



<http://www.diva-portal.org>

## Postprint

This is the accepted version of a paper published in *Developmental and Comparative Immunology*. This paper has been peer-reviewed but does not include the final publisher proof-corrections or journal pagination.

Citation for the original published paper (version of record):

Ericsson, L., Söderhäll, I. (2018)  
Astakines in arthropods-phylogeny and gene structure  
*Developmental and Comparative Immunology*, 81: 141-151  
<https://doi.org/10.1016/j.dci.2017.11.005>

Access to the published version may require subscription.

N.B. When citing this work, cite the original published paper.

This work is licensed under a Creative Commons License CC-BY-NC-ND

Permanent link to this version:

<http://urn.kb.se/resolve?urn=urn:nbn:se:uu:diva-347534>

Astakines in arthropods – phylogeny and gene structure

Lena Ericsson and Irene Söderhäll

Department of Comparative Physiology, Uppsala University, Norbyvägen 18A, 75326,  
Uppsala, Sweden

Corresponding author: Irene Söderhäll, [irene.soderhall@ebc.uu.se](mailto:irene.soderhall@ebc.uu.se)

The editorial process of this article was handled by Mirodrag Belosevic.

## Abstract

Astakine1 was isolated as a hematopoietic cytokine in the freshwater crayfish *Pacifastacus leniusculus*. In this study we detect and compare 79 sequences in GenBank, which we consider to be possible astakine orthologs, among which eleven are crustacean, sixteen are chelicerate and 52 are from insect species. Available arthropod genomes are searched for astakines, and in conclusion all astakine sequences in the current study have a similar exon containing CCXX(X), thus potentially indicating that they are homologous genes with the structure of this exon highly conserved. Two motifs, RYS and YP(N), are also conserved among the arthropod astakines. A phylogenetic analysis reveals that astakine1 and astakine2 from *P. leniusculus* and *Procambarus clarkii* are distantly related, and may have been derived from a gene duplication occurring early in crustacean evolution. Moreover, a structural comparison using the Mamba intestinal toxin (MIT1) from *Dendroaspis polylepis* as template indicates that the overall folds are similar in all crustacean astakines investigated.

Keywords: Astakine; Astakine-like; Prokineticin;

## 1. Introduction

Comparison of protein sequences can provide meaningful insights into how proteins function as well as how they have evolved (Ajawatanawong and Baldauf, 2013). During the past five years, the number of available annotated eukaryotic genomes has increased dramatically, from 40 in 2011 to 359 in 2016 ([https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/)). This increase of genomic data, combined with an even larger increase in reported protein sequences, has made it possible to perform comparisons between sequences from an increasing number of species.

In 2005, we published the first report of a hematopoietic cytokine in freshwater crayfish *Pacifastacus leniusculus* and named the protein astakine (Söderhäll et al., 2005). It was purified from plasma and sequenced by mass spectrometry, and the mRNA sequence was identified and characterized. Astakine was found to contain a prokineticin domain (pfam 06607) with 56% similarity and 31% identity to that of *Bombina variegata* Bv8 (GenBank accession no. [AAD45816](#)) and high similarity to other vertebrate prokineticins (Mollay et al., 1999). This astakine is now named as astakine 1 (Ast1). Further, we isolated cDNA for a second astakine from the penaeid shrimp *Penaeus monodon*, which had an insertion of 13 amino acids after amino acid 44 compared with *P. leniusculus* astakine1 (Söderhäll et al., 2005). . Later, we identified a similar, astakine2 in *P. leniusculus*, and this longer types of astakines is now named as astakine 2 (Ast2) (Lin et al., 2010). Since then, several astakine-like proteins have been described from different invertebrate animals, primarily arthropods (Hsiao and Song, 2010; Lin et al., 2010; Li et al., 2016; Shelby et al., 2015).

The prokineticin protein was originally isolated from black mamba venom (Boisbouvier et al., 1998) and then from skin secretions of frogs (Mollay et al., 1999). Prokinectins are 80-90 amino acids in length, and contain 10 cysteines forming 5 bridges. The amino-terminal

sequence in all vertebrate prokineticins is AVIT, and in addition to being present in snake venoms and frog skin secretions, these proteins are expressed in various tissues in mammals (Kaser et al., 2003). Vertebrate prokineticins are involved in not only angiogenesis and cancer (Monnier and Samson, 2010), immunity (Martucci et al., 2006) and hematopoiesis (LeCouter et al., 2004) but also reproduction (Wechselberger et al., 1999) pain regulation (Negri et al., 2009, 2002) and neural repair (Gordon et al., 2016). Moreover, two highly homologous G-protein coupled receptors for prokineticins have been identified (Lin, 2002).

A common trait of all arthropod astakines is that they lack the N-terminal sequence AVIT, which is a signature sequence for vertebrate prokineticins and is important for binding to their G-protein coupled receptors, PROKR1 and PROKR2 (Kaser et al., 2003). To date, no similar receptor has been detected for the invertebrate astakines, but binding studies have shown that *P. leniusculus* astakine binds to the beta subunit of ATP synthase (Lin et al., 2009) a finding later confirmed to also occur in shrimp (Liang et al., 2015).

Several arthropod protein sequences with similarity to that of crayfish astakine can be found in GenBank, but only a few studies about the function of this group of proteins have been published. An important role of *P. leniusculus* Ast1 in hemocyte proliferation and differentiation has been described (Lin et al., 2010) and more reports indicating roles of astakines in immunity and hematopoiesis have been published (Hsiao and Song, 2010; Jiravanichpaisal et al., 2007; Liang et al., 2015; Lin et al., 2008; Li et al., 2016; Shelby et al., 2015; Thomas et al., 2016; Wilson et al., 2015). However, knowledge of arthropod astakine functions remains scarce, and to date, no structure has been experimentally determined for any of these proteins, although we have performed homology modeling for *P. leniusculus* astakine1 and astakine2 by using mamba intestinal toxin 1 as a template (Lin et al., 2010).

Two structures have been determined experimentally for vertebrate prokineticins (Boisbouvier et al., 1998; Morales et al., 2010). A solution structure of the disulfide-bridge topology of mamba intestinal toxin 1 (MIT1), determined by NMR spectroscopy, reveals similarities with colipase (an enzyme secreted from pancreas). Both peptides show resistance to endoproteases, and the authors have suggested that exocrine glands such as the pancreas may have evolved into venom glands, owing to the structural similarities between colipase and Mamba intestinal toxin (MIT1) from *Dendroaspis polylepis* (Boisbouvier et al., 1998). Interestingly, several astakine-like sequences in arthropods show similarities with venom proteins in insects and spiders.

In the present study, we searched 27 arthropod genomes to find genes encoding astakine-like proteins and searched for additional astakine-like protein or cDNA sequences in GenBank. In total, we detected 79 sequences, which we deemed to be possible astakine orthologs. We compared these sequences, focusing on differences in putative indel sequences to identify possible structures that may be of interest for further functional studies. Further, we used the software Phyre2 to compare the putative 3-dimensional structures of crustacean astakines.

## 2. Materials and methods

### 2.1 Naming and definition of astakines

The astakine sequences were divided into two groups, astakine 1 (Ast1) and astakine 2 (Ast2), according to the naming in *P. leniusculus* where these molecules were first defined (Lin et al., 2010; Söderhäll et al., 2005). Ast1 contains a prokineticin domain with 10 cysteines, whereas Ast2 in addition has an insert of 10-20 amino acids containing the conserved YP(N/D) motif. Naming of the proteins was done as follows: the protein name begins with an abbreviation of the species followed by Ast1 or Ast2. When multiple copies were found, lower case letters were added in alphabetic order as additional identifier. For example, the two *Stegodyphus*

*mimosarum* Ast2 were named St-Ast2a and St-Ast2b. If several species had similar initials we named as in the following example: *Procambarus clarkii*= Pcl; *Polistes Canadensis*= Pca, or *Camponatus floridanus* = Ca.f; *Copidosoma floridanus*= Co.f.

## 2.2. BLAST search and sequence collection

Seventy-nine astakine-like sequences from different arthropods were investigated in this study. The sequence comparison was limited to the prokineticin domain, which was deemed to start with the second amino acid located to the N-terminal side of the first cysteine residue in the N-terminus and to end with the second amino acid after the tenth cysteine in the C-terminus. Sequences were numbered starting with 1 at the second amino acid preceding the first C and extending to the second amino acid after the tenth C (i.e., X<sub>1</sub>X<sub>2</sub>C<sub>3</sub> – C<sub>94</sub>X<sub>95</sub>X<sub>96</sub>). We searched for astakine-like sequences in several different ways. Astakine and astakine-like were used as keywords to search in GenBank at the NCBI web page (<https://www.ncbi.nlm.nih.gov/>). The protein and nucleotide sequences from *P. leniusculus* Ast1 and Ast2 (accessions AAX14635.1 and ABQ23256.1 respectively), were used as query sequences in BLAST searches (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) by using Protein BLAST, blastx and tblastn. The resulting astakine-like sequences from different arthropods were used for further searching via BLAST. Non-redundant protein sequences or specific arthropod genomes were selected as the BLAST databases. In total, 79 arthropod astakine-like sequences were found and used for further analysis (Supplementary table 1). Some of the astakine-like sequences were found only as nucleotide sequences, and in those cases, we used the ExPASy translate tool (<http://web.expasy.org/translate/>) to translate them into protein sequences (Artimo et al., 2012). As cutoff value for identity of 30 %, and moreover the conserved ten cysteine pattern was used as criteria for naming the sequence as astakine. In

some species, we found more than one astakine-like sequence, and in those cases, additional identifiers were assigned as described in section 2.1.

### *2.3. Investigation of the exon structure of arthropod astakines*

To explore the exon-intron structure of putative astakine genes, 33 different arthropods found in the GenBank genome assembly database

([https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/all/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/all/)) were analyzed with BLASTn for full-length astakine mRNA sequences from these species. BLAST hits, i.e., fragments of the mRNA sequences in the respective genomes, were downloaded. The nucleotide fragments were translated into amino acid sequences to determine the astakine protein sequences and thereby the exon structure and the length of the introns were estimated.

### *2.4 Multiple sequence alignments and phylogenetic analyses*

To compare the different astakine-like sequences, we performed multiple sequence alignment of their protein sequences by using the Muscle (Edgar, 2004) online tool from EMBL-EBI (<http://www.ebi.ac.uk/Tools/msa/muscle/>).

Phylogenetic trees were constructed for two different datasets of the astakine sequences. Two different methods were used for both datasets; the Bayesian method using MrBayes 3.2, (Ronquist et al., 2012) and the maximum likelihood (ML) method using the IQ-TREE-1.5.5 software (Nguyen et al., 2015). The best-fit evolution model for the ML analysis of the 79 sequences was WAG+I+G4 and for the ML analysis of the 33 sequences VT + I + G4. WAG combines two empirical models of protein evolution, Dayhoff and JTT, using an approximate maximum likelihood method (Whelan and Goldman, 2001). For the Bayesian analysis the best model was estimated for both datasets to be WAG. The rates were set to equal. IQ-TREE estimates the appropriate evolutionary model using Modelfinder (Kalyaanamoorthy et al., 2017). Ultrafast bootstrap approximation (Minh et al., 2013) was used to assess branch

support values. The number of replicates was set to 1000. Bayesian phylogenetic inference uses Markov chain Monte Carlo (MCMC) methods to produce the most likely phylogenetic tree for a given set of data. One of the dataset contained all the 79 astakine protein sequences in this study. The other dataset consisted of 33 astakine sequences based on sequences showing > 35 % identity to Pl-Ast2. All of the four trees were rooted by an astakine-like sequence from the collembolan hexapod *Folsomia candida*.

## 2.5. Detection of signal peptides and calculation of isoelectric points and molecular weights

The presence and locations of putative cleavage sites for signal peptides of the astakine sequences were predicted by the SignalP 4.1 server (Petersen et al., 2011) (<http://www.cbs.dtu.dk/services/SignalP/>). D-cutoff values were set to default (meaning the score above which the SignalP program will predict a cleavage site for a signal peptide for eukaryotes), and input sequences were allowed to include TM regions. After removal of the signal peptide sequences, the reduced astakine sequences were analyzed with the ExPASy compute pI/Mw tool web site ([http://web.expasy.org/compute\\_pi/](http://web.expasy.org/compute_pi/)), and average resolutions were used to calculate isoelectric points and molecular weights assuming no glycosylation or lipid binding of the proteins (Bjellqvist et al., 1993).

## 2.6. Structure prediction

Three-dimensional structures of full-length mature proteins (without signal peptide) of *L. vannamei*, *P. monodon*, *M. japonicus*, *P. clarkii*, *P. leniusculus* Ast2 and *P. leniusculus* Ast1 were predicted using the software Phyre2 (Kelley et al., 2015) at the Phyre2 web page (<http://www.sbg.bio.ic.ac.uk/phyre2/html/page.cgi?id=index>), using normal modeling mode. The three top-scoring models for each crustacean astakine were downloaded in Protein Data



Bank (PDB) format, and PyMol (Sanner, 1999) was used to display the models. Charge potentials for the protein models were computed with Adaptive Poisson-Boltzmann Solver (APBS) (Baker et al., 2001) and mapped to the surface with medium quality and a distance of 1.0 Å from the surface.

### 3. Results and Discussion

#### 3.1. Astakines or astakine-like sequences

In addition to *P. leniusculus* Ast1 and Ast2, we detected 77 putative astakine protein sequences in GenBank or by BLAST searches, and performed multiple sequence alignment of the conserved prokineticin-like domain of these sequences together with an astakine-like sequence from the primitive hexapod *Folsomia candida* (Springtail) as an outgroup sequence (Figure 1).

