_CDA1L2_Exon2.1 TCGAAGAACGAGTACATGTGGTTTCACATCAATCCCAACAGACGGCCACTGTTATCATTTGTGTAA
Lp175_CDA1L2_Exon2.2 TCGAAGAACGAGTACATGTGGTTTCACATCAATCCCAACAGACGGCCACTGTTATCATTTGTGTAA

```

C

```

Lp173_CDA1L2_Exon3.1 -----
Lp175_CDA1L2_Exon3.1 -----
Lp175_CDA1L2_Exon3.2 AGAAGACCCCGTGGATGCTTGACGCCAAGAGGCCCCCTTTTAAATTTTGTAGCTCTTCTATGGTACTTGTATCAGGTGCGGCACCGTACATTTTGTAAAGTATGTAGCGCATCG

Lp173_CDA1L2_Exon3.1 -----GTAACCC-CACACCACCCGACACAGCCAA
Lp175_CDA1L2_Exon3.1 -----
Lp175_CDA1L2_Exon3.2 TTTTAAATTTGTGTATGTAAGAAAGAAACCAAGTACAATGTTTCAATATCMAGTTTTTGGGTTAGATCGCGTAAAGTATAAATGTTGCGTGTAAACCCACACACCCGACACAGCCAA

Lp173_CDA1L2_Exon3.1 TGAGTAAAGGATTTGTCGTGTAACAAAACGCCAATGAGTATTTGAAATGTAATGTTGGATTTACTCTCCAGCACCCGATGGGCTGAAATAGTAACTAATAGCGTTTTGTTTAC
Lp175_CDA1L2_Exon3.1 TGAGTAAAGGATTTGTCGTGTAACAAAACGCCAATGAGTATTTGAAATGTAATGTTGGATTTACTCTCCAGCACCCGATGGGCTGAAATAGTAACTAATAGCGTTTTGTTTAC
Lp175_CDA1L2_Exon3.2 TGAGTAAAGGATTTGTCGTGTAACAAAACGCCAATGAGTATTTGAAATGTAATGTTGGATTTACTCTCCAGCACCCGATGGGCTGAAATAGTAACTAATAGCGTTTTGTTTAC

Lp173_CDA1L2_Exon3.1 GTGACGAATCCCTTACTAATTTGGCTGTGTGGGTTGTGGAGGTTACGACCAACATTTACTTACCGGATCTAACCCAAAATCTTTATCTGTAATAATTTTGTGCTTAAGCCTA
Lp175_CDA1L2_Exon3.1 GTGACGAATCCCTTACTAATTTGGCTGTGTGGGTTGTGGAGGTTACGACCAACATTTACTTACCGGATCTAACCCAAAATCTTTATCTGTAATAATTTTGTGCTTAAGCCTA
Lp175_CDA1L2_Exon3.2 GTGACGAATCCCTTACTAATTTGGCTGTGTGGGTTGTGGAGGTTACGACCAACATTTACTTACCGGATCTAACCCAAAATCTTTATCTGTAATAATTTTGTGCTTAAGCCTA

Lp173_CDA1L2_Exon3.1 TGCCTTAACTGATGATCGAACGCGGTTTCCGCTCCCCCGAGCCCGGGGAAGCTGTCCGAGTCCCGCAGCGGGGGCACGAGAGTTGGCATGACGATGATTTGCACCTGCCGCTTGAT
Lp175_CDA1L2_Exon3.1 TGCCTTAACTGATGATCGAACGCGGTTTCCGCTCCCCCGAGCCCGGGGAAGCTGTCCGAGTCCCGCAGCGGGGGCACGAGAGTTGGCATGACGATGATTTGCACCTGCCGCTTGAT
Lp175_CDA1L2_Exon3.2 TGCCTTAACTGATGATCGAACGCGGTTTCCGCTCCCCCGAGCCCGGGGAAGCTGTCCGAGTCCCGCAGCGGGGGCACGAGAGTTGGCATGACGATGATTTGCACCTGCCGCTTGAT

