# Supplementary Materials

## MHC in newts illuminates the evolutionary dynamics of complex regions in giant genomes

W. Babik, K. Dudek, G. Palomar, M. Marszałek, G. Dubin, M. H. Yun, M. Migalska

## 7 Supplementary methods

8 **Samples and RNAs.** Adult *Pluerodeles waltl* from a laboratory colony in TU Dresden were mated,
9 eggs and larvae developed at 19 °C and 12:12 L:D light regime. The following stages were
10 sampled and preserved in RNAlater: young larvae (stages 38-44 Shi and Boucaut 1995), older
11 larvae (stages 45-47), metamorphosing larvae (stages 48-49) and metamorphs (stage 50). Because
12 *P. waltl* larvae are larger than those of other species, starting from stage 46, we sampled separately
13 also intestine (middle part) and spleen. For all other species we sampled, during the breeding
14 season, two presumably inseminated females from natural populations (Table S3). The females
15 were kept in the laboratory until they deposited eggs, were then euthanized with MS222 and
16 heart, intestine (middle part), liver, lungs, spleen and tail tissue were preserved in RNAlater and
17 separate libraries were prepared for each tissue. Eggs and larvae developed in the laboratory at
18 room temperature under the natural light regime. Larvae were fed ad libitum with *Artemia* nauplii
19 and then with live chironomid larvae. Larvae were euthanized in MS222 and preserved in
20 RNAlater; depending on their size, they were cut into two (anterior/posterior) or three
21 (head/middle/tail) parts, preserved and further processed separately (Table S4).

22 **RNA extraction and Illumina RNA-seq**. Total RNA was extracted from tissues preserved in
23 RNAlater using RNAzol® RT (Sigma) or, if only a very small amount of tissues was available,
24 RNeasy kits (Qiagen). RNA-seq libraries were constructed using Novogene Plant and Animal
25 Eukaryotic Strand Specific mRNA (WOBI) approach and sequenced (2 x 150 bp) on the Illumina
26 NovaSeq platform by Novogene.

27 **Full-length transcripts.** PacBio HiFi reads were clustered into transcripts using SMRT Tools v.
28 13.1: lima, isoseq refine and isoseq cluster. Only high-quality (hq) Iso-Seq transcripts generated
29 by this pipeline were used in subsequent analyses. A transcript is labelled hq if: 1) it is supported
30 by full-length non-chimeric (FLNC) reads with both primers and poly(A) tail detected, 2) the
31 consensus sequence has been polished using the underlying FLNC reads until the predicted per-
32 base accuracy $\geq$ 99% (QV $\geq$ 30), and 3) it passes internal software filters against artifacts such as
33 chimeras or concatemers.

34 **Identification of MHC transcripts.** First, sequences of the axolotl MHC-I (GenBank
35 AAC60108) and MHC-II (XP_069491359.1, AAG42326, XP_069491364.1, XP_069491365.1)
36 proteins were used as queries to identify MHC-I and MHC-II sequences in *P. waltl* reference
37 transcriptome (Matsunami et al. 2019) using TBLASTN v. 2.15 (Camacho et al. 2009). Both
38 axolotl and *P. waltl* MHC protein sequences were used to TBLASTN search the Iso-Seq
39 transcripts from all seven newt species at the E-value threshold of 1e-30. TransDecoder v 5.7.1

40    (Haas et al. 2013) was used to identify Open Reading Frames (ORF) in Iso-Seq transcripts (both

41    complete and partial ORFs were allowed, single best ORF per transcript). Because a large

42    number of transcripts were identified (ca. 700 in the case of MHC-I) and visual inspection of

43    their protein alignments indicated that sequences were often highly similar but differed by the

44    presence/absence of entire exons, to reduce redundancy we clustered the proteins within each

45    species based on sequence identity. Clustering was performed using Clusterize() from the R

46    package DECIPHER (Wright 2016) at the protein sequence divergence threshold 0.1, ignoring

47    regions containing gaps (penalizeGapLetterMatches = FALSE). The longest protein was selected

48    as a representative of each cluster. Cluster representatives of the minimum length of 240 amino

49    acids (aa) were aligned using DECIPHER::AlignSeqs(). The resulting alignments were visually

50    inspected and sequences lacking more than half of any of α/β 1-3 domains, or showing extreme

51    divergence in a part of the alignment and high similarity to at least one other sequence in the

52    alignment were removed. While such sequences may represent genuine isoforms, they may also

53    be artefacts of a limited-cycle PCR which is a part of Iso-Seq library preparation protocol.

54    Because some Iso-Seq clusters were represented by few HiFi reads and because some divergent

55    MHC lineages contained representatives of only some species, we suspected that our Iso-Seq

56    dataset may have missed some MHC sequences, especially those poorly expressed in the

57    intestine. Therefore we supplemented the Iso-Seq dataset with sequences of transcripts obtained

58    from de novo assemblies of short RNA-seq reads. Separate assemblies were obtained for each of

59    five individuals per species, including adults, metamorphs, and, for some species, also advanced

60    larvae (Table S2, S4), reads from all individual's RNA-seq libraries were pooled and assembled

61    with Trinity v. 2.15 (Grabherr et al. 2011) using --SS_lib_type RF parameter suitable for our

62    stranded RNA-seq libraries, with other settings set to default values. Trinity assemblies were

63    searched for MHC sequences using TBLASTN with the axolotl and *P. waltl* MHC-I and MHC-II

64    proteins as queries; sequences with hits of min. 90 aa alignment length and E-value < 1e-40 were

65    retained. TransDecoder was used to identify the single best ORF per Trinity transcript and the

66    resulting proteins were clustered together with representative sequences of Iso-Seq clusters.

67    Although the majority of Trinity proteins fell into previously identified clusters, some new

68    clusters which did not include Iso-Seq sequences and had only Trinity sequences of at least 240

69    aa appeared. Representative (longest) sequences from such clusters were visually examined for

70    the signatures of chimaerism, and those that passed this check were added to the dataset. It is

71    important to note here that Trinity sequences, resulting from de novo assembly of short reads,

72    may not represent true MHC sequences because of the difficulties of reconstructing full-length

73    transcripts of similar genes from multigene families from short reads. We decided to include such

74  sequences because they represent, albeit imperfectly, extra MHC diversity present in individuals

75  or tissues for which Iso-Seq was not available or lowly expressed.

