

# **Phylogenomics of vertebrate serpins**

## **Dissertation**

zur Erlangung des akademischen Grades Doktor der  
Naturwissenschaften  
(Dr. rer. nat.)

vorgelegt von

**Abhishek Kumar, Master of Science**

geb. in Muzaffarpur (Indien)

Technische Fakultät der Universität Bielefeld  
2010

*Nothing in Biology Makes Sense  
Except in the Light of Evolution*

Theodosius Dobzhansky (1900-1975)

## List of publications and presentations from this work

### **Publications:**

1. **Kumar, A. & Ragg, H.** (2008). Ancestry and evolution of a secretory pathway serpin. *BMC Evolutionary Biology*, 8:250.
2. Ragg, H., **Kumar, A.** Köster, K. Bentele, C. Wang, Y. Frese, M.A. Prib, N. & Krüger, O. (2009). Multiple gains of spliceosomal introns in a superfamily of vertebrate protease inhibitor genes. *BMC Evolutionary Biology*, 9:208.

### **Poster:**

**Kumar, A. & Ragg, H.** (2008). Deciphering the phylogenetic history of neuroserpin orthologs across metazoans by analysis of synteny and rare genomic characters. German Conference on Bioinformatics (GCB) September 2008, Dresden, Germany.

### **Talk:**

**Kumar, A.** (2005). Delving into vertebrate serpins for understanding their evolution. Bioinformatics Research and Education Workshop (BREW), April 2005, Berlin. Germany.

# Contents

<b>1. Summary</b>	<b>1</b>
<b>2. Introduction</b>	<b>3</b>
<b>2.1 Overview of serpins</b>	<b>3</b>
2.1.1 Structure and mechanism of action	3
2.1.2 Physiological functions	6
2.1.3 Evolution of serpins	7
<b>2.2 Rare genomic changes</b>	<b>9</b>
2.2.1 Intron gain and loss	9
2.2.2 Rare indels	11
2.2.3 Gene duplications and fates of duplicated genes	11
2.2.4 Exonization of non-coding regions in genomes	13
<b>2.3 Intron evolution theories</b>	<b>15</b>
<b>2.4 Aim of this work</b>	<b>17</b>
<b>3. Materials</b>	<b>18</b>
<b>3.1 Genomes</b>	<b>18</b>
<b>3.2 Databases</b>	<b>18</b>
3.2.1 NCBI	18
3.2.2 RefSeq	19
3.2.3 Entrez	19
3.2.4 SWISS-PROT	19
3.2.5 UniProt	20
3.2.6 PROSITE	20
3.2.7 ENSEMBL	20
3.2.8 Serpin Databank	21
<b>3.3 Searching tools</b>	<b>21</b>
3.3.1 BLAST	21
3.3.2 PSI-BLAST	22
3.3.3 FASTA	22
3.3.4 Superfamily HMM library	23
<b>3.4 Multiple sequence analysis tools</b>	<b>23</b>
3.4.1 CLUSTAL	25
3.4.2 DIALIGN	26
3.4.3 MUSCLE	26
3.4.4 T-COFFEE	27
<b>3.5 Sequence editing tools</b>	<b>28</b>
<b>3.6 Phylogenetic tools</b>	<b>28</b>
3.6.1 MEGA 3.1	28
<b>3.6.2 PHYLIP</b>	<b>28</b>
<b>3.6.3 Phylodraw</b>	<b>29</b>



---

3.6.4	Geneious	29
3.6.5	Phylowin	29
3.6.6	TREEVIEW	29
3.6.7	NJPLOT	29
3.7	<b>Comparative genomics tools</b>	<b>30</b>
3.7.1	Genome browsing tools	30
3.7.2	GENELIGHT	31
4.	<b>Methods</b>	<b>33</b>
4.1	<b>Database searching and evaluations</b>	<b>34</b>
4.1.1	Searches with BLAST	34
4.1.2	Searches based on motifs	34
4.2	<b>Sequence alignment</b>	<b>35</b>
4.3	<b>Gene structure analysis</b>	<b>35</b>
4.3.1	Gene structure prediction	35
4.3.2	Mapping of intron positions	35
4.4	<b>Gene specific features</b>	<b>36</b>
4.5	<b>Orthology assignment</b>	<b>36</b>
4.5.1	Sequence identity and sequence similarity values from protein alignment	36
4.5.2	Rare indels	36
4.6	<b>Synteny analysis of vertebrate serpins</b>	<b>36</b>
4.6.1	Synteny analysis of group V1 serpins	36
4.6.2	Group V2 serpin synteny analysis	37
4.6.3	Synteny analysis of serpin groups V3-V6	37
4.7	<b>Analysis of the <i>Ciona intestinalis</i> genome</b>	<b>38</b>
4.7.1	Searching serpins in the <i>Ciona intestinalis</i> genome	38
4.7.2	Determination of the exon-intron structure	38
4.7.3	Synteny Building	38
4.7.4	Analysis of serine codon dichotomy at position 56	38
4.8	<b>Analysis of the <i>Branchiostoma floridae</i> genome</b>	<b>38</b>
4.9	<b>Analysis of the <i>Strongylocentrotus purpuratus</i> genome</b>	<b>38</b>
4.10	<b>Analysis of the <i>Nematostella vectensis</i> genome</b>	<b>39</b>
4.11	<b>Phylogenetic analysis and bootstrap analysis</b>	<b>39</b>
5.	<b>Results</b>	<b>41</b>
5.1	<b><i>Gallus gallus</i> and its serpins</b>	<b>42</b>
5.2	<b><i>Xenopus tropicalis</i> and its serpins</b>	<b>42</b>
5.3	<b><i>Danio rerio</i> and its serpins</b>	<b>43</b>
5.4	<b><i>Tetraodon nigroviridis</i> and its serpins</b>	<b>44</b>
5.5	<b><i>Fugu rubripes</i> and its serpins</b>	<b>45</b>
5.6	<b><i>Petromyzon marinus</i> and its serpins</b>	<b>47</b>
5.7	<b>Characterization of serpins from <i>Ciona intestinalis</i> genome</b>	<b>48</b>
5.8	<b><i>Branchiostoma floridae</i> and its serpins</b>	<b>53</b>
5.9	<b>The Sea urchin - <i>Strongylocentrotus purpuratus</i> and its serpins</b>	<b>53</b>

---

<b>5.10</b>	<b><i>Nematostella vectensis</i> and its serpins</b>	<b>55</b>
<b>5.11</b>	<b>Orthology analysis of group V1 serpins</b>	<b>55</b>
5.11.1	Gene structure of group V1 serpins	56
5.11.2	Synteny analysis of group V1 serpins	57
5.11.3	Sequence analysis of group V1 serpins	60
<b>5.12</b>	<b>Orthology analysis of group V2 serpin genes</b>	<b>63</b>
5.12.1	Gene structure of group V2 serpin genes	63
5.12.2	Synteny analysis of group V2 serpins in the $\alpha_1$ -antitrypsin cluster	65
5.12.3	Synteny analysis of the serpinA7 gene	68
5.12.4	Synteny analysis of the angiotensinogen (AGT) gene	68
5.12.5	Synteny analysis of the heparin cofactor II (HCII) gene	69
5.12.6	Genomic organization of fish specific group V2 serpins	70
5.12.7	Sequence analysis of group V2 serpins	70
<b>5.13</b>	<b>Orthology analysis of group V3 serpin genes</b>	<b>73</b>
5.13.1	Gene structure of group V3 serpins	73
5.13.2	Synteny analysis of PAI1 genes	75
5.13.3	Synteny analysis of GDN genes	76
5.13.4	Synteny analysis of E3 genes	76
5.13.5	Synteny analysis of neuroserpin-pancpin genes	77
5.13.6	Sequence analysis of group V3 serpins	79
<b>5.14</b>	<b>Orthology analysis of group V4 serpin genes</b>	<b>83</b>
5.14.1	Gene structure of group V4 serpin genes	83
5.14.2	Synteny analysis of PEDF and $\alpha_2$ -antiplasmin	84
5.14.3	Synteny analysis of the C1 inhibitor	85
5.14.4	Sequence comparisons of group V4 serpins	86
<b>5.15</b>	<b>Orthology analysis of group V5 serpin genes</b>	<b>89</b>
5.15.1	Gene structure of ATIII genes	89
5.15.2	Synteny analysis of ATIII genes	90
5.15.3	Sequence comparisons of ATIII genes	91
<b>5.16</b>	<b>Orthology analysis of group V6 serpin genes</b>	<b>92</b>
5.16.1	Gene structure of group V6 serpins	92
5.16.2	Synteny analysis of group V6 serpins	93
5.16.3	Sequence comparisons of group V6 serpins	94
<b>6.</b>	<b>Discussion</b>	<b>97</b>
<b>6.1</b>	<b>Overview of vertebrate serpins from fishes to mammals</b>	<b>97</b>
<b>6.2</b>	<b>Evolutionary history of group V1 serpins</b>	<b>98</b>
<b>6.3</b>	<b>Phylogenetic history of group V2 serpins</b>	<b>99</b>
<b>6.4</b>	<b>Evolution of group V3 serpins</b>	<b>100</b>
<b>6.5</b>	<b>Evolutionary history of group V4 serpins</b>	<b>104</b>
<b>6.6</b>	<b>Summary of group V5 serpins</b>	<b>104</b>
<b>6.7</b>	<b>Overview of group V6 serpins</b>	<b>105</b>
<b>6.8</b>	<b>Intron gain and loss in vertebrate serpins</b>	<b>106</b>
<b>6.9</b>	<b>Strength and weakness in this work</b>	<b>111</b>

---

<b>6.10 Outlook</b>	<b>111</b>
<b>7. References</b>	<b>112</b>
<b>8. Appendix</b>	<b>122</b>
<b>8.1</b> Highly conserved residues present in > 70 % of the serpins	122
<b>8.2</b> GENEDOC usage	123
<b>8.3</b> Alignments	124
8.3.1 Protein sequence alignment of serpins from <i>C. intestinalis</i>	125
8.3.2 Protein sequence alignment of serpins from <i>B. floridae</i>	129
8.3.3 Protein sequence alignment of serpins from <i>S. purpuratus</i>	131
8.3.4 Protein sequence alignment of serpins from <i>N. vectensis</i>	134
8.3.5 Alignment of MNEI from vertebrates	135
8.3.6 Alignment of PAI2	137
8.3.7 Alignment of SPB5 (maspin) protein sequences from vertebrates	138
8.3.8 Alignment of SPB6 orthologs and paralogs (pSPB6) from vertebrates	140
8.3.9 Alignment of group V1 serpin sequences from chicken genome	142
8.3.10 Alignment of group V1 serpin sequences from <i>Xenopus tropicalis</i> genome	145
8.3.11 Alignment of group V1 serpin sequences from <i>Danio rerio</i> genome	147
8.3.12 Alignment of AGT (serpinA8) protein sequences from vertebrates	150
8.3.13 Alignment of heparin cofactor II sequences from vertebrates	153
8.3.14 Alignment of ZPI (serpinA10) protein sequences from vertebrates	157
8.3.15 Alignment of $\alpha_1$ -antitrypsin like sequences – Spn_215c from <i>Fugu</i> and <i>Tetraodon</i>	160
8.3.16 Alignment of $\alpha_1$ -antitrypsin like sequences – Fru-Spn-17 and Tni-Spn-4 from <i>Fugu</i> and <i>Tetraodon</i> respectively	161
8.3.17 Alignment of $\alpha_1$ -antitrypsin like serpins from <i>Gallus gallus</i>	162
8.3.18 Alignment of $\alpha_1$ -antitrypsin like serpins from <i>Xenopus tropicalis</i>	165
8.3.19 Alignment of $\alpha_1$ -antitrypsin like serpins from <i>Danio rerio</i>	168
8.3.20 Alignment of $\alpha_1$ -antitrypsin sequences from vertebrates	171
8.3.21 Alignment of THBG (serpinA7) protein sequences from vertebrates	174
8.3.22 Alignment of PAI1 sequences from vertebrates	176
8.3.23 Alignment of GDN sequences from vertebrates	178
8.3.24 Alignment of serpinE3 sequences from vertebrates	180
8.3.25 Alignment of pancpin sequences from vertebrates	182
8.3.26 Alignment of neuroserpin sequences from vertebrates	184
8.3.27 Alignment of PEDF sequences from vertebrates	186
8.3.28 Alignment of $\alpha_2$ -antiplasmin from vertebrates	188
8.3.29 Alignment of protein sequences of C1 inhibitor and fish specific group V4 (FSG4)	192
8.3.30 Alignment of ATIII sequences from vertebrates	195
8.3.31 Alignment of HSP47 homologs from vertebrates	197
<b>8.4 List of marker genes flanking serpin genes</b>	<b>200</b>
8.4.1 Marker genes flanking <i>Ciona</i> serpins	200
8.4.2 Marker genes flanking group V1 serpin genes in vertebrates	200

---

8.4.3	Marker genes flanking group V2 serpin cluster	201
8.4.4	Marker genes flanking serpinA7 in mammals	201
8.4.5	Marker genes flanking group V3 serpins in vertebrates	201
8.4.6	Marker genes flanking group V4 serpins in vertebrates	202
8.4.7	Marker genes flanking group V5 serpins in vertebrates	202
8.4.8	Marker genes flanking group V6 serpins in vertebrates	202
<b>8.5</b>	<b>List of figures</b>	<b>203</b>
<b>8.6</b>	<b>List of tables</b>	<b>205</b>
<b>8.7</b>	<b>Abbreviations</b>	<b>206</b>
<b>9.</b>	<b>Acknowledgements</b>	<b>207</b>

---

## 1. Summary

The serpins constitute a superfamily of proteins that fold into a conserved tertiary structure and employ a sophisticated, irreversible suicide-mechanism of inhibition. More than 6000 serpins have been identified, occurring in all three forms of the life- the eukaryotes, the prokaryotes and the archaea. Vertebrate serpins can be conveniently classified into six groups (V1-V6), based on three independent biological features - gene organization, diagnostic amino acid sites and rare indels. In the present work, the phylogenetic relationships of serpins from *Nematostella vectensis*, *Strongylocentrotus purpuratus*, *Ciona intestinalis*, four fish species, frog, chicken and mammals were investigated, using gene architecture analyses and stringent criteria for identification of orthologs. With some deviations, all vertebrate serpin genes fit into one of the six exon/intron gene classes previously identified, dating the existence and maintenance of these gene organizations before or close to the divergence of fishes. Group V1 and V2 gene families underwent rapid adaptive radiation along the lineages leading to mammals as indicated by an up to nine-fold increased number of family members, accompanied by a rapid functional diversification. In contrast, gene groups V3 to V6 display a rather conservative evolution with little changes since the divergence of fishes and the other vertebrates. The orthology assessment indicates that all vertebrates are equipped with a subset of strongly conserved serpins with functions that can be clearly correlated with basic vertebrate-specific physiology.

None of serpin genes from *C. intestinalis* shares a common exon-intron architecture organisation with any of the vertebrate serpin gene classes, nor was it possible to identify orthologs of vertebrates. The lack of gene architecture similarity and the complete absence of orthology between urochordate and vertebrate serpins indicate that major changes with bursts of character acquisition must have occurred during evolution of serpins in the time interval separating urochordates from chordates, indicating massive intron gains or losses and events providing C and N-terminal sequence extensions characteristic for today's vertebrate serpins. Lancelets and sea urchin genomes, in contrast, share one orthologous serpin with vertebrates. Rare genomic characters are used to show that orthologs of neuroserpin, a prominent representative of vertebrate group V3 serpin genes, exist in early diverging deuterostomes and probably also in cnidarians, indicating that the origin of a mammalian serpin can be traced back far in the history of eumetazoans. A C-terminal address code assigning association with secretory pathway organelles is present in all neuroserpin orthologs, suggesting that supervision of cellular export/import routes by antiproteolytic serpins is an ancient trait.

Phylogenomic comparisons show that, after establishment of canonical exon-intron patterns in the serpin superfamily at the dawn of vertebrate evolution, multiple intron acquisition events have occurred during diversification of a lineage of actinopterygian fishes. The novel introns were acquired within a limited time interval (on an evolutionary timescale), and no such events were observed in other groups of vertebrates. Examination of the sequences flanking the intron insertion points revealed that the genetic requirements for acquisition of novel introns might be less stringent than previously suggested. Finally, we argue that genome

---

compaction, a phenomenon associated with the fish lineage depicting preferential intron gain, might promote intron acquisition.

## 2. Introduction

### 2.1 Overview of serpins

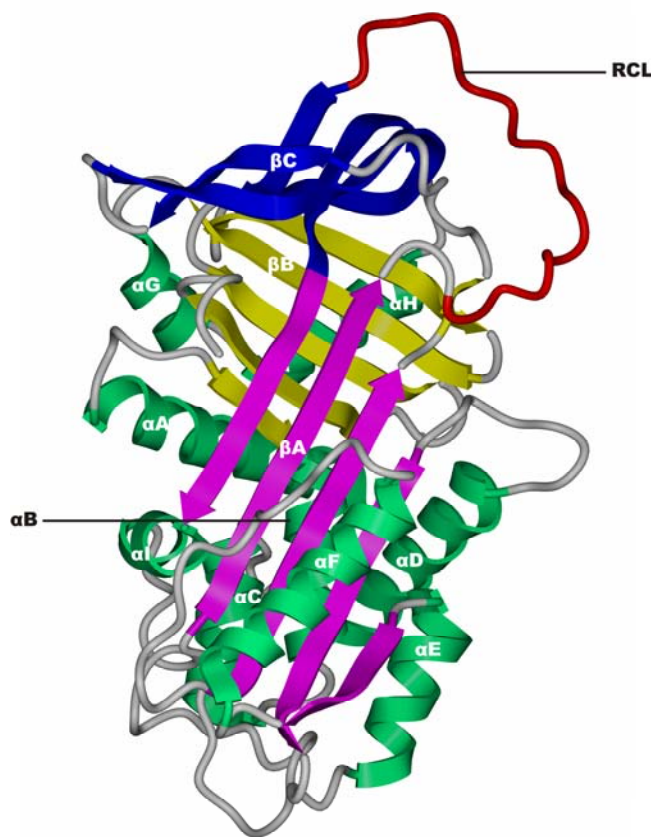
The serpins (serine protease inhibitors) are a superfamily of proteins that cover a highly divergent spectrum of functions. Serpins are primarily inhibitors of serine and/or cysteine proteases, but some family members have completely other tasks (Silverman *et al.*, 2001). Serpins are either classified into 16 different clades, designated A through H, based on sequence homology (Silverman *et al.*, 2001) or, based on intron-exon-structures, rare indels and diagnostic sites, they are categorized into six groups, V1-V6 (Ragg *et al.*, 2001). **Table 1** lists major vertebrate serpins.

**Table 1: Classification of vertebrate serpins.**

Groups	Clade	Serpins
V1	B	Ovalbumin, Gene Y Protein, Gene X Protein, Plasminogen activator inhibitor-2 or PAI-2, Squamous cell carcinoma antigen 1 or SCCA-1, SCCA-2, Protease Inhibitor 2 or PI-2, PI-6, PI-9, Bomapin, Headpin, Maspin, Megsin, Epipin, Yukopin
V2	A	$\alpha_1$ -antitrypsin, Corticosteroid-binding globulin, Protein C inhibitor, Angiotensinogen, $\alpha_1$ -antichymotrypsin, PI-4, Thyroxine-binding globulin
	D	Heparin cofactor II
V3	E	PAI-1, Nexin-1, SerpinE3,
	I	Neuroserpin, Pancpin
V4	F	$\alpha_2$ -antiplasmin or A2AP, Pigment epithelium derived factor or PEDF
	G	C1-inhibitor
V5	C	Antithrombin III or ATIII
V6	H	Heat shock protein 47kDa or HSP47

#### 2.1.1 Structure and mechanism of action

Serpins are single domain proteins with a conserved core of ~350-400 residues often possessing N- or C-terminal extensions, resulting in an overall molecular mass of ~40-60 kDa. N- and/or O-glycosylations are frequently observed in extracellular serpins (Gettins *et al.*, 1996; Gettins, 2002). The conserved three-dimensional structure of serpins is composed of three  $\beta$ -sheets ( $\beta$ A- $\beta$ C) and 8-9  $\alpha$ -helices ( $\alpha$ A- $\alpha$ I) (**Figure 1**). The hallmark of the serpin inhibitory mechanism is a large scale conformational change involving the reactive center loop (RCL). The RCL is an exposed flexible loop of about 17-20 residues, which interacts with a target protease (**Figure 1**). The RCL acts as a bait imitating a protease substrate that is cleaved between the positions P1 and P1'. Starting from the scissile bond, residues are designated P1, P2, P3,... and P1', P2', P3',... in the N- or C-terminal direction, respectively, according to the standard nomenclature (Schechter and Berger, 1967). Considering the interaction with serine proteases, composition and conformation of the RCL and especially the P1 residue are the major determinants of target specificity (Carrell and Travis, 1985).



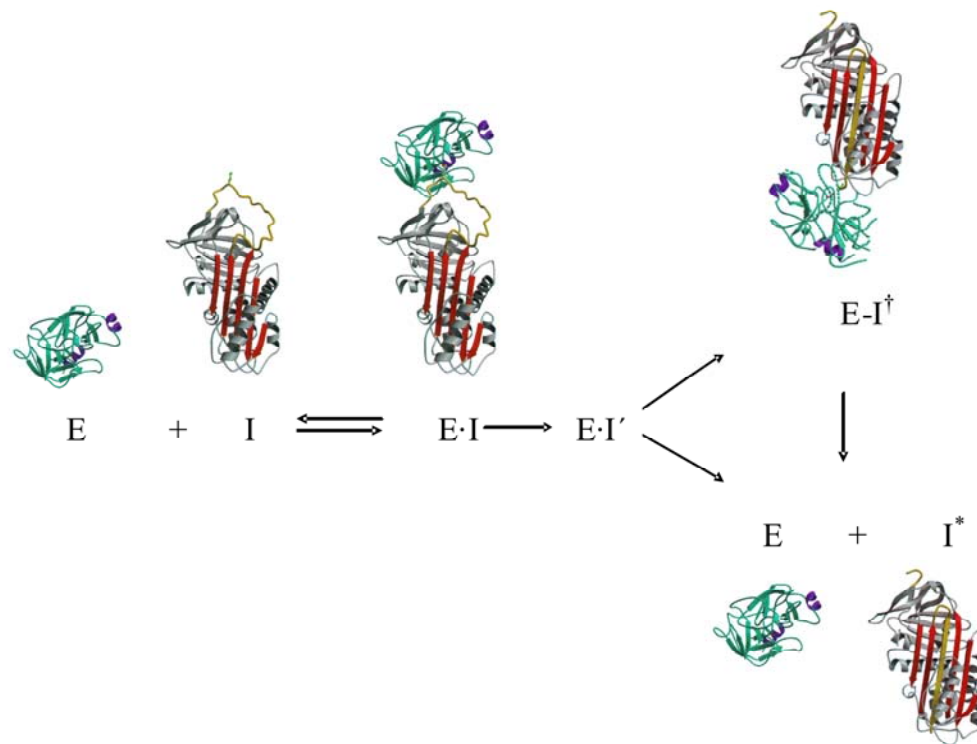
**Figure 1: Three-dimensional structure of uncleaved  $\alpha_1$ -antitrypsin (PDB code 1HP7), a member of the serpin family.** There are nine  $\alpha$ -helices ( $\alpha A$ - $\alpha I$ ) colored in green and three  $\beta$ -sheets ( $\beta A$ - $\beta C$ ) illustrated in magenta, yellow, and blue, respectively. The structure is visualized with YASARA<sup>1</sup>. RCL – Reactive Center Loop.

In the native state, serpins adopt a metastable, stressed conformation that can undergo substantial structural rearrangements upon cleavage by the target protease. The inhibition of proteases by serpins is described by the “branched pathway” mechanism (**Figure 2**), (Lawrence *et al.*, 2000; Gettins, 2002).

In the first step, the protease [**E**] recognizes the exposed RCL bait of the serpin [**I**] and forms a reversible, non-covalent Michaelis complex [**E**·**I**]. Serine and cysteine proteases are characterized by an active site that contains a nucleophilic serine or cysteine residue. The nucleophilic attack of the protease at the scissile P1–P1' bond of the serpin results in cleavage of the RCL and release of the C-terminal part of the serpin, followed by formation of a covalent acyl-enzyme intermediate [**E**–**I**], as described by the catalytic triad mechanism (Nelson and Cox, 2005). From this point on, the reaction can continue in two different directions. If the protease is able to fulfill its catalytic action, deacylation occurs (non-inhibitory pathway), leading to a release of the active protease [**E**] and the cleaved, inactive serpin [**I**\*]. In the inhibitory pathway, the serpin adopts its relaxed, thermodynamically favored conformation (Silverman *et al.*, 2001). The N-terminal part of the RCL (residues P1–P14) inserts into  $\beta$ -sheet A and extends the  $\beta$ -sheet structure to form a fully anti-parallel sheet with six instead of five  $\beta$  strands.

<sup>1</sup> YASARA webpage, [www.yasara.org](http://www.yasara.org)





**Figure 2: Branched pathway model.** The protease E (cyan) binds reversibly to the RCL of the serpin I (grey) and forms a non-covalent Michaelis complex [E-I]. Cleavage of the RCL results in formation of a covalent acyl-enzyme intermediate [E-I\*]. Insertion of the RCL into  $\beta$  sheet A leads to inactivation of the protease by deformation (E-I $^\dagger$ , inhibitory pathway), whereas deacylation produces inactive serpin I\* and active protease E. This figure is adopted from Huntington *et al.* (2000).

During this process, the covalently bound protease is translocated by 70 Å to the opposite pole of the serpin and compressed against the inhibitor body. According to the X-ray structure of the trypsin/ $\alpha_1$ -antitrypsin complex [E-I $^\dagger$ ] (Huntington *et al.*, 2000), the conformational change leads to a significant deformation of the protease and its catalytic center (**Figure 2**). As a result, deacylation rates are decreased by 6–8 orders of magnitude, kinetically trapping the acyl-enzyme intermediate, and inactivating both protease and serpin. In vitro, the enzyme-inhibitor complexes [E-I $^\dagger$ ] have half-lives between hours and weeks. In vivo, the complexes are recognized by receptors and cleared by proteolysis of both components (Silverman *et al.*, 2001; Gettins, 2002). The length and flexibility of the RCL, especially the hinge region (residues P15–P9), are important determinants for successful inhibition. Inhibitory serpins have a highly conserved hinge region with small residues to facilitate strand insertion. The positions P15 and P12–P9 are usually occupied by glycine and alanine residues, respectively. Mutations in this region result in a loss of the inhibitory function (Huber and Carrell, 1989). Partial insertion of the RCL, leads to an inactive, latent state of the serpin.

## 2.1.2 Physiological functions

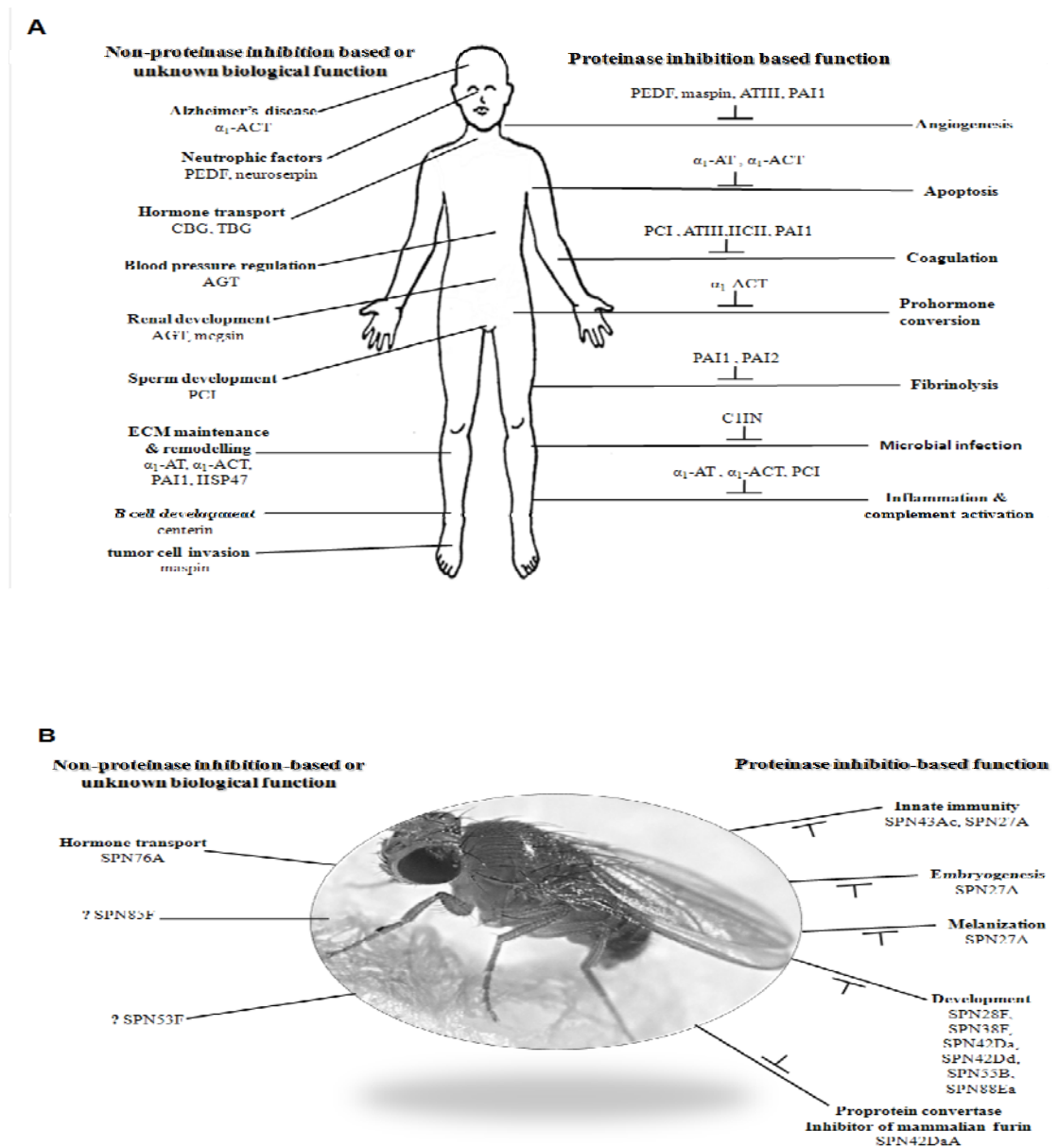


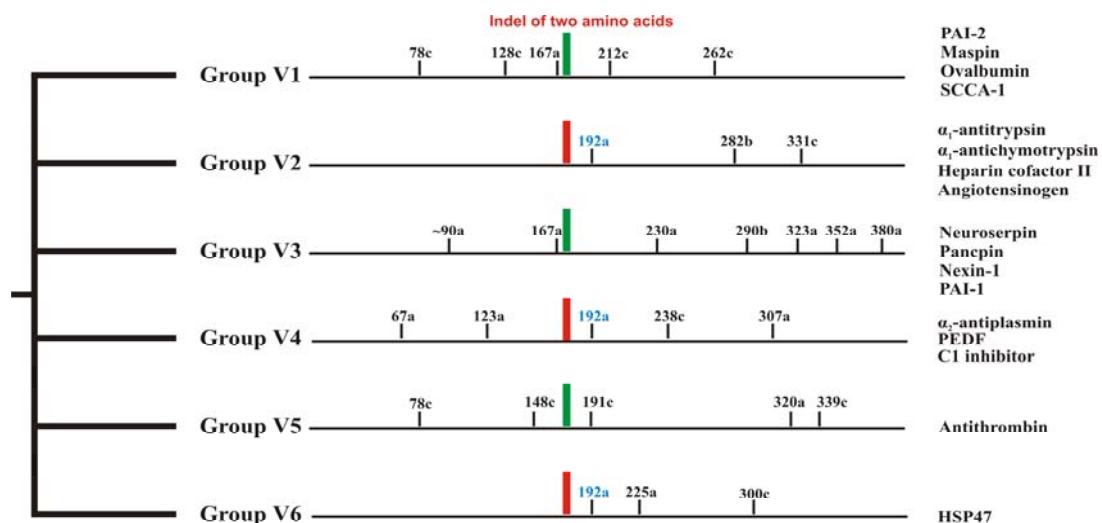
Figure 3: Physiological functions of selected serpins in vertebrates (A) or invertebrate model organisms (B).

Serpins are major factors in regulating proteolytic activities within our body to avoid excessive proteolysis. **Figure 3** depicts some important roles of human and *Drosophila melanogaster* serpins. A highly divergent functional spectrum is covered by serpins both with vertebrates (**Figure 3A**) and with invertebrates (**Figure 3B**). Serpins for example, regulate dorsal-ventral axis formation and immune regulation in insects such as *Drosophila* (Levashina *et al.*, 1999; Ligoxygakis *et al.*, 2003), embryo development in nematodes (Pak *et al.*, 2004), or proprotein convertases in lancelets (Bentele *et al.*, 2006).

### 2.1.3 Evolution of serpins

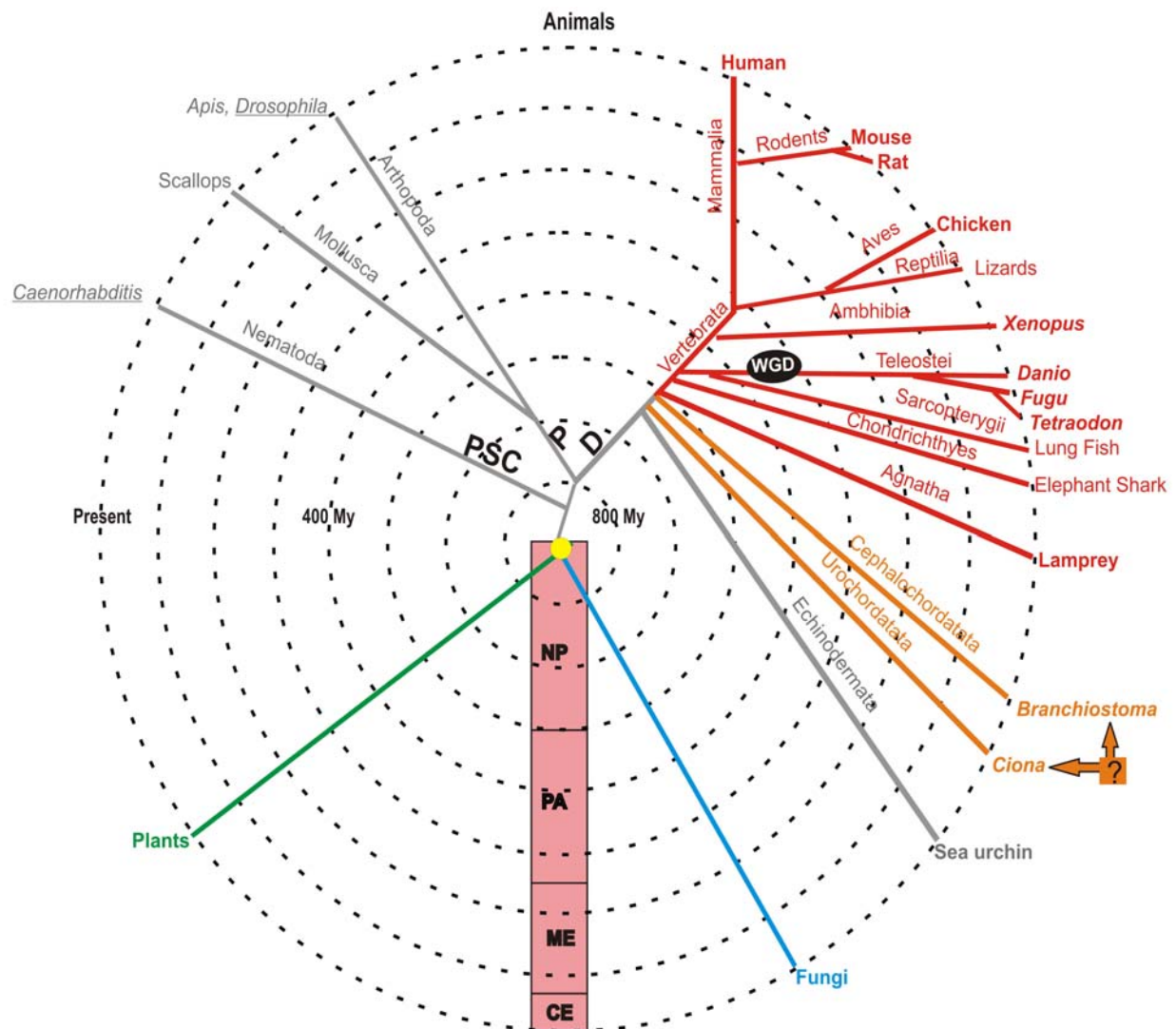
In metazoans, serpins have undergone divergent evolution over a period of about 650-700 million years (Kumar and Ragg, 2008). A number of phylogenetic studies have been undertaken using sequence analysis of the serpins. Early investigations suggested the establishment of this multigene family through inter- and intra-chromosomal gene duplications. Several gene clusters have arisen; encoding functionally diverse serpin proteins (see previous section). In metazoans, serpin genes display highly variable exon-intron patterns that, however, may be strongly conserved within some taxa. Gene architecture and other rare genetic characters constitute a robust basis to group vertebrate serpins. Based on number, positions, and phases of introns, serpins have been classified into six groups maintained at least since the fish/tetrapod split (**Figure 4**). Most known serpin genes contain a non-coding first exon and a partly non-coding last exon (Ragg *et al.*, 2001). However, the genes encoding  $\alpha_1$ -antitrypsin and heat shock proteins (HSP47) contain an alternatively spliced first exon. Computational analysis showed a strong similarity in the classification of vertebrate serpins either according to classical phylogenetic analysis of amino acid sequences or gene structure-based categorization (Atchley *et al.*, 2001; Ragg *et al.*, 2001). Vertebrate serpin genes with equivalent gene structures often tend to be organized in clusters (Benarafa and Remold-O'Donnell, 2005); however, close physical linkage is not always found.

Interestingly, none of altogether 24 intron positions mapping to the core domain of vertebrate serpins is shared by all of these six gene groups; however, characteristic amino acid indels provide some further cues for unraveling phylogenetic relationships (Ragg *et al.*, 2001). None of the group-specific vertebrate gene architectures is found in earlier diverging animal taxa, though a few vertebrate-specific intron positions are present in a scattered fashion in some basal metazoans.



**Figure 4: Gene structure-based phylogenetic classification of vertebrate serpins.** Positions of introns refer to the human  $\alpha_1$ -antitrypsin sequence. A two amino acid indel present between positions 173 and 174 ( $\alpha_1$ -antitrypsin numbering) suggests that groups V1, V3, and V5 (indicated by green bar) are more closely related to each other than to the other groups. Groups V2, V4, and V6 lack the 173/174 indel (marked by red bars) and depict an intron at position 192a, implying shared ancestry. Some group V1 members contain an additional intron at position 85c (not shown). This figure is based on three publications of Ragg's group (Atchley *et al.*, 2001; Ragg *et al.*, 2001; Kumar and Ragg, 2008).

In order to investigate the evolution of vertebrate serpin genes, I considered representative genomes of all classes of vertebrates with exception of reptiles, for which only a single initial genome draft version (lizard) is available. **Figure 5** shows the position and evolution of vertebrates within the Tree of Life (TOL), based on data from Kumar and Hedges (1998), Hedges (2002) and Ponting (2008). Together with urochordates and cephalochordates, vertebrates constitute the phylum Chordata.



**Figure 5: The phylogenetic tree of animal evolution.** The last common ancestor (yellow center) of multicellular life was split about 800 millions of years ago (Mya) into three main branches – animals, fungi, and plants. The evolution of animals started with branching out of Pseudocoelomata (PSC) followed by divergence into Proteostomia (P) and Deuterostomia (D). Higher invertebrates include echinoderms, cephalochordates, and urochordates. The position of cephalochordates and urochordates is still in debate (indicated by ?) culminating in the question which are being closer to vertebrates (connecting link between vertebrata and invertebrata). Vertebrates (red) arose about 500-520 Mya. About 336-404 Mya, a fish-specific whole genome duplication (WGD) is believed to have occurred. Genomes considered in this work are marked with bold letters. For simplicity, plants and fungi branches are not expanded here. Geological time periods are also shown. CE, Cenozoic, ME, Mesozoic, PA, Palaeozoic, NP, Neoproterozoic. Time lines and geological time periods are taken from Kumar and Hedges (1998), Hedges (2002) and Ponting (2008).

At the dawn and during evolution of chordates, genome duplications are believed to be responsible for bringing in diversities. During vertebrates evolution, whole genome duplication (WGD) events happened after separation of fishes from tetrapods (**Figure 5**) as proposed by Susumu Ohno (Ohno, 1970; Ohno, 1999).

## 2.2 Rare genomic changes

Rare genomic changes (RGC) are mutational changes that have occurred in the genomes of particular clades. These changes may serve as phylogenetic markers for characterization of particular clades (Rokas and Holland, 2000). **Table 2** gives an overview on RGCs – indicating types, taxonomic resolution, extent of homoplasy<sup>1</sup> and taxa in which RGCs are applicable.

**Table 2: Overview of rare genomic change (RGC) markers for phylogenetic purposes.** Modified from Rokas and Holand (2000). \$ = mitochondrial, # = chloroplast

Marker	Taxonomic resolution	Homoplasy	Taxa in which RGCs are applicable
Intron indels	Wide ranging	Low	Eukaryotes
Retroposons (SINEs and LINEs) <sup>2</sup>	Within orders	Zero to very low	Animals
Signature sequences	Wide ranging	Unknown	All branches of the life
mtDNA\$ genetic code variants	Phyla to classes	Low to moderate	Eukaryotes
Nuclear DNA genetic code variants	Phyla	Low to moderate	All branches of the life
mtDNA gene order	Wide ranging (phyla to families)	Low to moderate in animals, Higher plants, fungi and protists	Eukaryotes
cpDNA# gene order	Families	Low	Plants
Gene duplications	Wide ranging	Unknown	All branches of the life
Comparative cytogenetics	Within phyla	Unknown	All branches of the life
Overlapping genes	Wide ranging	Low	All branches of the life

### 2.2.1 Intron gain and loss

Gains or losses of introns are important evolutionary markers. The mechanisms of intron gain (**Figure 6**) and intron loss (**Figure 7**) have been reviewed in detail (Roy and Gilbert, 2006).

<sup>1</sup> Acquisition of the same character state in two taxa is not because of common descent.

<sup>2</sup> SINEs, short interspersed elements; LINEs, long interspersed elements.

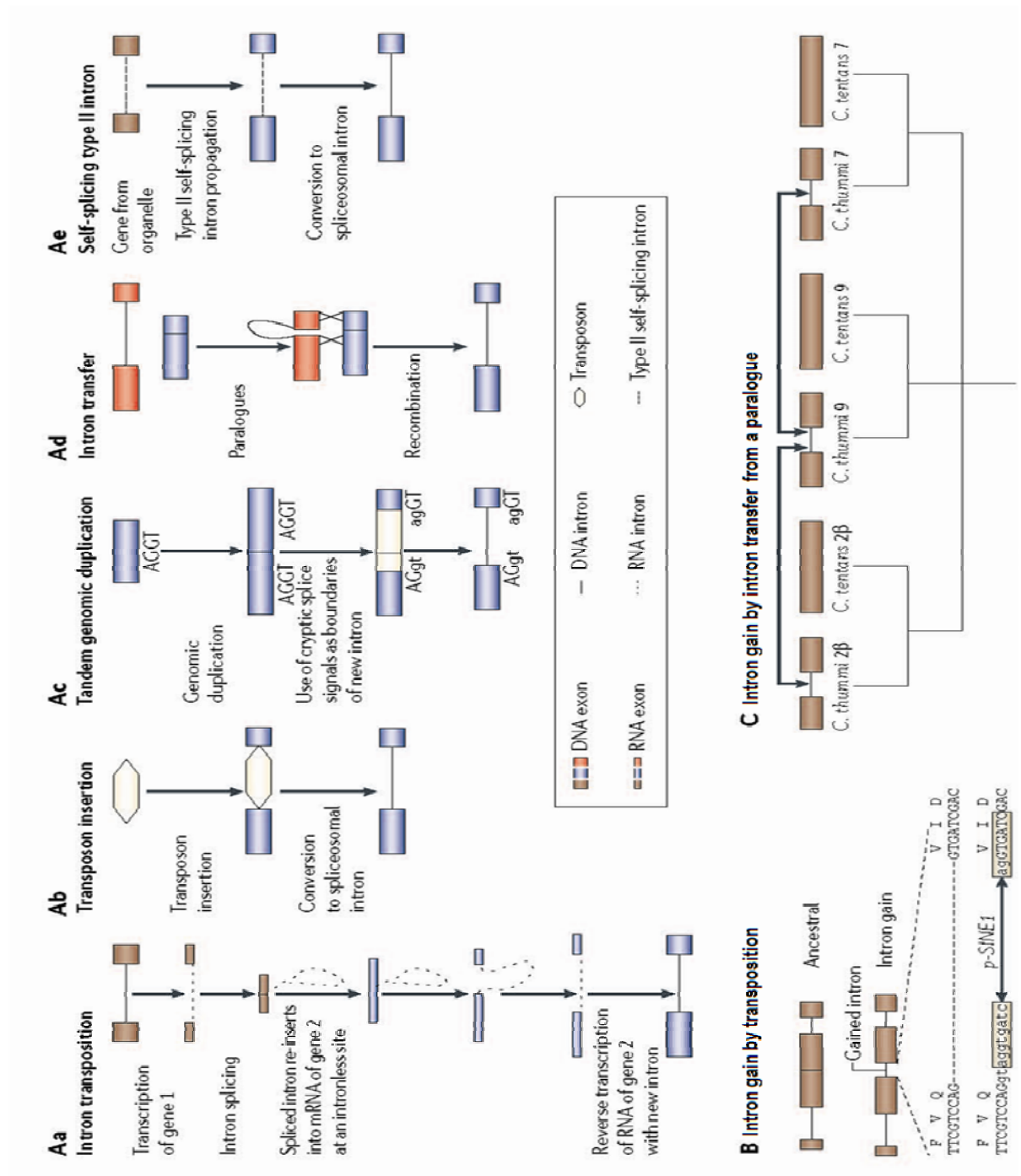


Figure 6: Models and examples of intron gain. Different models A (Aa-Ae) and examples B and C (taken from Roy and Gilbert, 2006).

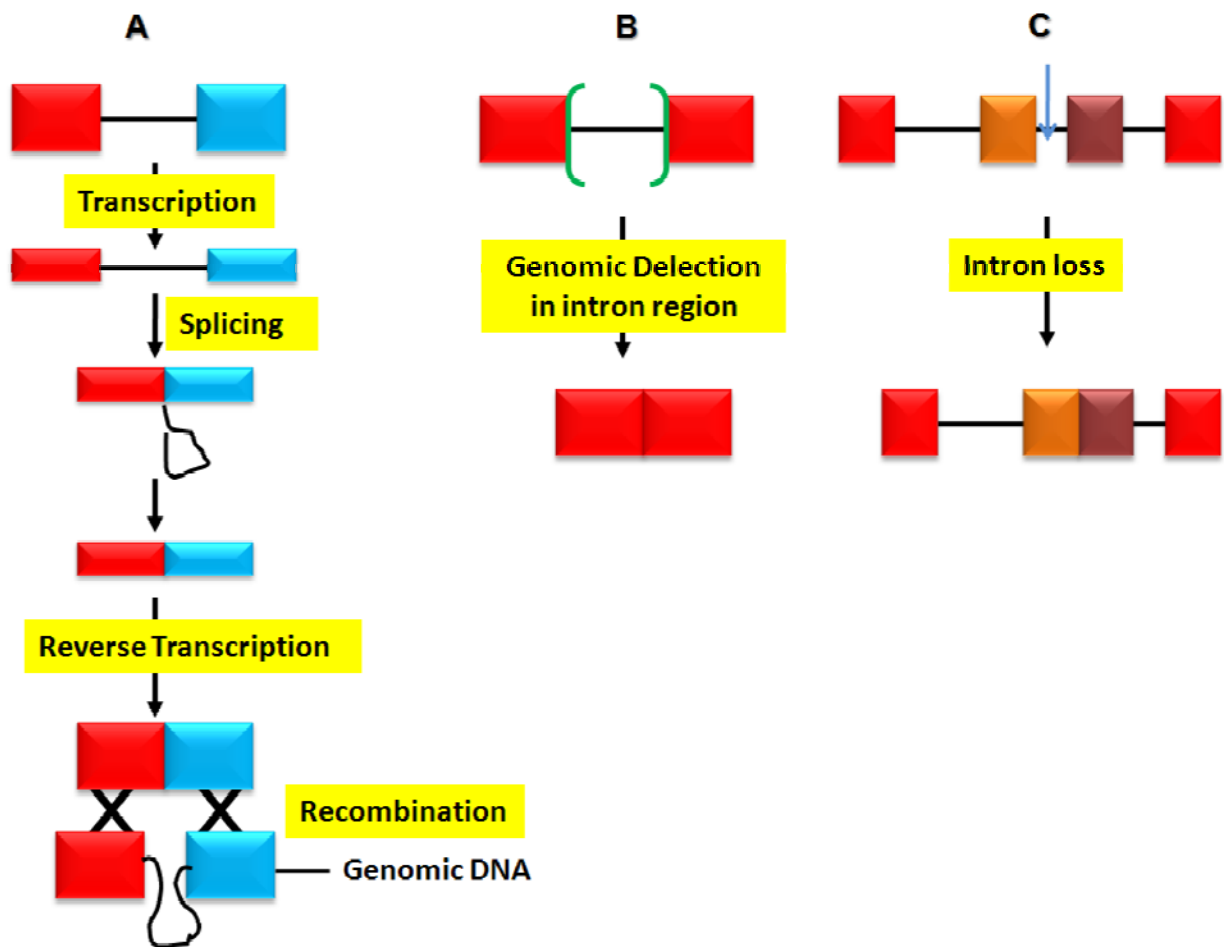


Figure 7: Models (A-B) and example (C) of intron loss (taken from Roy and Gilbert, 2006).

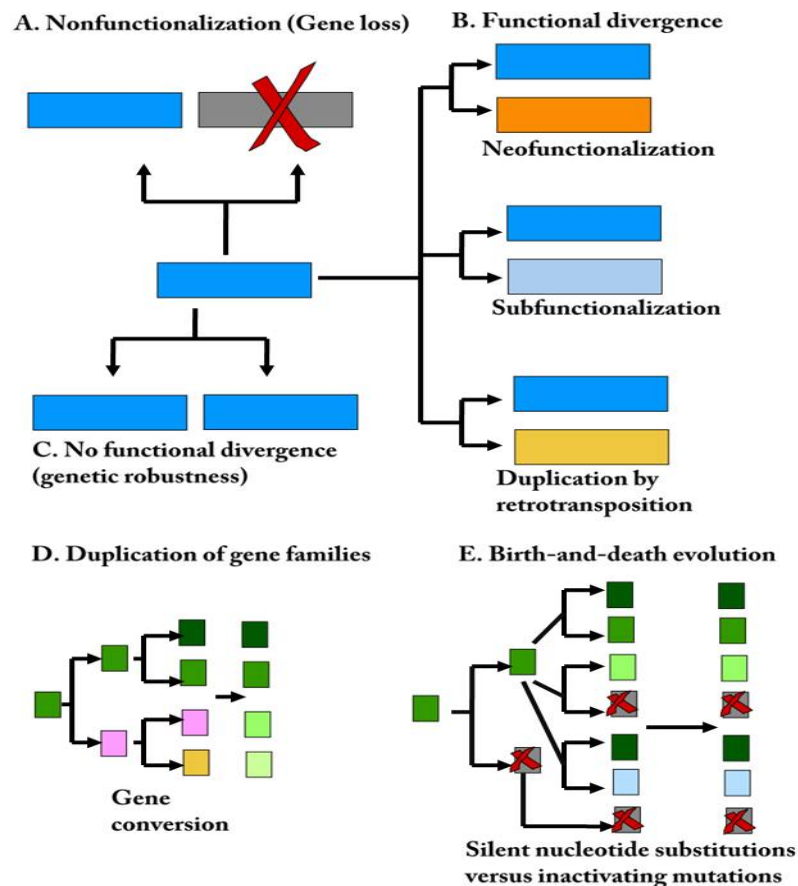
### 2.2.2 Rare indels

Events of insertion-deletion are well-commented examples of rare genomic changes. Indels can include gain/loss of few nucleotides and gain/loss of introns. For instance, vertebrate serpins are classified into six groups, based on sequence indels and intron indels (Ragg *et al.*, 2001)

### 2.2.3 Gene duplications and fates of duplicated genes

Gene duplications are considered to be major genetic basis for producing novel genetic variations. There are three types of gene duplications: whole genome, segmental and small scale duplications (Conrad and Antonarakis, 2007).





**Figure 8: Fate of duplicated single genes (A-C) and duplicated gene families (D-E).** Modified from Conrad and Antonarakis, 2007. Gene loss/inactivation is indicated by a red X.

About four decades ago, Susumu Ohno developed an insightful hypothesis arguing that gene duplication is a key factor shaping evolution. His model and its general predictions continue to attract much attention (Ohno, 1970; Ohno, 1999) in the post-genomic era with hundreds of genomes being available to test this hypothesis.

On an evolutionary scale, gene duplication may result in new functions via different scenarios (**Figure 8**) including: (i) Nonfunctionalization - predominant outcome is loss of function in one of the two gene copies (**Figure 8A**). (ii) Neofunctionalization - one gene copy may retain the original function while the other acquires a novel, evolutionarily advantageous/adaptive function (Force *et al.*, 1999). (iii) Subfunctionalization - after duplication, mutations may occur in both genes that specialize to perform complementary functions (Lynch and Conery, 2000; Lynch and Force, 2000).

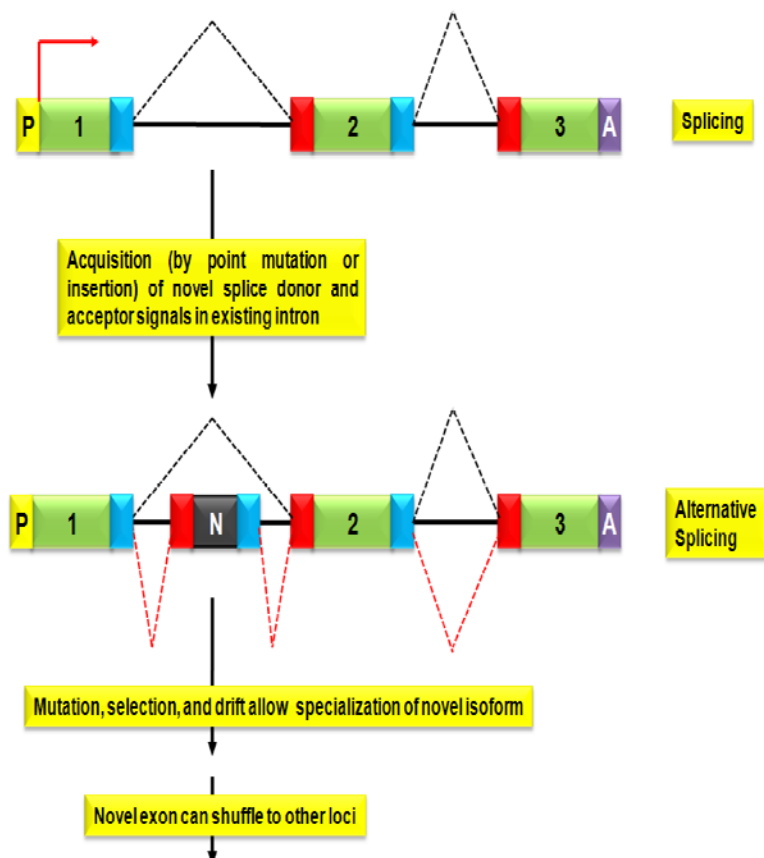
The question of how duplicate genes are retained in a population remains controversial. Classical duplication-degeneration-complementation/subfunctionalization models do not invoke positive selection, but stipulate a higher retention rate of duplicate genes in small rather than larger populations. Considerably more retentions and fewer losses of duplicate genes in rodents as compared with humans indicate that positive selection may play a more important role than originally anticipated (Shiu *et al.*, 2006). If two redundant gene copies were retained in the genome without significant functional divergence, the organism may



acquire increased genetic robustness against harmful mutations (**Figure 8C**). In multigene families descended from a common ancestor, individual genes in the group exert similar functions and have similar DNA sequences (Nei *et al.*, 2000; Nei and Rooney, 2005). One concept, concerted evolution, applies particularly to localized and typically tandem copies of a gene. The concept posits that all genes in a given group evolve coordinately, and that homogenization is the result of gene conversion (**Figure 8D**). For most multigene families, the currently favored model is birth-and-death evolution, according to which similarity in protein sequence among the members of a family is assured by strong purifying selection, such that individual genes evolve essentially via silent synonymous nucleotide substitutions (**Figure 8E**), (Nei *et al.*, 2000; Nei and Rooney, 2005).

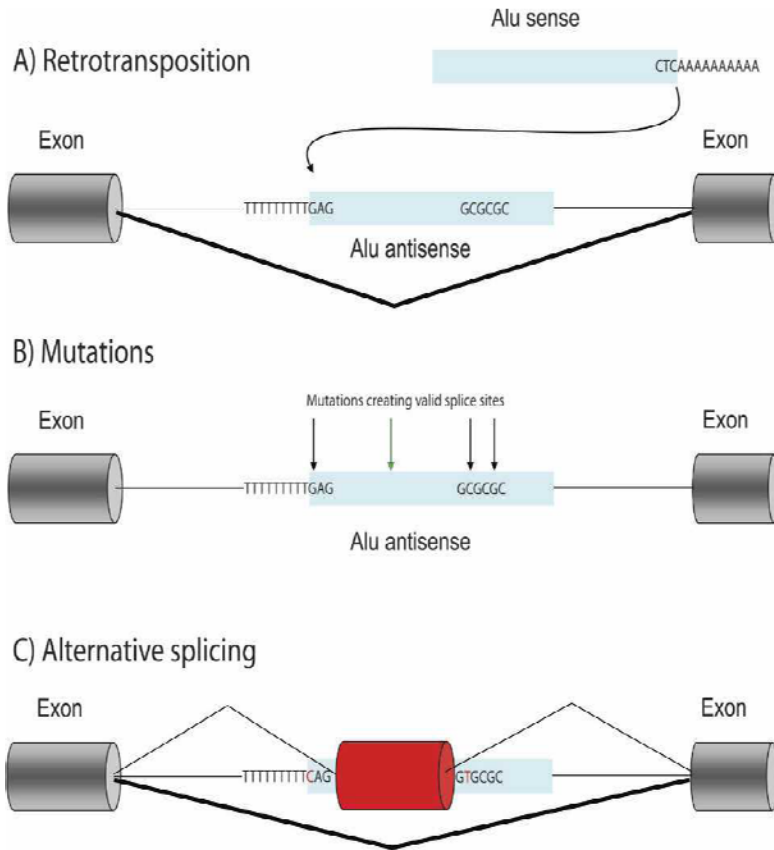
### 2.2.4 Exonization of non-coding regions in genomes

New exons are normally created by duplication of genes or exons in metazoan genomes. However, the most intriguing processes are exonization events, where intronic sequences are converted to *de novo* exons (**Figure 9**), (Schmidt and Davies, 2007).



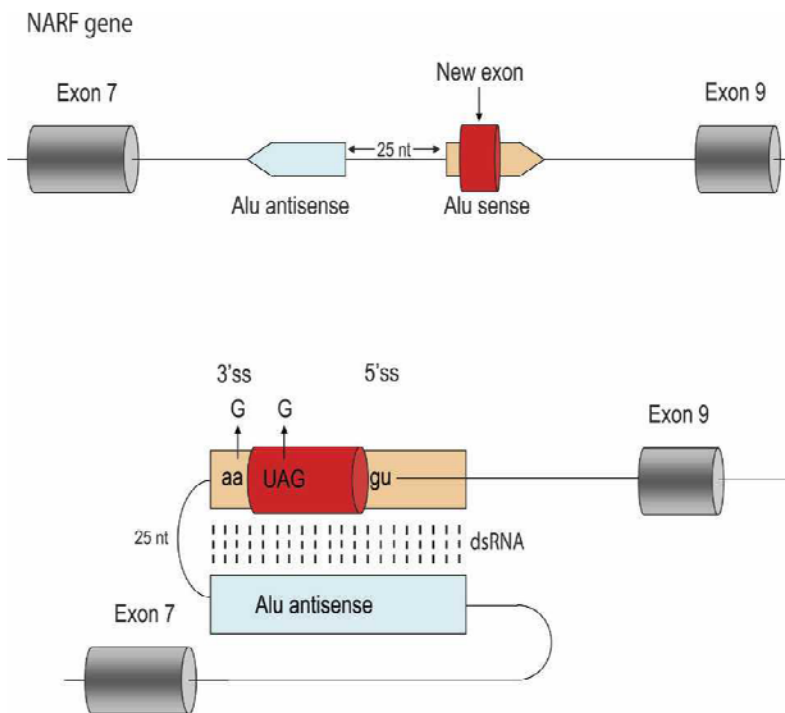
**Figure 9: Exonization of intron sequences to generate novel gene products.** A simple gene is depicted with three exons (green), introns (black lines) with splice donor (blue) and acceptor (red) signals, promoter P sequences (yellow), initiation site (red bent arrow), poly-A signals (violet), and splicing pattern (dashed lines). Acquisition of splice donor and acceptor signals within an existing intron can generate a novel exon (labeled “N”, black). The phase of the exon–exon junctions must be preserved. Alternative splicing (red bent lines) can produce either the original protein or the modified version with the novel polypeptide. This novel exon might subsequently be transferred to other parts of the genome. This figure is taken from Schmidt and Davies (2007).

Exonization of intronic sequences, particularly those originating from repetitive elements such as *Alu* repeats (**Figure 10**), are now widely documented in different vertebrate genomes (Sorek, 2007).



**Figure 10: Exonization of an *Alu* element.** (A) *Alu* element is inserted into introns of primate genes by retrotransposition. (B) During the course of evolution, mutations within pseudo-splice sites in the intronic *Alu* activate these sites (black arrows). Mutations changing splicing regulatory elements are also possible (arrow). (C) Following these mutations, part of the *Alu* sequence is recognized as a new exon (“exonized”), and spliced into the transcript. Typically, the *Alu*-containing transcript is the minor splice form, as in most cases the created splice sites are weak. Most exonizations involve the antisense orientation of the *Alu* sequence, presumably because of the preceding long poly-T that serves as a strong poly-pyrimidine tract necessary for the 3’SS recognition. This figure is taken from Sorek (2007).

It can be mediated by RNA-editing (Figure 11) (Lev-Maor *et al.*, 2007).



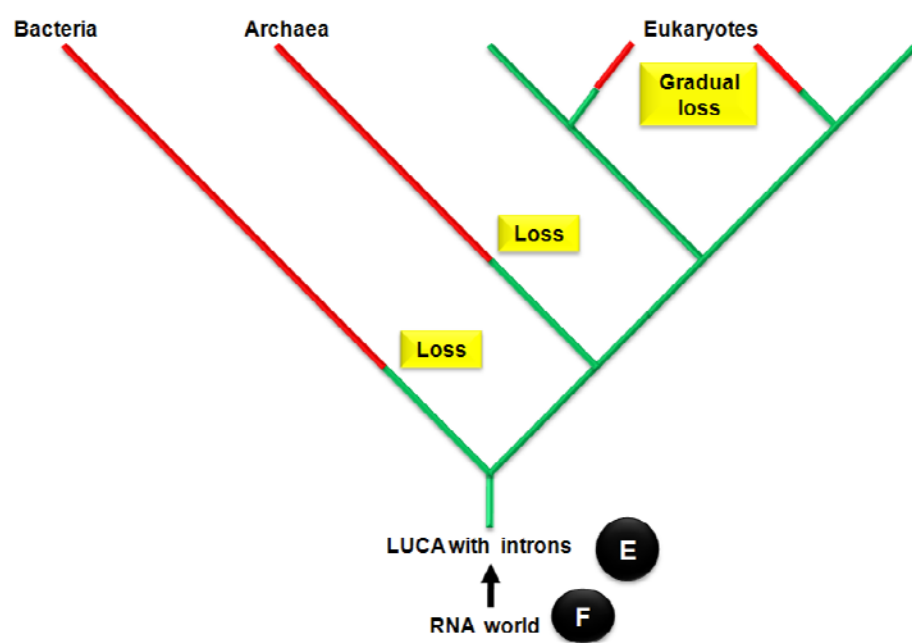
**Figure 11: Exonization through RNA editing** (Lev-Maor *et al.*, 2007; Sorek, 2007). Shown is a schematic illustration of the genomic region spanning exons 7–9 of the human NARF gene (not to scale). Exons are depicted as cylinders. The *Alu* element that is the source of the new exon is orange; an intronic, antisense orientation *Alu* sequence (light blue) is 25 bp upstream of the exonized *Alu*. Sense and antisense *Alus* fold to form a double-stranded RNA (dsRNA) secondary structure, thus allowing RNA editing to take place (lower panel). RNA editing changes an AA dinucleotide into a functional AG 3’ splice site and also changes a UAG stop codon into a UGG Trp codon. Thus, RNA editing leads to the creation of a new functional exon. This figure is taken from Sorek (2007).

Such *de novo* appearance of exons is very frequently associated with alternative splicing, with the new exon-containing variant typically being the rare one. This allows the new variant to

be evolutionarily tested without compromising the original gene product, and provides an evolutionary strategy for generation of novel functions with minimum damage to the existing functional repertoire. With multiple genomes available to study, it is becoming clear that exonizations of introns or intergenic sequences are an efficient way to produce novel gene products and are not as rare as expected before (Sorek, 2007).

### 2.3 Intron evolution theories

Spliceosomal introns are present in the nuclear genomes of all characterized eukaryotes. The discovery of introns and splicing in the 1970s led to debates about their origin. The most prominent hypotheses are the ‘Introns Early’ and the ‘Introns Late’ hypothesis (Jeffares *et al.*, 2006; Roy and Gilbert, 2006). However, there are compromised or mixed models. The ‘Introns Early’ theory proposed that introns were already present in the last universal common ancestor (LUCA) of prokaryotes and eukaryotes (‘E’ in **Figure 12**), where they were merely the genomic regions between genes (Darnell, 1978; Gilbert, 1978). These regions then suffered different fates in the different lineages: they were lost in all prokaryote lineages, while in eukaryotes they were maintained as introns by the appearance of the spliceosome. According to this theory, a modern protein is a concatenation of earlier, smaller proteins achieved by one of these two evolutionary processes. The ‘Intron first’ hypothesis proposes that introns

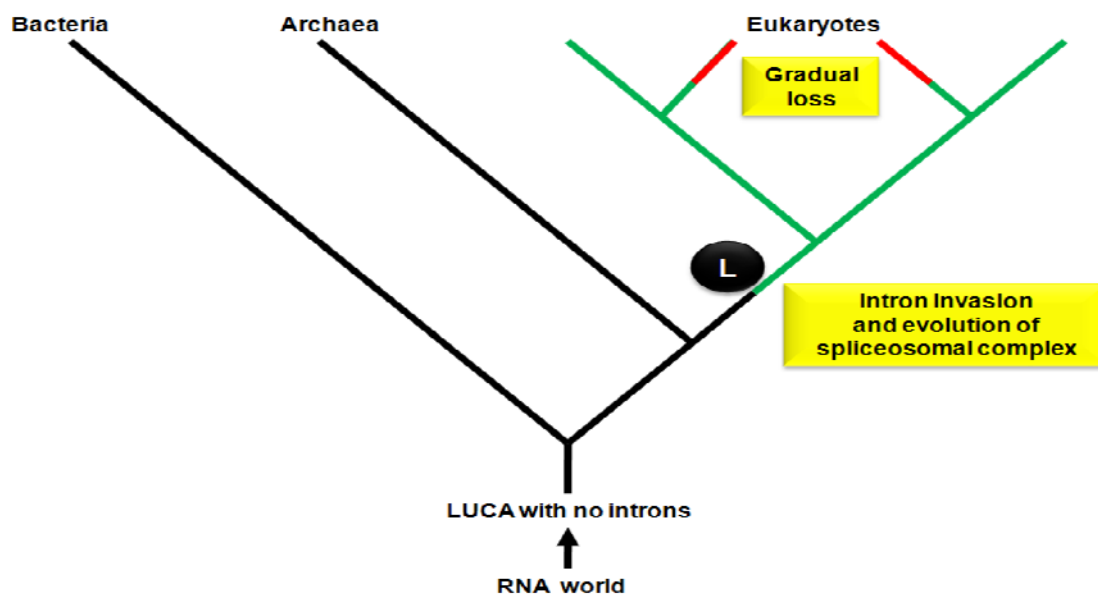


**Figure 12: Intron early (E) hypothesis and intron first (F) hypothesis during evolution of life.** The green branches indicate lineages containing introns; the black branches denote pre-intron stages and the red branches indicate secondary loss of introns. This figure is modified from Jeffares *et al.* (2006).

and the spliceosome are remnants from the RNA world (‘E’ in **Figure 12**) (Jeffares *et al.*, 1998; Poole *et al.*, 1998) and it is similar to intron-early model. This model was initiated from the observation that putatively ancient snoRNA genes are often encoded by introns. Because

RNAs were the only catalysts for the assembly of an all-RNA ribosome before the advent of proteins, snoRNAs must have been used for the assembly of the proto-ribosome as it evolved towards full protein producing capacity (Poole et al., 1999). Thus, the introns that contain snoRNAs pre-date the protein-coding exons that surround them. The splicing of snoRNA-encoding introns from transcripts without protein coding potential, and the processing of pre-rRNA and pre-tRNAs by RNase P are examples of how RNA processing might have occurred before proteins evolved (Jeffares et al., 1998; Poole et al., 1998).

In contrast to above hypotheses, the ‘Introns Late’ hypothesis proposes that spliceosomal introns only appeared in eukaryotes (‘L’ in **Figure 13**), where they were derived from self-splicing introns that invaded previously undivided genes, and that the spliceosome evolved as a way of removing them (Cavalier-Smith, 1991; Palmer and Logsdon, 1991; Boeke, 2003). Self-splicing introns/retrointrons are a type of genomic parasite: they insert themselves into the host genome and, when transcribed, their RNA catalyses its own excision – although sometimes assisted by a protein translated from sequences within the intron (Lambowitz and Zimmerly, 2004).



**Figure 13: Intron late (L) hypothesis during evolution of life.** The green branches indicate lineages containing introns; the black branches denote pre-intron stages and the red branches indicate secondary loss of introns. This figure is modified from Jeffares *et al.* (2006).

Mixed or compromised models of intron evolution include aspects of both the ‘Intron Early’ and the ‘Intron Late’ hypothesis. According to SW Roy, some introns are recent, most are ancient (Roy, 2003) whereas Rogozin *et al.* proposes that most introns are recent and some are ancient, but not necessarily very old (Rogozin *et al.*, 2003).

## 2.4 Aim of this study

Serpins are involved in a wide array of physiological processes amongst different taxa in the tree of life. Understanding evolutionary history of serpins is a challenging task and poses notorious problems in animal genomes. Notably, vertebrate serpins were classified into six groups (V1-V6) based on rare indels, diagnostic sites and gene structures (Ragg *et al.*, 2001). However, this classification was based on a limited set of genomic data, although it is a more reliable classification system than other sequence-based classification systems for serpins.

Therefore, the aim of present study is to examine an extended set of genomes from vertebrates of evolutionary importance in order to unravel whether this classification system holds in all vertebrates or whether during over 450-500 million years of vertebrate evolution deviations occurred. In order to extend our understanding of this classification to additional non-mammalian vertebrates, we chose the following evolutionary important genomes: i) *Gallus gallus* (bird), ii) *Xenopus tropicalis* (frog), and fish genomes - iii) *Fugu rubripes*, iv) *Tetraodon nigroviridis*, v) *Danio rerio*, vi) *Petromyzon marinus* (lamprey). The serpins from these genomes are to be characterized and compared with two more fish genomes - medaka and stickleback. Orthologs and paralogs of human serpins are to be assigned based of sequence features, indels, gene architectures, and syntenic analysis from above mentioned genomes.

A further aim of this study is to analyze intron gain/loss in different serpin genes in non-mammalian vertebrate genomes. There are 25 conserved intron positions as differentiating markers for six groups (V1-V6). An additional objective of this study is to extend this analysis to non-vertebrate model organisms such as *Branchiostoma floridae* (lancelet), *Ciona intestinalis* (sea squirt), *Strongylocentrotus purpuratus* (sea urchin) and *Nematostella vectensis* (sea anemone). This comparative analysis of serpins from metazoan genomes might provide some clues to the origin and ancestry of vertebrate serpin genes.

### 3. Materials

#### 3.1 Genomes

The genomes analyzed in our study are listed in the **Table 3**, which includes vertebrate genomes as well as the genomes of evolutionarily important animals.

**Table 3: Genomes analyzed.**

Genome	Major database used	Reference
<i>Homo sapiens</i>	<a href="http://www.ncbi.nlm.nih.gov/genome/guide/human/">http://www.ncbi.nlm.nih.gov/genome/guide/human/</a>	(Venter et al., 2001)
<i>Mus musculus</i>	<a href="http://www.ncbi.nlm.nih.gov/genome/guide/mouse/">http://www.ncbi.nlm.nih.gov/genome/guide/mouse/</a>	(Waterston et al., 2002)
<i>Rattus norvegicus</i>	<a href="http://www.ncbi.nlm.nih.gov/genome/guide/rat/">http://www.ncbi.nlm.nih.gov/genome/guide/rat/</a>	(Gibbs et al., 2004)
<i>Gallus gallus</i>	<a href="http://www.ncbi.nlm.nih.gov/genome/guide/chicken/">http://www.ncbi.nlm.nih.gov/genome/guide/chicken/</a>	(Hillier et al., 2004)
<i>Xenopus tropicalis</i>	<a href="http://genome.jgi-psf.org/Xentr4/Xentr4.home.html">http://genome.jgi-psf.org/Xentr4/Xentr4.home.html</a>	
<i>Fugu rubripes</i>	<a href="http://genome.jgi-psf.org/Takru4/Takru4.home.html">http://genome.jgi-psf.org/Takru4/Takru4.home.html</a>	(Aparicio et al., 2002)
<i>Tetraodon nigroviridis</i>	<a href="http://www.genoscope.cns.fr/externe/tetranew/">http://www.genoscope.cns.fr/externe/tetranew/</a>	(Jaillon et al., 2004)
<i>Danio rerio</i>	<a href="http://www.ensembl.org/Danio_rerio/index.html">http://www.ensembl.org/Danio_rerio/index.html</a>	(Birney et al., 2006)
<i>Petromyzon marinus</i>	<a href="http://pre.ensembl.org/Petromyzon_marinus/Info/Index">http://pre.ensembl.org/Petromyzon_marinus/Info/Index</a>	
<i>Branchiostoma floridae</i>	<a href="http://genome.jgi-psf.org/Brafl1/Brafl1.home.html">http://genome.jgi-psf.org/Brafl1/Brafl1.home.html</a>	(Putnam et al., 2008)
<i>Ciona intestinalis</i>	<a href="http://genome.jgi-psf.org/ciona4/ciona4.home.html">http://genome.jgi-psf.org/ciona4/ciona4.home.html</a>	(Dehal et al., 2002)
<i>Drosophila melanogaster</i>	<a href="http://www.fruitfly.org/">http://www.fruitfly.org/</a>	(Adams et al., 2000)
<i>Strongylocentrotus purpuratus</i>	<a href="http://www.hgsc.bcm.tmc.edu/projects/seaurchin/">http://www.hgsc.bcm.tmc.edu/projects/seaurchin/</a>	(Sodergren et al., 2006)
<i>Nematostella vectensis</i>	<a href="http://genome.jgi-psf.org/Nemve1/Nemve1.home.html">http://genome.jgi-psf.org/Nemve1/Nemve1.home.html</a>	(Putnam et al., 2007).