The alignment files were then examined manually, and numbering of the amino acids is as described in the method section. Ten highly conserved cysteine residues, two motifs (R<sub>22</sub>Y<sub>23</sub>S<sub>24</sub> and Y<sub>57</sub>P<sub>58</sub>) and a conserved proline residue (P<sub>80</sub> between C<sub>79</sub> and C<sub>81</sub>) were found in the alignment (numbers according to Figure 1). Seventy-nine sequences were defined as Ast1 or Ast2, among which 11 were crustacean, 16 were chelicerate and 52 were from Insecta (Figure 1, Supplementary table 1). In addition to the highly conserved amino acid motifs mentioned above, some amino acids were identified as being more or less conserved in specific groups of species. For crustaceans, the residues G<sub>11</sub>P<sub>12</sub> and P<sub>30</sub>L<sub>31</sub>G<sub>32</sub>D/E<sub>33</sub> and T<sub>95</sub>C<sub>96</sub>Q<sub>97</sub> were highly conserved among decapods.

Among the 52 insect sequences, including 13 diverse Hemiptera, 15 ants, and 21 other Hymenoptera, the highest similarities between sequences were detected among the ants. The amino acid P<sub>20</sub> and the motifs M<sub>29</sub>P<sub>30</sub>F/Y<sub>31</sub>Q<sub>32</sub>Q<sub>33</sub> and T<sub>49</sub>I<sub>50</sub>T<sub>51</sub>T<sub>53</sub>N<sub>54</sub>L<sub>55</sub>T<sub>56</sub> are highly conserved among ants but less conserved in other insect species. In bees, we also identified a

conserved Q<sub>27</sub>, which is also present in shrimp (Figure 1). Unfortunately, there have been few studies about the functions of the insect and chelicerate astakines, and therefore it is premature to draw any functional conclusions regarding this conservation of motifs at the organism level.

In conclusion, we could detect astakine sequences in chelicerats, crustaceans, and some insect orders. So far no astakine sequence was found in Myriapoda, which could be due to lack of sequence data for this group. Interestingly, most of the insect sequences belonged to insects of the orders Hemiptera and Hymenoptera, whereas no astakine-like sequences have been reported to date from Diptera, Coleoptera or Lepidoptera. According to several recent phylogenomic studies of insects, it seems clear that Hymenoptera is a basal order within the holometabolous group and is a sister group to Diptera, Coleoptera or Lepidoptera (Behura, 2015; Peters et al., 2017, 2014). Thus, the data presented in our study suggest that astakine genes may have been lost at the root of the clade Aparaglossata, which includes all holometabolous insects except Hymenoptera (Peters et al., 2014).

### 3.2 Signal peptides, isoelectric points and molecular weights in arthropod astakines

We could find cleavage sites for signal peptides in all astakine protein sequences. However, when we compared the *Atta cephalotes* (GenBank Accession number XP\_012063524.1) sequence with that of the close relative *Atta colombica* (GenBank Accession number KYM75707.1), these sequences were nearly identical, except for a longer N-terminal reported in *A. cephalotes*. This *A. cephalotes* sequence was predicted by an automatic analysis from genomic data and submitted as such to GenBank, whereas the *A. colombica* sequence was experimentally identified from transcriptomic sequences. This result indicates that in order to verify protein sequences, experimental confirmation is needed. Thus, the long deduced N-

terminal sequence in *A. cephalotes* reported as predicted astakine-like protein with Accession number XP\_012063524.1.

In the crustacean group, the isoelectric point for most of the investigated crustacean astakines in the current study varied between pI = 4.54 and pI = 5.13. However, two exceptions were Pl-Ast2 (pI = 7.04) and Pcl-Ast2 (pI = 7.69) (Supplementary table 1).

The isoelectric points for chelicerate astakines varied between pI = 3.92 and pI = 8.54. Most of the identified astakines had pI values higher than 4.50 and lower than 6.80, and hence they are negatively charged at neutral or physiological pH.

In summary, all astakines detected in our study have a predicted signal peptide and thus may be secreted proteins. That is similar to the prokineticins in vertebrates, which all are secreted. All astakines are small molecules with molecular mass between eight and fourteen kDa, and the prokineticin domain constitutes the main part of the mature protein. However, in contrast to vertebrate prokineticins, most of the arthropod astakines, with some exceptions have a pI below 7, meaning a negative charge at physiological pH. However, there are no studies published so far about experimentally determined structure for any arthropod astakine, and therefore it can not be concluded whether their surface charges are negative or positive with certainty.

### *3.3. Putative exon-intron structures of arthropod astakines*

We identified putative arthropod astakine sequences by searching in the arthropod genomes that are annotated at NCBI. Our deduced gene structures for the identified astakine genes among these arthropods showed some general similarities as well as some differences among the major classes. In all astakine genes investigated in the current study, the exon containing the 3' end of the prokineticin domain ends with CCXX or CCXXX (Figure 2 and

Supplementary table 2). Most of the genes in insects consist of two exons in total, with the first one ending as mentioned above (Figure 2). Four of the astakine genes from insects (*A. cephalotes*, *Linepithema humile*, *Solenopsis invicta* and *Vollenhovia emeryi*) have an additional exon in the N-terminal region, and in the *H. halys* astakine gene, there is an additional exon in the C-terminal region (Supplementary table 2). However, this conclusion must be considered carefully, because these sequences do not seem to have been confirmed experimentally; as mentioned above for *A. cephalotes*, the predicted first exon may not be expressed or may be inaccurate (Suen et al., 2011; Wurm et al., 2011).

Among the chelicerate astakine genes analyzed, most had an extra exon at the N-terminus encoding the signal peptide, and the second exon encoding the structure ending with CCXX or CCXXX (Figure 2). In *M. occidentalis*, an exception among the chelicerates, the PROK-domain is encoded by three exons, and the second exon ends with nucleotides encoding **CPCEG** (Supplementary table 2). No decapod genome is available to date, but the N-terminal exon has been found to share the same structure (ending with the nucleotides encoding CCXX or CCXXX) in the unpublished genome of the decapod marble crayfish, *Procambarus fallax* forma *virginialis* (Phattarunda Jaree, Frank Lyko and Julian Gutekunst, personal communication).

In summary, all astakine sequences in the current study have a similar exon encoding the structure ending with CCXX(X), thus potentially indicating that they are homologous genes with the structure of this exon highly conserved. In most of the astakines the PROK-domain is encoded by two exons, the first one ending as above and the second encoding the rest of the prokineticin domain. This result may indicate that these astakine sequences are more closely related to each other than the ones with the prokineticin domain encoded for by three exons. Some of the astakines appear to have an additional exon located at the N-terminus of the gene.

273 Because this part of the protein contains the signal peptide, it may be less conserved than  
274 other regions.

275

### 276 3.4 Multiple sequence alignment of arthropod astakines

277

278 The alignment of all astakine amino acid sequences showed that some of the residues are  
279 highly conserved. Ten cysteine residues are conserved in all sequences (numbering as in  
280 figure 1):

281 C<sub>3</sub> - (X<sub>5</sub>) - C<sub>10</sub> - (X<sub>4</sub>) - C<sub>15</sub> - C<sub>16</sub> - (X<sub>11</sub>) - C<sub>28</sub> - (X<sub>9</sub>) - C<sub>38</sub> - (X<sub>n</sub>) - C<sub>79</sub>(P)C<sub>81</sub> - (X<sub>5</sub>) - C<sub>89</sub> -  
282 (X<sub>n</sub>) - C<sub>96</sub>

283 Deviating from the structure above are two insertions of a P in Pl-Ast1 and Pcl-Ast1 between  
284 C<sub>3</sub> and C<sub>10</sub>. In Pcl-Ast1, there are also four deletions between C<sub>16</sub> and C<sub>28</sub>. Between C<sub>81</sub> and  
285 C<sub>89</sub> the two *Daphnia* sequences, Dm-Ast2 and Dp-Ast2, have two insertions consisting of an  
286 alanine and an asparagine (Figure 1).

287 Two motifs, RYS and YP(N/D), are conserved among the arthropod astakines. Only in  
288 PclAst1, the YS part of the RYS is missing in the alignment, and in the RYS motif, arginine is  
289 in some sequences substituted by the similar amino acid lysine and in one sequence, that of  
290 *Cimex lectularius* Cm-Ast2c, it is replaced by leucine. In 14 of the sequences, tyrosine is  
291 replaced by phenylalanine. The serine residue of the RYS is in seven of the insects replaced  
292 by an alanine, in Ast2 from *Polistes canadensis* and *V. emeryi* by valine and in four other  
293 insects by a threonine. In the YP(N/D) motif, the tyrosine is replaced by phenylalanine in *M.*  
294 *occidentalis* by tryptophan in *Copidosoma floridanum*, and by glutamine in *Diachasma*  
295 *alloeum*. It has previously been shown by mutant recombinant protein experiments that this  
296 motif is important for the function of Pl-Ast2 (Lin et al., 2010). In Pl-Ast1 and Pcl-Ast1, there

are 24 gaps between C<sub>38</sub> and C<sub>79</sub>. These two sequences also lack the YP(N/D) motif, which indicates different functions of Ast1 and Ast2.

There is also a proline residue in the position between C<sub>79</sub> and C<sub>81</sub>. This residue is conserved in all sequences except in Lhu-Ast2 in (replaced by serine), Tz-Ast2 (alanine), Lh-Ast2b and Zn-Ast2 (aspartic acid), Bg-Ast2 (glycine) and Hl-Ast2 (leucine), and if this finding is not due to sequencing errors, it may have implications for the function of these putative astakines (Figure 1).

Taken together, all the astakines have a conserved cysteine pattern with ten cysteines, and between the sixth and seventh cysteine there is an insertion of variable length in all Ast2 containing an YP(N/D) motif which for Pl-Ast2 is shown to be of importance for the function (Lin et al., 2010).

### 3.5. Phylogenetic analyses of some arthropod astakines

All of the arthropod astakine sequences were subjected to two different phylogenetic analyses, a maximum likelihood (ML) analysis by the IQ-Tree software (Figure 3a), and a Bayesian analysis by MrBayes software (Figure 3b). The study was restricted to these methods since distance based methods are less reliable when analyzing the high number of diverse sequences as in this study. An astakine-like sequence from the springtail *Folsomia candida* was used to root the tree. In the resulting IQ-TREE-file from the phylogenetic analysis there is a warning that deduction of the phylogeny should be done with caution. This was due to the larger number of parameters (branch lengths and model parameters) in relation to the sample size i.e. the length of the alignment. In order to improve the robustness in the phylogenetic estimation and avoid warnings, a second dataset was constructed of 33 sequences (Figure 4a-b). The sequences included were chosen by their percentage of identity to Pl-Ast2. Seven

sequences were crustacean, eight from chelicerates and 18 were insect astakines, and *F. candida* was used as root sequence. Since the number of parameters depends on the number of sequences there was no warning for the second analysis in IQ-Tree (Figure 4a).

Ast2 sequences from chelicerata and insecta are clustered in two distinct clades in all four analyses. The placement of the crustacean astakines seems to be more uncertain. In the smaller dataset, all the crustacean astakines are gathered in one clade with high support. However, in the analysis of the large dataset the topologies of the crustacean sequences are different depending on the phylogenetic method. Using the maximum likelihood method, the crustacean astakines Dm-Ast2 and Dp-Ast2 belonging to the Cladocera branches of early in a minor clade and differ from all the other astakines, which belongs to Decapoda. In contrast, in the Bayesian method, the crustacean sequences can not be fully resolved, and the decapod Cm-Ast2 is found outside all the other crustaceans in a polytomy. A comparison between the sequence structure of the crustacean astakines shows that Cm-Ast2 contains deviating amino acids in 21 positions. Nine of these positions contain amino acid residues not found in any other astakine sequence in this study. For the rest of the 21 positions identical residues have been found in some of the astakines from chelicerates and insects, and thus Cm-Ast2 is different from the other astakines of Decapoda.

In all phylogenetic analyses, Ast1 and Ast2 from *P. clarkii* and *P. leniusculus* were separated into different groups in the trees, indicating that these sequences are distantly related to each other. Gene duplication may have occurred in the crustacean astakines before the divergence of the species included in this analysis. Ast1 has to date been detected in only *P. leniusculus*, *P. clarkii* (Beltz and Brenneis, personal communication) and *P. fallax* forma *virginalis* (Jaree, Lyko and Jutekunst, personal communication), and it is possible that one of the variants has been lost during evolution in some groups (or has not yet been found).

347 All other sequences analyzed in this study belong to the Ast2 type. For some species more  
 348 than one sequence was found. For example, three duplicates of Ast2 from *L. polyphemus* were  
 349 grouped together, and thus are more closely related to one another than Ast1 and Ast2 in *P.*  
 350 *leniusculus*.

351 The dataset also contains two astakine sequences for the American house spider, *P.*  
 352 *tepidariorum*, and the African social velvet spider, *S. mimosarum*. Even if their placement is  
 353 somewhat uncertain, they are related to each other in the same order in both the ML and the  
 354 bayesian phylogenetic analyses (Figure 3a-b). Pt-Ast2a groups with Sm-Ast2a, and Pt-Ast2b  
 355 groups with Sm-Ast2b, thus indicating earlier gene duplication before speciation, compared  
 356 with the evolution of astakine duplicates in *L. polyphemus*.

357 The chelicerate astakines are clustered in almost the same way in all four analyses. However,  
 358 the phylogeny of De-Ast2 (*D. erythrina*), Sm-Ast2a and Pt-Ast2a could not be fully resolved  
 359 in this analysis, since there is a polytomy in both the ML and Bayesian analysis.

360

361 A large number of sequences in the current study belong to the insects. Several of the species  
 362 have more than one Ast2 sequence namely, *A pisum*, *L Hesperus*, *C lectularius*, *T pretiosum*,  
 363 *N vitripennis* and *A echinator*. In the insect clade, all astakines from hymenoptera except Da-  
 364 Ast2 are clustered in one large clade. Da-Ast2 is found in another clade together with  
 365 astakines from Phthiraptera, Hemiptera, Blattodea and Isoptera (Figure 3a-b). The topology of  
 366 this clade is similar but not identical in both trees, but the support values in this area of the  
 367 trees are lower in the Bayesian tree (Figure 3a). A comparison of the sequence structures of  
 368 the insect astakines shows that Da-Ast2 also contains several different amino acids, compared  
 369 to other sequences. Some of the residues are identical to the ones in Lhe-Ast2b. Therefore, it  
 370 is possible that Da-ast2 is another astakine variant than the other hymenoptera astakines.