76  **The genomic organisation of the MHC region.** For automatic annotation we used the

77  available gene predictions for *P. waltl* (Brown et al. 2025) and performed de novo prediction with

78  BRAKER3 (Gabriel et al. 2024) for the remaining three genomes. Because de novo gene

79  prediction is a complex and computationally intensive task, especially in large genomes, such as

80  those of newts, we limited our prediction to the MHC region(s) (both core and disparate). We

81  added other continuous genomic regions, totalling 950 Mb in *L. helveticus* and *L. vulgaris* and 700

82  Mb in *T. cristatus*, to reach a sufficient number of well-supported genes required to train

83  prediction algorithms in BRAKER3. The available version of *T. cristatus* genome assembly is

84  already softmasked for repeats, while in the case of *L. helveticus* and *L. vulgaris* assemblies, we

85  performed repeat identification and softmasking using RepeatModeler v. 2.0.6 (Flynn et al. 2020)

86  and RepeatMasker v. 4.1.7 (repeatmasker.org). Both RNA-seq data from target species and

87  vertebrate proteins from the OrthoDB database (Tegenfeldt et al. 2025) were used for gene

88  prediction. Predicted protein sequences were extracted from the genome assemblies using gffread

89  v 0.12.8. (Pertea and Pertea 2020) and annotated using eggNOG-mapper v2 (Cantalapiedra et al.

90  2021).

91  Because the initial mapping of cluster representatives to genomes indicated that automatic

92  predictions of MHC genes were in many cases not accurate, MHC-I and MHC-II genes were

93  annotated manually (also when present outside of the core MHC region). We used the following

94  sources of information: i) gene models from automatic prediction, ii) mapping of all RNA-seq

95  libraries from a target species to the genome with hisat2 v. 2.2.1 (Kim et al. 2019), iii) mapping of

96  all cluster representatives (obtained as described above) from the genus to the genome with

97  minimap2 v 2.28 (Li 2018), iv) mapping of all Iso-Seq transcripts from the species to the genome

98  with minimap2. Tracks representing all the data were visualised in IGV v. 2.18.4 (Robinson et al.

99  2011), coordinates of coding sequences were identified and exported as .bed files. The retrieved

100  sequences were checked for internal stop codons, and those that did not contain them were

101  considered as putative MHC genes, forming the dataset G.

102  To identify putative MHC pseudogenes or gene fragments, we extracted for each species exons >

103  100 bp from all annotated MHC genes, and used BLASTN to search the genome at the E-value

104  threshold of 1e-10 with scoring adjusted to divergent sequences (-reward 1 -penalty -1 -gapopen

105  3 -gapextend 2). Hits with at least 70% sequence identity covering at least 80% of the query

106  sequence located outside manually annotated MHC genes were considered to represent

107  pseudogenes or gene fragments. Complete annotations of the MHC genes and other genes in the

108  MHC region are provided as .gff files in Supplemental Data.

109  **MHC-I polymorphism.** We did not analyze polymorphism of MHC-I-like genes as the primers

110  used in previous studied did not amplify them. The alleles from Palomar et al. (2021) and

111  Gaczorek et al. (2023) were BLASTN-mapped to the respective genomes at the E-value

112  threshold of 1e-30. Because the majority of alleles produced multiple hits, we considered as valid

113  all hits within 0.9 of the bitscore of the best hit. Each hit was given the weight 1/n, where n was

114  the number of hits, allowing an allele matching multiple locations in the genome to be counted as

115  a fraction in each, avoiding overestimation of polymorphism for similar genes. The primers used

116  in Palomar et al. (2021) and Gaczorek et al. (2023) were developed from transcriptome data

117  before genome assemblies became available, and thus were not guaranteed to amplify all genes.

118  We therefore screened the primers against all gene sequences to identify potential mismatches in

119  binding sites. For *L. vulgaris* and *T. cristatus*, we designed additional primers (Table S9), combined

120  them with those previously published, and re-genotyped MHC-I in 24 individuals per species.

121  PCR, library preparation, sequencing, and genotyping followed Palomar et al. (2021) and

122  Gaczorek et al. (2023). In *L. vulgaris*, genotypes from the original and updated primer sets were

123  nearly identical (2/125 alleles, 1.5%, were new). In *T. cristatus*, 17/126 alleles (13.5%) were new,

124  yet estimates of polymorphism from the two primer sets were extremely highly correlated ($R^2 =$

125  99.8%; Fig. S9). For *P. waltl* we verified, using newly designed primers (Table S9), that *Plewa-06*

126  and *Plewa-08* genes that previously did not amplify are not polymorphic. We conclude thus that

127  our assessment of polymorphism of the MHC-I genes annotated in newt genomes is accurate.

128  **Phylogenies of MHC proteins.** The G+T dataset protein sequences, together with the axolotl

129  or *Andrias* as outgroups, were aligned with DECIPHER::AlignSeq(), separately for MHC-I,

130  MHC-IIA, MHC-IIB, MHC-DMA, MHC-DMB, as well as for MHC-I $\alpha 1$, $\alpha 2$, and $\alpha 3$ domains.

131  The amino-acid evolution models were identified by ModelTest-NG v. 0.1.7 (Darriba et al. 2020);

132  because for all alignments JTT + G4 or similar models were selected, we decided to use JTT +

133  G4 model for all phylogenetic analyses. Maximum likelihood (ML) phylogenies were

134  reconstructed in RAxML-NG 1.2.2 (Kozlov et al. 2019), and robustness of the obtained

135  topologies was tested with 100 bootstrap replicates.

136  **MHC-I and MHC-II expression through ontogeny and across tissues.** Expression was

137  estimated by mapping RNA-seq libraries to two types of references: (i) for the four species with

138  the available genome assemblies, RNA-seq reads were mapped to genomes with hisat2 and reads

139  mapped to the annotated MHC genes, pseudogenes and gene fragments were counted with

140     featureCounts v. 2.0.8 (Liao et al. 2014) – this analysis provided detailed expression information

141     for MHC sequences present in reference genomes, (ii) for all seven species, RNA-seq reads were

142     mapped to the cluster representatives using Bowtie2 v. 2.5.1 (Langmead and Salzberg 2012) and

143     the number of reads mapped to each reference was calculated with samtools idxstats – this

144     analysis provided expression estimates for all full length MHC coding sequences identified in our

145     samples. Expression was estimated separately for each library as Fragments Per Kilobase of

146     transcript length per Million reads mapped (FPKM); in the analysis of all seven species we

147     assumed mapping rate of 0.895 as this was the average for the four species with the available

148     genome assemblies.

149     **Gene tree – species tree reconciliation.** The analysis was performed using Notung v. 3.0_25

150     beta (Durand et al. 2006) on two datasets that used different gene trees: (i) G, and (ii) G+T.