#### 3.2 Databases

**Table 4: Major databases used.**

Database	URL	References
NCBI	<a href="http://www.ncbi.nlm.nih.gov/">http://www.ncbi.nlm.nih.gov/</a>	(Wheeler et al., 2006)
RefSeq	<a href="http://www.ncbi.nlm.nih.gov/RefSeq/">http://www.ncbi.nlm.nih.gov/RefSeq/</a>	(Pruitt et al., 2005)
Entrez	<a href="http://www.ncbi.nlm.nih.gov/Entrez/">http://www.ncbi.nlm.nih.gov/Entrez/</a>	(Maglott et al., 2005)
Swissprot	<a href="http://www.expasy.org">www.expasy.org</a>	(Bairoch et al., 2004; Schneider et al., 2004)
UniProt	<a href="http://www.uniprot.org">http://www.uniprot.org</a>	(Apweiler et al., 2004b; Wu et al., 2006)
PROSITE	<a href="http://www.expasy.org/PROSITE/">http://www.expasy.org/PROSITE/</a>	(Hulo et al., 2006)
ENSEMBL	<a href="http://www.ensembl.org">www.ensembl.org</a>	(Birney et al., 2006; Hubbard et al., 2007)
The Serpin Database	<a href="http://www-structmed.cimr.cam.ac.uk/serpins.html">http://www-structmed.cimr.cam.ac.uk/serpins.html</a>	

##### 3.2.1 NCBI

The National Center for Biotechnology Information [NCBI] provides analysis and retrieval resources for the data in GenBank (Pruitt et al., 2003; Pruitt et al., 2005) and other biological data (Wheeler et al., 2005; Wheeler et al., 2006). There are many databases and tools from NCBI which are extensively used in this work including- Entrez, My NCBI, PubMed, PubMed Central, Entrez Gene, the NCBI Taxonomy Browser, BLAST, BLAST Link (BLink), Electronic PCR, OrfFinder, Spidey, Splign, RefSeq, UniGene, HomoloGene, ProtEST, Entrez Genome, Genome Project and related tools, the Trace and Assembly Archives, the Map Viewer, the Conserved Domain Database (CDD) and the Conserved

---

Domain Architecture Retrieval Tool (CDART). There are many other databases and tools available from NCBI that are not related to our work and are not mentioned above.

### 3.2.2 RefSeq

The Reference Sequence database (RefSeq)<sup>1</sup> is maintained and curated at the NCBI. It aims to provide a non-redundant collection of reference protein sequences (Pruitt *et al.*, 2003; Pruitt *et al.*, 2005). RefSeq sequences exist for several species (Pruitt *et al.*, 2003; Pruitt *et al.*, 2005) including genomes analyzed in this work (**Table 3**). The main features of the RefSeq collection include non-redundancy, explicitly linked nucleotide and protein sequences, updates to reflect current knowledge of sequence data and biology, data validation and format consistency. In November 2006, the database contained 3,000,705 entries with approximately 40 % manually reviewed entries (Apweiler *et al.*, 2004a).

### 3.2.3 Entrez

NCBI's Entrez Protein<sup>2</sup> is another exhaustive sequence repository (Wheeler *et al.*, 2005; Wheeler *et al.*, 2006). The database contains sequence data translated from the nucleotide sequences of the DNA Data Bank of Japan [DDBJ] (Tateno *et al.*, 1998), the European Molecular Biology Laboratory [EMBL] Nucleotide Sequence Database (Stoesser *et al.*, 1997), GenBank database (Benson *et al.*, 2005; Benson *et al.*, 2006), as well as sequences from SWISS-PROT (Bairoch and Apweiler, 1996; Bairoch and Apweiler, 2000), the Protein Information Resource [PIR] (Barker *et al.*, 1987), RefSeq (Pruitt *et al.*, 2003; Pruitt *et al.*, 2005) and the Protein Data Bank [PDB] (Berman *et al.*, 2000). The entries list additional information that can be extracted from curated databases such as SWISS-PROT and PIR. Sequence collection in the database is redundant (Apweiler *et al.*, 2004a).

### 3.2.4 SWISS-PROT

SWISS-PROT<sup>3</sup> (Bairoch and Boeckmann, 1991) is an annotated protein sequence database established in 1986 and maintained collaboratively, since 1988, by the Department of Medical Biochemistry of the University of Geneva and the EMBL Data Library. The SWISS-PROT protein knowledgebase (Boeckmann *et al.*, 2003) represents carefully curated amino acid sequences providing an interdisciplinary overview of relevant information by bringing together experimental results and computed features. The SWISS-PROT database distinguishes itself from other protein sequence databases by three distinct criteria (Bairoch and Apweiler, 2000): (a) annotation, (b) minimal redundancy in the sequence data and (c) integration of other 66 databases with cross-referencing facilities. In this work, SWISS-PROT was extensively used because of these features, which helped us in understanding about orthologs.

---

<sup>1</sup> <http://www.ncbi.nlm.nih.gov/RefSeq/>

<sup>2</sup> <http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?db=Protein>

<sup>3</sup> <http://expasy.org/sprot/>

### 3.2.5 UniProt

Universal Protein Resource (UniProt<sup>1</sup>) is a comprehensive catalog of information on proteins (Apweiler *et al.*, 2004b; Wu *et al.*, 2006). It joins the information contained in Swiss-Prot, TrEMBL, and PIR. UniProt is comprised of three components, each optimized for different uses (Apweiler *et al.*, 2004a). The UniProt Knowledgebase [UniProt] is the central access point for extensive curated protein information. The UniProt Non-redundant Reference [UniRef] databases combine closely related sequences into a single record to speed searches. The UniProt Archive [UniParc] is a comprehensive repository, reflecting the history of all protein sequences. Protein sequences are retrieved from predominant publicly accessible resources. All new and updated protein sequences are collected and loaded daily into UniParc for full coverage (Leinonen *et al.*, 2004).

### 3.2.6 PROSITE

PROSITE<sup>2</sup> is a database of protein families and domains defined on the basis of signatures (Bairoch, 1991). From a multiple sequence alignment, it is possible to derive a signature for a protein family or domain, which distinguishes its members from all other unrelated proteins. Biologically significant patterns and profiles are formulated in such a way that with appropriate computational tools it can help to determine to which family of proteins (if any) a new sequence belongs, or which known domains are found in the new sequence (Hulo *et al.*, 2004; Hulo *et al.*, 2006). These signature sequences are regular expressions in pure computational sense and can be easily searched in the protein sequences using Unix grep or using simple perl script for searching regular expressions of a specific length. ScanPROSITE is a new and improved version of the web-based tool for detecting PROSITE signature matches in protein sequences using ProRul (Henikoff and Henikoff, 1991; de Castro *et al.*, 2006).

### 3.2.7 ENSEMBL

ENSEMBL<sup>3</sup> is a comprehensive database in the area of chordate comparative genomics with coverage of 33 different genomes (Birney *et al.*, 2006; Hubbard *et al.*, 2007). It includes facilities for annotation, synteny, and automatic orthology assignment. It has an excellent genome browser with the ability of aligning different genomes at a time. In this work, we have been using ENSEMBL as a platform for comparing the serpins from different vertebrates and for building synteny around different serpins.

---

<sup>1</sup> Uniprot website, <http://www.uniprot.org>

<sup>2</sup> Prosite website, <http://expasy.org/prosite/>

<sup>3</sup> Ensembl website, <http://www.ensembl.org/>



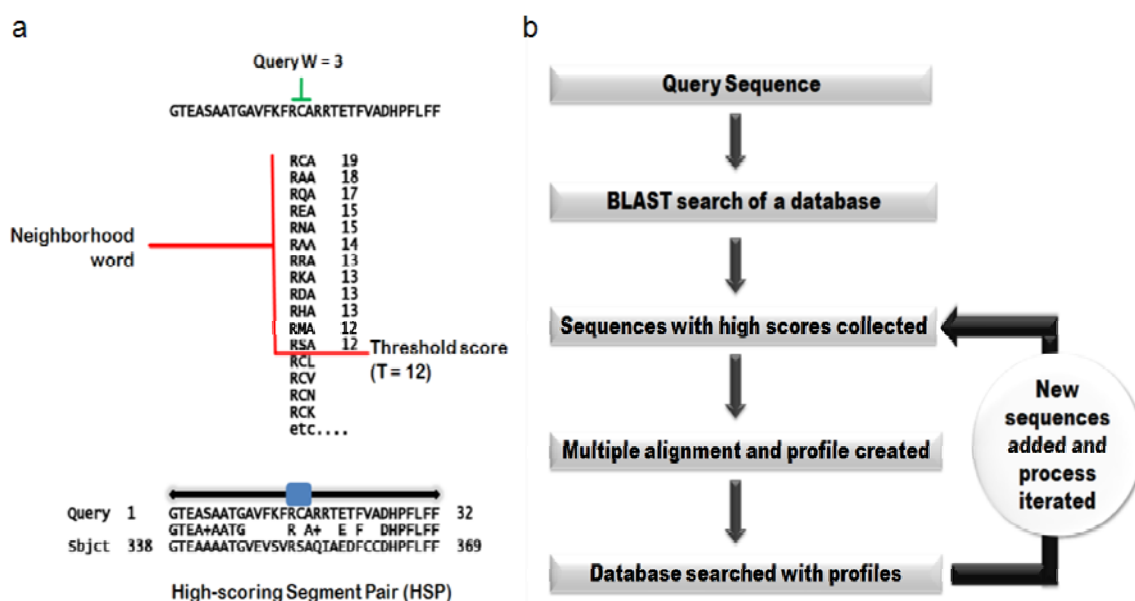
### 3.2.8 Serpin Database

The Serpin database<sup>1</sup> has information exclusively about serpins in terms of the sequences, known structures, and known mutations. This database is used in this work for gathering information about gene specific features.

## 3.3 Searching Tools

### 3.3.1 BLAST

BLAST [**B**asic **L**ocal **A**lignment **S**earch **T**ool] is a heuristic approach to find the highest scoring locally optimal alignments between a query sequence and sequences of a database (Altschul *et al.*, 1990). The overall approach of the BLAST algorithm is shown below with a random example of a reaction center loop region (RCL) in a serpin (**Figure 14a**). The heuristic search strategy of the BLAST is to find words of length  $W$  [e.g.,  $W = 3$  for proteins] that score at least  $T$  when aligned with the query and scored with a substitution matrix. The words in the database that score  $T$  or greater are extended in both directions in an attempt to find a locally optimal ungapped alignment called **high-scoring segment pair (HSP)** with a minimal score  $S$  or a minimal specified threshold  $E$ -value or a combination of the score  $S$  and  $E$ -value. The HSPs that meet these criteria are reported in BLAST output.



**Figure 14: BLAST algorithm.** (a) BLAST approach (Altschul *et al.*, 1990) for a reactive center loop (RCL) region of a randomly chosen serpin. (b) PSI-BLAST approach (Altschul *et al.*, 1997).

**Table 5** shows different versions of the BLAST approach which are used in this work.

<sup>1</sup> Serpin database website, <http://www-structmed.cimr.cam.ac.uk/serpins.html>

Table 5: Variants of BLAST suite.

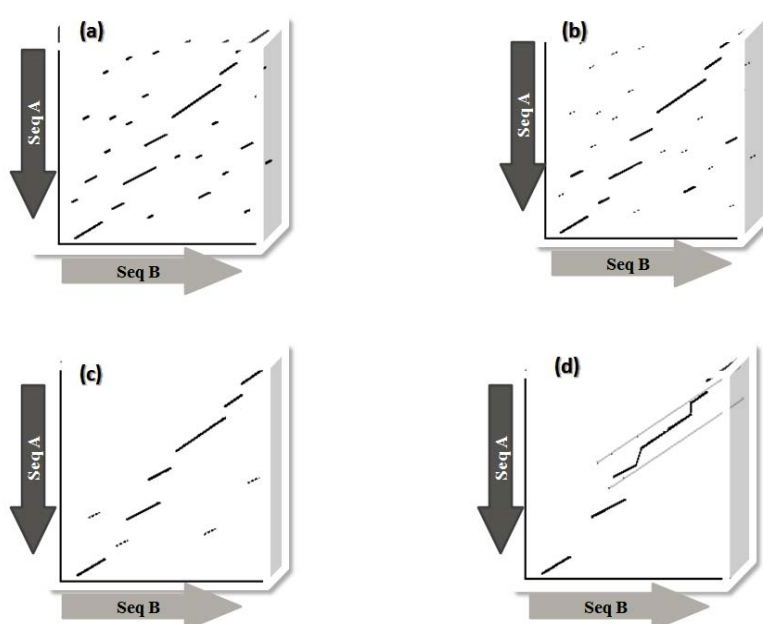
BLAST variant	Query sequences	Database
BLASTP	Protein	Protein
BLASTN	Nucleotide	Nucleotide
BLASTX	Translated nucleotide	Protein
TBLASTN	Protein	Translated nucleotide
TBLASTX	Translated nucleotide	Translated nucleotide
PSI-BLAST	Protein	Protein
MEGABLAST	Nucleotide	Nucleotide

### 3.3.2 PSI-BLAST

PSI-BLAST (**P**osition-**S**pecific **I**terated **B**LAST) was developed with three main goals - (a) speed, (b) simplicity and (c) automatic operation (Altschul *et al.*, 1997). The PSI-BLAST approach is summarized as follows (**Figure 14b**). The approach is basically a gapped BLAST of a protein sequence (Altschul *et al.*, 1997). From the gapped BLAST multiple alignments, the profiles of a length equals to the query length are created. From these profiles, the database is repeatedly searched until convergence. Since these steps are repeated or iterated and so, these steps are called as iteration [I]. In this work, we used PSI-BLAST with iteration  $I = 5$  because after 5 iterations, there was no significant change observed in the PSI-BLAST search. Unlike most profile-based search methods, PSI-BLAST runs as one program, starting with a single protein sequence and the intermediate steps of multiple alignment and profile construction are invisible to the user.

### 3.3.3 FASTA

Fasta compares one protein sequence to another protein sequence or to a protein database or a DNA sequence to another DNA sequence or a DNA library (Pearson and Lipman, 1988; Pearson, 1990). The algorithmic approach of FASTA (**Figure 15**) is a four-step process:



**Figure 15: Overview of FASTA algorithm.** The FASTA algorithm is four step process:

- Finding identities between two sequences A and B.
- Top scoring segments are selected based on a substitution matrix.
- Applying "joining threshold" to remove parts, which are not likely to be part of the alignment.
- Optimizing the alignment by joining top segments in a narrow band with help of dynamic programming.

- (a) Search: Finding identities between two sequences A and B (where, B = a sequence in a searching database).
- (b) Rescan: Top scoring segments are selected based on a rescanning using a substitution matrix. Now only the regions or segments of high density of identity are considered.
- (c) Join threshold: The “joining threshold” is applied to remove parts, which are not likely to be part of the alignment.
- (d) Optimization: The alignment is optimized by joining top segments in a narrow band with help of dynamic programming.

The FASTA variants used in this work are listed in **Table 6**.

**Table 6: Variants of FASTA.**

FASTA variant	Query sequences	Database
FASTP	Protein	Protein
FASTN	Nucleotide	Nucleotide
TFASTA	Protein	Translated nucleotide
FASTF	Protein Fragment	Protein
TFASTF	Protein Fragment	Translated nucleotide
FASTS	Protein Fragment	Protein
TFASTS	Protein Fragment	Translated nucleotide
FASTX	Translated nucleotide	Protein
FASTY	Translated nucleotide	Protein

### 3.3.4 Superfamily HMM library

The Superfamily HMM library was developed with the aim to provide structural and hence implied functional assignments to protein sequences at the superfamily level (Gough *et al.*, 2001). The online server and the software is available for local use from superfamily website<sup>1</sup> (Gough and Chothia, 2002). **Figure 16** shows the basic approach of a statistical model called HMM (**H**idden **M**arkov **M**odel) for sequence alignment (Krogh *et al.*, 1994) which is similar as that used in a sequence search of superfamily database.

## 3.4 Multiple sequence analysis tools

Multiple sequence alignments (MSA) of protein sequences are important in many applications, including phylogenetic tree estimation, structure prediction, and critical residue identification. About 30 different multiple sequence alignment tools are available which are used up till now as summarized in **Figure 17**. Traditionally, the most popular approach has been the progressive alignment method. A multiple alignment is built up gradually by aligning the closest sequences first and successively adding the more distant ones (Feng and Doolittle, 1987; Doolittle and Feng, 1990).

<sup>1</sup> Superfamily website, <http://supfam.mrc-lmb.cam.ac.uk/SUPERFAMILY/>

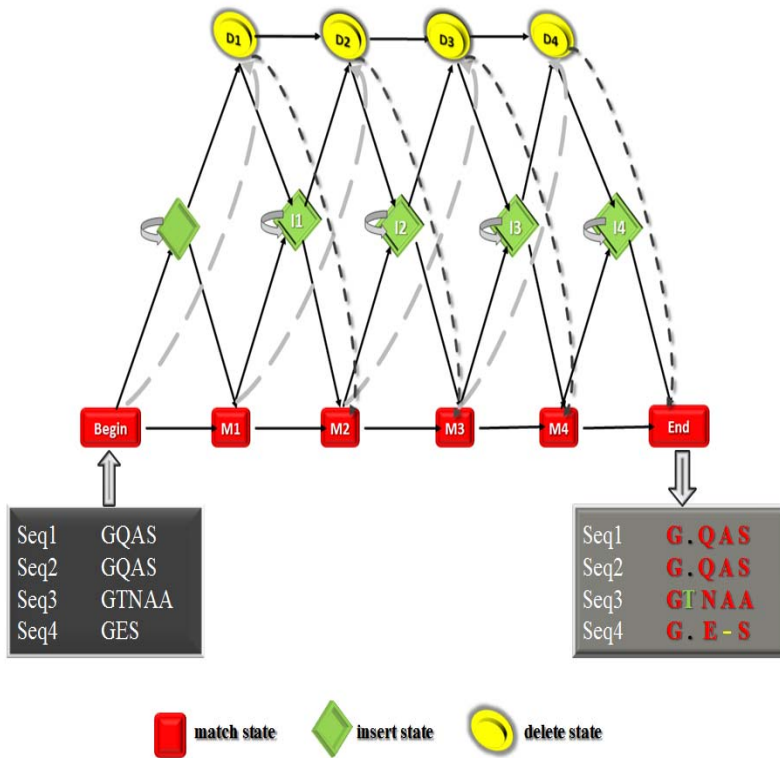


Figure 16: An example of the Hidden Markov Model for protein sequence alignment based on SAM (sequence alignment and modeling) (Krogh et al., 1994). The HMM consists of a series of states associated with the alignment probabilities. The match states are from begin to end, with in between M1-M4 (red squares), and are columns of the multiple sequence alignment. The “insert states” are insertions in the alignment (green diamonds). The delete states (yellow circles) are deletions or gaps marked D1-D4. Seq1 to Seq4 are input sequences, and the final output of the alignment of the same sequences is shown in the same colour as of all three states as match, insert and delete state, respectively.

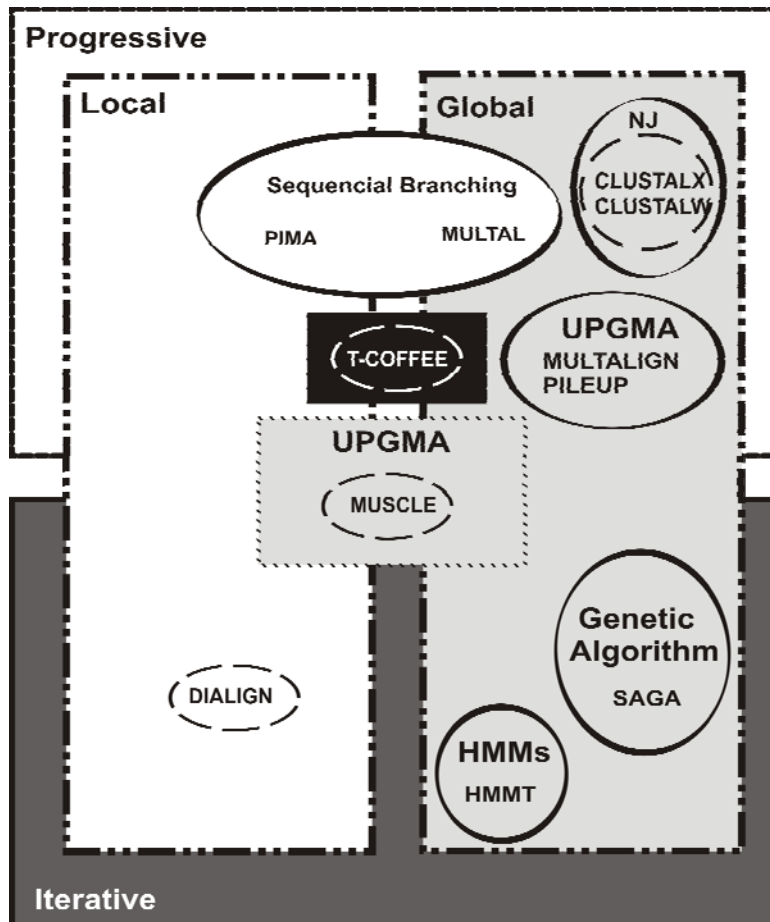


Figure 17: Summary of multiple sequence alignment algorithms (Thompson et al., 1999). The progressive alignment strategy is based on Feng and Doolittle (Feng and Doolittle, 1987; Doolittle and Feng, 1990) where first the closely related sequences are aligned and then the distant sequences. Re-aligning and improving is called iterative, used in DIALIGN (Morgenstern, 2000; Morgenstern, 2004). This figure is an updated modified of version published where T-Coffee (Notredame et al., 2000) and MUSCLE algorithms (Edgar, 2004a; Edgar, 2004b) are included and the broken oval circles indicates the tools used in this work. Alignment can be local or global.

There are many tools, which follow this approach, mainly differing in the method used to determine the order of alignment of the sequences. A common point of interest has been the application of iterative strategies to refine and improve the initial alignment and this is called iterative approach. **Table 7** summarizes the MSA tools used in this work.

**Table 7: Tools for multiple sequence alignment.**

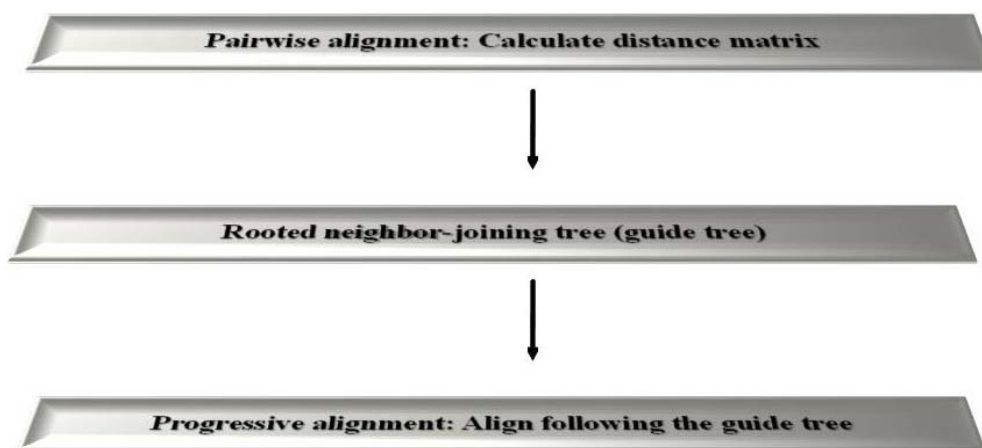
Tool	Approach	Sequence Type*	Alignment Type**	References
DIALIGN	Iterative alignment	PN	L	(Brudno <i>et al.</i> , 2004)
CLUSTALW	Progressive alignment	PN	L / G	(Thompson <i>et al.</i> , 1994)
MUSCLE	Progressive/iterative alignment	PN	L / G	(Edgar, 2004a; Edgar, 2004b)
T-COFFEE	More sensitive progressive alignment	PN	L / G	(Notredame <i>et al.</i> , 2000)
*Sequence type: protein(P) / nucleotide (N), both (PN)				
**Alignment Type: local(L) / global(G)				

### 3.4.1 CLUSTAL

CLUSTAL (CLUSTER ALIGNMENT) has been first developed in 1988 (Higgins and Sharp, 1988) and has been subsequently improved (Higgins *et al.*, 1996; Chenna *et al.*, 2003). CLUSTAL performs a global multiple alignment using following steps (**Figure 18**):

- (1) Perform pairwise alignment of all the sequences.
- (2) Use the alignment scores to produce the phylogenetic tree
- (3) Align the sequences sequentially, guided by the phylogenetic relationships indicated by the tree.

Thus, the most closely related sequences are aligned first, followed by additional sequences and groups of sequences are added guided by the initial alignments.



**Figure 18: Steps in CLUSTAL algorithm.**

---

CLUSTALW<sup>1</sup> (Thompson *et al.*, 1994) is the most recent version (where W stands for "weighing"), providing the ability of the program to provide weights to sequence and program parameters. The sensitivity of the CLUSTAL has been greatly improved for the alignment of divergent protein sequences using following steps with following four enhancement strategies:

- (a) Individual weights are assigned to each sequence in a partial alignment in order to down weigh near-duplicate sequences and up weigh the most divergent ones.
- (b) The amino acid substitution matrices are varied at different alignment stages according to the divergence of the sequences to be aligned.
- (c) The residue specific gap penalties and locally reduced gap penalties in hydrophilic regions encourage new gaps in potential loop regions rather than regular secondary structure.
- (d) The positions in early alignments where gaps have been opened receive locally reduced gap penalties to encourage the opening up of new gaps at these positions.

The graphical interface of CLUSTALW is called CLUSTALX (Thompson *et al.*, 1997). The CLUSTALX is a windows interface that makes it easy to use, provides an integrated system for performing multiple sequence, profile alignments, neighbor-joining tree building with bootstrapping facility and analyzing the results. A versatile sequence-coloring scheme allows the user to highlight conserved features in the alignment.

Overall, CLUSTALW and CLUSTALX are good tools for the multiple alignments performed in this work, because of the possibility to use diverse sequences and user-friendly graphical interfaces.

### 3.4.2 DIALIGN

DIALIGN<sup>2</sup> (DIagonal ALIGNment) is an automatic alignment tool that constructs pairwise and multiple alignments by comparing segment to segment of the sequences (Morgenstern, 1999; Morgenstern, 2004). DIALIGN's strength is in the comparison of sequences that share only local similarities.

### 3.4.3 MUSCLE

MUSCLE<sup>3</sup> (MULTI Sequence Comparison by Log-Expectation) is a multiple sequence alignment tool (Edgar, 2004a; Edgar, 2004b) that first makes a draft progressive alignment using a guided UPGMA (Sneath and Sokal, 1973) tree, further improvement using the another guided UPGMA (Sneath and Sokal, 1973) tree and finally refinement and re-alignment (**Figure 19**). MUSCLE provides a range of options that provide improved speed and / or alignment accuracy compared with CLUSTALW (Edgar, 2004a).

---

<sup>1</sup> ClustalW website <http://www.ebi.ac.uk/clustalw/>

<sup>2</sup> Dialign website, <http://bibiserv.techfak.uni-bielefeld.de/dialign/>

<sup>3</sup> <http://www.drive5.com/muscle/>

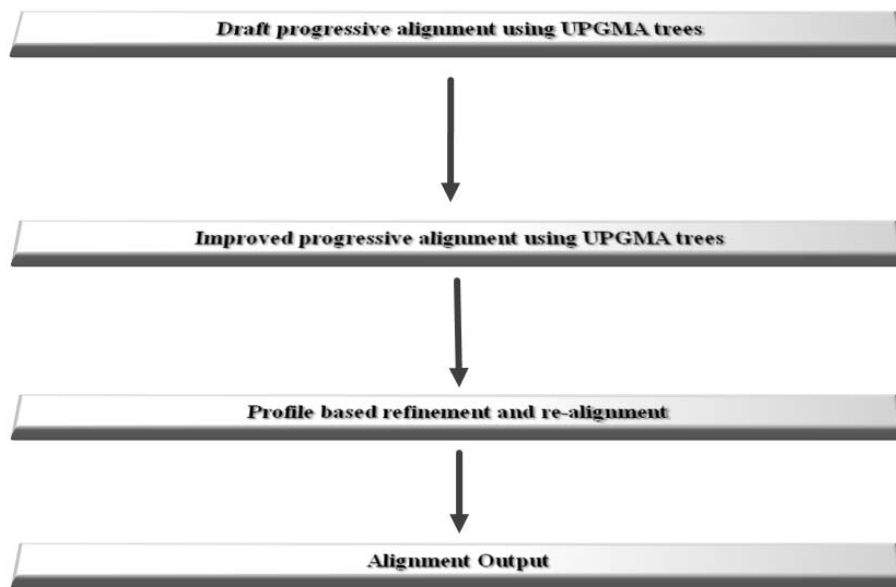


Figure 19: Algorithm of MUSCLE.

### 3.4.4 T-COFFEE

T-COFFEE<sup>1</sup> (Notredame *et al.*, 2000) is a sequence alignment package that allows the combination of a collection of multiple/pairwise, global or local alignments into a single model. It also enables estimation of the level of consistency of each position in the new alignment with the rest of the alignments (**Figure 20**). The strength of T-COFFEE is that it copes better with large gaps than CLUSTAL.

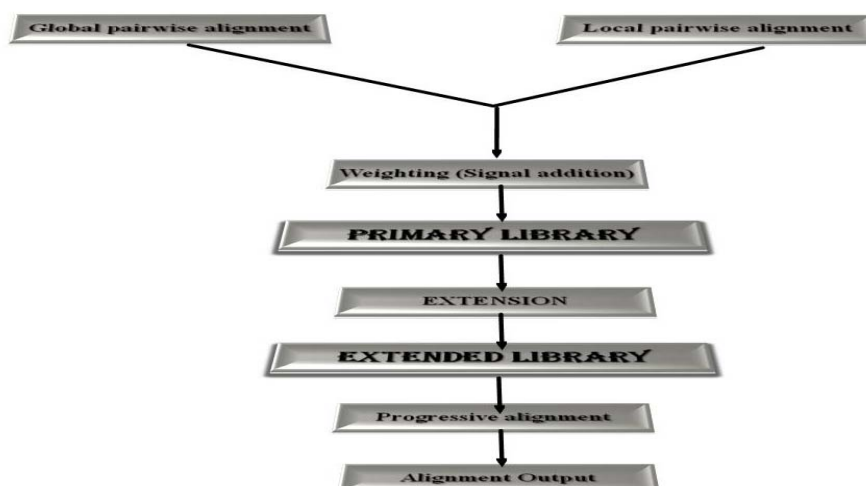


Figure 20: Algorithm of T-COFFEE.

<sup>1</sup> T-Coffee website, [http://www.igs.cnrs-mrs.fr/Tcoffee/tcoffee\\_cgi/index.cgi](http://www.igs.cnrs-mrs.fr/Tcoffee/tcoffee_cgi/index.cgi)

A new extension of T-Coffee called M-Coffee is a meta-method for assembling multiple sequence alignments (MSA) by combining the output of several individual methods into one single MSA (Wallace *et al.*, 2006).

### 3.5 Sequence editing tools

**Table 8** shows the multiple sequence alignment editing and representation tools used in this work.

**Table 8: Tools for multiple sequence alignment editing and representation.**

Tool	URL	Reference
GENEDOC	<a href="http://www.psc.edu/biomed/GENEDOC/">http://www.psc.edu/biomed/GENEDOC/</a>	(Nicholas <i>et al.</i> , 1997)
Jalview	<a href="http://www.jalview.org/">http://www.jalview.org/</a>	(Clamp <i>et al.</i> , 2004)
ESPrnt2.2	<a href="http://esprnt.ibcp.fr/ESPrnt/cgi-bin/ESPrnt.cgi">http://esprnt.ibcp.fr/ESPrnt/cgi-bin/ESPrnt.cgi</a>	(Gouet <i>et al.</i> , 1999)

### 3.6 Phylogenetic tools

There are many phylogenetic tools developed over the years in order to understand molecular evolution. The tools used in this work are summarized in **Table 9** and are described below.

**Table 9: Major phylogenetic tools.**

Tool	URL	Reference
MEGA3.1	<a href="http://www.MEGAsoftware.net/">http://www.MEGAsoftware.net/</a>	(Kumar <i>et al.</i> , 2001; Kumar <i>et al.</i> , 2004)
PHYLIP	<a href="http://evolution.genetics.washington.edu/PHYLIP.html">http://evolution.genetics.washington.edu/PHYLIP.html</a>	(Felsenstein, 1993; Felsenstein, 1996)
Phylodraw	<a href="http://pearl.cs.pusan.ac.kr/phylodraw/">http://pearl.cs.pusan.ac.kr/phylodraw/</a>	(Choi <i>et al.</i> , 2000)
Phylowin	<a href="http://pbil.univ-lyon1.fr/software/phylowin.html">http://pbil.univ-lyon1.fr/software/phylowin.html</a>	(Galtier <i>et al.</i> , 1996)
TREEVIEW	<a href="http://taxonomy.zoology.gla.ac.uk/rod/TREEVIEW.html">http://taxonomy.zoology.gla.ac.uk/rod/TREEVIEW.html</a>	(Page, 1996)
NJPLOT	<a href="http://pbil.univ-lyon1.fr/software/NJPLOT.html">http://pbil.univ-lyon1.fr/software/NJPLOT.html</a>	(Perriere and Gouy, 1996)

#### 3.6.1 MEGA 3.1

MEGA (*Molecular Evolutionary Genetics Analysis*) is a comprehensive tool for automatic and manual sequence alignment, inferring, editing and formatting phylogenetic trees, mining web-based databases, estimating rates of molecular evolution, and testing evolutionary hypotheses (Kumar *et al.*, 2001; Kumar *et al.*, 2004). MEGA3.1 is the most advanced version (Kumar *et al.*, 2004) which was extensively used in generating phylogenetic trees and editing for the purpose of the visualization of the trees in this work.

#### 3.6.2 PHYLIP

PHYLIP<sup>1</sup> (*PHYLogeny Inference Package*) contains programs for inferring phylogenies and is available for free (Felsenstein, 1993; Felsenstein, 1996). This package contains parsimony, distance matrix, and likelihood methods including bootstrapping, and consensus trees. It is a computational approach presented for minimizing the weighted sum of square of the

<sup>1</sup> Phylip website, <http://evolution.genetics.washington.edu/phylip/software.html>



differences between observed and expected pairwise distances between species, with the expectations are generated by an additive tree model. The method considers both Fitch and Margoliash criteria (Fitch and Margoliash 1967) along with Cavalli-Sforza and Edwards factor (Cavalli-Sforza and Edwards 1967). The parameters are weighted based on the least squares, with different weights. PHYLIP iterates lengths of adjacent branches in the tree three at a time. The weighted sum of squares never increases during the process of iteration, and in the iterative approach acquires a stationary point on the surface of the sum of squares. This approach makes it easy to maintain the constraint that branch lengths never become negative, although negative branch lengths can also be allowed. The PHYLIP approach is useful in studying the phylogenetic relationship among diverse sequences belonging to a particular family.

### 3.6.3 Phylodraw

Phylodraw<sup>1</sup> is a tree editor and manipulator (Choi *et al.*, 2000).

### 3.6.4 Geneious

Geneious<sup>2</sup> is a new bioinformatics tool which includes software for sequence analysis, phylogenetic methods, phylogenetic tree editing and literature mining (Drummond *et al.*, 2006).

### 3.6.5 Phylowin

Phylowin<sup>3</sup> is a graphical colour interface for molecular phylogenetic inference which performs neighbor-joining, parsimony and maximum likelihood methods and bootstrap (Galtier *et al.*, 1996).

### 3.6.6 TREEVIEW

TREEVIEW<sup>4</sup> (Page, 1996) is an useful tree editor which can read and manipulate many formats of the trees.

### 3.6.7 NJPLOT

NJPLOT<sup>5</sup> is a tree editor able to draw any phylogenetic tree. It allows zooming, branch swapping, display of bootstrap scores and printing in the PDF format (Perriere and Gouy, 1996).

---

<sup>1</sup> Phylodraw website, <http://pearl.cs.pusan.ac.kr/phylodraw/>

<sup>2</sup> Geneious website, <http://www.geneious.com/>

<sup>3</sup> Phylowin website, <http://pbil.univ-lyon1.fr/software/phylowin.html>

<sup>4</sup> Treeview website, <http://taxonomy.zoology.gla.ac.uk/rod/treeview.html>

<sup>5</sup> NJplot website, <http://pbil.univ-lyon1.fr/software/njplot.html>

### 3.7 Comparative genomics tools

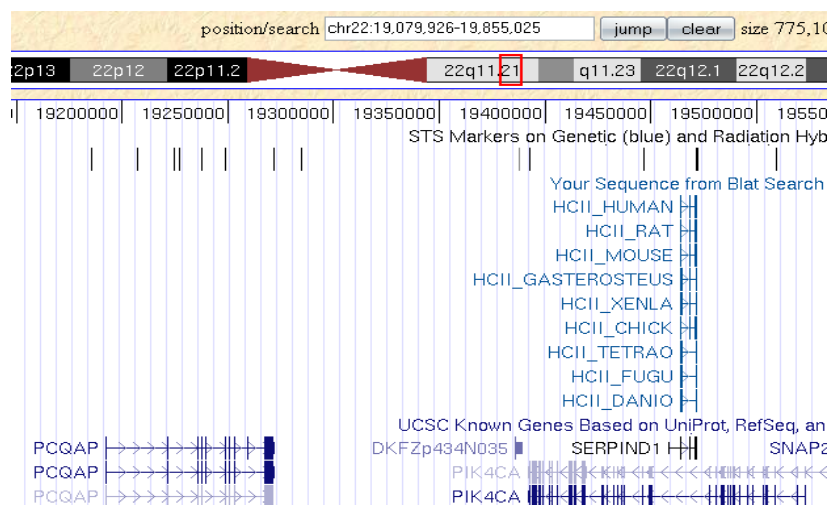
#### 3.7.1 Genome browsing tools

Genome browsers are tools for visualizing genomic regions on the genome. There are different genome browsers, which differ in way of visualization. **Table 10** lists different genome browsers, which we have used in this work.

**Table 10: Major genome browsers.**

Genome Browser	URL	Reference
ENSEMBL	www.ensembl.org	(Birney <i>et al.</i> , 2006; Hubbard <i>et al.</i> , 2007)
UCSC genome browser	http://genome.ucsc.edu/cgi-bin/hgGateway	(Kent <i>et al.</i> , 2002)
NCBI mapviewer	http://www.ncbi.nlm.nih.gov/mapview/	(Pruitt <i>et al.</i> , 2005)
JGI Fugu Genome Browser	http://genome.jgi-psf.org/cgi-bin/browserLoad/455749f979df66204c138bf7	
JGI Xenopus Genome Browser	http://genome.jgi-psf.org/cgi-bin/browserLoad/45574b9a2a40b9db41abc15f	
JGI Ciona Genome Browser	http://genome.jgi-psf.org/cgi-bin/browserLoad/45574d7a16b35bd737c994e2	
Tetraodon Genome Browser	http://www.genoscope.cns.fr/externe/tetranew/	

**Figures 21 to 23** show examples of the applications of genome browsers.



**Figure 21: The UCSC genome browser as seen in Firefox 2.0 web browser.** The structure of the heparin cofactor II gene on human chromosome 22 is illustrated.

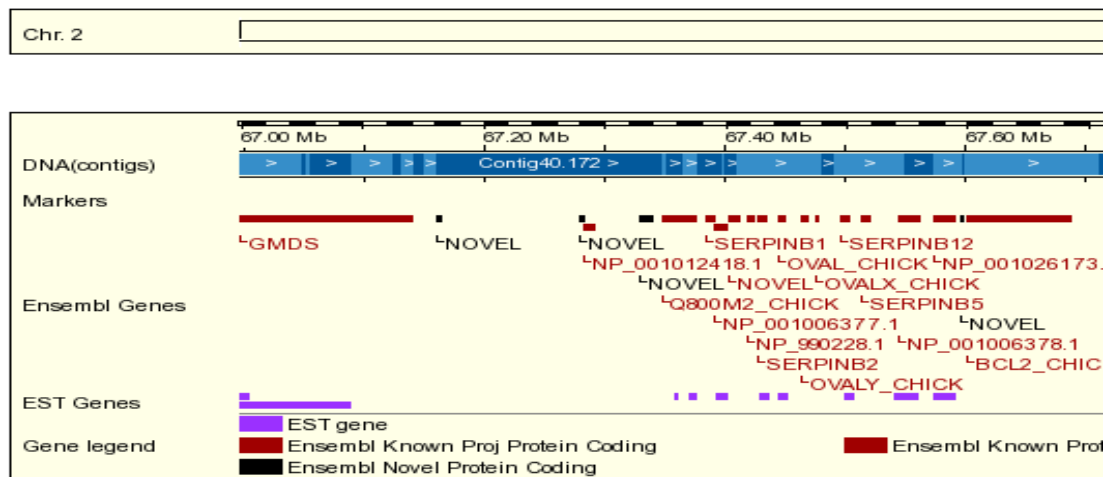


Figure 22: The Ensembl browser as seen in Firefox 2.0 web browser. The genomic organization of clade B serpins on chicken chromosome 2 is illustrated.

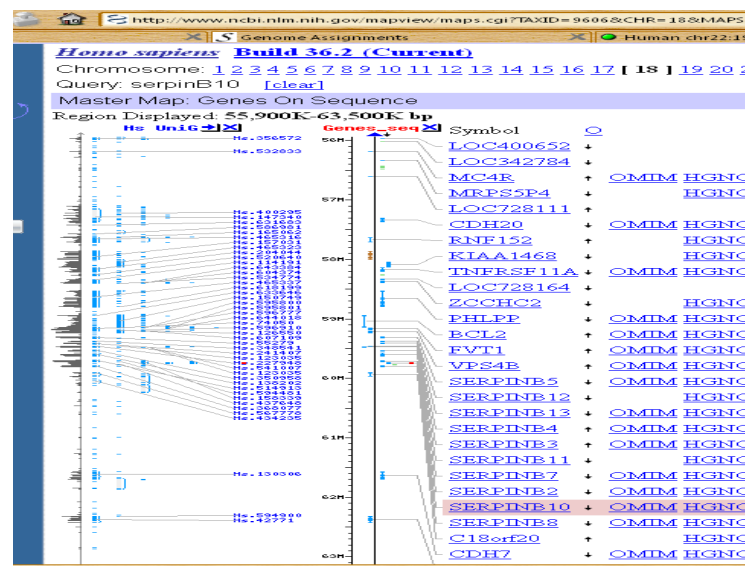


Figure 23: The NCBI mapviewer as seen in Firefox 2.0 web browser. The genomic organization of clade B serpins on human chromosome 18 is illustrated.

### 3.7.2 GENLIGHT

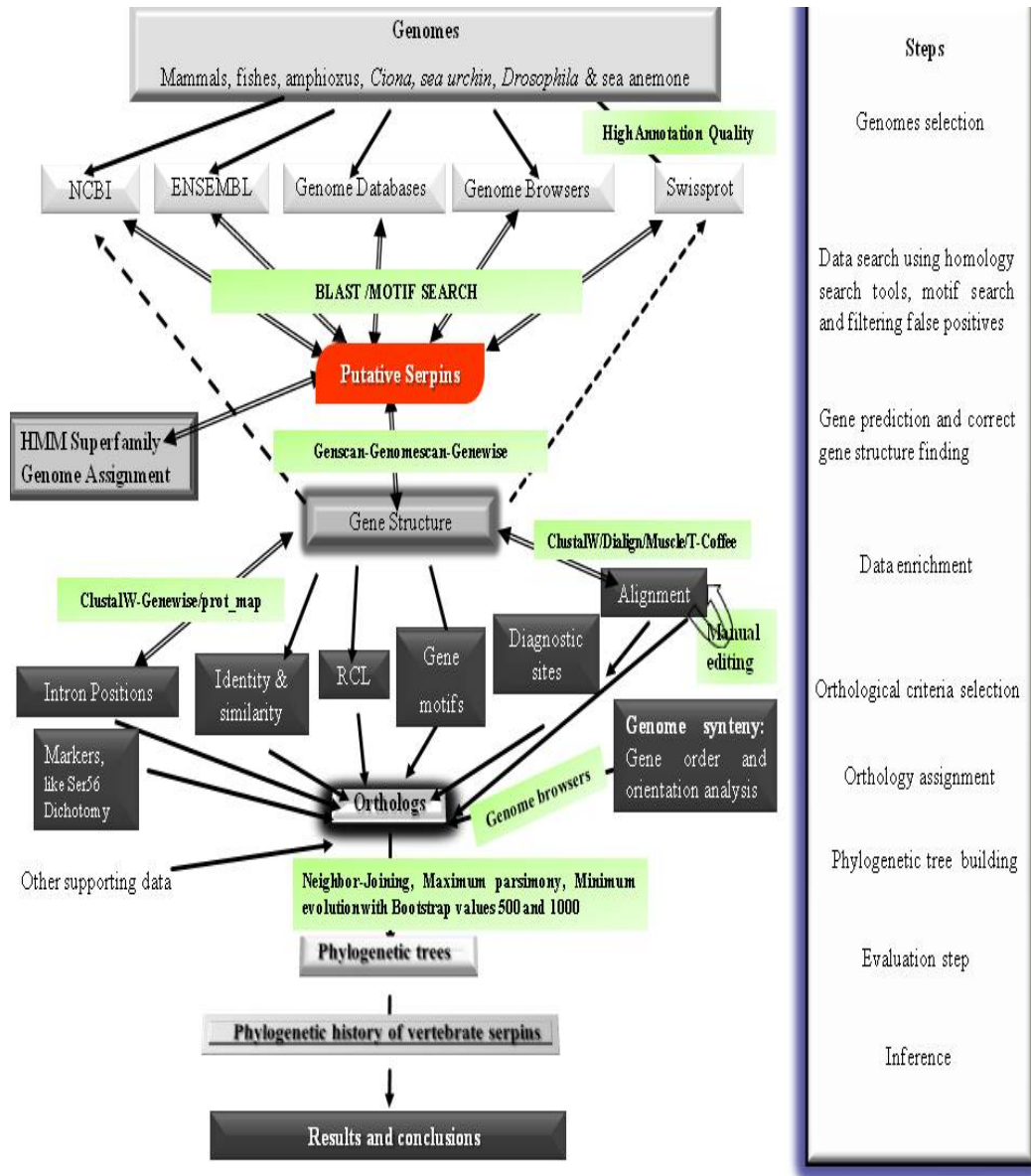
GENLIGHT is an extremely versatile system for comparative genomics and differential sequence analysis (Beckstette, 2004). It is a Client/Server based program suite build on the object relational database system on for large scale sequence analysis and comparative genomics and supports the management of nucleotide sequences as well as protein sequences. The assessment methods are complemented by a large variety of visualization methods for the evaluation of the results (Beckstette *et al.*, 2004). During this work, GENLIGHT was used extensively for benchmarking the sequence analysis using different filter parameters. The major usage of GENLIGHT in this work included using for quick BLAST and FASTA

---

searches, automatic orthology assignment using two way BLAST approach and for finding novel introns in different genomes using FASTA and SSHA searches provided in the GENLIGHT (Beckstette *et al.*, 2004).

## 4. Methods

An overview on the methods used in this work is shown in **Figure 24**.



**Figure 24: Protocol for the phylogenetic study of vertebrate serpins.** We selected different animal genomes and searched serpin genes from these genomes. The data was cross-validated from NCBI, ENSEMBL, SUPERFAMILY, SWISSPROT and UCSC genome browsers. The probable genomic sequences carrying serpin genes were collected and the gene structures were predicted using GENSCAN, GENOMESCAN and GENEWISE. Introns positions were mapped using GENEWISE or SoftBerry's PROT\_MAP tool and they were aligned with mature  $\alpha_1$ -antitrypsin using CLUSTALW. The alignments were edited with help of GENEDOC alignment editing software. Based on protein alignment percentage identity, reactive center loop characteristics, gene-specific motifs, diagnostic sites for group specific features of serpins were implemented and analyzed. The synteny amongst different genomes was build using different genome browsers. The phylogenetic trees were built. The final conclusions were made based on these trees and orthology assignment. The boxes show the steps used in this work in a generalised form. Double headed arrows indicate data verification stages while the backward and broken arrow indicate multiple repeating steps.

## 4.1 Database searching and evaluations

### 4.1.1 Searches with BLAST

The BLAST searches were performed with each of the genomes under consideration using (a) BLASTP with default expect value (E-value) 10 and word size 3 and repeated successively for E-value 0.01 and 0.001

(b) PSI-BLAST with E-value 0.01 and 0.001, word size 3 and number of iteration 5.

The different E-values were used to make the search more stringent. The use of multiple rounds of BLAST helped in avoiding false positives. During all BLASTP and PSI-BLAST searches,  $\alpha_1$ -antitrypsin was used as a standard protein for search. In some cases, a different serpin was used as standard serpin. For instance, for the search of group V1 serpins in a genome, MNEI was used. These searches were made both using standalone BLAST as well GENLIGHT incorporated BLAST (Beckstette *et al.*, 2004) where one can store and benchmark the data.

### 4.1.2 Searches based on motifs

Many serpins may be recognized by the presence of signature sequences. There are the following types of serpin signatures:

a) One type is based on serpin sequence alignments (Irving *et al.*, 2000; Ragg *et al.*, 2001).

The deduced serpin signature is shown in **Table 11**.

**Table 11: Locations of signature sequences in a typical serpin.** The positions refer to human  $\alpha_1$ -antitrypsin. The search patterns (amino acid codes) are shown in the AGREP-Notation: "." stand for a wildcard symbol; "#" stands for a variable number of positions.

Location	Position	Signature
s3A-breach-s4C	186 - 208	[NS]..[HYF]F[KR][GA].W...F...T...F
s5A-s4A	334 - 351	[HQ][KR]A...[DN][DE][DE]G[TS]EAA..[TS]
s1C-turn-s4B	364 - 386	F..[DN][HRK]PF.[FLV]#F.G

This signature is spread across three major regions of serpins; namely (i) region - s3A-breach-s4C; (ii)  $\beta$ -sheets s5A, s4A and its hinge region and (iii) between  $\beta$ -sheets s1C, s4B and their turn region - s1C-turn-s4B (**Figure 1**).

b) The PROSITE serpin signature: PROSITE<sup>1</sup> is a database of signature and profiles (Hulo *et al.*, 2006). The PROSITE serpin signature (id = PS00284) and is 11 amino acid long, as shown below:

[LIVMFY]-{G}-[LIVMFYAC]-[DNQ]-[RKHQS]-[PST]-F-[LIVMFY]-[LIVMFYC]-x-[LIVMFAH]

where,

[LIVMFY] = any one of the amino acid enclosed in the bracket [].

{G} = any amino acid but not G.

<sup>1</sup> Prosite website, <http://expasy.org/prosite/>

x = a position where any amino acid is accepted

The signatures were searched in newly identified sequences using following searching tools: AGREP 3.37 (Wu and Mander, 1992), DNA2AA with AGREP format (Krueger, 2003).

## 4.2 Sequence alignment

Protein alignments were generated with CLUSTALW/CLUSTALX 1.83 (Higgins *et al.*, 1996; Chenna *et al.*, 2003), or DALIGN 2.2.1 (Morgenstern, 1999; Morgenstern, 2004), or T-COFFEE (Notredame *et al.*, 2000), or MUSCLE (Edgar, 2004a; Edgar, 2004b) or a combinations of all these tools. The alignments were then visualized and edited using GENEDOC (Nicholas *et al.*, 1997) as described in **appendix 8.2**.

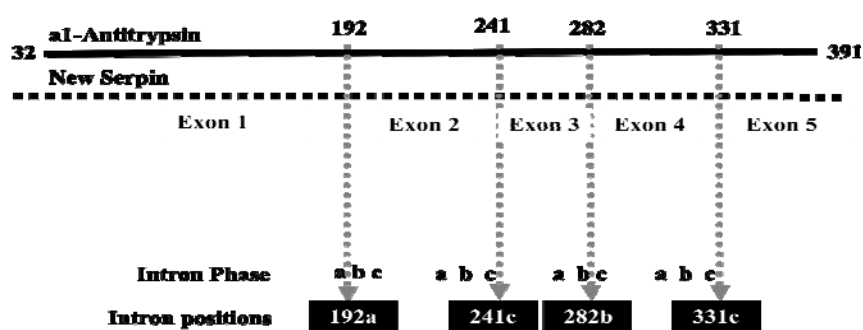
## 4.3 Gene structure analysis

### 4.3.1 Gene structure prediction

Gene structure prediction was done using GENSCAN<sup>1</sup> (Burge and Karlin, 1997; Burge and Karlin, 1998), GENOMESCAN<sup>2</sup> (Burge and Karlin, 1997; Burge and Karlin, 1998) and GENEWISE<sup>3</sup> (Birney *et al.*, 2004) or using a combination of all.

### 4.3.2 Mapping of intron positions

Mature human  $\alpha_1$ -antitrypsin was used as standard sequence for mapping of intron positions. Intron-exon structures were determined with the aid of GENEWISE (Birney *et al.*, 2004) and/or PROT\_MAP (softberry software<sup>4</sup>). The pairwise alignment of mature human  $\alpha_1$ -antitrypsin and putative serpins was created using CLUSTALW (Higgins and Sharp, 1988; Thompson *et al.*, 1994; Higgins *et al.*, 1996). Intron positions were marked semi-automatically, with manual inspection (**Figure 25**).



**Figure 25: Generalised scheme of intron mapping.** Intron positions are marked with respect to the conserved part (amino acids 32-391) of  $\alpha_1$ -antitrypsin. The intron position 192a means that the intron maps to amino acid 192 of the mature  $\alpha_1$ -antitrypsin and then its phasing is after the first base of the codon specifying amino acid 192.

<sup>1</sup>GENSCAN website, <http://genes.mit.edu/GENSCAN.html>

<sup>2</sup>GenomeScan website, <http://genes.mit.edu/genomescan.html>

<sup>3</sup>Wise2 website, <http://www.ebi.ac.uk/Wise2/>

<sup>4</sup>Softberry, <http://www.softberry.com/berry.phtml>

---

In some cases, there were novel introns inserted in the genes, such introns were analyzed with the aid of FASTA and SSHA searches incorporated in GENLIGHT (Beckstette *et al.*, 2004).

#### 4.4 Gene specific features

The alignment of serpins was built using CLUSTALW (Higgins *et al.*, 1996; Chenna *et al.*, 2003). Gene specific information was gathered from inspection of various publications and / or gathered from the serpin database<sup>1</sup>. Data were incorporated with help of the GENEDOC editor (Nicholas *et al.*, 1997).

#### 4.5 Orthology assignment

##### 4.5.1 Sequence identity and sequence similarity values from protein alignments

From gene specific sequence alignment the percentage of sequence identity and the percentage of sequence similarity values were calculated. This was one of the parameters used for orthology assignment.

##### 4.5.2 Group specific diagnostic sites

The group specific diagnostic sites were marked on the genes as described earlier for mammalian serpins (Ragg *et al.*, 2001).

##### 4.5.3 Rare indels

Rare indels were marked on the genes as described for mammalian serpins (Ragg *et al.*, 2001).

#### 4.6 Synteny analysis

##### 4.6.1 Synteny analysis of group V1 serpins

The group V1 synteny maps were built using following steps:

- i. Using NCBI mapviewer<sup>2</sup>, the human genome was scanned for the presence of group V1 serpins by zooming in and out. The genomic organization, location, and orientations of all group V1 serpins were marked. Some other genes were also marked as reference at the boundaries on the both sides of the clusters.
- ii. The genomic organization of serpins from the mouse and rat genomes was also built up with use of the corresponding NCBI mapviewer. Conserved markers were considered with respect to human genome.
- iii. The chicken genome was scanned for location and orientation of group V1 genes and for marker genes using the NCBI mapviewer. This analysis was repeated using the UCSC genome browser to confirm the accuracy of conservation.
- iv. The *Xenopus tropicalis* genome was scanned for group V1 genes using JGI Xenopus genome browser. The experiment was repeated using the ENSEMBL and the UCSC genome

---

<sup>1</sup> Serpin database website, <http://www-structmed.cimr.cam.ac.uk/serpins.html>

<sup>2</sup> NCBI Map Viewer, <http://www.ncbi.nih.gov/mapview>



browsers. The use of multiple genome browsers aided in assigning proper gene location and orientation.

- v. The *Fugu* genome was analyzed for group V1 genes using the JGI *Fugu* genome browser. The analysis was repeated using the ENSEMBL and the UCSC genome browsers. Group V1 serpin genes in the *Fugu* genome were found to be scattered on different scaffolds. The scaffolds were compared with human, chicken and frog.
- vi. The *Danio* genome was also scanned for group V1 genes and the marker genes, which formed the boundary of the cluster, using ENSEMBL as well as UCSC genome browsers.
- vii. The *Tetraodon* genome was searched for the group V1 genes using *Tetraodon* genome browser, the ENSEMBL and the UCSC genome browsers.

The tentative orthology of the marker genes was confirmed by bi-directional BLAST approach using the NR (**non-redundant**) database from NCBI. This step was considered because (i) this provided a confirmation of genes that are really conserved and (ii) there was no wrong annotation in NCBI mapviewer. Finally, the clusters from all vertebrates were compared with each other. The multiple genome synteny maps were repeatedly built using the ENSEMBL genome browser and the UCSC genome browser. This step was useful in resolving the problematic cases that arose during genome specific cluster building.

**Table 12: List of genomes and corresponding genome browsers used in building synteny maps of different serpins.**

Genomes	Genome browsers
<i>Homo sapiens</i>	NCBI mapviewer
<i>Mus musculus</i>	NCBI mapviewer
<i>Rattus norvegicus</i>	NCBI mapviewer
<i>Gallus gallus</i>	NCBI mapviewer, UCSC genome browser and ENSEMBL genome browser
<i>Xenopus tropicalis</i>	JGI <i>Xenopus</i> genome browser, ENSEMBL genome browser and UCSC genome browser.
<i>Fugu rubripes</i>	JGI <i>Xenopus</i> genome browser, ENSEMBL genome browser and UCSC genome browser
<i>Tetraodon nigroviridis</i>	<i>Tetraodon</i> genome browser, ENSEMBL genome browser and UCSC genome browser
<i>Danio rerio</i>	ENSEMBL genome browser and UCSC genome browser

#### 4.6.2 Group V2 serpin synteny analysis

Synteny of group V2 serpins was built as described in **section 4.6.1** using the genome browsers summarized in **Table 12**. There are three genes in the group V2 serpin genes, which are not located in the common cluster of the group V2 serpin gene, namely heparin cofactor II genes, serpinA7 genes and angiotensinogen genes. The synteny for these three genes was built using similar strategies as described in **section 4.6.1**.

#### 4.6.3 Synteny analysis of serpin groups V3-V6

The synteny of serpin groups V3-V6 was determined using the different genome browsers listed in **Table 12** are described in the **section 4.6.1**.

## 4.7 Analysis of the *Ciona intestinalis* genome

### 4.7.1 Searching serpins in the *Ciona intestinalis* genome

The *Ciona intestinalis* genome (version v1<sup>1</sup>) was searched for serpins with following homology search tools, using human  $\alpha_1$ -antitrypsin as standard:

- a) Using BLAST variants (Altschul and Lipman, 1990; Altschul *et al.*, 1997) as described in **section 4.1.1**.
- b) Using the Superfamily HMM library<sup>2</sup> (Gough and Chothia, 2002) search with default settings.

### 4.7.2 Determination of the exon-intron structure

Using human  $\alpha_1$ -antitrypsin as reference sequence, the exon-intron structures of all *Ciona* serpin genes were determined with help of GENEWISE and/or PROT\_MAP. The pairwise alignment of *Ciona* serpins and human  $\alpha_1$ -antitrypsin was created using CLUSTALW (Higgins *et al.*, 1996; Chenna *et al.*, 2003). Intron positions and phasing were assessed as described in **section 4.3.2** and **Figure 25**. The exon-intron structures were manually checked for accuracy.

### 4.7.3 Synteny Building

The *Ciona* serpins were collected on the scaffolds using the *Ciona* genome browser and the location of the serpin genes and their orientation were determined. The location and orientation of neighboring genes of *Ciona* serpin was marked. The data was then compared with that of human and fish serpin genes.

### 4.7.4 Analysis of serine codon dichotomy at position 56

The *Ciona* serpins were analyzed for serine codon dichotomy at position 56<sup>3</sup> (Krem and Di Cera, 2003) using CLUSTALW alignment (Higgins *et al.*, 1996; Chenna *et al.*, 2003) as well as the GENEWISE protein-nucleotide alignment (Birney *et al.*, 2004).

## 4.8 Analysis of the *Branchiostoma floridae* genome

The *Branchiostoma floridae* genome was analyzed for presence and characterization of serpin genes in a similar fashion as described in **section 4.7**, using *B. floridae* genome database<sup>4</sup>.

## 4.9 Analysis of the *Strongylocentrotus purpuratus* genome

The *Strongylocentrotus purpuratus* genome was analyzed for presence and characterization of

<sup>1</sup> *Ciona* genome v1 website, <http://genome.jgi-psf.org/ciona4/ciona4.home.html>

<sup>2</sup> Superfamily website, <http://supfam.mrc-lmb.cam.ac.uk/SUPERFAMILY/>

<sup>3</sup> The numbering in mature region of human  $\alpha_1$ -antitrypsin.

<sup>4</sup> *B. floridae* genome database website, <http://genome.jgi-psf.org/Braf11/Braf11.home.html>

serpin genes in a similar fashion as described in **section 4.7**, using *S. purpuratus* genome database<sup>1</sup> and NCBI sea urchin genome resources<sup>2</sup>.

#### 4.10 Analysis of the *Nematostella vectensis* genome

The *Nematostella vectensis* genome was analyzed for presence and characterization of serpin genes in a similar fashion as described in **section 4.7**, using *N. vectensis* genome database<sup>3</sup>.

#### 4.11 Phylogenetic analysis and bootstrap analysis

A phylogenetic tree is a two dimensional graph composed of branches and nodes which show evolutionary relationships between genes/proteins. Only one branch (or an edge) connects any two nodes. The nodes represent the taxonomic units called as taxa; the node is the intersection or terminating point of two or more branches. For instance, DNA/protein sequences are considered as taxons. An OTU (**O**perational **T**axonomic **U**nit) is an extant taxon present at an external node, or leaf: the OTUs are the available nucleic acid or protein sequences. There are two principal methods of making trees: character-based methods and distance-based methods. **Table 13** summarizes the different phylogenetic methods.

**Table 13: Summary of the phylogenetic methods.**

Phylogenetic methods	Summary
<b>Character based methods</b>	
<b>Maximum Likelihood (ML)</b>	The most likely output tree, given a probabilistic model of evolutionary changes in DNA or protein sequences.
<b>Maximum parsimony (MP)</b>	The minimum number of evolutionary steps required to generate the observed variation in a set of sequences, as found by comparison of the number of steps in all possible phylogenetic tree.
<b>Distance based methods</b>	
<b>Neighbor joining (NJ)</b>	Heuristic search algorithm that finds a minimum evolution tree from the distance between each pair of taxa in the tree (Saitou and Nei, 1987).
<b>Unweighted pair group method with arithmetic mean (UPGMA)</b>	A simple method for tree construction that assumes the rate of change along the branches of the tree is a constant and the distances are approximately ultrametric (Sneath and Sokal, 1973).

Bootstrapping is a statistical method for testing how well a particular data set fits a model. For instance, a sequence may be left out of an analysis to determine how much the sequence influences the results of that analysis. The phylogenetic analysis of serpins was done using following methods: Maximum parsimony (MP), Neighbor joining (NJ), and UPGMA [see Table 5] with bootstrapping values 500 and 1000, respectively, with the help of MEGA 3.1 (Kumar *et al.*, 2004) and/or PHYLIP<sup>4</sup> (Felsenstein, 1993; Felsenstein, 1996). Phylogenetic

<sup>1</sup> *S. purpuratus* genome database website, <http://www.hgsc.bcm.tmc.edu/projects/seaurchin/>

<sup>2</sup> NCBI sea urchin genome resources website, [http://www.ncbi.nlm.nih.gov/genome/guide/sea\\_urchin/](http://www.ncbi.nlm.nih.gov/genome/guide/sea_urchin/)

<sup>3</sup> *N. vectensis* genome database website, <http://genome.jgi-psf.org/Nemve1/Nemve1.home.html>

<sup>4</sup> Phylip website, <http://evolution.genetics.washington.edu/phylip/software.html>

---

trees were built at gene level, group level, within a genome, or sets of genomes. Different phylogenetic trees were built with serpins from selected genomes (**Table 3**). The original trees were extensively edited for visualisation using TREEVIEW<sup>1</sup> (Page, 1996), NJPLOT<sup>2</sup> (Perriere and Gouy, 1996), GENEIOUS tree editor<sup>3</sup> and MEGA 3.1 tree editor<sup>4</sup> (Kumar *et al.*, 2004). The individual cases are explained in the results section when visualising a tree.

---

<sup>1</sup> Treeview website, <http://taxonomy.zoology.gla.ac.uk/rod/treeview.html>

<sup>2</sup> NJplot website, <http://pbil.univ-lyon1.fr/software/njplot.html>

<sup>3</sup> Geneious software, <http://www.geneious.com/>

<sup>4</sup> Mega3.1 website, <http://www.megasoftware.net/>

## 5. Results

### 5.1. *Gallus gallus* and its serpins

The chicken is an important model organism for agriculture, biomedical research, developmental and aging research. Birds have been evolved separately from mammals for about 310 Mya (Hedges, 2002; Reisz and Muller, 2004). The chicken genomic sequence is released as draft version with the 6.6X coverage. The size of the chicken genome is about one third of mammalian genomes with lesser intergenic repeats, pseudogenes, and segmental duplications (Bourque *et al.*, 2005). Several rounds of homology searches revealed 27 serpin genes in the chicken genome, as listed in **Table 14**. These serpins are further characterized in **section 5.11. to 5.16.**

**Table 14: List of serpins of *Gallus gallus*.**

Gene Name	Accession Id	Clade	Protein length	Homology to known serpin in NR database
Gga-Spn-1	XP_418980 (LOC4)	B	378	MNE1
Gga-Spn-2	XP_418981	B	379	SPB6
Gga-Spn-3	XP_426040	B	378	BOMAPIN
Gga-Spn-4	NP_990228	B	410	MENT
Gga-Spn-5	XP_418982	B	412	SERPINB2
Gga-Spn-6	NP_990483	B	386	Ovalbumin
Gga-Spn-7	XP_418983	B	388	Gene Y protein
Gga-Spn-8	XP_418984	B	388	Gene X protein
Gga-Spn-9	XP_418985	B	422	SPB12
Gga-Spn-10	XP_418986	B	375	Maspin
Gga-Spn-11	XP_426460	A	374	$\alpha_1$ -AT
Gga-Spn-12	XM_001235489	A	419	$\alpha_1$ -AT
Gga-Spn-13	XP_421345	A	432	$\alpha_1$ -AT
Gga-Spn-14	XP_421344	A	437	$\alpha_1$ -AT
Gga-Spn-15	XP_421343	A	425	$\alpha_1$ -AT
Gga-Spn-16	XP_421342	A	425	$\alpha_1$ -AT
Gga-Spn-17	XP_421341	A	439	ZPI
Gga-Spn-18	XP_419584	A	464	AGT
Gga-Spn-19	AAC16324	D	489	HCII
Gga-Spn-20	gi:50730899	E	395	GDN
Gga-Spn-21	XM_417070	E	396	SERPINE3
Gga-Spn-22	gi:521387191	I	410	NEURO
Gga-Spn-23	gi:50758202	F	423	PEDF
Gga-Spn-24	XP_415807.2	F	514	A2AP
Gga-Spn-25	gi:50747972	G	448	C1IN
Gga-Spn-26	XP_422282	C	453	ATIII
Gga-Spn-27	gi:45384240	H	405	HSP47

## 5.2. *Xenopus tropicalis* and its serpins

The frog *Xenopus tropicalis* has the smallest genome of all known amphibians. It is a connecting link between mammals and fish. Therefore, it is of phylogenetic interest, in addition to its importance in early embryonic development and cell biology.



**Figure 26:** *Xenopus tropicalis*. This picture is taken from NCBI Xenopus genome resources<sup>1</sup>.

The *Xenopus tropicalis* genome assembly (release v4.1) was assembled using JAZZ, the JGI assembler, indicating a genome of approximately 1.5 Gb. The assembly contains 19,501 scaffolds with an average coverage of 7.65 fold. Roughly half of the genome is contained in 272 scaffolds, all at least 1.56 Mb in length. Gene models and associated transcripts/proteins are predicted or mapped using a variety of tools based on cDNA, protein homology and *ab initio* methods. The current release contains approximately 28,000 gene models, supported by EST and cDNA data of both *X. tropicalis* and the closely related species *X. laevis*. Homology searches detected 25 serpins in *X. tropicalis* genome, listed in **Table 15**. These serpins are further characterized in **section 5.11. to 5.16**.

**Table 15:** List of serpins from *Xenopus tropicalis* genome.

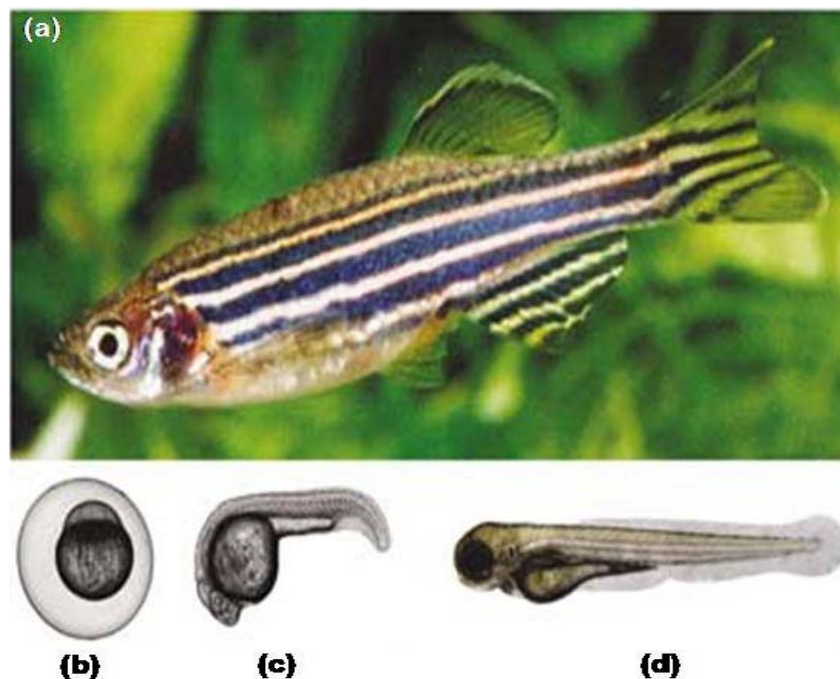
Name Given	Accession Id	Clade	Protein length	Homology to known serpin in NR database
Xtr-Spn-1	fgenes1_kg.C_scaffold_95000011	B	379	SPB5
Xtr-Spn-2	fgenes1_kg.C_scaffold_95000012	B	374	SPB12
Xtr-Spn-3	fgenes1_kg.C_scaffold_95000013	B	379	SPB6
Xtr-Spn-4	fgenes1_kg.C_scaffold_95000014	B	377	MNEI
Xtr-Spn-5	estExt_fgenes1_pm.C_2110010	B	370	MNEI
Xtr-Spn-6	fgenes1_pm_kg.C_scaffold_211000008	B	392	MNEI
Xtr-Spn-7	fgenes1_kg.C_scaffold_185000010	A	437	$\alpha_1$ -AT
Xtr-Spn-8	e_gw1.185.72.1	A	435	$\alpha_1$ -AT
Xtr-Spn-9	estExt_fgenes1_pg.C_1850041	A	435	$\alpha_1$ -AT
Xtr-Spn-10	estExt_fgenes1_pg.C_1850042	A	414	$\alpha_1$ -AT
Xtr-Spn-11	C_scaffold_185000011	A	434	$\alpha_1$ -AT
Xtr-Spn-12	e_gw1.185.79.1	A	413	$\alpha_1$ -AT
Xtr-Spn-13	e_gw1.185.80.1	A	384	$\alpha_1$ -AT
Xtr-Spn-14	e_gw1.49.222.1	A	390	ZPI
Xtr-Spn-15	fgenes1_pg.C_scaffold_2000123	A	458	AGT
Xtr-Spn-16	ENSXETP00000048524	D	484	HCI1
Xtr-Spn-17	estExt_Genewise1.C_7340032	E	356	PAI1

<sup>1</sup> NCBI Xenopus genome resources website, <http://www.ncbi.nlm.nih.gov/genome/guide/frog/>

Xtr-Spn-18	fgenes1_kg.C_scaffold_750000001	E	397	GDN
Xtr-Spn-19	e_gw1.233.93.1	E	404	SERPINE3
Xtr-Spn-20	ENSXETP00000049461	I	411	NEURO
Xtr-Spn-21	ENSXETP00000049481	I	410	PANC
Xtr-Spn-22	ENSXETP00000050413	F	409	PEDF
Xtr-Spn-23	ENSXETP00000029676	F	400	$\alpha_2$ -AP
Xtr-Spn-24	estExt_fgenes1_pm.C_10068	C	456	ATIII
Xtr-Spn-25	estExt_fgenes1_pg.C_2770030	H	425	HSP47

### 5.3. *Danio rerio* and its serpins

*D. rerio* is a blue spotted fresh water tropical fish (27a). The main habitat of the zebrafish is Southeast Asia. It serves as a model organism in developmental biology (Detrich *et al.*, 1999), embryogenesis (Driever and Fishman, 1996) and in genetics. *D. rerio* has a short life cycle and its transparent embryos and early adults (Figure 27b-d) can be used for light microscopy (Kari *et al.*, 2007).



**Figure 27:** *Danio rerio*. (a) Adult zebrafish. (b) embryo at one-cell stage, (c) embryo at 24 h post-fertilization and (d) embryo 3 days at post-fertilization.

The *D. rerio* genome sequencing project was started at Wellcome Trust Sanger Institute in 2001. About 73% of genome sequence was finished by May 2007 as reported on the *D. rerio* genome sequencing webpage<sup>1</sup>. The sixth assembly of the zebrafish genome (Zv6) was released in March 2006. The zebrafish genome assemblies are automatically annotated and

<sup>1</sup> *Danio rerio* genome sequencing webpage, [http://www.sanger.ac.uk/Projects/D\\_rerio/](http://www.sanger.ac.uk/Projects/D_rerio/)

are accessible from Ensembl<sup>1</sup>. **Table 16** lists 31 serpins as detected in the *D. rerio* genome after several rounds of homology searches.

**Table 16: List of serpins from *Danio rerio* genome.**

Gene Name	Accession id	Clade	Protein length	Homology to known serpin in NR database
Dre-Spn-1	CAI20749	B	380	MNEI
Dre-Spn-2	CAI20745	B	382	MNEI
Dre-Spn-3	AAH53300	B	380	SBP6
Dre-Spn-4	AAQ97848	B	384	MNEI
Dre-Spn-5	AAH66740	B	382	SBP6
Dre-Spn-6	AAH64292	B	433	MNEI
Dre-Spn-7	NP_001013277	A	429	$\alpha_1$ -AT
Dre-Spn-8	NP_001071226	A	429	$\alpha_1$ -AT
Dre-Spn-9	XP_001104678	A	372	$\alpha_1$ -AT
Dre-Spn-10	NP_001099059	A	372	$\alpha_1$ -AT
Dre-Spn-11	XR_029524	A	304	$\alpha_1$ -AT
Dre-Spn-12	XP_695000	A	372	$\alpha_1$ -AT
Dre-Spn-13	NP_001038536	A	391	ZPI
Dre-Spn-14	XP_001343164	A	396	ZPI
Dre-Spn-15	NP_932329	A	454	AGT
Dre-Spn-16	NP_878300	D	507	HCII
Dre-Spn-17	XP_690192	E	392	PAI1
Dre-Spn-18	Q7ZVL5	E	392	GDN
Dre-Spn-19	ENSDARP00000074162	E	407	SERPINE3
Dre-Spn-20	ENSDARP00000017430	I	412	NEURO
Dre-Spn-21	ENSDARP00000069366	F	406	PEDF
Dre-Spn-22	ENSDARP00000078640	F	480	A2AP
Dre-Spn-23	ENSDARP00000041512	G	403	C1IN
Dre-Spn-24	ENSDARG00000042684	C	452	ATIII
Dre-Spn-25	ENSDARP00000037780	H	405	HSP47
Dre-Spn-26	ENSDARP00000028177	H	403	HSP47
Dre-Spn-27	ENSDARP00000052941	H	414	HSP47
Dre-Spn-28	AAI53324	B	377	MNEI
Dre-Spn-29	AAI52147 (Zgc: 173729)	B	439	MNEI
Dre-Spn-30	XP_001331039	B	384	MNEI
Dre-Spn-31	XP_697505	B	440	MNEI

These serpins are further characterized in **section 5.11.** to **5.16.**

#### 5.4. *Tetraodon nigroviridis* and its serpins

*Tetraodon nigroviridis*, the green spotted freshwater pufferfish is highly popular as aquarium fish (**Figure 28**). It belongs to largest genus of the order Tetraodontiformes in the pufferfish family Tetraodontidae.

<sup>1</sup> The zebrafish genome at Ensembl website [http://www.ensembl.org/Danio\\_rerio/index.html](http://www.ensembl.org/Danio_rerio/index.html)





**Figure 28:** *Tetraodon nigroviridis*.

This pufferfish has a very small genome of approximately 350 Mb, consisting of 21 chromosomes. About 45,000 contigs were assembled in more than 12,000 scaffolds covering 332.5 Mb at a depth of about 8-fold (Jaillon *et al.*, 2004). The *T. nigroviridis* genome serves as a model system to study whole genome duplication (WGD) events and synteny to the genomes of mammals. **Table 17** lists 19 serpins detected in the *T. nigroviridis* genome after several rounds of homology searches. These serpins are further characterized in **section 5.11** to **5.16**.