Taken together, our phylogenetic trees give a hypothetic indication about the evolutionary relationship between astakine sequences, but it has to be taken into account that such tree analysis are limited by the number of sequences available. When all detected astakine sequences were used in one analysis the different trees were similar but several branches showed low values of support (Figure 3a-b). In contrast, the limited analysis in which sequences of high identity were used showed more robust trees with higher support values (Figure 4a-b).

### 3.6. *Structure prediction*

Six crustacean astakine sequences from Pm-Ast2, Lv-Ast2, Mj-Ast2, Pcl-Ast2, Pl-Ast2 and Pl-Ast2 were analyzed with Phyre2 for alignment (Kelley et al., 2015). These sequences were used for comparison in order to get an idea about what parts of the structure that is most likely to be of importance for functional difference between some marine and freshwater crustacean species. Phyre2 determines an evolutionary profile for the query sequence by heuristic searches in protein sequence databases. To search for the best templates, this profile, together with the secondary structure predicted in Phyre2, is scanned against a folding library containing proteins of known, experimentally determined structures. The best-scoring alignments between the query sequences and the library sequences are then used to build three-dimensional models of the query protein. The three top-scoring models from the Phyre2 results were identical for all analyzed sequences: mamba intestinal toxin 1 from *Dendroaspis polyepis*, PDB 1MIT (Boisbouvier et al., 1998); prokineticin Bv8 from *Bombina variegata* PDB 2KRA (Morales et al., 2010); and Dickkopf-related protein 1 (DKK1) from *Homo sapiens* PDB 3S8V (Cheng et al., 2011). These three templates gave different alignment with the astakines as shown in Supplementary figure 1. The confidence for the matches between the submitted astakine sequences and the models was between 98.5 and 100, thus indicating a

396 high percentage probability that the astakines and the models are homologous. The percentage  
397 identity between the astakines and the models was between 28% and 37%. In Phyre2, the  
398 proportion of disorder in secondary structures was predicted for all astakine sequences as  
399 reported in supplementary table 3. Thus, fairly large portions of these proteins probably lack  
400 fixed three-dimensional structures and are unstructured with conformational flexibility, owing  
401 to random coil structures. Structures with high proportions of disorder are more difficult to  
402 predict. The three models with high confidence for the matches, as mentioned above, in  
403 Phyre2, were then used as templates for the prediction of astakine structure, and the predicted  
404 structures were displayed in Phyton Molecule Viewer (Sanner, 1999). The overall folding and  
405 the core of the astakine structures for these models are shown in Figure 5a and Supplementary  
406 figures 2 and 3, with the ten cysteine residues, the RYS and YP(N/D) motifs and the indel  
407 regions marked (shown by yellow crosses, Figure 5a, Supplementary figures 2 and 3).

408 Figure 5a shows the structure of the astakines predicted with intestinal toxin 1 as template.  
409 The extension of the modeled astakines is between GXC<sub>3</sub> and C<sub>89</sub>XRXX, in the Pcl-Ast2  
410 model between HC<sub>3</sub> – C<sub>89</sub>SRTS (numbering according to Figure 1). The overall folding of the  
411 structures appears to be quite similar among the species. Most of the models contain four  
412 cysteine bridges, C<sub>3</sub> – C<sub>16</sub>, C<sub>10</sub> – C<sub>28</sub>, C<sub>15</sub> – C<sub>79</sub> and C<sub>38</sub> – C<sub>89</sub>, although the two cysteine  
413 residues C<sub>3</sub> and C<sub>16</sub> from Pcl-Ast2 and Pl-Ast1 models appear to be too distant from each  
414 other to form a bridge. In contrast, the secondary structures are less similar, especially in the  
415 region partly consisting of the indel region (Figure 5a). In this region, an alpha helix is found  
416 in Ast2 from *L. vannamei* and *P. monodon*, whereas Pl-Ast2 contains one helix and one beta  
417 sheet, and in Ast2 from *M. japonicus* and *P. clarkii*, the structures in this region consist only  
418 of coil structure. The predicted structure of Pl-Ast1, which has a deletion of 13 amino acid  
419 residues in this part of the structure, has a beta sheet. For the astakine structures with Bv8 as  
420 templates (Supplementary figure 2), the extension of the modeled residues is almost the same

421 as in the previous prediction. The only difference is one additional residue in the C-terminus.  
 422 All structures contain four cysteine bridges. The overall folds are similar in all crustacean  
 423 astakines, and the secondary structures are more similar within the two groups of astakines  
 424 compared with the structures with the intestinal toxin as template. Among the shrimp, an  
 425 alpha helix is found in the indel region, whereas only coil structures are found in the two Ast2  
 426 structures from *P. clarkii* and *P. leniusculus* (Supplementary figure 3).  
 427  
 428 In comparison, the predicted structures using the Dickkopf-related protein 1 (DKK1, 3S8V)  
 429 as template include all ten cysteine residues forming five cysteine bridges, including C<sub>81</sub> – C<sub>96</sub>  
 430 (Supplementary figure 3). Most of the modeled structures contain XC<sub>3</sub> in the N-terminus, but  
 431 the Pl-Ast1 structure contains GSC<sub>3</sub> and Pcl-Ast2 only C<sub>3</sub>. All of the structures end with C<sub>10</sub>Q,  
 432 except that of Pcl-Ast2, which ends with C<sub>96</sub>QL. Another difference is the overall folding of  
 433 the structures predicted from DKK1. These differences include, for example, the RYS motif  
 434 being located closer to the YP(N/D) residues. *L. vannamei* and *P. monodon* Ast2 have similar  
 435 secondary structures containing an alpha helix in the indel region, as does *P. clarkii*. In *M.*  
 436 *japonicus*, only a small helix is found, and in the two astakines from *P. leniusculus*, the indel  
 437 region consists only of coiled structure (Supplementary figure 3).  
 438  
 439 Two regions, the RYS motif and the YP(N/D) motif with an additional asparagine residue,  
 440 are conserved among the Ast2 and were investigated further. The charge potentials were  
 441 computed and mapped to the surface with 1MIT as a template (Figure 5b). The RYS and  
 442 YP(N/D) structures of the surfaces together with the structures of the residues in the overall  
 443 folds (Figure 5a and Supplementary figures 2 and 3) were compared across all the modeled  
 444 astakine structures. The structures of the residues in the RYS motif seem to be similar in the  
 445 models predicted with the intestinal toxin and prokineticin templates, although none of the

template contains this motif. In contrast, in the models predicted with DKK1 as template, the structure of the RYS motif is different in some astakines, whereas no model for these residues in Pcl-Ast2 could be predicted using this template.

All six crustacean structures modeled by Phyre2, except Pl-Ast1, contain the YP(N/D) motif, and this part of the molecule has previously been found to be important for function in the granular hemocyte lineage in *P. leniusculus* (Lin et al., 2010). This motif is near the variable indel region, and therefore it is possible that the structure around this motif might vary in the different models (Figure 5a, Supplementary figure 2-3).

In conclusion, the overall structure of the models predicted with intestinal toxin and prokineticin as templates seems to be similar for most of the predicted astakine structures. The secondary structure and the backbone of the proteins are also similar in most of the regions, although not in the region containing indels. The structure of the RYS motif also looks similar among the astakines modeled by using these two templates. These similarities may be because the two templates are similar, and the same four cysteine bridges are predicted in the modeled structures.

In contrast, in the models predicted with DKK1 as template, the overall folding and number of modeled cysteine residues and resulting number of bridges differ from the others. The structure of the RYS motif is more variable when this template is used, but DKK1 is the only template yielding structures with more similar YP(N/D) structures, possibly because these models contain five cysteine bridges, which may stabilize the structure, and especially the indel region and YP(N/D).

#### 4. Conclusions

Since the first reported astakine sequence in 2005, we could find 77 other arthropod astakines or astakine-like sequences in GenBank. A search in available genomes revealed a similar exon-intron structure among the arthropod astakines. Although all sequences are similar and contain the core astakine structure with ten cysteines, the RYS and YP(N/D) motifs, a phylogenetic analysis combining all arthropods were not fully resolved and gave trees with some polytomies. However, both the ML and the Bayesian method showed clearly separate crustacean, chelicerate and insect clades. In addition, the hymenopteran sequences all grouped together with one exception, and the hemiptera formed a common clade also with only one exception.

In crustaceans, the distance between Ast1 and Ast2 from *P. leniusculus* and *P. clarkii* indicates a gene duplication occurring early in crustacean evolution.

A structural comparison using the Phyre2 software gave some indication of a similar overall core structure, but since the available templates are fairly distant, such predictions has to be evaluated with care. Nevertheless, our structural comparison of five crustacean sequences could still show that the indel sequences following the preserved YP(N/D) motif is likely to give a specific surface structure that varies among species, and can be of specific interest to experimentally manipulate in order to reveal possible function.

#### Acknowledgements

We thank Professor B. S. Beltz and Dr. Georg Brenneis for information about *Procambarus clarkii* astakine 1 sequence, and Professor Frank Lyko and Dr. Julian Gutekunst for insight in *Procambarus clarkii* genome

#### References

Ajawatanawong, P., Baldauf, S.L., 2013. Evolution of protein indels in plants, animals and

496 fungi. BMC Evol. Biol. 13, 140. doi:10.1186/1471-2148-13-140  
 497 Artimo, P., Jonnalagedda, M., Arnold, K., Baratin, D., Csardi, G., de Castro, E., Duvaud, S.,  
 498 Flegel, V., Fortier, A., Gasteiger, E., Grosdidier, A., Hernandez, C., Ioannidis, V., Kuznetsov,  
 499 D., Liechti, R., Moretti, S., Mostaguir, K., Redaschi, N., Rossier, G., Xenarios, I., Stockinger,  
 500 H., 2012. ExPASy: SIB bioinformatics resource portal. Nucleic Acids Res. 40, W597–603.  
 501 doi:10.1093/nar/gks400  
 502 Baker, N.A., Sept, D., Joseph, S., Holst, M.J., McCammon, J.A., 2001. Electrostatics of  
 503 nanosystems: application to microtubules and the ribosome. Proc. Natl. Acad. Sci. U. S. A. 98,  
 504 10037–10041. doi:10.1073/pnas.181342398  
 505 Behura, S.K., 2015. Insect phylogenomics. Insect Mol. Biol. 24, 403–411.  
 506 doi:10.1111/imb.12174  
 507 Bjellqvist, B., Hughes, G.J., Pasquali, C., Paquet, N., Ravier, F., Sanchez, J.C., Frutiger, S.,  
 508 Hochstrasser, D., 1993. The focusing positions of polypeptides in immobilized pH gradients  
 509 can be predicted from their amino acid sequences. Electrophoresis 14, 1023–1031.  
 510 Boisbouvier, J., Albrand, J.P., Blackledge, M., Jaquinod, M., Schweitz, H., Lazdunski, M.,  
 511 Marion, D., 1998. A structural homologue of colipase in black mamba venom revealed by  
 512 NMR floating disulphide bridge analysis. J. Mol. Biol. 283, 205–219.  
 513 doi:10.1006/jmbi.1998.2057  
 514 Cheng, Z., Biechele, T., Wei, Z., Morrone, S., Moon, R.T., Wang, L., Xu, W., 2011. Crystal  
 515 structures of the extracellular domain of LRP6 and its complex with DKK1. Nat. Struct. Mol.  
 516 Biol. 18, 1204–1210. doi:10.1038/nsmb.2139  
 517 Edgar, R.C., 2004. MUSCLE: multiple sequence alignment with high accuracy and high  
 518 throughput. Nucleic Acids Res. 32, 1792–1797. doi:10.1093/nar/gkh340  
 519 Gordon, R., Neal, M.L., Luo, J., Langley, M.R., Harischandra, D.S., Panicker, N., Charli, A.,  
 520 Jin, H., Anantharam, V., Woodruff, T.M., Zhou, Q.-Y., Kanthasamy, A.G., Kanthasamy, A.,

521 2016. Prokineticin-2 upregulation during neuronal injury mediates a compensatory protective  
 522 response against dopaminergic neuronal degeneration. *Nat. Commun.* 7, 12932.  
 523 doi:10.1038/ncomms12932  
 524 Hsiao, C.-Y., Song, Y.-L., 2010. A long form of shrimp astakine transcript: molecular cloning,  
 525 characterization and functional elucidation in promoting hematopoiesis. *Fish Shellfish*  
 526 *Immunol.* 28, 77–86. doi:10.1016/j.fsi.2009.10.016  
 527 Jiravanichpaisal, P., Puanglarp, N., Petkon, S., Donnuea, S., Söderhäll, I., Söderhäll, K., 2007.  
 528 Expression of immune-related genes in larval stages of the giant tiger shrimp, *Penaeus*  
 529 *monodon*. *Fish Shellfish Immunol.* 23, 815–824. doi:10.1016/j.fsi.2007.03.003  
 530 Kalyaanamoorthy, S., Minh, B.Q., Wong, T.K.F., von Haeseler, A., Jermini, L.S., 2017.  
 531 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat. Methods* 14,  
 532 587–589. doi:10.1038/nmeth.4285  
 533 Kaser, A., Winklmayr, M., Lepperdinger, G., Kreil, G., 2003. The AVIT protein family.  
 534 Secreted cysteine-rich vertebrate proteins with diverse functions. *EMBO Rep.* 4, 469–473.  
 535 doi:10.1038/sj.embor.embor830  
 536 Kelley, L.A., Mezulis, S., Yates, C.M., Wass, M.N., Sternberg, M.J.E., 2015. The Phyre2 web  
 537 portal for protein modeling, prediction and analysis. *Nat. Protoc.* 10, 845–858.  
 538 doi:10.1038/nprot.2015.053  
 539 LeCouter, J., Zlot, C., Tejada, M., Peale, F., Ferrara, N., 2004. Bv8 and endocrine gland-  
 540 derived vascular endothelial growth factor stimulate hematopoiesis and hematopoietic cell  
 541 mobilization. *Proc. Natl. Acad. Sci. U. S. A.* 101, 16813–16818.  
 542 doi:10.1073/pnas.0407697101  
 543 Liang, G.-F., Liang, Y., Xue, Q., Lu, J.-F., Cheng, J.-J., Huang, J., 2015. Astakine LvAST  
 544 binds to the  $\beta$  subunit of F1-ATP synthase and likely plays a role in white shrimp  
 545 *Litopenaeus vannamei* defense against white spot syndrome virus. *Fish Shellfish Immunol.* 43,

546 75–81. doi:10.1016/j.fsi.2014.12.015

547 Lin, D.C.-H., 2002. Identification and Molecular Characterization of Two Closely Related G  
 548 Protein-coupled Receptors Activated by Prokineticins/Endocrine Gland Vascular Endothelial  
 549 Growth Factor. *J. Biol. Chem.* 277, 19276–19280. doi:10.1074/jbc.M202139200

550 Lin, X., Kim, Y.-A., Lee, B.L., Söderhäll, K., Söderhäll, I., 2009. Identification and properties  
 551 of a receptor for the invertebrate cytokine astakine, involved in hematopoiesis. *Exp. Cell Res.*  
 552 315, 1171–1180.