151     Notung analyses used ML trees constructed using protein sequences with RAxML-NG (see

152     above). Notung bootstrap edge weight threshold parameter, which indicates which branches are

153     "weak", i.e., weakly supported, was set to 50, the cost of duplication and cost of loss were kept at

154     the default values of 1.5 and 1.0. After reconciliation tree was adjusted to further optimise the

155     number of duplications and losses by rearranging "weak" edges, that can be swapped in their

156     placement without strongly contradicting the data.

157     **Testing for gene conversion.** Gene conversion was tested using coding DNA sequences of the

158     G+T dataset after removing sequences assembled by Trinity from short RNA-seq reads, as some

159     of them could have been computationally generated chimaeras that may produce a false

160     recombination signal. The analyses were performed separately for MHC-I and MHC-I-like

161     genes, because high sequence divergence between them makes conversion unlikely and because

162     including numerous sequences would reduce the power of the test by inflating the number of

163     pairwise comparisons that would necessitate stringent correction for multiple testing. MHC-I-like

164     sequences from all species were analysed together due to their relatively small number. MHC-I

165     sequences from each species were analysed separately as we were mostly interested in species-

166     specific gene conversion events. The analyses were performed using geneconv v. 1.81a (Padidam

167     et al. 1999) and were run both allowing no mismatches within conversion tracts ($g = 0$) and using

168     the more permissive setting allowing for some mismatches ($g = 1$) and thus enabling detection of

169     older conversion events.

170     **Selective pressures acting on MHC-I and MHC-I-like genes.** Tests of selective pressures

171     were performed applying methods implemented in HyPhy v. 2.5.8 (Kosakovsky Pond et al.

172     2020), using protein-guided codon alignment of the G+T dataset obtained with DECIPHER::

173     AlignTranslation() and approximately ML phylogeny inferred with FastTree v. 2.1.11 (Price et al.

174 2010). RELAX (Wertheim et al. 2015) was used to test whether a change in the strength of

175 selection occurred (i) along the long branches separating MHC-I-like1 and MHC-I-like2 from

176 other sequences in MHC-I phylogeny, (ii) within MHC-I-like1 and MHC-I-like2 clusters. FUBAR

177 (Murrell et al. 2013) was used to identify codons under positive/purifying selection, separately for

178 MHC-I, MHC-I-like1 and MHC-I-like2. FUBAR analysis was limited to α1–α3 domains.

179 **Structural modelling.** Genomic sequences of MHC-I-like1 molecules were obtained for *L.*

180 *vulgaris* and *P. waltl*, while MHC-I-like2 sequences were identified in *L. helveticus*, *L. vulgaris*, and *T.*

181 *cristatus*. A putative *P. waltl* MHC-I-like2 locus was also identified on a different chromosome

182 than the core MHC region. This locus contained an ORF spanning 20 exons, several of which

183 appeared to encode domains atypical of MHC molecules. Given its likely chimeric and non-

184 functional nature, this gene was excluded from structural modelling. Signal peptides were

185 predicted using SignalP 6.0 (Teufel et al. 2021) and removed prior to structural modelling. MHC-

186 I-like2 molecules contained short N-terminal sequences resembling leader peptides; however, not

187 confidently identified by SignalP. Therefore, these regions were manually trimmed based on

188 sequence alignment with classical MHC-I molecules. Genomic sequences of all MHC-I-like

189 molecules and one putative classical molecule per species (*L. helveticus*: *Lihe-01*, *L. vulgaris*: *Livu-03*,

190 *T. cristatus*: *Trcr-01*, *P. waltl*: *Plwa-01*), together with corresponding β2m, were modelled with the

191 Alphafold3 Server (Abramson et al. 2024; default settings, access April 2025). Structures were

192 visualized with UCSF ChimeraX 1.9 (Meng et al. 2023). A search for structural homology was

193 performed with DALI server (Holm 2022).

## Supplementary References

195 Abramson J, Adler J, Dunger J, Evans R, Green T, Pritzel A, Ronneberger O, Willmore L,
196       Ballard AJ, Bambrick J. 2024. Accurate structure prediction of biomolecular interactions
197       with AlphaFold 3. *Nature* **630**: 493–500.

198 Almeida T, Ohta Y, Gaigher A, Muñoz-Mérida A, Neves F, Castro LFC, Machado AM, Esteves
199       PJ, Veríssimo A, Flajnik MF. 2021. A highly complex, MHC-linked, 350 million-year-old
200       shark nonclassical class I lineage. *J Immun* **207**: 824–836.

201 Brown T, Mishra K, Elewa A, Iarovenko S, Subramanian E, Araus AJ, Petzold A, Fromm B,
202       Friedländer MR, Rikk L. 2025. Chromosome-scale genome assembly reveals how repeat
203       elements shape non-coding RNA landscapes active during newt limb regeneration. *Cell*
204       *genomics* **5**. doi: 10.1016/j.xgen.2025.100761.

205 Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, Madden TL. 2009.
206       BLAST+: architecture and applications. *BMC Bioinformatics* **10**: 1–9.

207 Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J. 2021. eggNOG-
208     mapper v2: functional annotation, orthology assignments, and domain prediction at the
209     metagenomic scale. *Mol Biol Evol* **38**: 5825–5829.

210 Darriba D, Posada D, Kozlov AM, Stamatakis A, Morel B, Flouri T. 2020. ModelTest-NG: a new
211     and scalable tool for the selection of DNA and protein evolutionary models. *Mol Biol Evol*
212     **37**: 291–294.

213 Durand D, Halldórsson BV, Vernot B. 2006. A Hybrid Micro–Macroevolutionary Approach to
214     Gene Tree Reconstruction. *J Comput Biol* **13**: 320–335.

215 Flynn JM, Hubley R, Goubert C, Rosen J, Clark AG, Feschotte C, Smit AF. 2020.
216     RepeatModeler2 for automated genomic discovery of transposable element families. *Proc*
217     *Natl Acad Sci USA* **117**: 9451–9457.

218 Gabriel L, Br\uuna T, Hoff KJ, Ebel M, Lomsadze A, Borodovsky M, Stanke M. 2024.
219     BRAKER3: Fully automated genome annotation using RNA-seq and protein evidence
220     with GeneMark-ETP, AUGUSTUS, and TSEBRA. *Genome Research* **34**: 769–777.

221 Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
222     Raychowdhury R, Zeng Q, et al. 2011. Full-length transcriptome assembly from RNA-
223     Seq data without a reference genome. *Nat Biotech* **29**: 644–652.