Lp173_CDA1L2_Exon3.1 GATCTGACGGTGGTGGAGTGCACGCCAGCAAGCAAGGACCTCCCGAAGCCAGCCGCCCCACGCTGCCAGGAAACGACAACAAGAAGACCCCGTGGATGCTTGAAGCCCAAG
Lp175_CDA1L2_Exon3.1 GATCTGACGGTGGTGGAGTGCACGCCAGCAAGCAAGGACCTCCCGAAGCCAGCCGCCCCACGCTGCCAGGAAACGACAACAAGAAGACCCCGTGGATGCTTGAAGCCCAAG
Lp175_CDA1L2_Exon3.2 GATCTGACGGTGGTGGAGTGCACGCCAGCAAGCAAGGACCTCCCGAAGCCAGCCGCCCCACGCTGCCAGGAAACGACAACAAGAAGACCCCGTGGATGCTTGAAGCCCAAG

Lp173_CDA1L2_Exon3.1 AGGGCCCCCTTTTAAATTTTGTACTCTTTCATGCGTACTTGTATCAGTGTGGGACCCGCTACATTTTGTAAACGATATGATCGGCATCGTTTTTAAATGTTGTGTATGTAAGAAAGAACCA
Lp175_CDA1L2_Exon3.1 AGGGCCCCCTTTTAAATTTTGTACTCTTTCATGCGTACTTGTATCAGTGTGGGACCCGCTACATTTTGTAAACGATATGATCGGCATCGTTTTTAAATGTTGTGTATGTAAGAAAGAACCA
Lp175_CDA1L2_Exon3.2 AGGGCCCCCTTTTAAATTTTGTACTCTTTCATGCGTACTTGTATCAGTGTGGGACCCGCTACATTTTGTAAACGATATGATCGGCATCGTTTTTAAATGTTGTGTATGTAAGAAAGAACCA

Lp173_CDA1L2_Exon3.1 AGTACAATTTGTCATAATCAAGGTTTTTGGATCAGATCGCGTTTATTTATAAATGTTGTCGTAACCCCTCCACACCCGACAGCCGCAACCGTAAATGTTGTACAGTTTTTCAAGGACA
Lp175_CDA1L2_Exon3.1 AGTACAATTTGTCATAATCAAGGTTTTTGGATCAGATCGCGTTTATTTATAAATGTTGTCGTAACCCCTCCACACCCGACAGCCGCAACCGTAAATGTTGTACAGTTTTTCAAGGACA
Lp175_CDA1L2_Exon3.2 AGTACAATTTGTCATAATCAAGGTTTTTGGATCAGATCGCGTTTATTTATAAATGTTGTCGTAACCCCTCCACACCCGACAGCCGCAACCGTAAATGTTGTACAGTTTTTCAAGGACA

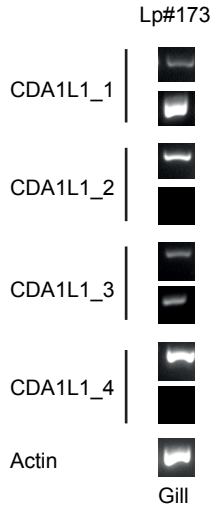
Lp173_CDA1L2_Exon3.1 AATTTATCACTTGTGTTAGCCACCCGGTGTGTTGAAGAGAGAATCAACAGCATTTT
Lp175_CDA1L2_Exon3.1 AATTTATCACTTGTGTTAGCCACCCGGTGTGTTGAAGAGAGAATCAACAGCATTTTAAATTTGAGTACAGCTTTCGCTGATAATTAACAATCAGCTTCATCGGGCAAAAATGGTGAAGC
Lp175_CDA1L2_Exon3.2 AATTTATCACTTGTGTTAGCCACCCGGTGTGTTGAAGAGAGAATCAACAGCATTTTAAATTTGAGTACAGCTTTCGCTGATAATTAACAATCAGCTTCATCGGGCAAAAATGGTGAAGC

Lp173_CDA1L2_Exon3.1 TCTCCTATATTGAATAATATTGGTCGCGAAAGCCCTACCGGTTAAAAAG
Lp175_CDA1L2_Exon3.1 TCTCCTATATTGAATAATATTGGTCGCGAAAGCCCTACCGGTTAAAAAG
Lp175_CDA1L2_Exon3.2 TCTCCTATATTGAATAATATTGGTCGCGAAAGCCCTACCGGTTAAAAAG