**Table 17:** List of serpins from *Tetraodon nigroviridis*. \* Partial sequence, # Due to missing sequence

Gene Name	Accession Id	Clade	Protein length	Homology to known serpin in NR database
Tni-Spn-1	GSTENP00015677001	B	380	MNEI
Tni-Spn-2	GSTENP00015675001	B	303	MNEI
Tni-Spn-3	GSTENP00007903001	A	401	$\alpha_1$ -AT
Tni-Spn-4	GSTENP00018460001	A	416	$\alpha_1$ -AT
Tni-Spn-5	GSTENP00008425001	A	394	$\alpha_1$ -AT
Tni-Spn-6	GSTENP00018459001	A	413	$\alpha_1$ -AT
Tni-Spn-7	GSTENP00031597001	A	201*	AGN
Tni-Spn-8	GSTENT00032260001	A	400	ZPI
Tni-Spn-9	GSTENP00028636001	D	504	HCII
Tni-Spn-10	GSTENP00026727001	E	397	GDN
Tni-Spn-11	GSTENP00034604001	I	374	NEURO
Tni-Spn-12	GSTENP00013159001	F	60*	PEDF
Tni-Spn-13	GSTENP00014689001	F	411	$\alpha_2$ -AP
Tni-Spn-14	GSTENP00009345001	G	593	C1IN
Tni-Spn-15	GSTENP00004792001	C	453	ATIII
Tni-Spn-16	GSTENP00006756001	H	287	HSP47
Tni-Spn-17	GSTENT00003787001	E	356	PAI1
Tni-Spn-18	GSTENT00016647001	B	405	SBP6
Tni-Spn-19	GSTENT00029213001	E	364	SERPINE3

### 5.5. *Fugu rubripes* and its serpins

*Fugu rubripes* is a poisonous marine fish (**Figure 29**). It has one of the smallest genomes out of all known vertebrates (390Mb) around one eighth the size of the human genome, but as a

vertebrate, it has a similar complement of genes to that of mammals (Elgar et al. 1996). Fugu genomic sequences are available in the form of 12,381 scaffolds, ranging in size from 657 to 2 kb with approximately 30,000 potential genes (Aparicio *et al.*, 2002).



Figure 29. *Fugu rubripes*.

Table 18 lists 21 serpins detected in the *F. rubripes* genome after several rounds of homology searches. These serpins are further studied in section 5.11 to 5.16.

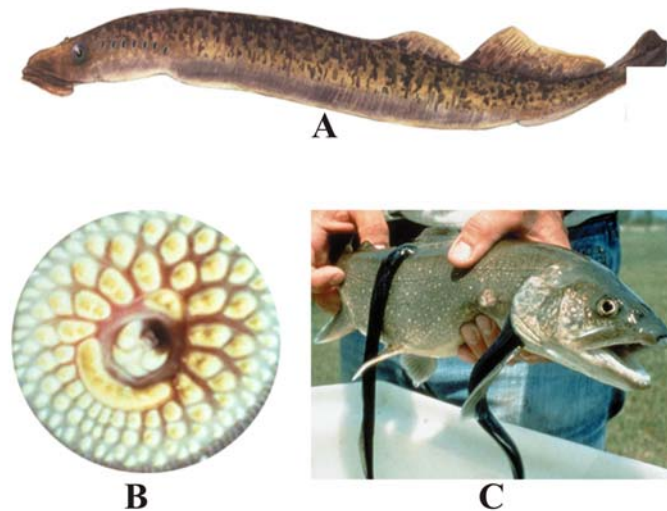
Table 18: List of serpins from *Fugu rubripes*.

Fugu Serpins	Gene Identifier in Fugu v4 [Fugu v3]	Scaffold Id	Clade	Protein Size	Homology to known serpin in NR database
Fru-Spn-1	e_gw2.131.10.1 [FRUP00000156735]	scaffold_131 [scaffold_2913]	H	405	HSP47
Fru-Spn-2	e_gw2.111.104.1 [FRUP00000155065]	scaffold_111 [scaffold_757]	A	423	$\alpha_1$ -AT
Fru-Spn-3	fgh5_pm.C_scaffold_488000001 [FRUP00000161527]	scaffold_488 [scaffold_2007]	I	407	NUERO
Fru-Spn-4	e_gw2.88.117.1 [FRUP00000132180]	scaffold_88 [scaffold_188]	A	397	PZI
Fru-Spn-5	e_gw2.123.110.1 [FRUP00000137160]	scaffold_123 [scaffold_417]	E	408	NEXIN
Fru-Spn-6	FRUP00000149263	scaffold_385	D	502	HCII
Fru-Spn-7	FRUP00000160285	scaffold_6239	A	413	$\alpha_1$ -AT
Fru-Spn-11*	FRUP00000140727	scaffold_508	A	462	AGT
Fru-Spn-12	FRUP00000141273	scaffold_1026	F	420	PEDF
Fru-Spn-14	FRUP00000162952	scaffold_154	F	435	$\alpha_2$ -AP
Fru-Spn-15	e_gw2.417.16.1	scaffold_417	F	455	$\alpha_2$ AP
Fru-Spn-17	FRUP00000155064	scaffold_111 [scaffold_757]	A	435	$\alpha_1$ -AT
Fru-Spn-18	FRUP00000146289	scaffold_641	A	392	$\alpha_1$ -AT
Fru-Spn-35	estExt_GW.C_1290044 [FRUP00000131353]	scaffold_129 [scaffold_110]	B	380	MNEI
Fru-Spn-36	FRUP00000163136	scaffold_2405	B	480	MNEI
Fru-Spn-37	e_gw2.671.2.1[FRUP00000138778]	scaffold_671 [scaffold_5139]	B	411	SBP6
Fru-Spn-38	e_gw2.269.120.1 [FRUP00000165249]	scaffold_1226	C	447	ATIII
Fru-Spn-39	FRUP00000133449	scaffold_3139	G	492	C1IN
Fru-Spn-40	e_gw2.275.54.1	scaffold_275	E	419	PAI
Fru-Spn-41	fgh5_pg.C_scaffold_186000009	scaffold_188	H	378	HSP47
Fru-Spn-42	FRUP00000142610	Scaffold_1209	E	201	SerpinE3

\* Numbering is not continuous since the same gene has two accession id detected.

### 5.6. *Petromyzon marinus* and its serpins

*Petromyzon marinus* is an aquatic eel-like, blood-sucking parasitic animal (**Figure 30**). It belongs to the most basal extant group of vertebrates and is supposed to have existed largely without any change for >500 million years. This organism serves as a model system for evolutionary biology since its study is expected to provide information on the early evolution of vertebrates. In addition, it serves as a model organism in developmental biology.



**Figure 30:** *Petromyzon marinus*. (A) Adult sea lamprey. (B) Sucking disc (mouth). (C) A lamprey attached to a trout<sup>1</sup>.

A preliminary 5.9-fold assembly of the sea lamprey genome (Feb 2007) is available via the Ensembl website<sup>2</sup>. Homology searches in this genomic assembly of the lamprey genome detected only serpins of groups V1, V2, V4, and V6 as summarized in **Table 19**. The inability to detect members of group V3 and V5 is most probably due to the incompleteness of this genomic assembly. These serpins are further analyzed in **section 5.11** to **5.12**.

**Table 19:** List of serpins from *Petromyzon marinus*.

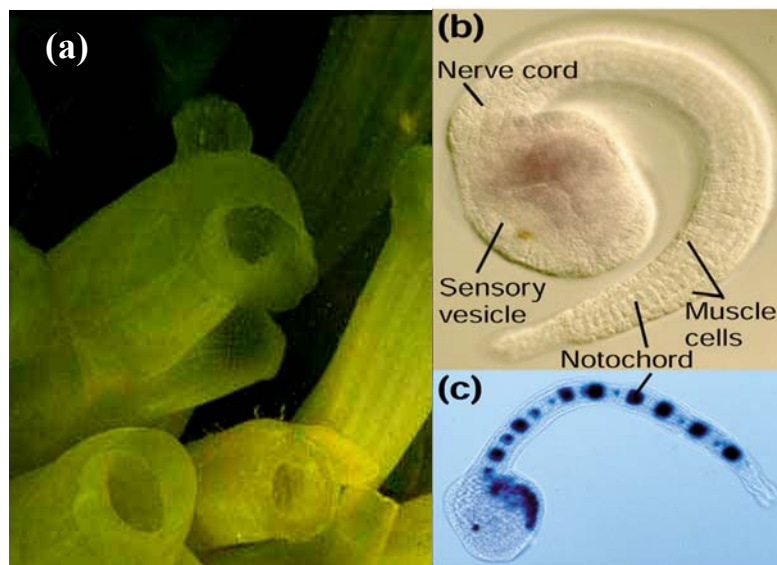
Name Given	Ensembl Accession id	Group	Clade	Protein Size	RCL P1-P1'
Pma-Spn-1	GENSCAN00000114312	V1	B	280	R-C
Pma-Spn-2	GENSCAN00000029305	V1	B	430	R-C
Pma-Spn-3	GENSCAN00000124947	V1	B	369	M-C
Pma-Spn-4	GENSCAN00000089208	V2	A	479	I-S
Pma-Spn-5	GENSCAN00000067410	V2	D	517	L-T
Pma-Spn-6	GENSCAN00000047295	V4	F	443	M-S
Pma-Spn-7	GENSCAN00000097429	V4	F	447	T-N
Pma-Spn-8	GENSCAN00000147606	V6	H	435	M-R

<sup>1</sup> Source: National human genome research institute [www.genome.gov](http://www.genome.gov).

<sup>2</sup> Ensembl website, [www.ensembl.org](http://www.ensembl.org).

### 5.7. Characterization of serpins from *Ciona intestinalis* genome

*Ciona intestinalis* (sea squirt) belongs to the urochordata (tunicates) in the class ascidiacea. It is a non-vertebrate chordate that diverged very early from the other chordates, namely cephalochordates and vertebrates, approximately 550 million years ago. Therefore, it is considered highly important for the understanding of the evolution of the vertebrates. *Ciona* species live in flat water areas of the oceans and go through two phases of the life cycle – an adult stage (**Figure 31a**) which metamorphoses from free swimming tadpole stage (**Figure 31b**). The tadpole is built of approximately 2500 cells, whose development can be observed easily under the microscope on the basis of the transparency of the larva (Corbo *et al.*, 1997) (**Figure 31c**). Additionally, this organism has the relatively short life cycle of approximately three months, making it a good system for developmental research.



**Figure 31:** *Ciona intestinalis*. (a) Adult. (b) Early tadpole larva stage. (c) a LacZ gene expressed in tadpole larva using electroporation technique for functional analysis (Corbo *et al.*, 1997). The figure is adopted from a review on *Ciona intestinalis* (Canestro *et al.*, 2003).

A whole genome shotgun assembly of *C. intestinalis* was released by the JGI with 11-fold coverage with estimated genome size of 173 Mb, which contains ~16,000 genes (Dehal *et al.*, 2002). Using *C. intestinalis* genomic assembly v1.95 (October 2002), serpins were searched and further analyzed for gene structure and synteny mapping as described in **section 4.7.1**.

This analysis and summary of *Ciona* serpins are similar to analysis of *Ciona* serpins carried out by Olaf Krüger (Krüger, 2003), with some exceptions as summarized below. There are eleven serpins in the *C. intestinalis* genome as summarized in **Table 20**.

The majority of *Ciona* serpins have AGY as codon for S56, except for Ci-Spn-1 and Ci-Spn-2 where it is TCN. Ci-Spn-4 and Ci-Spn-5 do not have serine at position 56. The presence of AGY codon in the majority of *Ciona* serpins indicates that it is an excellent case of TCN-

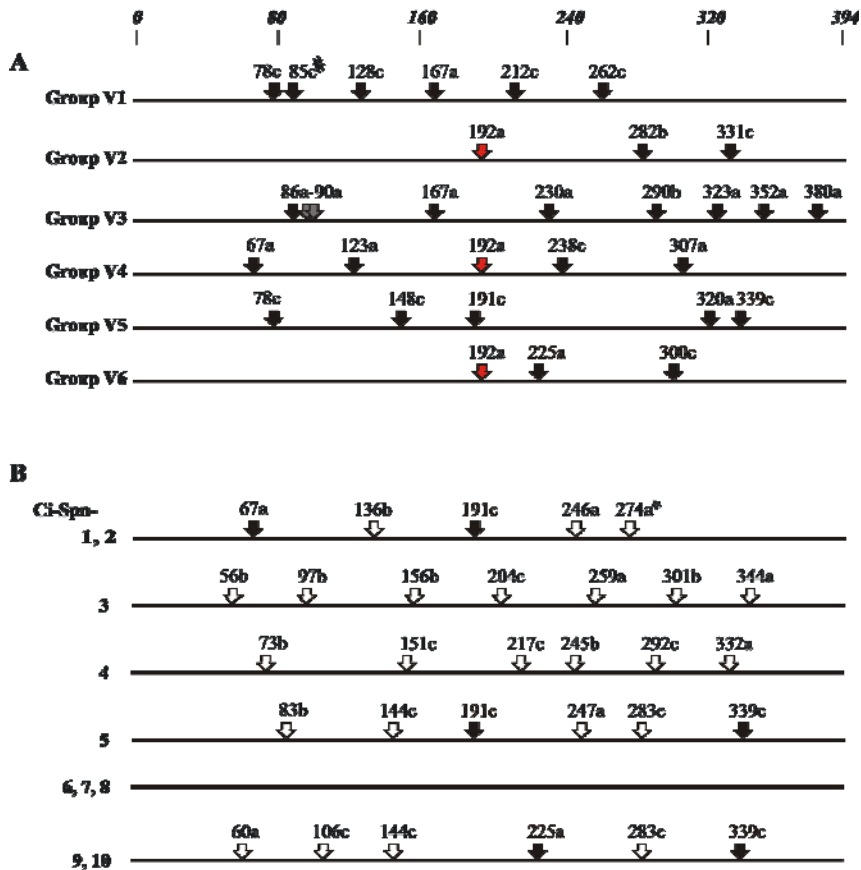
AGY usage dichotomy (Krem and Di Cera, 2003) being *Ciona* a deuterostome, with exception of two serpin genes.

**Table 20: List of serpins from *Ciona* genome draft version v1.95<sup>1</sup>.** The gene name with the prefix is enlarged Ci (for *Ciona intestinalis*) and Spn (for Serpin). Presence of Expressed tag sequence (EST) and S56 codon dichotomy, and P1-P1' residues in RCL is indicated. \$ indicates that ci0100146394 is accession id for Ci-Spn-6, Ci-Spn-7 and Ci-Spn-8 in database for *Ciona* genome draft version v1.95<sup>1</sup>. The serpin Ci-Spn-10 shows two variations of the RCL exon, named A and B, respectively.

Gene name	JGI protein accession id	Protein length	EST	S56 Codon dichotomy			RCL P1-P1'
				TCN	AGY	Comment	
Ci-Spn-1	ci0100132788	449	E	TCG			R-S
Ci-Spn-2	ci0100132818	412	E	TCG			R-S
Ci-Spn-3	ci0100134682	402	E		AGT		R-S
Ci-Spn-4	ci0100141118	441	E			No S56	R-S
Ci-Spn-5	ci0100143209	413	E			No S56	S-V
Ci-Spn-6	ci0100146394 <sup>\$</sup>	377	E		AGC		R-S
Ci-Spn-7	ci0100146394 <sup>\$</sup>	380	E		AGC		S-M
Ci-Spn-8	ci0100146394 <sup>\$</sup>	379	E		AGC		R-S
Ci-Spn-9	ci0100148346	409	E		AGT		D-S
Ci-Spn-10A	ci0100154072 <sup>#</sup>	408	E		AGT		R-S
Ci-Spn-10B	ci0100154072 <sup>#</sup>	407	E		AGT		P-L

*Ciona* serpins have unique gene structures as compared to vertebrate serpins (**Figure 32**). Some of the vertebrate intron positions were found to be conserved in *Ciona* serpins, such as intron positions 67c and 191c in Ci-Spn-1 and Ci-Spn-2. In case of Ci-Spn-5, introns at position 191c and 339c have counterparts in vertebrate group V5 (ATIII). The introns at positions 225a and 339c were present in Ci-Spn-9 and Ci-Spn-10, which tally with known introns of vertebrate group V6 and V5, respectively. But by and large, *Ciona* serpins have a different set of intron positions as compared to vertebrate serpins (Ragg *et al.*, 2001).

<sup>1</sup> Website of *Ciona* genome draft version v1, <http://genome.jgi-psf.org/ciona4/ciona4.home.html>



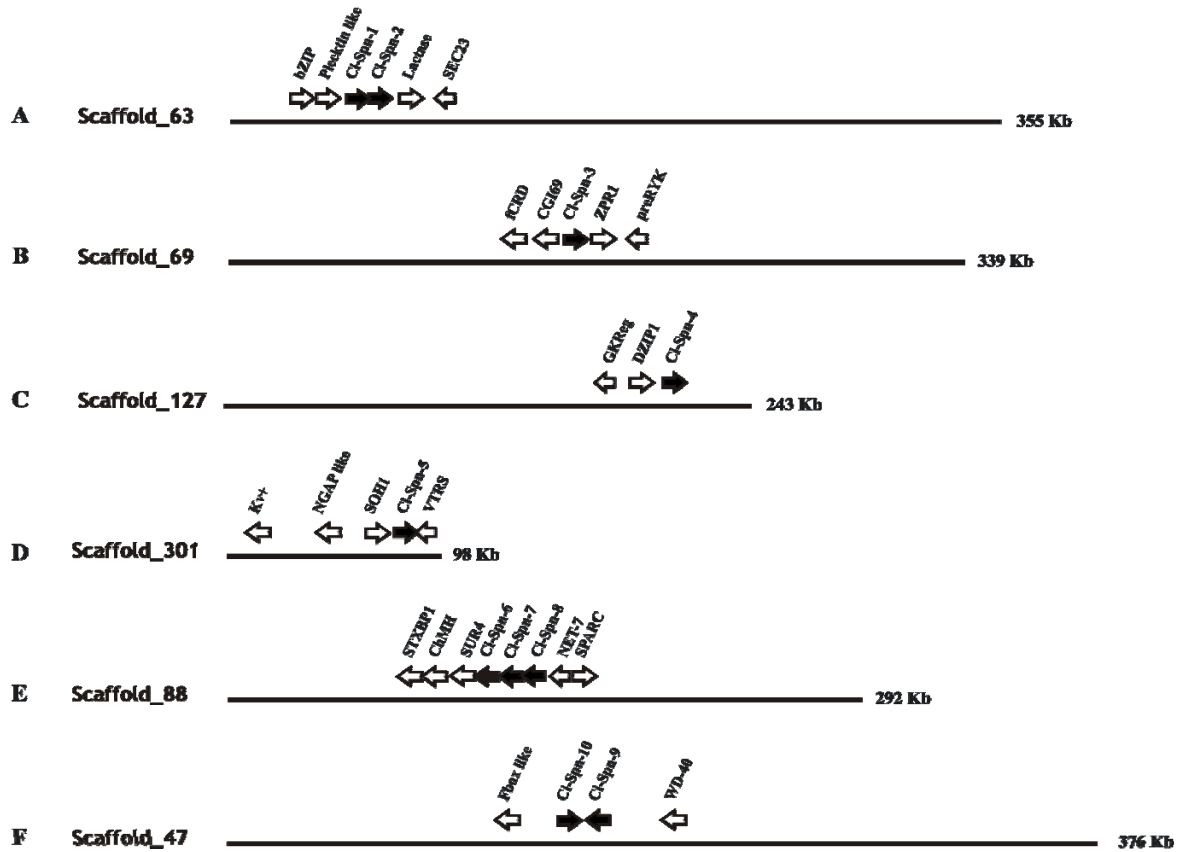
**Figure 32: Comparison of gene structures of vertebrate and *Ciona* serpins.** (A) Vertebrate serpins. (B) *Ciona serpins*. The black arrows indicate vertebrate specific intron positions marked with respect to  $\alpha_1$ -antitrypsin. White arrows indicate unique intron positions in *Ciona* serpins. # indicates 85c position that can differentiate between vertebrate serpin groups V1a and V1b. \* indicates that intron at position 274a is only found in Ci-Spn-1. The intron at position 192a is characteristic for group V2, V4 and V6 (marked by red color).

In order to understand orthology and divergence of *Ciona* serpins as compared to vertebrate serpins, synteny analysis was performed. Ci-Spn-1 and Ci-Spn-2 were found to be adjacent to each other and having the same orientation. These genes are located in scaffold\_63, flanked by bZIF and a Pleckstrin-like gene on one side, while on the other side the lactase and SEC23 genes are located (**Figure 33A**).

Ci-Spn-3 was located on scaffold\_69, surrounded by fCRD and CGI-69 on one side, and ZFR1 and preRYK are located on the other side (**Figure 33B**). Ci-Spn-4 was located on scaffold\_127 as a single serpin gene (**Figure 33C**). Ci-Spn-5 is located on scaffold\_301 as a single serpin gene flanked by Kv+-NGAP like gene-SOH1 on one side, while the VTRS gene is situated on the other side (**Figure 33D**). Ci-Spn-6, Ci-Spn-7, and Ci-Spn-8 are found on scaffold\_88 adjacent to each other in the same orientation (**Figure 33E**).

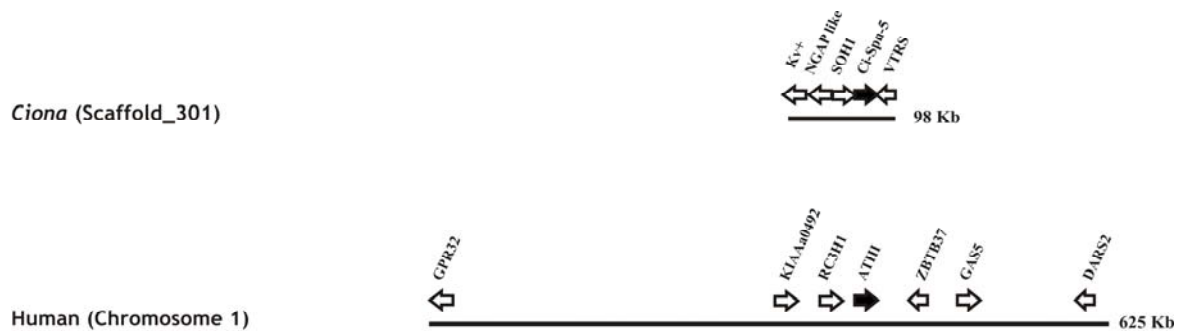
*Ciona* serpins Ci-Spn-9 and Ci-Spn-10 are found in opposite orientations on scaffold\_47 flanked by a Fbox like gene and WD-40 (**Figure 33F**). All marker genes are summarized in **appendix 8.4.1** with accession id. None of these gene organizations of *Ciona* serpins matched with that of the vertebrate serpins.





**Figure 33: Genomic organization of *Ciona* serpins.** *Ciona* serpins (black arrows) were identified on different scaffolds surrounded by marker genes (white arrows and appendix 8.4.1).

Ci-Spn-5 and human ATIII share two intron positions. To understand whether this is incidental or due to common ancestry, we compared the genomic locations of these two genes, but no common microenvironment was observed (**Figure 34**).



**Figure 34: Comparison of genomic organization of *Ciona* serpin Ci-Spn-5 and human ATIII.**

Also, the genomic locations of Ci-Spn-9 and Ci-Spn-10 and of human HSP47 were compared as these genes share an intron at position 225a (**Figure 32**) and a C-terminal endoplasmic reticulum retention signal in the respective protein. Again, these genes were not found to share a common genomic localization (**Figure 35**).

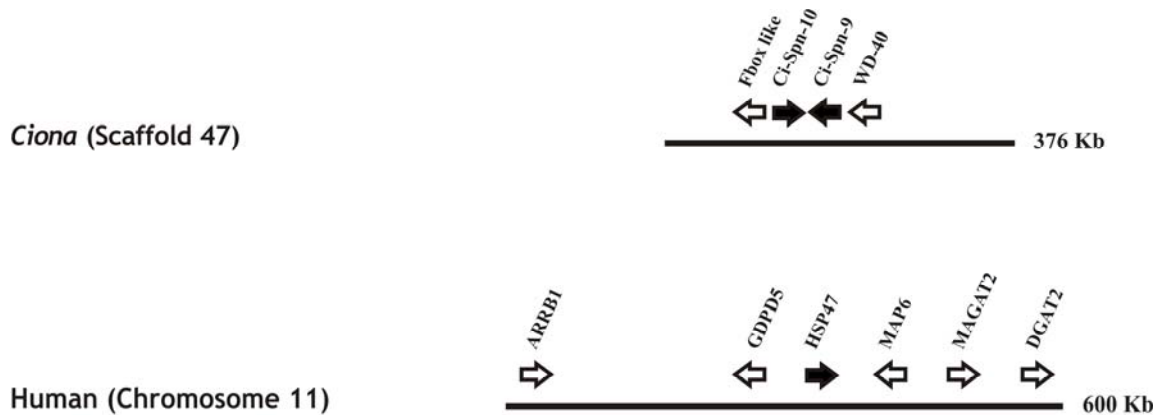


Figure 35: Comparison of genomic localization of *Ciona* serpins Ci-Spn-9 and Ci-Spn-10, and human HSP47.

The **appendix 8.3.1** summarizes the protein alignment of serpin sequences from *Ciona*. Furthermore, a phylogenetic tree of *Ciona* serpins (**Figure 36**) was generated based on the Neighbor-joining (NJ) method using Mega 3.1 (Kumar *et al.*, 2004). *Ciona* serpins that grouped in the same scaffold show clustering within this tree, supported by high bootstrap values.

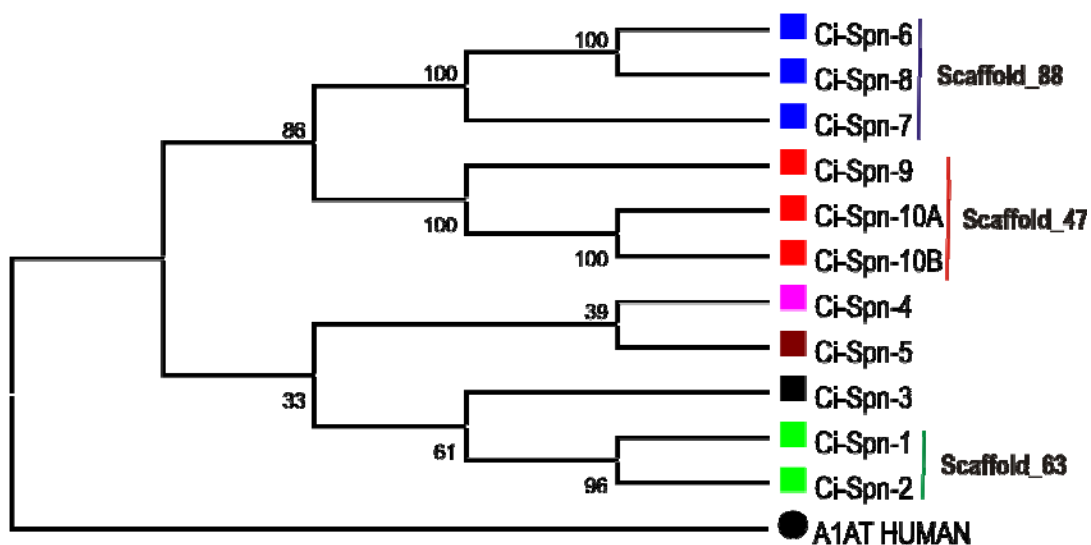


Figure 36: Phylogenetic tree of *Ciona* serpins. The tree was generated by Neighbor-joining (NJ) method using Mega 3.1 with bootstrap value = 1000. Human  $\alpha_1$ -antitrypsin was used as an outgroup. The color bars represent *Ciona* serpins in different scaffolds, which cluster in this phylogenetic tree together.

In summary, *Ciona* harbors serpins that are highly diverged and that have no orthologs in vertebrates based on analysis of synteny, gene structures, and protein sequences.



### 5.8. *Branchiostoma floridae* and its serpins

*Branchiostoma floridae* (amphioxus) belongs to cephalochordates, which are composed of 25-30 species of lancelets and which are small and fish-shaped creatures inhabiting shallow tropical and temperate oceans. The phylum Chordata is composed of the vertebrates, urochordates and cephalochordates that descended from a common ancestor that lived perhaps 550 million years ago (Putnam *et al.*, 2008). Thus, lancelets are important in evolutionary biology for understanding the origin of chordates.



**Figure 37:** *Branchiostoma floridae*. This figure is taken from JGI *B. floridae* genome database website<sup>1</sup>.

The draft assembly of the amphioxus genome sequence<sup>2</sup>, in which initial gene and protein predictions had been made using the JGI annotation pipeline, as described by Putnam *et al.* (2008). Nine serpins were detected from *B. floridae* genome as summarized in **Table 21**. The protein alignment of serpins from *B. floridae* is shown in **appendix 8.3.2**.

**Table 21:** List of serpins from *Branchiostoma floridae*.

Name	Gene Model@	Protein ID@	Protein size (aa)	S56 Codon dichotomy		RCL P1-P1'
				TCN	AGY	
Bfl-Spn-1	estExt_fgenes2_pg.C_4600026	130749	390		AGT	R-S
Bfl-Spn-2	fgenes2_pm.scaffold_460000005	62047	387		AGT	G-G
Bfl-Spn-3	fgenes2_pg.scaffold_11000109	66648	385		AGT	R-S
Bfl-Spn-4	estExt_fgenes2_pg.C_6170006	131993	377		AGC	R-S
Bfl-Spn-5	fgenes2_pg.scaffold_11000110	66649	377		AGT	R-S
Bfl-Spn-6	fgenes2_pg.scaffold_11000112	66651	416		AGT	M-S
Bfl-Spn-7	estExt_fgenes2_pg.C_200164	118881	309		AGT	C-A
Bfl-Spn-8	fgenes2_pg.scaffold_1013000002	112013	378		AGT	L-S
Bfl-Spn-9	estExt_fgenes2_pg.C_6170005	131992	416		AGT	R-S

@ Sequence information can be fetched from *Branchiostoma floridae* genome database using the following link:  
<http://genome.jgi-psf.org/cgi-bin/searchGM?db=Brafl1>

### 5.9. The sea urchin *Strongylocentrotus purpuratus* and its serpins

The purple sea urchin *Strongylocentrotus purpuratus* belongs to the phylum echinodermata. It serves as a research model system for molecular, evolutionary and cell biology (Sodergren *et*

<sup>1</sup> JGI *B. floridae* genome database website, <http://genome.jgi-psf.org/Brafl1/Brafl1.home.html>

<sup>2</sup> *B. floridae* genome website, <http://genome.jgi-psf.org/Brafl1/Brafl1.home.html>

*al.*, 2006). The genomic sequence of *S. purpuratus* (genome size - 814Mb with about 23000 genes) is available from Sea urchin genome database at Baylor College of Medicine<sup>1</sup>.



**Figure 38:** *Strongylocentrotus purpuratus*. This figure is taken from Monterey bay aquarium website<sup>2</sup>.

Using homology searches, 10 serpins were detected in this genome as summarized in **Table 22**. The protein alignment of serpins from *S. purpuratus* is shown in **appendix 8.3.3**.

**Table 22:** List of serpins from sea urchin - *Strongylocentrotus purpuratus*.

Name given	Accession ID@	Scaffold@	Protein (length)	Putative	S56 Codon dichotomy		Gene structure
				RCL	TCN	AGY	
Spu-spn-1	GLEAN3_28469	Scaffold49418	418	R-S		AGT	No intron <sup>§</sup>
Spu-Spn-2	GLEAN3_13378	Scaffold104538	393	G-C		AGC	No intron
Spu-Spn-3	GLEAN3_13377	Scaffold104538	393	R-C		AGC	No intron
Spu-Spn-4	GLEAN3_09346	Scaffold23825	413	C-L		AGC	No intron <sup>§</sup>
Spu-Spn-5	GLEAN3_18631	Scaffold85441	395	G-G		AGC	No intron
Spu-Spn-6	GLEAN3_24263	Scaffold21611	376	G-C		AGC	300c-intron
Spu-Spn-7	GLEAN3_18630	Scaffold85441	397	C-L		--	No intron
Spu-Spn-8	GLEAN3_18632	Scaffold85441	395	G-G		AGC	No intron
Spu-Spn-9	GLEAN3_04543	Scaffold60098	410#	M-M		AGC	No intron <sup>§</sup>
Spu-spn-10	GLEAN3_20278	Scaffold60098	196**	R-W			No intron

@ Sequence information can be fetched from Sea Urchin genome database using the following link:

<http://annotation.hgsc.bcm.tmc.edu/Urchin/cgi-bin/pubLogin.cgi>

# Low complexity regions are deleted from the sequence

\*Partial Sequence

§ No introns in conserved serpin domain (only one intron in signal peptide).

<sup>1</sup>Sea urchin genome database website, <http://www.hgsc.bcm.tmc.edu/project-species-o-Strongylocentrotus%20purpuratus.hgsc?pageLocation=Strongylocentrotus%20purpuratus>

<sup>2</sup> Monterey bay aquarium website, <http://www.montereybayaquarium.org/>

### 5.10. *Nematostella vectensis* and its serpins

The starlet sea anemone *Nematostella vectensis* (**Figure 39**) is a member of the oldest eumetazoan phylum, the Cnidaria that comprise anemones, corals, jellyfish and hydras.

*N. vectensis* is a simple eumetazoan, although recently available genome sequences suggest that its genome possesses more similarity to vertebrates than to flies or worms (Putnam *et al.*, 2007).



**Figure 39:** *Nematostella vectensis*. Taken from JGI *Nematostella vectensis* webpage<sup>1</sup>.

Three serpins were detected on BLAST searches against the *N. vectensis* genome as summarized in **Table 23**. The alignment of serpin sequences from *N. vectensis* is shown in **appendix 8.3.4**.

**Table 23:** List of serpins from *Nematostella vectensis*.

Name given	JGI accession id.	NCBI accession id.	Peptide Length	S56 Codon dichotomy		RCL
				TCN	AGY	P1-P1'
Nve-Spn-1	estExt_fgenesh1_pg.C_1860016	XP_001627732	397		AGC	R-S
Nve-Spn-2	e_gw.186.64.1	XP_001627750	374		AGC	R-C
Nve-Spn-3	estExt_GenewiseH_1.C_880258	XP_001632351	380		AGC	M-S

### 5.11. Orthology analysis of group V1 serpins

Group V1 vertebrate serpin genes depict five standard introns at positions 78c, 128c, 167a<sup>2</sup>, 212c, and 262c ( $\alpha_1$ -antitrypsin numbering) in their coding region (Ragg *et al.*, 2001). An additional intron is found in some group V1 members at position 85c, constituting group V1a, while members that lack intron at position 85c represent group V1b. These members are normally inhibitors of serine or cysteine proteases, but some of them are non-inhibitory (**Table 24**). Group V1 serpins are also named ov-serpins being closely related to ovalbumin; these serpins have been arranged in clade B in the clade-based classification system of serpins

<sup>1</sup> JGI *Nematostella vectensis* webpage, <http://genome.jgi-psf.org/Nemve1/Nemve1.home.html>

<sup>2</sup> This intron position is also shared by group V3 serpins

(Silverman *et al.*, 2001). These ov-serpins are primarily intracellular as they lack an N-terminal signal peptide. They also lack C-terminal extensions.

**Table 24: Physiological roles of group V1 serpins and associated diseases/syndromes.**

Group V1 serpins	Physiological Role(s)	Associated Disease(s)/Syndrome(s)
SERPINB1 (MNEI)	Inhibitor of neutrophil elastase	
SERPINB2 (PAI2)	Inhibitor of uPA	
SERPINB3 (SCCA1)	Inhibitor of cathepsin L and V	
SERPINB4 (SCCA2)	Inhibitor of cathepsin G and chymase	
SERPINB5 (Maspin)	Metastasis control by unknown mechanism, non-inhibitory.	Mouse knockouts are lethal
SERPINB6 (PI-6)	Inhibitor of cathepsin G	IgA nephropathy
SERPINB7 (Megsin)	Megakaryocyte maturation	
SERPINB8 (PI-8)	Inhibitor of furin	
SERPINB9 (PI-9)	Inhibitor of granzyme B	
SERPINB10 (Bomapin)	Inhibitor of thrombin and trypsin	
SERPINB11 (Epipin)	?	
SERPINB12 (Yukopin)	Inhibitor of trypsin	
SERPINB13 (Headpin)	Inhibitor of cathepsin L, protecting epithelial cells	

**Appendices 8.3.5 to 8.3.11** show protein alignments of group V1 serpins.

### 5.11.1. Gene structure of group V1 serpins

Since gene structure plays an important role in classifying groups V1-V6, the gene architectures of probable group V1 serpin homologs from different vertebrates were determined. **Table 25** shows that the vertebrate species investigated contain at least three genes that depict the basic exon-intron structure of group V1 serpins (for instance in *Fugu*, *Tetraodon*, and lamprey). Most of these genes contain the complete set of standard introns at the canonical positions 78c, 85c<sup>1</sup>, 128c, 167a<sup>2</sup>, 212c, and 262c. However, the SPB6 genes from *Fugu* and *Tetraodon* each possess two additional introns at positions 239c (novel) and 320a (**appendix 8.3.4**). Interestingly, group V5 serpins also contain an intron at position 320a. The SPB6 gene of *Tetraodon* has an intron at position ~85c (position cannot be exactly assigned, due to sequence ambiguity); this is a unique case since SPB6-like genes of all other vertebrates do not possess this intron. The structure of the Dre-Spn-5 gene cannot be determined, as only full-length cDNA information<sup>3</sup> of this gene is available from the zebrafish gene collection (ZGC).

<sup>1</sup> Only in serpin group V1a members.

<sup>2</sup> Also shared by group V3 serpins.

<sup>3</sup> cDNA id ZGC:76926.

Table 25: Intron positions of group V1 genes in different vertebrates. The presence (+) of intron positions is shown.

Group V1 serpin genes	Intron at position						Abnormalities in intron positions
	78c	85c	128c	167a	212c	262c	
MNEI_HSA (P30740)	+		+	+	+	+	
MNEI_MMU (Q5SUV7)	+		+	+	+	+	
MNEI_RNO (gi:72255515)	+		+	+	+	+	
MNEI_GGA (XP_418980)	+		+	+	+	+	
MNEI_XTR (fgenesh1_kg.C_scaffold_95000014)	+		+	+	+	+	
MNEI_FRU (FRUP00000131353)	+		+	+	+	+	
MNEI_TNI (GSTENP00015677001)	+		+	+	+	+	
MNEI_DRE (CAI20749)	+		+	+	+	+	
MNEIL_PMA (GENSCAN00000124947)	+		+	+	+	+	
PAI2_HSA (P05120)	+	+	+	+	+	+	
PAI2_MMU (P12388)	+	+	+	+	+	+	
PAI2_RNO (P29524)	+	+	+	+	+	+	
SPB5_HSA (P36952)	+		+	+	+	+	
SPB5_MMU (P70124)	+		+	+	+	+	
SPB5_RNO (P70564)	+		+	+	+	+	
SPB5_GGA (XP_418986)	+		+	+	+	+	
SPB5_XTR (fgenesh1_kg.C_scaffold_95000011)	+		+	+	+	+	
SPB6_HSA (P35237)	+		+	+	+	+	
SPB6_MMU (Q60854)	+		+	+	+	+	
SPB6_RNO (Q6P9U0)	+		+	+	+	+	
SPB6_GGA (XP_418981)	+		+	+	+	+	
SPB6_XTR (fgenesh1_kg.C_scaffold_95000013)	+		+	+	+	+	
pSPB6_FRU (e_gw2.671.2.1)	+		+	+	+	+	[+1239c, +]320a
pSPB6_TNI (GSTENT00016647001)	+	+	+	+	+	+	[+1239c, +]320a
pSPB6_DRE (AAH53300)	+		+	+	+	+	
SPB6_PMA (GENSCAN00000029305)	+		+	+	+	+	
Gga-Spn-3 (XP_426040/serpinB10)	+	+	+	+	+	+	
Gga-Spn-4 (NP_990228/MENT)	+	+	+	+	+	+	
Gga-Spn-5 (XP_418982)	+	+	+	+	+	+	
Gga-Spn-6 (ovalbumin)	+	+	+	+	+	+	
Gga-Spn-7 XP_418983/Gene Y protein)	+	+	+	+	+	+	
Gga-Spn-8 (XP_418984/Gene X protein)	+	+	+	+	+	+	
Gga-Spn-9 (XP_418985)	+	+	+	+	+	+	
Xtr-Spn-2 (fgenesh1_kg.C_scaffold_95000012)	+		+	+	+	+	
Xtr-Spn-5 (estExt_fgenesh1_pm.C_2110010)	+	+	+	+	+	+	
Xtr-Spn-6 (fgenesh_pm_kg.C_scaffold_211000008)	+		+	+	+	+	
Dre-Spn-2 (CAI20745)	+		+	+	+	+	
Dre-Spn-4 (AAQ97848)	+		+	+	+	+	
Dre-Spn-5 (AAH66740)#							.
Dre-Spn-6 (AAH64292)	+		+	+	+	+	
Dre-Spn-28 (AAI53324)	+		+	+	+	+	
Dre-Spn-29 (AAI52147)	+		+	+	+	+	
Dre-Spn-30 (XP_001331039)	+		+	+	+	+	
Dre-Spn-31 (XP_697505)	+		+	+	+	+	
Fru-Spn-36 (FRUP00000163136)	+		+	+	+	+	
Tni-Spn-2 (GSTENP00015675001)	+		+	+	+	+	
Pma-Spn-1 (GENSCAN00000114312)	+		+	+	+	+	

# Gene structure is not available.

### 5.11.2. Synteny analysis of group V1 serpins

To investigate orthology relationships of vertebrate group V1 serpins, the syntenic arrangements of these genes were analyzed in various vertebrates (**Figure 40**).

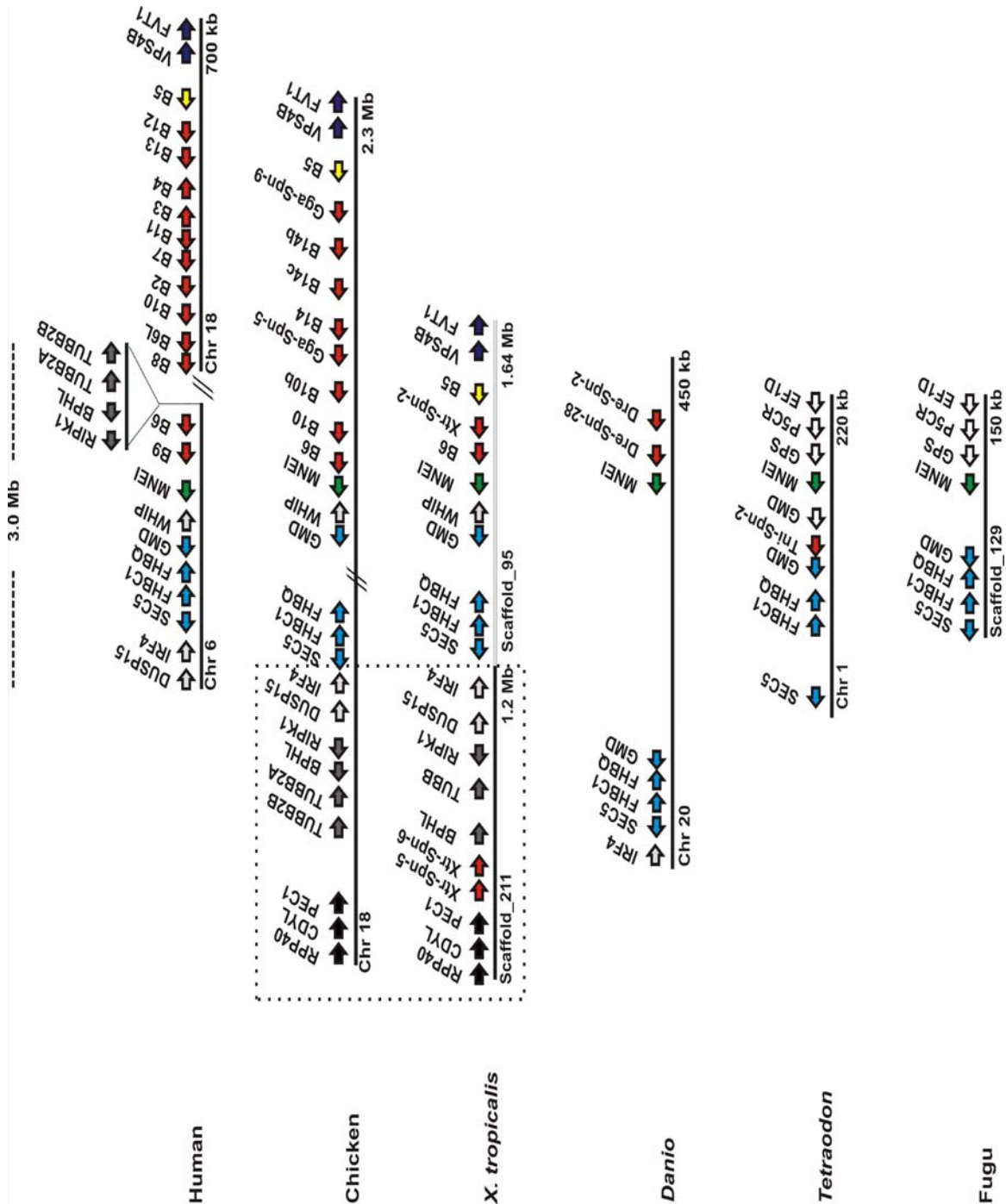


Figure 40: Synteny organization of group V1 serpins in vertebrates.

---

There are two clusters of group V1 serpins in the human genome, one on the chromosome 6 containing serpin-B1(MNEI), -B6 and -B9 in a 3 Mb region, that is flanked by markers RIPK1-BPHL-TUBB2A-TUBB2B<sup>1</sup> on one side and by WHIP-GMD-FHBQ<sup>1</sup> on the other side (**Figure 40**). The other cluster on the chromosome 18 encompasses serpin-B8, a -B6-like pseudo-gene, -B10, -B2, -B7, -B11, -B3, -B4, B13, B12 and -B5 within a 700 kb region, flanked by markers VPS4B-FVT1<sup>1</sup>. In the chicken genome, the group V1 serpins are organized in an uninterrupted, single cluster on chromosome 2, flanked by markers RIPK1-BPHL-TUBB2A-TUBB2B on one side, and markers VPS4B-FVT1 on the other side (Benarafa and Remold-O'Donnell, 2005).

In the *Xenopus tropicalis* genome, group V1 serpins are also organized in one cluster comprised of two scaffolds (scaffold\_211 and scaffold\_95) that are flanked by a series of similar marker genes as in the chicken genome. Interestingly, there are two unique serpins of group V1, named Xtr-Spn-5 and Xtr-Spn-6 in scaffold\_211, that are surrounded by markers RPP40-CDYL-PEC1 on the one side, and by marker triad BPHL-TUBB-RIPK1 on the other side. In the chicken genome, these extra group V1 serpins are not present, but the corresponding region has a set of conserved markers as in *Xenopus tropicalis* (black box in **Figure 40**). This indicates that the frog has a unique expansion of group V1 serpins adjacent to the main conserved cluster. Alternatively, these genes were lost in chicken.

In fish genomes, there is only one cluster that groups around serpinB1 and includes Dre-Spn-2 and Dre-Spn-28 in *Danio rerio* genome, whereas a pseudogene Tni-Spn-2 is present in the *Tetraodon* genome. A serpinB6-like gene is also present in these fish genomes, but this gene is found in another syntenic organization surrounded by conserved markers (**Figure 41**), suggesting that it is a paralogue, and not an ortholog of human serpinB6.

---

<sup>1</sup> See Appendix 8.4.2.

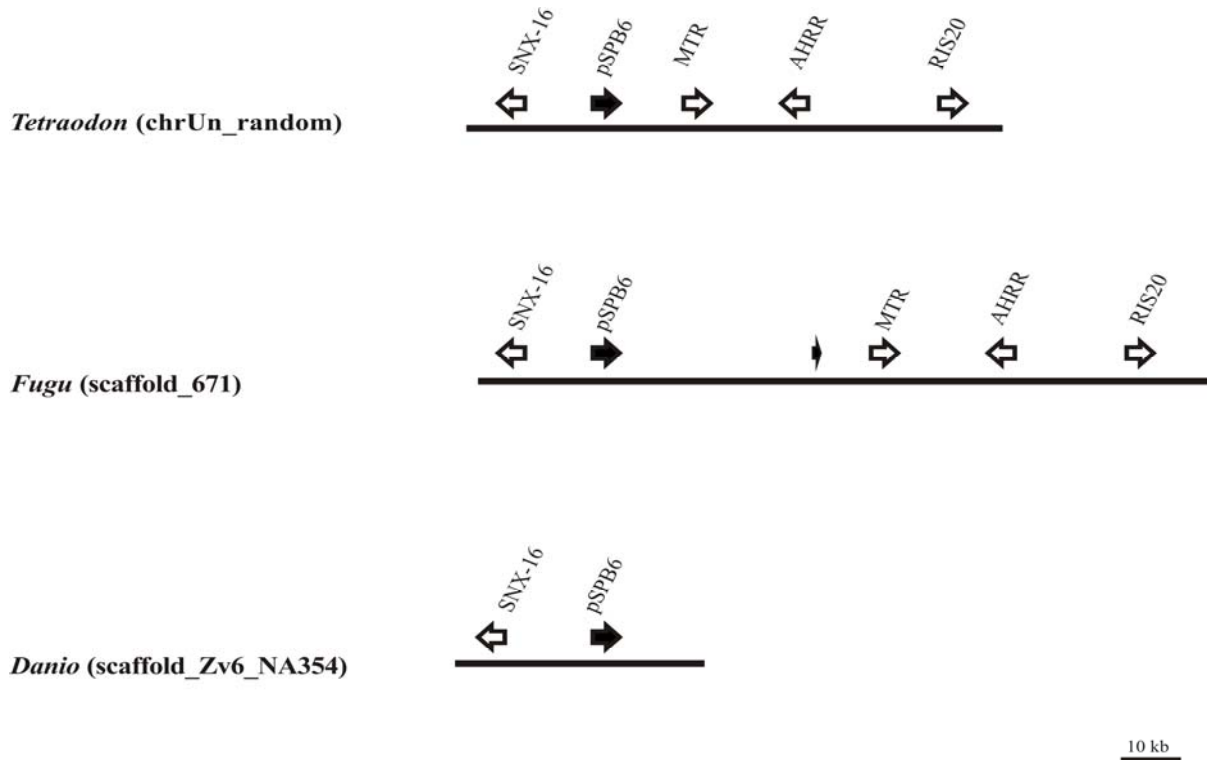


Figure 41: Genomic localization of serpinB6-like genes in fishes.

In summary, these data suggest that there is one orthologous cluster with group V1 serpin genes, conserved across different vertebrates since the fish/tetrapod split at around 450 million years (MY) ago. After separation of the chicken from mammals around 310 MY ago, this cluster bifurcated into two clusters by chromosomal breakage in mammals. There are some further fish-specific group V1 serpins with a unique genomic micro-environment as described in this section. Dre-Spn-4 of zebrafish is located on chromosome 19, flanked by a distinct set of markers (**appendix 8.4.1**), which assigns this gene to a unique micro-locus (**Figure 42**).



Figure 42: Genomic localization of Dre-Spn-4 from *Danio rerio*.

Dre-Spn-5 from *Danio rerio* cannot be located in the present assembly of genomic sequences, whereas Dre-Spn-6, Dre-Spn-29, Dre-Spn-30, and Dre-Spn-31 are localized on Contig: NW\_001884542.1 that is not assigned to any chromosome yet in the zebrafish genome (**Figure 43**).



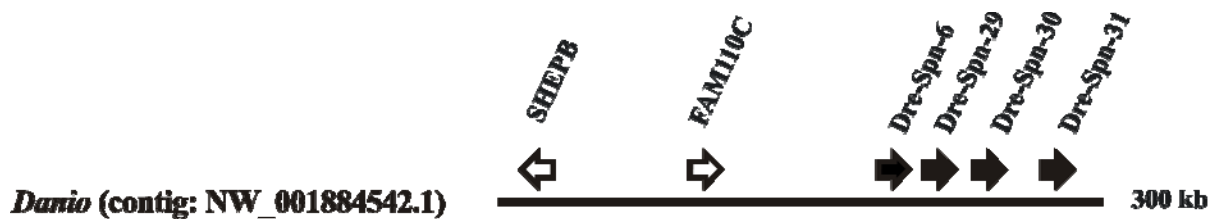


Figure 43: Genomic localization of four group V1 serpins from *Danio rerio*, namely Dre-Spn-6, Dre-Spn-29, Dre-Spn-30, and Dre-Spn-31.

### 5.11.3. Sequence analysis of group V1 serpins

To further delineate orthologs of group V1 serpins, sequence analyses were carried out. The major outcomes are reported below.

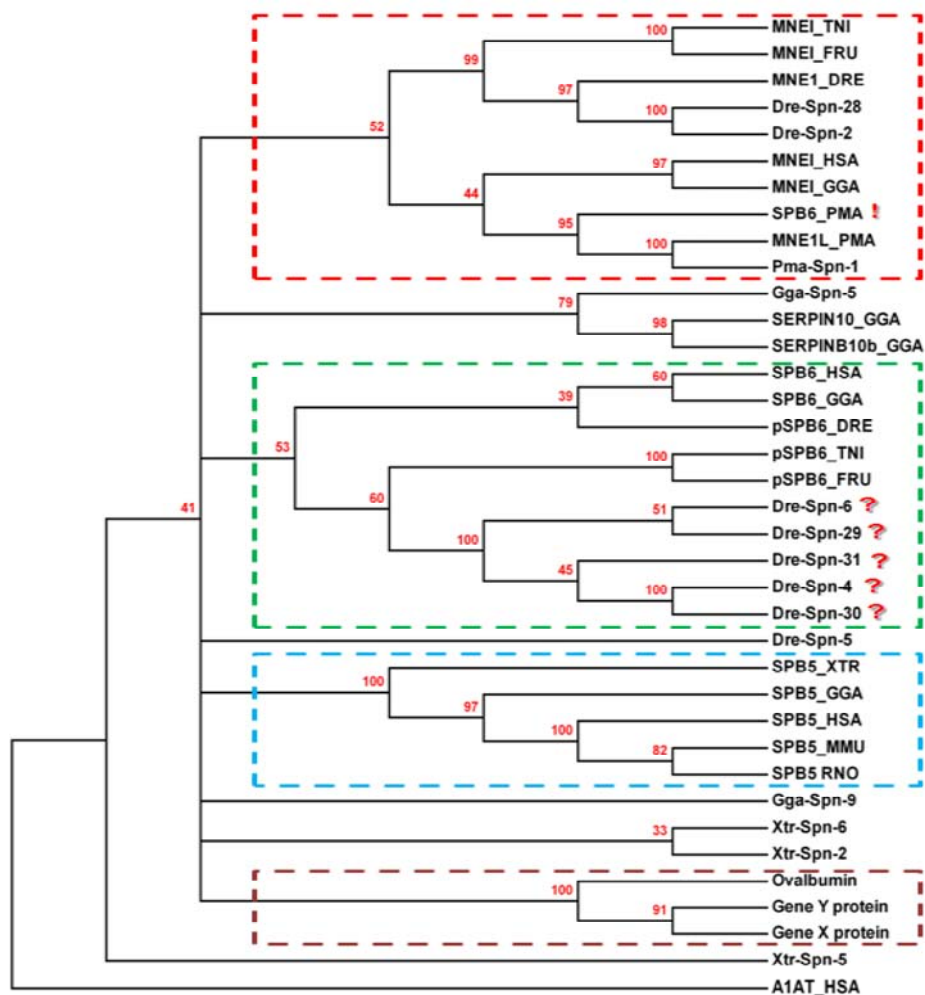
MNEI (serpinB1) is highly conserved in vertebrates, depicting 56-81% sequence identity and 72-92% sequence similarity on the amino acid level with human MNEI. The inhibitory RCL region is conserved, containing C-M at P1-P1' (**appendix 8.3.5**). PAI2 (serpinB2) is highly conserved in mammals, depicting 72-75% sequence identity and 86-87% sequence similarity on the amino acid level with human PAI2. The inhibitory RCL region is conserved exhibiting R-T at P1-P1' (**appendix 8.3.6**). PAI2 is further characterized by a conserved loop between helices C and D (CD loop), a cysteine disulfide bridge between C79-C161<sup>1</sup> and by absence of one amino acid insertion between positions 247/248. SPB5 (maspin/serpinB5) is conserved from frog to mammals with 53-89% sequence identity and 77-97% sequence similarity on the amino acid level with human SPB5 and is further characterized by non-inhibitory RCL (**appendix 8.3.7**). SPB6 from tetrapods and its paralogs in fishes (pSPB6) depict 39-75% sequence identity and 56-84% sequence similarity on the amino acid level with human SPB6 with a conserved inhibitory RCL region, containing R-C at P1-P1' (**appendix 8.3.8**). Chicken has ten group V1 serpins in a single cluster including orthologs of MNEI, SPB6 and SPB5 and other ovalbumin like genes (**appendix 8.3.9**). *Xenopus tropicalis* has six serpin genes into two clusters of group V1 serpins including orthologs of MNEI, SPB6 and SPB5; and several other group V1 serpins named as Xtr-Spn-2, Xtr-Spn-5 and Xtr-Spn-6 (**appendix 8.3.10**). *Danio* has ten group V1 serpins localized in different clusters including an ortholog of MNEI, a paralog of SPB6 and other eight ov-serpin genes (**appendix 8.3.11**).

To comprehend the relationship of these group V1 serpins, a phylogenetic tree (**Figure 44**) was created based on the NJ method (Saitou and Nei, 1987) with the help of MEGA4 (Tamura *et al.*, 2007). The phylogenetic tree of group V1 serpins from different vertebrate has three firmly established branches, MNEI-like, SBP6-like, and SBP5-like serpins, and out of these three classes, only first two are found in fishes. Dre-Spn-2 and Dre-Spn-28 are recent duplicates of MNEI\_DRE. Group V1 serpins from lamprey are clubbed with MNEI, corroborating basal nature of MNEI from where SBP6 originated by duplication. Thus, MNEI

<sup>1</sup> Numbering according to human PAI2 sequence.

is the ancestor of group V1 serpins found in early vertebrate originating at ~500 Mya. Furthermore, there are many paralogs of these genes in different organisms, making it difficult to decide orthology with human group V1 serpins. In many cases, orthologs such as ovalbumin, gene X protein, and gene Y proteins in chicken genome have no counterpart in humans. In addition, in *Xenopus tropicalis*, Xtr-Spn-5 and Xtr-Spn-6 are unique group V1 serpins with no orthologs in any other vertebrates.

Briefly, many different species possess species-specific group V1 serpins generated from tandem duplication events, apart from conserved members of group V1 serpins.



**Figure 44: Evolutionary tree of group V1 serpins from different vertebrates.** This tree has three major firmly supported branch - MNEI-like (red box), SPB5-like (blue box), and SPB6-like (green box). Group V1 serpins whose orthologs are not clear are named with species name such as Xtr-Spn-5 from *Xenopus tropicalis* (Xtr). The group V1 serpins from lamprey (PMA) are grouped with MNEI, and ! indicates origin of lamprey SBP6 as this gene is grouping with MNEI genes from other species, suggesting a basal nature of MNEI from which SBP6 originated by duplication. "?" indicates question that whether *Danio rerio* specific group V1 serpins are paralogue of SPB6 or not. This tree was created with the NJ method using MEGA4. Bootstrap values (in percentage) for 1000 replicates are shown. Branches corresponding to partitions reproduced in less than 30% bootstrap replicates are collapsed.

### 5.12. Orthology analysis of group V2 serpins

Group V2 of vertebrate serpins has been defined by a gene structure depicting three introns at homologous positions - 192a, 282b, and 331c ( $\alpha_1$ -antitrypsin numbering) in their coding region, and each member also has an intron mapping to the untranslated region (Ragg *et al.*, 2001). Group V2 is multi-membered, composed of  $\alpha_1$ -antitrypsin like serpins that are involved in different physiological roles, including inhibitors (like  $\alpha_1$ -antitrypsin or antichymotrypsin) and non-inhibitory members (like angiotensinogen) (**Table 26**).

**Table 26: Physiological roles of group V2 serpins and associated diseases/syndromes.**

Group V2 serpins	Physiological Role(s)	Associated Disease(s)/Syndrome(s)
SERPINA1 (A1AT)	Elastase inhibitor	Emphysema & serpinopathy
SERPINA2		
SERPINA3 (ACT)	Chymotrypsin inhibitor	Emphysema & serpinopathy
SERPINA4 (KALL)	Kallikrein inhibitor	
SERPINA5 (PCI)	Protein C inhibitor	
SERPINA6 (CBG)	Corticosteroid transporter	Chronic fatigue
SERPINA7 (THBG)	Thyroxine transporter	Hypothyroidism
SERPINA8 (AGT)	Blood pressure regulation	Hypertension
SERPINA9 (CEN)	B cell maintenance by inhibiting trypsin-like serine proteases	
SERPINA10 (ZPI)	Inhibitor of Factor Xa	Venous thromboembolic disease
SERPINA11		
SERPINA12 (VAS)	Adipokine with insulin-sensitizing effects	Metabolic syndrome
SERPINA13		
SERPIND1 (HCII)	Thrombin inhibitor	

The protein alignments of group V2 serpins are shown in **appendices 8.3.12 to 8.3.21**.

#### 5.12.1. Gene structure of group V2 serpins

To ensure group affiliation, the gene architectures of supposed group V2 serpin homologs were determined in different vertebrates. **Table 27** summaries that all vertebrates investigated contain several group V2 serpin genes with introns at the canonical positions 192a, 282b, and 331c. However, there are some genes with deviations from the standard structure. The intron at the position 331c in A1AT\_TNI cannot be assigned, probably due to sequencing errors around this position. AGT\_FRU has two additional introns at positions 77c and 233c. The 233c intron is also shared by AGT\_TNI. The presence of an intron at position 77c could not be identified in the AGT\_TNI gene, since there is a big gap in 5' part of this gene in the current version of the *Tetraodon* genome. HCII\_FRU and HCII\_TNI share an additional intron at position 241c. Furthermore, these genes have a common non-canonical intron in the non-conserved N-terminal domain, which can be assigned to position 85c (numbering according to HCII\_FRU, **appendix 8.3.13**). HCII\_PMA has additional introns in the serpin core at position 83c ( $\alpha_1$ -antitrypsin numbering). Similarly, HCII\_PMA also has two additional introns in its non-conserved 5' region (**appendix 8.3.13**). Fru-Spn-7 and TNI-Spn-3 share an extra intron at position 215c, whereas ZPI3\_FRU and ZPI3\_TNI each contain a non-standard intron at position 94a and both these two genes are renamed from here on as Spn\_215c and Spn\_94a

respectively. One of the group V2 serpins from *Danio*, Dre-Spn-11 is a pseudogene with the intron at position 331c not found due to a premature stop codon.

**Table 27: Intron positions of group V2 genes in different vertebrates.** The presence (+) of intron positions is shown. Abnormalities in intron positions within any group V2 serpin genes are also tabulated in last column. Note that only introns mapping to the serpin core are listed in this table.

Group V2 serpin genes	Intron at position			Abnormalities in intron positions
	192a	282b	331c	
A1AT_HSA (P01009)	+	+	+	
A1AT_MMU (P07758)	+	+	+	
A1AT_RNO (P17475)	+	+	+	
A1AT_GGA (XP_426460)	+	+	+	
A1AT_XTR (fgenesh1_kg.C_scaffold_185000010)	+	+	+	
A1AT_FRU (e_gw2.111.104.1)	+	+	+	
A1AT_TNI (GSTENP00018459001)	+	+	?	Gap in coding region
A1AT_DRE (NP_001013277)	+	+	+	
A2_HSA (P20848)	+	+	+	
A2_MMU (gi:20858201)	+	+	+	
A3_HSA (P01011)	+	+	+	
A3_MMU (Q9D490)	+	+	+	
A4_HSA (P29622)	+	+	+	
A4_MMU (P97569)	+	+	+	
A5_HSA (P05154)	+	+	+	
A5_MMU (Q5BKQ8)	+	+	+	
A5_RNO (Q66HL5)	+	+	+	
A6_HSA (P08185)	+	+	+	
A6_MMU (Q06770)	+	+	+	
A6_RNO (P31211)	+	+	+	
A7_HSA (P05543)	+	+	+	
A7_MMU (P61939)	+	+	+	
A7_RNO (P35577)	+	+	+	
AGT_HSA (P01019)	+	+	+	
AGT_MMU (P11859)	+	+	+	
AGT_RNO (P01015)	+	+	+	
AGT_GGA (gi:50741434)	+	+	+	
AGT_XTR (fgenesh1_pg.C_scaffold_2000123)	+	+	+	
AGT_FRU (FRUP00000140727)	+	+	+	[+177c, +1233c
AGT_TNI (GSTENP00031597001)	?	+	+	[?177c, +1233c
AGT_DRE (NP_932329)	+	+	+	
A9_HSA (Q86YP7)	+	+	+	
A9_MMU (Q9D7D2)	+	+	+	
A9_RNO (gi:56912218)	+	+	+	
ZPI_HSA (Q9UK55)	+	+	+	
ZPI_MMU (Q8R121)	+	+	+	
ZPI_RNO (Q62975)	+	+	+	
ZPI_GGA (XP_421341)	+	+	+	
ZPI1_XTR (e_gw1.49.222.1)	+	+	+	
ZPI1_FRU (e_gw2.88.117.1)	+	+	+	
ZPI_TNI (GSTENT00032260001)	+	+	+	
ZPI 1_DRE (NP_001038536)	+	+	+	
ZPI 2_DRE (XP_001343164)	+	+	+	
ZPI3_FRU (FRUP00000146289) /Spn_94	+	+	+	[+194a
ZPI3_TNI (GSTENP00008425001) /Spn_94	+	+	+	[+194a

A11_HSA (Q86YP6)	+	+	+	
A11_MMU (Q8CIE0)	+	+	+	
A11_RNO (gi:21717801)	+	+	+	
A12_HSA (Q8IW75)	+	+	+	
A12_MMU (Q8R4Z1)	+	+	+	
A12_RNO (Q6P6M3)	+	+	+	
HCII_HSA (P05546)	+	+	+	
HCII_MMU (P49182)	+	+	+	
HCII_RNO (Q64268)	+	+	+	
HCII_GGA (AAC16324)	+	+	+	
HCII_XTR (ENSXETP00000048524)	+	+	+	
HCII_FRU (FRUP00000149263)	+	+	+	[+]241c
HCII_TNI (GSTENP00028636001)	+	+	+	[+]241c
HCII_DRE (NP_878300)	+	+	+	
HCII_PMA (GENSCAN00000067410)	+	+	+	[+]83c
Gga-Spn-11 (XP_421342)	+	+	+	
Gga-Spn-12 (XP_421343)	+	+	+	
Gga-Spn-13 (XP_421344)	+	+	+	
Gga-Spn-14 (XP_421345)	+	+	+	
Gga-Spn-15 (XM_001235489)	+	+	+	
Xtr-Spn-8 (e_gw1.185.80.1)	+	+	+	
Xtr-Spn-9 (e_gw1.185.79.1)	+	+	+	
Xtr-Spn-10 (C_scaffold_185000011)	+	+	+	
Xtr-Spn-11 (estExt_fgenesh1_pg.C_1850042) EP45	+	+	+	
Xtr-Spn-12 (estExt_fgenesh1_pg.C_1850041)	+	+	+	
Xtr-Spn-13 (e_gw1.185.72.1)	+	+	+	
Fru-Spn-7 (FRUP00000160285) /Spn_215c	+	+	+	[+] 215c
Tni-Spn-3 (GSTENP00007903001) /Spn_215c	+	+	+	[+] 215c
Fru-Spn-17 (FRUP00000155064)	+	+	+	
Tni-Spn-4 (GSTENP00018460001)	+	+	+	
Dre-Spn-8 (NP_001071226)	+	+	+	
Dre-Spn-9 (NP_001104678)	+	+	+	
Dre-Spn-10 (NP_001099059)	+	+	+	
Dre-Spn-11 (XR_029524)	+	+	.*	
Dre-Spn-12 (XP_695000)	+	+	+	

\*pseudogene

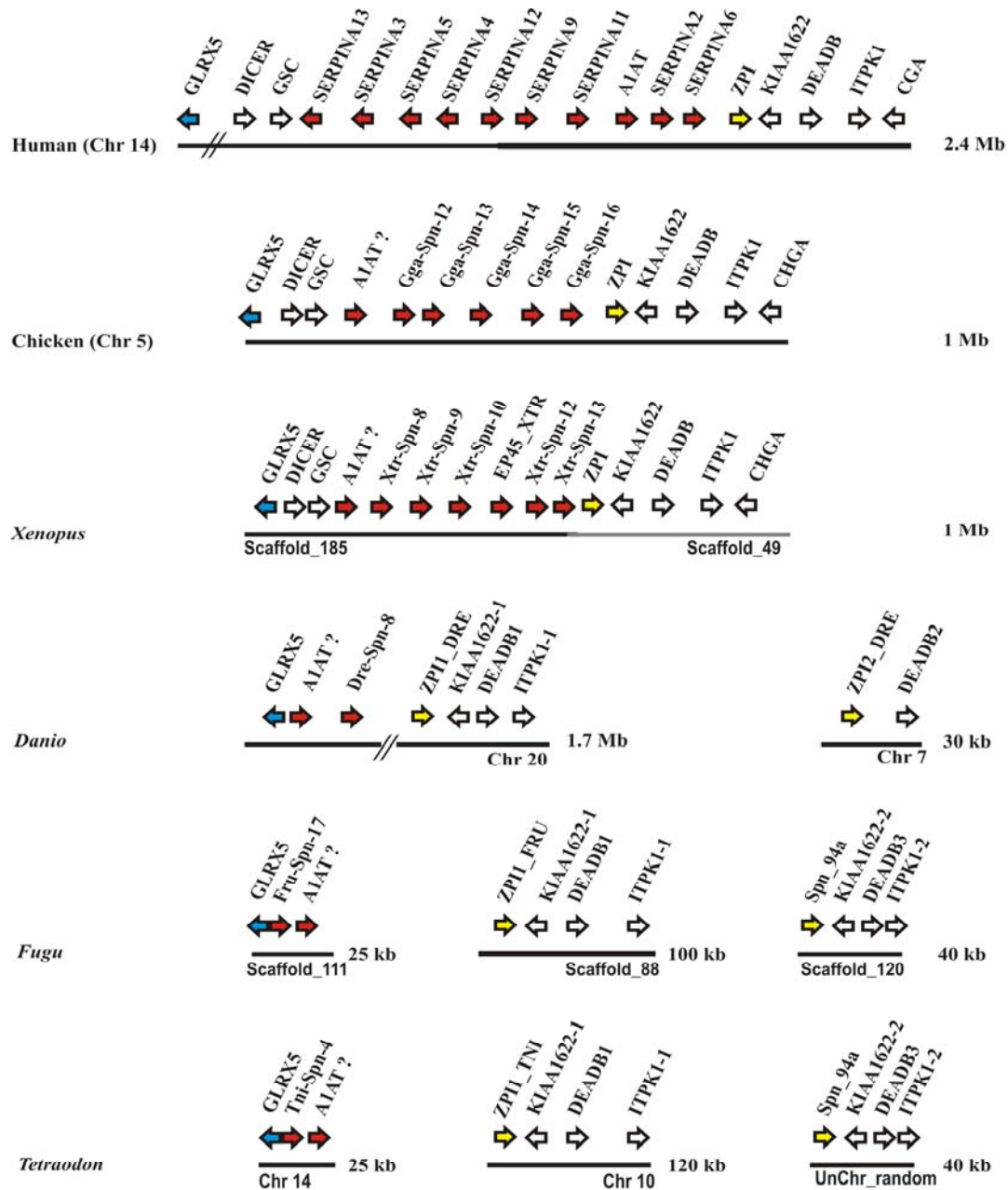
### 5.12.2. Synteny analysis of group V2 serpins in the $\alpha_1$ -antitrypsin cluster

To examine orthology of group V2 serpins, their chromosomal synteny in different vertebrates was investigated. In the human chromosome 14, a cluster of group V2 serpins ( $\alpha_1$ -antitrypsin like) containing serpins A13, A3, A5, A4, A12, A9, A11, A1, A2, A6 and A10 is present, flanked by markers GLRX5-DICER-GSC<sup>1</sup> on one side and by the triad (KIAA1622<sup>2</sup>-DEADB-ITPK1)<sup>1</sup> on the other side. A similar syntenic organization, containing a serpin gene cluster bounded by common marker sets was found to be conserved from fish to human (**Figure 45**). The ZPI gene is consistently found at the proximal end, adjacent to the (KIAA1622-DEADB-ITPK1)<sup>1</sup> cluster. The chicken chromosome 5 has seven group V2 serpins in this cluster, and these include A1AT like genes and a ZPI ortholog. The *Xenopus*

<sup>1</sup> See Appendix 8.4.3.

<sup>2</sup> KIAA1622 gene encodes a HEAT-like repeat-containing protein.

*tropicalis* cluster has eight group V2 serpin genes, including an ortholog of ZPI, the frog specific EP45 gene, and A1AT-like genes.



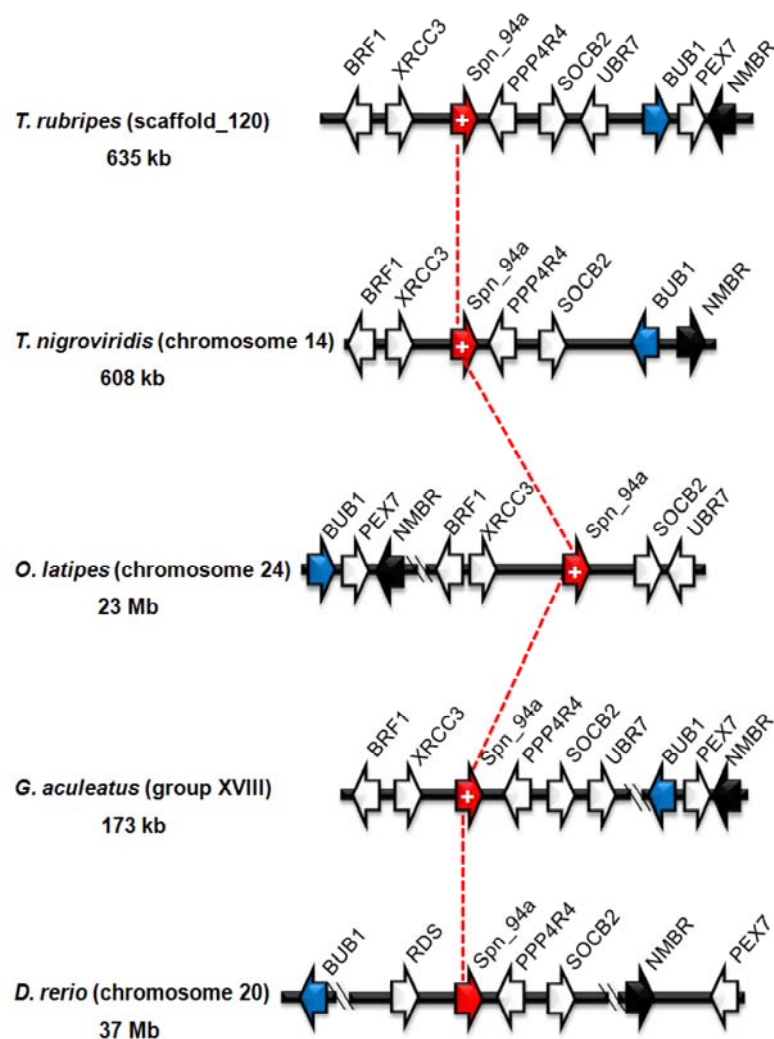
**Figure 45: Synteny of group V2 ( $\alpha_1$ -antitrypsin like) serpin genes in vertebrates.** Possible orthologs of human A1AT is marked by ? from different organisms have tentatively identified based only on RCL conservation (also see **section 5.12.7**).

In *Danio rerio*, A1AT-like genes and ZPI are found in a cluster on chromosome 20 but are separated by other markers in between. The genomes of *Fugu* and *Tetraodon* have one extra non-inhibitory serpin in this cluster, namely Fru-Spn-17 and Tni-Spn-4, respectively, along with A1AT-like gene flanked by marker GLRX5 (blue) on one side. These fish genomes have a paralogous genomic fragment a ZPI-like gene (ZPI2\_DRE) and a DEADB-like gene (DEADB2) can be located. This suggests that fishes have two types of ZPI like genes - ZPI1

is orthologous to human ZPI and is found in all three fishes. Furthermore, selected fishes have Spn\_94a serpin, which possess sequence similarity to ZPI.

The subset of group V2 serpins that is reversely oriented with respect to A1AT gene in the human serpin gene cluster is not found in any of non-mammalian vertebrates. This suggests that these serpins (A3-A5 and A13) are specific to mammals.

To understand origin of Spn\_94a genes in fishes, synteny analysis of selected fish genome was carried and orthologs were identified from selected fish (**Figure 46**). It becomes evident that this gene is found in different ray-finned fishes, however, extra intron at position 94a is found in all fishes except for *D. rerio*. It suggests that this intron is inserted after divergence of the *D. rerio* lineage.