553 Lin, X., Novotny, M., Söderhäll, K., Söderhäll, I., 2010. Ancient cytokines, the role of  
 554 astakines as hematopoietic growth factors. *J. Biol. Chem.* 285, 28577–28586.  
 555 doi:10.1074/jbc.M110.138560

556 Lin, X., Söderhäll, K., Söderhäll, I., 2008. Transglutaminase activity in the hematopoietic  
 557 tissue of a crustacean, *Pacifastacus leniusculus*, importance in hemocyte homeostasis. *BMC*  
 558 *Immunol.* 9, 58. doi:10.1186/1471-2172-9-58

559 Li, Y., Jiang, S., Li, M., Xin, L., Wang, L., Wang, H., Qiu, L., Song, L., 2016. A cytokine-  
 560 like factor astakine accelerates the hemocyte production in Pacific oyster *Crassostrea gigas*.  
 561 *Dev. Comp. Immunol.* 55, 179–187. doi:10.1016/j.dci.2015.10.025

562 Martucci, C., Franchi, S., Giannini, E., Tian, H., Melchiorri, P., Negri, L., Sacerdote, P., 2006.  
 563 Bv8, the amphibian homologue of the mammalian prokineticins, induces a proinflammatory  
 564 phenotype of mouse macrophages. *Br. J. Pharmacol.* 147, 225–234.  
 565 doi:10.1038/sj.bjp.0706467

566 Minh, B.Q., Nguyen, M.A.T., von Haeseler, A., 2013. Ultrafast approximation for  
 567 phylogenetic bootstrap. *Mol. Biol. Evol.* 30, 1188–1195. doi:10.1093/molbev/mst024

568 Mollay, C., Wechselberger, C., Mignogna, G., Negri, L., Melchiorri, P., Barra, D., Kreil, G.,  
 569 1999. Bv8, a small protein from frog skin and its homologue from snake venom induce  
 570 hyperalgesia in rats. *Eur. J. Pharmacol.* 374, 189–196.



571 Monnier, J., Samson, M., 2010. Prokineticins in angiogenesis and cancer. *Cancer Lett.* 296,  
 572 144–149. doi:10.1016/j.canlet.2010.06.011  
 573 Morales, R.A.V., Daly, N.L., Vetter, I., Mobli, M., Napier, I.A., Craik, D.J., Lewis, R.J.,  
 574 Christie, M.J., King, G.F., Alewood, P.F., Durek, T., 2010. Chemical synthesis and structure  
 575 of the prokineticin Bv8. *Chembiochem Eur. J. Chem. Biol.* 11, 1882–1888.  
 576 doi:10.1002/cbic.201000330  
 577 Negri, L., Lattanzi, R., Giannini, E., Canestrelli, M., Nicotra, A., Melchiorri, P., 2009.  
 578 Bv8/Prokineticins and their Receptors A New Pronociceptive System. *Int. Rev. Neurobiol.* 85,  
 579 145–157. doi:10.1016/S0074-7742(09)85011-3  
 580 Negri, L., Lattanzi, R., Giannini, E., Metere, A., Colucci, M., Barra, D., Kreil, G., Melchiorri,  
 581 P., 2002. Nociceptive sensitization by the secretory protein Bv8. *Br. J. Pharmacol.* 137, 1147–  
 582 1154. doi:10.1038/sj.bjp.0704995  
 583 Nguyen, L.-T., Schmidt, H.A., von Haeseler, A., Minh, B.Q., 2015. IQ-TREE: a fast and  
 584 effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Mol. Biol.*  
 585 *Evol.* 32, 268–274. doi:10.1093/molbev/msu300  
 586 Petersen, T.N., Brunak, S., von Heijne, G., Nielsen, H., 2011. SignalP 4.0: discriminating  
 587 signal peptides from transmembrane regions. *Nat. Methods* 8, 785–786.  
 588 doi:10.1038/nmeth.1701  
 589 Peters, R.S., Krogmann, L., Mayer, C., Donath, A., Gunkel, S., Meusemann, K., Kozlov, A.,  
 590 Podsiadlowski, L., Petersen, M., Lanfear, R., Diez, P.A., Heraty, J., Kjer, K.M., Klopstein, S.,  
 591 Meier, R., Polidori, C., Schmitt, T., Liu, S., Zhou, X., Wappler, T., Rust, J., Misof, B.,  
 592 Niehuis, O., 2017. Evolutionary History of the Hymenoptera. *Curr. Biol.* 27, 1013–1018.  
 593 doi:10.1016/j.cub.2017.01.027  
 594 Peters, R.S., Meusemann, K., Petersen, M., Mayer, C., Wilbrandt, J., Ziesmann, T., Donath,  
 595 A., Kjer, K.M., Aspöck, U., Aspöck, H., Aberer, A., Stamatakis, A., Friedrich, F., Hünefeld,

596 F., Niehuis, O., Beutel, R.G., Misof, B., 2014. The evolutionary history of holometabolous  
 597 insects inferred from transcriptome-based phylogeny and comprehensive morphological data.  
 598 BMC Evol. Biol. 14, 52. doi:10.1186/1471-2148-14-52  
 599 Ronquist, F., Teslenko, M., van der Mark, P., Ayres, D.L., Darling, A., Höhna, S., Larget, B.,  
 600 Liu, L., Suchard, M.A., Huelsenbeck, J.P., 2012. MrBayes 3.2: efficient Bayesian  
 601 phylogenetic inference and model choice across a large model space. Syst. Biol. 61, 539–542.  
 602 doi:10.1093/sysbio/sys029  
 603 Sanner, M.F., 1999. Python: a programming language for software integration and  
 604 development. J. Mol. Graph. Model. 17, 57–61.  
 605 Shelby, K.S., Perera, O.P., Snodgrass, G.L., 2015. Expression profiles of astakine-like  
 606 transcripts in the tarnished plant bug, *Lygus lineolaris*, exposed to fungal spores of *Beauveria*  
 607 *bassiana*. Insect Mol. Biol. 24, 480–490. doi:10.1111/imb.12175  
 608 Söderhäll, I., Kim, Y.-A., Jiravanichpaisal, P., Lee, S.-Y., Söderhäll, K., 2005. An ancient  
 609 role for a prokineticin domain in invertebrate hematopoiesis. J. Immunol. Baltim. Md 1950  
 610 174, 6153–6160.  
 611 Suen, G., Teiling, C., Li, L., Holt, C., Abouheif, E., Bornberg-Bauer, E., Bouffard, P., Caldera,  
 612 E.J., Cash, E., Cavanaugh, A., Denas, O., Elhaik, E., Favé, M.-J., Gadau, J., Gibson, J.D.,  
 613 Graur, D., Grubbs, K.J., Hagen, D.E., Harkins, T.T., Helmkampf, M., Hu, H., Johnson, B.R.,  
 614 Kim, J., Marsh, S.E., Moeller, J.A., Muñoz-Torres, M.C., Murphy, M.C., Naughton, M.C.,  
 615 Nigam, S., Overson, R., Rajakumar, R., Reese, J.T., Scott, J.J., Smith, C.R., Tao, S., Tsutsui,  
 616 N.D., Viljakainen, L., Wissler, L., Yandell, M.D., Zimmer, F., Taylor, J., Slater, S.C., Clifton,  
 617 S.W., Warren, W.C., Elsik, C.G., Smith, C.D., Weinstock, G.M., Gerardo, N.M., Currie, C.R.,  
 618 2011. The genome sequence of the leaf-cutter ant *Atta cephalotes* reveals insights into its  
 619 obligate symbiotic lifestyle. PLoS Genet. 7, e1002007. doi:10.1371/journal.pgen.1002007  
 620 Thomas, A., Sudheer, N.S., Kiron, V., Bright Singh, I.S., Narayanan, R.B., 2016. Expression

621 profile of key immune-related genes in *Penaeus monodon* juveniles after oral administration  
 622 of recombinant envelope protein VP28 of white spot syndrome virus. *Microb. Pathog.* 96, 72–  
 623 79. doi:10.1016/j.micpath.2016.05.002  
 624 Wechselberger, C., Puglisi, R., Engel, E., Lepperdinger, G., Boitani, C., Kreil, G., 1999. The  
 625 mammalian homologues of frog Bv8 are mainly expressed in spermatocytes. *FEBS Lett.* 462,  
 626 177–181.  
 627 Whelan, S., Goldman, N., 2001. A general empirical model of protein evolution derived from  
 628 multiple protein families using a maximum-likelihood approach. *Mol. Biol. Evol.* 18, 691–  
 629 699.  
 630 Wilson, W., Lowman, D., Antony, S.P., Puthumana, J., Bright Singh, I.S., Philip, R., 2015.  
 631 Immune gene expression profile of *Penaeus monodon* in response to marine yeast glucan  
 632 application and white spot syndrome virus challenge. *Fish Shellfish Immunol.* 43, 346–356.  
 633 doi:10.1016/j.fsi.2014.12.032  
 634 Wurm, Y., Wang, J., Riba-Grognuz, O., Corona, M., Nygaard, S., Hunt, B.G., Ingram, K.K.,  
 635 Falquet, L., Nipitwattanaphon, M., Gotzek, D., Dijkstra, M.B., Oettler, J., Comtesse, F., Shih,  
 636 C.-J., Wu, W.-J., Yang, C.-C., Thomas, J., Beaudoin, E., Pradervand, S., Flegel, V., Cook,  
 637 E.D., Fabbretti, R., Stockinger, H., Long, L., Farmerie, W.G., Oakey, J., Boomsma, J.J.,  
 638 Pamilo, P., Yi, S.V., Heinze, J., Goodisman, M.A.D., Farinelli, L., Harshman, K., Hulo, N.,  
 639 Cerutti, L., Xenarios, I., Shoemaker, D., Keller, L., 2011. The genome of the fire ant  
 640 *Solenopsis invicta*. *Proc. Natl. Acad. Sci. U. S. A.* 108, 5679–5684.  
 641 doi:10.1073/pnas.1009690108  
 642  
 643  
 644  
 645

646

647

648 Figure legends

649 Figure 1. Sequence comparison of the prokineticin domain of arthropod astakines.

650 Highly conserved residues are colored green. Other residues are colored according to their

651 chemical properties. Ten cysteine residues, a proline residue between C<sub>79</sub> and C<sub>81</sub> and two

652 motifs (R<sub>22</sub>Y<sub>23</sub>S<sub>24</sub> and Y<sub>57</sub>P<sub>58</sub>) are conserved among most of the astakine sequences.

653

654 Figure 2. Deduced exon-intron structure of the astakine gene of a representative each from the

655 insects, chelicerates and crustaceans. All astakine genes investigated in the current study

656 contain the exon colored yellow in the figure. This exon ends with CCXX or CCXXX.

657

658 Figure 3a. A phylogenetic tree of 79 astakine protein sequences from arthropods analyzed

659 using the maximum likelihood method with the IQ-TREE software. The astakine-like

660 sequence from the hexapod *F. candida* was used as root sequence.

661 Bootstrap values are given at the nodes, and light grey or dark grey shading indicates clades

662 of closely related taxa. The scale bar indicates substitutions per site.

663 Figure 3b. A phylogenetic tree of astakine protein sequences from arthropods analyzed using

664 the MrBayes software. The astakine-like sequence from the hexapod *F. candida* was used as

665 root sequence.

666 Node support values are given at the nodes, and light grey or dark grey shading indicates

667 clades of closely related taxa. The scale bar indicates substitutions per site.