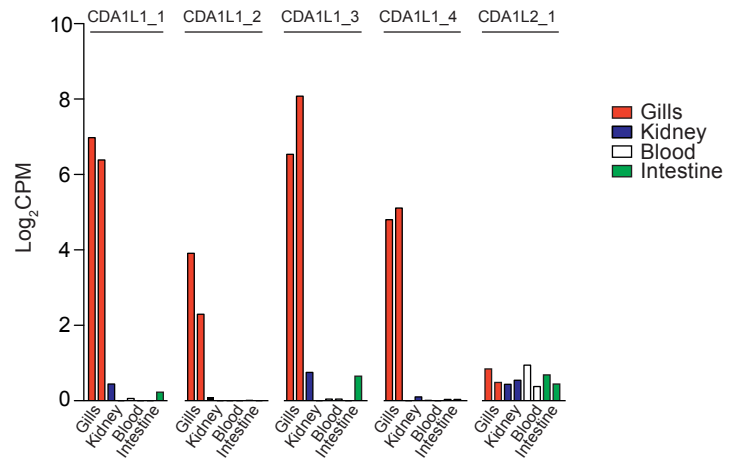
```

**Fig. S6.** Genomic organization of *CDAIL2* exons. (A) Schematic (left) and read coverage plots (right) for individuals containing either one (Lp#173) or two (Lp#175) *CDAIL2* genes. Green color indicates a coverage by >5 reads; blue, 2-5 reads, and red, single read. Orange bars correspond to region of contig containing open reading frame (ORF) of the gene. Due to the high degree of sequence similarity between exons 2.1 and 2.2 in Lp#175 whole genome sequences, only one single contig was assembled containing both exon variants. Orange asterisks indicate the position of SNPs contained within exon 2. (B, C) Nucleotide alignments of sequences containing exons 2 (B) and 3 (C) of *CDAIL2* genes. ORFs highlighted in red, SNPs in purple and splice signal site in orange.

A

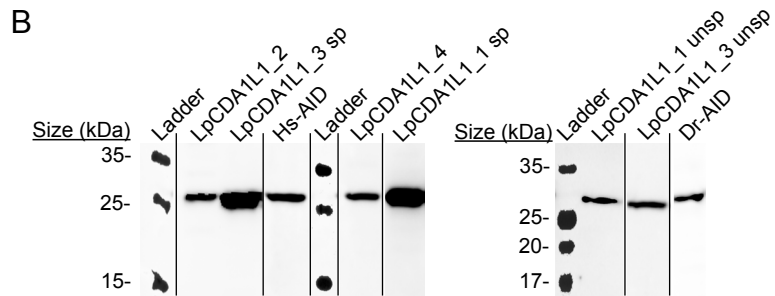
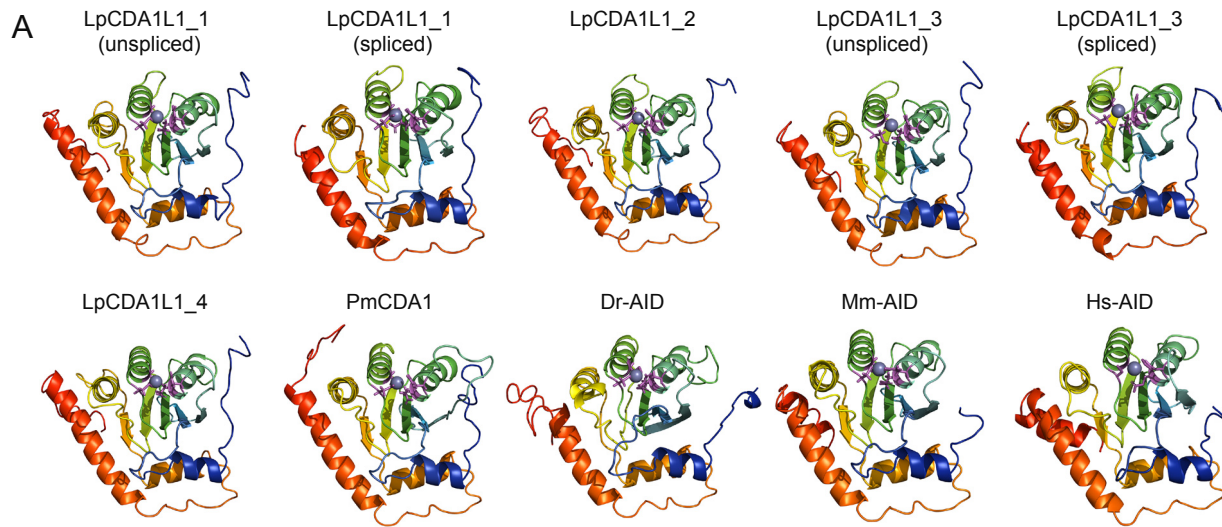


B

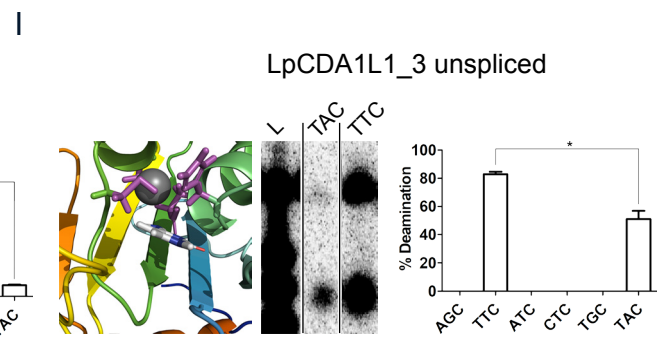
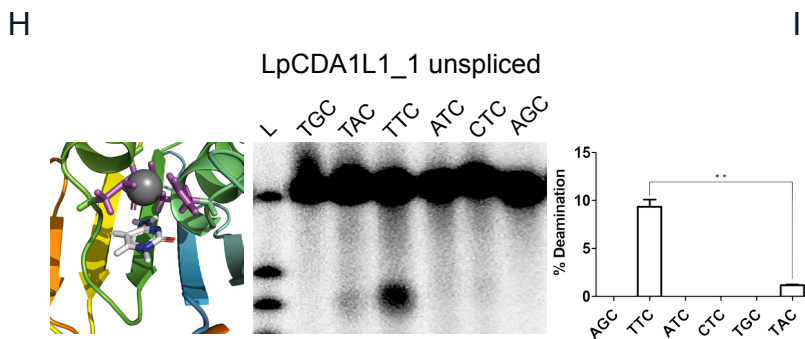
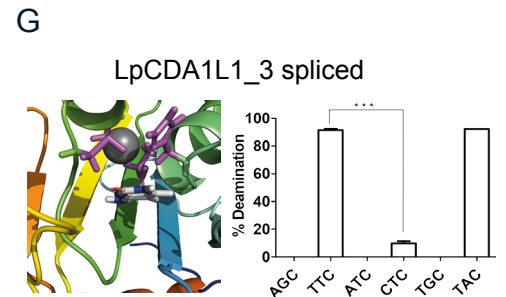
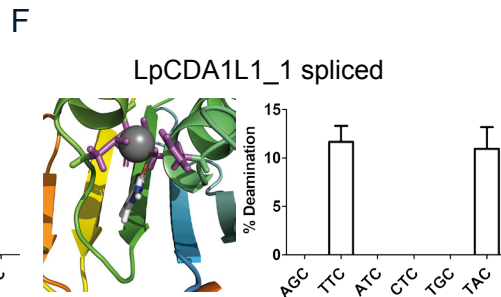
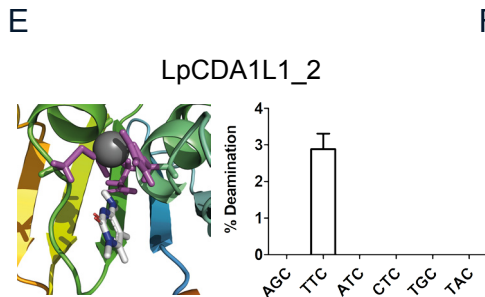
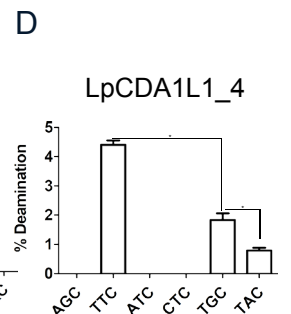
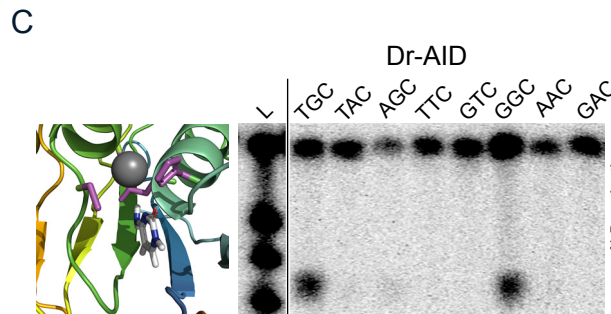
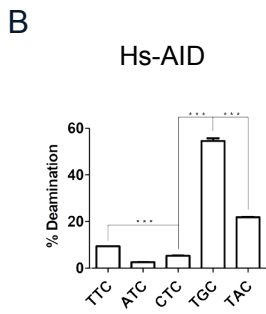
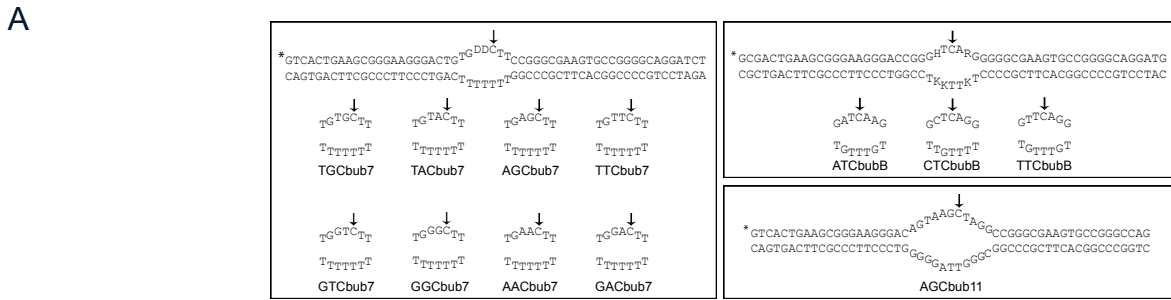


**Fig. S7.** Expression of *CDAI*-like genes. (A) RT-PCR analysis of *CDAILL1* gene expression using cDNA prepared from RNA of gills Lp#173, whose genome was confirmed to contain all 4 *CDAILL1* genes (c.f., Fig. S4). Gene-specific forward primers were used with universal reverse primers for unspliced (upper panels) or spliced (lower panels) gene products. Actin-specific primers were used as a control for cDNA integrity. (B) Expression of *CDAI*-like genes in immune organs of lamprey larvae. RNAseq-derived expression levels are expressed as  $\log_2$  of counts per million ( $\log_2$ CPM) in two *Lampetra planeri* individuals (Lp#131; Lp#132).