**Figure 46: *Spn\_94a* orthologs unraveled by chromosomal gene order from selected fishes.** With the exception of *Danio rerio*, all fishes investigated share a gene with an extra intron at position 94a (indicated by a plus sign). Chromosomal gene order corroborates that these genes, dubbed *Spn\_94a*, are orthologous. Intron gain hence took place after divergence of the *D. rerio* lineage.

### 5.12.3. Synteny analysis of the *serpinA7* gene

To detect the orthologs of human *serpinA7*, which located on the X-chromosome, the *serpinA7* micro-environment was investigated. In mammals, *serpinA7* is flanked by the *ILIRAPL2*-*NRK* (**appendix 8.4.4**) gene cluster on one side and by marker on the other side (**Figure 47**). This architecture is not found in any of non-mammalian vertebrates.



Figure 47: Synteny organization of *serpinA7* gene in mammalian genomes.

### 5.12.4. Synteny analysis of the angiotensinogen (AGT) gene

In humans and in chicken, the *AGT* (*serpinA8*) gene is flanked by the *COG2* marker on one side and *CAPN9* on the other side. A similar genomic architecture is maintained in frog and in fishes, but only *COG2* marker was found in these species (**Figure 48**).

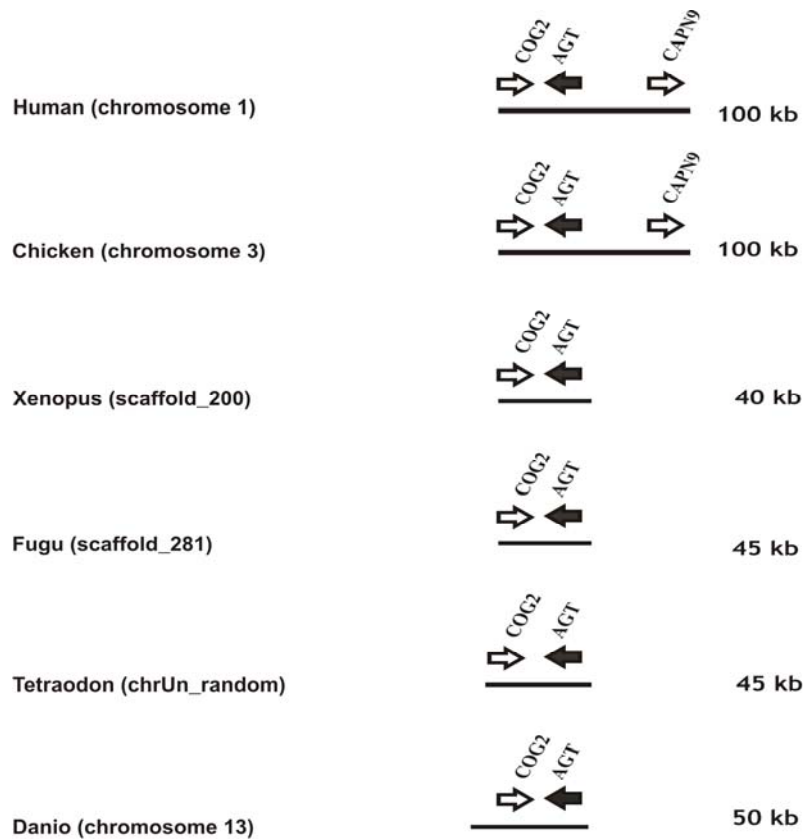


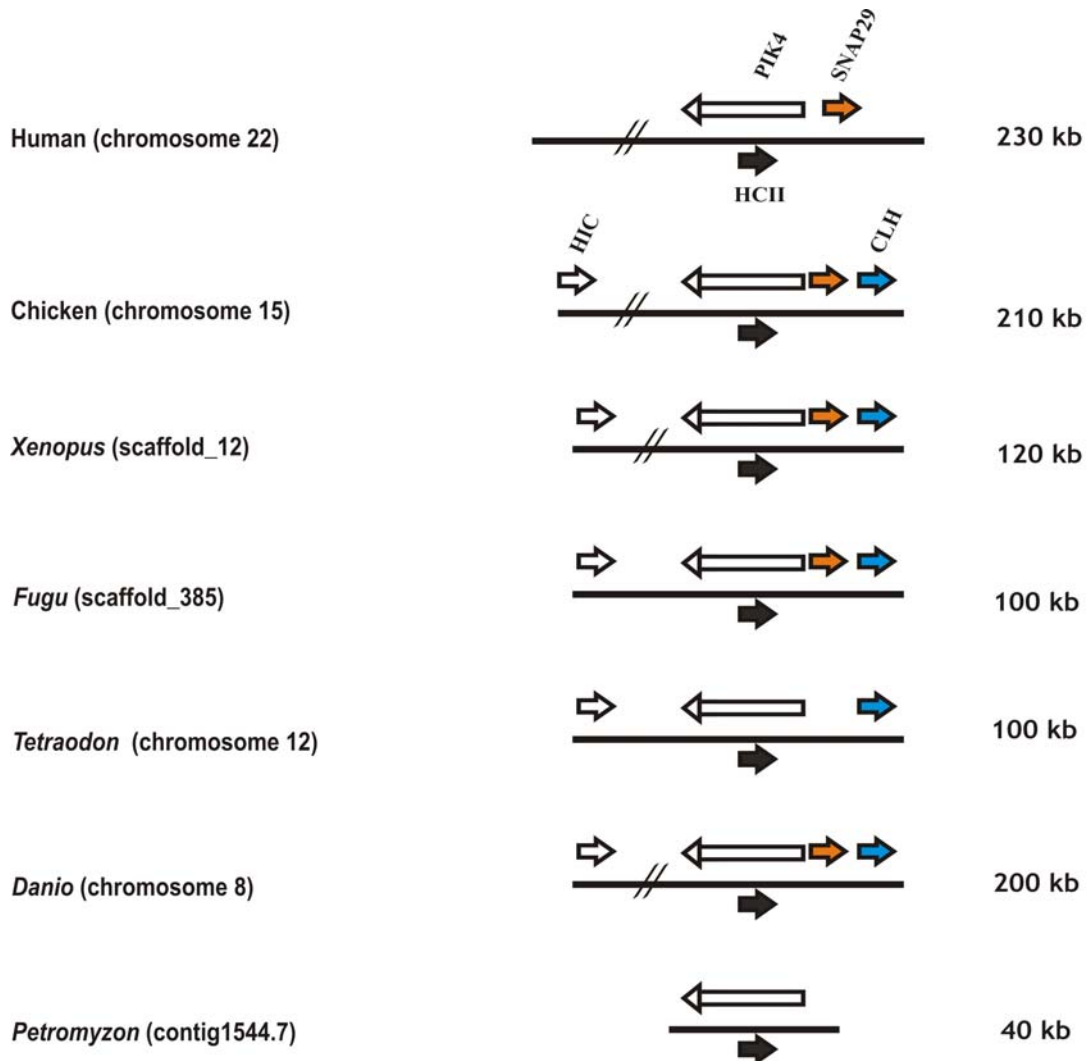
Figure 48: Synteny of the angiotensinogen (*AGT*) genes in vertebrate genomes.



These data suggest that orthologs of human AGT are maintained from fishes to mammals.

### 5.12.5. Synteny analysis of the heparin cofactor II (HCII) gene

To unravel orthologs of heparin cofactor II, the genomic environment of the HCII gene was compared (**Figure 49**).



**Figure 49: Synteny analysis of the heparin cofactor II (HCII) gene in vertebrates.** The HCII gene is consistently located within an intron of the PIK4 gene in reverse orientation, suggesting that HCII gene has been continuously maintained since divergence of lampreys.

The HCII gene in vertebrates is found to be conserved as *gene within gene* located in an intron of the PIK4 gene (in opposite orientation) and a common set of flanking markers. This conserved syntenic organization corroborates that orthologs of human HCII gene are found across vertebrates.

### 5.12.6. Genomic organization of fish specific group V2 serpins

There are two different fish specific genomic organizations of some serpin genes, not evident in any other vertebrates investigated. Based on sequence features, Fru-Spn-7 and Tni-Spn-3 are close homologs, and these genes share an additional intron at position 215c, thus these genes were renamed as Spn\_215c. Nevertheless, a syntenic localization could not be confirmed, due to lack of conserved markers in scaffold\_5339 of the *Fugu* genome (**Figure 50**).

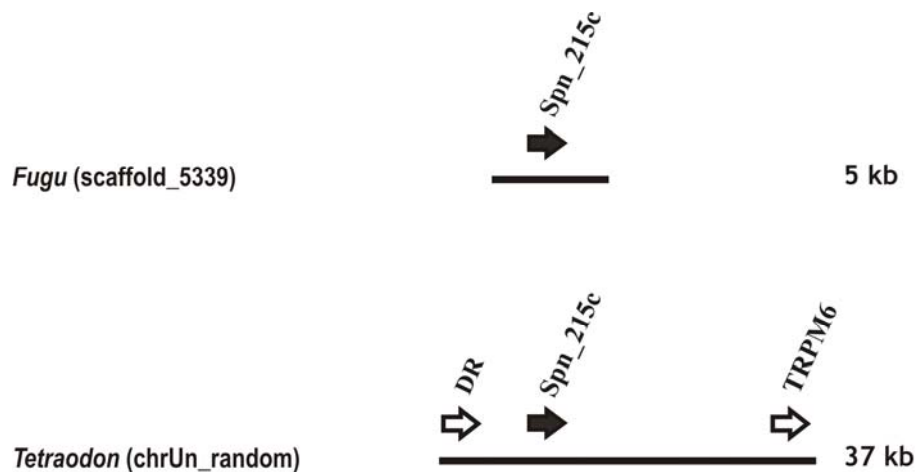


Figure 50: Genomic organization of the fish specific group V2 serpin – Spn\_215c

The *Danio* genome depicts four group V2 serpin genes in a unique syntenic organization on chromosome 5 (**Figure 51**).

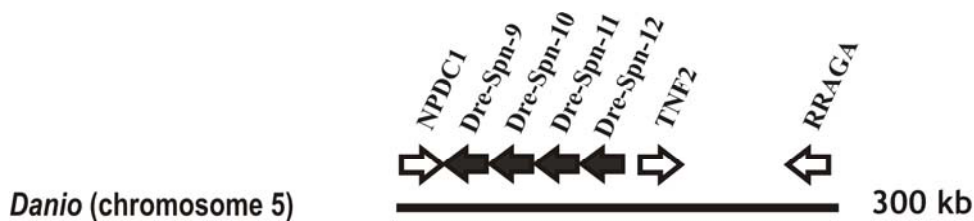


Figure 51: Genomic organization of the *Danio* specific group V2 serpin genes – Dre-Spn-9, Dre-Spn-10, Dre-Spn-11, and Dre-Spn-12.

### 5.12.7. Sequence analysis of group V2 serpins

In order to expand the understanding of orthologs and paralogs of different group V2 serpins, sequence analysis of group V2 serpins was carried out. The major findings are described as below.

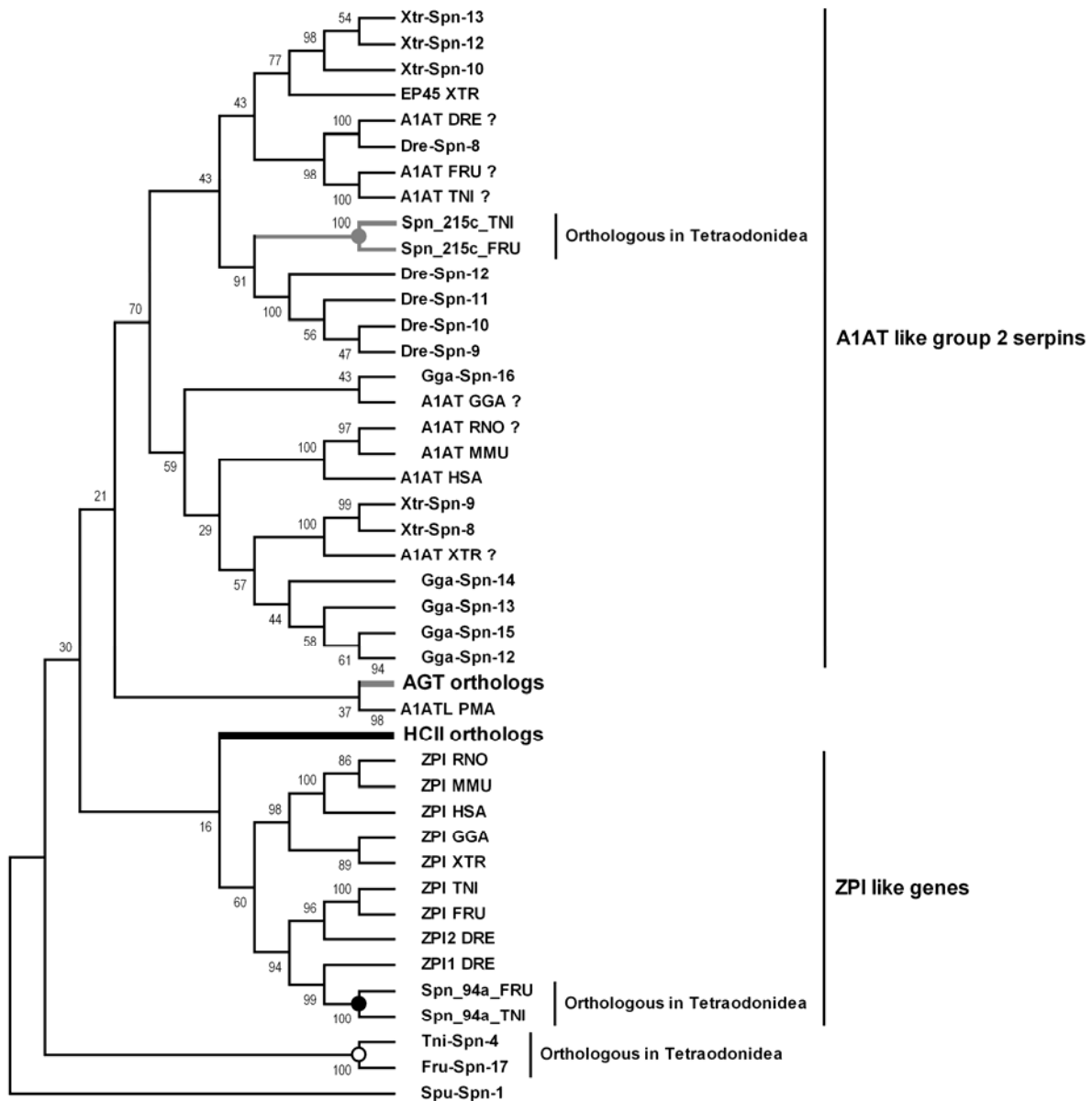
The AGT protein is maintained in vertebrates with 23-62% sequence identity and 42-75% sequence similarity between humans and fishes. *Tetraodon* AGT (AGT\_TNI) is only partially available due to the presence of a big gap in the genomic sequence after the intron at position 192a. The AGT proteins are characterized by the presence of the highly conserved angiotensin sequence close to the N-terminal end (cyan boxes in **appendix 8.3.12**) and the non-inhibitory RCL (red boxes in **appendix 8.3.12**).

The HCII protein is highly conserved in vertebrates with 39-81% sequence identity and 63-88% sequence similarity between humans and fishes. The heparin binding helix-D of HCII was found to be highly conserved (yellow boxes in **appendix 8.3.13**) which therefore represent a signature sequence for identifying HCII orthologs. The RCL region is highly conserved (red boxes in **appendix 8.3.13**).

The ZPI is conserved in vertebrates as orthologs and paralogs of human ZPI show 22-71% sequence identity and 44-81% sequence similarity at amino acid level with human ZPI from fish to mammals and possess inhibitory RCL, except in ZPI\_FRU and ZPI\_TNI (red boxes in **appendix 8.3.14**). This suggests that ZPI\_FRU and ZPI\_TNI acquired non-inhibitory function after duplication of micro-environment in these fish lineage.

Serpins possessing additional introns at position 215c – Spn\_94a\_FRU and Spn\_94a\_TNI share 73% sequence identity and 86% sequence similarity with each other and possess inhibitory RCL (**appendix 8.3.15**). Fru-Spn-17 and Tni-Spn-4 share 71% sequence identity and 79% of sequence similarity with each other and possess a non-inhibitory RCL (**appendix 8.3.16**). Unraveling orthology of the A1AT gene is a challenging task, since A1AT-like genes expanded by tandem duplication resulting in many A1AT-like genes in different organisms. Chicken, *Xenopus* and *Danio*, for instance, have six A1AT like genes (**appendices 8.3.17-8.3.19**). Based on RCL sequence conservation, A1AT proteins from each species were picked to examine orthology of human A1AT. However, orthology assignment based on RCL sequence alone, this makes a weak argument in favour. Therefore, these sequences are aligned in **appendix 8.3.20** as tentative A1AT orthologs. The thyroxine hormone binding globulin (THBG/serpinA7) is highly conserved in mammals and it show 75-76% sequence identity and 88% sequence similarity with human THB in addition, it possesses a non-inhibitory RCL (**appendix 8.3.21**).

To complement the understanding about group V2 serpins in different vertebrates, a phylogenetic tree (**Figure 52**) based on the NJ method (Saitou and Nei, 1987) was created with the help of MEGA4 (Tamura *et al.*, 2007). This phylogenetic tree has several major branches such as  $\alpha_1$ -antitrypsin like genes including  $\alpha_1$ -antitrypsin from different vertebrates, HCII orthologs, AGT orthologs, ZPI and its orthologs and paralogs; in addition, some genes are identified as *Danio* specific and *Tetraodontidae* specific. Presence of two group V2 serpins from initial version of lamprey genome, namely HCII and  $\alpha_1$ -antitrypsin like genes with group V2 specific gene structure, suggests that group V2 serpins arose at early vertebrate emergence.



**Figure 52: Phylogenetic tree of group V2 serpins from different vertebrates.**  $\alpha_1$ -antitrypsin like genes cluster together whereas ZPI like genes cluster in separate branch. There are three different tetraodontidae specific orthologous group V2 serpins, marked by circles of different colors. HCII and AGT branches are condensed to two thick single lines to make this tree simple. A1AT orthology assignment is challenging and therefore marked by ?. This tree was created with MEGA4 based on the NJ method. Bootstrap values (in percentage) for 1000 replicates are shown.

In summary, orthology of group V2 serpins can be assigned to ZPI, AGT and HCII genes across vertebrates whereas orthology assignment of A1AT is challenging, as different species contains many A1AT like genes. There is only one cluster of  $\alpha_1$ -antitrypsin like genes in mammals, chicken, and frog but fishes do possess additional clusters of  $\alpha_1$ -antitrypsin like genes. There are many group V2 serpin genes, which originated based on organism's requirements.

### 5.13. Orthology analysis of group V3 serpins

Group V3 of vertebrate serpins has been defined by a gene structure having seven introns at positions - 86a/88a or 90a<sup>1</sup>, 167a<sup>2</sup>, 230a, 290b, 323a, 352a and 380a ( $\alpha_1$ -antitrypsin numbering) in their coding regions<sup>3</sup> (Ragg *et al.*, 2001). The exact location of the first intron is uncertain in different group V3 family members due to alignment ambiguities. Group V3 has five inhibitory serpins as members that are involved in different physiological processes (Table 28).

Table 28: Physiological roles of group V3 serpins and associated diseases/syndromes.

Group V3 serpins	Physiological Role(s)	Associated Disease(s)/Syndrome(s)
SerpinE1 (Plasminogen activator inhibitor 1)	Inhibitor of plasminogen activation regulates tissue plasminogen inhibitor (tPA), urokinase plasminogen activator (uPA), and protein C.	Atherosclerosis, diabetes and Hypertension
SerpinE2 (Glia derived nexin /GDN)	Potent inhibitor of thrombin in central nervous system and in the vasculature.	
SerpinE3	Unknown	
SerpinI1 (Neuroserpin)	Inhibitor of tissue plasminogen activator in the nervous system.	Ischemia and FENIB <sup>4</sup>
SerpinI2 (Pancpin)	Inhibitor of cancer metastasis.	Cancer

The protein alignments of group V3 serpins are shown in **appendices 8.3.22 to 8.3.26**.

#### 5.13.1. Gene structure of group V3 serpins

Since gene structures are discriminatory features of group V1-V6 serpins, gene architectures of probable group V3 serpin homologs in different vertebrates were determined. **Table 28** shows that all vertebrates investigated contain at least three genes that depict the basic exon-intron structure of group V3 serpin genes. Introns at the canonical positions 86a/88a/90a, 167a, 230a, 290b, 323a, 352a, and 380a are conserved with some deviations. In the PAI1 genes of *Xenopus* (PAI1\_XTR) and *Tetraodon* (PAI1\_TNI) no intron at position 380a was found (indicated by ?), probably due to sequencing errors. PAI1 was not found in chicken, due to either a bird specific gene loss or alternatively, PAI1 in chicken escaped detection. Due to lack of data for the *Fugu* serpinE3 gene (E3\_FRU)<sup>5</sup>, not all introns can be assigned in this gene (indicated by ?) and for the same reasons, introns at positions 352a and 380a of the *Tetraodon* serpinE3 gene (E3\_TNI) cannot be assigned. Similarly, the intron at position 90a for NEURO\_TNI cannot be located, due to sequencing errors. Pancpin is only found in mammals and in the frog. There are two novel introns at positions 205b and 217a in PANC\_XTR. The gene structures of serpinE2 and serpinI1 perfectly correspond to the canonical group V3 gene structure.

<sup>1</sup> Tentative positions due to alignment ambiguities.

<sup>2</sup> Also shared by group V1 serpins.

<sup>3</sup> Out of seven intron positions, the last six are found at identical locations

<sup>4</sup> Familial encephalopathy with neuroserpin inclusion bodies.

<sup>5</sup> Mapping of intron positions was possible with use of trace archive data from Ensembl (accession id, SINFRUG00000134592.1).

**Table 28: Intron positions of group V3 genes in vertebrates.** The presence (+) or absence (-) of intron positions is shown. PAI1\_XTR and PAI1\_TNI lack the intron at position 380a, due to sequencing errors and similarly for the same reasons, the intron at position 90a in NEURO\_TNI (indicated by ?). Gene structures coding for E3\_FRU and E3\_TNI are incomplete (indicated by ?). PANC\_XTR has two novel introns at positions 205b and 217a (see in text).

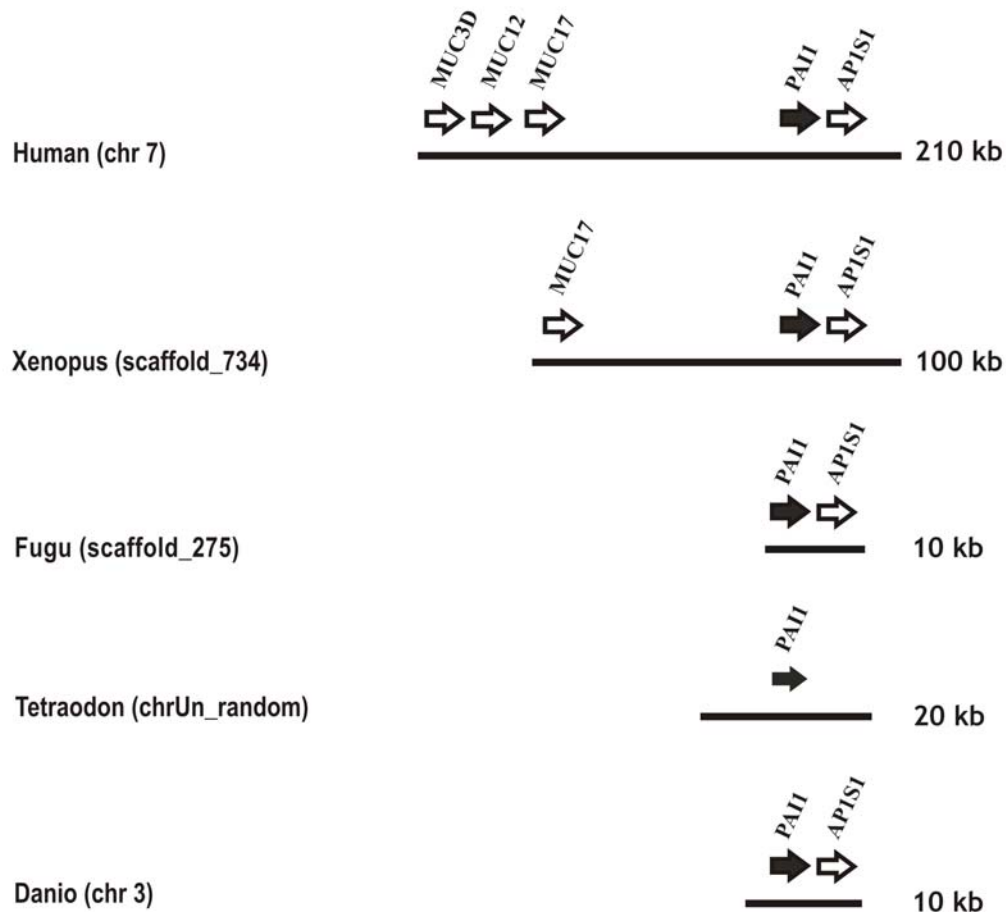
Group V3 serpin gene	Intron at position						
	86a <sup>1</sup>	167a	230a	290b	323a	352a	380a
<b>SerpinE1 (PAI1)</b>							
PAI1_HSA (P05121)	+	+	+	+	+	+	+
PAI1_MMU (P22777)	+	+	+	+	+	+	+
PAI1_RNO (P20961)	+	+	+	+	+	+	+
PAI1_XTR (estExt_Genewise1.C_7340032)	+	+	+	+	+	+	?
PAI1_FRU (e_gw2.275.54.1)	+	+	+	+	+	+	+
PAI1_TNI (GSTENT00003787001)	+	+	+	+	+	+	? <sup>2</sup>
PAI1_DRE (XP_690192)	+	+	+	+	+	+	+
<b>SerpinE2 (GDN)</b>	86a <sup>1</sup>	167a	230a	290b	323a	352a	380a
GDN_HSA (P07093)	+	+	+	+	+	+	+
GDN_MMU (Q07235)	+	+	+	+	+	+	+
GDN_RNO (P07092)	+	+	+	+	+	+	+
GDN_GGA (gi:50730899)	+	+	+	+	+	+	+
GDN_XTR (fgenes1_kg.C_scaffold_750000001)	+	+	+	+	+	+	+
GDN_FRU (e_gw2.123.110.1)	+	+	+	+	+	+	+
GDN_TNI (GSTENP00026727001)	+	+	+	+	+	+	+
GDN_DRE (Q7ZVL5)	+	+	+	+	+	+	+
<b>SerpinE3</b>	86a <sup>1</sup>	167a	230a	290b	323a	352a	380a
E3_HSA (XM_941682)	+	+	+	+	+	+	+
E3_MMU (AK053602)	+	+	+	+	+	+	+
E3_RNO (gi:109501642)	+	+	+	+	+	+	+
E3_GGA (XM_417070)	+	+	+	+	+	+	+
E3_XTR (e_gw1.233.93.1)	+	+	+	+	+	+	+
E3_FRU (FRUP00000142610)	?	?	+	+	+	+	?
E3_TNI (GSTENT00029213001)	+	+	+	+	+	?	?
E3_DRE (ENSDARP00000074162)	+	+	+	+	+	+	+
<b>SerpinI1 (Neuro)</b>	90a <sup>1</sup>	167a	230a	290b	323a	352a	380a
NEURO_HSA (Q99574)	+	+	+	+	+	+	+
NEURO_MMU (O35684)	+	+	+	+	+	+	+
NEURO_RNO (Q5M7T5)	+	+	+	+	+	+	+
NEURO_GGA (gi:521387191)	+	+	+	+	+	+	+
NEURO_XTR (ENSXETP00000049461)	+	+	+	+	+	+	+
NEURO_FRU (fgh5_pm.C_scaffold_488000001)	+	+	+	+	+	+	+
NEURO_TNI (GSTENP00034604001)	?	+	+	+	+	+	+
NEURO_DRE (ENSDARP00000017430)	+	+	+	+	+	+	+
<b>SerpinI2 (Panc)</b>	90a <sup>1</sup>	167a	230a	290b	323a	352a	380a
PANC_HSA (O75830)	+	+	+	+	+	+	+
PANC_MMU (Q9JK88)	+	+	+	+	+	+	+
PANC_RNO (gi:16758618)	+	+	+	+	+	+	+
PANC_XTR (ENSXETP00000049481)	+	+	+	+	+	+	+

<sup>1</sup> Tentative position, due to sequence ambiguities.

<sup>2</sup> 3' end of PAI1\_TNI is not present in databases.

### 5.13.2. Synteny analysis of PAI1 genes

To unveil orthology of vertebrate PAI1 genes, the syntenic arrangements in vertebrates were analyzed (**Figure 53**).



**Figure 53: Genomic localization of PAI1 genes in vertebrates.**

In humans, the PAI1 gene is found on chromosome 7, flanked by AP1S1<sup>1</sup> on one side and a gene cluster (MUC3D, MUC12 and MUC17)<sup>1</sup> on the other side. In *Xenopus tropicalis*, a similar syntenic organization is evident. In fishes, only the AP1S1<sup>2</sup> marker is found to flank PAI1. Together these data suggest that orthologs of human PAI1 are retained from fish to mammals.

<sup>1</sup> Appendix 8.4.5.

<sup>2</sup> AP1S1 marker is not detectable in *Tetraodon*, due to sequencing errors.

### 5.13.3. Synteny analysis of GDN genes

In humans, GDN is flanked by AP1S3<sup>1</sup> on one side and CUL3<sup>1</sup> on the other side (**Figure 54**). A similar syntenic organization is found in chicken and in the frog. In fishes, the linkage of GDN and AP1S3 is maintained, but instead of CUL3, the S28<sup>1</sup> gene is found on one side as marker.

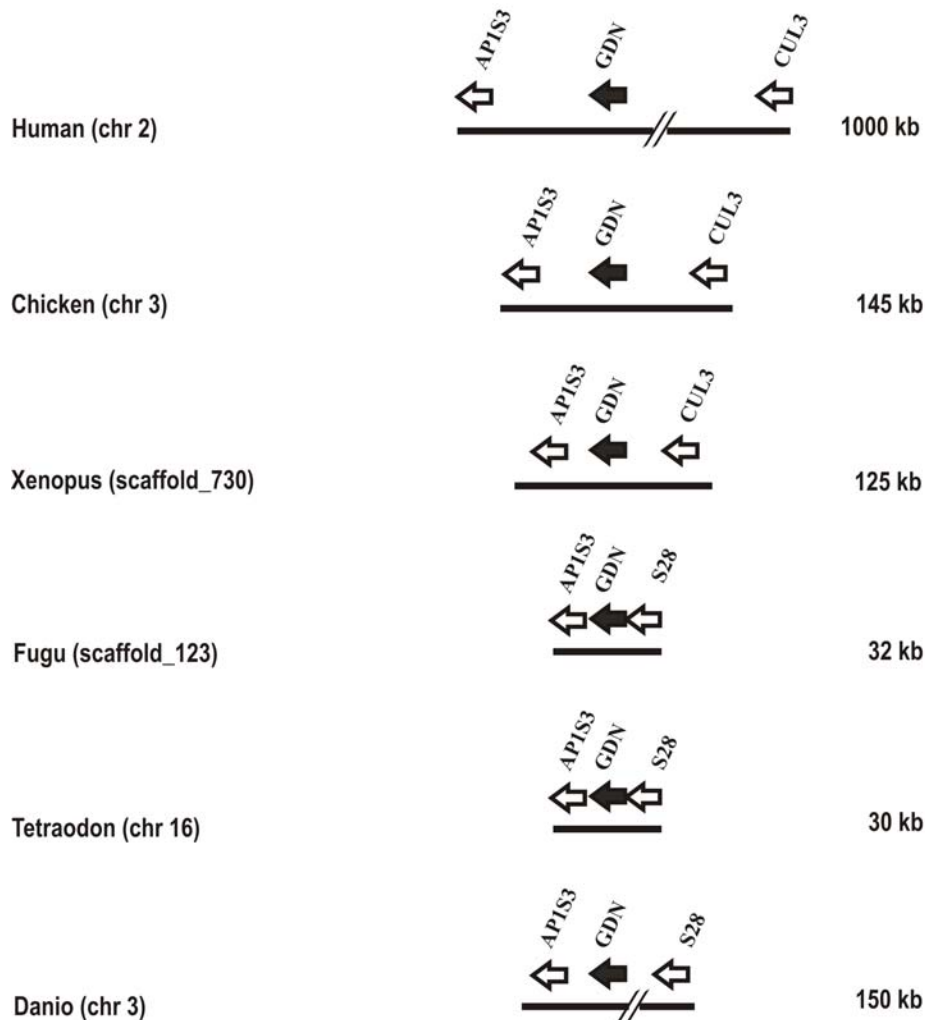


Figure 54: Genomic localization of GDN genes in vertebrates.

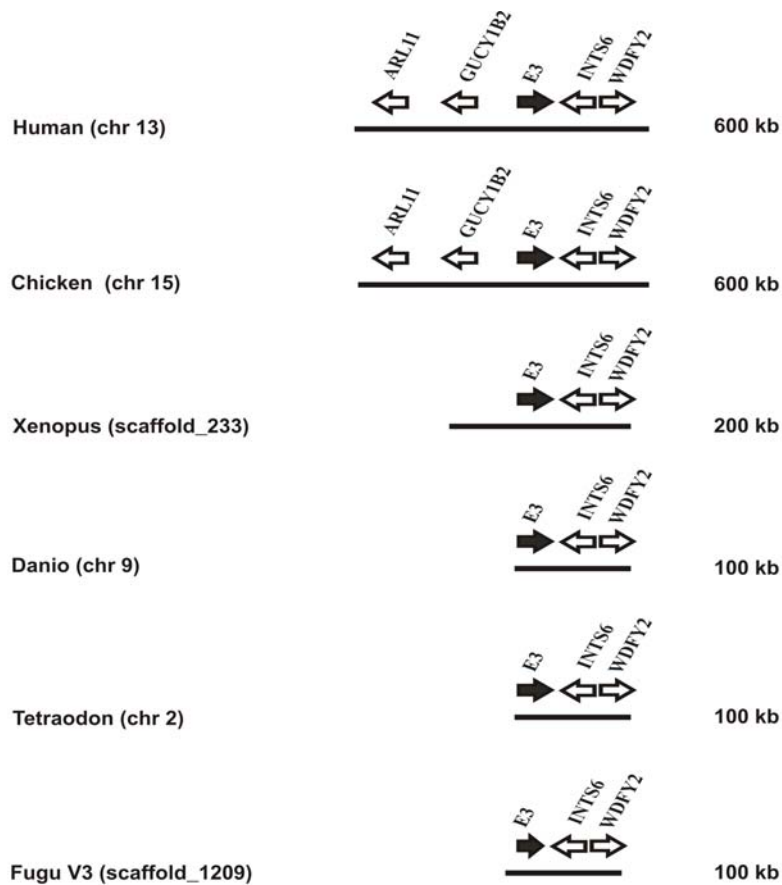
### 5.13.4. Synteny analysis of serpinE3 genes

To unravel the orthology of the serpinE3 gene, its genomic micro-environment was investigated (**Figure 55**). In humans and in chicken, the serpinE3 gene is flanked by the ARL11-GUCY1B2<sup>1</sup> cluster on one side and by INTS6-WDFY2<sup>1</sup> on the other side. In the frog and in fishes, a similar syntenic architecture is maintained, but only the marker genes INTS6-WDFY2 are found to be conserved on one side. Due to sequencing errors in current versions of genomic sequences of *Fugu* (versions V3 and V4) and *Tetraodon* (version V7), complete serpinE3 gene sequences cannot be located. However, SerpinE3 can be identified in two other

<sup>1</sup> Appendix 8.4.5.



fish genomes - Medaka<sup>1</sup> and stickleback<sup>2</sup> where it is arranged in a similar syntenic organization (not shown). This supports that serpinE3 is conserved in different vertebrates.



**Figure 55: Genomic localization of serpinE3 genes.** Two versions (V3 and V4) of Fugu genomic sequences were used to deduce the synteny.

### 5.13.5. Synteny analysis of neuroserpin and pancpin genes

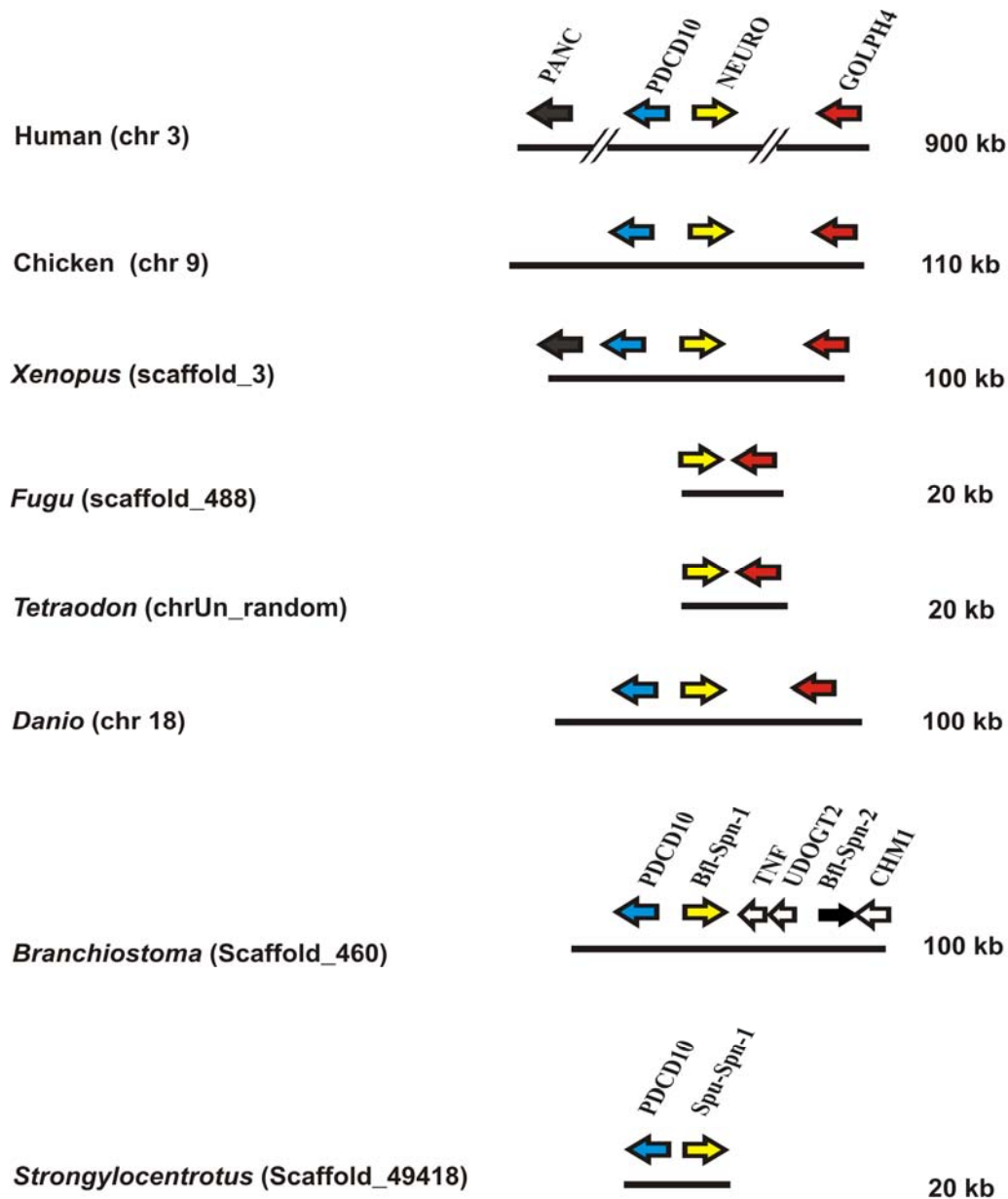
**Figure 56** shows the syntenic architectures of neuroserpin and pancpin genes across vertebrates. The genes coding for neuroserpin and pancpin are found in a cluster on chromosome 3 in humans separated by 261 kb. They are flanked by marker genes PDCD10 on one side and GOLPH4 on the other side. This synteny is found in all vertebrates investigated with some deviations. The pancpin gene is missing in chicken and in fishes, and the PDCD10 marker was not found in *Fugu* and *Tetraodon*. However, it is present in the *Danio* genome.

The detection of possible ancestors of vertebrate serpins is a major aim of this work. Importantly, Bfl-spn-1 from lancelet and Spu-spn-1 from sea urchin, respectively, show an arrangement comparable to that of vertebrate neuroserpin-pancpin cluster. The highly

<sup>1</sup> *Oryzias latipes* [Ensembl peptide id ENSORLP00000012629]

<sup>2</sup> *Gasterosteus aculeatus* [Ensembl peptide id ENSGACP00000000316]

conserved PDCD10 marker gene is found in a head-to-head orientation to these serpin genes in a similar fashion as in the vertebrate genomes. These relationships are investigated further in **section 5.13.6**.



**Figure 56: Genomic localization of neuroserpin and pancpin genes in vertebrates and comparative analysis of micro-synteny with serpins of higher invertebrates – Bfl-spn-1 (*B. floridae*) and Spu-spn-1 (*S. purpuratus*).** The neuroserpin gene in vertebrates is consistently found associated with the PDCD10 gene, which is highly conserved in all eukaryotes. In lancelets and sea urchins, the PDCD10 gene is found adjacent to Bfl-spn-1 and Spu-spn-1, respectively. The distance between neuroserpin and pancpin genes in the human genome is 261 kb.

### 5.13.6. Sequence analysis of group V3 serpins

PAI1 is conserved in vertebrates, depicting 38-80% sequence identity and 59-95 % sequence similarity on the amino acid level with human PAI1. The inhibitory RCL region is conserved containing R-M at P1-P1' (**appendix 8.3.22**).

GDN is also highly conserved in vertebrates and it shows 51-84% sequence identity and 70-93% sequence similarity with human GDN. The helix-D region is highly conserved among GDN orthologs of different vertebrates and an N-glycosylation site<sup>1</sup> (positions 163-165<sup>2</sup>) is conserved. The inhibitory RCL region is also strongly conserved (**appendix 8.3.23**).

SerpineE3 is maintained in vertebrates, show 27-64% sequence identity, and 37-74 % sequence similarity on the amino acid level with human serpinE3. The inhibitory RCL region is conserved and contains a cluster of hydrophobic amino acid preceding the presumptive P1 position (**appendix 8.3.24**).

Pancpin orthologs are only found in mammals and in *Xenopus*, showing 49-76% sequence identity and 68-88% sequence similarity on the amino acid level. The C-terminal end is strongly maintained (**appendix 8.3.25**).

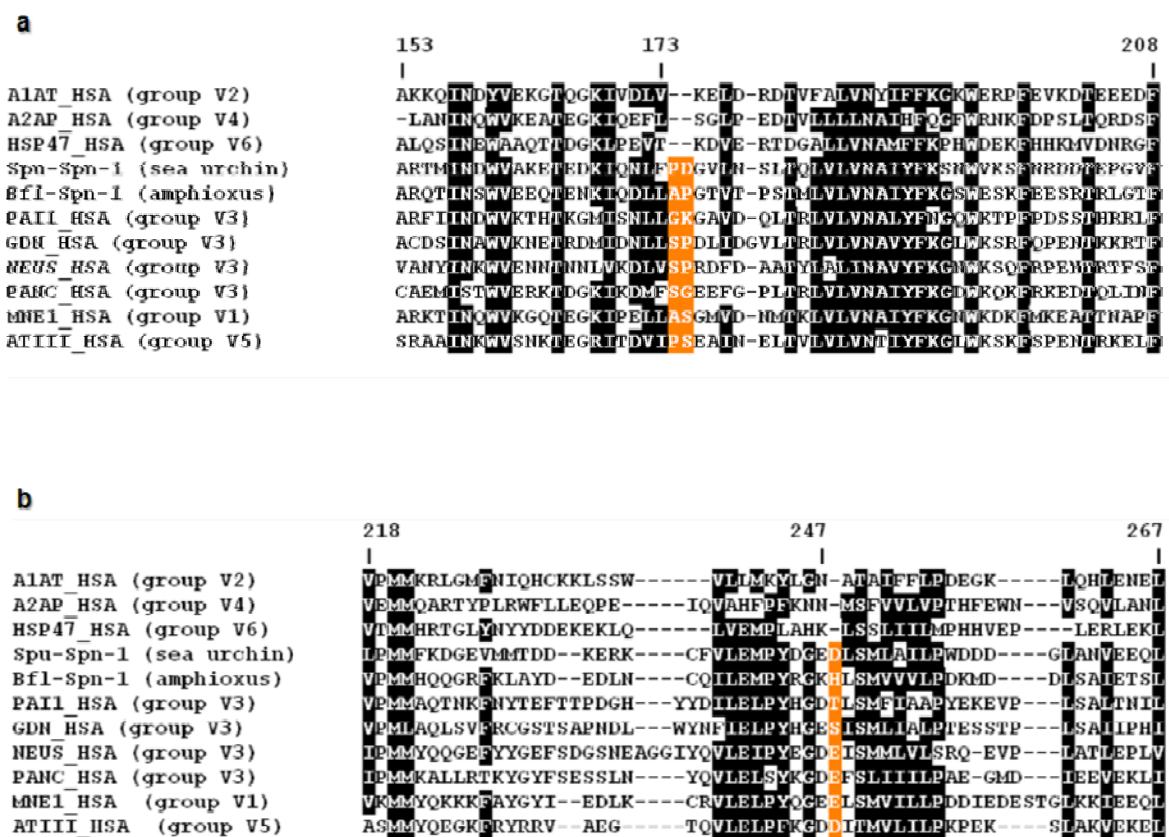
The neuroserpin gene is highly conserved in vertebrates, and the protein shows 47-81% sequence identity and 65-95% sequence similarity with the human ortholog. The inhibitory RCL region always contains an R at P1. An N-glycosylation signal (residues 163-165) is conserved. A C-terminal extension shown to direct neuroserpin to the regulated secretory pathway (Ishigami *et al.*, 2007) is strongly conserved (**appendix 8.3.26**).

An ancestor of neuroserpin is found in sea urchins based on synteny analysis (**section 5.13.4**). To explore these relationships further sequence comparisons were carried out. The *Spu-spn-1* gene has no introns in the conserved part of the serpin domain and it contains a single intron in the signal peptide. The *Bfl-spn-1* gene has only two introns at positions 75c and 174a.

There are three discriminating indels shared between group V1 and group V3 serpins, namely (a) two amino acids between position 171/172 or alternatively 173/174 based on serpin sequences used for protein alignment, (b) one amino acid position 247/248 and (c) an intron present at position 167a (Ragg *et al.*, 2001). First two of these indels are also shared by group V5 and group V6 serpins share second of these discriminating indels. These are maintained in all group V3 serpins from fish to mammals (indicated by \* in **appendices 8.3.22 to 8.3.26**).

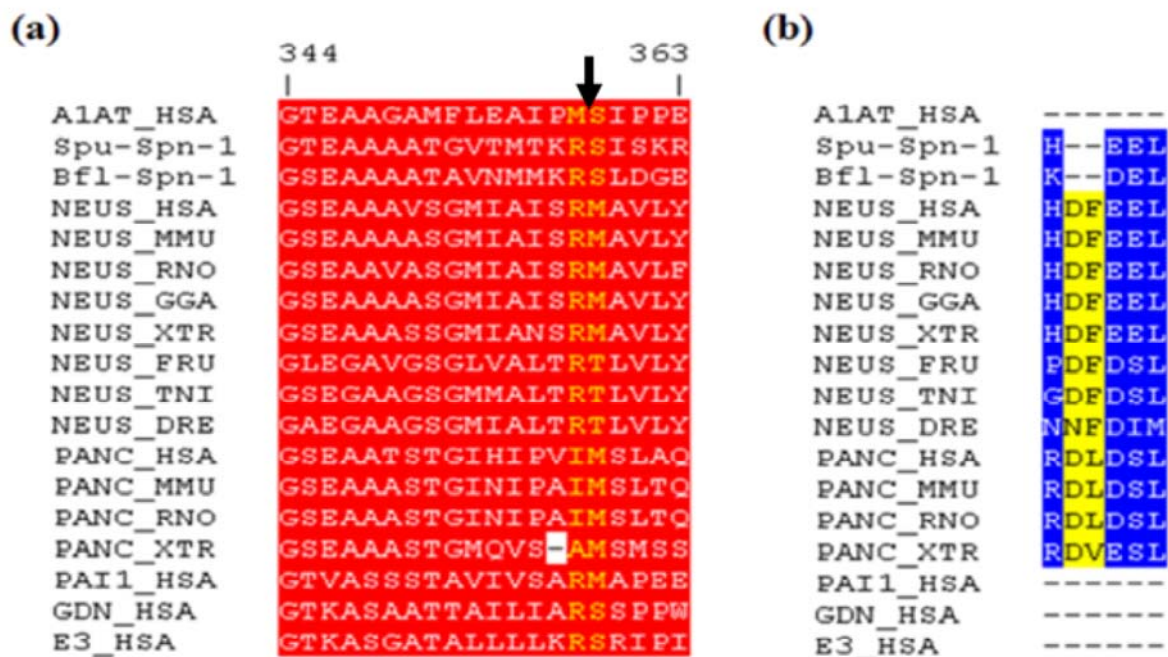
<sup>1</sup> N-glycosylation site, NX[ST], where X = any amino acid except P.

<sup>2</sup>  $\alpha_1$ -antitrypsin numbering.



**Figure 57: Comparison of discriminating amino acid indels among selected human serpins and invertebrate serpins.** (a) Two amino acids between positions 173/174 are found in vertebrate groups V1, V3, and V5. (b) Characteristic insertion of one amino acid position between positions 247/248 is maintained in vertebrate serpin groups V1, V3, V5, and V6. These two discriminatory indels are also found in serpins - Bfl-spn-1 and Spu-spn-1 from lancelet and sea urchins, respectively. This supports these invertebrate serpins are closely related to group V3 serpins of vertebrates as evident from synteny analysis (section 5.13.5). The numbering of amino acids refers to mature human A1AT.

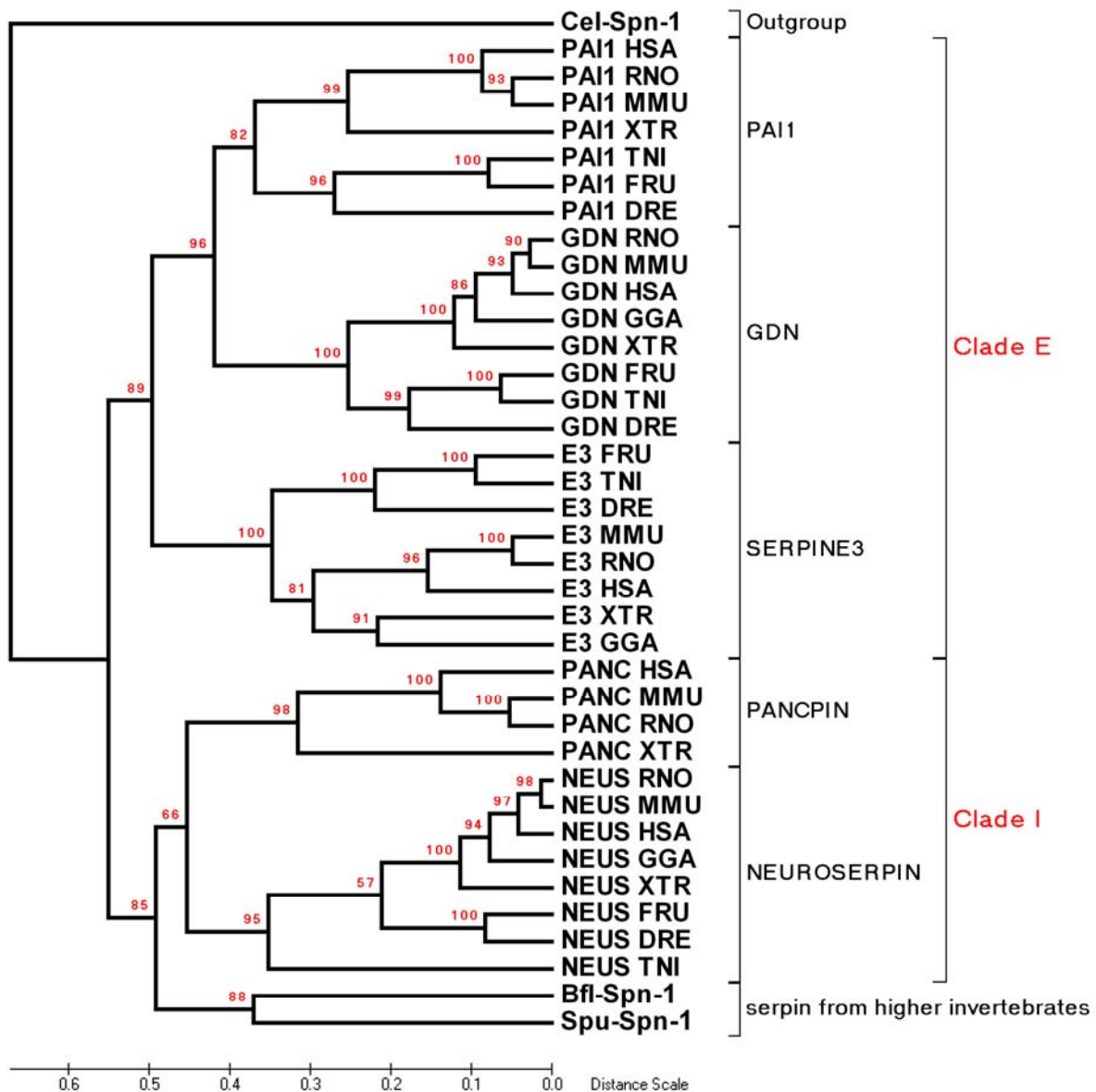
Investigating the presence of these indels in Bfl-spn-1 and Spu-spn-1, it was found that the two amino acid insertions between positions 173/174 (Figure 57a) and the insertion of one amino acid between positions 247/248 (Figure 57b) are maintained in these invertebrate serpins. Additionally, there are two common features conserved from sea urchin to humans – P1-P1' positions in the RCL (Figure 58a) and a conserved C-terminal extension (Figure 58b) on the sequence level.



**Figure 58: Sequence comparisons among selected group V3 serpins and invertebrate serpins.** (a) RCL is inhibitory among different group V3 serpins and Bfl-spn-1 and Spu-spn-1 from lancelet and sea urchins, respectively. (b) C-terminal ends of neuroserpins, pancpins, Bfl-spn-1 and Spu-spn-1 share a conserved extension. This supports these invertebrate serpins are closely related to group V3 serpins of vertebrates as evident from synteny analysis (section 5.13.5) and indel analysis. P1-P1' positions are marked in yellow and arrow indicates the cleavage site. The numbering of amino acids refers to mature human A1AT.

This suggests that the sea urchin genome harbors a close relative of the possible ancestor of modern day group V3 vertebrate serpins and this ancestor serpin can be dated back about 550 million years (Figure 6), when echinoderms separated (Sodergren *et al.*, 2006). This neuroserpin-like gene (Spu-spn-1) has no introns in conserved part of the serpin domain. Possibly, present day vertebrate group V3 serpins were created by massive intron insertion events at the time point of vertebrate emergence. The orthologous serpin gene (Bfl-Spn-1) has only two introns that, however, do not match with group V3 serpins gene architecture. Alternatively, introns might have been lost in the sea urchin or in the lancelet.

To complete orthology analysis of group V3 serpins, a phylogenetic tree (Figure 59) based on the UPGMA method (Sneath and Sokal, 1973) was constructed with help of MEGA4 (Tamura *et al.*, 2007). Group V3 serpins cluster into two major branches, constituting clades E and I as suggested (Silverman *et al.*, 2001). The serpins of invertebrates (Bfl-spn-2 and Spu-spn-1) group in the branch of clade I (neuroserpin-pancpin). This corroborates that present day clade I serpins are derived from the PDCD10-serpin locus of invertebrates.



**Figure 59: Evolutionary tree of group V3 serpins and related serpins from lancelets and sea urchins.** Group V3 serpins cluster into two major branches (clade E and clade I). *Bfl-spn-1* and *Spu-spn-1* from amphioxus and sea urchin, respectively, are grouped in the branch of clade I, supporting that these sequences are closely related to clade I serpins. The outgroup is *C. elegans* serpin 1 (Genbank, gi:2435565). This tree was created with MEGA4 based on the UPGMA method. Bootstrap values (in percentage) for 1000 replicates are shown (red color). A distance scale is shown below the tree.

In summary, most group V3 serpins are found from fishes to mammals. A serpin resembling the ancestor gene of group V3 serpins is unveiled in the lancelet and in sea urchins based on synteny and sequence related features. This suggests that the original locus of vertebrate group V3 serpins dates back at least to the time point of echinoderm separation about 550 My ago (**Figure 6**).

Furthermore, a serpin gene (JGI id - estExt\_fgenes1\_pg.C\_1860016/ NCBI id - XP\_001627732) is detected in *N. vectensis* genome as a probable neuroserpin ortholog based on sequence features (Kumar and Ragg, 2008).

### 5.14. Orthology analysis of group V4 serpins

Group V4 of vertebrate serpins has been defined by a gene structure depicting a conserved set of five introns at positions 67a, 123a, 192a<sup>1</sup>, 238c and 307a ( $\alpha_1$ -antitrypsin numbering) in the coding region (Ragg *et al.*, 2001). In mammals, group V4 serpins consists of three genes - pigment epithelium derived factor (PEDF/serpinF1),  $\alpha_2$ -antiplasmin ( $\alpha_2$ -AP/serpinF2) and C1 inhibitor (C1IN/serpinG1). These group V4 serpin genes are involved in very different physiological functions. PEDF is a non-inhibitory serpin that possesses neuroprotective and antiangiogenic functions (Steele *et al.*, 1993; Sawant *et al.*, 2004; Tombran-Tink, 2005).  $\alpha_2$ -antiplasmin is an inhibitor of plasmin and its fibrin bound form is a major regulator of blood clot lysis (Coughlin, 2005). C1 inhibitor is the primary inhibitor of two serine proteases (C1s and C1r) that, together with C1q, constitute the C1 complex of the classical pathway of complement (Cooper, 1985; Lener *et al.*, 1998). The protein alignments of group V4 serpins are shown in **appendices 8.3.27 to 8.3.29**.

#### 5.14.1. Gene structure of group V4 serpins

Since gene structure is a primary distinguishable parameter for classifying a new vertebrate serpin, the gene architectures of probable group V4 serpin homologs in different vertebrates were determined. **Table 30** shows that all vertebrate investigated contain at least two genes that depict the basic exon-intron structure of group V4 serpin genes. Introns at the canonical positions 67a, 123a, 192a, 238c, and 307a are conserved with some deviations. Due to lack of data for the *Tetraodon* PEDF gene (PEDF2\_TNI), introns cannot be assigned in this gene (indicated by?). Currently only the region spanning the C-terminal part with intact RCL is found (**appendix 8.3.23**). In *Fugu*, there are two A2AP like genes (A2AP2\_FRU and A2AP2\_FRU). The A2AP1\_FRU lacks the intron at position 123a. Lamprey (*Petromyzon marinus*) has two A2AP like group V4 members as A2APL1\_PMA and A2APL2\_PMA<sup>2</sup>. From the rather fragmented lamprey genome data, gene structure can be deduced for A2APL2\_PMA, but not for A2APL1\_PMA (indicated by ?). There is another type of group V4 serpins in fishes that possess similarity with the C1 inhibitor, which has two extra Ig domains in N-terminal part and this group V4 serpins are unique and are only found in fishes up to now, therefore these serpins are named as fish-specific group V4 (FSG4) serpins.

<sup>1</sup> also shared by group V2 and V6.

<sup>2</sup> A2AP like genes of lamprey are indexed with L1 and L2 to differentiate from other vertebrates A2AP genes since its orthology cannot be assigned yet and only tentative name is assigned based on sequence comparisons.



**Table 30: Intron positions of group V4 genes.** The presence (+) or absence (-) of intron positions is shown. In *Tetraodon* PEDF gene (PEDF2\_TNI), introns cannot be assigned (indicated by ?). There is loss of one intron at position 123a in *Fugu* A2AP gene 1, A2AP1\_FRU. Lamprey (*Petromyzon marinus*) has two A2AP like genes, A2APL1\_PMA and A2APL2\_PMA. From currently available lamprey genome data, only the gene structure of A2APL2\_PMA can be deduced.

Group V4 serpin gene	Intron at position				
	67a	123a	192a	238c	307a
PEDF_HSA (P36955)	+	+	+	+	+
PEDF_MMU (P97298)	+	+	+	+	+
PEDF_RNO (Q80ZA3)	+	+	+	+	+
PEDF_GGA (gi:50758202)	+	+	+	+	+
PEDF_XTR (ENSXETP00000050413)	+	+	+	+	+
PEDF2_FRU (FRUP00000141273)	+	+	+	+	+
PEDF2_TNI (GSTENP00013159001)	?	?	?	?	?
PEDF2_DRE (ENSDARP00000069366)	+	+	+	+	+
A2AP_HSA (P08697)	+	+	+	+	+
A2AP_MMU (Q61247)	+	+	+	+	+
A2AP_RNO (Q68FT8)	+	+	+	+	+
A2AP_GGA (XP_415807.2)	+	+	+	+	+
A2AP_XTR (ENSXETP00000029676)	+	+	+	+	+
A2AP1_FRU (FRUP00000162952)	+	-	+	+	+
A2AP2_FRU (e_gw2.417.16.1)	+	+	+	+	+
A2AP2_TNI (GSTENP00014689001)	+	+	+	+	+
A2AP2_DRE (ENSDARP00000078640)	+	+	+	+	+
A2APL1_PMA (GENSCAN00000047295)	?	?	?	?	?
A2APL2_PMA (GENSCAN00000097429)	+	+	+	+	+
C11N_HSA (P05155)	+	+	+	+	+
C11N_MMU (P97290)	+	+	+	+	+
C11N_RNO (NP_954524.1)	+	+	+	+	+
C11N_GGA (gi:50747972)	+	+	+	+	+
FSG4_FRU (FRUP00000133449)	+	+	+	+	+
FSG4_TNI (GSTENP00009345001)	+	+	+	+	+
FSG4_DRE (ENSDARP00000041512)	+	+	+	+	+

#### 5.14.2. Synteny analysis of PEDF and $\alpha_2$ -antiplasmin

In humans, the group V4 genes PEDF and  $\alpha_2$ -AP are found in a cluster on chromosome 17 and this arrangement has been retained in chicken and *Xenopus*. This PEDF-A2AP cluster is flanked by marker genes SCF-WDRD<sup>1</sup> on the one side and by a set of three genes (RPA1-RTN4R-DPH1<sup>1</sup>) on the other side (**Figure 60**). A similar genomic organization is found in *Fugu* (scaffold\_156), however, only the  $\alpha_2$ -AP gene 1 (A2AP1\_FRU) is present, whereas the PEDF gene is lacking. This syntenic arrangement suggests that A2AP1\_FRU is a fish orthologue of mammalian  $\alpha_2$ -AP.

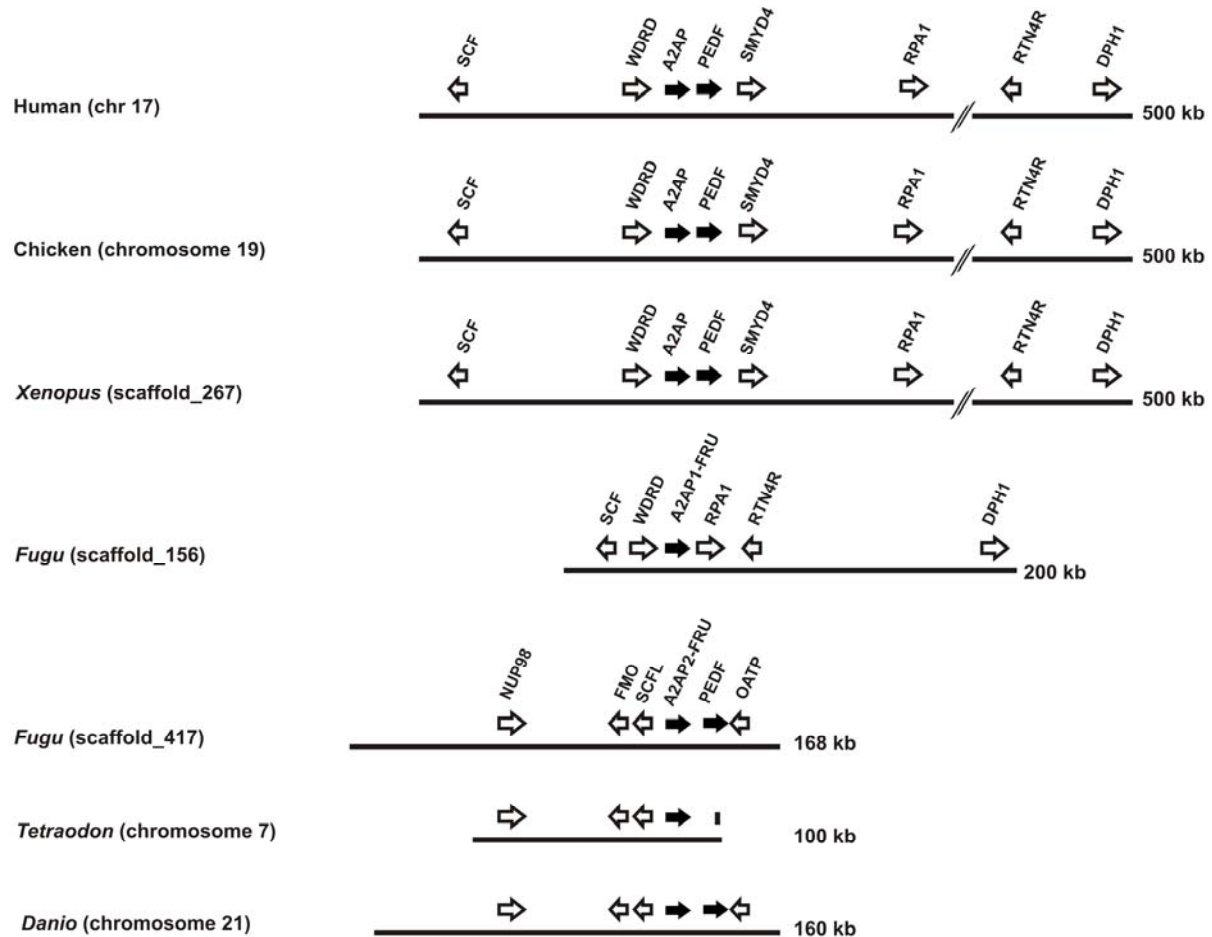
In contrast, no such syntenic arrangement is found in *Danio* and *Tetraodon*. This may be due to loss of the genomic fragment containing these genes. Alternatively, these genomic fragments have yet escaped detection.

Interestingly, another genomic locus with an A2AP-like gene (A2AP2\_FRU) and a PEDF-like gene is retained in *Fugu* (scaffold\_417) that is flanked by different sets of markers. This

<sup>1</sup> Appendix 8.4.6.



organization is shared by *Danio* and *Tetraodon*, suggesting that all fishes have acquired another genomic fragment carrying paralogs of representing A2AP and PEDF genes. Consequently, these genes have been named with index 2 (e.g. PEDF2\_FRU or A2AP2\_TNI).



**Figure 60: Synteny of the group V4 genes,  $\alpha_2$ -AP, and PEDF.** In most vertebrates,  $\alpha_2$ -AP and PEDF are clustered and flanked by a set of common marker genes. *Fugu* has two  $\alpha_2$ -AP like genes, one matching to the mammalian cluster and other one showing fish specific cluster. This suggests presence of a mammalian ortholog cluster and a paralog cluster in *Fugu*.

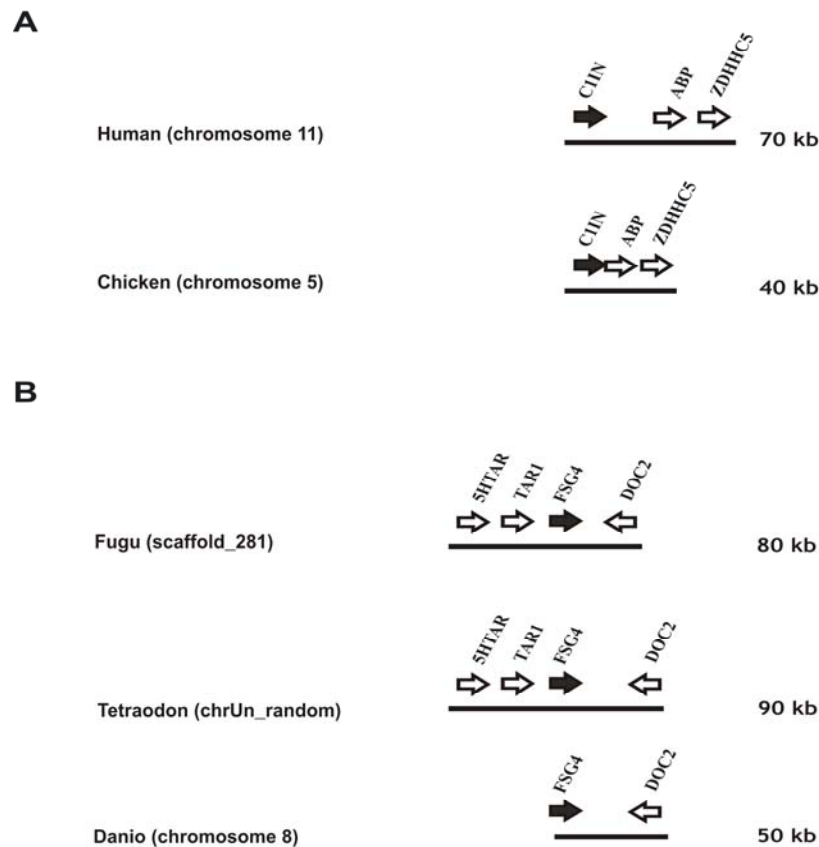
This suggests that out of fishes investigated, only *Fugu* has mammalian of ortholog of A2AP (A2AP1\_FRU) and none of these fishes has mammalian ortholog of PEDF gene.

### 5.14.3. Synteny analysis of the C1 inhibitor

Synteny analysis of the C1 inhibitor (C1IN) gene and a fish specific group V4 (FSG4) gene is shown in **Figure 61**. In higher vertebrates, the C1IN is flanked consistently by ABP-ZDHHC5<sup>1</sup> genes on one side (**Figure 61A**). A C1 inhibitor like gene is not detected in frog.

<sup>1</sup> Appendix 8.4.6.

This syntenic organization is not found in fishes, suggesting that fishes do not have human C1IN ortholog. Nevertheless, another gene FSG4 (with similarity to C1 inhibitor) is found in a distinct syntenic organization (**Figure 61B**), flanked by DOC2B<sup>1</sup> on one side and 5HTAR-TAR1 markers (not found in *Danio*) on the other side.



**Figure 61: Genomic localization of the C1 Inhibitor gene and a fish specific group V4 (FSG4) gene. (A)** Genomic localization of C1 inhibitor (C1IN) genes in human and chicken. **(B)** Genomic localization of FSG4 in different fishes. FSG4 gene is flanked by a distinct set of markers, which do not match with markers flanking C1 inhibitor genes in higher vertebrate. This difference in genomic localizations suggests that C1 inhibitor genes of higher vertebrates and FGS4 genes of fishes do not share orthology.

#### 5.14.4. Sequence comparisons of group V4 serpins

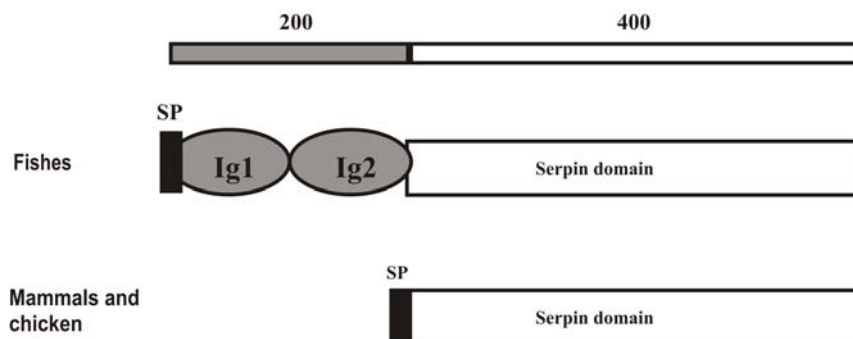
To further investigate orthology and paralogy of group V4 serpin genes in vertebrates, sequence comparisons of group V4 serpins were carried out.

Vertebrate PEDFs and its paralogs in fishes (PEDF2) show 34-85% sequence identity and 55-90% sequence similarity on the amino acid level with human PEDF (**appendix 8.3.27**). The RCL is probably non-inhibitory in all species due to bulky residues in the hinge region of PEDF protein (red boxes in **appendix 8.3.27**). The PEDF homologs are further characterized by the presence of a nuclear localization signal (NLS) (Tombran-Tink, 2005; Tombran-Tink *et al.*, 2005) (brown boxes in **appendix 8.3.27**).

$\alpha_2$ -AP from vertebrates and its orthologs and paralogs in fishes (A2AP2) show 25-74% sequence identity and 41-86% sequence similarity on the amino acid level to human  $\alpha_2$ -AP

(**appendix 8.3.28**). The inhibitory RCL of human  $\alpha_2$ -AP has two overlapping reactive sites within RCL i.e. R-M for inhibition of plasmin and trypsin, respectively, and M-S for inhibition of chymotrypsin (Potempa *et al.*, 1988). These two reactive sites are fully conserved (R-M-S) in mammalian  $\alpha_2$ -AP. However, there are variations in these sequences for  $\alpha_2$ -AP like genes from non-mammals.  $\alpha_2$ -AP-like genes are further characterized by N-terminal and C-terminal extensions.

The protein sequences of C1IN and of FSG4 share 20-68% sequence identity and 38-80% sequence similarity with human C1IN (**appendix 8.3.29**). The RCL is inhibitory, displaying residues R-[TSNI] and R-[TS] at positions P1 and P1' of C1IN and FSG4, respectively (red boxes in appendix 9.3.16). FSG4 from *Fugu*, *Tetraodon*<sup>1</sup>, and *Danio* carry two immunoglobulin (Ig) like domains (200 amino acids long) in the N-terminal region as predicted by the SMART program<sup>2</sup> (Schultz *et al.*, 1998; Letunic *et al.*, 2006) (**Figure 62**).



**Figure 62: Domain architecture comparisons of C1 inhibitor and a fish specific group V4 (FSG4) serpin.** FSG4 of fishes has two immunoglobulin like domains (Ig1 and Ig2) (after signal peptide [SP]) and a serpin domain, whereas another type N-terminal extension plus a serpin domain is found in C1 inhibitor genes of higher vertebrates. This difference in domain organizations indicates that C1 inhibitor and FSG4 of fishes do not share orthology. SP - signal peptide.

This unique group V4 serpin - FSG4 - is also found in trout<sup>3</sup> (Wang and Secombes, 2003) and the Japanese flounder<sup>4</sup> (Inoue *et al.*, 1997).