668

669 Figure 4a. A phylogenetic tree of 33 astakine protein sequences with highest similarity to *P.*  
670 *leniusculus* astakine 2 analyzed using the maximum likelihood method with the IQ-TREE  
671 software. The astakine-like sequence from the hexapod *F. candida* was used as root sequence.  
672 Bootstrap values are given at the nodes, and the scale bar indicates substitutions per site.  
673  
674 Figure 4b. A phylogenetic tree of 33 astakine protein sequences with highest similarity to *P.*  
675 *leniusculus* astakine 2 analyzed with MrBayes software. The astakine-like sequence from the  
676 hexapod *F. candida* was used as root sequence.  
677 Node support values are given at the nodes, and the scale bar indicates substitutions per site.  
678  
679 Figure 5a. Overall fold and the core of six crustacean astakines, determined by using Mamba  
680 intestinal toxin 1 as the template. The highly conserved cysteine residues and the RYS and  
681 YP(N/D) motifs are colored, and the indel regions are shown by yellow crosses.  
682  
683 Figure 5b. RYS and YP(N/D) structures, showing the charge potential mapped to the surface,  
684 of six crustacean astakines, determined by using Mamba intestinal toxin 1 as the template.  
685

[illegible]

INSECTA



*Megachile rotundata* (Leafcutter Bee)

MTPIFVTLFLLFVLSCSSRAQTNRPDYIQCQSNAECD SGYCCNI  
GPLRYSIPQCKVMAEGEICRPGSTSPTNMTLGYPDGALVTLTN  
VHYILCPCANGLTCDTKEGICKDTGEGHDTNRLFEEHKRHD

CHELICERATA



*Parasteatoda tepidariorum* (American House Spider)

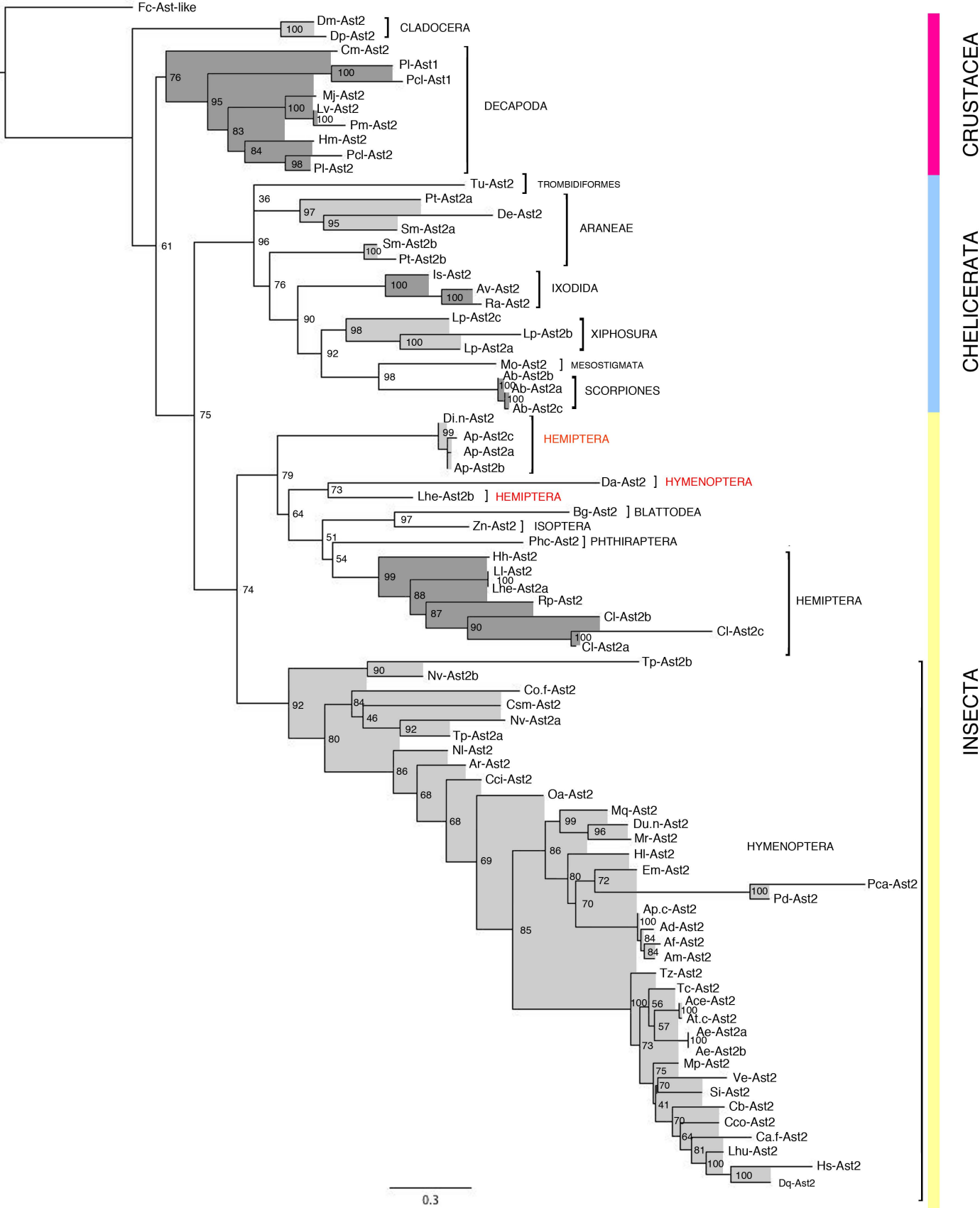
FTLSIVVSLLFQ  
VCICNTPRECSSKRDCGPNECCVV  
GRTRYSIPECKPNGRVGNTCLRGAESEDLTLYYPNGQRELEG  
VYTLFCPCDQNLVCKSNRCTV

CRUSTACEA

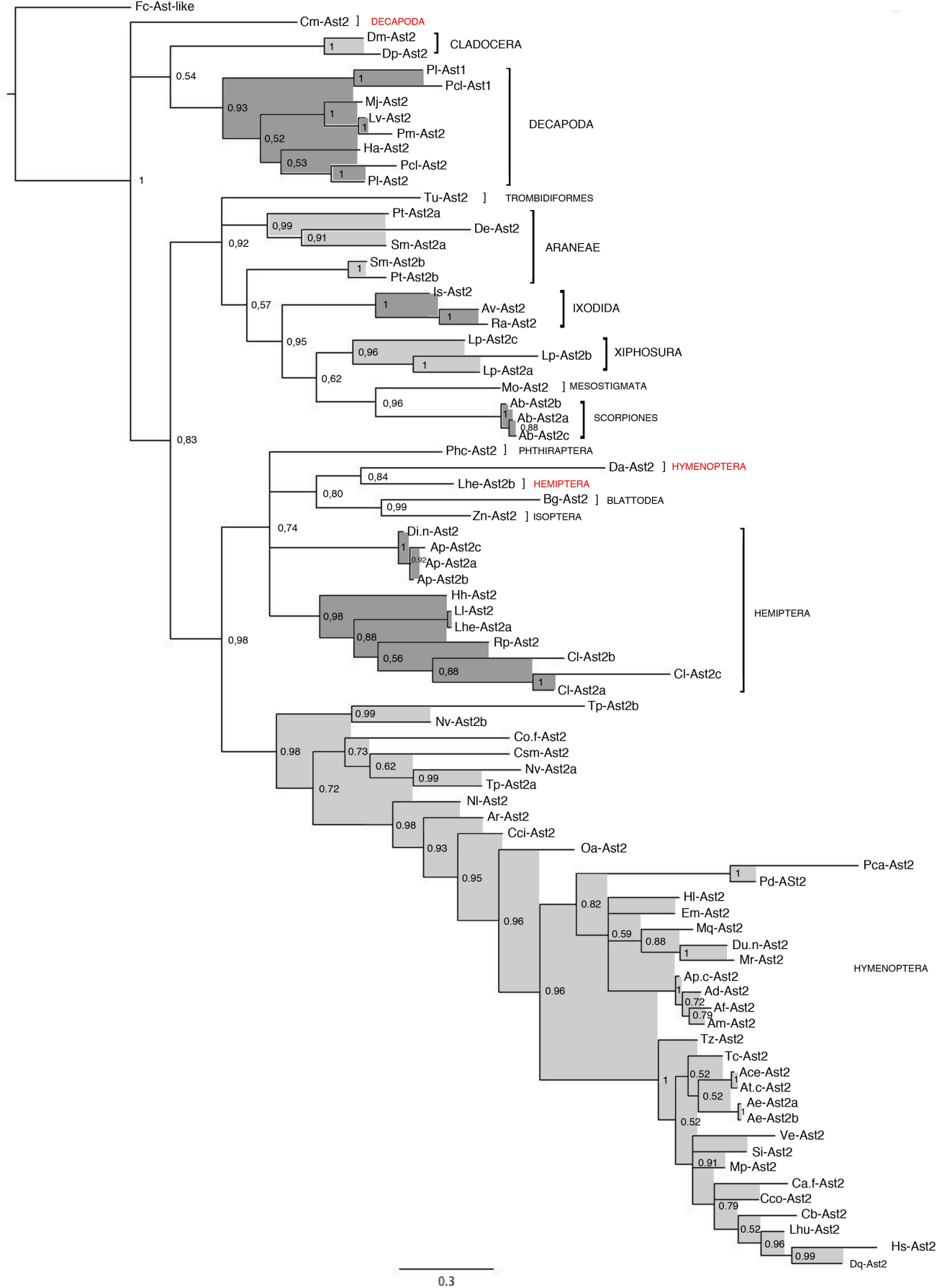


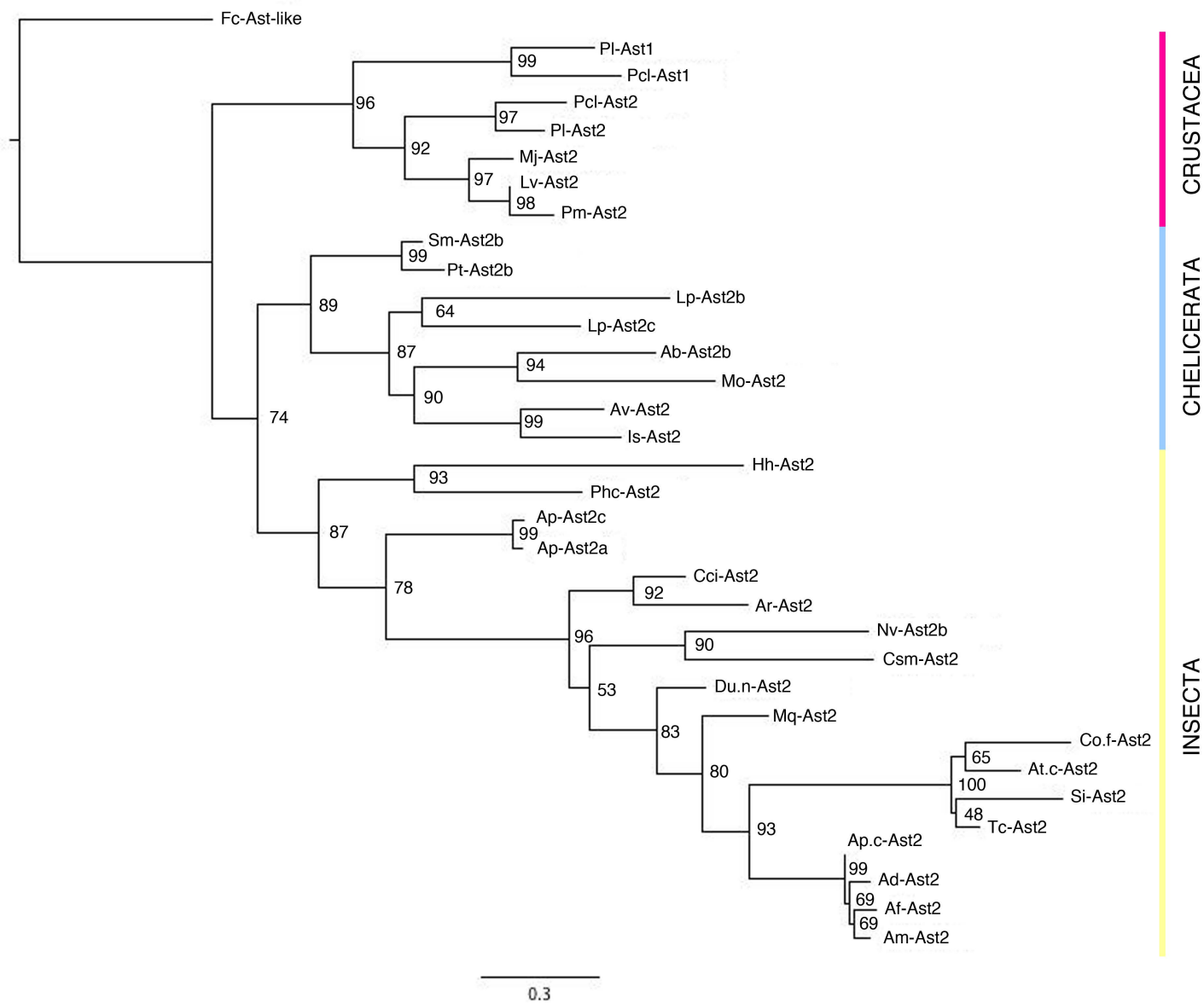
*Daphnia magna* (Water Flea)

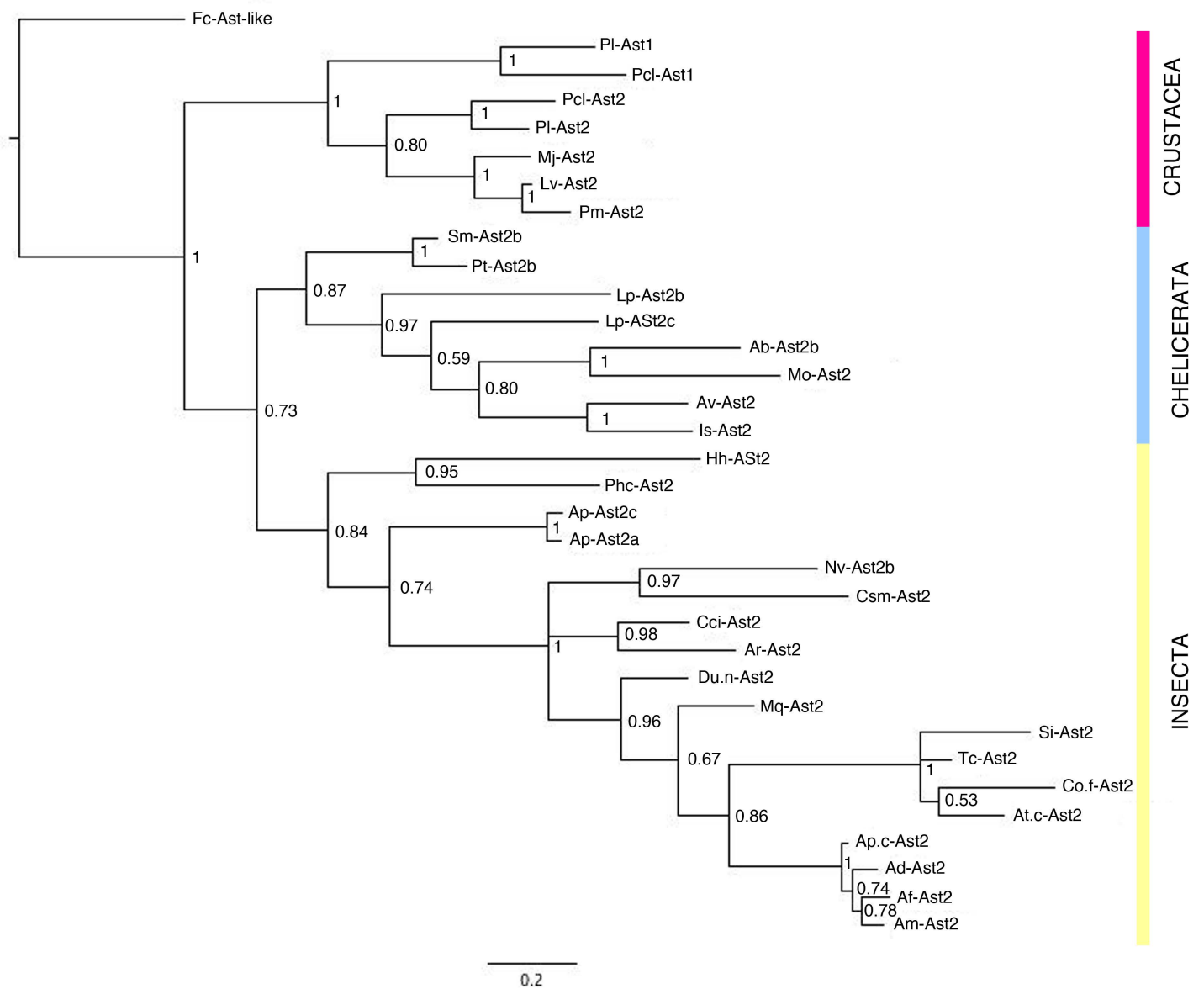
MLKECSLLFVCWTTLALTATLQPLPSYGVTGDCRSSEDCGPSSCCLL  
GMMRYSTPWCAPLLKLGDECRPSSHQLINRTLSYPGGLEIFLKDAHQV  
LCPCDANEGLVCSPLKGTVCYDVANDITPL