**Fig. S8.** Predicted structures of CDA1L1 proteins. (A) Representative ribbon models of LpCDA1L1 variants, PmCDA1 and DrAID, MmAID, and HsAID. All structures were predicted to have the same overall conserved tertiary structure apart from the C-terminal  $\alpha 7$  domain of HsAID, which was not predicted in the LpCDA1L1 variants or PmCDA1. N- to C-terminal progression is shown from blue to red, with the catalytic residues colored purple and the coordinated Zn denoted by a grey sphere. (B) Western blot analysis confirming the expression of CDA variants in 293T cells. SDS-PAGE electrophoresis on a 10% polyacrylamide gel was conducted on lysates of 293T cells transfected with expression constructs encoding each variant. Proteins were transferred onto nitrocellulose membranes and probed with polyclonal rabbit anti-V5 Tag antibody. Expression was confirmed by a specific band present at the expected size ~27kDa. sp= spliced, unsp= unspliced.



**Fig. S9.** Cytidine deamination activity of CDA1L1 variants. (A) Oligonucleotide substrates for cytidine deamination by alkaline cleavage experiments. The target (d)C is indicated by an arrow. The size of cleaved product at the target (d)C indicative of cytidine deamination activity is 28 nt. Left panel, DD(d)Cbub7 set, where D= A or T or G. All substrates in this set are identical except for the two nucleotides immediately upstream of the target (d)C in the bubble. Top right panel, bubB set. Substrates in this set have identical double-stranded arms, with different 7 nucleotide single-stranded bubbles. H= A or C or T, R= A or G, K= T or G. Bottom right panel, AGCbub11. This substrate has an AGC motif in the center of an 11-nucleotide bubble. (B) Substrate specificity profile of Hs-AID measured in the alkaline cleavage assay. (C) Catalytic pocket of Dr-AID docked with (d)C in a deamination-feasible configuration, showing the interactions between the