To explore orthology and relationships of different group V4 serpins further, a phylogenetic tree (**Figure 63**) was constructed based on the Maximum Parsimony method with help of MEGA4. The inhibitory and non-inhibitory members separate into distinct branches in this phylogenetic tree. The lamprey group V4 members are tentatively assumed to be A2AP-like genes because of the presence of terminal extensions and similarities at RCL regions. These genes are clustering with the inhibitory branch of group V4 serpins. The A2AP and C1 inhibitor sequences also separate into distinct sub-branches and the fish specific group V4 serpin - FSG4 cluster into a separate sub-branch together with orthologs of C1IN of higher vertebrates. The C1IN gene of higher vertebrates and the FSG4 gene of fishes differ as assessed by several criteria, suggesting that orthologs of human C1IN have been not detected

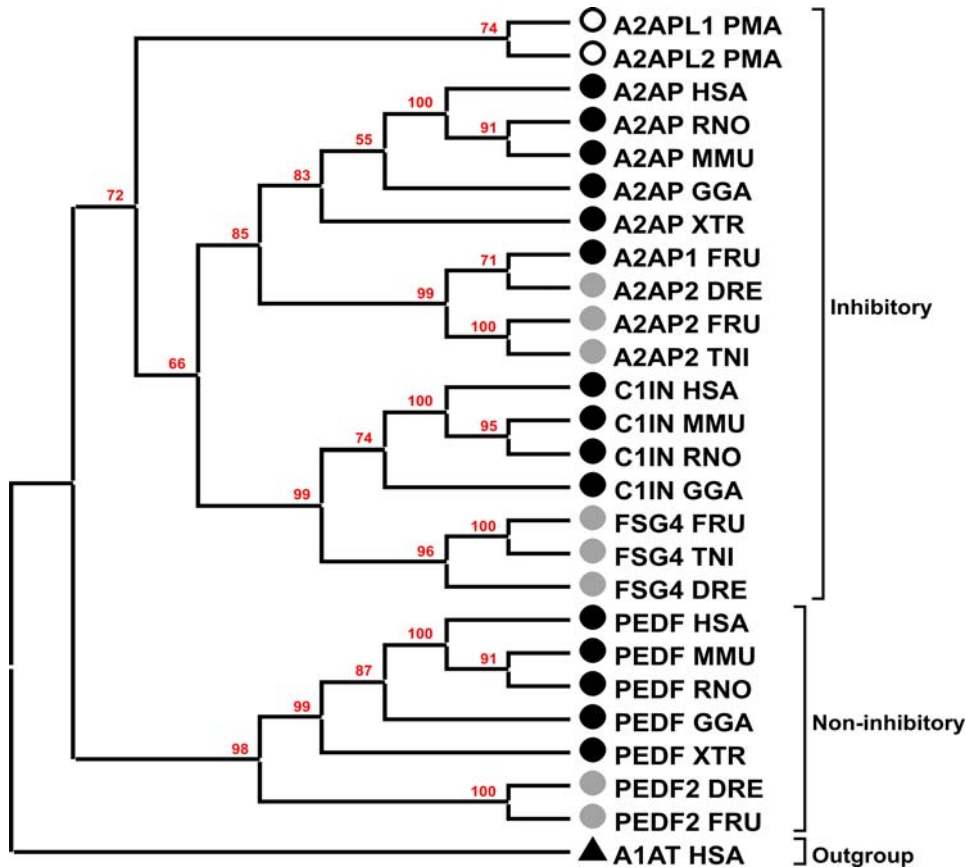
<sup>1</sup> *Tetraodon* [cDNA, GenBank CR656519]

<sup>2</sup> SMART website, <http://smart.embl-heidelberg.de/>

<sup>3</sup> *Oncorhynchus mykiss* [cDNA, GenBank AJ519930]

<sup>4</sup> *Paralichthys olivaceus* [cDNA, GenBank BN000290 and EST, GenBank C23239, C23240]

in fishes up to now. Possibly, the C1IN gene was lost in these fishes. Instead, another gene FSG4, with Ig like extra domains may have been acquired, suggesting that FSG4 of fishes possibly functions differently (neofunctionalization) as compared to C1IN genes in higher vertebrates.



**Figure 63: Phylogenetic tree of group V4 serpins based the Maximum Parsimony method.** There are two major branches separating non-inhibitory (PEDF) and inhibitory group V4 serpins. Inhibitory group V4 members constitute three distinct sub-branches separating  $\alpha_2$ -AP, C1IN and FSG4, and  $\alpha_2$ -AP like genes from lamprey A2APL1\_PMA and A2APL2\_PMA. The outgroup is human  $\alpha_1$ -AT (black triangle). Bootstrap values (in percent) for 1000 replicates are shown (red color). Orthologs and paralogs<sup>1</sup> of human group V4 serpins are depicted by black and grey circles, respectively. Orthology of lamprey group V4 members is still open (white circles).

Interestingly, the majority of group V4 serpins in fishes do not have human orthologs (grey circles in **Figure 63**).

In summary, orthologs of most of human group V4 serpins are lost in fishes or cannot be found in current genomic sequence versions, with exception of A2AP in *Fugu* (A2AP1\_FRU). If undetected this may indicate that during evolution, fishes lost the orthologs of the higher vertebrate group V4 serpin loci. Instead, they have paralogs due to genome duplication and diversification. The syntenic divergence of group V4 serpins in fishes thus

<sup>1</sup> Paralogs of human group V4 serpins are indexed by 2 after their names.

provides a rudimentary insight into whole genome duplication event in fishes and subsequent gene diversification events.

### 5.15. Orthology analysis of V5 serpin genes

Group V5 consists of a single member - antithrombin III (ATIII). Its gene encompasses seven exons and six introns with conserved intron positions based on gene structure analysis of several mammalian serpins (Ragg *et al.*, 2001). In the human genome, the ATIII gene is located on chromosome 1q23–q25. ATIII is the major thrombin inhibitor in the blood coagulation cascade (Jordan, 1983), requires heparin for activation and has potent anti-angiogenic activity in certain conformations (Gettins *et al.*, 1996). The alignment of ATIII homologs of different vertebrates is available in **appendix 8.3.30**.

#### 5.15.1. Gene structure of ATIII genes

Since gene structure plays an important role in distinguishing group V1-V6, the exon-intron structures of ATIII orthologs were determined. It was found that the gene structure was conserved in ATIII of different vertebrates with group V5 specific introns maintained at positions 78c, 148c, 191c, 320a, and 339c with some variations as shown in **Table 31**.

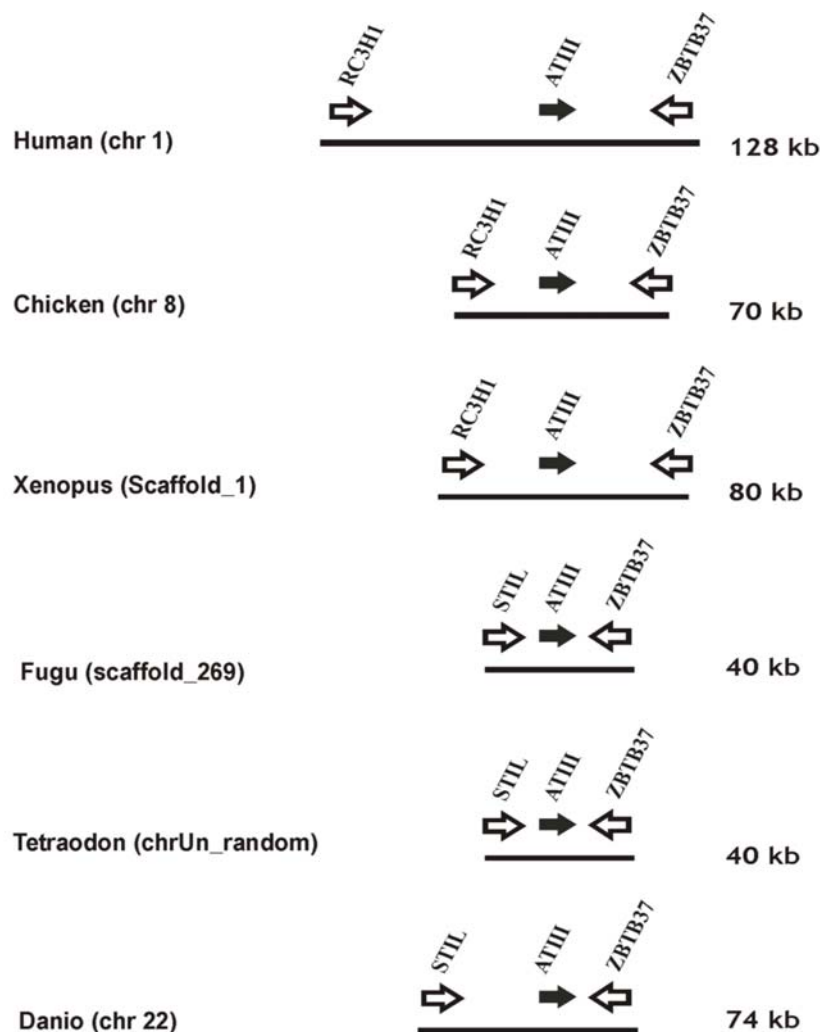
**Table 31: Intron positions of the ATIII gene in different vertebrates.** The presence (+) or absence (-) of intron positions is shown. The novel intron at position 262c of ATIII genes from fishes is also found in group V1 serpins. There were gaps in the genomic sequence of chicken. Hence, introns at positions 320a and 339c for ATIII\_GGA could not be detected (indicated by ?).

	Intron at Position					
	78c	148c	191c	262c	320a	339c
ATIII_HSA (P01008)	+	+	+	-	+	+
ATIII_MMU (P32261)	+	+	+	-	+	+
ATIII_RNO (Q5M7T5)	+	+	+	-	+	+
ATIII_GGA (XP_422282)	+	+	+	-	?	?
ATIII_XTR (estExt_fgensch1_pm.C_10068)	+	+	+	-	+	+
ATIII_DRE (ENSDARG00000042684)	+	+	+	+	+	+
ATIII_TNI (GSTENP00004792001)	+	+	+	+	+	+
ATIII_FRU (e_gw2.269.120.1)	+	+	+	+	+	+

In fishes, a novel intron at position 262c was found. This intron position is normally characteristic for group V1 serpins. Due to gaps in genomic region containing the chicken ATIII, the introns at positions 320a and 339c could not be identified.

### 5.15.2. Synteny analysis of ATIII genes

In order to investigate orthology of ATIII genes further, an analysis of the ATIII locus in different vertebrates was carried out (**Figure 64**). The ATIII gene in the human genome is surrounded by the marker genes RC3H1<sup>1</sup> (same orientation) on one side and the ZBTB37 gene (opposite orientation of ATIII gene) on the other side. Similar synteny arrangements were found in chicken and in *Xenopus*. In fishes, the ATIII-ZBTB37 synteny is conserved, but on the other side, another marker gene – STIL<sup>1</sup> is adjacent to the ATIII gene.



**Figure 64: Synteny comparison of ATIII genes in different vertebrate genomes.** The ATIII gene (black arrow) is flanked from fish to mammals by marker gene ZBTB37 (**appendix 8.4.6**) on one side. On the other side, either the RC3H1 marker (**appendix 8.4.6**) (from mammals to *Xenopus*) or the STIL gene (**appendix 8.4.6**) is found (in fishes).

These data document that the ATIII gene synteny is conserved in different vertebrates.

<sup>1</sup> See appendix 8.4.7.

### 5.15.3. Sequence comparisons of ATIII genes

To investigate the orthology using sequence comparisons, the protein sequences of ATIII genes from vertebrates were analyzed. The ATIII gene is highly conserved in vertebrates and the sequences show 50-87% sequence identity and 67-97% sequence similarity on amino acid level with human ATIII from fish to mammals. From this alignment, several signature sequences have been deduced: The helix D region of the ATIII, which is involved in heparin binding (Gandrille *et al.*, 1990) was found to be highly conserved (yellow boxes in appendix 11.3.17). No other vertebrate serpin has these specific arrangements of basic residues in the helixD region (**appendix 8.3.30**). There are eight basic residues reported to be important in heparin binding – four in N-terminal part of the ATIII molecule (positions K11, K13, R46, and R47), and four in the helix D region (molecule positions K126, R129, R132, and R133) in mature human ATIII<sup>1</sup> (Gandrille *et al.*, 1990; Backovic and Gettins, 2002). The majority of these residues are conserved in vertebrates with the exception of R46, which is only found in mammalian ATIII genes (orange boxes in **appendix 8.3.30**). The inhibitory RCL region (red boxes in **appendix 8.3.30**) is also highly conserved with P1-P1' position (R-S) maintained in all vertebrates as shown in **Figure 65**.



**Figure 65: Sequence logo of RCL region of ATIII from different vertebrates.** Most positions of the RCL region are highly conserved and the cleavage site between P1-P1' is marked with an arrow. This logo was created using weblogo<sup>2</sup> (Schneider and Stephens, 1990; Crooks *et al.*, 2004).

The hinge region residues P14-P15 (G-S) are highly conserved in all vertebrate ATIII protein, whereas in majority of other serpins these residues are G-T. It has been reported that three pairs of disulfide bridges are required in human ATIII in order to bind heparin with high affinity and to inhibit proteinases (Longas *et al.*, 1980; Ferguson and Finlay, 1983). The ATIII from all vertebrates investigated has maintained the six cysteines constituting these three pairs of disulfide bridges (marked C1, C2, and C3 pairs in **appendix 8.3.30**). From the serpin specific conserved 51 amino acid positions (summarized in **appendix 8.1**), 35 are found to be fully conserved in vertebrate ATIII genes (black boxes in **appendix 8.3.30**). There are four N-glycosylation site<sup>3</sup> in mature human ATIII (positions N96, N135, N155, and N192; cyan boxes in **appendix 8.3.30**) (Backovic and Gettins, 2002). These sites are found to be conserved in ATIII with some exceptions, like N-glycosylation sites at N96 and N135, which are not found in ATIII of chicken and in ATIII of fishes, respectively. The N-glycosylation site at N155 is not found in ATIII of *Fugu* and *Tetraodon*. These fishes have acquired a

<sup>1</sup> Numbering is based on mature human ATIII.

<sup>2</sup> Weblogo website, <http://weblogo.berkeley.edu/>

<sup>3</sup> N-glycosylation site, NX[ST], where X = any amino acid except P.

different N-glycosylation site at N160, which is also present in chicken ATIII. In summary, gene structure and synteny conservation, presence of conserved helix D, RCL, three pairs of disulfide bridges and conserved basic residues, together with sequence identity and sequence similarity, suggests the presence of ATIII orthologues from fish to human.

### 5.16. Orthology analysis of group V6 serpins

Group V6 of vertebrate serpins has been defined by a gene structure depicting three introns at positions 192a, 225a and 300c in their coding regions (Ragg *et al.*, 2001). These genes code for heat shock protein 47 kDa (HSP47), which possesses a C-terminal endoplasmic reticulum (ER) retention signal<sup>1</sup> (Pelham, 1990). HSP47 is a non-inhibitory serpin that is found in the ER of collagen producing cells where it is involved in the correct folding of procollagen triplet helices. Furthermore, it assists in transport of procollagen from the ER to the Golgi complex (Nagata, 1996; Lamande and Bateman, 1999; Hendershot and Bulleid, 2000; Sauk *et al.*, 2005). The alignment of HSP47 homologs of different vertebrates is available in the **appendix 8.3.31**.

#### 5.16.1. Gene structure of group V6 serpins

Since gene structure is a primary discriminatory factor for classification as a prospective member of groups V1-V6, the exon-intron structures of suspected HSP47 homologs in different vertebrates were determined. **Table 32** shows that all vertebrates investigated contain at least one gene that basically depicts the exon-intron structure of group V6 serpins. Introns at positions 192a, 225a, and 300c are conserved with some deviations. In contrast to humans and other vertebrates, which contain a single group V6 gene (HSP47), there are two or three group V6 homologs in *Fugu* and *Danio*, respectively. The *Fugu* HSP47 gene 1 (HSP47\_1\_FRU) has two additional unique introns at positions 36b and 102c. In the *Tetraodon* HSP47 gene (HSP47\_TNI) the intron at position 192a was not identified, probably due to sequencing errors in the coding region of this gene.

**Table 32: Intron positions of group V6 genes in different vertebrates.** The presence (+) or absence (-) of intron positions is shown. Unique introns are found in *Fugu* HSP47 gene 1 (HSP47\_1\_FRU) at positions 36b and 102c. In the *Tetraodon* HSP47 gene (HSP47\_TNI), the presence or absence of the intron at position 192a cannot be confirmed since there are sequencing errors in the coding region of this gene (indicated by ?).

Group V6 serpin gene	Intron at position				
	36b	102c	192a	225a	300c
HSP47_HSA (P29043)	-	-	+	+	+
HSP47_MMU (P97290)	-	-	+	+	+
HSP47_RNO (NP_954524)	-	-	+	+	+
HSP47_GGA (gi:45384240)	-	-	+	+	+
HSP47_XTR (estExt_fggenesh1_pg.C_2770030)	-	-	+	+	+
HSP47_1_FRU (e_gw2.131.10.1)	+	+	+	+	+
HSP47_2_FRU (fgh5_pg.C_scaffold_186000009)	-	-	+	+	+
HSP47_TNI (GSTENP00006756001)	-	-	?	+	+

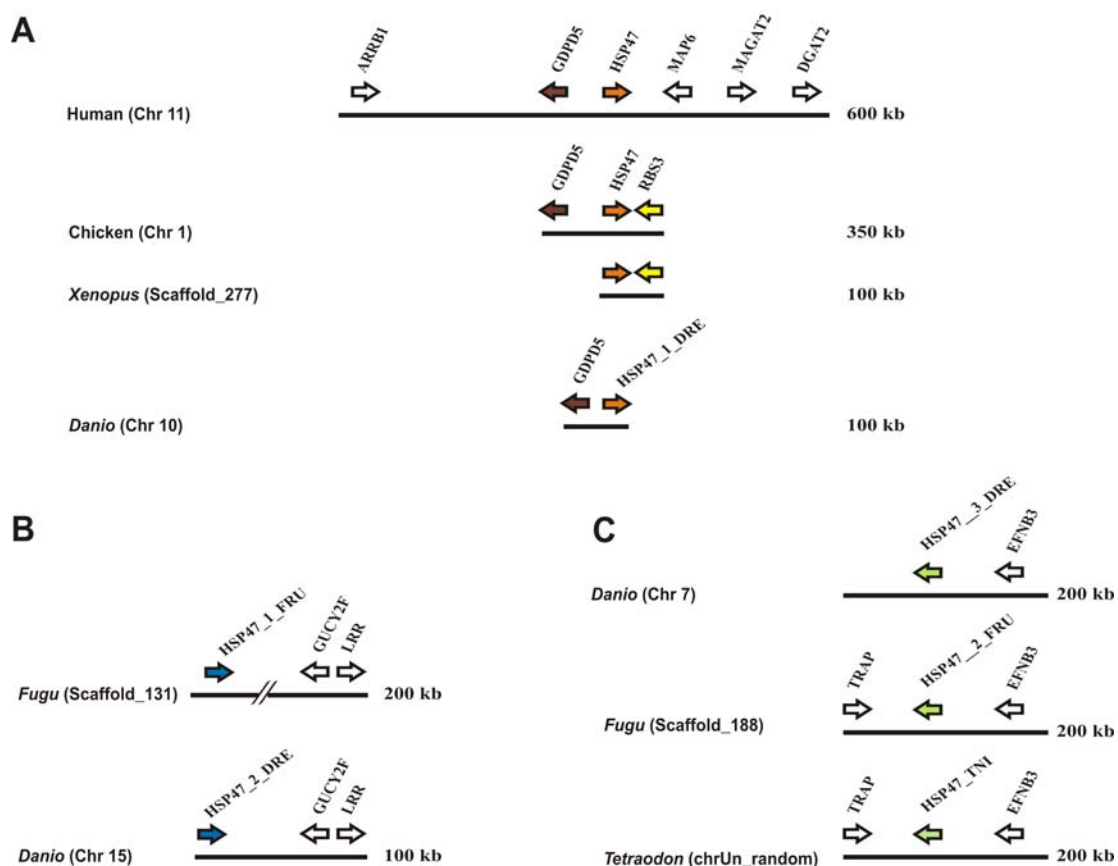
<sup>1</sup> ER-retention signal, [RKH]DEL (Pelham, 1990).



HSP47_1_DRE (ENSDARP00000037780)	-	-	+	+	+
HSP47_2_DRE (ENSDARP00000028177)	-	-	+	+	+
HSP47_3_DRE (ENSDARP00000052941)	-	-	+	+	+
HSP47_PMA (GENSCAN00000147606)	-	-	+	+	+

### 5.16.2. Synteny analysis of group V6 serpins

Since there were three or two group V6 genes in *Danio* and *Fugu*, respectively, the genome micro-synteny was analyzed to resolve orthology with mammalian HSP47 (**Figure 66**). The HSP47 gene in the human genome is found on chromosome 11, flanked by marker gene GDPD5<sup>1</sup>(brown) on one side, and the MAP6-MAGAT2-DGAT2 gene cluster on the other side (**Figure 66A**).



**Figure 66: Genomic localization of HSP47 homologs in different vertebrate genomes. (A)** Syntenic arrangement of human HSP47 orthologs (orange). **(B)** Syntenic arrangement of group V6 homologs HSP47\_1\_FRU and HSP47\_2\_DRE (blue). **(C)** Syntenic arrangement of HSP47\_2\_FRU, HSP47\_3\_DRE, and HSP47\_TNI (light green). Obviously, there are three sets of HSP47 homologs in fishes.

In chicken and frog, the HSP47 gene is surrounded by GDPD5 (brown) and RBS3 (yellow) markers. HSP47\_1\_DRE<sup>2</sup> and marker gene GDPD5 were found to be syntenic on

<sup>1</sup> Appendix 8.4.8.

<sup>2</sup> cDNA available from the Zebrafish Model Organism Database ([www.zfin.org](http://www.zfin.org)), ZFIN ID: ZDB-GENE-990415-93.

chromosome 10 in *Danio*, unveiling this gene as a true orthologue of mammalian HSP47. This mammalian HSP47 ortholog is also found in two other fishes – Medaka<sup>1</sup> and Sicklefis<sup>2</sup> with EST evidence (data not shown), which suggests that fishes generally possess a true ortholog of mammalian HSP47. Consequently, HSP47\_2\_DRE and HSP47\_3\_DRE<sup>3</sup> are paralogs of mammalian HSP47. Based on synteny analysis, orthologs of mammalian HSP47 in *Fugu* and *Tetraodon* were not identified. Probably the true HSP47 orthologs have been overlooked in these organisms.

In *Danio* and *Fugu*, HSP47\_2\_DRE and HSP47\_1\_FRU respectively, are flanked by markers GUCY2F and LRR<sup>4</sup> (**Figure 66B**). Similarly, HSP47\_3\_DRE of *Danio*, HSP47\_2\_FRU of *Fugu*, and HSP47\_TNI of *Tetraodon* are flanked by marker genes EFNB3 and TRAP<sup>2</sup> (not present in *Danio*), revealing these genes as orthologs (**Figure 66C**). This synteny is also found in Medaka<sup>5</sup> (data not shown), advocating the presence of this group V6 gene in different fishes.

The micro environment of the single lamprey HSP47 gene cannot be depicted using the current version of genomic sequences (version PMA3).

### 5.16.3. Sequence comparisons of group V6 serpins

To further unravel the relationships of group V6 genes, sequence based comparisons were carried out. HSP47-like genes are conserved from lamprey to mammals and these genes show 22-96% sequence identity and 37-98% sequence similarity with human HSP47, respectively. The HSP47\_TNI protein is highly diverged from standard HSP47 protein as well as from all other serpin sequences (**Table 33**).

**Table 33: Sequence comparisons of HSP47 homologs in vertebrates.** Percentage sequence identity (SI) and percentage sequence similarity (SS) values are shown as compared to HSP47\_HSA and A1AT\_HSA. Synteny based clustering divides group V6 genes into three sets: set I – true mammalian HSP47 orthologs (orange), set II - fish specific paralogs as compared to **Figure 66B** (blue) and set III as in **Figure 66C** (light green). Orthology of lamprey 6 group gene, HSP47\_PMA (grey) cannot be decided on this basis.

Human Serpins	Values (%)	HSP47_MMU	HSP47_RNO	HSP47_GGA	HSP47_XTR	HSP47_1_FRU	HSP47_2_FRU	HSP47_TNI	HSP47_1_DRE	HSP47_2_DRE	HSP47_3_DRE	HSP47_PMA
		HSP47_HSA	SI	96	96	76	70	63	29	22	65	64
	SS	98	98	88	83	82	46	37	83	82	52	65
A1AT_HSA	SI	23	23	25	24	24	18	14	25	24	17	23
	SS	45	45	45	46	44	35	26	45	46	37	41

All HSP47 homologs appear to be non-inhibitory (red boxes in **appendix 8.3.31**), since they contain bulky amino acids in the hinge region. All vertebrate group V6 members have an ER

<sup>1</sup> Mammalian HSP47 ortholog in Medaka, Ensembl accession Id - ENSORLG00000014312.

<sup>2</sup> Mammalian HSP47 ortholog in sicklefish, Ensembl accession Id - ENSGACG00000006375.

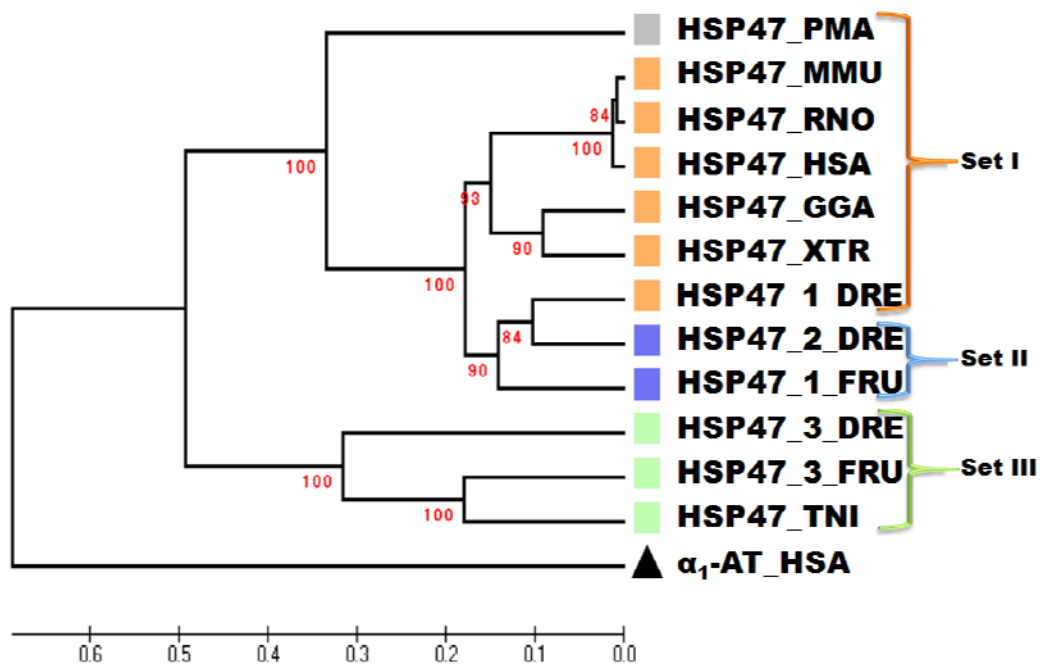
<sup>3</sup> cDNA available from the Zebrafish Model Organism Database (www.zfin.org), ZFIN ID: ZDB-GENE-050417-12.

<sup>4</sup> Appendix 8.4.8.

<sup>5</sup> The group V6 gene in Medaka, Ensembl accession Id - ENSORLG00000003689.

retention signal ([RH]DEL) at the C-terminus (**appendix 8.3.31**). Out of 51 amino acid positions conserved in the majority of serpins (**appendix 8.1**), 31 residues are fully conserved (black boxes in **appendix 8.3.31**).

To understand orthology of the group V6 genes further, a phylogenetic tree was constructed (**Figure 67**) using the UPGMA method (Sneath and Sokal, 1973). The mammalian orthologs of human HSP47 gene (set I) cluster in one branch of the phylogenetic tree (orange squares). In fishes, a recent branching has created set II genes (blue squares). Set III genes – the second cluster of paralogues are divided into a distinct branch (light green squares) which comprises genes with lower sequence identities (**Table 33**).



**Figure 67: Evolutionary tree of HSP47 homologs from lamprey to human created with the UPGMA method, using MEGA4.** Three distinct sets of HSP47 homologs are colored according to syntenic arrangements (**Figure 66**). The lamprey HSP47 (grey square), whose syntenic arrangement is not known, clusters with mammalian HSP47 sequences. The outgroup is human  $\alpha_1$ -AT (black triangle). Bootstrap values (1000 replicates) are shown in percentage (red color) and a distance scale is shown below the tree.

Phylogenetic analysis suggests that the single group V6 gene from lamprey (HSP47\_PMA) is an ortholog of mammalian HSP47. However, this needs to be confirmed by synteny analysis, which cannot be carried out with current version of genomic sequences (version PMAL3). The *Tetraodon*, HSP47\_TNI gene is problematic, because of sequencing errors in the coding region. Since a low complexity region was found at the intron at the position 300c, this region was deleted from the HSP47\_TNI protein sequence in **appendix 8.3.27**. The issue whether this gene has been pseudogenized or carries exonized intron sequences generating novel polypeptide domains (Schmidt and Davies, 2007) remains an open question.

In summary, orthologues of human HSP47 gene have been found from mouse to *Danio* as well as in some other fishes. In *Fugu* and *Tetraodon*, this gene might have been lost and only paralogues of HSP47 genes have been retained in these fishes. This suggests that the

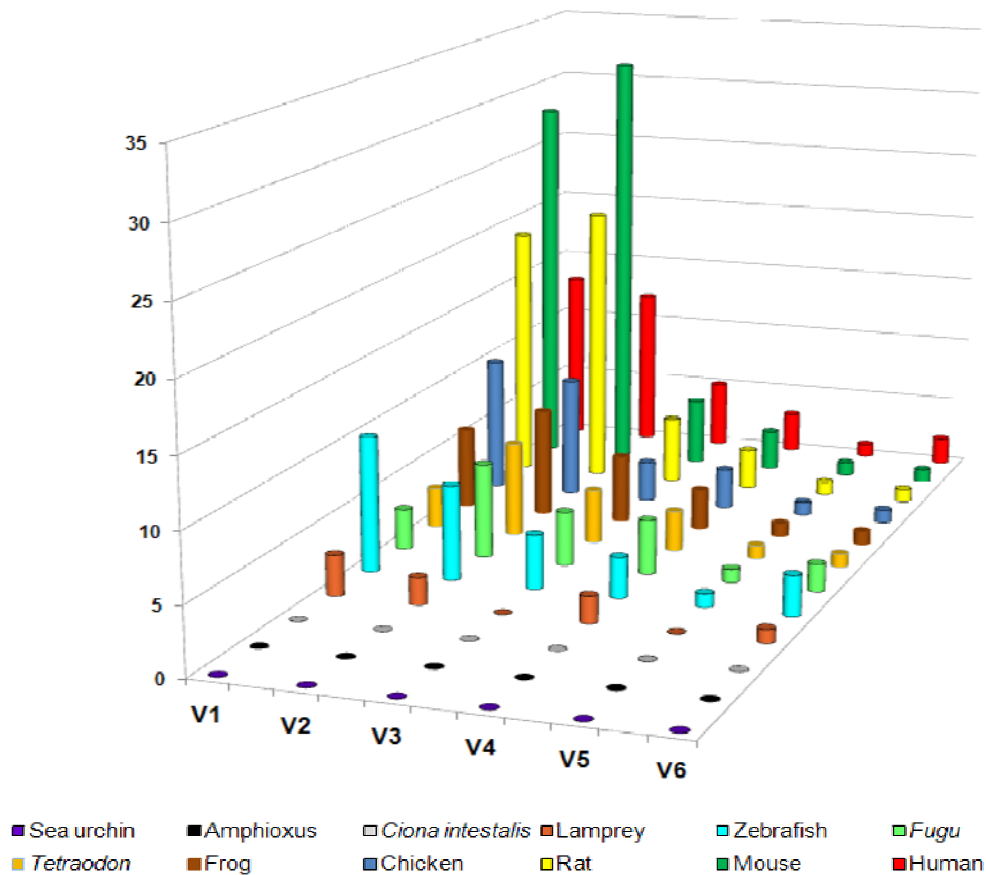
---

mammalian HSP47 orthologue has been lost in the pufferfish family (*Tetraodontidae*) and only paralogues have been retained, though it cannot be excluded that the true HSP47 orthologue in these fishes remained undetected up to now. Set III of HSP47 genes of fishes might represent a class of ancestor group V6 genes as is evident from the branching in the phylogenetic tree (**Figure 67**).

## 6. Discussion

### 6.1 Overview of vertebrate serpins from fishes to mammals

Figure 68 shows the group specific distribution of serpins in different metazoans.



**Figure 68:** Distribution of vertebrate serpins based on their intron-coded classification into six groups (V1-V6). The number of serpin genes from lamprey will probably increase, since the genome project is in its initial draft stage.

The number of group V1 and group V2 serpins varies considerably in different organisms. An expansion is evident from fish to mammals. Strikingly, mice and rats have more members of group V1 and V2 than human. In contrast, group V3 serpins are retained from fish to mammals without marked expansions. However, some exceptions are found such as PAI1 and serpinE3, which are missing in the chicken genome. Pancpin was not found in any fish genome analyzed. Based on sequence analysis alone, group V4 members appears to be conserved in vertebrates, but on analyzing synteny, we found that there are difference in the distribution of orthologs/paralogs (**Figure 60**). The only member of group V5 – the ATIII gene - is conserved across all vertebrates. The HSP47 gene of group V6 is conserved in most vertebrates, but there is a varying numbers of group V6 paralogs in fishes. Based simply on analysis of gene architectures, no serpin genes were found in sea squid, amphioxus, and sea urchin that share the gene organization of their vertebrate counterparts. However, by

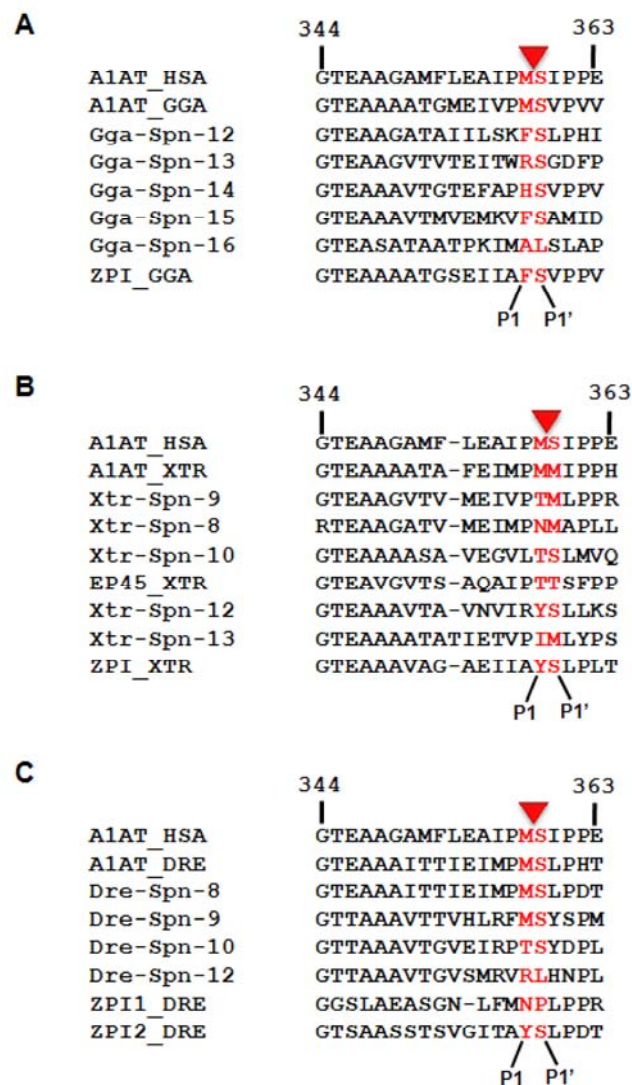
combining syntenic information and sequence-specific features, it was possible to trace orthologs of neuroserpin in amphioxus and sea urchin (Kumar and Ragg, 2008; section 6.4).

## 6.2 Evolutionary history of group V1 serpins

Group V1 serpin genes are found from lampreys to human. The human genome has two clusters of group V1 serpins that are located on chromosomes 6 and 18, respectively. In contrast, the chicken has only one such cluster and therefore it is argued that there was a split after mammal/bird divergence at around 310 Mya (Benarafa and Remold-O'Donnell, 2005; Kaiserman and Bird, 2005; Izuhara *et al.*, 2008). A chicken-type genomic organization of group V1 serpins is also found in frogs and in fishes (**Figure 40**). Fishes, in addition, possess some paralogous clusters of serpin genes (**Figures 41–43**). In frog, an additional cluster containing two serpins (with EST evidences) adjacent to the conserved orthologous cluster is found. The serpins SPB1/SPB6 of group V1 are probably conserved descendants of the ancestor of all group V1 serpins, since these genes are found in lampreys and other fishes and are also conserved across other vertebrate taxa. The group V1 serpins may be classified into sub-groups V1a and V1b, since these differ by one intron. Some scholars have argued that a serpin gene of group V1b (7 exons) is the ancestor of group V1a (8 exons) that has emerged in birds after divergence of frogs (Benarafa and Remold-O'Donnell, 2005; Kaiserman and Bird, 2005; Izuhara *et al.*, 2008). Their first argument coincides with our data, suggesting that group V1a serpins are derived from 7-exon genes such as MNEI/SPB6. However, the argument that 8-exon genes first arose in chickens does not hold, since Xtr-Spn-5 in *X. tropicalis* and pSPB6 in *T. nigroviridis*, are group V1a members having the 8 exons / 7 introns architecture. However, due to sequence alignment problems, the position of the extra intron at position 85c in *T. nigroviridis* is ambiguous. Therefore, we propose that group V1b is ancestral to all group V1 serpins and group V1a is suggested to have arisen independently several times in different vertebrates from fishes to mammals. The ancestor of group V1 serpins appears to have been generated during the emergence of vertebrates, and the oldest group V1 serpins are SPB1/SPB6 orthologs that are present in lamprey. An ancestor of serpinB6 was claimed to be present in urochordates (Kaiserman and Bird, 2005), however, using synteny, gene structures, and sequence motif for analysis, I did not find any evidence suggesting close relationships between any of *Ciona* serpins and group V1 serpins. BLAST searches using human serpinB1 or serpinB6 sequence for querying organisms such as insects (*Drosophila* and *Anopheles*), worms (*C. elegans*), sea urchin, and amphioxus also provided no clear evidence for direct ancestor/offsprings relationships of group V1 serpins. In contrast, using rare indels and synteny analyses, we have identified an ortholog of neuroserpin in deep-branching metazoans. This clearly shows that inclusion of synteny and indel analysis may facilitate kinship recognition. The complete genomic sequences of lamprey and hagfish will shed further light on this issue. In conclusion, an expansion of group V1 serpins was found from fish to mammals that, as previously reported, is particularly evident within mammalian genomes (Kaiserman *et al.*, 2002). To understand this expansion in detail, further comparative genomic studies including basal mammals such as marsupials and Platypus are essential.

### 6.3 Phylogenetic history of group V2 serpins

From fish to human, group V2 comprises multiple paralogs of  $\alpha_1$ -antitrypsin like genes. Genuine orthologs of angiotensinogen and HCII were identified from fish to human, using synteny and signature sequences. Concerning the other genes of group V2, one-to-one orthology allocation proved to be difficult, since in most genomes the clusters containing group V2 genes are derived from recent duplications resulting in proteins with high sequence similarities, often even within the usually hypervariable RCL region (**Figure 69**). The orthologs of the ZPI gene were identified by considering the syntenic conservation of marker genes (**Figure 45**). In fishes, the common microenvironment of the Spn\_94a gene (named due to a novel intron at position 94a) corroborates its fish specific ortholog and a paralog of human ZPI gene.



**Figure 69: Comparison of reactive center loops (RCL) of selected group V2 serpins in (A) chicken, (B) *X. tropicalis* and (C) zebrafish.** The proposed cleavage site and P1-P1' residues are marked (in red color). The RCL region of human  $\alpha_1$ -antitrypsin (A1AT\_HSA) is included.

Both *Fugu* and *T. nigroviridis* possess one more group V2 gene with an additional intron at position 215c (Spn\_215c), suggesting that they are orthologs. The origin of these genes, however, is unclear. No orthologs of the hormone binding serpins (THBG/CBG) were detected in non-mammalian vertebrates. In short, the conserved set of group V2 comprises only orthologs of angiotensinogen and HCII. In contrast, some fish-specific group V2 genes and the  $\alpha_1$ -antitrypsin-like genes are differentially expanded in vertebrates, particularly in mammalian lineages, such as rodents (Forsyth *et al.*, 2003) and cattle (Pelissier *et al.*, 2008). The expansion of group V2 members should be explored further by analyzing marsupials and Platypus, which branched out early in mammalian evolution. The presence of group V2 members in the lamprey genome suggests that this group originated during emergence of vertebrates. Further investigation of group V2 members in the hagfish genome and the complete lamprey genome will shed more light on this issue.

#### 6.4 Evolution of group V3 serpins

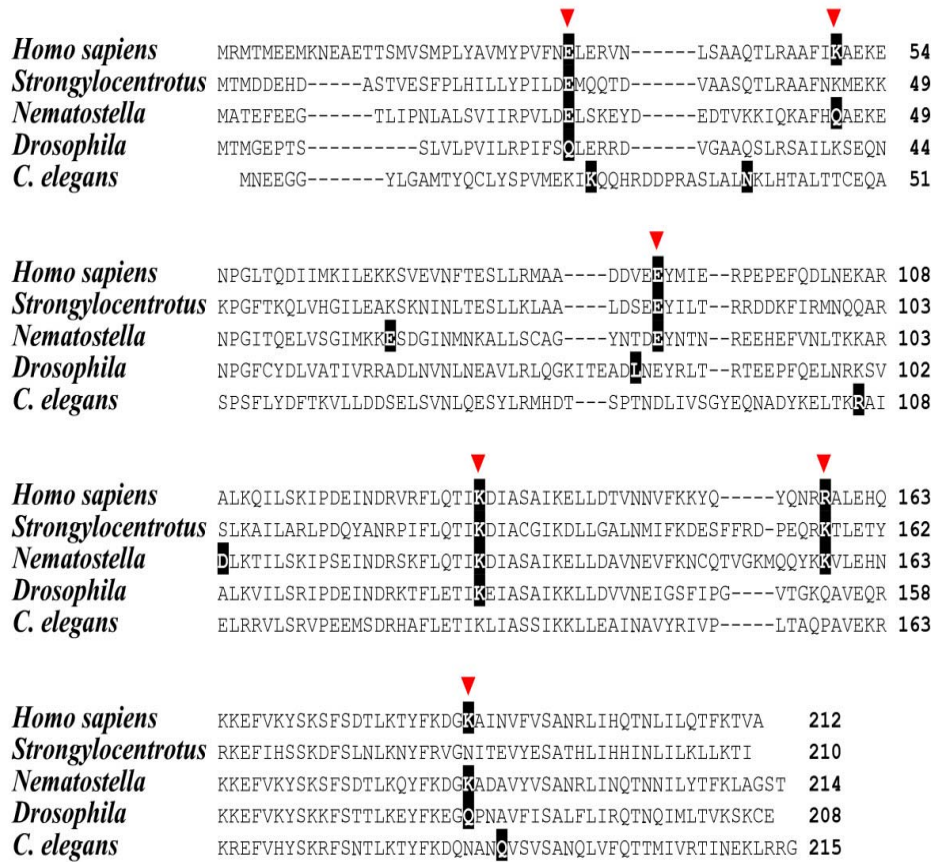
Group V3 encompasses five highly conserved inhibitory members: SerpinE1/plasminogen activator inhibitor 1 (PAI1), SerpinE2/glia derived nexin (GDN), SerpinE3, SerpinI1/neuroserpin, and SerpinI2/pancpin. While studying group V3 serpins, comprehensive insight into the phylogenetic history of neuroserpin was unraveled by combining discriminatory data from the genomic, gene and protein level. With aid of these data, the previously unknown origin of neuroserpins during metazoan evolution was settled. Synteny analysis proved to be very instrumental in this respect, demonstrating that rare genomic characters can provide very useful information for decoding of bonds in protein families with intricate evolutionary history. The strongly conserved syntenic association of PDCD10 and neuroserpin orthologs during diversification of deuterostomes is unraveled here. The conserved close linkage of expression of these two head-to-head oriented genes may have been caused by a bi-directional and asymmetrical promoter region inserted within the ~0.9 kb intergenic region separating the coding regions (Chen *et al.*, 2007). Dependence from a common regulatory unit may have forced the maintenance of this linkage. The rapidly increasing flood of data from genome sequencing projects (with rapid change in genome sequencing technologies) will certainly continue to provide further discriminatory markers, such as codon usage dichotomy (Krem and Di Cera, 2003), to enable robust classification of other metazoan serpins.

A C-terminal, KDEL-like motif deters secretion of soluble endoplasmic reticulum (ER) – resident proteins (Lewis *et al.*, 1990; Semenza *et al.*, 1990; Raykhel *et al.*, 2007). There are 24 possible variants of ER retention signals listed as a PROSITE motif - [KRHQSA]-[DENQ]-E-L in the PROSITE database (Hulo *et al.*, 2004; Hulo *et al.*, 2006). In addition, there are some ER retention signals that do not fit into the PROSITE motif (Raykhel *et al.*, 2007). A few serpins that are apparently engaged in secretory pathway are possessing such peptide sequences at their C-terminal ends (Ragg, 2007), distributed in organisms of wide evolutionary spectrum. In early diverging deuterostomia, neuroserpin orthologs like Spu-spn-1 of *Strongylocentrotus* and the Spn-1 gene of lancelets contain HEEL and KDEL, respectively, at their C-terminal end. Furthermore, the putative neuroserpin ortholog, Nve-



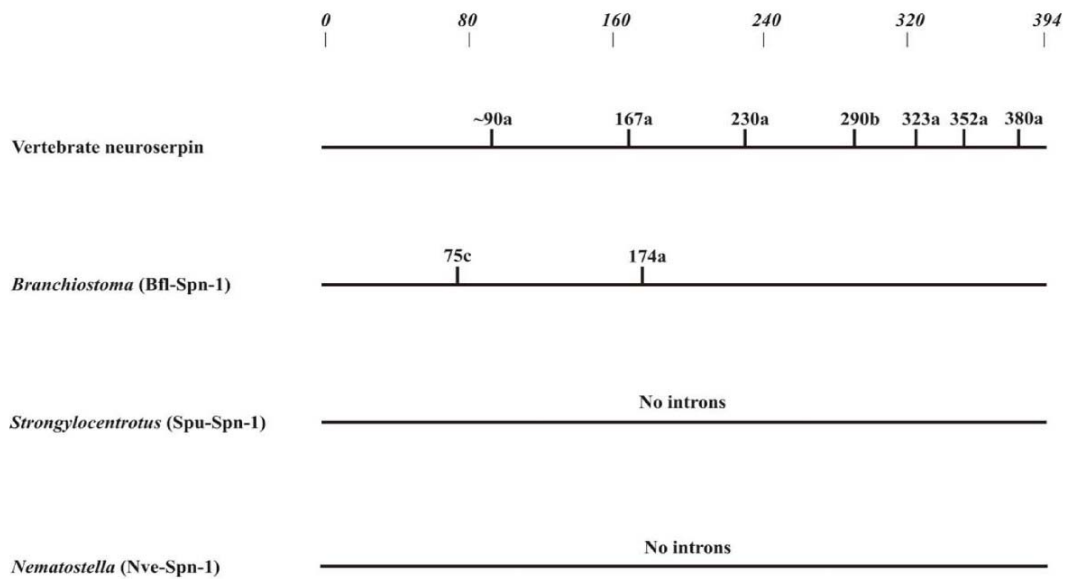
Spn-1 sequence from the sea anemone (*N. vectensis*) has SDEL at the C-terminus, which fits as one variant of ER retention/retrieval signal. Thus, it is clear that the neuroserpin ortholog, as corroborated by synteny analysis, from these animal species possesses one of the above mentioned 24 variants of ER-retention signals. In contrast, the C-terminal end of neuroserpin from tetrapods is HDFEEL (**Figure 58b**). In HeLa cells that express three different KDEL receptors with overlapping, but differential passenger specificities, the FEEL sub-sequence targets attached passenger proteins primarily to the Golgi, though one-fourth of cells depict ER localization (Raykhel *et al.*, 2007), whereas, in transfected COS cells, intracellular neuroserpin localizes to either the ER or Golgi (Ishigami *et al.*, 2007). In cells with a regulated secretory pathway, however, neuroserpin resides in large dense core vesicles, mediated by a C-terminal extension encompassing the last 13 amino acids (ETMNTSGHDFEEL) including the FEEL sequence (Ishigami *et al.*, 2007). Collectively, these data suggest that in orthologs of neuroserpin from deep-branching metazoans, a two amino acid insertion “FE” constitutes (in combination with additional residues?) a modified sorting signal attributing a more specialized subcellular localization. The surveillance of the secretory pathway routes by serpins is an ancient and conserved trait in eukaryotes as indicated by the putative neuroserpin ortholog present in the sea anemone genome. It will be interesting to investigate experimentally, whether the C-terminal extensions of neuroserpin orthologs from fishes are functional and mediate differential localization in a similar fashion as mammalian neuroserpin. Due to variations in their RSL region, ER-localized serpins may work differently in the secretory pathway. In vitro, neuroserpin from vertebrates inhibits tissue-type plasminogen activator (tPA) using the Arg residue at the P1 position in the RSL region (Osterwalder *et al.*, 1998). The cleavage site of Bfl-sp-1 is preceded by the dipeptide motif Lys-Arg (KR), a discernable feature for substrates and inhibitors of proprotein convertases (PCs). Similar sequences were found for Bla-Spn-1 from *B. lanceolatum* (Bentele *et al.*, 2006). Since the serpins Bfl-sp-1 of *B. floridae*, Spu-sp-1 of the sea urchin, and Nve-Spn-1 of the sea anemone also possess the Lys-Arg dipeptide motif (KR), a similar physiological role of these serpins has to be expected. Questions remain open concerning the presence of a neuroserpin ortholog in the arthropod lineage. Several labs have investigated a serpin acting as furin inhibitor - *Spn4* equipped with a classical ER targeting signal (HDEL) in *D. melanogaster* (Oley *et al.*, 2004; Osterwalder *et al.*, 2004; Richer *et al.*, 2004) and a homologous gene - *SRPN10* in *Anopheles gambiae* (Danielli *et al.*, 2003). However, the orthology of these genes to neuroserpin is unclear; because of following reasons:- (i) homoplasy due to convergent evolution and (ii) recombination events in protein coding regions especially in the RSL coding region (Börner and Ragg, 2008). Thus, a meticulous investigation will be needed in order to establish the relationships among the *Spn4* gene from *D. melanogaster* or *SRPN10* from *A. gambiae* and neuroserpin orthologs from deuterostomes. On comparing exon-intron structures of neuroserpin and PDCD10 genes, these two closely associated genes have very different fates in terms of intron patterns over deep-animal evolution. PDCD10 orthologs have undergone few changes in the exon-intron architecture since divergence of lineages leading to sea anemones and vertebrates (**Figure 70**). In PDCD10 genes, six out of eight intron positions occurring in humans or in the cnidarian are conserved. This is in favour of previous reports adducing that the majority of genes from early

diverging present-day eumetazoans are intron rich with most introns apparently maintained since ancient times (Raible *et al.*, 2005; Putnam *et al.*, 2007).



**Figure 70: Intron positions of *PDCD10* genes in metazoans.** Intron positions (white-on-black printing, phasing not indicated) were identified with GENEWISE and mapped onto the protein sequences. Intron positions conserved in at least two species are marked with an arrowhead. Accession numbers for *PDCD10* sequences: AAH16353 (*H. sapiens*); XP\_001186662 (*S. purpuratus*); EDO34838 (*N. vectensis*); AAF55190 (*D. melanogaster*); CAA90115 (*C. elegans*). Adopted from Kumar and Ragg (2008).

In contrast, the circumstances appear to be reverse for serpin genes. The putative sea anemone *neuroserpin* ortholog *Nve-Spn-1* and the sea urchin *neuroserpin* ortholog *Spu-spn-1* possess no intron within their serpin core domains (**Figure 71**). *Neuroserpin* orthologs from two lancelets - *B. floridae* and *B. lanceolatum* (Bentele *et al.*, 2006) - demonstrates two introns mapping to identical sites within the serpin body with no matching intron position with any known intron positions from vertebrate serpin genes (Ragg *et al.*, 2001).



**Figure 71: Exon-intron organisation of the neuroserpin gene lineage.** The *N. vectensis* serpin gene Nve-Spn-1 is included, though orthology with the deuterostome counterparts is currently only supported by protein-based signature sequences. Specifications for intron positions and their phasings refer to mature human  $\alpha$ 1-antitrypsin. Only introns mapping to the serpin core domain (residues 33 to 394 of the reference) are considered. Adopted from Kumar and Ragg (2008).

Thus, it is noteworthy that a substantial fraction of introns is recent in the serpin lineage leading to mammalian *neuroserpin* and they may have been inserted during metazoan evolution. Nevertheless, an alternate explanation cannot be snubbed, namely that large intron loss events are responsible for gene architecture of present-day serpin genes from cnidarians (and in *neuroserpin* orthologs from sea urchins and lancelets), whereas most other introns that have survived in these animals. Intron insertion is possibly not as rare as sometimes believed (Zhuo *et al.*, 2007), however, it could be confined to certain gene families and/or to discrete evolutionary phases (Babenko *et al.*, 2004), for as yet unexplored reasons.

The pancpin gene is localized in close proximity to the neuroserpin gene (**Figure 56**). Pancpin also possesses a C-terminal extension and indels like neuroserpin, suggesting its close relatedness to these proteins. Its lack in fishes led us to conclude that the pancpin gene might have originated by tandem duplication of neuroserpin after separation of tetrapods from the fish lineage. From fishes to human, the remaining group V3 members such as PAI1 (**Figure 53**), GDN (**Figure 54**) and serpinE3 (**Figure 55**) are present at various genomic locations. This suggests that they evolved independently since the origin of vertebrates. Since we were unable to trace even a single member of group V3 in the initial version of the lamprey genome, it is not clear whether they were present in basal vertebrates, such as lampreys and hagfishes. Thus, it should be interesting to investigate when group V3 genes originated. Since there is also no evidence for orthologs of these genes in invertebrate model genomes, they might have been originated during the emergence of vertebrates.

### 6.5 Evolutionary history of group V4 serpins

Group V4 has three members in mammalian genomes – pigment epithelium derived factor (PEDF),  $\alpha_2$ -antiplasmin, and C1-inhibitor. I have identified a fish-specific group V4 (FSG4) gene. The FSG4 protein has a serpin core domain, in addition to two immunoglobulin domains. Similar proteins were identified from other fishes (Inoue *et al.*, 1997; Wang and Secombes, 2003) other than analyzed fish genomes in this study. The well-conserved Ig domains in FSG4 proteins in several fishes imply that they must have functional significance, probably by supplying binding sites for the extracellular matrix and plasma proteins to enable its function to be strictly regulated. Indeed, the function of the mammalian C1IN can be regulated by ligand binding of heparin, type IV collagen, lamin, and entactin (Patston and Schapira, 1997; Bos *et al.*, 2002). The function of the Ig domains in the FSG4 is an issue of further investigations. However, during this study, evidence for syntenic conservation of C1IN from tetrapods and FSG4 of fishes was not found, suggesting that their origin is independent. During diversification events, fishes may have lost the original C1IN loci, and the duplicated C1IN loci diverged and during this process, one of the copies might have been fused with Ig-like domains. This speculation suggests that the FSG4 gene possibly functions differently as compared to C1IN genes of higher vertebrates (neofunctionalization).

On the protein sequence level, it looks that PEDF and  $\alpha_2$ -AP-like genes are conserved throughout vertebrates. However, on scrutiny of group V4 serpins, orthologs of most human group V4 serpins other than A2AP1\_FRU in *Fugu* cannot be found in current genomic sequence versions of fish genomes (**Figure 60**). It appears that fishes lost the orthologs of group V4 serpin loci present in tetrapods. Instead, they have paralogs, probably due to fish-specific genome duplications and diversifications. Moreover, in the draft version of the lamprey genome, two members of group V4 were detected resembling  $\alpha_2$ -AP like genes named as A2APL1\_PMA and A2APL2\_PMA (**Figure 63**). This suggests that group V4 exists since the beginning of vertebrates.

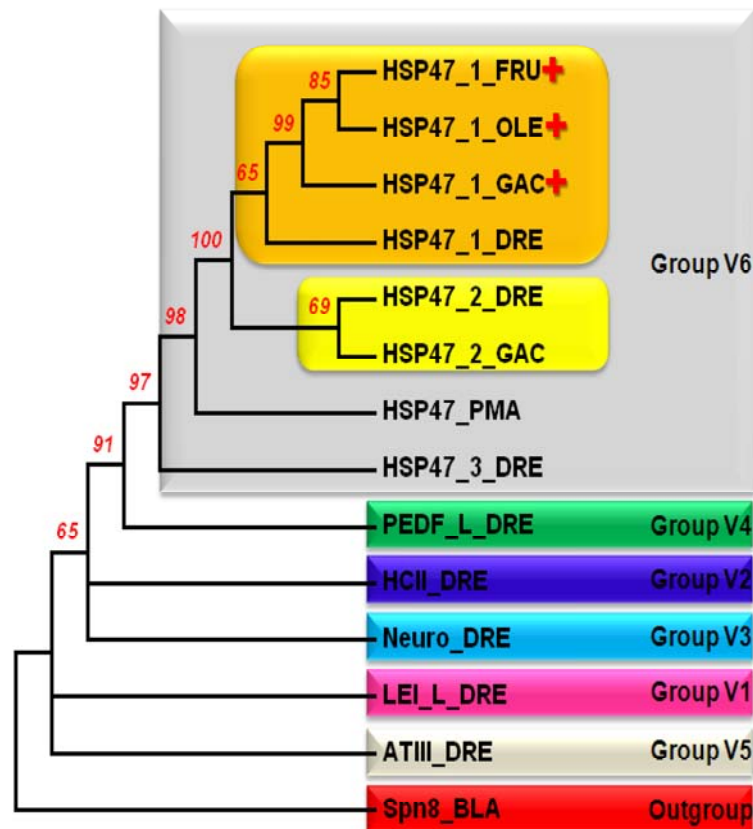
### 6.6 Summary of group V5 serpins

Group V5 consists of a single member - antithrombin III (ATIII). Grounded on gene structure and synteny conservation, the presence of the conserved helix D sequence, RCL-sequence, three pairs of disulfide bridges and other features, the ATIII gene is found to be maintained from fish to human. A remarkable difference between fish ATIII and the orthologs from tetrapods was observed with regard to the intron at position 262c (also a characteristic of group V1). This finding suggests that group V1 and group V5 are closely related. Furthermore, it is suspected that this intron was lost in ATIII from tetrapods. The inability to identify ATIII gene in current genomic assembly of lamprey hinders further tracing of the 262c intron. Additionally, the intron at position 339c of the ATIII gene is found in several serpins from an array of evolutionary distant organisms, such as *C. elegans* (Zang and Maizels, 2001), *B. malayi* (Zang *et al.*, 1999), lancelets, and *C. intestinalis*. It was not possible to unravel evolutionary history of the 339c intron from the datasets used for this study. It

would be interesting to unravel whether intron 339c is ancestral or has independently emerged multiple times.

### 6.7 Overview of group V6 serpins

Group V6 comprises the HSP47 gene and its paralogs in different vertebrates. Tetrapods have a single copy of the HSP47 gene, while there are two or three HSP47-like genes in some fishes as demonstrated in **Figure 72**.



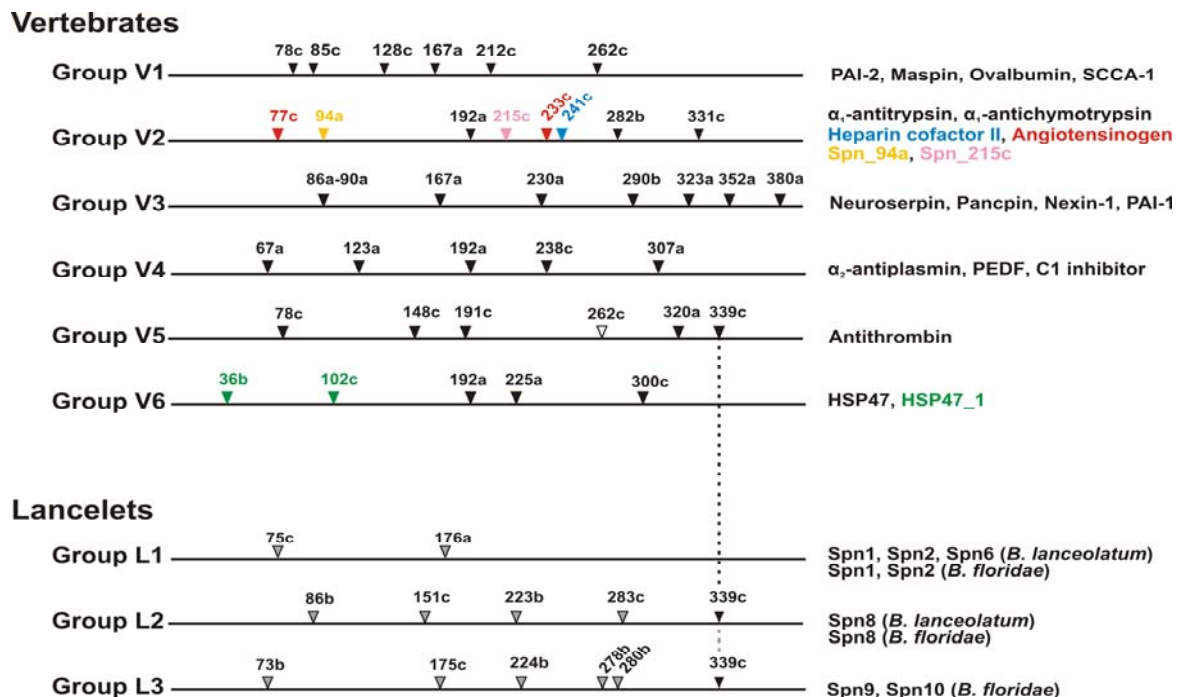
**Figure 72: Summary of evolutionary history of HSP47-related serpins from fishes.** The phylogram is based on the Neighbor-Joining method and includes one representative each of *serpin* groups V1 to V5 from *Danio rerio* (DRE). Within group V6, HSP47\_1 (orange box) and HSP47\_2 genes (yellow box) constitute separate sub-trees, supported by their genomic localization. + indicates the presence of two novel introns at positions 36b and 102c, respectively. The percentage of replicate trees, in which serpins clustered together in the bootstrap test (1000 replicates), is indicated. The outgroup in this tree is Spn8\_BLA (Spn8 gene from *Branchiostoma lanceolatum*, Genbank accession id. - FM242707). PMA: lamprey, GAC: stickleback, OLE: Japanese medaka.

Intron patterns and syntenic organizations of HSP47-like genes shed light on the evolution of HSP47 genes in fishes. *D. rerio* has three HSP47-related genes (named HSP47\_1\_DRE, HSP47\_2\_DRE and HSP47\_3\_DRE) containing the standard introns of group V6. Orthologs of HSP47\_1\_DRE were identified in *Fugu*, stickleback and medaka (HSP47\_1\_FRU, HSP47\_1\_GAC, and HSP47\_1\_OLE, respectively) by analysis of syntenic conservation. All orthologs with the exception of HSP47\_1\_DRE have novel introns at positions 36b and 102c, respectively. HSP47\_2\_DRE and HSP47\_2\_GAC genes cluster together in the phylogenetic

tree (yellow box in **Figure 72**) and they are orthologs, since they share a similar genomic microenvironment. Both genes possess the standard introns of group V6, like HSP47 from lamprey, tetrapods, and all three HSP47-like genes of *D. rerio*. These findings suggest that introns 36b and 102c of the HSP47\_1 orthologs have been gained in selected ray-finned fishes. The alternative possibility, intron loss in HSP47\_1\_DRE, however, cannot be excluded. This issue is discussed in detail in the next section.

### 6.8 Intron gain and loss in vertebrate serpins

Gain or loss of spliceosomal introns are rare events in evolution, which can serve as markers for phylogenetic analysis. Intron gain has been reported to be very rare in many metazoan lineages, including mammals and other vertebrates (Coulombe-Huntington and Majewski, 2007; Loh *et al.*, 2008). However, during this study, I found several instances of newly acquired introns in a single vertebrate protein superfamily, the serpins, while a single apparent intron loss event in specific ray-finned fishes became evident. Combining data from cDNA and gene sequences of serpin genes from *L. fluviatilis*, and *B. lanceolatum* obtained from other members of AG Zelluläre Genetik, Bielefeld and study of genomic sequences of stickleback, medaka and *Petromyzon marinus*, this finding is further authenticated (Ragg *et al.*, 2009).



**Figure 73: Gene structure comparisons between vertebrate (V1-V6) and lancelet (L1-L3) serpin groups.** Novel intron positions are marked in different colors with corresponding serpin genes.

The angiotensinogen gene provides a good example of intron gain. The introns at positions 77c and 233c found in some orthologs of this gene appear after the split of the *D. rerio* lineage from the other actinopterygians, whereas lampreys, tetrapods and *D. rerio* depict the standard exon/intron pattern of group V2. Similarly, all other non-standard introns found in

genes of the serpin superfamily (**Table 34**) are also confined to selected ray-finned fishes. Intron insertions have been proposed to occur primarily at “proto-splice sites” with the consensus sequence MAG↑R, where M is A/C, R is A/G, and the arrow (↑) represents the intron insertion site (Dibb and Newman, 1989). These sites are considered as “hot spot” for intron acquisition events (Coghlan and Wolfe, 2004; Qiu *et al.*, 2004; Sadusky *et al.*, 2004; Tordai and Patthy, 2004). The listed novel introns of specific serpins (**Table 34**) are characterized by following features: (i) canonical proto-splice site with some exceptions (marked in white-on-black printing in **Table 34**). (ii) intron sizes ranging from 68-178 base pairs with pre-dominant intron phasing c (five out of seven listed introns and one each for phase a and b, respectively). These novel intron positions listed for group V2 serpins are neither found in any other paralogs of vertebrate group V2 nor in any other vertebrate serpin genes reported so far. These introns are most likely acquired *de novo* rather than inherited from a common ancestor.

**Table 34: Sequences flanking the insertion points of novel introns in vertebrate *serpin* genes.** Adopted from Ragg *et al.* (2009).

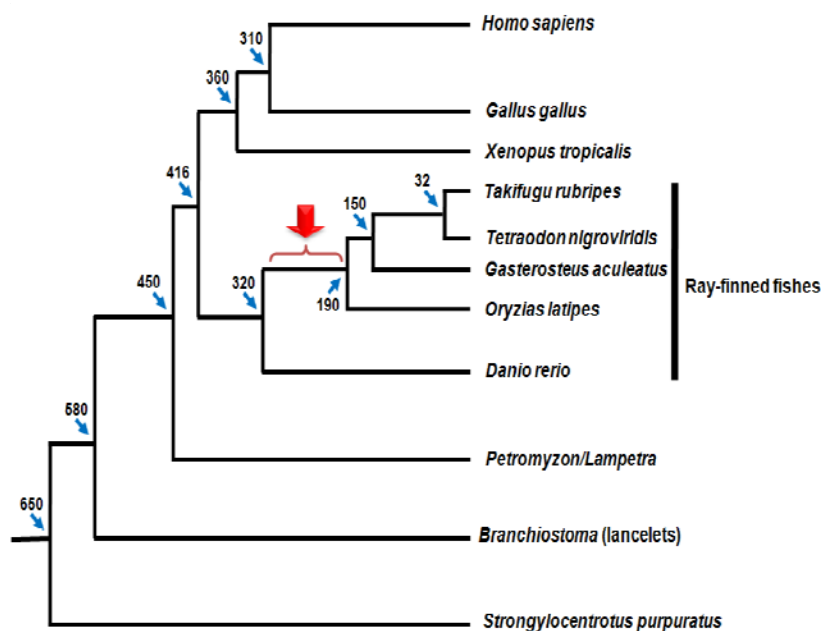
Species	Gene	Intron (Intron Size)	Flanking Sequences
<i>T. rubripes</i>	Angiotensinogen	77c (75)	CCAG↑TCTC
<i>G. aculeatus</i>	Angiotensinogen	77c (140)	CCAG↑TACC
<i>O. latipes</i>	Angiotensinogen	77c (82)	TCTG↑CGTC
<i>T. rubripes</i>	Angiotensinogen	233c (80)	TAAG↑GTTC
<i>G. aculeatus</i>	Angiotensinogen	233c (112)	TAAG↑GTAC
<i>O. latipes</i>	Angiotensinogen	233c (80)	TAAG↑TTGA
<i>T. rubripes</i>	HCII	241c (75)	ACAG↑CTCC
<i>T. nigroviridis</i>	HCII	241c (70)	ACAG↑CTCC
<i>G. aculeatus</i>	HCII	241c (82)	ACAG↑CTCC
<i>O. latipes</i>	HCII	241c (98)	ACAG↑CTCC
<i>T. rubripes</i>	HSP47_1	36b (178)	TCAG↑CCTC
<i>G. aculeatus</i>	HSP47_1	36b (141)	TCAG↑CCTC
<i>O. latipes</i>	HSP47_1	36b (100)	TTAG↑CCTT
<i>T. rubripes</i>	HSP47_1	102c (88)	TCAG↑TTGA
<i>G. aculeatus</i>	HSP47_1	102c (123)	CGAG↑GTGA
<i>O. latipes</i>	HSP47_1	102c (97)	TCAG↑GTGA
<i>T. rubripes</i>	Spn_94a	94a (68)	CCAG↑AGCT
<i>T. nigroviridis</i>	Spn_94a	94a (68)	CCAG↑ATCT
<i>G. aculeatus</i>	Spn_94a	94a (74)	CCAG↑ATCT
<i>O. latipes</i>	Spn_94a	94a (111)	CCAG↑ATCT
<i>T. rubripes</i>	Spn_215c	215c (76)	CAAG↑GTTC
<i>T. nigroviridis</i>	Spn_215c	215c (68)	CAAG↑GTCC

Arrows indicate the intron insertion points. Bases deviating from the proto-splice site sequence (MAG↑R) (Dibb and Newman, 1989) are printed in white-on-black.

No novel intron appears to have been acquired at the expense of adjacent introns as no losses of standard introns are found in these serpin genes. There are no non-standard introns in *serpin* genes from other vertebrate taxa. Thus, these novel introns appear to have been gained



during radiation of actinopterygians. Similar findings were reported for several of the few other well-documented cases of novel introns in vertebrates (Figueroa *et al.*, 1995; Venkatesh *et al.*, 1999; Schioth *et al.*, 2005; Moriyama *et al.*, 2008). Vertebrate group V6 serpins provide a clear picture about time period and processes possibly associated with intron gain. During evolution of ray-finned fishes, group V6 was split into three lineages as found in *D. rerio*, probably a consequence of whole genome and/or large fragment duplications. The extra introns at positions 36b and 102c however, are only unraveled in one lineage, the *HSP47\_1* orthologs from *Fugu*, stickleback, and medaka. In contrast, the second lineage comprised of *HSP47\_2* genes from stickleback and zebrafish depicts the standard intron pattern of group V6 as found in other HSP47-like genes from *D. rerio*, lampreys and tetrapods. These findings corroborate that intron gain was not associated with the fish-specific genome duplication events. Moreover, they are supportive of the concept that newly gained introns are maintained by the existence of several gene copies. The estimated timing of intron gains can be confined to the period before or during emergence of the stickleback/Japanese medaka/pufferfish lineage at about 320-190 MYA (Figure 74) by using phylogenetic timescale information (Ponting, 2008).



**Figure 74: Phylogenetic tree of vertebrates emphasizing timescale and lineages displaying intron gain in *serpin* genes.** The estimated divergence times (in mya) were taken from Ponting, 2008 and are marked with blue arrows. The time interval of intron gains is indicated (red arrow). Adopted from Ragg *et. al.* 2009.

With exception of intron 94a in the *Spn\_94a* gene, all novel introns have exact insertion points without any deletions / insertions at the borders. Indels are tolerated by many serpins, largely in the regions of non-secondary structural elements as reported in various sequence alignments without any functional implications and thus, it is important to understand when considering mechanisms of intron gain.

The nucleotides flanking the novel introns correspond to the proto-splice site as previously proposed (Dibb and Newman, 1989; Sverdlov *et al.*, 2004). However, a considerable fraction of the affiliated 3'-exons starts with pyrimidine rather than purine. The 3'-side flanking insertions thus appears to be less stringent than insinuated and may rather resemble MAG↑N, at least in vertebrates. There are several mechanism supposed to be responsible for birth of



introns (Roy and Gilbert, 2006; Irimia *et al.*, 2008). Duplication events operating on pre-existing sequences at some stage are vital to these mechanisms. The activity of transposons is believed to be responsible for intron insertion (Roy and Gilbert, 2006). Fish genomes are characterized by their diversity of retrotransposable elements, especially retrotransposons. Selected retrotransposons are reported to be active in recent times in fish lineages (Aparicio *et al.*, 2002; Volff *et al.*, 2003; Volff, 2005). No significant similarity to known repetitive elements was detected in the non-standard introns (Ragg *et al.*, 2009). Hence, the involvement of duplication dependent transposons in these intron gain events becomes unlikely, though preferential loss of transposons from newly inserted introns cannot be ignored. The source of these newly acquired introns remains open for investigations. Largely, due to the fact that during my searches using different homology search suites (BLAST suite or FASTA suite), these novel introns do not have significant homology either to flanking sequence within the locus or anywhere else in overall sequenced parts of these fish genomes (analysis data stored in GENLIGHT, not shown here as it do not provide any significant information). Every genome sequencing project faces some problems in sequencing process, which result in unavailability of small proportion of genomic sequences. Similarly, fractions of genomic sequences from these selected fish genomes are likely to remain unsequenced, so the remote possibility exists that the novel introns are derived from some unsequenced portion of these fish genomes. A similar finding is reported for a documented case of strain specific intron insertion in *Daphnia pulex* (Omilian *et al.*, 2008).

Several different types of processes are believed to be responsible for intron births. Nevertheless, these processes not necessarily need to be related with the events responsible for the primordial emergence of spliceosomal introns. Excision of intron sequences, probably created by expansion of short simple repeats or more complex repetitive elements (Figuroa *et al.*, 1995; Zhuo *et al.*, 2007) or by intronization of exonic sequences (Irimia *et al.*, 2008), manifests that the spliceosome will operate as long as the essential splice signals are present, not matter how the introns were generated. It is difficult to find out some clear hints supporting currently discussed intron gain mechanisms. To find some solutions on this issue, genome size of selected animals were compared (**Table 35**) as proposed by Ragg *et al.* (2009) using the Animal genome size database (Gregory, 2008), which is a comprehensive database of genome size studies. Every eukaryotic species has a characteristic amount of genomic DNA. The amount of this DNA in the haploid cell of a species is called C-value. C-value is expressed in base pairs or picogram or molecular weight (daltons). One picogram of DNA corresponds to approximately 1Gb. The lack of co-relation between genomic size (C-value) and biological complexity is called as C-value paradox (Hartl, 2000). Selected ray-finned fishes exhibit considerable reduction in genomic contents (marked in bold in **Table 35**) as compared to lampreys and zebrafish. Apparently, after the fish-specific whole genome duplication; compaction processes have led to a considerable reduction of genome sizes in many actinopterygians (Hinegardner, 1968; Aparicio *et al.*, 2002; Vandepoele *et al.*, 2004; Gregory, 2008). Reduction in genomic DNA size may affect three levels: (a) whole genes. (b) intergenic regions and (c) intronic sequences. *D. rerio* possesses considerably larger introns than pufferfishes. It is a fascinating quest whether this provides some clues to the mechanisms of intron insertion.

**Table 35: Genome size of selected animals.** A depletion of size in selected ray-finned fishes is evident (marked in bold). Source: Animal genome size database (Gregory, 2008).

Selected organism	Species	Genome size - Total haploid DNA content C-value in pictogram (pg)
Human	<i>H. sapiens</i>	3.50
Chicken	<i>G. domesticus</i>	1.25
Frog	<i>X. laevis</i>	3.69
Zebrafish	<i>D. rerio</i>	1.78
<b>Fugu</b>	<i>T. rubripes</i>	<b>0.40</b>
<b>Tetraodon</b>	<i>T. nigroviridis</i>	<b>0.35</b>
<b>Stickleback</b>	<i>G. aculeatus</i>	<b>0.70</b>
<b>Medaka</b>	<i>O. latipes</i>	<b>0.75</b>
Sea lamprey	<i>P. marinus</i>	2.44
European river lamprey	<i>L. fluviatilis</i>	1.45
Sea squid	<i>C. intestinalis</i>	0.20
Lancelet	<i>B. lanceolatum</i>	0.59
Sea urchin	<i>S. purpuratus</i>	0.89
Fly	<i>D. melanogaster</i>	0.18
Worm	<i>C. elegans</i>	0.10
Sea anemone	<i>N. vectensis</i>	0.23

Depletion of genomic contents is considerably associated with events of DNA breakage and recombination that further require DNA repair and recombination events. Intron acquisition events probably are associated with such changes in DNA contents of these fishes (Ragg *et al.*, 2009).

Formations of new genes either by whole genome duplication or by tandem duplication of paralogs, might conceivably favor intron gain and maintenance of novel introns, since an unaffected gene copy remains within in the genome. A direct co-relation between intron insertions in serpin genes and events of gene/genome duplications cannot be established, since *D. rerio* does not possess any of these novel introns even though, the time of divergence of *D. rerio* is closer to the fish-specific genome duplication event (Meyer and Van de Peer, 2005) than that of the other fishes investigated. However, conservation of novel introns could indeed be favored by the co-existence of paralogs. Other well-documented cases of intron gain apply to multi-membered gene families (Figuroa *et al.*, 1995; Schioth *et al.*, 2005), however, association of preferred intron gain with multi-copy gene families is still controversial.

In conclusion, a group of ray-finned fishes exhibits multiple intron insertions in selected serpins that are not shown by any other vertebrates. Depletion in genomic contents of these fishes may have played a crucial role in these intron acquisitions. Fishes exhibit a high diversity after separation from last common ancestor of tetrapods/fishes lineage and these diversities can be explained by rapid change in DNA contents by processes such as whole genome duplication and genome compaction. Losses/gains in gene contents, introns, and intergenic regions are crucial to these events.