224 Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, Couger MB, Eccles D,
225     Li BO, Lieber M. 2013. De novo transcript sequence reconstruction from RNA-seq using
226     the Trinity platform for reference generation and analysis. *Nature protocols* **8**: 1494–1512.

227 Holm L. 2022. Dali server: structural unification of protein families. *Nucl Acid Res* **50**: W210–
228     W215.

229 Kaufman J, Salomonsen J, Flajnik M. 1994. Evolutionary conservation of MHC class I and class
230     II molecules - different yet the same. *Sem Immun* **6**: 411–424.

231 Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. 2019. Graph-based genome alignment and
232     genotyping with HISAT2 and HISAT-genotype. *Nat Biotech* **37**: 907–915.

233 Kosakovsky Pond SL, Poon AF, Velazquez R, Weaver S, Hepler NL, Murrell B, Shank SD,
234     Magalis BR, Bouvier D, Nekrutenko A. 2020. HyPhy 2.5—a customizable platform for
235     evolutionary hypothesis testing using phylogenies. *Mol Biol Evol* **37**: 295–299.

236 Kozlov AM, Darriba D, Flouri T, Morel B, Stamatakis A. 2019. RAxML-NG: a fast, scalable and
237     user-friendly tool for maximum likelihood phylogenetic inference. *Bioinformatics* **35**: 4453–
238     4455.

239 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Meth* **9**: 357–359.

240 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**: 3094–3100.

241 Liao Y, Smyth GK, Shi W. 2014. featureCounts: an efficient general purpose program for
242     assigning sequence reads to genomic features. *Bioinformatics* **30**: 923–930.

243 Matsunami M, Suzuki M, Haramoto Y, Fukui A, Inoue T, Yamaguchi K, Uchiyama I, Mori K,
244     Tashiro K, Ito Y. 2019. A comprehensive reference transcriptome resource for the

245          Iberian ribbed newt *Pleurodeles waltl*, an emerging model for developmental and
246          regeneration biology. *DNA Research* **26**: 217–229.

247   Meng EC, Goddard TD, Pettersen EF, Couch GS, Pearson ZJ, Morris JH, Ferrin TE. 2023.
248          UCSF CHIMERAX : Tools for structure building and analysis. *Protein Sci* **32**: e4792.

249   Murrell B, Moola S, Mabona A, Weighill T, Sheward D, Kosakovsky Pond SL, Scheffler K. 2013.
250          FUBAR: a fast, unconstrained bayesian approximation for inferring selection. *Mol Biol*
251          *Evol* **30**: 1196–1205.

252   Padidam M, Sawyer S, Fauquet CM. 1999. Possible emergence of new geminiviruses by frequent
253          recombination. *Virology* **265**: 218–225.

254   Pertea G, Pertea M. 2020. GFF utilities: GffRead and GffCompare. *F1000Research* **9**: ISCB-
255          Comm.

256   Price MN, Dehal PS, Arkin AP. 2010. FastTree 2–approximately maximum-likelihood trees for
257          large alignments. *PloS one* **5**: e9490.

258   Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP.
259          2011. Integrative genomics viewer. *Nat Biotech* **29**: 24–26.

260   Sammut B, Du Pasquier L, Ducoroy P, Laurens V, Marcuz A, Tournefier A. 1999. Axolotl MHC
261          architecture and polymorphism. *Eur J Immunol* **29**: 2897–2907.

262   Shi DL, Boucaut JC. 1995. The chronological development of the urodele amphibian *Pleurodeles*
263          *waltl* (Michah). *Int J Dev Biol* **39**: 427–441.

264   Tegenfeldt F, Kuznetsov D, Manni M, Berkeley M, Zdobnov EM, Kriventseva EV. 2025.
265          OrthoDB and BUSCO update: annotation of orthologs with wider sampling of genomes.
266          *Nucl Acid Res* **53**: D516–D522.

267   Teufel F, Armenteros JJA, Johansen AR, Gíslason MH, Pihl SI, Tsirigos KD, Winther O, Brunak
268          S, Von Heijne G, Nielsen H. 2021. SignalP 6.0 achieves signal peptide prediction across
269          all types using protein language models. *BioRxiv* 2021–06.

270   Wertheim JO, Murrell B, Smith MD, Kosakovsky Pond SL, Scheffler K. 2015. RELAX: detecting
271          relaxed selection in a phylogenetic framework. *Mol Biol Evol* **32**: 820–832.

272   Wright ES. 2016. Using DECIPHER v2. 0 to analyze big biological sequence data in R. *R Journal*
273          **8**: 352–359.

274

275

## Supplementary Tables

Are in a separate Excel workbook.

**Fig. S1. Phylogeny of MHC-I and MHC-I-like sequences.** The RAxML-NG maximum likelihood tree constructed from protein sequences under the JTT+G4 amino-acid substitution model. The tree was rooted with *Andrias* sequence (GenBank AGY55973, not shown). Support values for clades with the minimum bootstrap support of 70% (100 bootstrap replicates) are shown. The tree contains sequences from the G+T dataset (as described in Materials and Methods) obtained from PacBio Iso-Seq, de novo assembly of RNA-seq and protein sequences predicted for manually annotated MHC genes in genome assemblies of *L. helveticus*, *L. vulgaris*, *P. waltl* and *T. cristatus*. For sequences predicted in genomes the number of protein-coding exons is shown. Circles at tippoints show maximum relative expression – the maximum fraction of total MHC-I expression (FPKM) attributable to a sequence across all RNA-seq libraries from postmetamorphic individuals within a species. Triangles show N conserved AA – the number of residues important for anchoring the termini of antigenic peptides containing amino acids that are conserved in classical MHC-I of most taxa. Sequences are color-coded according to the species and species phylogeny is in Fig. 1.

# DMA



# DMB



# MHC-IIA



# MHC-IIB

**Fig. S2. Phylogenies of MHC-II sequences.** Separate trees were constructed for α and β chains of DM and other class II proteins. The RAxML-NG maximum likelihood trees were constructed from protein sequences under the JTT+G4 amino-acid substitution model. The trees were rooted with *Andrias* or axolotl sequences (not shown). Support values for clades with the minimum bootstrap support of 70% (100 bootstrap replicates) are shown. The tree contains sequences from the G+T dataset (as described in Materials and Methods) obtained from PacBio Iso-Seq, de novo assembly of RNA-seq and protein sequences predicted for manually annotated MHC genes in genome assemblies of *L. helveticus*, *L. vulgaris*, *P. waltl* and *T. cristatus*. Sequences are color-coded according to the species and species phylogeny is in Fig. 1.