Zn-coordinating triad and target (d)C; the coordinated Zn is depicted as a grey sphere, with the Zn-coordinating and catalytic glutamic acid residues colored purple (left), representative alkaline cleavage gel (middle), and substrate specificity profile graphs (right). (D) Substrate specificity profile for LpCDA1L1\_4. (E-G) Substrate specificity profiles for LpCDA1L1\_2, LpCDA1L1\_1, and LpCDA1L1\_3 proteins, the latter two in spliced configuration (c.f., Fig. 3); catalytic pockets of proteins docked with (d)C in a deamination-feasible configuration are shown in the left panels. (H,I). Catalytic pockets of unspliced LpCDA1L1\_1, and LpCDA1L1\_3 proteins docked with (d)C in a deamination-feasible configuration (left panels), representative alkaline cleavage gels (middle panels), and substrate specificity profile graphs (right panels), \*P ≤ 0.05, \*\*P < 0.01, \*\*\*P < 0.001. Error bars represent standard error of the mean. Some display items were combined from different parts of the same autoradiographic films; the splice sites are indicated by solid lines. Ladder size as indicated in Figure 3.

Table S1: Summary of lamprey usage and tissue sources used in this study

Animal ID	Species	Stage	River	gDNA				RNA	
				Genome Sequencing	Southern Blot	PFGE	PCR	Transcriptome	RT-PCR
Lp#8a	<i>Lampetra planeri</i>	Larva	La Sélune				WB		
Lp#8b	<i>Lampetra planeri</i>	Larva	La Sélune				WB		
#159	<i>Lampetra planeri</i>	Larva	Rhine			B			
#130	<i>Lampetra planeri</i>	Larva	Rhine						K
#131	<i>Lampetra planeri</i>	Larva	Rhine					K G B T	G
#132	<i>Lampetra planeri</i>	Larva	Rhine					K G B T	
#144	<i>Petromyzon marinus</i>	Larva	Rhine	B					
#173	<i>Lampetra planeri</i>	Larva	Rhine	K G B T					G
#175	<i>Lampetra planeri</i>	Larva	Rhine	K G B T					
#196	<i>Lampetra planeri</i>	Larva	Rhine		WB		WB		
#236	<i>Lampetra planeri</i>	Larva	La Sélune	WB	WB				
#242	<i>Lampetra planeri</i>	Larva	Rhine	WB					