### 6.9 Strength and weakness in this work

This work has extended the intron-coded group V1-V6 serpin classification system across different vertebrates from lamprey to human. Prior to this work, this classification system based on exon-intron architecture was established mainly from mammals (Ragg *et al.*, 2001). During this work, this serpin classification system has been validated over different types of vertebrates including birds, frogs, and fishes. Orthologs and paralogs of human serpins were also determined during this work based on exon-intron architecture, micro-synteny analysis, and sequence motifs. This work has highlighted the weakness of solely sequence-based methods for evaluating orthology/paralogy such as bidirectional BLAST, since these computational tools can only provide evidence of homology; and only limited insight into the origin of a gene. Synteny analysis, coupled with gene architecture and motif features, provide a better solution for assigning orthology of genes. During this work, many serpins were clearly classified as orthologs or paralogs of their human counterparts, especially in case of fish serpins which are enigmatic due to fish-specific whole genome duplication events (Ohno, 1970) that led to many paralogs. During this work, the deep evolutionary roots of mammalian neuroserpin (a secretory-pathway associated serpin) was analyzed and resolved at least since the emergence of deuterostomes and most probably even since divergence of *Bilateria* from eumetazoans using synteny, rare indels and sequence motif data (Kumar and Ragg, 2008). Furthermore, this work was instrumental in unraveling the intron gains in specific serpin genes in selected ray-finned fishes (Ragg *et al.*, 2009). Overall, a validated classification system for vertebrate serpins now exists, into which serpin family members from newly sequenced genomes of vertebrates can be easily incorporated. The proposed neofunctionalization or subfunctionalization in serpins needs to be experimentally validated especially in the case of  $\alpha_2$ -AP and HSP47 of fishes.

### 6.10 Outlook

By the time this work was finished, the lamprey genomic sequence was in an initial stage. It would be interesting to investigate the lamprey genome more closely in order to unravel vertebrate group V3 and group V5 serpins since by now members of these groups were not identified in this species. Similarly, improvements in genomic assemblies of the sea urchin and the sea anemone may provide new insight into deep metazoan evolution of serpins. As by now, only evidence for neuroserpin orthologs in these animals is available. The genomic sequences of hagfish, representing a basal vertebrate will be helpful to unravel the structure and function of serpins at the origin of vertebrates in more details. Group V1 and group V2 have multiple paralogs in different mammalian genomes, a study of basal mammals such as marsupials and Platypus will help in understanding the molecular mechanisms of their expansions. Furthermore, it is recognized that the intron at position 339a is found in some serpin genes from various organisms (that are important in metazoan evolution). It would be interesting to investigate whether this intron is ancestral or has emerged multiple times independently.

---

## 7. References

- Adams, M.D., *et al.* (2000) The genome sequence of drosophila melanogaster. *Science*, **287**, 2185-2195.
- Altschul, S.F., *et al.* (1990) Basic local alignment search tool. *J Mol Biol*, **215**, 403-410.
- Altschul, S.F. and Lipman, D.J. (1990) Protein database searches for multiple alignments. *Proc Natl Acad Sci U S A*, **87**, 5509-5513.
- Altschul, S.F., *et al.* (1997) Gapped blast and psi-blast: A new generation of protein database search programs. *Nucleic Acids Res*, **25**, 3389-3402.
- Aparicio, S., *et al.* (2002) Whole-genome shotgun assembly and analysis of the genome of fugu rubripes. *Science*, **297**, 1301-1310.
- Apweiler, R., Bairoch, A. and Wu, C.H. (2004a) Protein sequence databases. *Curr Opin Chem Biol*, **8**, 76-80.
- Apweiler, R., *et al.* (2004b) Uniprot: The universal protein knowledgebase. *Nucleic Acids Res*, **32**, D115-119.
- Atchley, W.R., *et al.* (2001) Phylogenetic analyses of amino acid variation in the serpin proteins. *Mol Biol Evol*, **18**, 1502-1511.
- Babenko, V., Rogozin, I., Mekhedov, S. and Koonin, E. (2004) Prevalence of intron gain over intron loss in the evolution of paralogous gene families. *Nucleic Acids Res*, **32**, 3724 - 3733.
- Backovic, M. and Gettins, P.G. (2002) Insight into residues critical for antithrombin function from analysis of an expanded database of sequences that includes frog, turtle, and ostrich antithrombins. *J Proteome Res*, **1**, 367-373.
- Bairoch, A. (1991) Prosite: A dictionary of sites and patterns in proteins. *Nucleic Acids Res*, **19 Suppl**, 2241-2245.
- Bairoch, A. and Apweiler, R. (1996) The swiss-prot protein sequence data bank and its new supplement trembl. *Nucleic Acids Res*, **24**, 21-25.
- Bairoch, A. and Apweiler, R. (2000) The swiss-prot protein sequence database and its supplement trembl in 2000. *Nucleic Acids Res*, **28**, 45-48.
- Bairoch, A. and Boeckmann, B. (1991) The swiss-prot protein sequence data bank. *Nucleic Acids Res*, **19 Suppl**, 2247-2249.
- Bairoch, A., Boeckmann, B., Ferro, S. and Gasteiger, E. (2004) Swiss-prot: Juggling between evolution and stability. *Brief Bioinform*, **5**, 39-55.
- Barker, W.C., *et al.* (1987) Protein sequence database of the protein identification resource (pir). *Protein Seq Data Anal*, **1**, 43-98.
- Beckstette, M., *et al.* (2004) Genlight: An interactive system for high-throughput sequence analysis and comparative genomics. . *Journal of Integrative Bioinformatics*, **1(1)**, 79-94.
- Beckstette, M.S., A. Selzer, PM. (2004) Genlight: An interactive system for high-throughput sequence analysis and comparative genomics. . *German Conference on Bioinformatics, GCB 2004*. GI-Edition, **Bielefeld University, Bielefeld, Germany**, pp. P-53.
- Benarafa, C. and Remold-O'Donnell, E. (2005) The ovalbumin serpins revisited: Perspective from the chicken genome of clade b serpin evolution in vertebrates. *Proc Natl Acad Sci U S A*, **102**, 11367-11372.
- Benson, D.A., *et al.* (2005) Genbank. *Nucleic Acids Res*, **33**, D34-38.
- Benson, D.A., *et al.* (2006) Genbank. *Nucleic Acids Res*, **34**, D16-20.

- Bentele, C., *et al.* (2006) A proprotein convertase-inhibiting serpin with an endoplasmic reticulum targeting signal from branchiostoma lanceolatum, a close relative of vertebrates. *Biochem J*, **395**, 449-456.
- Berman, H.M., *et al.* (2000) The protein data bank. *Nucleic Acids Res*, **28**, 235-242.
- Birney, E., *et al.* (2006) Ensembl 2006. *Nucleic Acids Res*, **34**, D556-561.
- Birney, E., Clamp, M. and Durbin, R. (2004) Genewise and genomewise. *Genome Res*, **14**, 988-995.
- Boeckmann, B., *et al.* (2003) The swiss-prot protein knowledgebase and its supplement trembl in 2003. *Nucleic Acids Res*, **31**, 365-370.
- Boeke, J.D. (2003) The unusual phylogenetic distribution of retrotransposons: A hypothesis. *Genome Res*, **13**, 1975-1983.
- Börner, S. and Ragg, H. (2008) Functional diversification of a protease inhibitor gene in the genus drosophila and its molecular basis. *Gene*, **415**, 23 - 31.
- Bos, I.G., Hack, C.E. and Abrahams, J.P. (2002) Structural and functional aspects of c1-inhibitor. *Immunobiology*, **205**, 518-533.
- Bourque, G., *et al.* (2005) Comparative architectures of mammalian and chicken genomes reveal highly variable rates of genomic rearrangements across different lineages. *Genome Res*, **15**, 98-110.
- Brudno, M., Steinkamp, R. and Morgenstern, B. (2004) The chaos/dialign www server for multiple alignment of genomic sequences. *Nucleic Acids Res*, **32**, W41-44.
- Burge, C. and Karlin, S. (1997) Prediction of complete gene structures in human genomic DNA. *J Mol Biol*, **268**, 78-94.
- Burge, C.B. and Karlin, S. (1998) Finding the genes in genomic DNA. *Curr Opin Struct Biol*, **8**, 346-354.
- Canestro, C., Bassham, S. and Postlethwait, J.H. (2003) Seeing chordate evolution through the ciona genome sequence. *Genome Biol*, **4**, 208.
- Carrell, R. and Travis, J. (1985)  $A_1$ -antitrypsin and the serpins: Variation and countervariation. *Trends in Biochemical Science*, **10**, 20-24.
- Cavalier-Smith, T. (1991) Intron phylogeny: A new hypothesis. *Trends Genet*, **7**, 145-148.
- Chen, P.Y., *et al.* (2007) Two non-homologous brain diseases-related genes, serpinil and pdcd10, are tightly linked by an asymmetric bidirectional promoter in an evolutionarily conserved manner. *BMC Mol Biol*, **8**, 2.
- Chenna, R., *et al.* (2003) Multiple sequence alignment with the clustal series of programs. *Nucleic Acids Res*, **31**, 3497-3500.
- Choi, J.H., Jung, H.Y., Kim, H.S. and Cho, H.G. (2000) Phylodraw: A phylogenetic tree drawing system. *Bioinformatics*, **16**, 1056-1058.
- Clamp, M., Cuff, J., Searle, S.M. and Barton, G.J. (2004) The jalview java alignment editor. *Bioinformatics*, **20**, 426-427.
- Coghlan, A. and Wolfe, K.H. (2004) Origins of recently gained introns in caenorhabditis. *Proc Natl Acad Sci U S A*, **101**, 11362-11367.
- Conrad, B. and Antonarakis, S.E. (2007) Gene duplication: A drive for phenotypic diversity and cause of human disease. *Annu Rev Genomics Hum Genet*, **8**, 17-35.
- Cooper, N.R. (1985) The classical complement pathway: Activation and regulation of the first complement component. *Adv Immunol*, **37**, 151-216.
- Corbo, J.C., Levine, M. and Zeller, R.W. (1997) Characterization of a notochord-specific enhancer from the brachyury promoter region of the ascidian, ciona intestinalis. *Development*, **124**, 589-602.
- Coughlin, P.B. (2005) Antiplasmin: The forgotten serpin? *Febs J*, **272**, 4852-4857.
- Coulombe-Huntington, J. and Majewski, J. (2007) Characterization of intron loss events in mammals. *Genome Res*, **17**, 23-32.

- Crooks, G.E., Hon, G., Chandonia, J.M. and Brenner, S.E. (2004) Weblogo: A sequence logo generator. *Genome Res*, **14**, 1188-1190.
- Danielli, A., Kafatos, F. and Loukeris, T. (2003) Cloning and characterization of four anopheles gambiae serpin isoforms, differentially induced in the midgut by plasmodium berghei invasion. *J Biol Chem*, **278**, 4184 - 4193.
- Darnell, J.E., Jr. (1978) Implications of rna-rna splicing in evolution of eukaryotic cells. *Science*, **202**, 1257-1260.
- de Castro, E., *et al.* (2006) Scanprosite: Detection of prosite signature matches and prorule-associated functional and structural residues in proteins. *Nucleic Acids Res*, **34**, W362-365.
- Dehal, P., *et al.* (2002) The draft genome of ciona intestinalis: Insights into chordate and vertebrate origins. *Science*, **298**, 2157-2167.
- Detrich, H.W., 3rd, Westerfield, M. and Zon, L.I. (1999) Overview of the zebrafish system. *Methods Cell Biol*, **59**, 3-10.
- Dibb, N.J. and Newman, A.J. (1989) Evidence that introns arose at proto-splice sites. *EMBO J*, **8**, 2015-2021.
- Doolittle, R.F. and Feng, D.F. (1990) Nearest neighbor procedure for relating progressively aligned amino acid sequences. *Methods Enzymol*, **183**, 659-669.
- Driever, W. and Fishman, M.C. (1996) The zebrafish: Heritable disorders in transparent embryos. *J Clin Invest*, **97**, 1788-1794.
- Drummond, A., *et al.* (2006) Geneious v2.5, available from <http://www.Geneious.Com/>.
- Edgar, R.C. (2004a) Muscle: A multiple sequence alignment method with reduced time and space complexity. *BMC Bioinformatics*, **5**, 113.
- Edgar, R.C. (2004b) Muscle: Multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*, **32**, 1792-1797.
- Felsenstein, J. (1993) Phylip (phylogeny inference package) version 3.5c; distributed by the author.
- Felsenstein, J. (1996) Inferring phylogenies from protein sequences by parsimony, distance, and likelihood methods. *Methods Enzymol*, **266**, 418-427.
- Feng, D.F. and Doolittle, R.F. (1987) Progressive sequence alignment as a prerequisite to correct phylogenetic trees. *J Mol Evol*, **25**, 351-360.
- Ferguson, W.S. and Finlay, T.H. (1983) Localization of the disulfide bond in human antithrombin iii required for heparin-accelerated thrombin inactivation. *Arch Biochem Biophys*, **221**, 304-307.
- Figueroa, F., *et al.* (1995) Evidence for insertion of a new intron into an mhc gene of perch-like fish. *Proc Biol Sci*, **259**, 325-330.
- Force, A., *et al.* (1999) Preservation of duplicate genes by complementary, degenerative mutations. *Genetics*, **151**, 1531-1545.
- Forsyth, S., Horvath, A. and Coughlin, P. (2003) A review and comparison of the murine alpha1-antitrypsin and alpha1-antichymotrypsin multigene clusters with the human clade a serpins. *Genomics*, **81**, 336-345.
- Galtier, N., Gouy, M. and Gautier, C. (1996) Seaview and phylo\_win: Two graphic tools for sequence alignment and molecular phylogeny. *Comput Appl Biosci*, **12**, 543-548.
- Gandrille, S., *et al.* (1990) Important role of arginine 129 in heparin-binding site of antithrombin iii. Identification of a novel mutation arginine 129 to glutamine. *J Biol Chem*, **265**, 18997-19001.
- Gettins, P.G. (2002) Serpin structure, mechanism, and function. *Chem Rev*, **102**, 4751-4804.
- Gettins, P.G.W., Patston, P.A. and Olson, S.T. (1996) *Serpins: Structure, function and biology*. R. G. Landes Co., Austin.

- Gibbs, R.A., *et al.* (2004) Genome sequence of the brown norway rat yields insights into mammalian evolution. *Nature*, **428**, 493-521.
- Gilbert, W. (1978) Why genes in pieces? *Nature*, **271**, 501.
- Gooptu, B., *et al.* (2000) Inactive conformation of the serpin alpha(1)-antichymotrypsin indicates two-stage insertion of the reactive loop: Implications for inhibitory function and conformational disease. *Proc Natl Acad Sci U S A*, **97**, 67-72.
- Gouet, P., Courcelle, E., Stuart, D.I. and Metzoz, F. (1999) Esprict: Analysis of multiple sequence alignments in postscript. *Bioinformatics*, **15**, 305-308.
- Gough, J. and Chothia, C. (2002) Superfamily: Hmms representing all proteins of known structure. Scop sequence searches, alignments and genome assignments. *Nucleic Acids Res*, **30**, 268-272.
- Gough, J., Karplus, K., Hughey, R. and Chothia, C. (2001) Assignment of homology to genome sequences using a library of hidden markov models that represent all proteins of known structure. *J Mol Biol*, **313**, 903-919.
- Gregory, T.R. (2008) Animal genome size database, <http://www.Genomesize.Com>.
- Hartl, D.L. (2000) Molecular melodies in high and low c. *Nat Rev Genet*, **1**, 145-149.
- Hedges, S.B. (2002) The origin and evolution of model organisms. *Nat Rev Genet*, **3**, 838-849.
- Hendershot, L.M. and Bulleid, N.J. (2000) Protein-specific chaperones: The role of hsp47 begins to gel. *Curr Biol*, **10**, R912-915.
- Henikoff, S. and Henikoff, J.G. (1991) Automated assembly of protein blocks for database searching. *Nucleic Acids Res*, **19**, 6565-6572.
- Higgins, D.G. and Sharp, P.M. (1988) Clustal: A package for performing multiple sequence alignment on a microcomputer. *Gene*, **73**, 237-244.
- Higgins, D.G., Thompson, J.D. and Gibson, T.J. (1996) Using clustal for multiple sequence alignments. *Methods Enzymol*, **266**, 383-402.
- Hillier, L.W., *et al.* (2004) Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature*, **432**, 695-716.
- Hinegardner, R. (1968) Evolution of cellular DNA content in teleost fishes. *American Naturalist*, **102**, 517-523.
- Hubbard, T.J., *et al.* (2007) Ensembl 2007. *Nucleic Acids Res*, **35**, D610-617.
- Huber, R. and Carrell, R.W. (1989) Implications of the three-dimensional structure of alpha 1-antitrypsin for structure and function of serpins. *Biochemistry*, **28**, 8951-8966.
- Hulo, N., *et al.* (2006) The prosite database. *Nucleic Acids Res*, **34**, D227-230.
- Hulo, N., *et al.* (2004) Recent improvements to the prosite database. *Nucleic Acids Res*, **32**, D134-137.
- Huntington, J.A., Read, R.J. and Carrell, R.W. (2000) Structure of a serpin-protease complex shows inhibition by deformation. *Nature*, **407**, 923-926.
- Inoue, S., Nam, B.H., Hirono, I. and Aoki, T. (1997) A survey of expressed genes in japanese flounder (*paralichthys olivaceus*) liver and spleen. *Mol Mar Biol Biotechnol*, **6**, 376-380.
- Irimia, M., *et al.* (2008) Origin of introns by 'intronization' of exonic sequences. *Trends Genet*, **24**, 378-381.
- Irving, J.A., Pike, R.N., Lesk, A.M. and Whisstock, J.C. (2000) Phylogeny of the serpin superfamily: Implications of patterns of amino acid conservation for structure and function. *Genome Res*, **10**, 1845-1864.
- Ishigami, S., *et al.* (2007) Identification of a novel targeting sequence for regulated secretion in the serine protease inhibitor neuroserpin. *Biochem J*, **402**, 25-34.
- Izuhara, K., *et al.* (2008) Recent progress in understanding the diversity of the human ov-serpin/clade b serpin family. *Cell Mol Life Sci*.

- Jaillon, O., *et al.* (2004) Genome duplication in the teleost fish tetraodon nigroviridis reveals the early vertebrate proto-karyotype. *Nature*, **431**, 946-957.
- Jeffares, D.C., Mourier, T. and Penny, D. (2006) The biology of intron gain and loss. *Trends Genet*, **22**, 16-22.
- Jeffares, D.C., Poole, A.M. and Penny, D. (1998) Relics from the rna world. *J Mol Evol*, **46**, 18-36.
- Jordan, R.E. (1983) Antithrombin in vertebrate species: Conservation of the heparin-dependent anticoagulant mechanism. *Arch Biochem Biophys*, **227**, 587-595.
- Kaiserman, D. and Bird, P. (2005) Analysis of vertebrate genomes suggests a new model for clade b serpin evolution. *BMC Genomics*, **6**, 167.
- Kaiserman, D., *et al.* (2002) Comparison of human chromosome 6p25 with mouse chromosome 13 reveals a greatly expanded ov-serpin gene repertoire in the mouse. *Genomics*, **79**, 349-362.
- Kari, G., Rodeck, U. and Dicker, A.P. (2007) Zebrafish: An emerging model system for human disease and drug discovery. *Clin Pharmacol Ther*, **82**, 70-80.
- Kent, W.J., *et al.* (2002) The human genome browser at ucsc. *Genome Res*, **12**, 996-1006.
- Krem, M.M. and Di Cera, E. (2003) Conserved ser residues, the shutter region, and speciation in serpin evolution. *J Biol Chem*, **278**, 37810-37814.
- Krogh, A., *et al.* (1994) Hidden markov models in computational biology. Applications to protein modeling. *J Mol Biol*, **235**, 1501-1531.
- Krueger, O. (2003) Dna2aa, doctoral thesis. University of Bielefeld, Bielefeld.
- Krüger, O. (2003) Eine phylogenetische analyse von serpin-genen aus eukaryotischen genomprojekten. *Cellular Genetics, Technical Faculty*. University of Bielefeld, Bielefeld, Vol. Ph.D.
- Kumar, A. and Ragg, H. (2008) Ancestry and evolution of a secretory pathway serpin. *BMC Evol Biol*, **8**, 250.
- Kumar, S., Tamura, K., Jakobsen, I.B. and Nei, M. (2001) Mega2: Molecular evolutionary genetics analysis software. *Bioinformatics*, **17**, 1244-1245.
- Kumar, S., Tamura, K. and Nei, M. (2004) Mega3: Integrated software for molecular evolutionary genetics analysis and sequence alignment. *Brief Bioinform*, **5**, 150-163.
- Lamande, S.R. and Bateman, J.F. (1999) Procollagen folding and assembly: The role of endoplasmic reticulum enzymes and molecular chaperones. *Semin Cell Dev Biol*, **10**, 455-464.
- Lambowitz, A.M. and Zimmerly, S. (2004) Mobile group ii introns. *Annu Rev Genet*, **38**, 1-35.
- Lawrence, D.A., *et al.* (2000) Partitioning of serpin-proteinase reactions between stable inhibition and substrate cleavage is regulated by the rate of serpin reactive center loop insertion into beta-sheet a. *J Biol Chem*, **275**, 5839-5844.
- Leinonen, R., *et al.* (2004) Uniprot archive. *Bioinformatics*, **20**, 3236-3237.
- Lener, M., *et al.* (1998) Molecular cloning, gene structure and expression profile of mouse c1 inhibitor. *Eur J Biochem*, **254**, 117-122.
- Letunic, I., *et al.* (2006) Smart 5: Domains in the context of genomes and networks. *Nucleic Acids Res*, **34**, D257-260.
- Lev-Maor, G., *et al.* (2007) Rna-editing-mediated exon evolution. *Genome Biol*, **8**, R29.
- Levashina, E.A., *et al.* (1999) Constitutive activation of toll-mediated antifungal defense in serpin-deficient drosophila. *Science*, **285**, 1917-1919.
- Lewis, M.J., Sweet, D.J. and Pelham, H.R. (1990) The erd2 gene determines the specificity of the luminal er protein retention system. *Cell*, **61**, 1359-1363.
- Ligoxygakis, P., Roth, S. and Reichhart, J.M. (2003) A serpin regulates dorsal-ventral axis formation in the drosophila embryo. *Curr Biol*, **13**, 2097-2102.



- Loh, Y.H., Brenner, S. and Venkatesh, B. (2008) Investigation of loss and gain of introns in the compact genomes of pufferfishes (fugu and tetraodon). *Mol Biol Evol*, **25**, 526-535.
- Longas, M.O., Ferguson, W.S. and Finlay, T.H. (1980) A disulfide bond in antithrombin is required for heparin-accelerated thrombin inactivation. *J Biol Chem*, **255**, 3436-3441.
- Lynch, M. and Conery, J.S. (2000) The evolutionary fate and consequences of duplicate genes. *Science*, **290**, 1151-1155.
- Lynch, M. and Force, A. (2000) The probability of duplicate gene preservation by subfunctionalization. *Genetics*, **154**, 459-473.
- Maglott, D., Ostell, J., Pruitt, K.D. and Tatusova, T. (2005) Entrez gene: Gene-centered information at ncbi. *Nucleic Acids Res*, **33**, D54-58.
- Meyer, A. and Van de Peer, Y. (2005) From 2r to 3r: Evidence for a fish-specific genome duplication (fsgd). *Bioessays*, **27**, 937-945.
- Morgenstern, B. (1999) Dialign 2: Improvement of the segment-to-segment approach to multiple sequence alignment. *Bioinformatics*, **15**, 211-218.
- Morgenstern, B. (2000) A space-efficient algorithm for aligning large genomic sequences. *Bioinformatics*, **16**, 948-949.
- Morgenstern, B. (2004) Dialign: Multiple DNA and protein sequence alignment at bibiserv. *Nucleic Acids Res*, **32**, W33-36.
- Moriyama, S., et al. (2008) Gene structure and functional characterization of growth hormone in dogfish, *squalus acanthias*. *Zoolog Sci*, **25**, 604-613.
- Nagata, K. (1996) Hsp47: A collagen-specific molecular chaperone. *Trends Biochem Sci*, **21**, 22-26.
- Nei, M., Rogozin, I.B. and Piontkivska, H. (2000) Purifying selection and birth-and-death evolution in the ubiquitin gene family. *Proc Natl Acad Sci U S A*, **97**, 10866-10871.
- Nei, M. and Rooney, A.P. (2005) Concerted and birth-and-death evolution of multigene families. *Annu Rev Genet*, **39**, 121-152.
- Nelson, D.L. and Cox, M.M. (2005) *Lehninger principles of biochemistry*. W. H. Freeman and Company, New York.
- Nicholas, K.B., Nicholas H.B. Jr. and Deerfield, D.W.I. (1997) Genedoc: Analysis and visualization of genetic variation. *EMBNEW.NEWS*, **4**, 14.
- Notredame, C., Higgins, D.G. and Heringa, J. (2000) T-coffee: A novel method for fast and accurate multiple sequence alignment. *J Mol Biol*, **302**, 205-217.
- Ohno, S. (1970) *Evolution by gene duplication*. NY: Springer Verlag, New York.
- Ohno, S. (1999) Gene duplication and the uniqueness of vertebrate genomes circa 1970-1999. *Semin Cell Dev Biol*, **10**, 517-522.
- Oley, M., Letzel, M. and Ragg, H. (2004) Inhibition of furin by serpin spn4a from drosophila melanogaster. *FEBS Lett*, **577**, 165 - 169.
- Omilian, A.R., Scofield, D.G. and Lynch, M. (2008) Intron presence-absence polymorphisms in daphnia. *Mol Biol Evol*, **25**, 2129-2139.
- Osterwalder, T., et al. (1998) The axonally secreted serine proteinase inhibitor, neuroserpin, inhibits plasminogen activators and plasmin but not thrombin. *J Biol Chem*, **273**, 2312 - 2321.
- Osterwalder, T., et al. (2004) Drosophila serpin 4 functions as a neuroserpin-like inhibitor of subtilisin-like proprotein convertases. *J Neurosci*, **24**, 5482 - 5491.
- Page, R.D. (1996) Treeview: An application to display phylogenetic trees on personal computers. *Comput Appl Biosci*, **12**, 357-358.
- Pak, S.C., et al. (2004) Srp-2 is a cross-class inhibitor that participates in postembryonic development of the nematode *caenorhabditis elegans*: Initial characterization of the clade I serpins. *J Biol Chem*, **279**, 15448-15459.

- Palmer, J.D. and Logsdon, J.M., Jr. (1991) The recent origins of introns. *Curr Opin Genet Dev*, **1**, 470-477.
- Patston, P.A. and Schapira, M. (1997) Regulation of c1-inhibitor function by binding to type iv collagen and heparin. *Biochem Biophys Res Commun*, **230**, 597-601.
- Pearson, W.R. (1990) Rapid and sensitive sequence comparison with fastp and fasta. *Methods Enzymol*, **183**, 63-98.
- Pearson, W.R. and Lipman, D.J. (1988) Improved tools for biological sequence comparison. *Proc Natl Acad Sci U S A*, **85**, 2444-2448.
- Pelham, H.R. (1990) The retention signal for soluble proteins of the endoplasmic reticulum. *Trends Biochem Sci*, **15**, 483-486.
- Pelissier, P., *et al.* (2008) An original serpin3 gene cluster: Elucidation of genomic organization and gene expression in the bos taurus 21q24 region. *BMC Genomics*, **9**, 151.
- Perriere, G. and Gouy, M. (1996) Ww-query: An on-line retrieval system for biological sequence banks. *Biochimie*, **78**, 364-369.
- Ponting, C.P. (2008) The functional repertoires of metazoan genomes. *Nat Rev Genet*, **9**, 689-698.
- Poole, A., Jeffares, D. and Penny, D. (1999) Early evolution: Prokaryotes, the new kids on the block. *Bioessays*, **21**, 880-889.
- Poole, A.M., Jeffares, D.C. and Penny, D. (1998) The path from the rna world. *J Mol Evol*, **46**, 1-17.
- Potempa, J., Shieh, B.H. and Travis, J. (1988) Alpha-2-antiplasmin: A serpin with two separate but overlapping reactive sites. *Science*, **241**, 699-700.
- Pruitt, K.D., Tatusova, T. and Maglott, D.R. (2003) Ncbi reference sequence project: Update and current status. *Nucleic Acids Res*, **31**, 34-37.
- Pruitt, K.D., Tatusova, T. and Maglott, D.R. (2005) Ncbi reference sequence (refseq): A curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res*, **33**, D501-504.
- Putnam, N., *et al.* (2007) Sea anemone genome reveals ancestral eumetazoan gene repertoire and genomic organization. *Science*, **317**, 86 - 94.
- Putnam, N.H., *et al.* (2008) The amphioxus genome and the evolution of the chordate karyotype. *Nature*, **453**, 1064-1071.
- Qiu, W.G., Schisler, N. and Stoltzfus, A. (2004) The evolutionary gain of spliceosomal introns: Sequence and phase preferences. *Mol Biol Evol*, **21**, 1252-1263.
- Ragg, H. (2007) The role of serpins in the surveillance of the secretory pathway. *Cell Mol Life Sci*, **64**, 2763 - 2770.
- Ragg, H., *et al.* (2009) Multiple gains of spliceosomal introns in a superfamily of vertebrate protease inhibitor genes. *BMC Evol Biol*, **9**, 208.
- Ragg, H., *et al.* (2001) Vertebrate serpins: Construction of a conflict-free phylogeny by combining exon-intron and diagnostic site analyses. *Mol Biol Evol*, **18**, 577-584.
- Raible, F., *et al.* (2005) Vertebrate-type intron-rich genes in the marine annelid platynereis dumerilii. *Science*, **310**, 1325 - 1326.
- Raykhel, I., *et al.* (2007) A molecular specificity code for the three mammalian kdel receptors. *J Cell Biol*, **179**, 1193-1204.
- Reisz, R.R. and Muller, J. (2004) Molecular timescales and the fossil record: A paleontological perspective. *Trends Genet*, **20**, 237-241.
- Richer, M., *et al.* (2004) The spn4 gene of drosophila encodes a potent furin-directed secretory pathway serpin. *Proc Natl Acad Sci*, **101**, 10560 - 10565.

- Rogozin, I.B., *et al.* (2003) Remarkable interkingdom conservation of intron positions and massive, lineage-specific intron loss and gain in eukaryotic evolution. *Curr Biol*, **13**, 1512-1517.
- Rokas, A. and Holland, P.W. (2000) Rare genomic changes as a tool for phylogenetics. *Trends Ecol Evol*, **15**, 454-459.
- Roy, S.W. (2003) Recent evidence for the exon theory of genes. *Genetica*, **118**, 251-266.
- Roy, S.W. and Gilbert, W. (2006) The evolution of spliceosomal introns: Patterns, puzzles and progress. *Nat Rev Genet*, **7**, 211-221.
- Sadusky, T., Newman, A.J. and Dibb, N.J. (2004) Exon junction sequences as cryptic splice sites: Implications for intron origin. *Curr Biol*, **14**, 505-509.
- Saitou, N. and Nei, M. (1987) The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Mol Biol Evol*, **4**, 406-425.
- Sauk, J.J., Nikitakis, N. and Siavash, H. (2005) Hsp47 a novel collagen binding serpin chaperone, autoantigen and therapeutic target. *Front Biosci*, **10**, 107-118.
- Sawant, S., *et al.* (2004) Regulation of factors controlling angiogenesis in liver development: A role for pedf in the formation and maintenance of normal vasculature. *Biochem Biophys Res Commun*, **325**, 408-413.
- Schechter, I. and Berger, A. (1967) On the size of the active site in proteases. *Biochem. Biophys. Res. Commun.*, **27**, 157-162.
- Schioth, H.B., Haitina, T., Fridmanis, D. and Klovins, J. (2005) Unusual genomic structure: Melanocortin receptors in fugu. *Ann N Y Acad Sci*, **1040**, 460-463.
- Schmidt, E.E. and Davies, C.J. (2007) The origins of polypeptide domains. *Bioessays*, **29**, 262-270.
- Schneider, M., Tognolli, M. and Bairoch, A. (2004) The swiss-prot protein knowledgebase and expasy: Providing the plant community with high quality proteomic data and tools. *Plant Physiol Biochem*, **42**, 1013-1021.
- Schneider, T.D. and Stephens, R.M. (1990) Sequence logos: A new way to display consensus sequences. *Nucleic Acids Res*, **18**, 6097-6100.
- Schultz, J., Milpetz, F., Bork, P. and Ponting, C.P. (1998) Smart, a simple modular architecture research tool: Identification of signaling domains. *Proc Natl Acad Sci U S A*, **95**, 5857-5864.
- Semenza, J.C., Hardwick, K.G., Dean, N. and Pelham, H.R. (1990) Erd2, a yeast gene required for the receptor-mediated retrieval of luminal er proteins from the secretory pathway. *Cell*, **61**, 1349-1357.
- Shiu, S.H., *et al.* (2006) Role of positive selection in the retention of duplicate genes in mammalian genomes. *Proc Natl Acad Sci U S A*, **103**, 2232-2236.
- Silverman, G.A., *et al.* (2001) The serpins are an expanding superfamily of structurally similar but functionally diverse proteins. Evolution, mechanism of inhibition, novel functions, and a revised nomenclature. *J Biol Chem*, **276**, 33293-33296.
- Sneath, P.H. and Sokal, R.R. (1973) *Numerical taxonomy*. Freeman, San Fransico, CA.
- Sodergren, E., *et al.* (2006) The genome of the sea urchin *Strongylocentrotus purpuratus*. *Science*, **314**, 941-952.
- Sorek, R. (2007) The birth of new exons: Mechanisms and evolutionary consequences. *RNA*, **13**, 1603-1608.
- Steele, F.R., Chader, G.J., Johnson, L.V. and Tombran-Tink, J. (1993) Pigment epithelium-derived factor: Neurotrophic activity and identification as a member of the serine protease inhibitor gene family. *Proc Natl Acad Sci U S A*, **90**, 1526-1530.
- Stoesser, G., *et al.* (1997) The embl nucleotide sequence database. *Nucleic Acids Res*, **25**, 7-14.

- Sverdlov, A.V., Rogozin, I.B., Babenko, V.N. and Koonin, E.V. (2004) Reconstruction of ancestral protosplice sites. *Curr Biol*, **14**, 1505-1508.
- Tamura, K., Dudley, J., Nei, M. and Kumar, S. (2007) Mega4: Molecular evolutionary genetics analysis (mega) software version 4.0. *Mol Biol Evol*, **24**, 1596-1599.
- Tateno, Y., *et al.* (1998) DNA data bank of japan at work on genome sequence data. *Nucleic Acids Res*, **26**, 16-20.
- Thompson, J.D., *et al.* (1997) The clustal\_x windows interface: Flexible strategies for multiple sequence alignment aided by quality analysis tools. *Nucleic Acids Res*, **25**, 4876-4882.
- Thompson, J.D., Higgins, D.G. and Gibson, T.J. (1994) Clustal w: Improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Res*, **22**, 4673-4680.
- Thompson, J.D., Plewniak, F. and Poch, O. (1999) A comprehensive comparison of multiple sequence alignment programs. *Nucleic Acids Res*, **27**, 2682-2690.
- Tombran-Tink, J. (2005) The neuroprotective and angiogenesis inhibitory serpin, pedf: New insights into phylogeny, function, and signaling. *Front Biosci*, **10**, 2131-2149.
- Tombran-Tink, J., *et al.* (2005) Pedf and the serpins: Phylogeny, sequence conservation, and functional domains. *J Struct Biol*, **151**, 130-150.
- Tordai, H. and Patthy, L. (2004) Insertion of spliceosomal introns in proto-splice sites: The case of secretory signal peptides. *FEBS Lett*, **575**, 109-111.
- Vandepoele, K., *et al.* (2004) Major events in the genome evolution of vertebrates: Paraneome age and size differ considerably between ray-finned fishes and land vertebrates. *Proc Natl Acad Sci U S A*, **101**, 1638-1643.
- Venkatesh, B., Ning, Y. and Brenner, S. (1999) Late changes in spliceosomal introns define clades in vertebrate evolution. *Proc Natl Acad Sci U S A*, **96**, 10267-10271.
- Venter, J.C., *et al.* (2001) The sequence of the human genome. *Science*, **291**, 1304-1351.
- Volff, J.N. (2005) Genome evolution and biodiversity in teleost fish. *Heredity*, **94**, 280-294.
- Volff, J.N., Bouneau, L., Ozouf-Costaz, C. and Fischer, C. (2003) Diversity of retrotransposable elements in compact pufferfish genomes. *Trends Genet*, **19**, 674-678.
- Wallace, I.M., O'Sullivan, O., Higgins, D.G. and Notredame, C. (2006) M-coffee: Combining multiple sequence alignment methods with t-coffee. *Nucleic Acids Res*, **34**, 1692-1699.
- Wang, T. and Secombes, C.J. (2003) Complete sequencing and expression of three complement components, c1r, c4 and c1 inhibitor, of the classical activation pathway of the complement system in rainbow trout *oncorhynchus mykiss*. *Immunogenetics*, **55**, 615-628.
- Waterston, R.H., *et al.* (2002) Initial sequencing and comparative analysis of the mouse genome. *Nature*, **420**, 520-562.
- Wheeler, D.L., *et al.* (2006) Database resources of the national center for biotechnology information. *Nucleic Acids Res*, **34**, D173-180.
- Wheeler, D.L., *et al.* (2005) Database resources of the national center for biotechnology information. *Nucleic Acids Res*, **33**, D39-45.
- Whisstock, J.C., *et al.* (2000a) Conformational changes in serpins: Ii. The mechanism of activation of antithrombin by heparindagger. *J Mol Biol*, **301**, 1287-1305.
- Whisstock, J.C., Skinner, R., Carrell, R.W. and Lesk, A.M. (2000b) Conformational changes in serpins: I. The native and cleaved conformations of alpha(1)-antitrypsin. *J Mol Biol*, **295**, 651-665.
- Wu, C.H., *et al.* (2006) The universal protein resource (uniprot): An expanding universe of protein information. *Nucleic Acids Res*, **34**, D187-191.

- 
- Wu, S. and Mander, U. (1992) Fast text searching allowing errors. *Communications of the ACM*, **35**, 83 - 91.
- Zang, X. and Maizels, R.M. (2001) Serine proteinase inhibitors from nematodes and the arms race between host and pathogen. *Trends Biochem Sci*, **26**, 191-197.
- Zang, X., *et al.* (1999) A novel serpin expressed by blood-borne microfilariae of the parasitic nematode *Brugia malayi* inhibits human neutrophil serine proteinases. *Blood*, **94**, 1418-1428.
- Zhuo, D., Madden, R., Elela, S. and Chabot, B. (2007) Modern origin of numerous alternatively spliced human introns from tandem arrays. *Proc Natl Acad Sci USA*, **104**, 882 - 886.

## Appendix 8.1: Highly conserved residues present in > 70% of the serpins.

Residue <sup>1</sup>	%	Location	Comment
<b>Phe33</b>	79	middle of hA	shutter, packs against conserved position 54
<b>Asn49</b>	87	start of s6B	gate, extensive hydrogen bond network of C-terminal residues (389–393)
<b>Ser53*</b>	93	end of s6B	shutter, forms hydrogen bond of backbone of conserved positions 56 and 383
<b>Pro54*</b>	90	start of hB	shutter, forms tight turn
<b>Ser56**</b>	72	hB	shutter, makes hydrogen bond of side chain to conserved position 186
<b>Leu61</b>	75	hB	shutter, buried hydrophobic residue, packs against conserved positions 80, 184, 299, 303, and 312
<b>Gly67*</b>	80	end of hB	forms tight turn, packs against conserved position 130
<b>Thr72</b>	87	start of hC	makes hydrogen bonds to loop between hI and s5A
<b>Leu80</b>	75	end of hC	shutter, buried hydrophobic, packs against conserved position 61
<b>Phe130</b>	75	start of hE	packs against conserved position 67
<b>Phe147</b>	84	start of hF	packs into interface between hF and the A $\beta$ -sheet
<b>Ile157</b>	83	hF	shutter, packs into interface between hF and the A $\beta$ -sheet
<b>Asn158*</b>	94	hF	shutter, forms hydrogen bonds to loop joining hF to s3A
<b>Val161</b>	78	hF	shutter, packs into interface between hF and the A $\beta$ -sheet
<b>Thr165</b>	89	end of hF	shutter, inserts into A $\beta$ -sheet in $\delta$ conformation (Gooptu <i>et al.</i> , 2000)
<b>Gly167</b>	75	end of hF	shutter, inserts into A $\beta$ -sheet in $\delta$ conformation (Gooptu <i>et al.</i> , 2000)
<b>Ile169</b>	84	loop between hF/s3A	shutter, inserts into A $\beta$ -sheet in $\delta$ conformation (Gooptu <i>et al.</i> , 2000)
<b>Thr180</b>	75	loop between hF/s3A	hydrogen bonding stabilizes turn into s3A
<b>Leu184</b>	74	s3A	shutter, buried hydrophobic, packs against conserved position 61
<b>Asn186</b>	85	s3A	shutter, hydrogen bond to conserved position 334, Ser 56 and P8 of RCL in cleaved form (Whisstock <i>et al.</i> 2000a)
<b>Phe190</b>	95	s3A	breach, buried hydrophobic, packs against conserved position 244
<b>Lys191</b>	78	s3A	breach, makes salt bridge to Asp 341 and hydrogen bonds to uninserted RCL
<b>Gly192</b>	74	end of s3A	breach, mobile region where sheet swings open to accept RCL during loop insertion
<b>Trp194</b>	94	end of s3A	breach, buried hydrophobic, packs against conserved positions 198 and 244
<b>Phe198</b>	95	s4C	breach, buried hydrophobic, packs against conserved positions 194 and 221
<b>Thr203</b>	84	s4C	gate, hydrogen bonds to conserved position 342 (Whisstock <i>et al.</i> 2000b) (Whisstock <i>et al.</i> , 2000a)
<b>Phe208</b>	98	s4C	gate, buried hydrophobic, packs against conserved positions 218, 369, and 370
<b>Val218</b>	80	s3C	gate, buried hydrophobic, packs against conserved positions 208, 220, 289, and 391
<b>Met220*</b>	84	s3C	gate, buried hydrophobic, packs against conserved positions 218 and 289
<b>Met221</b>	86	s3C	breach/gate, buried hydrophobic, packs against conserved positions 289, 198, and 342
<b>Tyr244</b>	76	s2B	breach, packs against conserved positions 190 and 194. Makes hydrogen bonds to P14 of RCL in inserted form
<b>Leu254</b>	80	s3B	gate, buried hydrophobic, packs against s1C and conserved position 370
<b>Pro255</b>	93	s3B	gate, buried hydrophobic, packs against conserved position 370
<b>Pro289</b>	96	start of s6A	gate, buried hydrophobic, packs against conserved positions 208, 218, 220, and 370
<b>Lys290</b>	72	start of s6A	gate, makes salt bridge to conserved position 342
<b>Leu299</b>	79	start of hI	buried hydrophobic, packs against conserved positions 61, 303, and 334
<b>Leu303</b>	90	hI	buried hydrophobic, packs against conserved positions 299 and 61
<b>Gly307</b>	83	end of hI	forms tight turn at end of hI
<b>Phe312</b>	90	loop between hI/s5A	buried hydrophobic, packs underneath A $\beta$ -sheet and against conserved position 61
<b>Ala316</b>	80	loop between hI/s5A	buried hydrophobic, packs underneath A $\beta$ -sheet
<b>Leu327</b>	72	loop between hI/s5A	buried hydrophobic, packs underneath A $\beta$ -sheet
<b>His334*</b>	78	s5A	shutter, H-bond to conserved position 186 (Whisstock <i>et al.</i> , 2000b), packs against conserved position 299
<b>Glu342*</b>	91	top of s5A	breach, H-bond bond to conserved position 203, salt bridge to conserved position 290, packs against conserved position 221
<b>Gly344</b>	89	RCL	hinge region (breach when RCL inserted)
<b>Ala347</b>	79	RCL	hinge region (shutter when RCL inserted)
<b>Pro369*</b>	96	start of s4B	gate, forms tight turn, packs against conserved position 208
<b>Phe370*</b>	97	s4B	gate, buried hydrophobic packs against conserved positions 208, 254, 255, and 289
<b>Leu383*</b>	80	s5B	shutter, buried hydrophobic, forms $\beta$ -bulge in s5B, packs against conserved position 384
<b>Phe384</b>	94	s5B	shutter, buried hydrophobic, forms -bulge in s5B, packs against conserved positions 190 and 383
<b>Gly386*</b>	89	s5B	shutter
<b>Pro391*</b>	95	C terminus	gate, buried hydrophobic; packs against conserved positions 208 and 218

<sup>1</sup> Numbering of residue is based on mature region of human  $\alpha$ -antitrypsin.

## Appendix 8.2: GENEDOC Usage

The alignments were visualised in the better way using GENEDOC sequence alignment editor. The following steps were followed

- ❖ After sequences of interest were aligned with sequence aligning tools and the output were saved the file into \*.aln format.
  - ❖ The alignment file (\*.aln) were imported into the GENEDOC software sequence alignment editor with import command then used following parameters:
    - The layer 1 was used for no color alignment with crossing the button I.
    - Go to project configure then From the configure under project button was opened and following parameters were selected:
      - Font Settings:
        - Points 10
        - Normal
      - Gap Ind: Dash
      - Seq Loc Ind:
        - On
        - After Name
        - After Seq
      - Residues
        - Normal
      - Seq Blocked Sizing:
        - Fixed = 40 (40 characters per line)
      - Project type
        - Protein
      - DNA Ambiguity
        - Disabled
      - Make Backups –yes
      - Show `~`as `~` – no
      - Show Man Shade – yes
      - Show comments – yes
      - Consensus Line: No consensus.
      - Summary Disp:
        - Sum Cols/Inch 24:24
        - 2 Col 1/1
      - Pict File Adjust
        - Width = 0
        - Height = 0
        - Ascent = 0
      - Scoring:
        - Pair Wise.
      - Marker Line:
        - Enabled: No
      - Cons Gap Sys: No
      - Max NameLen: 10
    - Name Separation: “ ” (One Gap character)
      - Indicator Separator: “ ” (One Gap character)
- ❖ 51 residues for characteristic serpins [Appendix 8.1] were marked in black background color and white text colors.
- ❖ The RCL region was marked in red background color and white text color, and residues P1-P1' were marked in black background color and yellow text color.
- ❖ Introns were marked in gray background color and black text colors. On top, intron positions are indicated and novel intron positions are marked by \* notation following the intron position.
- ❖ Gene specific characters and other features (if required) were marked in other colors as specified.

---

### 8.3. Alignments.

Common notations for all alignments in this section are described here and will be followed constantly in all appendices in this section:

(a) Intron features:

The conserved intron positions [number on top]

Novel intron position insertion [number on top with !]

Novel intron position insertion non-conserved part [number on top with #]

Novel intron position loss [number on top with !!]

Novel intron insertion/loss possessing serpin sequence [grey box at the position]

(b) Sequence based features:

Signal peptide [green box],

RCL [red box] with P1-P1' [yellow font in red box]

serpin specific conserved 51 amino acid position [black boxes] as summarized in **Appendix 8.1**.

Furthermore, marked sequence features are explained in the individual alignments.



**Appendix 8.3.1: Protein sequence alignment of serpins from *C. intestinalis*.** This alignment depicts signal peptide (green) RCL region (red color), amino acid residues conserved above 70% (black shade) C-terminal ER-retention signals (blue), intron position of each gene is marked with number and grey color in corresponding sequence.

ALAT_HSA	1	-----EDPQGDAAQKTDTS	SHHDQDHPTFNK	<b>K</b>	ITPNL-----	30	
Ci-Spn-1	1	<b>MYTASAYGVFLCLAIYQVGAT</b>	<b>KHLQAEFDYGEYEDDANSWDP</b>	<b>RQ</b>		45	
Ci-Spn-2	1	-----MLLV	IACMLSAAFNGAVGEP	YAPTNAFREP	VAH-----	33	
Ci-Spn-3	1	-----MKLLICS	LLLLVIATGYCQ	RRWINHFTD	NQ-----	31	
Ci-Spn-4	1	---MFLKQVLVLCVFF	FMTSSAF	YMPMVRTHPPQMDMP	PAYCAEVV	42	
Ci-Spn-5	1	-----MRFIFLCFVLLV	SAGFNEAK	RTRVISKWRLTA	IAN----	35	
Ci-Spn-6	1	-----		MAFCKVAAAK	-----	10	
Ci-Spn-7	1	-----		MAFCKVAAAK	-----	10	
Ci-Spn-8	1	-----		MAFCKVAAAK	-----	10	
Ci-Spn-9	1	-----MQFLYAI	VMLVLDANAK	I	IDTSEHVEKLSEAN	33	
Ci-Spn-10A	1	-----MQFLYAI	VMLVLDANAK	I	IDTSEHVEKLSEAN	33	
Ci-Spn-10B	1	-----MQFLYAI	VMLVLDANAK	I	IDTSEHVEKLSEAN	33	
ALAT_HSA	31	-----		AEFAFS	<b>S</b> LYRQ <b>L</b> A-HQ	44	
Ci-Spn-1	46	ASITKIGKMDGLT	IDQVELPPFEKPAARLVN	NFAFKLLNE	IA-SD	89	
Ci-Spn-2	34	-----		ALYDFGMDMY	<b>N</b> Q <b>L</b> EP <b>S</b> W	50	
Ci-Spn-3	32	-----		NTFSGS	<b>S</b> LYW <b>A</b> IS-KE	45	
Ci-Spn-4	43	NATRVFSG	-----	FILNAATHANS	ATDEY	66	
Ci-Spn-5	36	-----		KLFAHR	LFMEVARTT	50	
Ci-Spn-6	11	-----		TDFALGLY	<b>K</b> EL <b>S</b> -QK	24	
Ci-Spn-7	11	-----		TDFALGLY	<b>K</b> EL <b>S</b> -QK	24	
Ci-Spn-8	11	-----		TDFALGLY	<b>K</b> EL <b>S</b> -QK	24	
Ci-Spn-9	34	-----		IEFTLNLY	<b>K</b> N <b>L</b> I-EG	47	
Ci-Spn-10A	34	-----		IEFTLNLY	<b>K</b> N <b>L</b> I-EG	47	
Ci-Spn-10B	34	-----		IEFTLNLY	<b>K</b> N <b>L</b> I-EG	47	
			56b	60a	67a	73b	83b
ALAT_HSA	45	SNSTN	I	I	I	I	I
		IFFSPV	S	I	A	F	A
		SI	A	F	A	M	L
		S	L	G	T	K	A
		D	H	D	E	I	L
		E	L	E	G	L	E
		---					
		NFNLT					
		85					
Ci-Spn-1	90	N-EDNVV	F	S	P	L	S
		I	F	T	S	L	A
		T	L	R	P	A	L
		N	G	T	S	L	E
		Q	L	N	D		
		---					
		VTGLD					
		127					
Ci-Spn-2	51	RPTEN	I	V	I	S	P
		M	S	M	Y	A	I
		L	S	I	L	L	P
		G	L	N	G	A	S
		H	D	Q	V	Y	N
		---					
		ALRMT					
		89					
Ci-Spn-3	46	KPNKN	V	L	F	S	P
		I	S	V	S	Q	T
		L	G	M	V	L	A
		G	A	M	G	N	T
		Y	D	E	I	T	R
		---					
		ALQMT					
		84					
Ci-Spn-4	67	VAERN	V	F	S	P	F
		G	A	A	N	V	V
		G	I	L	R	L	A
		S	A	G	R	T	R
		E	Q	F	D	G	L
		P					
		---					
		LFSSI					
		107					
Ci-Spn-5	51	PEQBN	F	F	I	S	P
		Y	A	V	S	A	G
		L	S	M	P	L	Y
		G	A	H	S	T	A
		R	E	I	M	D	T
		L	G	Y	T	Q	L
		S	T	S			
		---					
		95					
Ci-Spn-6	25	E-DGNL	F	F	S	P	Y
		S	I	S	T	A	L
		M	M	T	L	L	G
		S	K	E	K	T	R
		E	E	M	L	D	
		---					
		VLGLK					
		62					
Ci-Spn-7	25	G-DGNL	F	F	S	P	Y
		S	I	S	T	A	L
		M	M	T	L	L	G
		S	K	E	K	T	R
		E	E	M	L	D	
		---					
		VLGLK					
		62					
Ci-Spn-8	25	E-DGNL	F	F	S	P	Y
		S	I	S	T	A	L
		M	M	T	L	L	G
		S	K	E	K	T	R
		E	E	M	L	D	
		---					
		VLGLK					
		62					
Ci-Spn-9	48	DPMKN	V	M	F	S	P
		I	T	T	A	L	A
		I	A	H	L	G	A
		K	G	N	T	A	K
		Q	I	D	D		
		---					
		AFMFS					
		86					
Ci-Spn-10A	48	DPMKN	V	M	F	S	P
		I	S	A	A	L	A
		M	T	H	L	G	A
		K	G	K	T	A	K
		Q	I	D	D		
		---					
		AFMFS					
		86					
Ci-Spn-10B	48	DPMKN	V	M	F	S	P
		I	S	A	A	L	A
		M	T	H	L	G	A
		K	G	K	T	A	K
		Q	I	D	D		
		---					
		AFMFS					
		86					

		97b	106c		
ALAT_HSA	86	EIPEAQIHEGFQ	ELLRTLNOQDSQ	-----	LQITTGNGLEFLSE 122
Ci-Spn-1	128	TIRESDMNDMYDGIF	KKSSS	-----	YKIKQASRIYVDR 160
Ci-Spn-2	90	NLPRNGVDAESAMCS	KIFQINPN	-----	YDLTRANRIEFGDR 125
Ci-Spn-3	85	DLTPSRIHTLMRKTR	NNVVMRPNG	-----	QTVKLANSVFVIGS 121
Ci-Spn-4	108	LQHNDRFMRGFTQ	TLRSLILSVTSFPGTNP	GGQSDKDLLNSGVETSR 152	
Ci-Spn-5	96	NFNQAKVPRLYQKML	HQVHQKDHG	-----	FELTSVNRMFGEES 132
Ci-Spn-6	63	DLNESDINSGLFLQ	ILHHLRSSRGD	-----	VVLEMANKLFPEA 99
Ci-Spn-7	63	DLNESDINSGLFLQ	ILHHLRSSKGD	-----	VVLEMANKLFPEA 99
Ci-Spn-8	63	DLNESDINSGLFLQ	ILHHLRSSRGD	-----	VVLEMANKLFPEA 99
Ci-Spn-9	87	KIEDGRFHSAFGE	LHGLLFDKASEK	-----	VTAKSSNRVVFADK 124
Ci-Spn-10A	87	KIEDGRFHSAFGE	LHGLLFDKASDN	-----	VTVKSSNRVVFADK 124
Ci-Spn-10B	87	KIEDGRFHSAFGE	LHGLLFDKASDN	-----	VTVKSSNRVVFADK 124
		136b	144c	151c	156b
ALAT_HSA	123	GLKLVDKLELDV	KKLYHSEAF	TVNF	FG-DTEEAKKQINDYVEKGTQ 166
Ci-Spn-1	161	GIRLSRSYRTDLYR	MKISRARRLD	FERRAPEESRNKINKYVKKRTR 205	
Ci-Spn-2	126	TLTFKKS	YKNETSWHHKAAHKKVDF	OHYPNRARRKMRYVSKMTD 170	
Ci-Spn-3	122	NYPVVQOYIDLLRQ	NKSSVFPVNFH	-NSNAAANMINEWVSNMTE 165	
Ci-Spn-4	153	WLYLQTRFISDAR	NFYKAVVASVDFS	-DPELASSHINMWINARTQ 196	
Ci-Spn-5	133	RNIFVPSYVKGVEH	FYGAKLKKVDF	RRNPERARQEINTWVEEVIN 177	
Ci-Spn-6	100	IYKLEKDELSKCKE	FYETEIQALDFKGN	PDASREAINAWAEKETS 144	
Ci-Spn-7	100	TYKLEDELSKCKQ	FYETEIQALDFKGN	PDASREAINVWAEKETS 144	
Ci-Spn-8	100	TYKLEKDELSKCKE	FYETEIQALDFKGN	PDASREAINAWAEKETS 144	
Ci-Spn-9	125	HITVFEDYQD-SL	SVYSATVSVDFK	-MPKS	AVKKINDWSSDATN 167
Ci-Spn-10A	125	KRKVLEDYKN-AL	TVYGAKLENVDFK	-TFSNAVKQINDWASDATN 167	
Ci-Spn-10B	125	KRKVLEDYKN-AL	TVYGAKLENVDFK	-TFSNAVKQINDWASDATN 167	
				191c	204c
ALAT_HSA	167	GKIVDLVK--ELDR	DTVFALVNYIFK	GKWERPEEVKDINEEDEH 209	
Ci-Spn-1	206	KLIKELVPVGAISS	ATMMYLVAIYLKAK	WDIPFQKSLTRMRFR 250	
Ci-Spn-2	171	GEIQQLIPREAVTT	DRIFLVNAIAFKAA	WKSSEIKDATTLTNEH 215	
Ci-Spn-3	166	DKIRELVDPSSITA	FTRMILVNAVYFQAD	WAISEFRRIPTKQN-EF 209	
Ci-Spn-4	197	RKITKIVSPSDLSP	TTLVTVENTLFEAL	WKHPETTGRISNSTEV 241	
Ci-Spn-5	178	GTIREALPPNSVTA	EULLVDMSTLYFK	GLWEKPF EIN--LRSTFY 220	
Ci-Spn-6	145	GKIKDLLPSGSID	SLVRLVLNAVYFKG	SWLHKFKQQTIMKDFH 189	
Ci-Spn-7	145	GKINDLLPNGSIN	SLVRLVLNAVYFKG	SWLHKFKDYDSIESNEH 189	
Ci-Spn-8	145	GKIKDLLPSGSID	SLVRLVLNAVYFKG	SWLHKFKQQTIMKDFH 189	
Ci-Spn-9	168	GVIKSMLEEDGVN	NDTALLIINALYFR	GNWDEYEFDEGRITKRRPFY 212	
Ci-Spn-10A	168	GKISNMLQDDAVD	SNTALIVNAVYFRG	DWHSKENEMQTERRAFY 212	
Ci-Spn-10B	168	GKISNMLQDDAVD	SNTALIVNAVYFRG	DWHSKENEMQTERRAFY 212	

			217c		225a				
ALAT_HSA	210	VDQVTTVKVPMMKRL	-----	GMFNIQHCKKLS	SSWVLLM		242		
Ci-Spn-1	251	VSNNESIRVETMISK	-----	NTFC	TRVNNRDLQASVTV		283		
Ci-Spn-2	216	VSPTKVKQAATMYTSS	-----	AVCF	HQSHDAQLES	DLIVL	250		
Ci-Spn-3	210	LSNGTTVQVPEFMVRW	-----	EAVV	KSYNYRDKIEFFFI		242		
Ci-Spn-4	242	LANGTPVLTPEMMEVTANHFLHYSGEFCQLF	SMR	RCHP	NTPDIVIL		286		
Ci-Spn-5	221	TTNNEQYQTDQVQQT	-----	MFAL	HSFSEQFOAHIVEL		253		
Ci-Spn-6	190	IRENKVEKVMFMFK	-----	RKFR	FNFDQSLGLQVVEI		222		
Ci-Spn-7	190	VKEGTTTQVKMMNQK	-----	EWEN	FKTDPDLGLKIAEL		222		
Ci-Spn-8	190	IRENKVEKVMFMFE	-----	RKFR	FNFDQSLGLQVVEI		222		
Ci-Spn-9	213	VSKDKAVETSEMFQN	-----	EHEF	KYAYINELTLQVLEM		245		
Ci-Spn-10A	213	VSHYKIVETPEMFQR	-----	GHEF	KYAYISELTLQVLEM		245		
Ci-Spn-10B	213	VSHYKIVETPEMFQR	-----	GHEF	KYAYISELTLQVLEM		245		
			246a	247a	259a	274a			
ALAT_HSA	243	KYLC	-----	NATAIFFLP	---DEGKI	QHLENELTHDIIITKFL	ENE	279	
Ci-Spn-1	284	LSLGG	-----	SFSFVIMSPH	---SAGNF	SRFYDDGVTTMQEK	MTRAF	322	
Ci-Spn-2	251	PFKGA	-----	KTTMVFIVPI	---VAGNF	GPLKGAVGASKISQALDRY		289	
Ci-Spn-3	243	RYKTTSN	QNTYFV	VGLPG	---DNYML	QQFSREAAQILSRFRTINK		284	
Ci-Spn-4	287	PYKGE	-----	RRQMI	VLIEN	---QNTIREIERQFGTNLEKWRSSLV		325	
Ci-Spn-5	254	PFKTSSSRYKVMVQLILPE	SRGADN	LNLI	EDQF	DEENFD	FATEDQ	298	
Ci-Spn-6	223	PYIGN	-----	KLSMV	VFLPT	---ERFAL	NKIENALTTEKLHGLLAGL	261	
Ci-Spn-7	223	FYKGG	-----	DYSMV	VLLPD	---EKYGL	NKCLEKLTSEKLQHISSGM	261	
Ci-Spn-8	223	PYIGN	-----	KLSMV	VFLPT	---ERFAL	NKIENALTTEKLHGLLAGL	261	
Ci-Spn-9	246	DYAGT	-----	DYSMV	LMP	---ENFDL	LAKVEANLNHANLTKWLSAL	283	
Ci-Spn-10A	246	DYAGK	-----	DYSMV	LMP	---ENFDL	LAKVEANLNHANLTKWLSAL	283	
Ci-Spn-10B	246	DYAGK	-----	DYSMV	LMP	---ENFDL	LAKVEANLNHANLTKWLSAL	283	
			245b	283c	292c	301b			
ALAT_HSA	280	-----	DRRSASLHLPKLS	ITGTYD	---	LKSVL	QGLGITKVF	S	313
Ci-Spn-1	323	NKIWTRRGNRQQQLCSVKLPKFKVDYAEN	---	IKEVL	KGLGIRDIFS			366	
Ci-Spn-2	290	WTG-YRNMP	IPMRVCEVRMPKFK	I	THSVDDIMGAMRAMNVTDIFS			333	
Ci-Spn-3	285	-----	MFRITHFKMPLIELSHKTD	---	VKEVL	QTLGVVDLFD		318	
Ci-Spn-4	326	-----	DGNVELHHPKFE	LKSNLD	---	LKSVL	RSEGLTEPEN	358	
Ci-Spn-5	299	-----	ENISVTIRLPKFRLE	YETD	---	LKETL	YNMGIQSLFS	332	
Ci-Spn-6	262	-----	WEETLMLSLPRMKFE	QDFD	---	LGGVL	KKMGMDAFD	295	
Ci-Spn-7	262	-----	MRTELALSLPHMKFE	QLD	---	LVGSL	KKLGLVDLFN	295	
Ci-Spn-8	262	-----	WEETLMLCLPRMKFE	QDFQ	---	LGEVL	KKMGMDAFS	295	
Ci-Spn-9	284	-----	EHESVDLTIPKFKLE	ETLQ	---	LQEVLP	PKMGVTDLFD	317	
Ci-Spn-10A	284	-----	KYKSVDLSVPKFKLE	ETLQ	---	LQEVLP	PKMGVTDLFD	317	
Ci-Spn-10B	284	-----	KYKSVDLSVPKFKLE	ETLQ	---	LQEVLP	PKMGVTDLFD	317	

			332a	339c	344a	
ALAT_HSA	314	NGAD-LSG-VTEEAPLKLSKAVHKAVLTIDEK				GTEAAGAMFLEAI 356
Ci-Spn-1	367	INADFSRLSVRNNRELYVSEARHSVAVLSADEA				GVEAAGATAFGIS 411
Ci-Spn-2	334	TEADFSPM---TPELVYVTDMRHKAVIKVNEQ				GVKATAATSIGLT 375
Ci-Spn-3	319	SGASNLTG-IISTVEQLYVSEFTQKAYINVNENG				TVAAAASAATVQ 362
Ci-Spn-4	359	RTTADYST--MTRQLAISKLFQTASISMDET				GVRATSTTAFFE 401
Ci-Spn-5	333	RGEADLSG-IISTNGDLSLGSAAHKTFIQVDES				GTTAGASYAQQGF 376
Ci-Spn-6	296	ERAANFEA-ISGSRDLVISKVVHKAFIEVNEE				GSEAAAATAVVM 339
Ci-Spn-7	296	GNKSNLRG-ISDDGDLFVSQVAHKAFIEVNET				GTEAAAATAMIAM 339
Ci-Spn-8	296	KGAAANFEA-ISGSRDLVISKVVHKAFIEVNEE				GSEAAAATAVVVK 339
Ci-Spn-9	318	RQACDLTG-IANRNNLFVDQIVHKTVLDVNEQ				GSEAAATTSVRTQ 361
Ci-Spn-10A	318	RQACDLTG-IKSKDLNVDQIVHKTVLEVDEQ				GSEAAATTTVRIQ 361
Ci-Spn-10B	318	RQACDLTG-IKSKDLNVDQIVHKTVLEVEEN				GGAVPQERADANQ 361

ALAT_HSA	357	PMSIPPE-----VKFNKPFVFLMIEQNTKSPLEFMGKVVNPTQK--				394
Ci-Spn-1	412	LRSTSLQ-----VTVNKPFIFALRHDP SGALIFVVKIVRPSVG--				449
Ci-Spn-2	376	GRSLPIR-----VEINRPFMYMIRHEPTGALLFLGRVVDPTK---				412
Ci-Spn-3	363	GRSLSIPRQ---VTVDRPFFIGVYQEKSNSEFLFLGKVENPLEN--				402
Ci-Spn-4	402	LRSFFFRTR---INANKPFLFIIEDIHTRTPFLFLGRVTDPRPL--				441
Ci-Spn-5	377	RSVSDLD-----LVFNHPFIVIIIREKYTQMPMFMGRVAREPMY---				413
Ci-Spn-6	340	LRSMAPPVVM--VNC DHPFLFLIRHNQTKTILFLGRFSGP-----				377
Ci-Spn-7	340	QSMAMP SVPPVQFNCDHPFLFLIKHNPTNSVLFGLGRCS DPS----				380
Ci-Spn-8	340	ARSMPCLP EM--VNC DHPFLFLIQHNETKTILFLGRFSGPSI---				379
Ci-Spn-9	362	CDSVAFNPIS--FVADHPFLWAI RHRQSELLIFMGRFSRPEGPLL				404
Ci-Spn-10A	362	ARSLNSRPS---FVADHPFLWAI RHRQSELLIFMGRLSRPEGPLL				403
Ci-Spn-10B	362	TPALDRPV----VYVDHPFLIIIVRGRANNAFHFLFGAYKRPA G KIR				402

ALAT_HSA	-	-----	-
Ci-Spn-1	-	-----	-
Ci-Spn-2	-	-----	-
Ci-Spn-3	-	-----	-
Ci-Spn-4	-	-----	-
Ci-Spn-5	-	-----	-
Ci-Spn-6	-	-----	-
Ci-Spn-7	-	-----	-
Ci-Spn-8	-	-----	-
Ci-Spn-9	405	GHDEF	409
Ci-Spn-10A	404	DHDEF	408
Ci-Spn-10B	403	SHDEL	407

**Appendix 8.3.2: Protein sequence alignment of serpins from *B. floridae*.** This alignment depicts RCL region (red color), amino acid residues conserved above 70% (black shade), sequence indels (orange), and C-terminal ER-retention signals (blue).