α1

α2

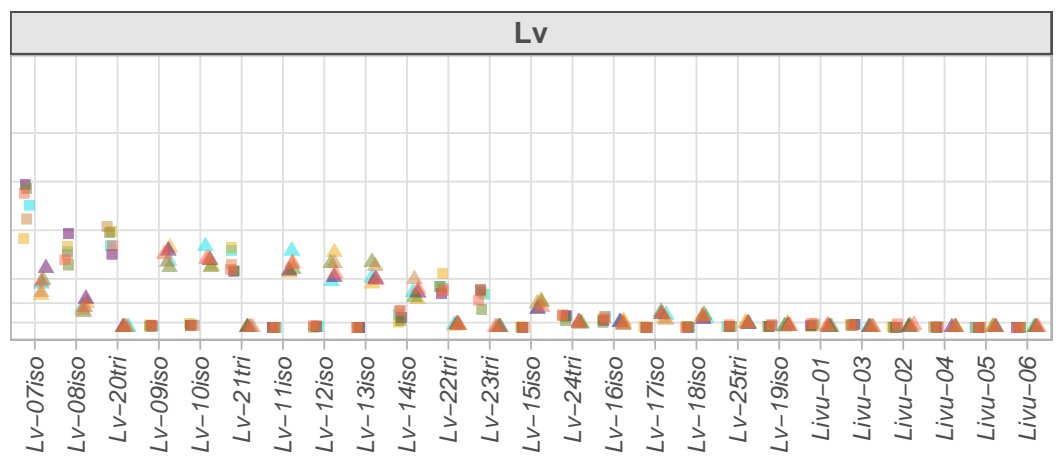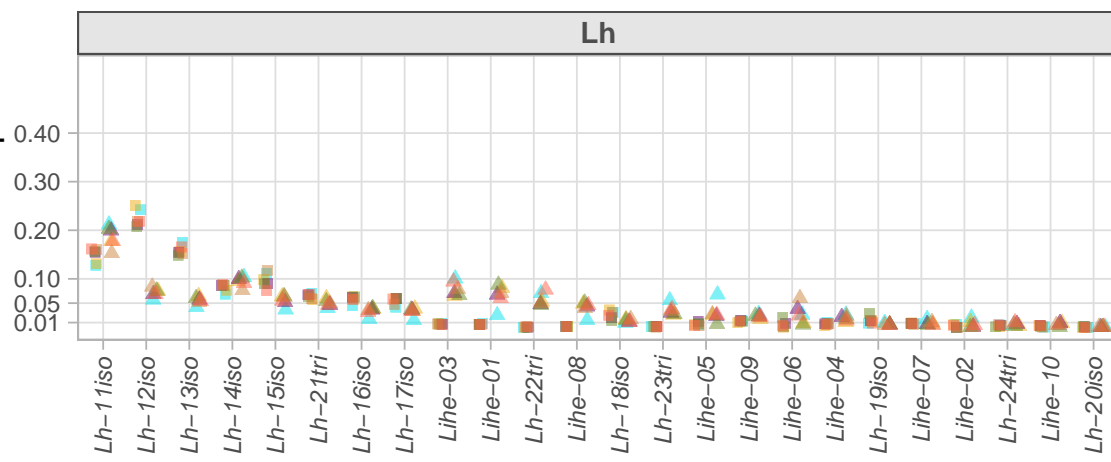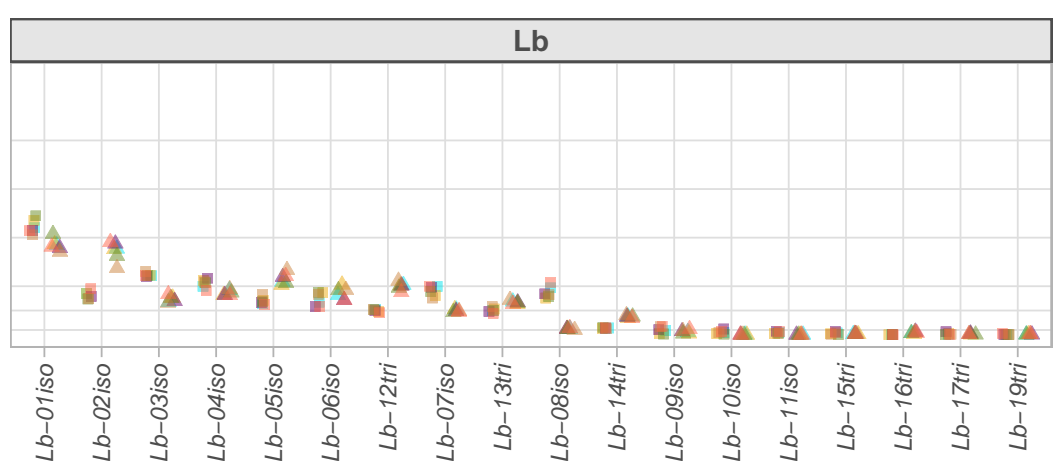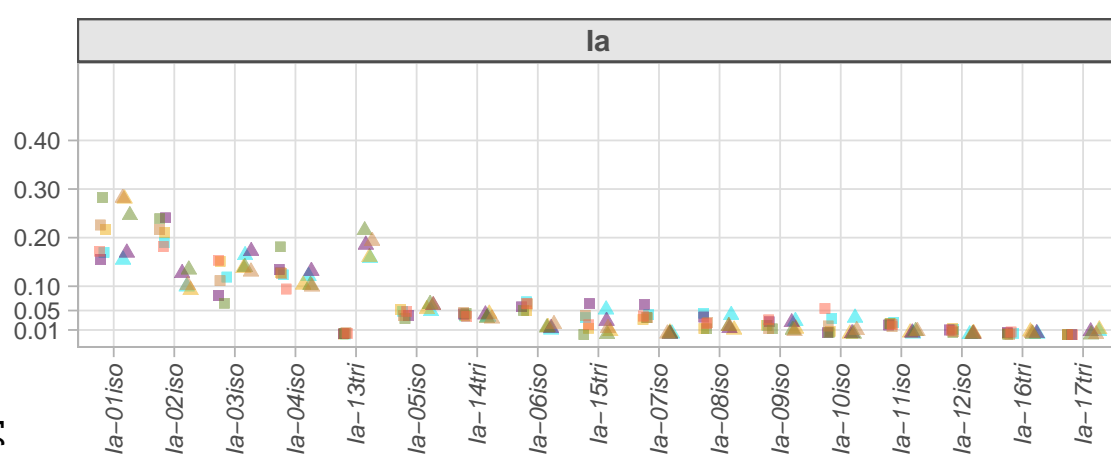**Fig. S3-S5. Phylogenies of MHC-I and MHC-I-like domains α1, α2, and α3.** The RAxML-NG maximum likelihood tree constructed from protein sequences under the JTT+G4 amino-acid substitution model. The tree was rooted with *Andrias* sequence (GenBank AGY55973, not shown). Support values for clades with the minimum bootstrap support of 70% (100 bootstrap replicates) are shown. The tree contains sequences from the G+T dataset (as described in Materials and Methods) obtained from PacBio Iso-Seq, de novo assembly of RNA-seq and protein sequences predicted for manually annotated MHC genes in genome assemblies of *L. helveticus*, *L. vulgaris*, *P. waltl* and *T. cristatus*. For sequences predicted in genomes the number of protein-coding exons is shown. Circles at tippoints show maximum relative expression – the maximum fraction of total MHC-I expression (FPKM) attributable to a sequence across all RNA-seq libraries from postmetamorphic individuals within a species. Triangles show N conserved AA – the number of residues important for anchoring the termini of