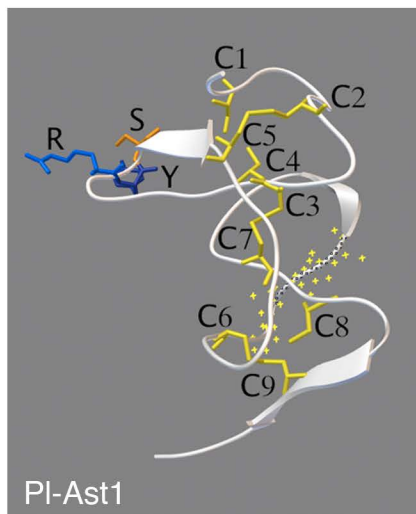
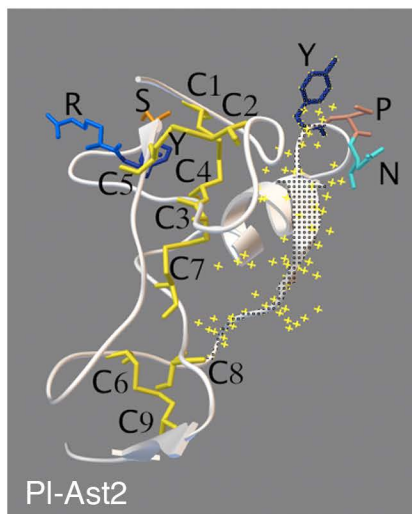
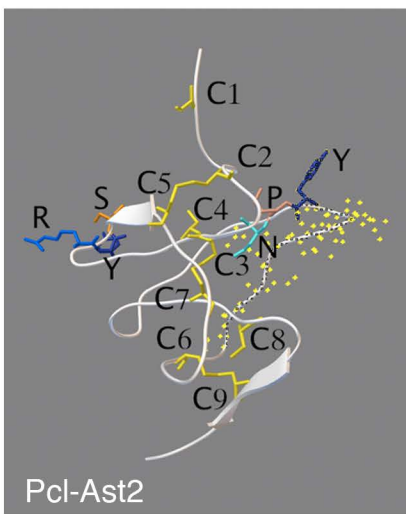
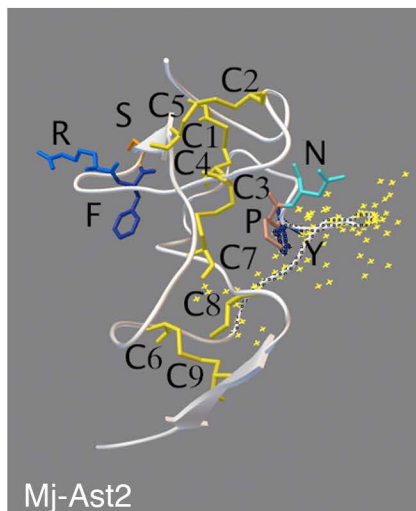
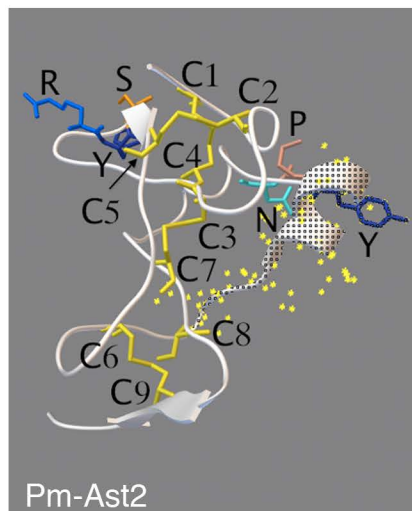
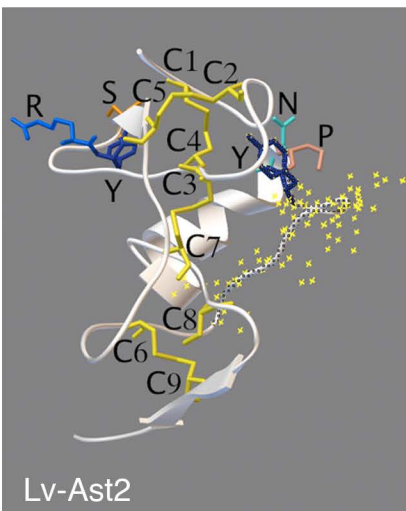


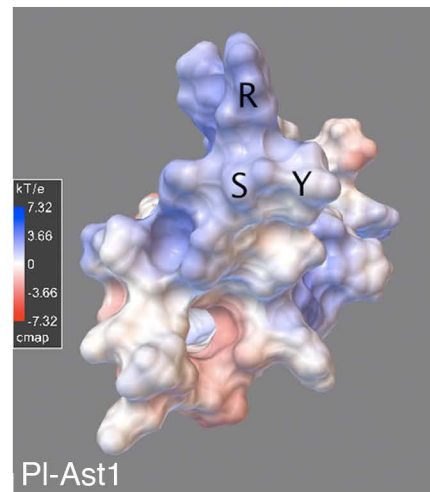
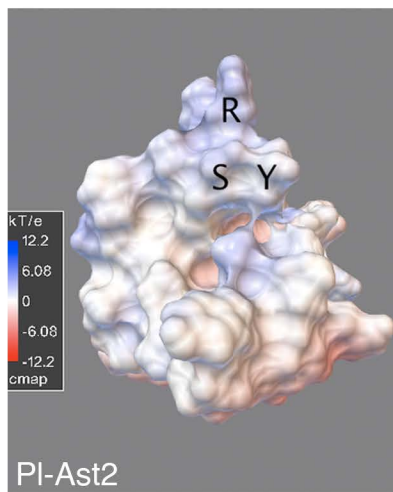
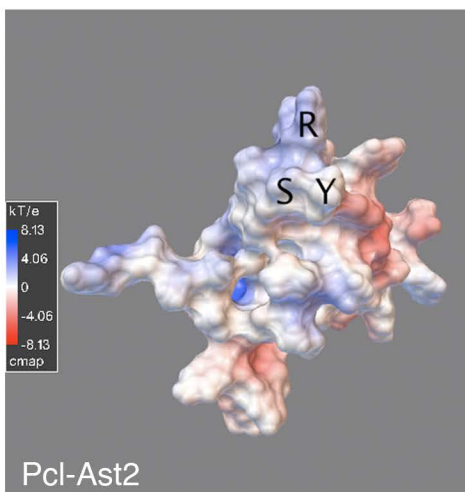
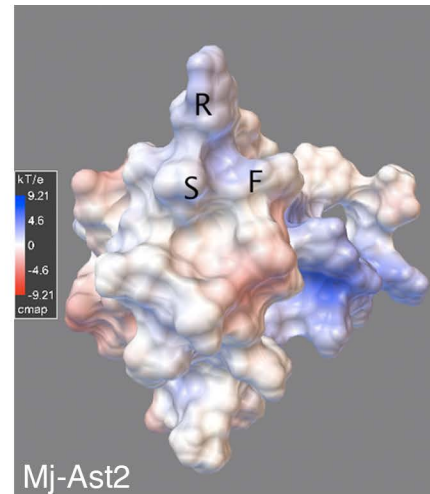
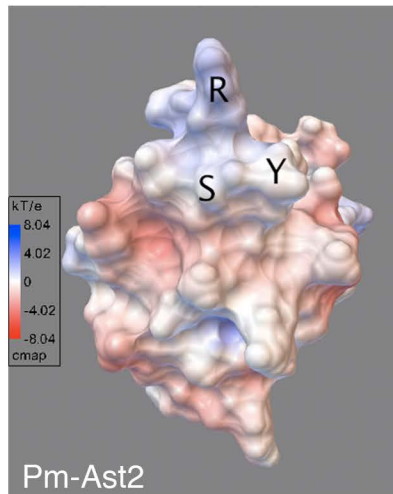
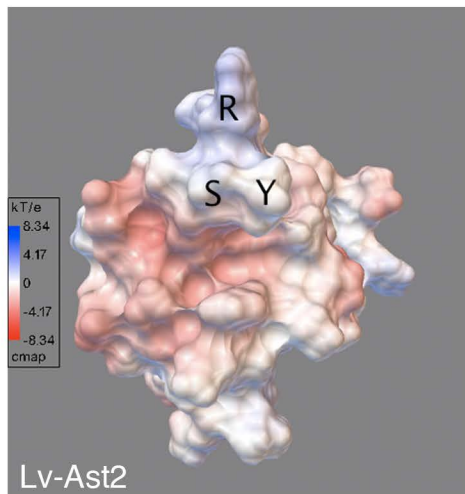