ALAT_HSA	1	-----EDPQGDAAQKTDTS	SHHDQDHPTFNK	ITPNLAEFAFSLYRQLAH	43	
Bfl-Spn-1	1	-----MRS	STSQESTPLADIN	SEFALELYKALHK	29	
Bfl-Spn-2	1	-----MLCPHYWLLGLLVA	IATAENPFPQ	EPTSLADVNSEFALELYKALHK	46	
Bfl-Spn-3	1	-----XNPLKTLVA	AANGKFA	LDDLYKKLTS	24	
Bfl-Spn-4	-	-----	-----	-----	-	
Bfl-Spn-5	1	MGGKKKIKNRRKPPARLR	GRERTVSASTGV	TTSWEAKSSLSSLVEANS	SAFALGLFRRLCD	60
Bfl-Spn-6	1	-----MSEDEVAEANS	SAFALS	LYRQLSQ	23	
Bfl-Spn-7	1	-----MSEDDVAECNS	SAFALS	LYRQLSQ	23	
Bfl-Spn-8	1	-----MGQSAIEFALELY	KVLHK	18		
Bfl-Spn-9	1	MGGKKKIKNRRKPPARLR	GREKTVSASTGV	TTSWEAKSSLSSLVEANS	SAFALGLFRRLCD	60

ALAT_HSA	44	QNSTN	IFFSPVSTATAFAMLSL	GTKADTHDEILEG	INFNLTEIPEAQIHEGFQ	ELLRTL	103
Bfl-Spn-1	30	DHP-EN	IFFSPFSISTCLAMT	YLGARNDTAQ	QMSRVILRFHKM-D-ASD	DHMLFHDLLTQL	86
Bfl-Spn-2	47	DHP-EN	IFFSPFSISTCLAMT	YLGARNDTAQ	QIRQVILRFNKS-N-QT	DHDFRDLAQL	103
Bfl-Spn-3	25	QSD-GN	MVFSPLSISTALAMT	YLAAGKTA	EQMKTMHFDDL-S-EL	TLHKTFAKLTET	81
Bfl-Spn-4	1	-----MT	ISSSLRCCF	-----	FFDRYKD--GGI	22	
Bfl-Spn-5	61	STD-GN	IVFSPLSISAAMAMT	YICARGNTRY	QMERILRFHYFQN-ED	DLHSTFSAIEDVI	118
Bfl-Spn-6	24	RTD-GN	IFFSPYSISAALAMT	YMGARHTTAA	QMAEVLHLLT-----	EGDFHQAFSNLS--R	75
Bfl-Spn-7	24	RTD-GN	IFFSPYSISAALAMT	YMGARHTTAA	QMAEVLHLLT-----	EGDFHQAFSNLS--R	75
Bfl-Spn-8	19	DHP-EN	IFFSPFSISTCLAMT	YLGARNDTAQ	QMSRVILRFNKL-N-QT	DHDFRDLAQL	75
Bfl-Spn-9	61	STD-GN	IVFSPLSISAAMAMT	YICARGNTRY	QMERILRFHYFQN-ED	DLHSTFSAIEDVI	118

ALAT_HSA	104	NQPDS--	QLQLTTGNGLFL	SEGLKLVDFK	LEDVKKLYHSEAF	VNF-GDTEEAKKQ	IINDY	160	
Bfl-Spn-1	87	HHSDR--	PYILK	TANRLFGQNS	SFEFVQKFLAET	SRHYRAQLAP	VDFHG	NTEGARQTINSW	144
Bfl-Spn-2	104	HHSDR--	PYILK	TANRLFGQNS	SFTFVQKFLDET	SRHYGADLAP	VDFHG	DTEGARQTINSW	161
Bfl-Spn-3	82	STNMT--	SYTILSMANRLE	VQEDFDV	LQSYIDGMKQHY	GAEVGRVDF-G	DSKVASD	MINIW	138
Bfl-Spn-4	23	GEVSD--	KYTLQ	TANRLYGE	QYTSFLQDFL	DATNKNYGAELAA	VDFKGA	AEQVRGTINQ	80
Bfl-Spn-5	119	STSG	READYTFVQANR	LFQAGMS	FRHDFLMDT	SRHYHSSLAT	VEFS	SDEEM-ARLAINSW	177
Bfl-Spn-6	76	TMFG	NLKKHTLVEANK	LFQOG	MKLEDDFLSGT	SRYYNARMEK	VDFDEER-S	RSRINSW	134
Bfl-Spn-7	76	AMFG	NLKKHTLVEANK	LFQOG	MKLEDDFLSGT	SRYYNARMEK	VDFDEER-S	RSRINSW	134
Bfl-Spn-8	76	HHSDR--	PYILK	TANRLFGQNS	SFKFVQKFLDET	SRHYGADLAP	VDFHG	NSEGARQTINSW	133
Bfl-Spn-9	119	STSG	READYTFVQANR	LFQAGMS	FRHDFLMDT	SRHYHSSLAT	VEFS	SDEEM-ARLAINSW	177

173

ALAT_HSA	161	VEKGT	QGGKIVDLV--	KELDRD	TVFALVNYIF	FKGK	WERP	FEVKDTE	EEDEHVDQVTT	TKV	218							
Bfl-Spn-1	145	VEEQ	ENKIQD	LLAPG	TVTP	STML	VLVN	AIYFKG	SWESKF	ESRIRL	GTGFHISRDEK	VEV	204					
Bfl-Spn-2	162	VEEQ	DNKIQD	IMAP	GSVP	PETLL	VLVN	AIYFKG	WESQ	FYS	SDIMLR	PFHVSP	EEVQV	221				
Bfl-Spn-3	139	VEEK	TQKIQD	ISE	DLNDL	TRL	VLVN	AIYFK	AKW	NEFN	PFDD	DRPE	FRTEEDS	VDV	198			
Bfl-Spn-4	81	VEEQ	QKNKIK	DLIP	AGAV	DAM	TRL	VLVN	AIYFK	GNW	DEQF	DANM	TRDR	ENINN	NEKVKV	140		
Bfl-Spn-5	178	VAGR	TGGK	VKGV	IPQ	LLKPL	TKL	VLVN	AVYF	AGK	WRTE	FDP	QLIN	MADE	FFIGPER	AVKV	237	
Bfl-Spn-6	135	VSTQ	TKRK	INDL	PKD	VLN	ALTR	LVN	AVYF	KG	WQ	TQ	DPRE	TYDR	KEF	FASSGNH	VTT	194
Bfl-Spn-7	135	VSTQ	TKRK	INDL	PKD	VLN	ALTR	LVN	AVYF	KG	WQ	TQ	DPRE	TYDR	KEF	FASSGNH	VTT	194
Bfl-Spn-8	134	VEEQ	ENKIQD	IMAP	GSVP	PETLL	VLVN	AIYFKG	WESQ	FYS	SDIMLR	PFHV	NHEEK	VQV	193			
Bfl-Spn-9	178	VAGR	TGGK	VKGV	IPQ	LLKPL	TKL	VLVN	AVYF	AGK	WRTE	FDP	QLIN	MADE	FFIGPER	AVKV	237	



247  
|

ALAT_HSA	219	PMMKRLGMEFIQHCKKLSWVLLMKYLGK-ATAFFFLPDE-GKLOHLENEIETHDIITKFL	276
Bfl-Spn-1	205	PMMHQQCRFKLAYDEDLNCQILEMPYRGRKLSMNVVVLDPKMDLISAETSLITPDLRHRW	264
Bfl-Spn-2	222	PMMYQDGIKFLGRDDDLNCSL-----STQLYLRDALQG-KALKYGGSL-----LTRRN	267
Bfl-Spn-3	199	PMHRSGNHHILFDPEVGCVLELPPYKQKDLMLVIVPTEKEGLRQVEDKTIIMDTLRGWR	258
Bfl-Spn-4	141	KMMRKEANFNHYGVFEDLKCVRVLELPPYVEKELSMILIFLPDAAATLQTVS-----	188
Bfl-Spn-5	238	PIMQLSGEFNVTEDEPILDCAVVELPPYSNIEIVMDIVLNPQRDGLERLQGLTIRRALNRI	297
Bfl-Spn-6	195	PTMHQRGKFERMADLPNLRCRMLELPPYAGDELAMFVILPKQMFGLKDVAVLITSEALLD	254
Bfl-Spn-7	195	PTMHQRGKFERMADLPNLRCRMLELPPYAGDELAMFVILPKQMFGLKDVAVLITSEALLD	254
Bfl-Spn-8	194	PMMYQDGIKFLGRDDDLNCSLLEIPYKGRKLSMNVVVLDPDEIGGLKTITETSLITPEVLQ	253
Bfl-Spn-9	238	PIMQLSGEFNVTEDEPILDCAVVELPPYSNIEIVMDIVLNPQRNGLERLQGLTIRRVLNRI	297

ALAT_HSA	277	ENEDRRSAS---LHLPKLSITGTYDLKSVLGQLGITKVFSL-NGADLSCVVEEAPLKLKSKA	332
Bfl-Spn-1	265	KSMSEESTM---VQIPKFKVEQDFLLKKEKLAEKMGMDLFSMADADLSCITGSRDLHVSHV	321
Bfl-Spn-2	268	RRTTEHRR-----CFESGRNFRLEKLESEMGMDDLFG-TDADLSCMTGSRDLHVDAL	318
Bfl-Spn-3	259	NALNDTFSL---VYLPKFKLEYSVSLTEHLKQMGMEDLFDLADLSCITGSRDLHVSVQV	315
Bfl-Spn-4	189	SKMYSTKVN---LLLPKFKLEQDFGLGDTLKKMGMGEAFS-DAADFSGMSGKDLFISAV	244
Bfl-Spn-5	298	RRYLPLEG---SVLPPKFKLETEEFSLKAOQLTAMGMDLFSQNRADLSCMTGQPGMHVSDA	354
Bfl-Spn-6	255	RSKSLQEVRSLDVALPKFRLTHALSILKNOQLTALGMDLFSMETADLSCVITGEKGLHVSE	314
Bfl-Spn-7	255	RSKSLQEVRSLDVALPKFRLTHALSILKNOQLTALGMDLFSMETADLSCVITGEKGLHVSE	314
Bfl-Spn-8	254	KSMVKEDVG---ILMPKFKLEQDFGLSEKLESEMGMDDLFG-TDADLSCMTGSRDLHVDAL	309
Bfl-Spn-9	298	RRYLPLEG---SVLPPKFKLETEEFSLKAOQLTAMGMDLFSQNRADLSCMTGQPGMHVSDA	354

P1P1'  
||

ALAT_HSA	333	VHKAVLTIDEK <b>GTEAAGAMFLEAIPMSIP</b> -----PEVKFNKPFVFLMTEQNIKSPLEFMG	386
Bfl-Spn-1	322	VHKAFVEVNEE <b>GSEAAAAATAVNMMKRSL</b> ---DGE---MFFADHPFLFLIRDNDSNSVLFGLG	376
Bfl-Spn-2	319	VHKAFVEVNEK <b>GTEAAAAATAG-IIVLGG</b> --PAH---EFAADHPFLFLFKDNETNSILFMG	372
Bfl-Spn-3	316	VQKAFVEVNEK <b>GSEAAAAATGVVIRLMSG</b> NFWLETPTVVRADRPFLFLIRDNRSNSVLFGLG	375
Bfl-Spn-4	245	VHKAFVEVNEE <b>GTEAAAAATGVVMLCAL</b> ---DLEGPPEFVADHPFLFLIRDNRSNSVLFGLG	302
Bfl-Spn-5	355	LHKAVTEVSEE <b>GTGAAAPAAAITDRS</b> -----RRGFEEFRADHPFLFLIRDKRTGSVLFLG	409
Bfl-Spn-6	315	LHKAFVEVNEE <b>GSEAAAAATAVVMRGRSG</b> ----NFGRFEMVNRPFLEFLQHKPTGTILFLG	370
Bfl-Spn-7	315	LHKAFVEVNEE <b>GSEAAAAATAVVMRGRSG</b> ----NFGRFEMVNRPFLEFLQHKPTGTILFLG	370
Bfl-Spn-8	310	VHKAFVEVNEK <b>GTEAAAAATAG-LILLSG</b> --PTH---EFAADHPFLFLFKDNETNSILFMG	363
Bfl-Spn-9	355	LHKAVTEVSEE <b>GTGAAAPAAAITDRS</b> -----RRGFEEFRADHPFLFLIRDKRTGSVLFLG	409

ALAT_HSA	387	KVVNETQK-----	394
Bfl-Spn-1	377	RLVREEGHTT- <b>KDEL</b>	390
Bfl-Spn-2	373	RLVREEGTT <b>KDEL</b>	387
Bfl-Spn-3	376	RVADPTGGKE-----	385
Bfl-Spn-4	303	RMYKPEA-----	309
Bfl-Spn-5	410	RLVDPERN-----	416
Bfl-Spn-6	371	RVTNENE-----	377
Bfl-Spn-7	371	RVTNENE-----	377
Bfl-Spn-8	364	RLVREEGTT <b>KDEL</b>	378
Bfl-Spn-9	410	RLVDPERN-----	416

**Appendix 8.3.3: Protein sequence alignment of serpins from *S. purpuratus*.** This alignment depicts RCL region (red color), amino acid residues conserved above 70% (black shade), sequence indels (orange) and C-terminal ER-retention signals (blue).

```

ALAT_HSA      1  -----EDPQGDAAQKTDTSSHHDQDHPTFNKITTPNLAEFAFSLY  38
Spu-Spn-1    1  MARSKGMAFHITMVTLLTLGMMGSEVKAAEVSEHTQQLARANNAFALHLY  50
Spu-Spn-2    1  -----MAFYKQODISEQLVQLSSANTRFALDLY  28
Spu-Spn-3    1  -----MAFSKQODISGQLVQLSSANTRFALHLY  28
Spu-Spn-4    1  ---MFQFVQIIEYKYGQKYIRNMAFSKQODISGQLVQLSSANTRFALHLY  47
Spu-Spn-5    1  -----MAFSEQODISGQLVQLSSANIGFALDLY  28
Spu-Spn-6    1  -----MAFSKQODISGKLVQLSSANMGFALDLY  28
Spu-Spn-7    1  -----MAFSKQODISGQLMQLSSANIGFALDLY  28
Spu-Spn-8    1  -----MAFSKQODISGQLMQLSSANIGFALDLY  28
Spu-Spn-9    1  -----MAFSKQODISGQLVQLSSANIGFALDLY  28
Spu-Spn-10   -  -----

```

```

ALAT_HSA      39  RQLAHQSNSTNIFFSPVSIATAFAMLSLGTKADTHDEILEGLNFNLTETIP  88
Spu-Spn-1    51  SATRATHPDQNLFFSPLSVSTALGMTHLGARGTSSAQMSEVLRFNLLLEE-  99
Spu-Spn-2    29  QTEQDERRGINLFFSPLSISTALAMTQLGACGDTATQIADVFRNQVDQ-  77
Spu-Spn-3    29  QTLQDERRGINLFFSPLSISTALAMTQLGACGDTATQIADVFRNQVDQ-  77
Spu-Spn-4    48  QTLQDERKGKNLFFSPLSISTLLAMTQLGACGDTATQIADVFRNQVDQ-  96
Spu-Spn-5    29  QTLQDERRGINLFFSPLSISTALAMTQLGARGDTATQIADVFRNQVDQ-  77
Spu-Spn-6    29  QTLQDERRGINLFFSPLSISTALAMTQLGARGDTATQIADVFRNQVDQ-  77
Spu-Spn-7    29  QTLQDERRGKNLFFSPLILSTAMAMTQLGARGDTATQIADVFRNQVDQ-  77
Spu-Spn-8    29  QTLQDERRGKNLFFSPLSISTALAMTQLGARGDTATQIADVFRNQVDQ-  77
Spu-Spn-9    29  QTLQDERRGINLFFSPLSISTALAMTQEGARGDTATQIADVFRNQVDQ-  77
Spu-Spn-10   -  -----

```

```

ALAT_HSA      89  EAQTHEGFQELLRLTNQPDSQLQLTTGNGLFLSEGLKLVDKFLEDVKKLY  138
Spu-Spn-1    100 -EHLHASFKQLNALLYGSSNKYTLKSANKLFGKAGADFLQEFLDNTGNFY  148
Spu-Spn-2    78  -DQLHGTFKELNLLYQTNSGYKLHAANRLYGKSGYNFVQSFLEGTASY  126
Spu-Spn-3    78  -DQLHGTFKELKNLLYQTDSGYKLHAANRLYGKSGYNFVQSFLEGTASY  126
Spu-Spn-4    97  -DQLHGTFKELNLLYQTDSGYKLHAANRLYGKSGYNFVQSFLEGTATYY  145
Spu-Spn-5    78  -DQLHGTFKELNLLYQTDSGYKLHSANRLYGKSGYSFVQSFLEGTIVFYY  126
Spu-Spn-6    78  -DQLHGTFKELNKLLYQTNSGYKLHSANRLYGKSGYNFVQSFLEGTATYY  126
Spu-Spn-7    78  -DQLHGTFKELNLLYQTDSGYKLHSANRLYGKSGYNFVQSFLEESASY  126
Spu-Spn-8    78  -DQLHGTFKELNLLYQTDSGYKLHSANRLYGKSGYNFVQSFLEGTIVFYY  126
Spu-Spn-9    78  -DQLHGTFKELNLLYQTNYDSQIHCAYRLYGKSGYNFVQPLEGTASY  126
Spu-Spn-10   -  -----

```

173

```

ALAT_HSA      139 HSEAFTV-NFGDTEEAKKQINDYVEKGTQGKIVDLV--KELDRDTVFALV  185
Spu-Spn-1    149 NSAFEAVKDFA-APEARTMINDWVAKETEDKIQNLFPDGVLNSSLTQLVLV  197
Spu-Spn-2    127 GAAIEAVDNFA-APTATHSINDWVSKQTEDKITNLIPPGILNDLTRLVLV  175
Spu-Spn-3    127 GAAIEAVDNFA-APIATKSINDWVSKQTEDKITNLIAPGILNDLTRLVLV  175
Spu-Spn-4    146 GAAIEAVNNFA-APIATKSINDWVSKQTENKIKNLIAPGILNDLTRLVLV  194
Spu-Spn-5    127 GAAIEAVNDFA-SPIIATOSINDWVSKQTEGKIKNLIAPGILNDLTRLVLV  175
Spu-Spn-6    127 GAAIEAVNNFA-APTATOSINDWVSKQTEGKIKNLIAPGILNDLTRLVLV  175
Spu-Spn-7    127 GAAIEAVNNFA-SPIIVTOSINDWVSEQTEDKIKNLIAPGILNDLTRLVLV  175
Spu-Spn-8    127 GAAIEAVNDFA-SPIIATOSINDWVSKQTEGKIKNLIAPGILNDLTRLVLV  175
Spu-Spn-9    127 GAAIEVVNNNFA-APIIATOSINDWVSKQTEGKIKNLIAPGILNELTRLVLV  175
Spu-Spn-10   -  -----

```

```

ALAT_HSA      186  NYIFFKGNWERPFEVKDTEEE--DFHVDQVT-----TVKVPMMKRLGM 226
Spu-Spn-1    198  NATYFKSNWVKSFNRRDTEEG--VFRMKDET-----AHLPMMFKDGGE 237
Spu-Spn-2    176  NATYFKGNWESKFRAENTTQE--TFKVLDER-----KKVPVSLMIQK GK 217
Spu-Spn-3    176  NATYFKGNWESKFRAENTTQE--TFKVLDER-----KKVPVSLMIQK GK 217
Spu-Spn-4    195  NATYFKANWDSKFLAENTTKD--KFKVIGKR-----KKVRVSLMSQEDR 236
Spu-Spn-5    176  NATYFKANWKSIFYAENTTHDCTFKVFEDER-----YKVPVSLMSQKGR 219
Spu-Spn-6    176  NATYFKGNWESKFRAHNTTQE--TFKVFEDER-----KKVPVSLMSQEGS 217
Spu-Spn-7    176  NATYFKANWHSFRAHNTTKD--KFKVFEDERKKVFKREKVPVSLMSQGGP 223
Spu-Spn-8    176  NATYFKANWKSIFYAENTTHDCTFKVFEDER-----YKVPVSLMSQKGR 219
Spu-Spn-9    176  NATYFKANWKATFCAHNTTKD--KFKVFDMRKKVFKRKKVPVSLMSLKGK 223
Spu-Spn-10   1    -----MGVGLGVTN--QK GK 13

```

```

ALAT_HSA      227  FNIQHCKKLSWVLLMKYLGK--ATAIFFLP-DEGKLOHLENELTHDIT 273
Spu-Spn-1    238  VMITDDKERKCFVLEMPYDGED-LSMLAIALPWDDGLANVEEQLSMEVLD 286
Spu-Spn-2    218  FALAVDKTNDCLVLEMPYQGRN-LSLLIALPVKDDGLGNLQTKLSADILR 266
Spu-Spn-3    218  FALAVDNTNDCLVLEMPYQGRN-LSLLIALPVKDDGLGQLETKLSADILQ 266
Spu-Spn-4    237  FELVVDKTNDCLVLEMPYERFD-LSLLIALPVKDDGLGQLETKLSADILQ 285
Spu-Spn-5    220  FALTVDNTNDCLVLELPYESHN-LSLLIALPTKDDGLGQLETKLSVDVLQ 268
Spu-Spn-6    218  FALAVDNTNDCLVLELPYQGHN-LGLLIALPVKDDGLGQLETKLSADVLQ 266
Spu-Spn-7    224  FALAVDKTNDCLVLEMPYQGHN-MSLLIALPVKDDGLGQLETKLSVDVLQ 272
Spu-Spn-8    220  FALTVDNTNDCLVLELPYESHN-LSLLIALPTKDDGLGQLETKLSVDVLQ 268
Spu-Spn-9    224  FDLAVDKTNDCLVLELPYEGHN-LSLLIALPTKDDGLGQLETKLSVDVLQ 272
Spu-Spn-10   14  FSLAVDKTNDCLVLEMPYQGRNYLSLLIALPTKDDGLGQLETKLSADILQ 63

```

```

ALAT_HSA      274  KFLENEDRRSASLHLPKLSITGTYDLKSVLGGQLGITKVFNSNG-ADLSGVT 322
Spu-Spn-1    287  FWDSDLEPENAMYWVPRFKLEDTFSLSSILQSMGMADAFDATKADFSGMT 336
Spu-Spn-2    267  SWDAGLKSROVNVLLPKFKLEAEFQFQKEVLRMGMPDAFNDGKANFEGIS 316
Spu-Spn-3    267  SWDAGLKSROVNVLLPKFKLEAQFQFQKEFLORMGMSDAFDEDNANFEGIS 316
Spu-Spn-4    286  SWDAGLEWRKVKVLLPKFKLEAEFQFQKEVLRMGMPDAFNDGLANFEGIS 335
Spu-Spn-5    269  SWDVGLKSRRVNVLLPKFKLEATFQFQKEVLRMGMPDAFDEDNANFKGIS 318
Spu-Spn-6    267  SWDAGLKS-----EVLKRMGMPDAFDEDNANFEGIL 299
Spu-Spn-7    273  SWDAGLEWREVDVLLPKFKLEATYQFQKEVLRMGMPDAFDEDNANFKGIS 322
Spu-Spn-8    269  SWDVGLKSRRVNVLLPKFKLEATFQFQKEVLRMGMPDAFDEDNANFKGIS 318
Spu-Spn-9    273  SWDAGLEWREVDVLLPKFKLEATFQFQKEVLRMGMPDAFDEDNANFEGIS 322
Spu-Spn-10   64  SWDTGLKSREVDVLLPKFKLEATFQFQKEVLRMGMPDAFDEDNANFEGIS 113

```

```

ALAT_HSA      323  -EEAPLKLSKAVHKAVLITIDEKGTEAAGAMFLEAIPMSIPPE----- 363
Spu-Spn-1    337  -GDQSLHISEV IHKAFVNEEGTEAAAATGVTMTKRSISKR-----Y 378
Spu-Spn-2    317  -GDRELYISAV IHKAFVDINEEGSEAAAATAVAVKLGCARPR--EPEKPT 363
Spu-Spn-3    317  -GDRELHISAV IHKAFVDVNEEGSEAAAATAVVMRRRCAPPR--EPEKPI 363
Spu-Spn-4    336  -DDRELHISAV IHKAFVSVNEEGSEAAAATAVVMGNCLSMRRGKP-EKP I 383
Spu-Spn-5    319  SEEREFYISAV IHKAFVNEEGSEAAAATAVVMAGGPPRIER---EKP I 365
Spu-Spn-6    300  -GDRELHISAV IHKAFVDVNEEGSEAAAATAVVMAFGCSLPRER--EKP I 346
Spu-Spn-7    323  -DERELHISAV IHKAFVNEEDGTEAAAATGRLLGGCLPEDR---EKP I 367
Spu-Spn-8    319  SEEREFYISAV IHKAFVNEEGSEAAAATAVVMAGGPPRIER---EKP I 365
Spu-Spn-9    323  -GDFHISAV IHKAFVDINEEGSEAAAATAFGMNAMMSLCPPGEREKPI 371
Spu-Spn-10   114 -GDRELHISAV IHKAFVDVNEEGSEAAAATAVEMRRRWGCAPPSEREKPI 162

```



---

```

ALAT_HSA      364 -VKFNKPFVFMIEQNTKSPILFMGKVVNPTQK----- 394
Spu-Spn-1    379 RLRFDHPFLELIRDRRTKAVLFLGRLVDEPHDTRVNHEEEL 418
Spu-Spn-2    364 LFRADHPFLEMIRHRPTKSVLFMGRMMDES----- 393
Spu-Spn-3    364 LFRADHPFLEMIRHRPTKSVLFMGRMMDES----- 393
Spu-Spn-4    384 LFRADHPFLEMIRHQATKSVLFMGRMMDES----- 413
Spu-Spn-5    366 LFRADHPFLEMIRHRSTKSVLFMGRMMDES----- 395
Spu-Spn-6    347 LFRADHPFLEMIRHRSTKSVLFMGRMMDES----- 376
Spu-Spn-7    368 LFRADHPFLEMIRHRSTKAVLFMGRMMDES----- 397
Spu-Spn-8    366 LFRADHPFLEMIRHRSTKSVLFMGRMMDES----- 395
Spu-Spn-9    372 LFRADHPFLEMIRHRPTKSVLFMGRMMDESYILTPMSTL- 410
Spu-Spn-10  163 LFRADKRFPFPPREARNFWCFLLLYRKEDDSYHF----- 196

```

**Appendix 8.3.4: Protein sequence alignment of serpins from *N. vectensis*.** RCL region is marked in red color, amino acid residues conserved above 70% is marked by black shade.

AlAT_HSA	1	EDPQGDAAQKTDTS	SHHDQDHP	TFNKITPNLAEFAFS	SYRQLAHQSNSTNI	50
Nve-Spn-1	1	-----	-----	MASAVVHGSNDFALRI	LQSL--KSSASNV	28
Nve-Spn-2	1	-----	-----	MIISVCYQGK	-----GDAAKNT	17
Nve-Spn-3	1	-----	-----	MAQSPAGLSTNAFALDI	HRVLTAQDQGQTNL	30
AlAT_HSA	51	FFSPVSIATAFAMLSL	GPKADTHDE	IIEGLNFNLT	EIPAQIHEGFQELL	100
Nve-Spn-1	29	FFSPLSMSMALGLVYL	GSRGTTAIQIANIF	GWK--ESEFEET	HRTFKQFH	76
Nve-Spn-2	18	FYSPASTCVALGMVYA	GARGETADEMATA	MHWEGHKPMLP	SKHQEHKELS	67
Nve-Spn-3	31	FYSPASIVVALAMTYL	GARGNTATQMTKTF	HFP---TDVPEK	FHDFLQAL	77
AlAT_HSA	101	RTLNQPDS-QLQLTT	GNGLFLSEGLKLV	DKFL	EDVKKLVHSEAF	TVNFG-
Nve-Spn-1	77	EALLTSDLGYGEIQ	LVNKLWGHDFE	IL	EEFLHGTREFYH	SEMAQVDFVN
Nve-Spn-2	68	VALNPGA-TNEMSI	ANNLELQKDFS	SILKEFTD	ICQKYDADISL	VDYKT
Nve-Spn-3	78	NASNSDG---NQIL	MANRLEFAQMGFE	IL	EEFKKASKE	SFSASMALVDYVK
AlAT_HSA	149	DTEEAKKQINDYVE	KGTQKIVDLVKE	--LDRD	IVFALVNYIF	FKGKWER
Nve-Spn-1	127	KAFDARKEVNAWVH	QQTKGNTKELIP	HGVINSL	TRILLVNAVY	FKGVWKK
Nve-Spn-2	117	DFEGARKHVNQWVE	ERTKKKICDLTAP	GVFNM	LTRLVNAVY	FKGMWDK
Nve-Spn-3	125	NSNGARDTVNRWVE	QKTKDKIKNL	IPEGM	FNKDTILL	CLNAVYFKGSMK
AlAT_HSA	197	PFEVKDTEEEDEH	VDQVTT--VKVP	MMKRLGMENIQ	HCKKLS	SWVLLMKY
Nve-Spn-1	177	EFGEENTFHAAFF	VPESSHE	SKIEVEM	TRKMKVNFY	YDADIKCRVVELPY
Nve-Spn-2	167	PFKKEHSHSSEFR	TSSNE--VEVEM	FQKSKFKYLH	SDKYKCKL	LELPY
Nve-Spn-3	175	HENRNATQSGKEK	TTPSOE--IQVQ	FMYQSSE	FYLESS	TLCQIVELPY
AlAT_HSA	245	LG-NATAIFFFLP	DEGK-LQHL	ENELTHD	IITKFL	ENEDR---RSASLHLP
Nve-Spn-1	227	SGDDTAMVILPE	EPSSGIF	SILEKS	IDVEIME	KWRRMLN---TTVEVSIP
Nve-Spn-2	215	VDTQLSMVLVLP	DETEGLARFE	QDLTHDK	MTDIFNSV	SSQRPADVEVYIP
Nve-Spn-3	223	AGEKLSMVVLLP	NEVDGLGKLES	SLNKET	LQEA	MTSLRNSHP
Nve-Spn-3	272	EEVEVTLP				
AlAT_HSA	290	KLSITGTYDLKSVL	GQLGITKVF	SNG-ADL	SGVTE	EAP-DKLSKAVHKAV
Nve-Spn-1	274	KFRLSQKLELRSL	LQDLGVSDF	DSRKADL	SGISA	AKG-LVVS
Nve-Spn-2	265	KFKMTSEFKLINE	ALQELGMKMF	DQAAAD	FTGISL	PEHLEFVSAVLHKAF
Nve-Spn-3	273	KFTLTQEFSLGET	LKGMGASDL	FSPGKADL	SGISA	AAP--LVVSEVVHKAF
AlAT_HSA	338	LTIDEK	GTEAAGAMFLEA	IPMSIPPE	---VKF-NK	PFVFLMIEQNTKSPL
Nve-Spn-1	323	IEVNER	GTVA AATTGV	VMAKRSLDM	NE--VFYAD	HPFLFSIHKPSSAIL
Nve-Spn-2	315	VEVNEE	GTEAAAATAA	IMMRC	AIMREPLV	FRADHPFLFLIQHCKSKCVL
Nve-Spn-3	321	VEVNEE	GTIAAAATG	VGIMLSM	MPMPN	P--VFYANHPFLFLIRHNDTGAVL
AlAT_HSA	384	FMGKVVNETQK	-----			394
Nve-Spn-1	371	FLGKVMQPTRV	GKVS	PHSDKPL	SDEL	397
Nve-Spn-2	365	FMGRVMNPVE	-----			374
Nve-Spn-3	369	FMGRLVVPDKDN	-----			380

**Appendix 8.3.5: Alignment of MNEI from vertebrates.** Gene specific features include an inhibitory RCL (red box). Conserved intron positions are indicated above the alignment. Group V1 specific sequence indels and intron indels are marked by \* and #, respectively.

```

ALAT_HSA 1 EDPQGDAAQKTDTSHHDDQDHPFTFNKITPNLAEEAFSLYRQLAHQSNSTNIFFPSPVSIATA 60
MNEI_HSA 1 -----MEQLSSANTRFALDLFLALSENNPAGNIFISPFSSISSA 38
MNEI_MNU 1 -----MEQLSSANTLFALELFQTLNESSTPTGNIFFPSPSSISSA 38
MNEI_RNO 1 -----MEQLSSANSLEFALELFHTLSESSPTGNIFFPSPSSISSA 38
MNEI_GGA 1 -----MESLSNANSRFALDLFRKVNENPNPSGNIFFPSPSSISSA 38
MNEI_XTR 1 -----MENLSSACTHFSFDLFRKINENNATGNVFPSPSSISSA 38
MNEI_FRU 1 -----MAAISGSNTAFALDELLRRLSQQNPSGNIFVSPSSISSA 38
MNEI_TNI 1 -----MAAISSSNTAFALDLRLSQQNPSGNIFMSPSSISSA 38
MNEI_DRE 1 -----MEGVSRRANSLFALDLYRALSSASSAEGNIFFPSPSSISAA 38

```

78c#

```

ALAT_HSA 61 FAMLSLGTKADTHDEILEGLINFNLTEIPEAQIHEGFQELLRTLNQPDSQ-LQLTTGNGLF 119
MNEI_HSA 39 MAMVFLGTRGNTAAQLSKTFHFNTVEE----VHSRFQSLNADINKRGAS-YILKLANRLY 93
MNEI_MNU 39 LAMVILGAKGSTAAQLSKTFHFFDSVED----IHSRFQSLNAEVSKRGAS-HTLKLANRLY 93
MNEI_RNO 39 LAMVFLGTKGTTAAQLSKTFHFFDSVED----VHSRFQSLNAEVSKRGAS-HTLKLANRLY 93
MNEI_GGA 39 LAMVLLGSRGNTETQVLKTFHFDEVEN----IHSRFRALTADINRRDSS-CLLRIANRLY 93
MNEI_XTR 39 LAMVLLGARGNTAAQISRILHFDAVKD----LHSNFQTLNAEINKKNVSSYALNLANRLF 94
MNEI_FRU 39 LAMVYLGAKGETAAQMAQALSFSSGKD----VHADFQTLNGEINSPSAS-YTLKLANRLY 93
MNEI_TNI 39 LAIVYLGAKGDIAAQMAQALSFNSGHD----VHADFQTLNGEINSPSAS-YILRLANRLY 93
MNEI_DRE 39 LSMVYLGARGDTAGEMEKVLSFSSVSD----VHSHFESLISSINSPSAS-YILRLANRLY 93

```

128c

167 a#

\*\*

```

ALAT_HSA 120 LSEGLKLVDKFLEDVKKLYHSEAFTVNFG-DTEEAKKQINDYVEKGTQKIVDLVK--EL 176
MNEI_HSA 94 GEKTYNFLPEFLVSTQKYGADLASVDFQHASEDARKTINQWVKGQEGKIPELLSAGMV 153
MNEI_MNU 94 GEKTYNFLPEFLVSTQKMYGADLAPVDFQHASEDARKEINQWVKGQEGKIPELLSVGVV 153
MNEI_RNO 94 GEKTYNFLPEFLTSTQKMYGADLAPVDFQHASEDARKEINQWVKGQEGKIPELLAVGVV 153
MNEI_GGA 94 GEKSYSFLLEFLTNTQKLYGADLAAVDFLHAYGEARKEINQWVEEKEEGKIPDLLSEGSV 153
MNEI_XTR 95 GEKSFKFLPDFLSSVKKQYNADLGTVDFISAAEDARKEINTWVSEQIKGKIPEVLSAGAV 154
MNEI_FRU 94 GESTANFLSEFLDATQKYYHADLKAIDFIGATEECRABEINSWVEEQENKIKDLLKPGTV 153
MNEI_TNI 94 GETTSNFLSEFLKATQKHYHADLRAVDFIGAPEECRABEINTWVEQQENKIKDVLKPGSV 153
MNEI_DRE 94 GEKTFSFLPEYLSSSLNLYHADLQAVDFIGASEQSRQLNKWVEEQENKIRDLLKPGMV 153

```

212c

```

ALAT_HSA 177 DRDITVFALVNYIFFKGKWERPEFVKDTEEEDEHVDQVTTVKVPMMKRLGMFNIQHCKKLS 236
MNEI_HSA 154 DNMTKLVLVNAIYFKGNWKDKEMKEAITNAPERLNKKDRKTVKMYQKKFAYGYIEDLK 213
MNEI_MNU 154 DSMTKLVLVNAIYFKGMWEEKFMTEDITDAPERLSKKDTKTVKMYQKKFPFGYISDLK 213
MNEI_RNO 154 DSMTKLVLVNAIYFKGMWEEKFMQDITDAPERLNKKNTKSVKMYQKKFFFGYISDLK 213
MNEI_GGA 154 NSMTKLVLVNAIYFKGNWAEKFEEANTADMPERLNKNERKTVKMYQKKFRFGYISDMK 213
MNEI_XTR 155 NSFTKLVLVNAIYFKGDWAKKEKAEHTKDMPEFQLNKKEQKTVKMYQMEKLPFNYIPEIN 214
MNEI_FRU 154 STMTRLALVNAIYFKGNWMHREANTIKEMPEKVNQNESKPVQMYQMKKLPNYNIPEHG 213
MNEI_TNI 154 NTMTRLALVNAIYFKGNWMHPENEAFTIKEMPEKINQNESKPVQMYQMKKLPNYNIPDHS 213
MNEI_DRE 154 TGMTRLALVNAIYFKGNWLORENAODIKEMPEKINQKENRPVQMYQKKFPFNYIDHR 213

```

262c  
|

\*

```

ALAT_HSA 237 SWVLLMKYLGN-ATAIFFLPPDEGK-----LQHLENELTHDIITKFLENEDRRS---ASLH 287
MNEI_HSA 214 CRVLELPYQGEELSMVILLPPDDIEDESTGLKKIEEQITLEKLEHWTKPENLDF-IEVNVS 272
MNEI_MNU 214 CKVLEMPYQGGELSMVILLPKDIEDESTGLKKIEKQITLEKLEWTKRENLEF-IDVHVK 272
MNEI_RNO 214 CKVLEMPYQGGELSMVILLPEDIEDESTGLKKIEEQITLEKLEWTKRENLEN-IDVHVK 272
MNEI_GGA 214 TRVLELPYDEREFSMIILLPPDDIEDDSTGLQKLEQQITLEKLEWTRPEHLYS-TDVHVH 272
MNEI_XTR 215 CRVLELPYVDYELSMVIVLPPDNINDDTTGLQQLEKELSLEKINEWT--ENMMP-IDVHVH 271
MNEI_FRU 214 VQILELPYVEEELSMFILLPEETTDGSPLLKLENELTREKLEWNTRENMDVHSEVLVH 273
MNEI_TNI 214 LQILELPYAQEELSMFILLPEETTDGSPNLLKLENELTREKLEWNTRENMDVSSEVRVH 273
MNEI_DRE 214 VQVLELPYVKEELSMLILLPEETTDGSDPLLKLESELTIDKLEHWNTNRNMDTQTDIVH 273

ALAT_HSA 288 LPKLSITGTIDLKSVLGGQLGITKVFNSNG-ADLSGVTEEAPIKLSKAVHKAVLTIDEKGTE 346
MNEI_HSA 273 LPRFKLEESYTLNSDLARLGVQDLFNSSKADLSGMSGARDIFISKIVHKSFVEVNEEGTE 332
MNEI_MNU 273 LPRFKIEESYTLNSNLGRLGVQDLFSSSKADLSGMSGSRDLFISKIVHKSFVEVNEEGTE 332
MNEI_RNO 273 LPRFKIEESYILNSNLGRLGLQDLFNSSKADLSGMSGSRDLFISKIVHKAFVEVNEEGTE 332
MNEI_GGA 273 LPKFKLEESYDLKSDLSAMGLLDIFDSAKADLSGMSGAHDLFISKIVHKAFVEVNEEGTE 332
MNEI_XTR 272 LPKFKLEDSYKLSQIAGMCMADLFEAGSADLSGMSGSDIYLSEVIHKSFVEVNEEGTE 331
MNEI_FRU 274 LPKFKLEEDYEMNEALAKLGMTDVFCAAKADLSGMNGDGGDLFLSTVAHKAFVEVNEEGTE 333
MNEI_TNI 274 LPKFKLEENYEMKEALAKLGMTDVFCAKADLSGMNSDGGDLFLSTVAHKAFVEVNEEGTE 333
MNEI_DRE 274 LPRFKLEIESSLVEILMGMGMSVVFQEGKADLTGMTGHGGDLFLSAVAHKAFVDVNEEGTE 333

ALAT_HSA 347 AAGAMFLEAIPMSIPPEVKFN--KPFVFLMIEQNTKSPLFMGKVVNPTQK 394
MNEI_HSA 333 AAAATAGIATFCMLMPEENFTADHPFLFFIRHNSSGSILFLGRFSSP--- 379
MNEI_MNU 333 AAAATGGIATFCMLLPEEEFTVDHPFIFFIRHNPTSNVFLGRVCSP--- 379
MNEI_RNO 333 AAAATAGIATFCMLLPEEEFTADHPFIFFIRHNPTANVFLGRVCSP--- 379
MNEI_GGA 333 AAAATAGIAMLCMVI-EEDFNADHPFLFFLRHNPTKSIVFFGRYASP--- 378
MNEI_XTR 332 AAAASAGIAMMCLMR-EEEFNANHPFLFFIRHNATKSIFFGRYSSP--- 377
MNEI_FRU 334 AAAATAGMVAFCLR-EEHFTADHPFLFFIRHNKTKSILFLGRYSSPQ-- 380
MNEI_TNI 334 AAAATVAMVTFCLR-EEHFTADHPFLFFIRHNKTKSILFLGRYSSPQ-- 380
MNEI_DRE 334 AAAATAAIVAFCLR-EEHFMDADHPFLFYIRHNPTNSILFFGRFRGSS-- 380

```



**Appendix 8.3.7: Alignment of SPB5 (maspin) protein sequences from vertebrates.** Gene specific features includes a non-inhibitory RCL (red box). Conserved intron positions are indicated above the alignment. Group V1 specific sequence indels and intron indels are marked by \* and #, respectively.

ALAT_HSA	1	EDPQGDAAQKTDTS	SHHDQDHPTFNK	KITPNLAE	FAFSLYRQLAHQSNSTN	IFFSPVSI	IATA	60																																																			
SPB5_HSA	1	-----	MDALQLANS	AF	AVDLFKQLCEKEPLGNVLF	SPICL	STS	38																																																			
SPB5_MMU	1	-----	MDALRLANS	AF	AVDLFKQLCERDPAGNILE	SPICL	STS	38																																																			
SPB5_RNO	1	-----	MDALRLANS	AF	AVELFKQLCEKEPAGNILE	SPICL	STS	38																																																			
SPB5_GGA	1	-----	MDALQLANT	AF	AVDMFKKLCEKDRTANIVF	APLCT	STS	38																																																			
SPB5_XTR	1	-----	MDALRLANT	ALAVD	IFKKLCEKSATDNFVCS	PLCISS		38																																																			
78c#																																																											
ALAT_HSA	61	FAMLSL	CTKAD	THDE	I	LEGD	INFNLTE	IPEAQIHEGFQELLRTL	NQ	PDS	Q	L	Q	L	T	T	G	N	G	L	F	L	120																																				
SPB5_HSA	39	LSLAQV	CAKGD	THANE	IGQV	TH	FENVKD	----	IPFGFQ	T	V	T	S	D	V	N	K	L	S	S	F	Y	S	L	K	L	K	R	L	Y	V	94																											
SPB5_MMU	39	LSLAQV	CAKGD	THANE	IGQV	TH	FENVKD	----	V	P	F	G	F	Q	T	V	T	S	D	V	N	K	L	S	S	F	Y	S	L	K	L	V	K	R	L	Y	I	94																					
SPB5_RNO	39	LSLAQV	CAKGD	THANE	IGQV	TH	FENVKD	----	V	P	F	G	F	K	P	I	T	S	D	V	N	K	L	S	S	F	Y	S	L	K	L	K	R	L	Y	I	94																						
SPB5_GGA	39	LALAYK	ATKGD	THADQ	M	K	V	D	H	L	Q	D	V	K	D	----	V	S	F	G	F	Q	T	V	T	A	D	V	S	K	L	T	S	F	F	A	L	K	M	V	K	R	L	F	V	94													
SPB5_XTR	39	LSLIRK	CSQGN	THASE	L	E	K	A	TH	F	E	K	V	K	D	----	P	D	F	G	F	Q	L	L	S	S	D	I	S	K	I	S	S	A	N	S	L	K	L	L	K	R	V	Y	V	94													
128c																																																											
ALAT_HSA	121	SEGLKLV	DKF	LEDV	K	K	L	Y	H	S	E	A	F	T	V	N	E	G	-	D	T	E	E	A	K	Q	I	N	D	Y	V	E	K	G	T	Q	G	K	I	V	D	L	V	K	--	ELD	177												
SPB5_HSA	95	DKSLNL	STEF	IS	ST	K	R	P	Y	A	K	E	L	T	V	D	E	K	D	K	L	E	E	T	K	G	Q	I	N	S	I	K	D	L	T	D	G	H	F	E	N	I	L	A	D	N	S	V	N	154									
SPB5_MMU	95	DKSLNP	STEF	IS	ST	K	R	P	Y	A	K	E	L	T	V	D	E	K	D	K	L	E	E	T	K	G	Q	I	N	S	S	I	K	E	L	T	D	G	H	F	E	D	I	L	S	E	N	S	I	S	154								
SPB5_RNO	95	DKSLNL	STEF	IS	ST	K	R	P	Y	A	N	E	L	T	V	D	E	K	D	K	L	E	E	T	K	G	Q	I	N	S	S	I	K	E	L	T	D	G	H	F	E	D	I	L	P	E	N	S	I	S	154								
SPB5_GGA	95	DKSLSP	TTDF	V	N	S	T	K	R	P	F	P	S	E	L	E	L	V	E	F	F	K	E	T	E	T	R	Q	K	I	N	K	S	L	S	E	L	T	D	G	K	M	E	N	I	L	N	E	D	S	V	S	154						
SPB5_XTR	95	DNSIE	CKDF	I	N	S	A	K	P	Y	P	L	E	L	T	I	D	E	K	S	Q	A	E	E	A	R	T	Q	I	N	S	S	V	K	E	L	T	D	G	N	F	E	T	V	L	N	E	G	S	D	154								
212c																																																											
ALAT_HSA	178	RDIVF	ALV	NYI	F	F	K	G	K	T	E	R	P	E	F	V	K	D	I	E	E	E	D	F	H	V	D	Q	V	T	T	V	K	V	P	M	M	K	R	L	G	M	F	N	I	Q	H	C	K	L	S	237							
SPB5_HSA	155	DQTKIL	VVNA	AYF	V	G	K	M	K	K	F	P	E	S	E	T	K	E	C	P	F	R	L	N	K	T	D	T	K	P	V	Q	M	M	N	L	E	A	T	F	C	M	G	N	I	D	S	I	N	C	214								
SPB5_MMU	155	DQTKIL	VVNA	AYF	V	G	K	M	K	K	F	P	E	S	E	T	K	E	C	P	F	R	I	S	K	T	D	T	K	P	V	Q	M	M	N	L	E	A	T	F	C	L	G	N	I	D	D	I	S	C	214								
SPB5_RNO	155	DQTKIL	VVNA	AYF	V	G	K	M	K	K	F	P	E	S	E	T	K	E	C	P	F	R	I	N	K	T	D	T	K	P	V	Q	M	M	N	L	E	A	T	F	C	L	G	N	I	D	D	I	N	C	214								
SPB5_GGA	155	DQIQIL	VVNA	AYF	V	T	N	M	K	K	F	P	E	A	E	I	K	E	C	P	F	K	V	N	K	T	E	T	K	P	V	Q	M	M	N	L	E	A	T	F	C	L	G	Y	V	K	E	L	N	V	214								
SPB5_XTR	155	ENTKI	I	M	L	G	A	S	F	K	G	K	V	Y	T	E	N	K	S	E	T	K	E	M	D	F	H	I	N	K	K	E	T	K	P	V	Q	M	M	H	L	E	A	R	L	S	I	G	Y	I	N	E	L	K	T	214			
262c																																																											
ALAT_HSA	238	WVLLM	K	Y	L	G	N	-	A	T	A	I	F	F	L	P	----	D	E	G	K	L	Q	H	L	E	N	E	L	T	H	D	I	I	T	K	F	L	--	N	E	D	R	R	S	A	S	L	H	L	P	289							
SPB5_HSA	215	KIIE	L	P	F	Q	N	K	H	L	S	M	L	I	L	L	P	K	D	V	E	D	E	S	T	G	L	E	K	I	E	Q	L	N	S	E	S	L	S	Q	W	T	N	P	S	T	M	A	N	A	K	V	K	L	S	L	P	274	
SPB5_MMU	215	KIIE	L	P	F	Q	N	K	H	L	S	M	L	I	V	L	P	K	D	V	E	D	E	S	T	G	L	E	K	I	E	Q	L	N	P	E	T	L	L	Q	W	T	N	P	S	T	M	A	N	A	K	V	K	L	S	L	P	274	
SPB5_RNO	215	KIIE	L	P	F	Q	N	K	H	L	S	M	L	I	V	L	P	K	D	V	E	D	E	S	T	G	L	E	K	I	E	Q	L	N	P	E	T	L	L	Q	W	T	N	P	S	T	M	A	N	A	K	V	K	L	S	L	P	274	
SPB5_GGA	215	AILE	L	P	C	L	N	K	H	I	S	M	L	I	L	L	P	K	D	I	E	D	E	T	T	G	L	E	K	L	E	K	A	L	T	P	E	T	L	L	Q	W	T	N	P	S	M	M	A	N	T	K	V	N	V	F	L	P	274
SPB5_XTR	215	MVLE	M	P	F	Q	S	K	H	F	S	M	L	I	L	L	P	K	D	I	E	D	D	S	T	G	L	K	K	L	E	Q	D	M	T	F	E	K	Y	T	H	W	T	N	P	S	M	M	A	N	S	K	V	K	V	S	L	P	274

---

```

ALAT_HSA 290 KLSITGTYDIKSVLGLGKITKVFNSN-GADLSGVTEEAPLKLSKAVHKAVLTIDEKGTEAA 348
SPB5_HSA 275 KFKVEKMIDPKACIENLGLKHIFSEDTSDFSGMSETKGVALSNIHKKVCLEITEDGGDSI 334
SPB5_MMU 275 KFKVEKMIDPKASILESLGLKSLFNESTSDFSGMSETKGVLSNVIHRVCLEITEDGGESI 334
SPB5_RNO 275 KFKVEKMIDPKASILESLGLKSLFNESTSDFSGMSETKGVSVSNVIHRVCLEITEDGGDSI 334
SPB5_GGA 275 KFSVEGDYDIKPLLESLGNTNVEFNEASDFSEMCECTKGVVLSKIIHKVSLEVNEQGGESL 334
SPB5_XTR 275 KFKMENSYDIKDMKSLGINDAFNEEASDFSEMTESKGISISQAIQKACIEVDEDGTESA 334

```

```

ALAT_HSA 349 GAMFLEAIPMSIPPEVKFNK--PFVFLMIEQNTKSPLFMGKVVNPTQK- 394
SPB5_HSA 335 EVPGA----RILQHKDELNADHPPFIYIIRHNKTRNIIFFGKFCSP---- 375
SPB5_MMU 335 EVPGS----RILQHKDEFNADHPPFIYIIRHNKTRNIIFFGKFCSP---- 375
SPB5_RNO 335 EVPGS----RILQHKDEFKADHPPFLFIVRHNKTRNIVFLGKFSSP---- 375
SPB5_GGA 335 EVPGY----RILQHKDEFKADHPPFIIFLFRHNKTRNVILSGRFCSP---- 375
SPB5_XTR 335 DVSME----RRIMNKEEFLADHPPFIYIIRHNKTRTIIMLGRYCGPSEAS 379

```



**Appendix 8.3.8: Alignment of SPB6 orthologs and paralogs (pSPB6) from vertebrates.** Gene specific features includes an inhibitory RCL (red box). Conserved intron positions are indicated above the alignment. pSPB6 from *Fugu* and *Tetraodon* have common additional intron at 238c (novel), and 320a (feature of group V5). pSPB6 of *Tetraodon* has intron at position 85c with CD loop, not found in SPB6 of any other species reported till date. Presence of additional intron positions are marked by !. Group V1 specific sequence indels and intron indels are marked by \* and #, respectively.

```

AIAT_HSA      1  ED PQG DAA QK TDT SH HD Q DH P T FN K I TP N LA E F A F S L Y R Q L A H Q S N S T M I F F S P V S I A T A      60
SPB6_HSA      1  -----MDV LAE AN G T F A L N L L K T L G K D N - S K N V F F S P M S M S C A      37
SPB6_MMU      1  -----MD P L Q E AN G T F A L N L L K I L G E D S - S K N V F L S P M S I S S A      37
SPB6_RNO      1  -----MD H L Q E G N G T F A L K L L K T L S E D S - S N N I F L S P I S I S A A      37
SPB6_GGA      1  -----MD S L S A A N S T F A L D L L R E L R E K S S T K N L F F S P F S I S S A      38
SPB6_XTR      1  -----MD S L S A A N G T F A I N F L K K I N E S N K T G N I F V S P L S I S S A      38
pSPB6_DRE     1  -----ME P L S A A H A R F C L S L F Q K I S D G D S S Q N V F F S P L S I S S A A      38
pSPB6_FRU     1  -----MA A P S P L C K A N T S F S L A L F R K L S D N D T T A N I F Y S P F S I S S A      41
pSPB6_TNI     1  -----MA S P S P L S K A N T S F S L A L F R E L G D N D R T A N I F Y S P F S I S S A      41
SPB6_PMA      1  -----M O R R S A S P D      9

```

```

                                     78c#      85c!
                                     |          |
AIAT_HSA      61  F A M L S L G T K A D T H D E I L E G L N F N -----L T E I      87
SPB6_HSA      38  L A M V Y M G A K G N T A A Q M A Q I L S -----F N K S G G -      64
SPB6_MMU      38  L A M V F M G A K G T T A S Q M A Q A L A -----L D K C S G N      65
SPB6_RNO      38  L T M V F M G A K G M T A S Q M V Q T L S -----L D K C S G N      65
SPB6_GGA      39  L S M I L L G S K G D T E A Q I A K -----V L S L N      61
SPB6_XTR      39  L S M V L L G A K G N T A T Q M S Q L L K E L A I F D -----Y E F P L S I N E F K Q V L K L D      82
pSPB6_DRE     39  L S M L S L G A A G N T K D Q M S Q T L H F D G A E S -----      65
pSPB6_FRU     42  L A M V L L G A R G N T A A Q M S E V H H S N P A A S -----L K I K G      73
pSPB6_TNI     42  L A M V L L G A G G N T A T E M S E V L C F T E A E K P K D V E E Q Q Q Q Q Q L Q Q Q H H L R L P D F L K K C L K T E G      101
SPB6_PMA      10  L R -----D L E G D T W V Q H V N -----      23

```

```

                                     128c
                                     |
AIAT_HSA      88  P E A Q I H E G F Q E L L R T L N Q P D S Q L Q L T T G N G L F L S E G L K L V D K F L E D V K K L Y H S E A F T V N F      147
SPB6_HSA      65  -G G D I H Q G F Q S L L T E V N K T G T Q Y L L R V A N R L F G E K S C D F L S S F R D S C Q K F Y Q A E M E E L D F      123
SPB6_MMU      66  G G G D V H Q G F Q S L L T E V N K T G T Q Y L L R T A N R L F G D K T C D L L A S F K D S C L K F Y E A E L E E L D F      125
SPB6_RNO      66  G G G D V H Q G F Q S L L A E V N K T G T Q Y L L K T A N R L F G E K T C D L L A S F K D A C R K F Y E A E M E E L D F      125
SPB6_GGA      62  K A E D A H N G Y Q S L L S E I N N P D T K Y I L R T A N R L Y G E K T F E F L S S F I E S S Q K F Y H A G L E Q T D F      121
SPB6_XTR      83  K V D D A H C N F Q S L I S E I N K S G T N Y L L R T A N R L Y G E K S Y T F L E E F L G S T Q K H Y H A D L K A V D F      142
pSPB6_DRE     66  ---Q I H A G F T K L L T E M N R A G A P H T L S L A S R L Y G E Q S C R F Q E T F L S D T R R L Y G A E L Q P L D F      122
pSPB6_FRU     74  L E D D V H V S F S Q L L N E L H K E N A P Y A L S V A N R L Y G E Q S Y Q F V E D F L G S T K K H Y R A E L S V D F      133
pSPB6_TNI     102 C Q D D I H T S F S Q L L D E L H K K N A P Y A L S V A N R L Y G -----K H Y R A E L S V D F      146
SPB6_PMA      24  -----Y V G W G K M V H R F D -----F L D S S A K F Y R A E L A A V N F      53

```

```

                                     167 a#
                                     |          **
AIAT_HSA      148  G - D T E E A K K Q I N D Y V E K G T Q G K I I V D L V K --E L D R D I V F A L V N Y I F F K G K W E R P F E V K D T E      204
SPB6_HSA      124  I S A V E K S R K H I N T W V A E K T E G K I A E L L S P G S V D P L I R L V L V N A V Y F R G N W D E Q E D K E N T E      183
SPB6_MMU      126  Q G A T E E S R Q H I N T W V A K K T E D K I K E V L S P G T V N S D T S L V L V N A Y F K G N W E K Q E N K E H T R      185
SPB6_RNO      126  K G D T E Q S R Q R I N T W V A K K T E D K I K E L L A P G I V D P D T V L V L V N A Y F K G N W D K Q E N K E H T R      185
SPB6_GGA      122  K N A S E D S R K Q I N G W V E E K T E G K I Q K L L A E G I I N S M I K L V L V N A Y F K G N W E E K E D K E R I K      181
SPB6_XTR      143  S R K A E S R G E I N E W V A Q K T E G K I K D L L S G S V D S L I R L V L V N A Y F K G N W A N K E N P D H T H      202
pSPB6_DRE     123  I S Q P E A S R G I I N R W V E Q Q T H E K I R D L L A E G S V D S L S R L V L V N A V Y F K S S W E R K E L E E H H H      182
pSPB6_FRU     134  R A A A E T S R S N I N S W V E K Q T E G K I K D L L S D D V T G D T R L V L V N A Y F K G N W N E Q E K E N A T R      193
pSPB6_TNI     147  Q S A A E A S R I H I N S W V E K Q T E G K I K D L L V Q G I V S S D T R L V L V N A Y F K G R W N K Q E K E E A T R      206
SPB6_PMA      54  K G A F E E A R K E I N A W V E G Q T E G K I Q D L L A S G V V N S L I R L V L V N A V Y F K G S W D A K E D P E V I R      113

```



		212c		238c!		*		
AIAT_HSA	205	EEDFHVDQVTTVKVPMKRLGMFNIOHCKKLSWVLLMKVLGN-ATAIFFLE					-----	DEG 258
SPB6_HSA	184	ERLFKVKSKNEEKPVQMMFKQSTFKKTYIGEIFTQILVLPYVGKELNMIIMLP					-----	DETT 239
SPB6_MMU	186	EMPFKVKSKNEEKPVQMMFKKSTFKMTYIGEIFTKILLLPYVSELNMIIMLP					-----	DEHV 241
SPB6_RNO	186	EKPFKVKSKTEEKPVQMMFMKSTFKMTYIGEIFTKILLLPYAGNELNMIIMLP					-----	DEHI 241
SPB6_GGA	182	EMPFKINKNETKPVQMMFRKGYNMTYIGDLETKILEIPYIGNELSMIVLIPDAIQDEST						241
SPB6_XTR	203	ESPFRLNKNETKPVQMMFKKAKFPMTYIGELFTKVVEIPYVDNELSMIILLPDDINDGTT						262
pSPB6_DRE	183	EQQFRTSRNESKPVQMMFQKGRFPLAFIPDVNCQILELPYAGKELSMIVLIPNAMEDDGT						242
pSPB6_FRU	194	DATFHISKNSKPVQMMNQTSKFPFVFISEANCOVLQLPYVGKELSMILIFLPNQIEDSTT						253
pSPB6_TNI	207	DAQFNVTKNSKPVQMMHQTSKFPFTFIPEAKCOILEMPYIGEELSMILIFLPYQMEDSST						266
SPB6_PMA	114	DAEFKINKNEKPVQMMYKKAKYNFSHVEELNVNIVELPYEGHKLSMIVLIPLATEDETT						173
		262c		290c				
AIAT_HSA	259	KLQHLENELTHDIIITKFLENE--DRRSASLHLPKLSITGTYYDLKSVLGGQITKVFVS-NG						315
SPB6_HSA	240	DLRTVEKELTYEKFVEWTRLDMMDEEEVEVSLPRFKLEESYDMSVLRNLGMDAFELGK						299
SPB6_MMU	242	ELSTVEKEVITYEKFIEWTRLDKMDDEEEVEVFLPKFKLEENYHMDALYKLGMDAFG-GR						300
SPB6_RNO	242	ELKTVEKELTYEKFIEWTRLDMLDEEEVEVFLPRFKLEENYDMKVVLGKLGMDAFMEGR						301
SPB6_GGA	242	GLEKLERELTYEKLMDWINPEMMDSTEVRSLPRFKLEENYDLKPILSHMGMRDAFDLRM						301
SPB6_XTR	263	GLEALEKELTYEKFLLKWTNPEMMDITEMELSLPRFKLEDDYDLESFLSTMGMSDAFDQRR						322
pSPB6_DRE	243	GLEKLERALTLETITLDWTRSDMMDVLEVEVSLPRLRVEERLELKLPLVELGMPDAFDQRR						302
pSPB6_FRU	254	GLEKLEKLLTYDNFMEWTRPETMKEVEVQVGLPRFKMEEKCNMKNILVSMGMVDAFNEAA						313
pSPB6_TNI	267	GLEKLEKLLTYDKFMEWTRPDMDSVEVQVGLPRFKLEEKFNMKNVLVMGMVEAFDVAT						326
SPB6_PMA	174	GLEKLESALTLLKSLRQWTSPEINMSKLEVEHLPRFRLEKSYTLNEHLQRLGMSVFTQGE						233
		320a!						
AIAT_HSA	316	ADLSGVTEEAPLKLKSKAVHKAVLTIDEK	GTEAAGAMFLEATPMS	IPPE	VK	----	FNKPFV	371
SPB6_HSA	300	ADFSGMS-QTDLSLSKVVHKSFVEVNEE	GTEAAAAATAA	IMMRCARFV	PR	--	FCADHPFL	356
SPB6_MMU	301	ADFSGMSKQGLFSLKVVHKAFVEVNEE	GTEAAAAATAGM	MTVRCMRE	TPR	--	FCADHPFL	358
SPB6_RNO	302	ADFSGIASKQGLFSLKVIHKAFVEVNEE	GTEAVAATGST	ITMRCLE	TPR	--	FLADHPFL	359
SPB6_GGA	302	ANFSGISSGNELVLSEVHKSFVEVNEE	GTEAAAAATAGVM	LRCAMIV	P	--	DFTADHPFL	359
SPB6_XTR	323	ADFSGMS SANDLFLSKVLHKSFVDVNEE	GTEAAAAATAA	IMMLRCAMIT	P	--	PRIVVTCDHPFL	382
pSPB6_DRE	303	ADFSGVCAGGELLSTVHVHQSFLVEVNEE	GTEAAAAATAAVM	TRCLMRAE	R	--	FCADHPFL	360
pSPB6_FRU	314	SDFSGISPANDLFLSDVHVHKAFFVEVNEE	GTEASAATGAV	FKFRCA	RRTET	--	FVADHPFL	371
pSPB6_TNI	327	SNFSMSPANDLFLSEVHKAFFVEVNEE	GTEAAAAATGA	IMMLRCAR	P	--	FYADHPFL	384
SPB6_PMA	234	ADFSGINGARDLYVSHVAHKAFVEVNEE	GTEAAAAATAIV	MRCARMG	P	--	VRADHPFI	291
AIAT_HSA	372	FLMIEQNTKSPLEMGKVVNDTQK						394
SPB6_HSA	357	FFIQHKTNGILFCGRFSSP---						376
SPB6_MMU	359	FFIHHVKTNGILFCGRFSSP---						378
SPB6_RNO	360	FFIQHVKTNGILFCGRFSSP---						379
SPB6_GGA	360	FFIRHNKTSSILFCGRYCSP---						379
SPB6_XTR	383	FFIMHRQTRSILFCGRFSSP---						402
pSPB6_DRE	361	MLIRHNP TGSLIF YGRVCNP---						380
pSPB6_FRU	372	FFIRHNP SRNIFLAGRYCFPE--						392
pSPB6_TNI	385	FFIRHNP SMSILFAGRYCSPE--						405
SPB6_PMA	292	FFIRENSSGSVLF LGRFASP---						311

Appendix 8.3.9: Alignment of group V1 serpin sequences from chicken genome. Group V1 specific sequence indels and intron indels are marked by \* and #, respectively.

AlAT_HSA	1	EDPQGDAAQKTDTSHHDQDHPTFNKITPNLAEF*	SLYRQLAHQSNSTNIF	51	
MNEI_GGA	1	-----	ME SLSNANSRFALDLFRKVNENP SGNIF	29	
SPB6_GGA	1	-----	MDSL SAANSTFALDLLRELREKSSTKNLF	29	
SPB10_GGA	1	-----	MERLSASTNSFTLDLYKKLDVTSKGNIF	29	
SPB10b_GGA	1	-----	MEQVSASIGNFTVDFLNKLNENRDKNIF	29	
Gga-Spn-5	1	-----	MEALNKANTSFALDFFKHECQEDDNKNIL	29	
SPB14_GGA	1	-----	MGSIGAASMEFCF*	FDVFNEMKVVHANE NIF	29
SPB14b_GGA	1	-----	MDSISVTNAKFCF*	FDVFNEMKVVHVNENIL	29
SPB14c_GGA	1	-----	MGSISAANAFCF*	FDVFNELKVVQHTNENIL	29
Gga-Spn-9	1	-----	MGSISRMIIEFCF*	LDLYNKLNRTAKGQ NIV	29
SPB5_GGA	1	-----	MDALQLANTAFVDMFKKLCEKDRTANIV	29	

78c

|

AlAT_HSA	52	FSPVSIATAFAMLSLGT	TKADTHDEILEGLNFN	-----	83
MNEI_GGA	30	FSPLSISTALAMVLLG	SRGNETQVLKTFH	-----	59
SPB6_GGA	30	FSPFSISSALSMILLG	SKGDTAQAQIAKVLIS	-----	59
SPB10_GGA	30	FAPWSIATALAMVYL	GAKGDTATQMAKGL	-----	58
SPB10b_GGA	30	FSPWSISSALALTYLA	AKGSTAREMAEVILH	FTAV-----RAESSSV	71
Gga-Spn-5	30	FSPLSISSALATVYL	GAKGNTADQMAKVL	YFNEAEGARNITTTIRMQVYSR	80
SPB14_GGA	30	YCP IAIMSALAMVYL	GAKDSTRTQINKVVR	FDKLPGFG-----	67
SPB14b_GGA	30	YCP LSI L TALAMVYL	GARGNTESQMKKVILH	FDSITGAG-----	67
SPB14c_GGA	30	YSP LSI I VALAMVYM	GARGNTEYQMEKALH	FDSIAGLG-----	67
Gga-Spn-9	30	FSPMSISTSLGLILL	CARNNTAAQIEEVILH	VSHATGTT SLES--ELEGAVP	78
SPB5_GGA	30	FAPLCTSTSLALAYK	ATKGD TADQMKKVILH	-----	59

85c

|

AlAT_HSA	84	-----	LTEIPEAQIHEGFQELLR	TNLNQPDSQLQLTTGNGLF	119
MNEI_GGA	60	-----	FDEVENI-HSRFRALTAD	INRRDSSCLLR IANRLY	93
SPB6_GGA	60	-----	LNK AEDA-HNGYQSLLSE	INNPDTKYILRTANRLY	93
SPB10_GGA	59	-----	EYEETENI-HSGFKELLS	AINKPGNTYLLKSANQLF	93
SPB10b_GGA	72	ARPSRGRPKRRRMDPE	HEQAENI-HSGFKELLT	AFNKPRNNYSLRSANRIY	121
Gga-Spn-5	81	TDERLSNHRACFQKTE	IGKSGNI-HAGFKALNLE	INOPTKSYLLRSINQLY	130
SPB14_GGA	68	-----	DSIEAQCGTSVNV-HS	SLRDILNQITKPNVDVYSFSLASRLY	107
SPB14b_GGA	68	-----	STTDSQCGSSEYV-HN	LFKELLSEITRPNATYSLEIADKLY	107
SPB14c_GGA	68	-----	GSTQTKCGKSVNI-HL	LFKELLSDITASKANYSLRIANRLY	107
Gga-Spn-9	79	ENKSELSQERESSPSL	CNTDGNLNHEAFHALLL	QLQLNLGKDYVLSLANS LF	129
SPB5_GGA	60	-----	LQDVKDV-SFGFQTVT	ADVSKLTSFFALKMVKRLF	93

128c

|

167a

|

AlAT_HSA	120	LSEGLKLVDKFLEDVKK	LYHSEAFVNFVGDTEE	AKKQINDYVEKGTQG--	167
MNEI_GGA	94	GEKSYSFLLEFLTNTQ	KLYGADLAAVDFLHAYGE	ARKEINQWVEEKTEG--	142
SPB6_GGA	94	GEKTFEFLSFFIESSQ	KFYHAGLEQTD FKNASE	SRKQINGWVEEKTEG--	142
SPB10_GGA	94	EDKTYPLLPKFLQLIT	RYYYQAKPQAVNFKTD	AEQARAQINSWVENETER--	142
SPB10b_GGA	122	VEKTYALLPTYLQLS	KKYKAEQPQVNFKTAPE	QSRKEINTWVEKQTES--	170
Gga-Spn-5	131	GEKSLPFSKEYLQLAK	KYYSAEPQSVDFVGAAN	AIRREINSTVEHQTEG--	179
SPB14_GGA	108	AEERYPILPEYLQCVK	EYLRGGLEPINFQTAADQ	ARELINSWVESQING--	156
SPB14b_GGA	108	VDKTFSVLPEYLSAR	KFYTGVEEVNFKTA	AEERQ LINSWVEKETING--	156
SPB14c_GGA	108	AEKSRPILPYLKC	VKKLYRAGLETVNFKTA	ASDQARQLINSWVEKQTEG--	156
Gga-Spn-9	130	IQQGFEPHQKYL	MC SKELYRAALETVDFQ	RALEASRLKINDWVSE	TQKGT 180
SPB5_GGA	94	VDKSLSP TDFVN	STKRFPFSELELVEFKE	KTEETROKINKSLSEL	LDG-- 142

		212c	
		**	
AlAT_HSA	168	---KIVD--LVKELDRDVFALVNYIFFKGRWERPFVEVKDIEEEDFHVDQ	212
MNEI_GGA	143	---KIPDLLSEGSVNSMFKLVLVNATYFKGNWAEKFEEANTADMPFRLNK	189
SPB6_GGA	143	---KIQKLLAEGIINSMFKLVLVNATYFKGNWEEKFDKERTKEMPFKINK	189
SPB10_GGA	143	---KIQNLLPAGSLDSDTFLVLVNATYFKGNWEEKRFLEKDTSEMPFRLSK	189
SPB10b_GGA	171	---KIKNLLSSDDVKATFRLILVNAIYFKAEWVKFQAEKTSIQPFRLSK	217
Gga-Spn-5	180	---KIKSLLPPGSIDSLFRLVLVNATYFKGNWATKFDADDTQRPFRLSK	226
SPB14_GGA	157	---IIRNVLQPSVSDSQTAMVLVNATYFKGLTEKTFKDEDITQAMPFRVTE	203
SPB14b_GGA	157	---QIKDLLVSSSIDFGTTFVFIINTIYFKGTFKIAFNTEDTREMPFRLSK	203
SPB14c_GGA	157	---QIKDLLVSSSTDLDITLVLVNAIYFKGTFKTAFAAEDTREMPFRLSK	203
Gga-Spn-9	181	EQNCKIKELFAPGVIDSHTILVLVNVYFKASWEHKFEEKNTVQRDFKLNQ	231
SPB5_GGA	143	---KMNILNEDSVSDQITQILLVNAAYEVTNWMKKFPEAEIKECFKVNK	189
*			
AlAT_HSA	213	VTTVKVPMKRLGMFNIQHCKKLSWVLLMKV LGN-ATAIFFLDEGK---	259
MNEI_GGA	190	NERKTVKMMYQKKKFRFGYISDMKTRVLELPDREFSMIILLDDIEDDS	240
SPB6_GGA	190	NETKPVQMMFRKGYKNTYIGDLETKILEIPYIGNELSMIVLLPDAIQDES	240
SPB10_GGA	190	TKTKAVQMMFLRDTFLMLHEQTMFKIIELEPVENELSMFVLLPDDISDNT	240
SPB10b_GGA	218	NKSKPVKMMYMRDTPFVLIMEKMFKMIELPVPKRELSMFILLPDDIKDGT	268
Gga-Spn-5	227	HTTKPVPIHLSDKFNWTVESAQIDVLELPVNNELSMFILLPREI----	273
SPB14_GGA	204	QESKPVQMMYQIGLFRVASMASEKMKILELPFASGTM SMLVLLPDEVS---	251
SPB14b_GGA	204	EESKPVQMMCMNNSFNVAATLPAEKMKILELPFASGDL SMLVLLPDEVS---	251
SPB14c_GGA	204	EESKPVQMMCMNNSFNVAATLPAEKMKILELPFASGDL SMLVLLPDEVS---	251
Gga-Spn-9	232	NERKPVQMMYQKGTFLGYIEELGTQVLELPYQKLLSMIILLGETADGS	282
SPB5_GGA	190	TETKPVQMMNLEATFCLGYVKELNVAILELPCLNKHISMLILLKDIET	240
262c			
AlAT_HSA	260	---LQHLENELTHDIITKFLNEDR--RSASLHLPKLSITGTYDIKSVLGG	305
MNEI_GGA	241	-TGLQKLEQQLTLEKLQEWTRPEHLYSTDVHVHLPKFKLEESYDIKSDLSA	290
SPB6_GGA	241	-TGLEKLERELTYEKLMDWINPEMMDSTEVRLSLPRFKLEENYDKPILSN	290
SPB10_GGA	241	-TGLELVERELTHEKLAEWSNSARMKVEVELYLPKPKIEENYDLTSTLSN	290
SPB10b_GGA	269	-TGLEQLERELTYERLSEWADSKMMTETLVLDLHLPKFSLEDRIDLRDLRN	318
Gga-Spn-5	274	-TGLQKLINELTFEKLSAWT SPELMEKMKMEVYLP RFTVEEKYDIKSTLSK	323
SPB14_GGA	252	--GLEQLESIIINFEKLTETWSSNVMEERKIKVYLP RPKMEEKYNLTSVLMA	300
SPB14b_GGA	252	--GLERIEKTINFDKLREWTSTNAMAKKSMKVYLP RPKMEEKYNLTSVLMA	300
SPB14c_GGA	252	--GLERIEKTINFDKLREWTSTNAMAKKSMKVYLP RPKMEEKYNLTSVLMA	300
Gga-Spn-9	283	PSGLEQTESTMTYENMLWFSSSEHMFETVVEVYLP RPKFKLEGTFNMEVLKA	333
SPB5_GGA	241	-TGLEKLEKALTPETLLQWTNPSMMANTKVNVLPRPKFSVEGDYDIKPLLES	290
AlAT_HSA	306	LGITKVF SNG-ADLSGVTEEAPLKLKSKAVHKAVLTIDEKGT <b>TEAAGAMFLEA</b>	355
MNEI_GGA	291	MGLLDIFDSAKADLSGMSGAHDLFLSKIVHKAFVEVNEEG <b>TEAAAAATAGIA</b>	341
SPB6_GGA	291	MGM RDAFDLRMANFSGISSGNELVLSEVVHKSFVEVNEEG <b>TEAAAAATAGVM</b>	341
SPB10_GGA	291	MGIQNAFDPVQADFTRMSAKKDFFLSKVIHKAFVEVNEEG <b>TEAAAAATGVLV</b>	341
SPB10b_GGA	319	MGMTTAF T-TNADFRGMTDKKDLAISKVIHQSFVAVDEKGT <b>TEAAAAATAVII</b>	368
Gga-Spn-5	324	MGIEDAFTTEGQADFRGMS ENADLFLSQVFHKCYVEVNEEG <b>TEAAAASSASL</b>	374
SPB14_GGA	301	MGITDVFSSS-ANLSGISSAESLKISQAVHAAHAEINEAG <b>REVVGSAEAGV</b>	350
SPB14b_GGA	301	LGMTDLEFSRS-ANLTGISSVDNLMISDAVHGVFMEVNEEG <b>TEATGSTGAIG</b>	350
SPB14c_GGA	301	LGMTDLEFIPS-ANLTGISSAESLKISQAVHGA FMELSEDG <b>IEMAGSTGVIE</b>	350
Gga-Spn-9	334	MGMTDIFSESKADLSALSSEKSLVLSNIVHKAYVEVNEEG <b>TTAAAAATGATI</b>	384
SPB5_GGA	291	LGMTNVFNE SASDFSEM CETKGVVLSKIIHKVSLVNEEG <b>GESLEVPGYRI</b>	341

---

```

A1AT_HSA 356 IPMSIPPE----VKFNKPFVFLMIEQNTKSPLEFMGKVVNPTQK 394
MNEI_GGA 342 MLCMVIEE---DFNADHPFLFFLRHNPTKSIVFFGRYASE--- 378
SPB6_GGA 342 VLRCAMIV--PDFTADHPFLFFIRHNKTSSEILFCGRYCSP--- 379
SPB10_GGA 342 LRSRTPRV---TFKADHPFLFFIRHNKSKTILFFGRLCSP--- 378
SPB10b_GGA 369 SFTTSVINHVLKFKVDHPFHFFIRHNKSKTILFFGRFCQVE- 410
PAI2_GGA 375 ASRTLGAT--VIFVADHPFLFIIRHNKTKCILEFLGRFCSP--- 412
SPB14_GGA 351 DAASVSEE----FRADHPFLFCIKHIATNAVLEFFGRCVSP--- 386
SPB14b_GGA 351 NIKHSLEL--EEFRADHPFLFFIRYNPTNAILEFFGRYWSP--- 388
SPB14c_GGA 351 DIKHSPEL--EQFRADHPFLFLIKHNPTNTIVYFGRYWSP--- 388
SPB12_GGA 385 VRRSLPLI--EVFIADRPFLFFIRHNPTSTILFFGKFCSP--- 422
SPB5_GGA 342 LQHKDEFK-----ADHPFLFLFRHNKTRNVILSGRFCSP--- 375

```

Appendix 8.3.10: Alignment of group V1 serpin sequences from *Xenopus tropicalis* genome. Group V1 specific sequence indels and intron indels are marked by \* and #, respectively.

ALAT_HSA	1	EDPQGDAAQKTDTSHHDDHPTFNKITPNLAEFAFSLYRQLAHQSNSTNI	50
MNEI_XTR	1	-----MENLSSACTHFSFDLFRKINENNATGNV	28
SPB6_XTR	1	-----MDSLSAANGTFAINFLKKINE SNKTGNI	28
SPB5_XTR	1	-----MDALRLANTALAVDIFKKLCEKSATDNF	28
Xtr-Spn-2	1	-----MD-ICTANNEFTIDVLR EISKTAAGQNV	27
Xtr-Spn-5	1	-----MSSLKSFSEFSLDLCKELKKNP EKKNI	28
Xtr-Spn-6	1	-----MESINKSINEFSLDIFKELNSSCENKNI	28
		78c#	
			85c
ALAT_HSA	51	FFSPVSIATAFAMLSLGTKADTHDEILEGLNFNL-----TEI	87
MNEI_XTR	29	FFSPISISTALAMVLLGARGNTAQQISRILHFDAVKDLHSN-----	69
SPB6_XTR	29	FVSPLSISSALSMVLLGAKGNTATQMSQVLKLDKVDDAHCHN-----	69
SPB5_XTR	29	VCSPLCISSSLSLIRKGSQGN TASELEKALHFEKVKDPDFG-----	69
Xtr-Spn-2	28	VFSSMSIMTSLAMVYLGAGNTAADMGKALHFDEVEDVHAQ-----	68
Xtr-Spn-5	29	LFSPLSICSAMGLVLLGSKGDTAAEIEKVFHFPAAAGSRSSKPS CQQQTC	78
Xtr-Spn-6	29	FFSPMSISAALYLLHLGSR EDTATQIQKVSECGKVS DAHSK-----	69
ALAT_HSA	88	PEAQIHEGFQELLRTL NQPDSQ-LQLTTGNGLFLSEGLKLV DKEFLEDVKK	136
MNEI_XTR	70	-----FQTLNAEINKKNVSSYALNLANRLFGEKSFKFLPDEFLSSVKK	111
SPB6_XTR	70	-----FQSLISEINKSGTN-YL LRTANRLYGEKSYTFLEEF LGSTQK	110
SPB5_XTR	70	-----FQLSSDISKISSA-NSL KLLKRVYVDNSIECKDFEINS AKK	110
Xtr-Spn-2	69	-----FRVLLKELMKN GND-YTLTFTVNKLFGEKKY YFLPTFLKAINA	109
Xtr-Spn-5	79	QAQGVHLLFKDLFSALNKP NDH-YELSIANRAYGEKSF PFSEQYLLCIEQ	127
Xtr-Spn-6	70	-----FHALLSKLTEDPKG-VEL QIANGMFAQMNF PFLQQYLECAQA	110
		167a#	
			**
ALAT_HSA	137	LVHSEAF TVNFGDT-EEAKKQINDYVEKGTQGGKIIVD--LVKELDRDVFVA	183
MNEI_XTR	112	QYNADLGTVD FISAEDARKEINTWVSEQTKGKIPEVLSAGAVNSFKLV	161
SPB6_XTR	111	HYHADLKAVDF SRKAEESRGEIN EWVAQKTEGKIKDLLSSG SVDSLRLV	160
SPB5_XTR	111	PMPLELETID FKSQAEEARTQINS SVKELLDGNFETVLNEGSCDENIKII	160
Xtr-Spn-2	110	FYGAPLEKVD FSSNPEATRSYINAWIQEKTGKTIQNL LPENSISPNTVLM	159
Xtr-Spn-5	128	LVNATLESVD FKTADDVIQQINAWVESKTGKTIQNL FAKGSLDSTTALA	177
Xtr-Spn-6	111	LVNAKLQNVDF EK--DETRENIN SWVESKTGKTIKDLFEKNSLDKRIALV	158
		212c	
ALAT_HSA	184	LVNYIFFKGRTERPFEVKDTEEDFHVDQVTTVKVPMKRLGMFN IQHCK	233
MNEI_XTR	162	LVNAIYFKGDFAKKFKAEHTK DMPFQLNKKEQKTVKMMYQMEKLPFNYIP	211
SPB6_XTR	161	LVNAIYFKGNFANKFNPDH THESPERLNKNETKPVQMMFKKAKFPMTYIG	210
SPB5_XTR	161	MLGAASEFKGRVYTFNKSET KEMDFHINKKETKPVQMMHLEARLSIGYIN	210
Xtr-Spn-2	160	VANTLYFLANWTTQFSEHAT SKAPFTLITNEQIKVNM MATMNTFNMKRIK	209
Xtr-Spn-5	178	LVNAVYFKGSSWKKQFKKENT DAPFFLNKNDKTSVKMMSQKGKYKLG SNP	227
Xtr-Spn-6	159	LVNAIYFKGDFSNPFQEVHTK DAPFYVSKDVVKSVPMMYQS QKFNLGAIK	208

			*?	262c			
ALAT_HSA	234	KLSSWVLLMKY	LGN--ATAIFFLP	-----DEGKI	QHLENELTHDIITKFL	276	
MNEI_XTR	212	EINCRVLELEPY	YVD-YELSMVIVLP	DNINDDTTGL	QQLEKELSLEKINEWT	260	
SPB6_XTR	211	ELFTKVVEIPY	VD-NELSMIILLP	DDINDGTTGLE	ALEKELTYEKFLKWT	259	
SPB5_XTR	211	ELKTMVLEMPF	QS-KHFSMLILLP	KDIEDDSTGL	KKLEQDMTFEKYTHWT	259	
Xtr-Spn-2	210	NPGMSVLELEPY	GDTKDLSMVLMLP	----DNSTV	LTKVDREISYENLSKWT	255	
Xtr-Spn-5	228	ELKCRILKLEPY	EE-G-FSMKIILP	----DDIDGL	AELETHLTYETFTKLM	271	
Xtr-Spn-6	209	ELNAQILELEPY	QL-GALSMFILLT	----NEKFG	LQKIEQQLSWNYLAKGM	253	
ALAT_HSA	277	ENEDR--RSASLHLP	KLSITGTYDLK	SVLGQLGITKVF	SNG-ADLSGVTE	323	
MNEI_XTR	261	EN-MM-PIDVHVHLP	KFKLEDSYK	LKSQLAGMG	MADLFEAGSADLSGMSG	308	
SPB6_XTR	260	NPEMMDITEMELSLP	KFKLEDDYDLES	FLSTMGM	SDAFDQRRADFSGMSS	309	
SPB5_XTR	260	NP SMMANSKVKVSLP	KFKMENSYDLK	DMKSLGINDAF	NEEASDFSEMTE	309	
Xtr-Spn-2	256	RSENMS SNYLAVYLP	RFRMEKSFSL	KKVLS SLGMSS	SAFSQSRANFSGMGK	305	
Xtr-Spn-5	272	DLQRTREVQVVVKLP	QFKFGETYS	LTEVILQSM	GMTSAFHG--ANLSGISD	319	
Xtr-Spn-6	254	SNMEN--TKLDVYIP	RFRLEESLDL	GSHLINMGM	VDAFSEAKANLSGISD	301	
ALAT_HSA	324	EAPLKLKSKAVH	KAVLTIDEK	GTEAAGAMFL	---EAIPMSIPPE	VKFNKPF	370
MNEI_XTR	309	SNDLYLSEVTHK	SFVEVNEE	GTEAAAASAG	-IAMM-CLMREEE	FNANHPF	356
SPB6_XTR	310	ANDLFLSKVLH	KSFVDVNEE	GTEAAAATAA	-IMMLRCAMIIP	RIVCDHPF	358
SPB5_XTR	310	SKGISISQAIQ	KACIEVDED	GTESADV	----MERRLMNKEE	FLADHPF	354
Xtr-Spn-2	306	QKQLYVSDVH	HKTFIEVNEK	GTEAASATGS	-VMSIRSLAN-	EEFKADHPF	353
Xtr-Spn-5	320	KAGLAISTVVH	KSYIEVNEE	GTEAAAATG	GIGITVTSAPL	PPQEFIVDRPF	369
Xtr-Spn-6	302	VP-LYVSKIVH	KAFVEVNEE	GTVA A AATGV	QIAPKMAVIP-	RVFKADHSE	349
ALAT_HSA	371	VFLMIEQNTKSP	LFMCKVVND	TQK-		394	
MNEI_XTR	357	LFFIRHNATKS	ILFFGRYSSP	----		377	
SPB6_XTR	359	LFFILHRPSQS	ILFCGRFALP	----		379	
SPB5_XTR	355	IYILRHNKTRT	IIMLGRYCGP	SEAS		379	
Xtr-Spn-2	354	HFFIRHNKTNC	ILLYGKFYSP	----		374	
Xtr-Spn-5	370	LFCIEHISTKSL	LFYGRFTFPEI	--		392	
Xtr-Spn-6	350	LFFIKDNPNDT	ILFFGKYESP	----		370	



Appendix 8.3.11: Alignment of group V1 serpin sequences from *Denio rerio* genome. Group V1 specific sequence indels and intron indels are marked by \* and #, respectively. X indicates yet non-identifiable amino acids in respective serpins of *Danio*.