antigenic peptides containing amino acids that are conserved in classical MHC-I of most taxa. Sequences are color-coded according to the species and species phylogeny is in Fig. 1. Note, that in Fig. S5 (α3) tree, a single *T. marmoratus* MHC-I-like2 sequence did not cluster together with the remaining MHC-I-like2 sequences, so the labelled MHC-I-like2 clade does not include this sequence.
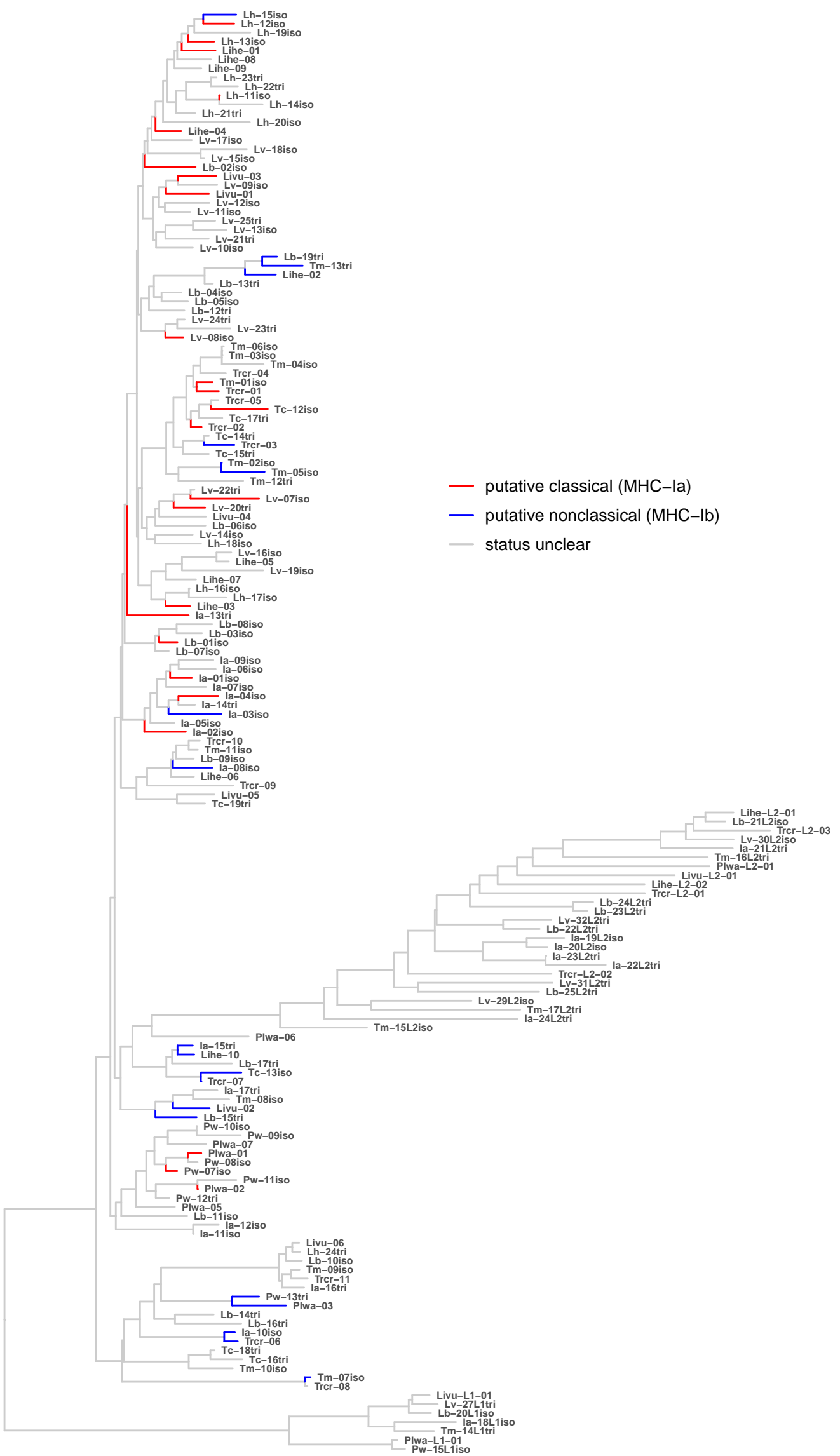
**Fig. S6. Expression of the G+T dataset MHC-I sequences through ontogeny.** Relative expression, i.e., FPKM for a particular cluster representative divided by the total MHC-I FPKM in a library. Boxplots show medians, interquartile and total ranges. Grey stripes indicate putative nonclassical (MHC-Ib) sequences identified as described in the text. Note, that for genes *Livu-02* and *Trcr-07* the pattern of higher relative expression in larval stages, visible when mapping RNA-seq data to genomes and used to classify them as MHC-Ib, is not clearly visible when mapping was performed to the broader set of G+T reference sequences. Species abbreviations: Ia – *I. alpestris*, Lb – *L. boscai*, Lh - *L. helveticus*, Lv - *L. vulgaris*, Pw – *P. waltl*, Tc – *T. cristatus*, Tm - *T. marmoratus*.