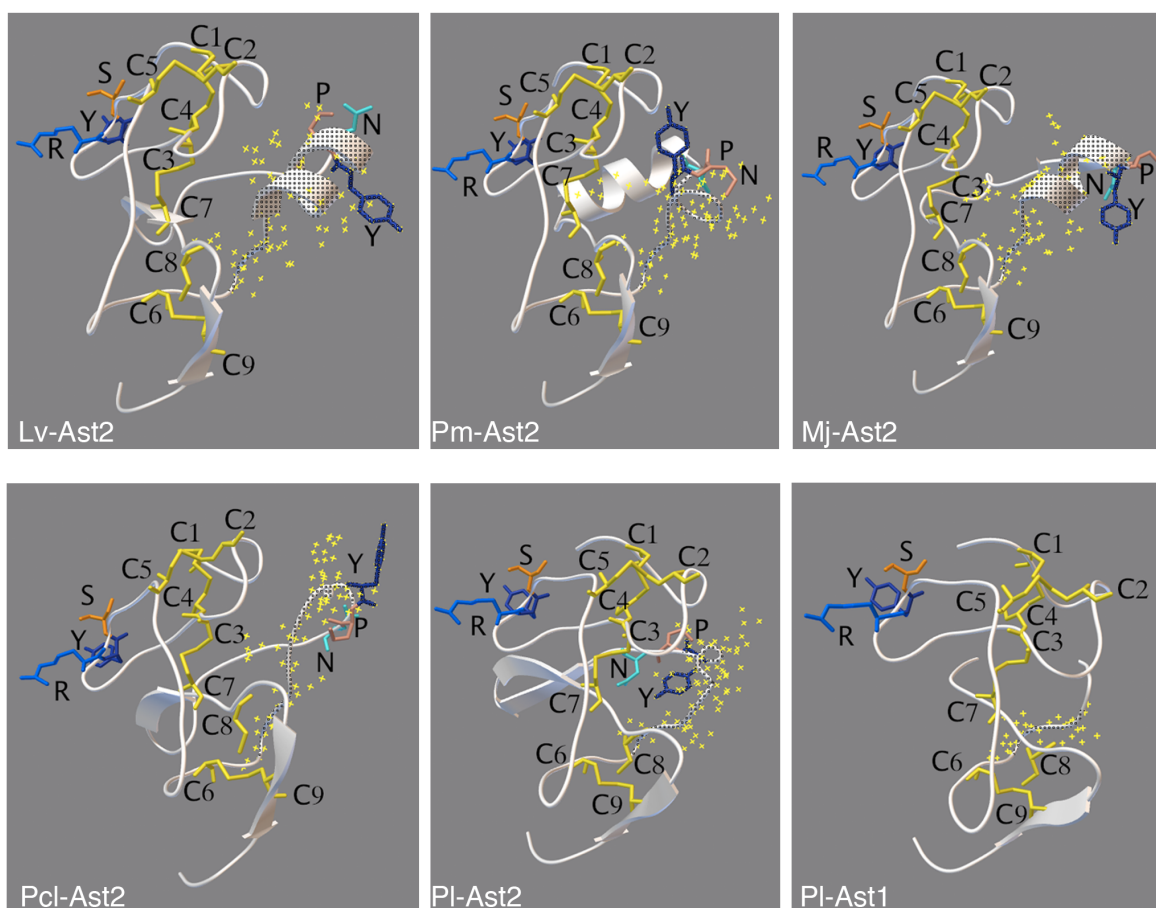






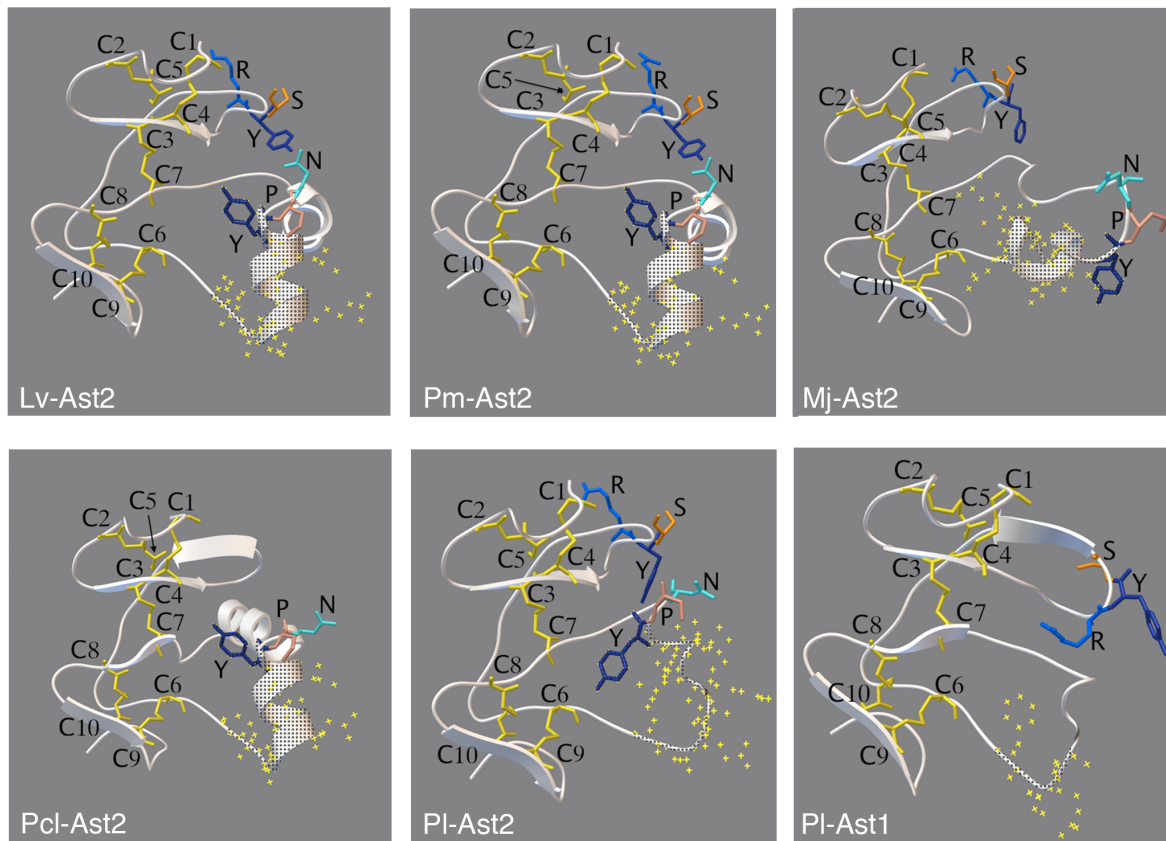






Supplementary Figure 2. The overall fold and the core of six crustacean astakines using the prokineticin Bv8 from *Bombina variegata* as model. The highly conserved cysteine residues and the RYS and YP(N/D) motifs are colored and the indels regions are shown by yellow crosses.





Supplementary Figure 3. The overall fold and the core of six crustacean astakines using Dickkopf-related protein 1 from *Homo sapiens* as model. The highly conserved cysteine residues and the RYS and YP(N/D) motifs are colored and the indels regions are shown by yellow crosses.



Supplementary Table 1. List of astakines included in the study.

Arthropods	Protein name	Species	pI	Mw	Accession number
<b>Crustacea</b>					
	Mj-Ast2	<i>Marsupenaeus japonicus</i>	4.88	11.28	<b>BAJ34645.1</b>
	Lv-Ast2	<i>Litopenaeus vannamei</i>	4.68	11.32	<b>ADM53424.1</b>
	Pm-Ast2	<i>Penaeus monodon</i>	5.13	11.30	<b>AAX14636.1</b>
	Pcl-Ast2	<i>Procambarus clarkii</i>	7.69	10.60	<b>AEC50077.1</b>
	Pl-Ast2	<i>Pacifastacus leniusculus</i>	7.04	11.19	<b>ABQ23255.1</b>
	Ha-Ast2	<i>Homarus americanus</i>	4.84	10.28	<i>FE535609</i>
	Pl-Ast1	<i>Pacifastacus leniusculus</i>	5.00	8.73	<b>AAX14635.1</b>
	Pcl-Ast1	<i>Procambarus clarkii</i>	5.08	8.15	*
	Cm-Ast2	<i>Carcinus maenas</i>	5.07	10.25	<i>DW585080</i>
	Dm-Ast2	<i>Daphnia magna</i>	4.85	11.48	<b>KZS05559.1,</b> <i>LRGB01002901.1</i>
	Dp-Ast2	<i>Daphnia pulex</i>	4.54	11.95	<i>FE329237</i>
<b>Chelicerata</b>					
	Pt-Ast2a	<i>Parasteatoda tepidariorum</i>	7.62	9.27	<b>XP_015917287.1,</b> <i>XM_016061801.1</i>
	De-Ast	<i>Dysdera erythrina</i>	5.31	10.18	<i>CV178181</i>
	Sm-Ast2a	<i>Stegodyphus mimosarum</i>	8.54	9.85	<b>KFM62184</b>
	Sm-Ast2b	<i>Stegodyphus mimosarum</i>	3.93	10.66	<b>KFM69031.1</b>
	Pt-Ast2b	<i>Parasteatoda tepidariorum</i>	3.92	11.02	<b>XP_015920703</b>
	Lp-Ast2b	<i>Limulus polyphemus</i>	6.07	10.20	<b>XP_013776682.1</b> <i>XM_013921228.1</i>
	Lp-Ast2a	<i>Limulus polyphemus</i>	6.70	10.40	<b>XP_013775587.1</b> <i>XM_013920133.1</i>
	Lp-Ast2c	<i>Limulus polyphemus</i>	5.09	9.61	<b>XP_013785672.1</b>
	Is-Ast2	<i>Ixodes scapularis</i>	4.60	11.29	<i>EW845057</i>
	Av-Ast2	<i>Amblyomma variegatum</i>	4.92	11.73	<b>DAA34752.1</b>
	Ra-Ast2	<i>Rhipicephalus appendiculatus</i>	5.06	11.62	<i>CD794853</i>
	Ab-Ast2b	<i>Androctonus bicolor</i>	5.31	8.74	<b>AIX87718.1</b>
	Ab-Ast2a	<i>Androctonus bicolor</i>	5.31	8.78	<b>AIX87717.1</b>
	Ab-Ast2c	<i>Androctonus bicolor</i>	5.31	8.81	<b>AIX87719.1</b>
	Tu-Ast2	<i>Tetranychus urticae</i>	7.66	11.46	<b>XP_015792820.1,</b> <i>XM_015937334.1</i>
	Mo-Ast2	<i>Metaselulus occidentalis</i>	4.59	13.76	<b>XP_003743670,</b> <i>XM_003743622.1</i>
<b>Insecta</b>	Da-Ast2	<i>Diachasma alloeum</i>	6.94	10.98	<b>XP_015127605.1</b>
	Hs-Ast2	<i>Harpegnathos saltator</i>	5.15	9.29	<b>EFN86043.1,</b> <i>GL447755.1</i>
	Ve-Ast2	<i>Vollenhovia emeryi</i>	5.81	14.95	<b>XP_011861926.1,</b>

					<i>XM 012006536.1</i>
	Ca.f-Ast2	<i>Camponotus floridanus</i>	7.63	9.63	<b>XP_011266500.1,</b> <i>XM 011268198.1</i>
	Dq-Ast2	<i>Dinoponera quadriceps</i>	5.11	9.19	<b>XP_014486827.1</b>
	Lhu-Ast2	<i>Linepithema humile</i>	9.07	16.07	<b>XP_012228085.1,</b> <i>XM 012372662.1</i>
	Cco-Ast2	<i>Cyphomyrmex costatus</i>	4.60	9.38	<b>KYN05408.1</b>
	Cb-Ast2	<i>Cerapachys biroi</i>	6.05	9.69	<b>EZA54888.1, KK107235.1</b>
	Si-Ast2	<i>Solenopsis invicta</i>	4.53	14.80	<b>EFZ14050.1, GL767121.1</b>
	Mp-Ast2	<i>Monomorium pharaonis</i>	7.43	9.84	<b>XP_012539106.1,</b> <i>XM 012683652.1</i>
	Tz-Ast2	<i>Trachymyrmex zeteki</i>	4.49	12.59	<b>KYQ52101.1</b>
	Ace-Ast2	<i>Atta cephalotes</i>	9.27	17.89	<b>XP_012063524.1,</b> <i>XM 012208134.1</i>
	At.c-Ast2	<i>Atta colombica</i>	8.57	9.92	<b>KYM75707.1</b>
	Tc-Ast2	<i>Trachymyrmex cornetzi</i>	6.48	9.71	<b>KYN16810.1</b>
	Ae-Ast2a	<i>Acromyrmex echinator</i>	5.49	9.87	<b>EGI67870.1, GL888084.1</b>
	Ae-Ast2b	<i>Acromyrmex echinator</i>	8.67	17.00	<b>XP_011050611.1</b>
	Cl-Ast2b	<i>Cimex lectularius</i>	8.58	10.75	<b>XP_014256495.1</b>
	Cl-Ast2c	<i>Cimex lectularius</i>	7.50	10.08	<b>XP_014256561</b>
	Cl-Ast2a	<i>Cimex lectularius</i>	4.94	11.47	<b>XP_014256494.1</b>
	Ll-Ast2a	<i>Lygus lineolaris</i>	4.19	12.39	<b>AJR27902.1</b>
	Lhe-Ast2a	<i>Lygus hesperus</i>	4.19	12.39	<b>JAG12052.1</b>
	Hh-Ast2	<i>Halyomorpha halys</i>	5.53	12.79	<b>XP_014294052.1,</b> <i>XM 014438566.1</i>
	Rp-Ast2	<i>Rhodnius prolixus</i>	4.33	11.80	<b>JAA75272.1</b>
	Phc-Ast2	<i>Pediculus humanus corporis</i>	4.27	10.73	<b>XP_002431167.1,</b> <i>XM 002431122.1</i>
	Lhe-Ast2b	<i>Lygus hesperus</i>	5.25	11.95	<b>JAG7549.1</b>
	Bg-Ast2	<i>Blatta germanica</i>	4.46	9.41	<b>FG125716</b>
	Zn-Ast2	<i>Zootermopsis nevadensis</i>	5.13	11.26	<b>KDR15950.1</b>
	Ap-Ast2c	<i>Acyrtosiphon pisum</i>	4.24	10.15	<i>EX635976.1</i>
	Ap-Ast2a	<i>Acyrtosiphon pisum</i>	4.16	10.18	<i>FF336283.1</i>
	Ap-Ast2b	<i>Acyrtosiphon pisum</i>	4.02	10.02	<i>EX610914.1</i>
	Di.n-Ast2	<i>Diuraphis noxia</i>	4.29	12.41	<b>XP_015368848.1</b>
	Tp-Ast2b	<i>Trichogramma pretiosum</i>	6.63	10.87	<b>XP_014230684.1,</b> <i>XM 014375198.1</i>
	Nv-Ast2b	<i>Nasonia vitripennis</i>	5.16	11.02	<b>XP_001605660.1,</b> <i>XM 001605610.3</i>
	Csm-Ast2	<i>Ceratosolen solmsi marchali</i>	5.02	11.94	<b>XP_011499199.1,</b> <i>XM 011500897.1</i>
	Co.f-Ast2	<i>Copidosoma floridanum</i>	7.63	9.63	<b>XP_011266500.1,</b> <i>XM 011268198.1</i>
	Nv-Ast2a	<i>Nasonia vitripennis</i>	5.07	13.41	<b>XP_008213060.1,</b> <i>XM 008214838</i>
	Tp-Ast2a	<i>Trichogramma pretiosum</i>	4.29	12.59	<b>XP_014235233.1,</b> <i>XM 014379747.1</i>
	Af-Ast2	<i>Apis florea</i>	7.64	9.74	<b>XP_012345115.1,</b> <i>XM 012489661.1</i>
	Ad-Ast2	<i>Apis dorsata</i>	6.86	9.65	<b>XP_006617146.1,</b> <i>XM 006617083.1</i>
	Ap.c-Ast2	<i>Apis cerana</i>	5.56	9.71	<b>XP_016909725.1,</b>

					<i>XM_017054236.1</i>
	Am-Ast2	<i>Apis mellifera</i>	6.71	9.84	<b>XP_003250271.1,</b> <i>XM_003250223</i>
	Em-Ast2	<i>Eufriesea mexicana</i>	4.72	10.50	<b>OAD53282.1</b>
	Hl-Ast2	<i>Habropoda laboriosa</i>	5.10	11.85	<b>KOC63709.1</b>
	Mq-Ast2	<i>Melipona quadrifasciata</i>	7.99	14.09	<b>KOX80272.1</b>
	Du.n-Ast2	<i>Dufourea novaeangliae</i>	4.87	11.97	<b>KZC05910.1</b>
	Mr-Ast2	<i>Megachile rotundata</i>	5.19	12.02	<b>XP_003708605.1,</b> <i>XM_003708557.1</i>
	Oa-Ast2	<i>Orussus abietinus</i>	6.50	17.02	<b>XP_012283423.1,</b> <i>XM_012428000.1</i>
	Cci-Ast2	<i>Cephus cinctus</i>	4.84	14.50	<b>XP_015602971.1</b>
	Ar-Ast2	<i>Athalia rosae</i>	5.21	12.19	<b>XP_012261298.1,</b> <i>XM_012405875.1</i>
	Nl-Ast2	<i>Neodiprion lecontei</i>	4.86	12.51	<b>XP_015520701.1</b>
	Pca-Ast2	<i>Polistes canadensis</i>	6.98	10.71	<b>XP_014616409.1,</b> <i>XM_014760923.1</i>
	Pd-Ast2	<i>Polistes dominula</i>	5.18	10.56	<b>XP_015184135.1,</b> <i>XM_015328649.1</i>
Root sequence	Fc-Ast-like	<i>Folsomia candida</i>	5.08	12.81	<b>XP_021968286.1</b>