Lf#29	<i>Lampetra fluviatilis</i>	Larva	La Sélune				F		
Lf#33	<i>Lampetra fluviatilis</i>	Larva	La Sélune				F		
Pm#1	<i>Petromyzon marinus</i>	Larva	Rhine	WB					
Lj#1	<i>Lethenteron japonicum</i>	Adult	Ishikari	Te					

B=Blood; F=Fin Clip G=Gills; K=Kidney; T=Typhlosole (Intestine); Te=Testes; WB=Whole Body

Table S2: Summary of *de novo* transcriptome assembly statistics

<i>L. planeri</i>	Tissue	Total # Read Pairs	Total # Reads after trimming		Total # of ORFs	# Unique ORFs
			R1	R2		
#131	Blood	186590513	183786365	155576954	675747	<b>153032</b>
	Gills	202230155	199814264	173904091		
	Kidney	209200560	206317660	180153510		
	Intestine	169020866	166889768	138656896		
	Total	767042094	756808057	648291451		
#132	Blood	173750048	171219322	148579163	645503	<b>152284</b>
	Gills	207350277	204725620	177506558		
	Kidney	199057420	196587887	169331118		
	Intestine	132513609	131013341	103726338		
	Total	712671354	703546170	599143177		

Table S3: Summary of CDA1L1 genotype of different lamprey species.

WGS or PCR	P	W	P	W	P	W	W	P	P	W	P	W	W
CDA1L1_1											?		
CDA1L1_2											?		
CDA1L1_3											?		
CDA1L1_4											?		
CDA1													
	#196	#173	#8a	#175	#8b	#236	#242	#29	#33	Lj#1	Lj#2	Pm#1	#144
	<i>Lp</i>							<i>Lf</i>		<i>Lj</i>		<i>Pm</i>	

*Lp* = *Lampetra planeri*; *Lf* = *Lampetra fluviatilis*; *Lj* = *Lethenteron japonicum*; *Pm* = *Petromyzon marinus*

## References

1. Smith JJ, et al. (2013) Sequencing of the sea lamprey (*Petromyzon marinus*) genome provides insights into vertebrate evolution. *Nat Genet* 45: 415-421.
2. Smith JJ, et al. (2018) The sea lamprey germline genome provides insights into programmed genome rearrangement and vertebrate evolution. *Nat Genet* 50: 270-277.
3. Rougemont Q, et al. (2015) Low reproductive isolation and highly variable levels of gene flow reveal limited progress towards speciation between European river and brook lampreys. *J Evol Bio.* 28: 2248-2263.