**Fig. S7. Expression of MHC-I sequences in adult tissues.** Relative expression, i.e., FPKM for a particular cluster representative divided by the total MHC-I FPKM in a library. RNA-seq libraries from different tissues are indicated with different colors and individuals with different shapes. Species abbreviations: Ia – *I. alpestris*, Lb – *L. boscai*, Lh - *L. helveticus*, Lv - *L. vulgaris*, Pw – *P. waltl*, Tc – *T. cristatus*, Tm - *T. marmoratus*. Sequence names correspond to those in Fig. 4, S1-S5. Note, that the cluster represented by the gene *Plwa-03* included also another gene annotated in the genome – *Plwa-04*.

putative classical (MHC−Ia)

putative nonclassical (MHC−Ib)

status unclear

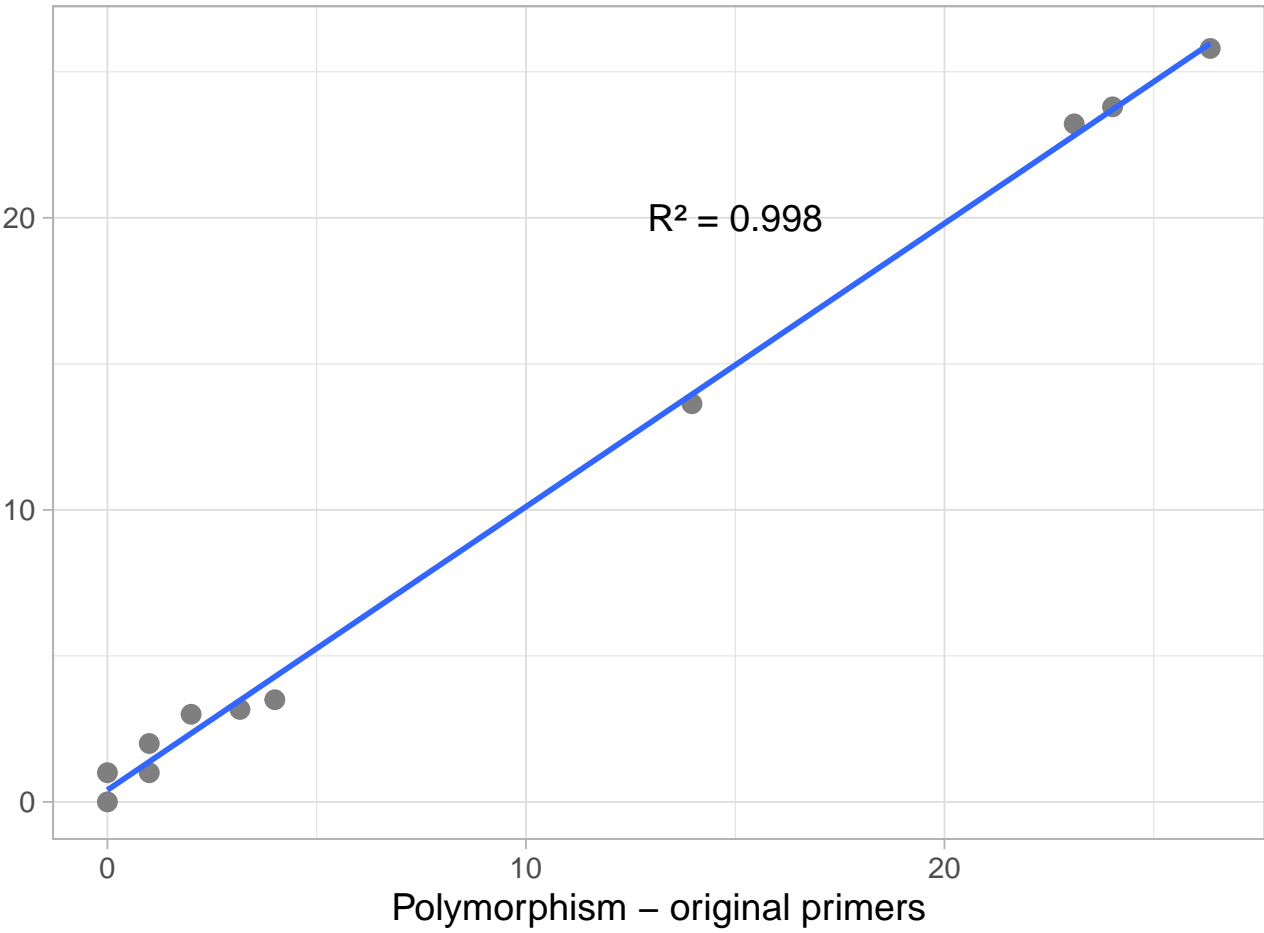**Fig. S8. Position of putative classical and nonclassical MHC-I sequences on phylogeny from Fig. S1.** Putative classical (MHC-Ia) and nonclassical (MHC-Ib) sequences were identified as described in the text and their branches are marked red and blue, respectively. The remaining sequences which status is unclear, as well as MHC-I-like sequences, are marked in gray.
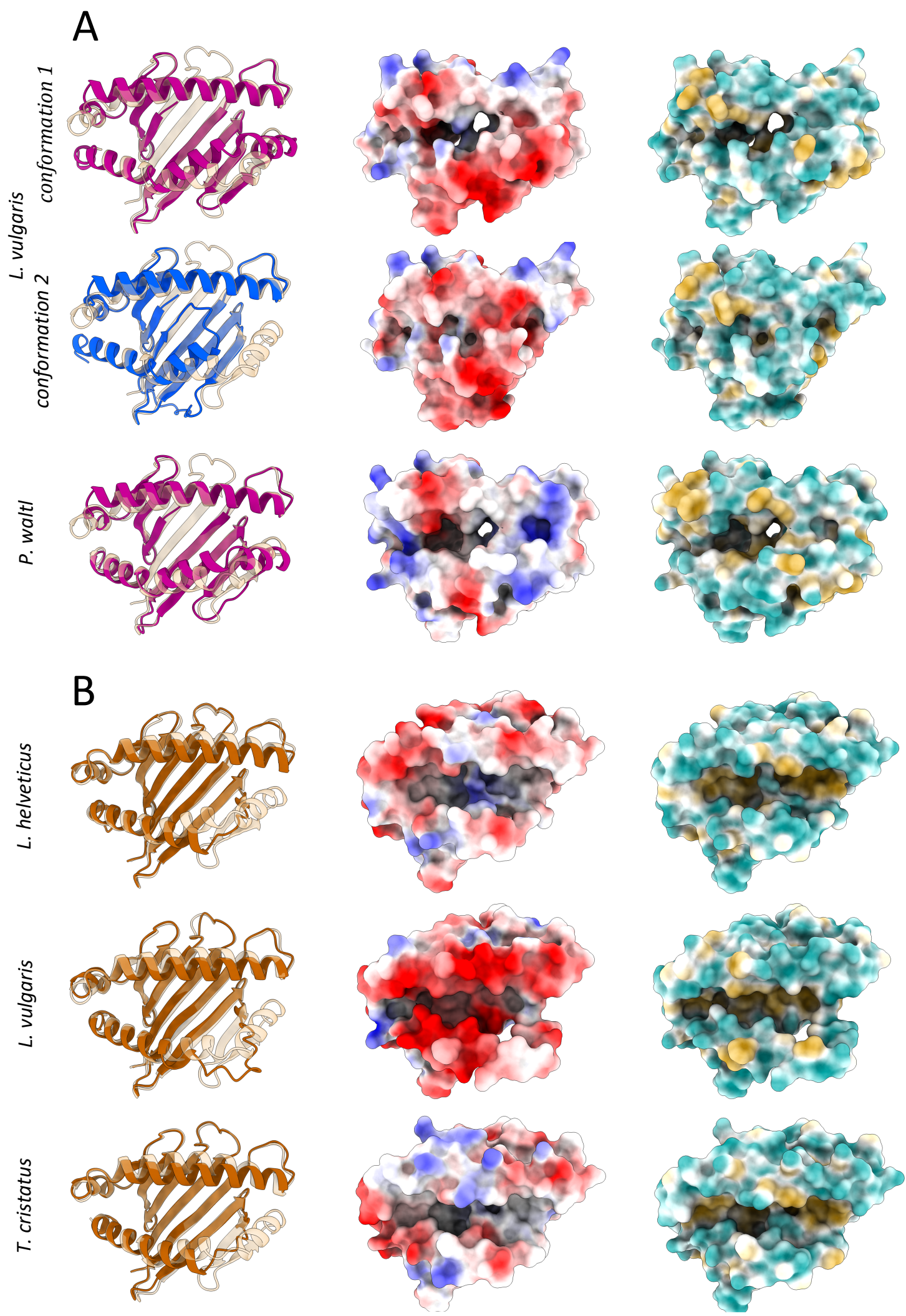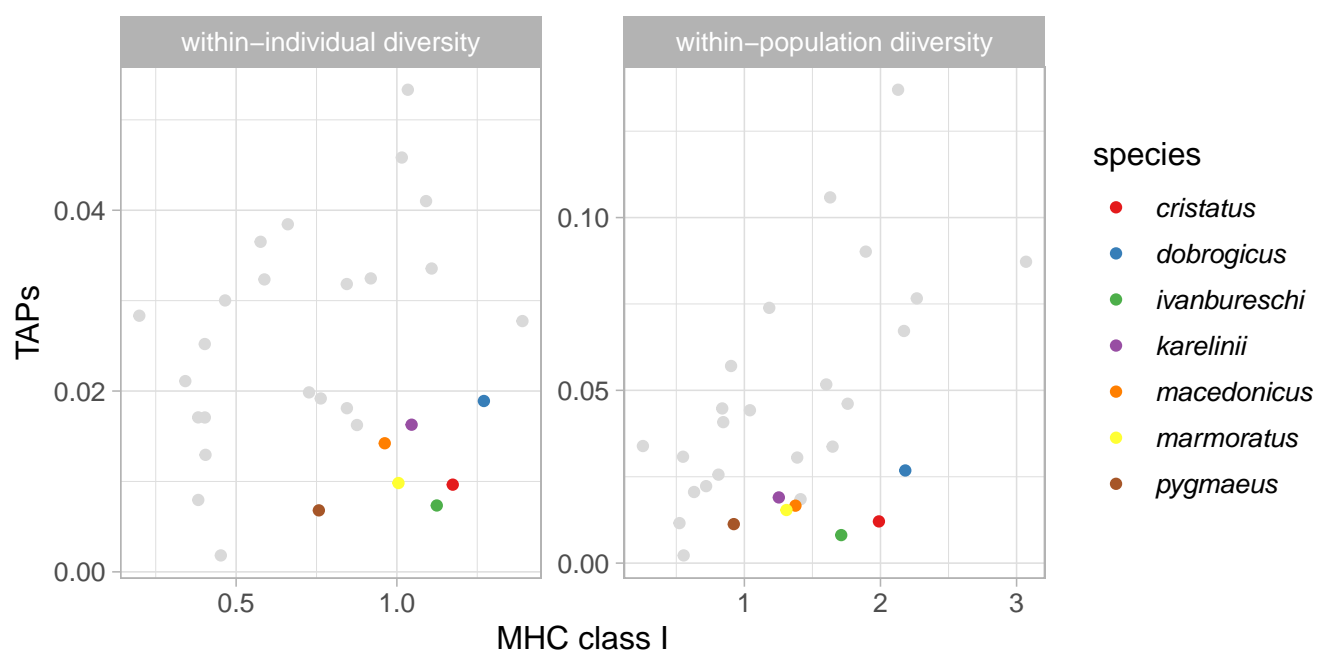
**Fig. S9. Polymorphism of MHC-I genes in *T. cristatus* estimated with previously used and updated primers.** Twenty four individuals that were analysed previously by Palomar et al. (2021, original primers) were genotyped again using primers updated so that they would match exon 2 sequences of all genes annotated in the genome (updated primers). Polymorphism of each gene in both datasets was estimated as described in Materials and Methods and in Supplementary Materials.

A

*L. vulgaris*
conformation 1

*L. vulgaris*
conformation 2

*P. waltl*

B

*L. helveticus*

*L. vulgaris*

*T. cristatus*

**Fig. S10. Top views of α1-α2 domains of a) MHC-I-like1 and b) MHC-I-like2 molecules**, generated with Alphafold3. Left hand-side structures: ribbon models of MHC-I-like molecules (magenta – top conformation of MHC-I-like1, blue – alternative MHC-I-like1 conformation, brown - top conformation of MHC-I-like2), overlaid on classical molecules of respective species. Middle structure: Coulombic electrostatic potential mapped on the molecular surface, with red for negative and blue for positive potential. Right hand-side structures: lipophilicity potential mapped on the molecular surface, with cyan most hydrophilic and goldenrod most lipophilic.

**Fig. S11. The relationship between MHC-I and *TAP* (*TAP1* + *TAP2*) diversity in *Triturus* newts.** Estimates of phylogenetic alpha (within-individual) and gamma (within-population) diversity for 30 salamander species are from Palomar et al. (2021). There is a tendency for *Triturus* species (color-coded) to show low *TAP* diversity and little evidence for a correlation between MHC-I and *TAP* diversity, as opposed to the remaining 23 species representing 15 genera and six families.