Supplementary Table 2. Arthropod astakine gene structures.  
([https://www.ncbi.nlm.nih.gov/genome/annotation\\_euk/all/](https://www.ncbi.nlm.nih.gov/genome/annotation_euk/all/))

1) Insecta

*Harpegnathos saltator* (taxonomy ID: 610380) Hs-Ast2

MSSILSILLIITVGLVFSSNGQCTNNADCLSDECCLL

Intron – 71 bp –

GPMRYSTPTCIPYQKKGDQCRVNAEFVTTNLTYPNNSHLEVKNVSYILCPCVKETSCNKETGICD

*Camponotus floridanus* (taxonomy ID: 104421) Ca.f-Ast2

MSLMSNVLLLITLAGIVPAFPFNSFKNCKTDLECSSLNLCCLL

Intron – 225 bp –

LLGPTRYAIPTCMPFQQKGEOQCRVNADTITANLTYPNNLQLEIRNANFILCPCANGLFCERGICN

*Cerapachys biro* (Taxonomy ID: 443821) Cb-Ast2

MSSILGLLLLISIAVAVPTSRTQQCVTNSDCPSNHCCLL

Intron – 183 bp –

GPSRYATPACMPFQQRGEQCRVNADTISTNLTYPDDSRIEVESIHYILCPCADGLSCNFKKGICN

*Monomorium pharaonis* (taxonomy ID: 307658) Mp-Ast2

MSPISGILIFVISIVATSNIGSVTSSSQDCATNSECKSNSCCLL

Intron – 236 bp –

GPSRYAIPTCMPFQQKGEOQCRVNAKTITTTTLFYPDGSQVEVKDIHSILCPCADGLSCDPKRGICK

*Nasonia vitripennis* (taxonomy ID: 7425) Nv-Ast2a

MKMIMRLGLLLLCAMVINTKALARFPRNWNSHIECTNSLQCAPGHCCTI

Intron – 128 bp –

STERYSYPRCQKLHEVGDYCRAEGPLLTNGNMTYPDGSKPSDVHLEDVYLLFCPCAPGLVCDSDERI  
CRQPDMKDFNYLKEQETGSNKSDD

*Nasonia vitripennis* (taxonomy ID: 7425) Nv-Ast2b

MTSAVLLLSLMIGSLYAAAEQSIIPPSWVECTSHLDCRPGSCCTI

Intron – 238 bp –

GQQRYSIPMCSPQPTLGEQCRPNPRLTNTTLGYPDGSTILIKDAYLMLCPCSTGLSCDSRIGLCQV  
KQQA EA

*Trichogramma pretiosum* (taxonomy ID: 7493) Tp-Ast2a

MLAARIKFSLSLLLLVGAVNAATLHEQTGLLADEIECTSNVDCAPGYCCTI

Intron – 86 bp –

SHERYSYPRCQKFQDLGDFCRPGGPLTTSGDRYYPDGTSIQLDEIYMQFCPCGPGLLCDRGEQVCRD  
ASDFNSVQLNQSGKSDD

*Trichogramma pretiosum* (taxonomy ID: 7493) Tp-Ast2b

MCRLFSFLFLSFATIIHLIDANDYRSGVIECINHSDCGPDSCCTIS

Intron — 523 bp —

MDRYSKPRCSKRPLEGEFCHPYLHKIENHNFWYPHNNNIFVEEAHYLSCPCFSGLRCDVEKAICKSK  
IA

*Ceratosolen solmsi marchali* (taxonomy ID: 326594) Csm-Ast2

MSRLLLLLLLLFHITAILIRDVRAGHPSWIHCTSNLHCAPGYCCRM

Intron — 1686 bp —

GFQRYISIPSCPEVLKDGEPCKPGEPFITNGTRGYPDGTTIELEDVYVMFCPCAIGLACDREAFVCRD  
ASEMKDFNHLSGKSDKTDD

*Apis mellifera* (taxonomy ID: 7460) Am-Ast2

MMTSIFEISFLFLVLAYPECHAQNDYIKQTSSECQPNHCCTL

Intron — 96 bp —

LGSVRYISIPQCKPMQKGGEVCRPTNSSTTFNVTLGYPDGSSLKIEDVYFIFCPCIDGLSCEKGICKE  
KN

*Apis dorsata* (taxonomy ID: 7462) Ad-Ast2

MMTSIFVISFLFLVLAYPECHAQNDYIKQTSSECQPNHCCTL

Intron — 110 bp —

LGSARYSIPQCKPMQKGGEVCRPTNEPFNVTLGYPDGSLKIEDVHFIFCPCIDGLSCEKGICKEKN

*Apis florea* (taxonomy ID: 7463) Af-Ast2

MMTSIFVISFLFLVLAYPECHTQNDYIKQTSSECRPNHCCTL

Intron — 95 bp —

GPIRYSIPQCKPMQKGGEVCRPTNVTNVTLGYPDGSLKIEDVHFIFCPCIDGFSCEKGVCKEKN

*Megachile rotundata* (taxonomy ID: 143995) Mr-Ast2

MTPIFVTLFLLFVLSCSSRAQTNRPDYIQQSNAECDSGYCCNI

Intron — 95 bp —

GPLRYSIPQCKVMQAEGEICRPGSTSPNTMTLGYPDGLVTLTNVHYILCPCANGLTCDTKEGICKD  
TGEHDTNRLFEEHKRHD

*Acyrtosiphon pisum* (taxonomy ID: 7029) Ap-Ast2c

MNTNIMGKLSLVAVTILVATVAAYPSPKPSFLGCQSSQDCGMNECCVLG

Intron — ? bp —

GMMRYSVPTCRPLGEEGDTICIPNSGDVQPQNVTVTPDGSSADLYVHTMLCPCVSGLECSGDMSCGTG  
LNGAGKLMPAPRLHGRG

*Orussus abietinus* (taxonomy ID: 222816) Oa-Ast2

MTQKLILAMILAIAGMSGPIIAGVASNRPSHVQCVSNSECLRGSCCTI

Intron — 946 bp —

APYKFSVPQCQSMQEEGAQCRPMGHETINTTLTPDGSLELKGVHYILCPCDYGLTCDPKDGICRD  
VSQRRDFNHLQNEAIAHED

*Athalia rosae* (taxonomy ID: 37344) Ar-Ast2

MTQSLTFAAIVVIVGICQVGYTSARVTMRPPYIQCSNSECLPGNCCSI

Intron — 69 bp —

GQNRFSIPQCKPMQDQGGVCRPRGPMTSNTTLVYPDGSQVQLVEVHIGFCPCGYGLTCNP EEGLCRD  
PSQRRGFNSLLDEASVQDD

*Polistes canadensis* (taxonomy ID: 91411) Pca-Ast2

MSSTISTSFTILLLLGLVSFLFAASIKQEPPDVQCRNDKECPDDHCCVI

Intron — 81 bp —

INGGRYVIPQCRPLLKKTETCKGDDRLFNNTLYPNDDKLTISGVHVFVLCPCHEGLICGLKEKVCIS  
NN

*Polistes dominula* (taxonomy ID: 743375) Pd-Ast2

MSSTISTSFTILLLLGLVSFLFAASIKQEPPANVQCHNNKECPSDHCCVL

Intron — 68 bp —

GGGRYTIPQCSPLLEEAATCRPNNELNMTLHYPNDTQLKISDVYHILCPCNEGLICDRKEGVCINN  
N

*Halyomorpha halys* (taxonomy ID: 2867706) Hh-Ast2

MSMTLFQLGAIASIFLTVYAMPNDRPGYIDCLDSSECGRDKCCSI

Intron — 2710 bp —

CSIGMGYSIPMCYAKGNIGDKCIPNNTLQKMTSLSYPDGTSINLTNFYFYHACPLDNLICSKDTE  
TCEDPLF

Intron — 1105 bp —

NYDFRGRYYQHTMGRF

*Atta Cephalotes* (taxonomy ID: 12957) Ace-Ast2

MCQNAERTQRYKTMRSSVRAIIRRAXCTASSGPRFKTT

Intron — 1108 bp —

EIIKSKQAIMPLMLNVVILITSIVVFPNIDPVTSSSPVFQKNCTTNTECEPNSSCCLLG

Intron — 1242 bp —

LGPMRYSIPTCMPFRQKGELCRVNAETITTNLTYPNTLEIKVKDIHYILCPCADGLSCNPKRGICK

*Linepithema humile* (taxonomy ID: 83485) L.hu-Ast2

MKWRYKMMRLSLRTVIRRAVCDKCGSRIALQL

Intron — 518 bp —

QLVKSKQTAIMSPILVALLFISLAIAAPPLIPSEQCTTDSECP SDFCCLL

Intron — 139 bp —

GPSRYAMPACMPYQQKGEOQCRVNAKTITTNLTYPDNSQLEVKNINFILCSCADGLSCNKKTGICN

*Solenopsis invicta* (taxonomy ID: 13686) Si-Ast2

MRSSVGAIIRRAXCAANS GPRI RND

Intron — 1049 bp —

RDYQTQTAIMTPISGILILVISIMATSSISSVPSYQKCN TNEDCKSSSCCLL

Intron — 189 bp —

GPSRYALPSCMPYQQKGEOQCRMNADTITTNLSYPDNSQIEVKDIHLILCPCADGLSCDFGICEEDA

*Vollenhovia emeryi* (taxonomy ID: 411798) Ve-Ast2

MRSSVRAIIHASCVAAGSGPRARND

Intron — 1225 bp —

SGYQAQTAIMSPMPGVLLFISIVTLPLNISSIPLSSENCVTNSECQTDSCCVL

Intron — 254 bp —

GASRYVIPTCMPFQQIGETCRVNAATITTNLSYPDNSQLEVTAVHFILCPAAGLSCDSKHGTCE

## 2) Chelicerate

### *Parasteatoda tepidariorum* (taxonomy ID: 114398) Pt-Ast2a

FTLSIVVSLLFQ

Intron — 1913 bp —

VCICNTPRECSSKRDCGPNECCVVG

Intron — 3377 bp —

GRTRYSIPECKPNRVRGNTCLRGAESDLTLYYPNGQRELEGVYTLFCPCDQNLVCKSNRCTV

### *Parasteatoda tepidariorum* (taxonomy ID: 114398) Pt-Ast2b

MGTPVHMYAFLAAMLVCCFSQ

Intron — 6785 bp —

QVSSYTLATSECRSQADCGPGECCVLG

Intron — 7829 bp —

MMRY SMAQCMLPGQVEDYCRDDNPPENRTLYYPNGEPVEVYEIYTHVCPDES LQCTDNFCAMDESY  
ENNYLY

### *Tetranychus urticae* (taxonomy ID: 32264) Tu-Ast2

MLCYSTKFIIIFALM

Intron — 549 bp —

MVTVSGRTWNFFALNSPKPCQSSDDCRRGECCAIG

Intron — 240 bp —

GFARFSVPMCKPMGRINDWCYPDNEPENMTLHYPYGSEAYTNVHRNFCPCCKQPLTCEHNICKFERFN  
YY

### *Limulus polyphemus* (taxonomy ID: 6850) Lp-Ast2a

MRTLVAIIILVAQ

Intron — 4625 bp —

MAQSFPGFRCRSQQDCDPGSCCVV A?

Intron — 429 bp —

MERFSTPRCQKLSQQGEYCRPRNSALNTSLSYPNGILDVTNLYTVLCPCDVGLICEQAMCQPNTFLQ  
SNHLA

### *Limulus polyphemus* (taxonomy ID: 6850) Lp-Ast2b

MKTIVCVFLILLELQ

Intron — ? bp —

VILSFPGFDGCRSPQDCDKSSCCVI T

Intron — 429 bp —

MEKYSVPHCRKLGNKEEYCRTNSAQNMTLNYPNGSVDVFGVYRILCPCNDGLECVQSVCQLLHSDT  
IL

*Limulus polyphemus* (taxonomy ID: 6850) Lp-Ast2c

MIMRPEITLFIIFTIIM

Intron — 2450 bp —

TIIMLAIGVPYFYGCKSPADCEPGECCVIG

Intron — 921 bp —

MNRYSFPRCEKFGQKNDCLPSNTPQNKTLYYPNGAVDFSNIYMLFCPCDTGFICYQAHCESA

*Metaseiulus occidentalis* (taxonomy ID: ?) Mo-Ast2

MRSSMERLSLLLLISLSVFLEFAACEDVEVRSCRKPSDCDPGYCCRI

Intron — 1268 bp —

GMRFSQPFQKFGTVGDTCRMGAEPEDKILWFPGGLTFDVGVRQFCPCEG

Intron — 153 bp —

GGLACKEAMCQPESAKIAAPTKKYNIDDFDYESLDNRAKSDNFAEFDI

### 3) Crustacea

*Daphnia pulex* (taxonom ID: 6669) Dp-Ast2

MKECGLLFVCWATVVLGILQPLPSHSAMGDCRSNEDCGPNRCCLLG

Intron — 74 bp —

GMMRYSTPWCAPLLNLGEDCRPTSSNEPSITNRTLVPYGGLEIFLKDAYQ

Intron — 75 bp —

ILCPCDANQGLVCSHLSGACISDESNDISPL

*Daphnia magna* (taxonomy ID: 35525) Dm-Ast2

MLKECSLLFVCWTTLALTATLQPLPSYGVGTGDCRSSEDCGPSSCCLL

Intron — 74 bp —

GMMRYSTPWCAPLLKLGDECRPSSHQLINRTLSYPGGLEIFLKDAHQV

Intron — 74 bp —

LCPCDANEGLVCSPLKGTCVYDVANDITPL



Supplementary Table 3. Prediction of proportion of disorder in Phyre2 for the five sequences studied.

Astakine sequence	Proportion of disorder (%)
Pl-Ast1	49
Pl-Ast2	61
Pcl-Ast2	50
Mj-Ast2	62
Pm-Ast2	64
Lv-Ast2	66