ALAT_HSA	1	EDPQGDAAQKTDTS	SHHDQDHPTFNKITPNLAEFAFSLTYRQLAHQSNSTNI	50
MNEI_DRE	1	-----	MEGVSRANSLFALDLYRALSSASSAEGNI	28
pSPB6_DRE	1	-----	MEPLSAAHARFCLSLFQKISDGDSSQNV	28
Dre-Spn-2	1	-----	MEGVSRANSLFALDLYRALSSASSAEGNI	28
Dre-Spn-4	1	-----	MESLSAANTQFSLNLFKKISGGNASGNV	28
Dre-Spn-5	1	-----	MEPVIAANTKFSLDLLKQLCQ-KSKDNV	27
Dre-Spn-6	1	-----	MESLSAANTQFSLNLFKKISGGNASGNV	28
Dre-Spn-28	1	-----	NSLFALDLYQALSASSAEGNI	21
Dre-Spn-29	1	-----	MESLSAANTQFSLNLFKKISGGNASGNV	28
Dre-Spn-30	1	-----	MESLSAANTQFSLNLFKKISGGNASGNV	28
Dre-Spn-31	1	-----	MESLSAANTQFSLNLFKKISGGNASGNV	28
			78c#	
ALAT_HSA	51	FFSPVSIATAFAMLSL	GLTKADTHDEILEGLNFNL-----	84
MNEI_DRE	29	FFSPLSISAALSMVYL	GARGDTAGEMEKVLSFSS-----	62
pSPB6_DRE	29	FFSPLSISAALSMLSL	GAAAGNTKDOMSQTLHFDG-----	62
Dre-Spn-2	29	FFSPLSISAALSMVYL	GARGDTAGEMEKVLCFSS-----	62
Dre-Spn-4	29	FYSPVSISSALAMVSL	GAKGNTADQMFKVLGFN-----	61
Dre-Spn-5	28	LFSPVSISSALGVL	LGAKGETADQMYKVLQFFK-----	61
Dre-Spn-6	29	FYSPVSISSALAMVSL	GAKGNTADQMFKVLGFNPPKPGGATPTP--AQA	76
Dre-Spn-28	22	FFSPLSISAVLSMVYL	GARGDTAAEMERVLSS-----	55
Dre-Spn-29	29	FYSPVSISSALAMVSL	GAKGNTADQMFKVLGFNPPKPGGATPTP--AQA	76
Dre-Spn-30	29	FYSPVSISSALAMVSL	GAKGNTADQMFKVLGFN-----	61
Dre-Spn-31	29	FYSPVSISSALAMVSL	GAKGNTADQMFKVLGFNNLPKSAGATPEAHQSM	78
			85c	
ALAT_HSA	89	-----	TEI-P----EAQIH	93
MNEI_DRE	63	-----	VSDVH	67
pSPB6_DRE	63	-----	AESQIH	68
Dre-Spn-2	63	-----	VSDFH	67
Dre-Spn-4	62	SQAHQP-----	VEQIH	72
Dre-Spn-5	62	-----	ET	63
Dre-Spn-6	77	TQKPQITCGVKSQHEP	QALQQPQKFELPADLKKCP-AQPVPGQKAEQIH	125
Dre-Spn-28	56	-----	VSDVH	60
Dre-Spn-29	77	TQKPQITCGVKSQHEP	QALQQPQKFELPADLKKCP-AQPVPGQKAEQIH	125
Dre-Spn-30	62	SQAHQP-----	VEQIH	72
Dre-Spn-31	79	QQAQPKPSGVKDQHG	QAMMQOTQKIDIPAEKVCNQC SAVPGQKAEQIH	128
			128c	
ALAT_HSA	94	EGFQELLRTLNPDS	QLQLTTGNGLEFLSEGLKLVDKFIEDVKKLYHSEAF	143
MNEI_DRE	68	SHEESLISSINSP	SASYILRLANRLYGEKTF SFLPEYLSSSLNLYHADLQ	117
pSPB6_DRE	69	AGFTKLLTEMNR	AGAPHTLSLASRLYGEQSCR FQETFLSDTRRLYGAELQ	118
Dre-Spn-2	68	AHEFKTLISSIN	SPSASYILRLANRLYGEKTF SFLPMYVDSTMKLYHAEPQ	117
Dre-Spn-4	73	SNEFKLMRELNK	PGAPYVLSLANRLYGEQTYQLIEKFLNDTKRYDAGLE	122
Dre-Spn-5	64	IIMHDLYEAIN	QGGKNRKLLKLYNRMFGDRT IDFLDGYLDKCEEWCFAGIR	113
Dre-Spn-6	126	SSENKFMSELNK	PGAPYVLSLANRLYGEQTYQFLEKYLSDAKKYAAGLE	175
Dre-Spn-28	61	SHEESLISSIN	SPSASYILRLANRLYGEKSF SFLPECLDSTMKLYHAELQ	110
Dre-Spn-29	126	SSENKFMSELNK	PGAPYVLSLANRLYGEQTYQFLEKFLSDAKTYAAGLE	175
Dre-Spn-30	73	SNEKFMSELNK	PEAPYVLSLANRLYGEQTYQLIEKFLNDTKRYDAGLE	122
Dre-Spn-31	129	SNEKFMSELNK	PGAPYVLSLANRLYGEQTYQFVEKFLSDAKRYEAGLE	178

		167 a#	
		**	
ALAT_HSA	144	TVNFG-DTEEAKKQINDYVEKGT-----	QGKIVD--LVKE 175
MNEI_DRE	118	AVDFIGASEQSRQLINRWVEEQT-----	ENKIRDLLKPGM 152
pSPB6_DRE	119	PLDFISQPEASRGILNRWVEQQT-----	HEKIRDLLAEGS 153
Dre-Spn-2	118	TVDFIRAADDSRQFINRWVEKQT-----	ENQIKDLLQPGV 152
Dre-Spn-4	123	KVDFINKSEDARVNIINTWVEKNT-----	QEKIKDLLPSGA 157
Dre-Spn-5	114	NVDFKTNPEAARAQIINTWVKNKTKVDCKPDADQKCDPYKNS IENLLGKED	163
Dre-Spn-6	176	KVDFKNKSEASRVNINKWVEKNT-----	QEKIKDLLPSGA 210
Dre-Spn-28	111	TVDFIGASEGSRQLINKWVEKQT-----	ENKIRDLLKPGM 145
Dre-Spn-29	176	KVDFKNKSEASRVNINKWVEKNT-----	QEKIKDLLPSGA 210
Dre-Spn-30	123	KVDFINKSEDARVNIINTWVEKTH-----	KVRXXXXXXXXXX 157
Dre-Spn-31	179	KVDFKNKSEARVNIINTWVEKNT-----	QEKIKDLLPSGA 213

		212c	
ALAT_HSA	176	LDRTVFALVNIYIFFKGKWERPFVVDTEEDFHVDQVTTVKVPMMKRLG	225
MNEI_DRE	153	VTGMTRLALVNAIYFKGNWLQRFNAQDTKEMPFKINQKENRPVQMMYQKK	202
pSPB6_DRE	154	VDSLRLVVLVNAVYFKSSWERKFLEEHTHEQQFRTSRNESKPVQMMFQKG	203
Dre-Spn-2	153	VNEMTRLLLVAIYFKGNWMHTFDAHATKEMPFKINQNESRPVQMMDQVE	202
Dre-Spn-4	158	IDAMTRLVVLVNAIYFKGNWEEKFPKEATRDGVFRLNKNQTKPVKMMHQKA	207
Dre-Spn-5	164	VSKDSVLALISVMHFEARWAQSFEP LHT-----NKDKAONGQMMQTTQ	206
Dre-Spn-6	211	IDAMTRLVVLVNAIYFKGNWEEKFPKEATKDGQFKLNKNQTKPVKMMHQKA	260
Dre-Spn-28	146	VTTMTRLALVNAIYFKGKWTHTFQAKYTREMAFKINQKESHPVRMMHQLN	195
Dre-Spn-29	211	IDAMTRLVVLVNAIYFKGNWEEKFTKEATRDGQFKLNKNQTKPVKMMHQKA	260
Dre-Spn-30	158	XXXNQT KP VKMMHQKA	207
Dre-Spn-31	214	IDAMTRLVVLVNAIYFKGNWERKFPKEATNDGQFKLNKNQTKPVKMMYQKA	263

		*	
ALAT_HSA	226	MFNIQHCKKLS SWVLLMKYLG N-ATAIFFLPDEGK-----LQHLENEILTH	269
MNEI_DRE	203	KFPFNYIYDHRVQVLELPIYVKEELSMLILLPEETQDGS DPLLKLESEIT I	252
pSPB6_DRE	204	RFPLAFIPDVNCQILELPIYAGKELSMLVLLPNAMEDDGTGLEKLERAILL	253
Dre-Spn-2	203	NFPYRCIPEYKLVLELPIYQOELSMLILLPDEIKYGSDPLLKLESEINL	252
Dre-Spn-4	208	EFPSGYIEEMKSHVLELPIYAGKNLSMLIILPDEIEDATTGLEKLERALTY	257
Dre-Spn-5	207	IFPLGEIPDADSKMLELPIYENS DV SFLLIILPNQN----DGLPKLLDSMTH	252
Dre-Spn-6	261	QFPFVVIPEINSQILELPIYVGKNLSMLIILPDEIEDATTGLQKLEKALTY	310
Dre-Spn-28	196	KLPFRCLPEYKLVLELPIYIQOELSMLILLPDETKDGS DPLLKLEKELTL	245
Dre-Spn-29	261	RFSLASIPEMNSQVLELPIYAGKNLSMLIILPDQIEDATTGLQKLEKALTY	310
Dre-Spn-30	208	EFPSGYIEEMKSHVLELPIYAGKNLSMLIILPDEIEDETTGLQKLERALTY	257
Dre-Spn-31	264	HFPLASIPEMNSQVLELPIYVGKNLSMLIILPDQIEDATTGLEKLEKALTY	313

ALAT_HSA	270	DIITKFLNEDRRSA---SLHLPKLSITGT YDLKSVLGGQLGITKVF S-NG	315
MNEI_DRE	253	DKLHEWTNRNMDTQTDIIIVHLPRFKLEIES SILVEIIMGMGMS SVFQEGK	302
pSPB6_DRE	254	ETLFDWTRSDMDV L-EVEVSLPRLRVEERLELKP LLLVBLGMPDAFDPQR	302
Dre-Spn-2	253	QKLLDWT SRGKMDTWRKIIIVRLPKFKLEIESCLSETLEKMGMS SVFQETK	302
Dre-Spn-4	258	EKIMEWTKPEVMHQR-EVQVSLPKFKTEQTYDMKSLLVSMGMEDVDFDPQK	306
Dre-Spn-5	253	EKILMWTQSHWMTPT-EVTVDMP IFQLKEKYDLEEAMKALGMDVDFG-DS	300
Dre-Spn-6	311	EKIMQWTK--VMRQQ-EVQVSLPKFKTEQTYDMKSLLVSMGMEDVDFDPLK	357
Dre-Spn-28	246	EKILLDWTNRDKMDTQGAIVHLPKFKLEIESCLSETLEKMGMS SVFQETK	295
Dre-Spn-29	311	EKIMEWTKPSMCCQ-EVQVSLPKFKTEQTYDMKSLLVSMGMEDVDFDPQK	359
Dre-Spn-30	258	EKIMEWTKPEVMHQR-EVQVSLPKFKTEQTYDMKSLLVSMGMEDVDFDPQK	306
Dre-Spn-31	314	EKIMEWTKPEVMRQQ-EVQVSLPKFKMEQTYDMKSLLVSMGMEDVDFDPQK	362



ALAT_HSA	316	ADLSGVTEEAPLKLSKAVHKAVLTIIDEK	GTEAAGAMFLEAIPMSIPP	---	362
MNEI_DRE	303	ADLTGMTGHGGLFLLSAVAHKAFVDVNEE	GTEAAAATAAIVAFCLRE	---	349
pSPB6_DRE	303	ADFSGV CAGGELLSTV VHQSFLEVNEE	GTEAAAATAAVMMTRCL	---	MR 349
Dre-Spn-2	303	ADLTGMSSNGGLFLLSAVIHKAFVEVNEE	GTEAAAATALLLPISACQG	---	A 350
Dre-Spn-4	307	VNLTGMSSSNDLVLSKVIHKAFVEVNEE	GTEAAAATAAIKFLCY	---	IP 353
Dre-Spn-5	301	CDLSGMAS-GKLLKLSKV VHGCSVNVDEK	GINADTGNGGVVKS LCK	---	VP 347
Dre-Spn-6	358	VNLTGMSSSNDLVLSKVIHKAFVEVNEE	GTEAAAATGAVVSI RIL	----	402
Dre-Spn-28	296	ADLTGMSSNGGLFVSAVIHKAFVDVSEE	GTEAAAATCVYIITSYVPR	EP	345
Dre-Spn-29	360	VNLTGMSSSNDLVLSKVIHKAFVEVNEE	GTEAAAATGV IATLTSMP	L	SP 408
Dre-Spn-30	307	VNLTGMSSSNDLVLSKAIHKAFVEVNEE	GTEAAAATAAIEKLMCY	---	IP 353
Dre-Spn-31	363	VNLTGMSSSNDLVLSKVIHKAFVEVNEE	GTEAAAATGAIMMLR	CT	RL 409

ALAT_HSA	363	--E V K F N K P F V F L M I E Q N I K S P L F M G K V V N P T Q K-	394
MNEI_DRE	350	-E H F M A D H P F L F Y I R H N P I N S I L F F G R F R G P S---	380
pSPB6_DRE	350	A E R F C A D H P F L M L I R H N P I G S L L F Y G R V C N P----	380
Dre-Spn-2	351	F H D F I A D H P F M F F I R H N P I N S I L F L G R F R A P S---	382
Dre-Spn-4	354	P V S F N A D H P F L F F I R H N P I K S I L F Y G R F C S P----	384
Dre-Spn-5	348	T N R F V A D R P F L F F I R H N P I K S I L F W G R F N P Q G P V N	382
Dre-Spn-6	403	A Q I F N A D H P F L F F I R H N P I N T I L F Y G R F C S P----	433
Dre-Spn-28	346	R Y Y F T A D H P F M F F I R H N P S N N I L F L G R Y R S P S---	377
Dre-Spn-29	409	P K T F T A D H P F I F F I R H N P I N A I L F Y G R F S S P----	439
Dre-Spn-30	354	P L S F N A D H P F L F F I R H N P I K S I L F Y G R L C S P----	384
Dre-Spn-31	410	P Q S F N A D H P F L F F I R H N P I K S I L F Y G R F C S P----	440

Appendix 8.3.12: Alignment of AGT (serpinA8) protein sequences from vertebrates. Gene specific features include conserved angiotensin (cyan) and non-inhibitory RCL (red boxes). Conserved intron positions are indicated above the alignment and additional introns are marked with \*, found in AGT of *Fugu* (two) and *Tetraodon* (one).

ALAT_HSA	-	-----	-----	-		
AGT_HSA	1	MRKRAPQSEMAPAGVSLRATILCLLAWAGLAAG	DRVYIHPFHL	VIHNEST 50		
AGT_MMU	1	-----	MPPTGAGLKATIFCILTWVSLTAG	DRVYIHPFHL	LYHNKST 41	
AGT_RNO	1	-----	MPPTGAGLKATIFCILTWVSLTAG	DRVYIHPFHL	LYYSKST 41	
AGT_GGA	1	-----	MKLAAGLLCLLLCFTAVGC	DRVYVHPFSL	NAINESA 36	
AGT_XTR	1	-----	MNIQRIWLCLTVICIGYSL	NRVYIHPFNL	FAYNKSE 36	
AGT_FRU	1	-----	MQLLQPLLPALLLCCYLSP	QANRVYVHPFSL	FAAENVS 39	
AGT_TNI	-	-----	-----	-----	-	
AGT_DRE	1	-----	MKMFLAFLFLSCFAMART	NRVYVHPFNL	FSSNIS 35	
ALAT_HSA	1	-----	-----	EDPQGDAQKTDTS	SHHDQ 18	
AGT_HSA	51	CEQLAKANAGKPKDPTFIPAP	IQAKTSPVDEKALQDQLV	LVAAKLDTEDK	100	
AGT_MMU	42	CAQLENP SVETLPESTFEPVP	IQAKTSPVNEKTLHDQLV	LAAEKLEDEDR	91	
AGT_RNO	42	CAQLENP SVETLPEPTFEPVP	IQAKTSPVDEKTLRDKLV	LATEKLEAEDR	91	
AGT_GGA	37	CEELERLAQEG--KTFVPAS	IESQTPAYEEDVKDEVRLD	SPSLSVRGR	84	
AGT_XTR	37	CEKVEKQNHTI--EALFTPVS	IEVNIISP--EETLGST-VQ	SKLLGIVER	81	
AGT_FRU	40	CESLQQT SKPLQTI PVAPLET	DVLTDPD-----	SKDVVKIEGQR	DIVT 82	
AGT_TNI	-	-----	-----	-----	-	
AGT_DRE	36	CEVIQSEEHKPLETVHPLP	PLPGSTDPD-----	PRTASAAESLKN-LT	77	
ALAT_HSA	19	DHPTFNKITPNLAE	FAFSLYRQLAHQSNST	NIFFSPVSI	ATAFAMLSLGT 68	
AGT_HSA	101	LRAAMVGMLANFLG	FRIYGMHSELWGVVHGAT	VLSPPTAVFGT	LASLYLGA 150	
AGT_MMU	92	KRAAQVAMIANFVG	FRMYKMLNEAGSGASGA-	ILSPPALFGT	LVSFYLGS 140	
AGT_RNO	92	QRAAQVAMIANFMG	FRMYKMLSEARGVASGA-	VLSPPALFGT	LVSFYLGS 140	
AGT_GGA	85	QKLIYLKDFVHVLG	MRFYNLQREAR--QGQNV	LLSPPTSLYGS	SLASFYLGA 132	
AGT_XTR	82	QRVSIPLSLVNDAG	ERSFNGWRKTH--KDDSI	LMSEFTNLF	GSLSVSFYLGA 129	
AGT_FRU	83	ERTMALAGLVNVLG	LGRMY--EALS-KRHST	NLLSPVST	CGTLVNFYLGA 129	
AGT_TNI	-	-----	-----	-----	-	
AGT_DRE	78	QRTAVLAELQNSL	GGRMY--QTLRSTQKHT	NLLSPLNA	FGALVTLYLGA 125	
		77c*				
ALAT_HSA	69	KADTHDEILEGLN	FN-----	LTEIPEAQIHEGFQ	ELLRTLNQPD	SQ- 109
AGT_HSA	151	LDHTADRLQAILG	VVPWKDN--CTSR	LDAHKVLSALQ	AVQGLLVAQ	GRAD 198
AGT_MMU	141	LDPTASQLQTL	LDVPVKEGD--CT	SRLDGHKVLA	ALRAVQGLL	VTOGGSS 188
AGT_RNO	141	LDPTASQLQVLL	GVPVKEGD--CT	SRLDGHKVLT	ALQAVQGLL	VTOGGSS 188
AGT_GGA	133	SNQTAADLQGLL	GFVPPSGDSNCT	SRVDGRKLL	ESLRTIESL	VKTQDEE- 181
AGT_XTR	130	STHTSADLQAF	LGFQAHQSGDQDCV	SKVDALKV	ISTLKHIDN	RFLSKDNS- 178
AGT_FRU	130	SKKTAASSFQSL	LGLSRSDGEDCV	SLMDGHKV	LKTLQININ	SLVDDGPKD- 178
AGT_TNI	-	-----	-----	-----	-----	-
AGT_DRE	126	SKKTAISYQQL	LGLNLESEQTD	CAYFVDGHT	VLRTLQAI	SAHVDESRK-- 173

ALAT_HSA	110	----LQLTTGNGLFLLSEGLKLVDKFLEDVVKLYHS-EAFTVNF-GDTEEA	153
AGT_HSA	199	SQAQLLLSTVVGVTAPGLHLKQPFVQGLALYTPVVLPRSLDF-TELDVA	247
AGT_MMU	189	SQTPLLQSIMVGLFTAPGFRLLKHSFVQSLALFTPALFPRSLDLSTDPVLA	238
AGT_RNO	189	SQTPLLQSTVVGLFTAPGLRLKQPFVESLGPFTPAIFPRSLDLSTDPVLA	238
AGT_GGA	182	----LLFSKVFCFLFSAPGILLSQQFVHNLLPSADAFYTRAVDF-TNPSEA	226
AGT_XTR	179	----VESLKMTCFLVSKHVSLSSETFIQNLIP SADKFYVRGVDF-TNSAKA	223
AGT_FRU	179	-----EITTHVWTFTRPQIQLSDFVQGTKDFS DASFIRSVNF-SSPEVA	222
AGT_TNI	-	-----	-
AGT_DRE	174	-----ELRTLWTFVNSDADLSKEFLRGTQDFSDDSFVRSVDF-SQAKDA	217

## 192a

ALAT_HSA	154	KKQINDYVEKGTQCKIVDLVKELDRDITVFALVNYIFKCKKWERPFEVKDT	203
AGT_HSA	248	AEKIDRFMQAVTGWKTGCSLMGASVDSTLAFNTYVHFOCKMKGFSLLAEP	297
AGT_MMU	239	TEKINRFKAVTGWKMNLPLEGVSTDSTLLFNTYVHFOGTMRGFSQLPGV	288
AGT_RNO	239	AQKINRFVQAVTGWKMNLPLEGVSTDSTLFFNTYVHFOCKMRGFSQLTGL	288
AGT_GGA	227	TKQINAFVEAKSKGQSKHLLTDLDPITDLLVAVDVRLAANAKKASWLKEP	276
AGT_XTR	224	VELINEYLNTRSTKKSTYISTPVDDSVNLFSTYIHFCKGTLKNSYLIPEP	273
AGT_FRU	223	ELTVNMFVEKTSDEKVKSAFKNLNSSSNLLFLT SFNFQGSWRTAFQPERT	272
AGT_TNI	1	-----SWRTAFQPOST	11
AGT_DRE	218	EVEVNFVFIQKTS DNKVKSMFKGVTPKITDLLFASSVHFKGNWKTAFQPEAT	267

## 233c\*

ALAT_HSA	204	EEEDFHVDQVTTVKVPMKRLGMFNIOHCKKLS SWVLLMKYLG NATAIFF	253
AGT_HSA	298	--QEFWVDNSTSVSVPMLSGMGTFOHWSDIQDNFSVTQVPFTESACLLLI	345
AGT_MMU	289	--HEFWVDNSISVSVPMISGTGNFQHWSDAQNNFSVTCVPLGERATLLLI	336
AGT_RNO	289	--HEFWVDNSTSVSVPMLSGTGNFQHWSDAQNNFSVTRVPLGESVTLILLI	336
AGT_GGA	277	--QEFWVDNRAISVPMLSVTGMFKYMTDTSETFSATEIPVGNVLLVLL	324
AGT_XTR	274	--QDFWIEPGKKIATPMISLSGMFNKYKHDINMSQLIVKVPLGENDFMILLI	321
AGT_FRU	273	SMQEFHINETTTVKAPIMTHTGQYHYLNDKVHRCTIVKLP LSKRSSMLLV	322
AGT_TNI	12	SVEEFHTNDT VTMAPIMTHTGQYHYLKDQVHRCTVVKLP LSKRSSMLLV	61
AGT_DRE	268	SDQDFWTQKNSSVQVPMHTGDYKYLDDAGRKCSIVRLGLSKRTFMILLV	317

## 282b

ALAT_HSA	254	LPDEG-KLQHLLENLTHDIITKFLNEDRRSASLHLPKLSITGTYDIKSV	302
AGT_HSA	346	QPHYASDLKVEGLTFQONSLNWMKLSPTIHLTMPQLVLQGSYDLQDL	395
AGT_MMU	337	QPHCTSDLDRVEALIFRNDLLTWIENPPPRAIRLTLPQLEIRGSYNLQDL	386
AGT_RNO	337	QPQCASDLDRVEVLVFOHDFLTWIKNPPPRAIRLTLPQLEIRGSYNLQDL	386
AGT_GGA	325	QPINGNDLDKVEAKLPQS-SAWLENLSPRKIKLTLPDEFRIEDSSDLQEF	373
AGT_XTR	322	QPINGNTLENMESSLSWDTFLKWLENLSSRYINLSLPMKMEIESSYDIQEI	371
AGT_FRU	323	LPHERTDLHNVE SKLPKNIISDWIQNLSEGTLELTPKFSMSSVHDMRDL	372
AGT_TNI	62	LPHQGSSSLQEI ESKLAKNIMSDWVQNLSEGTLELTPKFSMSSVHHMQDL	111
AGT_DRE	318	LPHEGASLQDIEKPLLT-VIPTWLRHLKEKYLELSLPMKFSLTAVTDLRSV	366

331c

```

ALAT_HSA 303 LGQLG---ITKVFENSGADLSGVTEEAPIKLSKAVHKAVLTIDEKGTEAAG 349
AGT_HSA 396 LAQAE---LPAILHTELNLOKLSNDR-IRVGEVLNSIFFELEAD-ERE-P 439
AGT_MMU 387 LAEDK---LPTLLGAEANLSNIGDTN-PRVGEVLNSILLELKA-GEEEQP 431
AGT_RNO 387 LAQAK---LSTLLGAEANLGKMGDTN-PRVGEVLNSILLELQA-GEEEQP 431
AGT_GGA 374 LADMK---LPALLGKEADLSKISDTH-LTVGKIMNKAFFKLSSDATHQPE 419
AGT_XTR 372 LSDME---LPYLLGKKADLSKISNAE-LTVGKVINKVHFELKESGEDTDI 417
AGT_FRU 373 LANMNPEIEAKLLGSQAQFSQLGNTKPFNVQVINKVIFEMSEEGTEVQE 422
AGT_TNI 112 LANMNPEIEAQLLGSQAQFSLGLPQPFSTDQVINKVIFEMAE-GGAEVQ 160
AGT_DRE 367 LSEMA--VEKYLMSDASFRMSKENFTVDKVLNKVVFEMTE-GGSEVQ 413

```

```

ALAT_HSA 350 AMFLEAIPMSIPPEVKFNKPEVFLMIEQNTKSPLEEMKVVNETQK-- 394
AGT_HSA 440 TESTQQLNKPEVLEVTLNRPFLFAVYDQSATALHELGRVANPLSTA- 485
AGT_MMU 432 TTSVQQPGSPEALDVTLSSPFLFAIYEQDSGTLHFLGRVNNPQSVV- 477
AGT_RNO 432 TESAQQPGSPEVLDVTLSSPFLFAIYERDSGALHFLGRVDNPQNVV- 477
AGT_GGA 420 DATAQ-EEDSVPQEVMLNKPELLAVFEAKSRAMLFLGRVTNPLQEV- 464
AGT_XTR 418 PLN--EKDQEPLEIKFQKPELFVVFEGTKALLEIGRVKSPLN--- 458
AGT_FRU 423 SVE----GPSPLKLSEFNRPFFFCVSEANSNAIMLLKITNPTE--- 462
AGT_TNI 161 GSAR----GAGSPLKVSEFNRPFFFCVSEANSNAILLLKITDPTL--- 201
AGT_DRE 414 NRTD----DGRAPHKVTFNRPFFFAVVEGNSNAIMLLKIINPTA--- 454

```

**Appendix 8.3.13: Alignment of heparin cofactor II sequences from vertebrates.** Gene specific features include an inhibitory RCL (red box), acidic repeats (blue boxes), Heparin binding residues (orange boxes). Conserved intron positions are indicated above the alignment. Novel intron positions are marked with \* in conserved part and in non conserved part are marked with HCII of the species and corresponding HCII sequence numbering.

		PMA-37c*		
ALAT_HSA	-	-----	-----	-
HCII_HSA	1	MKHSLNALLIFLIITSAWGGSKGPLDQLEKGGETAQASADPQWEQLNNKN-		49
HCII_MMU	1	MKHPLCT-LLSLITFMCIG-SKG-----	LAEQLTNEN-	30
HCII_RNO	1	MKHPAYTLLLSLIMSMCAG-SKG-----	LAEQLTKEN-	31
HCII_GGA	1	MKFLFPLLALAVIITSTFCG IKDFSDHFES-----	LKDAHTHENG	40
HCII_XTR	1	-MKLLHLATIFLLIHATLGGVKDLQEH-----	FEDTSTGIN-	35
HCII_FRU	1	MWVISLVCVAYLMVSPSLAGNIDLSSAFSDPKP-----	DPRGFEG--A	41
HCII_TNI	1	MWVISLVCVAWLMASPSLAETKHPSSPLSDPKP-----	DPRGFEG--T	41
HCII_GAC	1	MWVLGAVSVACLIVAPSLAEIKDLGSHFADP-----	EPRGFVEVG-G	40
HCII_DRE	1	MWLVPVIVVACLINSPALAGVKDLSSHSTLEKEKTV---	DARGLSPGGE	47
HCII_PMA	1	MFLYGLIFALSVLQWVEGQDQTKTAVDVN-----	KDAQLSFS--	37
FRU-85c*				
ALAT_HSA	-	-----	-----	-
HCII_HSA	50	-LSMPLLPADFHKENTVTNDWIPEGEED--	DDYLDLEKIFS--EDDDYI	93
HCII_MMU	31	-LTSFLPANFHKENTVTNDWIPEGEED--	EDYLDLEKLLG--EDDDYI	74
HCII_RNO	32	-LTVSLLPPNFHKENTVTNDWIPEGEED--	DDYLDLEKLLS--EDDDYI	75
HCII_GGA	41	TYNMPLPLEFHRENTITNDLIPEEEEE--	EDYLDLDKILG--EDD-YS	84
HCII_XTR	36	PRGSQTQAVENLLDDTVTNDLSTEGEDE--	EDYLDLFDKIFG--EDEDYI	80
HCII_FRU	42	AVDIEAIPLEFHKENTVTTEILFDGFED--	EDYIDFDKILS--SDEYFEG	87
HCII_TNI	42	EMDIEALPLEFHKENTVTKEIIFDGFED--	EDYIDFDKILAEGSDDYSDG	89
HCII_GAC	41	GADMEAVPLEFHKENTVTNDLLEFDGFED--	DDYIDFDKILAAGSDDYTEG	88
HCII_DRE	48	NTDME SIPLDFHRENTVTNDLP-EGQDD--	EDYVDFDKILG--EDDYSEG	92
HCII_PMA	38	--RYPNKPSDSLMDDTLALELDGFTDEDSLE	EDYIDFDKLLN--EDDDYP	82
PMA-118c*				
ALAT_HSA	1	-----	EDPQGDAAQKTDTSHH	16
HCII_HSA	94	DIVDSLSVSPD-----	SDVSAGNILQL	116
HCII_MMU	75	YIID--AVSPD-----	SESSAGNILQL	95
HCII_RNO	76	YVVD--AVSPD-----	SESSAGNILQL	96
HCII_GGA	85	DIID--AAPHV-----	SEIQQGNILEL	105
HCII_XTR	81	DIID--AAPEIKN-----	SETQQGNIFEL	102
HCII_FRU	88	DNIDEIATPAPDIDIFAEPD-----	PKIRRARLLRL	119
HCII_TNI	90	DNIDEIATPAPDIDIFAEPD-----	PKIRRARLLRL	121
HCII_GAC	89	DEIDEIATPAPDIDIFAEPD-----	PKIRRARLLRL	120
HCII_DRE	93	DHIDEISTPAPDLDFYEPD-----	PKIRRARLLRL	124
HCII_PMA	83	DEIDDINEDGSTGVTVDAEKVGLLHFTLSFSTEIKNLVDASFNKKLFLRR		132

ALAT_HSA	17	DQDHPTFNKITPNLAEF	AFSLYRQLAHQSN-STN	IFFSPVSIATAFAMLS	65
HCII_HSA	117	FHGKSRIQRLNILNAKFA	FNLYRVLKDQVNTFDN	IFIAIPVGISTAMGMIS	166
HCII_MMU	96	FQGKSRIQRLNILNAKFA	FNLYRVLKDQATTS	DNLFIAIPVGISTAMGMIS	145
HCII_RNO	97	FQGKSRIQRLNILNAKFA	FNLYRVLKDQATSS	SDNLFIAIPVGISTAMGMIS	146
HCII_GGA	106	FQGKTRIQRLNILNANFG	FNLYRSVADKANS	SDMILMAPVGISTAMAMIS	155
HCII_XTR	103	FHGKTRVQRLNIINANFG	FNLYRAIKNNTDASE	NILLAPVGISTAMATIS	152
HCII_FRU	120	FQGQSRLQRLNIINARFG	FNLYRSLRNTVNQSD	NILLAPAGISIAIGMMS	169
HCII_TNI	122	FHGQSRLQRLNIVNAHFG	FNLYRSLRNTVNQSD	NILLAPAGISIAMGMMS	171
HCII_GAC	121	FHGRSRLQRLNIVNAHFG	FNLYRSIRNDVNQSD	NILLAPAGISVAMGMMS	170
HCII_DRE	125	FHGQTRLQRLINVVNARFG	FLRYRKLRLNQLT	DNILLAPVGISIAMGMMG	174
HCII_PMA	133	FQGKTRIQRLSIVNSDF	AFNLYRSVSESTP	SGENLLLAPLGISSTLG	MIA 182

83c\*

|

ALAT_HSA	66	LGTKADTHDEILEG	INFN-----L	TEIPEAQIHEGFQELL	RTLNQPDSQL 110
HCII_HSA	167	LGLKGETHEQVHSIL	HFKDFVNASSKYE	ITTIHNLFRKLR	HLFRFRNFGY 216
HCII_MMU	146	LGLRGETHEEVH	SVLHFRDFVNASS	KYEVTTIHNLFRK	LRHLFRFRNFGY 195
HCII_RNO	147	LGLRGETHEEVH	SVLHFRDFVNASS	KYEVTTIHNLFRK	LRHLFRFRNFGY 196
HCII_GGA	156	LGLKGQTQQEVLS	VLGFEDFINASAK	YELMTVHNLFRK	LRHLFRFRNFGY 205
HCII_XTR	153	LGTGQTTLEQVLL	TLGFKDFLNASS	KYELTLHNVFRK	LRHLFRFRNFGY 202
HCII_FRU	170	LGTGAGTHDQIY	KAMGFSEFVNASH	HYDNTTVHKLFR	KLRHLFRFRNFGY 219
HCII_TNI	172	LGAGAETQDQIY	KAMGFSEFVNASH	CHYDNTTVHKLFR	KLRHLFRFRNFGY 221
HCII_GAC	171	LGAGSGTRDQIY	GALGFADFNASH	HYDNTTVHKLFR	KLRHLFRFRNFGY 220
HCII_DRE	175	LGVGPNTQEQLF	QTVGFAEFVNASH	HYDNSTVHKLFR	KLRHLFRFRNFGY 224
HCII_PMA	183	LGANGGTHKEIY	KALGFESLVDSS	SKYNI	STVHKLFRHLNHLFRFRNFGY 232

ALAT_HSA	111	QLTTGNGLFLSEGL	KLVDFLELVKKLYH	SEAFTVNEGDTEE	AKKQINDY 160
HCII_HSA	217	TLRSVNDLYIQKQ	FPILLDFKTKVRE	YFAEAQIADFS	-PAFISKTNNH 265
HCII_MMU	196	TLRSVNGLYIQKQ	FPIREDFKAAMRE	FYFAEAQEANFP	D-PAFISKANNH 244
HCII_RNO	197	TLQSVNDLYIQKQ	FPIREDFKAAMRE	FYFAEAQEADFS	-PAFISKANSH 245
HCII_GGA	206	TLRSVNDLYIRK	DFSILNDFRNMM	KTYFADAQPA	DFSD-PNFITKTNER 254
HCII_XTR	203	TLRSVNDIYV	KRDFLIREPFKNN	LKNYYFAEAQT	VDFGN-KDFLTKANKR 251
HCII_FRU	220	KLRAVNDVYV	KDVAVKDVFRAE	TKAYYFAEPQSV	NFRD-PGFLDKANSR 268
HCII_TNI	222	NLRAVNDVYIK	DVAVKDAFRAE	TKAYYFAEPQSV	NFRD-PAFLDKANSR 270
HCII_GAC	221	TLRSVNDVYV	KREVAVKDAFRAE	TKAFYFAEPQSV	DFGD-PAFLDKANSR 269
HCII_DRE	225	TLRSVNDLYV	KRVQIQDSFRAD	AKTYFAEPQSV	DFAD-PAFLVKANQR 273
HCII_PMA	233	TLKSASALYLQ	RWRPLLSYQQCL	RKTYFAEAHTV	DFKD-PATVQRINRW 281



## 192a

ALAT_HSA	161	VEKGTQCKTIVDLVKELDRDITVVFALVNYIFFKGGKWERPFEEVKDTEEEDEHV	210
HCII_HSA	266	IMKLTAKGLIKDALENIDPATQMMILNCIYFKGSGVVKFPEVEMTHNHNFERL	315
HCII_MMU	245	ILKLTAKGLIKEALENIDPATQMLILNCIYFKGTFVVKFPEVEMTHNHNFERL	294
HCII_RNO	246	ILKLTAKGLIKEALENTDSATQMMILNCIYFKGATMVKFPEVEMTHNHNFERL	295
HCII_GGA	255	ILKLTAKGLIKEALVNVPITILMILNCLYFKGTWENKFPVEMTKRSEFRL	304
HCII_XTR	252	IQQLTAKGLIKEALTNVDPALLMLLVNCIYFKGTWENKFPVEYTONMNFRL	301
HCII_FRU	269	ILKLTAKGLIRQPLKSIDPNMVLMLLNYLYFKGTWEQKFPKENTHYRNFVRV	318
HCII_TNI	271	ILKLTAKGLIRQPLKSIDPNMVLMLLNYLYFKGTWEQKFPKESHYRNFVRV	320
HCII_GAC	270	ILKLTAKGLIKEPLKSVDPNMVLMLLNYLYFKGAWEQKFPKEMHYRNFVRV	319
HCII_DRE	274	IQKLTAKGLIKEPLKSVDPNMAVMLLNYLYFKGTWEQKFPKELTHHRQFRV	323
HCII_PMA	282	VSSATKGTIISDAVTNIDPSTVFLVINSVYFKGPWEIKFASKHQISVRSERL	331

## 241c\*

ALAT_HSA	211	DQVTTVKVPMKRLGMFNIQHCKKLSVWVLLMKYLG NATAIFFLPD--EG	258
HCII_HSA	316	NEREVVKVSMQTKGNFLAANDQELDCDILQLEYVGGISMLIVVPHKMSG	365
HCII_MMU	295	NEREVVKVSMQTKGNFLAANDQELDCDILQLEYVGGISMLIVVPRKLSG	344
HCII_RNO	296	NEREVVKVSMQTKGNFLAANDQELDCDILQLEYVGGISMLIVIPRKLSG	345
HCII_GGA	305	NEKQTIKVPMQTKGNFLA AADPELDCGV IQLPFVGNISMLIVLPHKLSG	354
HCII_XTR	302	NEKELVKVPMKTKGNFLVAADPELDCAVLQLPVVGNISMLIVLPHKLSG	351
HCII_FRU	319	NEKTSVRVPMINKGNFLA AADHELECDILQLPYTGNISMLIALPRKITG	368
HCII_TNI	321	NEKTQVRVPMINRGNYL AADHDLDLDCDILQLPYRGNISMLIALPRKITG	370
HCII_GAC	320	NEKTNVRVPMTNKGNFLA AADHELEQCDILQLPYTGEISMLIALPSKING	369
HCII_DRE	324	NEKKQVRVILMQNKGSYLA AADHELNC DILQLPYAGNISMLIAVPOKLSG	373
HCII_PMA	332	NDKETVKVQMMQTKASF LVTTDHLELGC DILQLAYQGNVSMILAVPHKLKG	381

## 282b

ALAT_HSA	259	KLQHLENELTHDIITKFLENEDR-RSASLHLPKLSITGTYYDKSVLGLQLG	307
HCII_HSA	366	-MKTLEAQLTPRVVERWQKSMTN-RTREVLLPKFKLEKKNYNI VESLKLIMG	413
HCII_MMU	345	-MKTLEAQLTPQVVERWQKSMTN-RTREVLLPKFKLEKKNYNI VEV LKSMG	392
HCII_RNO	346	-MKTLEAQLTPQVVERWQKSMTN-RTREVLLPKFKLEKKNYNI VEV LKSMG	393
HCII_GGA	355	-MKALEKQITPQVVEKWQKSMTNSRTREV VLPKFKLEKKNYNI IGF LRSMG	403
HCII_XTR	352	-MKLLEKQISPQVVERWQNI MTN-RTREVFLPRFKLEKSYDILQKVL SMMG	399
HCII_FRU	369	-MRTLEQEI SPTVVSKWFKNM TN-RTREVVIPR FKLEQSYDIENLQELG	416
HCII_TNI	371	-MRTLEQDISPTVVSKWLKNM TN-RTREVVLP R FKLEQSYDIENLQKLG	418
HCII_GAC	370	-MRTLEQEI SPTVVKWLKNM TN-RTREVAIP R FKLEQNYDIENL KEMG	417
HCII_DRE	374	-MRSLEQEI SPTLVNKWLSNM TN-RTREV VFP R FKLEQNYDIENL KEMG	421
HCII_PMA	382	GLKTLERALSFDLLEKWLQAM TN-RTRDVIIP R FNLQOKYNIENNLKELG	430

## 331c

|

ALAT_HSA	308	ITKVF	SNGAD	DL	SGV	T	E	E	A	P	L	K	L	S	K	A	V	H	K	A	V	L	T	I	D	E	K	G	T	E	A	A	G	A	M	F	L	E	A	I	P	357							
HCII_HSA	414	IRMLF	D	K	N	G	N	M	A	G	I	S	-	D	Q	R	I	A	I	D	L	F	K	H	O	G	T	I	T	V	N	E	E	G	T	Q	A	A	A	V	T	T	V	G	F	M	P	462	
HCII_MMU	393	ITKLF	N	K	N	G	N	M	S	G	I	S	-	D	Q	R	I	A	I	D	L	F	K	H	O	S	T	I	T	V	N	E	E	G	T	Q	A	A	A	V	T	T	V	G	F	M	P	441	
HCII_RNO	394	ITKLF	N	K	N	G	N	M	S	G	I	S	-	D	Q	R	I	I	I	D	L	F	K	H	O	S	T	I	T	V	N	E	E	G	T	Q	A	A	A	V	T	T	V	G	F	M	P	442	
HCII_GGA	404	IEELF	S	E	K	G	N	Y	C	G	V	S	-	E	E	K	V	S	I	D	R	F	N	H	O	G	T	I	T	V	N	E	E	G	T	E	A	G	A	I	T	N	V	G	F	M	P	452	
HCII_XTR	400	ATDLF	T	-	H	G	D	F	S	G	V	S	-	D	K	D	I	N	I	G	L	F	Q	H	O	G	T	I	T	V	N	E	E	G	T	E	A	A	A	V	T	V	G	F	M	P	447		
HCII_FRU	417	LTDLF	F	K	D	S	G	D	F	S	E	M	T	-	S	E	K	V	S	M	N	W	L	K	H	O	G	T	I	T	V	N	E	E	G	T	E	A	A	A	L	T	Q	V	G	F	M	P	465
HCII_TNI	419	LTDLF	F	S	S	G	G	F	S	E	M	T	-	S	E	K	V	S	M	N	W	L	K	H	O	G	T	I	T	V	N	E	E	G	T	E	A	A	A	L	T	Q	V	G	F	M	P	467	
HCII_GAC	418	LTDLF	F	Q	E	S	G	D	F	S	A	M	T	-	S	D	K	V	H	M	S	W	L	K	H	O	G	T	I	T	V	N	E	E	G	T	E	A	A	A	L	T	Q	V	G	F	M	P	466
HCII_DRE	422	MTDLF	T	E	K	G	D	F	S	P	M	T	-	S	E	K	V	I	I	N	W	F	K	H	O	G	S	I	T	V	N	E	E	G	T	E	A	A	A	M	T	H	I	G	F	M	P	470	
HCII_PMA	431	VTELF	Q	A	N	A	D	L	S	G	M	T	G	A	K	D	V	Q	V	S	S	F	Q	H	O	G	F	I	K	I	D	E	E	G	S	E	A	A	A	V	T	T	V	G	F	T	P	480	

ALAT_HSA	358	M	S	I	P	P	E	V	K	F	N	K	P	F	V	F	L	M	I	E	Q	N	T	K	S	P	L	F	M	G	K	V	V	N	P	T	Q	K	394
HCII_HSA	463	L	S	T	Q	V	R	F	T	V	D	R	P	F	L	F	L	V	Y	E	H	R	T	S	C	L	L	F	M	G	R	V	A	N	P	S	R	S	499
HCII_MMU	442	L	S	T	Q	V	R	F	T	V	D	R	P	F	L	F	L	V	Y	E	H	R	T	S	C	L	L	F	M	G	K	V	T	N	P	A	K	S	478
HCII_RNO	443	L	S	T	Q	V	R	F	T	V	D	R	P	F	L	F	L	V	Y	E	H	R	T	S	C	L	L	F	M	G	R	V	A	N	P	A	K	S	479
HCII_GGA	453	L	S	T	Q	I	R	F	I	V	D	R	P	F	L	F	L	V	Y	E	H	R	T	N	C	L	L	F	M	G	R	V	V	N	P	A	K	P	489
HCII_XTR	448	L	S	T	Q	A	R	F	V	A	D	R	P	F	L	F	L	V	Y	E	H	R	T	N	C	L	V	F	M	G	R	V	A	N	P	T	K	S	484
HCII_FRU	466	L	S	S	Q	I	R	F	T	V	D	H	P	F	L	F	L	V	Y	E	H	R	T	D	C	L	V	F	I	G	R	V	S	N	P	S	Q	S	502
HCII_TNI	468	L	S	S	Q	I	R	F	T	V	D	H	P	F	L	F	L	V	Y	E	H	R	T	D	C	L	V	F	I	G	R	V	V	D	P	S	Q	S	504
HCII_GAC	467	L	S	S	Q	I	R	F	T	A	D	R	P	F	L	F	L	V	Y	E	H	R	T	D	C	L	V	F	M	G	R	V	V	N	P	S	Q	N	503
HCII_DRE	471	L	S	T	Q	T	R	F	I	V	D	R	P	F	L	F	L	V	Y	E	H	R	T	G	C	V	V	F	M	G	R	V	V	D	P	S	Q	T	507
HCII_PMA	481	L	T	S	H	N	R	F	V	A	D	R	P	F	V	F	I	V	Y	E	H	H	T	M	S	V	L	F	L	G	Q	V	S	N	P	A	K	N	517



**Appendix 8.3.14: Alignment of ZPI (serpinA10) protein sequences from vertebrates.** Gene specific features include an inhibitory RCL (red box), except ZPI3\_FRU and ZPI3\_TNI. Conserved intron positions are indicated above the alignment. There is a unique intron 94a position in fish specific Spn\_94a\_FRU and Spn\_94a\_TNI.

ALAT_HSA	1	-----EDPQGDAAQKTD-----	13
ZPI_HSA	1	MKVVP SLLS VLLAQVWLVPLA SPQSPETPAPQNTSRVQAPKEEED	51
ZPI_MMU	1	MRVASSLFLP VLLTEVWLVTSFNLS SHSPEASVHLESQDYENQTWEEYTRT	51
ZPI_RNO	1	MRVSSSLFLP VLLAEVWLVSSFNLS SHTPEAP IRLVSQDYENQTWEEYEWA	51
ZPI_GGA	1	MKTR--IYLLLLCELCFEISKADIKPKSPKDKRLNFLGRNKNVSISEEW	49
ZPI_XTR	1	-----MNSSLPSGN-----	9
ZPI1_FRU	1	-----MKMGF IFFFISAFICAH-----	18
ZPI1_TNI	1	-----MKIGLVFFVTSMFICAH-----	18
ZPI1_DRE	1	-----MKMGFFTLLIEASLLSV-----	18
ZPI2_DRE	1	-----MEFRLLLVFISACFLCSA-----	18
Spn_94a_FRU	1	-----MTPLLSLLLPGFLLGL-----	17
Spn_94a_TNI	1	-----MTPLFSLLLAGFLFLDL-----	17

ALAT_HSA	14	-----SHHDQDHPTFNKITPNLAEFAFSLYRQLAHQSNST	48
ZPI_HSA	52	EQEASEEKASE-----EEKAWLMASRQQLAKETSNFGFSLLRKISMRHDG	96
ZPI_MMU	52	DPREEEEEEEEEKEEGKDEEYWLRS-QQLSNETSSFGFNLLRKISMRHDG	100
ZPI_RNO	52	DPRD-----DNEYWLRS-QQLSNETSSFGFSLLRKISMRHDG	88
ZPI_GGA	50	QHKNDHKPL-----EEQSFEELTLHNFTEKTANFGFNLYRKIAMKLDN	92
ZPI_XTR	10	-----EEQIPTVLSFANVSQMSDFGFNLYRKIANKHDN	43
ZPI1_FRU	19	-----ORVQLPSATISDLSFKNMFAMNLYRTISSFHDK	52
ZPI1_TNI	19	-----BRLHLP SATISDLSFKNVD FAMNLYRKISSFHDK	52
ZPI1_DRE	19	-----VLGQTTD--VEELAIKNADFATRLYSKIASSDD	50
ZPI2_DRE	19	-----BHEELRTPDISDLAFRNTDFAINLYRKISSLHDR	52
Spn_94a_FRU	18	-----ASPATTDG SLEKLTNGNTDFAAKLYQAVASRTDD	51
Spn_94a_TNI	18	-----VSAQIPDGSVENLASRNVDFAARLYQAVASRTDD	51

94a!

|

ALAT_HSA	49	NIFFSPVSIATAFAMLSLGTKADTHDEILEGLNFNLTE-IPEAQIHEGFQE	98
ZPI_HSA	97	NMVFSPFGMSLAMTGLMGATGPIETQIKRGLHLQALK-PTKPGLLPSLFK	146
ZPI_MMU	101	NVIFSPFGLSVAMVNMLGAKGETKVQIENGLNLQALS-QAGPLILPALFK	150
ZPI_RNO	89	NVIFSPFGLSVAMVNMLGAKGETKVQVENGLNLQALS-QAGPLILPALFK	138
ZPI_GGA	93	NILISPLSVTTIMATYLLAAEGETHRQIAKALNLHSLKDRDR-HYLPALFK	142
ZPI_XTR	44	NIFFSPFSVSLGSSLLGTRGNTYDQLLHGLNYPFKDQENPYLLPELLK	94
ZPI1_FRU	53	NIFFSPLSISTSFAALLMASDGVITYKEILEGLNLHQLQAGQLDIPGLFK	103
ZPI1_TNI	53	NIFFSPLSISASFAALLMSDGVITYEEILEGLNLHQLQAGQDQPEVIPGLFK	103
ZPI1_DRE	51	NVAVSTLGATLALATLAAAGGATQSELLQIGVDSMVKDGQERIQNILQ	101
ZPI2_DRE	53	NVVFSPLSVSTCF SALLLAAQGSRTTEILKGLNLEALD-GGDSRRVPELFO	102
Spn_94a_FRU	52	NVCLSTFALSTALSALLSATSGPTQEQLLQGL--GLTGLDAQMLPELFOQL	100
Spn_94a_TNI	52	NVCLSTFALSSVLSALLSATSGPTREQLSQGL--ALTGLDPQTLPELFOQL	100

ALAT_HSA	99	LLRTLNPQDSQLQLTTGNGLFLSEGLKLVDFLEDVKKLYHSEFTVNEGD	149
ZPI_HSA	147	GLRETLNRNLELGLTQGSFAFIHKDFDVKETFNLSKRYFDTECVPMNFRN	197
ZPI_MMU	151	KVKETFSSNRDLGLSQGSFAFIHKDFDIKETYFNLSKKYFDIEYVSNFQN	201
ZPI_RNO	139	RVKETFSSNKKLGLTQGSFAFIHKDFEIKKTYFNLSMIFYDTEYVPTNFRN	189
ZPI_GGA	143	QLKDNITNEELLFVQGLSFIQKDFTVREAFNLNSKQYDFDMEFLCVDFQN	193
ZPI_XTR	95	TIKEKIAKNEELVNLIGLSFLHETFSMKDEFVNLTKKYFDMEYELIDFH-	144
ZPI1_FRU	104	LLMNNITQNGSLRLDQGMALFMHPKFRVEKTFQDQLKTTFFDADIKSVNFTN	154
ZPI1_TNI	104	LLSDNITQNGSLQLEQGMALFINTDFMVEKTFNEQLKTTFFDADIKSANFAD	154
ZPI1_DRE	102	QLREDA-----AQIPATGLFIKQDVKADDSFSNQVKQYYNADVQNVNYAN	146
ZPI2_DRE	103	QLHQNIS----LQMEQGTALFLDQHFHLQTNFSQQIQRFFNAEVLRVDFSK	149
Spn_94a_FRU	101	RTAIQPGN--ITNLKQAVALLPSHNFEVSASLRELVQTKFGGYMPSLKYTD	149
Spn_94a_TNI	101	RTTTQQGNT-ATCLKQAVAVLPSHNFEVSASERQLVQTKFGGYIPNVRYSD	150

		192a	
ALAT_HSA	150	TEEAKKQINDYVEKGIQCKIVDLVKEIDRDIVFALVNYIFFKGGKWERPFEV	200
ZPI_HSA	198	ASQAKRLMNHYNKEIEKQIPKLFDEINPEIKLILVDYILFKGKWLTPFDP	248
ZPI_MMU	202	SSQARGLINHCIVKEIEKQIPKLFDEINPEIKLILVDYVLFKGGKWLTPFDP	252
ZPI_RNO	190	SSQARGLMNHYNKIEIEKQIPKLFDEINPEIKLILVDYILFKGKWLTPFDP	240
ZPI_GGA	194	STQAKFVINQNIKQRIKCKISELFEVDRHSKLLLLDYIFFKGGKWLTPFNS	244
ZPI_XTR	145	SSKAKNEINAYVEKLIKGLISNFYDFIDPQIKLILLDYIFFKGGKWLTPFNP	195
ZPI1_FRU	155	TRGSVKFINEYIKRKSHDKISNMVSSIDPMTGLMLTNTIFFQGSWELPFNP	205
ZPI1_TNI	155	SMGTVKLINEYFKSKIHDKISDVVSSVDAMIQLVLTNTIFFQGSWELPFNP	205
ZPI1_DRE	147	GQQAKGSINDYVVRGRIEKEKVRDVENVDPOSMALLISAFFETGQWLQPFNA	197
ZPI2_DRE	150	PAVCRSLINEFVSRKIGRKVLEMLESVPELITQMLLNTIFYKGDWERPFNP	200
Spn_94a_FRU	150	QAEAISTINRWAQDQIQGDIQQVVTAVDAQTELLLATVSYKYQTQFSPLFNA	200
Spn_94a_TNI	151	QAEAISTINRWAQDQIQGDKIQQVFTALDAQIQLLLLATVSYKYQTQFSPLFNA	201

ALAT_HSA	201	KDTEEEDEHVDQVTVKIVPMMKRLGMFNIQHCKKLSWVLLMKYILGNATAI	251
ZPI_HSA	249	VFTEVDTEHLDKYKTIKVPMYAGAGKFASTFDKNFRCHVLKLPYQGNATML	299
ZPI_MMU	253	SFTEADTEHLDKYRAIKVPMYREGNFTSTFDKKFRCHILKLPYQGNATML	303
ZPI_RNO	241	IFTEADTEHLDKYKAVKVPMMYREGNFSTFDKKFRCHILKLPYQGNATML	291
ZPI_GGA	245	EFTEIETEHINKYRSVQVPMFKSDKVNSTYDENLRCNVIKLPYKGGKAYML	295
ZPI_XTR	196	ALTEVDSEFIDKYNVTVPMYKTDKVASVFDKDLSCVFKLPYRGNHML	246
ZPI1_FRU	206	NITVNAPEYIDNYSVVQVPMFLEDKFYMMVDTFLGVKVLKLPYKEGVSMML	256
ZPI1_TNI	206	NFTVIAPPEYIDNYSVQVPMFLEDKFYMTMDKNLGVNVKLPYKEGVSMML	256
ZPI1_DRE	198	TFTIQEDRFYVKNYNIVQVPMMLRSKGYLAYDPTFKVGLKLPCENGIAML	248
ZPI2_DRE	201	NNTEKSRFYVDKYNIVQVPMMLLEKFSVVEDRDLRARVLRLLPYRGGASML	251
Spn_94a_FRU	201	SLTQDERFYVKNYVVMVPMFRADKYFLAYDRSLKVGVLKLPMSDGMAML	251
Spn_94a_TNI	202	SLTQDERFYVKNYAVVMVPMFRADKFFLAYDPLLVKVGVLKLPMSDGTAML	252

		282b	
ALAT_HSA	252	FFLEDE-GKQLHLENELTHDIITKFLNEDRRSASLHLPKLSITGTYDLKS	301
ZPI_HSA	300	VVILMEKMGDHLALEDYLTDDLVEVTLWRNMKTRNMEVFFPKFKLDQKYEMHE	350
ZPI_MMU	304	VVILMEKTGDYLALEDYLTVDLVEVTLWQNMKTRKMEVFFPKFKLNQRYEMHE	354
ZPI_RNO	292	VVILMEKSGDHLALEDYLTDDLVEVWLQDMKTRKMEVFFPKFKLNQRYEMHE	342
ZPI_GGA	296	IVIDPEKGEDYVSLDHLTMELVESWLANMKSRNMDISEFPKFKLEQKYKMKK	346
ZPI_XTR	247	IIKPEKEGDFGILEDHITKELINSWQAKMQSRKTDIFFPKFKLDQKYKLS	297
ZPI1_FRU	257	IVLBNKNVDYTTDDEITADRIFRWTKMLRKVKLEVHLPKFKMEHSYSLHE	307
ZPI1_TNI	257	IVLBNKNVDYTEIDDEITADRIFRWTKRLQKTKLEVNLKPKFKMEQSYRLHE	307
ZPI1_DRE	249	VLLPDEDVDYTYVDESMTGEVFRGWVAKLKKTKLEIQLPFRFSLKQSNLSV	299
ZPI2_DRE	252	ILLPNADADYTAIEDEISAERLHGWIKNMRRIKTYTPHTETHNTHAHTCHV	302
Spn_94a_FRU	252	VVLPDEDVDIIVVEEKVTGKIRGWIRQLKKTKEVQLPFRFMLEKSYALRD	302
Spn_94a_TNI	253	VVLPDEDVDIIDVEEKMTGKIRAWIRQLKKTKEVQFPFRFLEKSYMMGD	303

		331c	
ALAT_HSA	302	VLGQLGITKVFSSNGADLSGVTEE-APLKLSKAVHKAVLTIDEK <b>GTEAAGAM</b>	351
ZPI_HSA	351	LLRQMGIRRIEFPFADLSELSATGRNLQVSRVLQRTVIEVDER <b>GTEAVAGI</b>	401
ZPI_MMU	355	LLKQMGIRRLFSTADLSELSAMARNLQVSRVLQSSVLEVDER <b>GTEAVSGT</b>	405
ZPI_RNO	343	LLKQVGIIRRIEFPFADLSELSAVARNLQVSKVVQSSVLEVDER <b>GTEVSGT</b>	393
ZPI_GGA	347	LLYALGIKNLFARTADLSHLTDQ-KHLTVSQVVQKAVIEVDEK <b>GTEAAAAT</b>	396
ZPI_XTR	298	SLNELGIKELFTGKANLTDLTEE-RNLMLTEITQAMIEVDER <b>GTEAAAVA</b>	347
ZPI1_FRU	308	IILPDMGMASVEDDS-ANLSKLRCLTFYISLKVLIHKAVIEVDET <b>GTTAAAAT</b>	357
ZPI1_TNI	308	IILPDVGMASIEDNS-ANLTKLSKDQGLKVSEVLIHKAVIEVDET <b>GTIAAAVT</b>	357
ZPI1_DRE	300	SILPSLGVKEIFGSTANLTGISSS-EGCLKLSEVVQKVAVDVDES <b>GGSLAEAS</b>	349
ZPI2_DRE	303	HTCAIKHICFFCQVNESVCGLAQALTLTSLVSETKAVIEVYEQ <b>GTSAASST</b>	353
Spn_94a_FRU	303	VILQTLNMNKMFDQDADIEMGS--KGPKLTQVYQKSAVAVGDS <b>RDEVATAG</b>	351
Spn_94a_TNI	304	FLQTLNVTMVFDQDAGEIREMGA--KGPRLTQVYQTSALSVRDS <b>SEEVMTGG</b>	352

---

ALAT_HSA	352	<b>F</b> <b>L</b> <b>E</b> <b>A</b> <b>I</b> <b>P</b> <b>M</b> <b>S</b> <b>I</b> <b>P</b> <b>P</b> <b>E</b>	VKFNKPPVFLMIEQNTKSPLEMGKVVNP	TQK	394
ZPI_HSA	402	<b>L</b> <b>S</b> <b>E</b> <b>I</b> <b>T</b> <b>A</b> <b>Y</b> <b>S</b> <b>M</b> <b>P</b> <b>P</b> <b>V</b>	IKVDRPPEHFMIYEETSGMLLFLGRVNP	PTLL	444
ZPI_MMU	406	<b>L</b> <b>S</b> <b>E</b> <b>I</b> <b>I</b> <b>A</b> <b>Y</b> <b>S</b> <b>M</b> <b>P</b> <b>P</b> <b>A</b>	IKVNRPEHFIIYEEMSRMLLFLGRVNP	PTVL	448
ZPI_RNO	394	<b>V</b> <b>S</b> <b>E</b> <b>I</b> <b>T</b> <b>A</b> <b>Y</b> <b>C</b> <b>M</b> <b>P</b> <b>P</b> <b>V</b>	IKVDRPPEHFIIYEEMSRMLLFLGRVNP	PTVL	436
ZPI_GGA	397	<b>G</b> <b>S</b> <b>E</b> <b>I</b> <b>I</b> <b>A</b> <b>F</b> <b>S</b> <b>V</b> <b>P</b> <b>P</b> <b>V</b>	LKVDREPELFMIFEETFKTLLFLGRVVD	PTET	439
ZPI_XTR	348	<b>G</b> <b>A</b> <b>E</b> <b>I</b> <b>I</b> <b>A</b> <b>Y</b> <b>S</b> <b>L</b> <b>P</b> <b>L</b> <b>T</b>	IRVNRPELFMIFEAYQSLLFLGRVMD	PTKL	390
ZPI1_FRU	358	<b>I</b> <b>I</b> <b>G</b> <b>I</b> <b>T</b> <b>P</b> <b>F</b> <b>S</b> <b>L</b> <b>P</b> <b>R</b> <b>T</b>	FTVNRPEFFFIYHEKTNCLMFMGRV	IDPT--	398
ZPI1_TNI	358	<b>A</b> <b>V</b> <b>G</b> <b>I</b> <b>T</b> <b>P</b> <b>F</b> <b>S</b> <b>L</b> <b>P</b> <b>R</b> <b>T</b>	FFVNRPEFFFIYHEKTNCLMFMGRV	IDPTKD	400
ZPI1_DRE	350	<b>G</b> <b>N</b> <b>L</b> <b>F</b> <b>M</b> <b>N</b> <b>P</b> - <b>L</b> <b>P</b> <b>P</b> <b>R</b>	LTFNRPELFVVYHEVTKCLLYIGRV	VDPTKN	391
ZPI2_DRE	354	<b>S</b> <b>V</b> <b>G</b> <b>I</b> <b>T</b> <b>A</b> <b>Y</b> <b>S</b> <b>L</b> <b>P</b> <b>D</b> <b>T</b>	FIINRPEFFFLYHEETASLLMFMGRV	IDPTLS	396
Spn_94a_FRU	352	<b>G</b> <b>A</b> <b>S</b> <b>T</b> <b>F</b> <b>S</b> <b>Y</b> - <b>P</b> <b>P</b> <b>P</b> <b>R</b>	LFINRPELFIIYHQTIGSVLEMGRV	TNP	392
Spn_94a_TNI	353	<b>G</b> <b>A</b> <b>S</b> <b>M</b> <b>F</b> <b>S</b> <b>D</b> - <b>P</b> <b>P</b> <b>P</b> <b>R</b>	LFINRPEVFLIYHOMSGIVLLIGRV	SDPTLQ	394

**Appendix 8.3.15: Alignment of  $\alpha_1$ -antitrypsin like sequences – Spn\_215c from *Fugu* and *Tetraodon*.** Gene specific features include an inhibitory RCL (red box). Conserved intron positions are indicated above the alignment. There is a unique intron 215c position in these fish specific group V2 serpins.

```
ALAT_HSA      1  -----EDPQG-DAAQKTD--TSHHDQDHPTFNKITPNL      30
Spn_215c_FRU  1  MNATRCVWILSITICVARGHVGNDIGQNQKEQDTSDNSTESLSLVTAAN  50
Spn_215c_TNI  1  -----MLLVARGHQG-DGAEKLEGQQSSAANS SAGVPLLLTAAN    37

ALAT_HSA      31  AEFATSLYRQLAHQSNS--TNIFFSPVSIATAFAMLSLGTKADTHDEILE  78
Spn_215c_FRU  51  REFAFRLYRSLAANPDSQGNIFFSVSVSVALAALAVGARGETHRQLFR  100
Spn_215c_TNI  38  REFAFRLYRSLAAQPD SRGNVFFSPLSVSVALAALAVGARGETHTQQLFR  87

ALAT_HSA      79  GLNFNLTEIPEAQIHEGFQELLRTLNQPDSQLQLTTGNGFLFSEGLKLV  128
Spn_215c_FRU  101  GLAFNSTWLSQTDVDFQAFQSLFEKTKKASNEVTS-EGTAVFMDNLFKPQP  149
Spn_215c_TNI  88  GLGLSNTSLSQAQVDQAFQSLFEQTRRTSSQVTR-EGTAVFVDHLFKAQP  136

ALAT_HSA      129  KPELEDVKKLYHSEAFVNFVGDTEEAKQINDYVEKGTQCKIIVDLVKELDR  178
Spn_215c_FRU  150  EFLDTLKKSYFADGFNVDFTKSSSEANTINKYVEEKTSGKIDKLVESLDP  199
Spn_215c_TNI  137  GELHTLKQSYFADGFVDFSKSSESTDTINKYVKEKTSGKIDKLVKDLDP  186

                      192a      215c!
                      |      |
ALAT_HSA      179  DTVFALVNYIFFKGRTERPFEVKDTEEEDEFHVDQVTTVKVPMMKRLGMFN  228
Spn_215c_FRU  200  TTVMYLLISYLYFKGRTERPEFDPDLLKEDELFMVDEKTKVPVQMMNIEKRFE  249
Spn_215c_TNI  187  STVMYLLISYLYFKGRWSEPEFDPDLLQEDVFTVDEETKVPVQMMNLERRFE  236

ALAT_HSA      229  IQHCKKLSWVLLMKYLGNATAIFPLPDEGKLQHLENELTHDIIITKFLEN  278
Spn_215c_FRU  250  TYRDQMFNTSVLHLFPNSSHSMLLLEDD--MSKLEN AISAAHVTKWLKW    297
Spn_215c_TNI  237  TYHDQTVNTSVLRLFPNSSHSMLLLEPEH--MAQLEQALSPAHI SKWLKW    284

                282b
                |
ALAT_HSA      279  EDRRSASLHLPKLSITGTYDKSVLGQLGITKVFNSGADLSGVTEEAPLK    328
Spn_215c_FRU  298  MKYRKYSVYIPKFSIKTSSYIKTVLTEMGMVDMFGDRADLSGIAEGQQLA  347
Spn_215c_TNI  285  MKSRTFNVYVPKFSIKTSASIKDVLTEMGMADMFGDRADLTGISEGGRLS  334

                331c?
                |
ALAT_HSA      329  LSKAVHKAVLTIDEKGTAAAGAMFLEAIPMSIPPE-V-KFNKPFVFLMIE  376
Spn_215c_FRU  348  VSEVVHQATLDVDEAGATAAATGIAITLFSYNYVPVLKFNRPFMVIITD  397
Spn_215c_TNI  335  VSEVVHQATLDVDEAGATAAATGIGITLFSFHHVPVLKFDRPFMVIITE  384

ALAT_HSA      377  QNTKSPLEFMGKVVNEITQK      394
Spn_215c_FRU  398  HSSDNILFMGKITNENI-      414
Spn_215c_TNI  385  HSTESI L L G K I T N P K I -      401
```

**Appendix 8.3.16: Alignment of  $\alpha_1$ -antitrypsin like sequences – Fru-Spn-17 and Tni-Spn-4 from *Fugu* and *Tetraodon* respectively.** Gene specific features include a non-inhibitory RCL (red box), and G-rich region from 108-165 (blue boxes). Conserved intron positions are indicated above the alignment.

ALAT_HSA	1	-----EDPQGDAAQKTDTSHHDDQDHPTFNKI	26
Fru-Spn-17	1	<b>MSLRCSFSCLTALTQLML</b> --- <b>CVAH</b> GTFGPRAEKEFAGPLPQEAHQHLNT	47
Tni-Spn-4	1	<b>MSLRRSFWCLTALTQLMLQTL</b> <b>CGAH</b> GTFGPTAQKDFARPGEASRHLNT	50
ALAT_HSA	27	TPNLAE <b>F</b> AFSLYRQLAH-----QSNSTNI <b>F</b> FPV <b>S</b> IATAFAMLS-- <b>LG</b>	67
Fru-Spn-17	48	SALNTLLAFEPYHGLASRVSTEPEAQQR-NILFS <b>P</b> LGLASAVVLLSRVSR	96
Tni-Spn-4	51	SALNTLLAFEPYQALAS---AEPEAPQ <b>Q</b> PNILFS <b>P</b> LGLASAAALLSRV <b>S</b> <b>G</b>	97
ALAT_HSA	68	TKAD <b>T</b> HDEILEG <b>I</b> NFNLTEIPEAQIHEGFQELLRTL <b>N</b> Q <b>P</b> DSQ-LQLTTGN	116
Fru-Spn-17	97	SESRSQALELLGLAANSTERSVEDAVSS <b>L</b> TDLLHNLTL <b>P</b> <b>E</b> <b>G</b> <b>R</b> <b>G</b> <b>G</b> <b>G</b> <b>R</b> <b>G</b> <b>A</b>	146
Tni-Spn-4	98	PERRSQALTLLGLAAGSPEQSVEDT <b>L</b> SALT <b>N</b> LLHNLTL <b>P</b> <b>E</b> <b>E</b> ----- <b>G</b> <b>E</b> <b>G</b> <b>A</b>	142
ALAT_HSA	117	GLFLSEG-----LKLVDK <b>F</b> LEDVKKLYHSEAF <b>T</b> -V <b>N</b> <b>F</b> GDTEEA	153
Fru-Spn-17	147	<b>G</b> SEAG <b>A</b> GT <b>T</b> AG <b>D</b> <b>G</b> <b>D</b> <b>G</b> <b>G</b> SDA <b>G</b> CRADAAEA <b>G</b> TH <b>A</b> <b>G</b> SOLKV <b>V</b> <b>S</b> <b>G</b> L <b>H</b> AG <b>G</b> <b>N</b> QSD	196
Tni-Spn-4	143	<b>G</b> AW <b>G</b> <b>G</b> <b>A</b> <b>G</b> -----D <b>N</b> <b>C</b> SSTDAEA <b>E</b> A <b>A</b> TH <b>A</b> <b>G</b> SOLKV <b>V</b> <b>S</b> <b>R</b> LQ <b>A</b> <b>D</b> <b>C</b> <b>K</b> <b>Q</b> <b>A</b> <b>D</b>	182
		192a	
ALAT_HSA	154	KKQ <b>I</b> NDY <b>V</b> EKG <b>I</b> Q <b>K</b> I <b>V</b> DLVKELDRD <b>T</b> VFALV <b>N</b> Y <b>I</b> FF <b>K</b> <b>G</b> <b>K</b> <b>W</b> ER <b>P</b> <b>E</b> <b>V</b> <b>K</b> <b>D</b> <b>T</b>	203
Fru-Spn-17	197	YQ <b>S</b> FLSE <b>G</b> <b>W</b> <b>S</b> <b>G</b> <b>F</b> NY <b>T</b> FD <b>T</b> LQ <b>K</b> DLES <b>S</b> DELE <b>L</b> N <b>M</b> Y <b>A</b> <b>F</b> <b>K</b> <b>G</b> --R <b>L</b> <b>P</b> <b>F</b> ERR <b>H</b> <b>T</b>	244
Tni-Spn-4	183	HQ <b>S</b> F <b>P</b> <b>S</b> <b>G</b> <b>N</b> Q <b>R</b> <b>G</b> --SS <b>N</b> AS <b>D</b> T <b>S</b> Q <b>V</b> SD <b>K</b> L <b>N</b> L <b>R</b> S <b>Y</b> A <b>F</b> <b>K</b> <b>G</b> --R <b>L</b> <b>P</b> <b>F</b> ERR <b>H</b> <b>T</b>	228
ALAT_HSA	204	EEED <b>F</b> HVDQ <b>V</b> TT <b>V</b> <b>K</b> <b>V</b> <b>P</b> <b>M</b> <b>M</b> <b>K</b> <b>R</b> L <b>G</b> <b>M</b> <b>F</b> <b>N</b> I <b>Q</b> H <b>C</b> <b>K</b> <b>K</b> <b>L</b> <b>S</b> <b>S</b> <b>W</b> <b>V</b> <b>L</b> <b>M</b> <b>K</b> <b>Y</b> <b>L</b> <b>G</b> <b>N</b> --A <b>T</b> <b>A</b> <b>I</b>	251
Fru-Spn-17	245	V <b>P</b> <b>R</b> <b>S</b> <b>F</b> <b>Q</b> L <b>N</b> <b>A</b> <b>T</b> <b>A</b> <b>S</b> <b>L</b> <b>E</b> <b>V</b> <b>A</b> <b>M</b> <b>M</b> <b>F</b> <b>R</b> <b>D</b> <b>D</b> <b>S</b> <b>S</b> <b>D</b> <b>V</b> <b>M</b> <b>L</b> <b>Y</b> <b>D</b> <b>T</b> <b>N</b> <b>C</b> <b>S</b> <b>A</b> <b>T</b> <b>V</b> <b>V</b> <b>Q</b> <b>L</b> <b>A</b> <b>Q</b> <b>S</b> <b>E</b> <b>R</b> <b>L</b> <b>A</b> <b>W</b> <b>L</b>	294
Tni-Spn-4	229	V <b>R</b> <b>R</b> <b>S</b> <b>F</b> <b>Q</b> L <b>N</b> <b>A</b> <b>T</b> <b>A</b> <b>S</b> <b>V</b> <b>E</b> <b>V</b> <b>A</b> <b>M</b> <b>M</b> <b>F</b> <b>R</b> <b>D</b> <b>D</b> <b>S</b> <b>S</b> <b>D</b> <b>V</b> <b>R</b> <b>M</b> <b>L</b> <b>Y</b> <b>D</b> <b>T</b> <b>N</b> <b>C</b> <b>S</b> <b>A</b> <b>T</b> <b>V</b> <b>V</b> <b>Q</b> <b>L</b> <b>A</b> <b>Q</b> <b>S</b> <b>E</b> <b>R</b> <b>L</b> <b>A</b> <b>W</b> <b>L</b>	278
		282b	
ALAT_HSA	252	FF <b>L</b> <b>P</b> <b>D</b> <b>E</b> <b>G</b> <b>K</b> <b>L</b> <b>Q</b> <b>H</b> <b>L</b> <b>E</b> <b>N</b> <b>E</b> <b>L</b> <b>T</b> <b>H</b> <b>D</b> <b>I</b> <b>I</b> <b>T</b> <b>K</b> <b>F</b> <b>L</b> <b>E</b> <b>N</b> <b>E</b> <b>D</b> <b>R</b> <b>R</b> <b>S</b> <b>A</b> <b>S</b> <b>L</b> <b>H</b> <b>L</b> <b>P</b> <b>K</b> <b>L</b> <b>S</b> <b>I</b> <b>T</b> <b>G</b> <b>T</b> <b>Y</b> <b>D</b> <b>L</b> <b>K</b> <b>S</b>	301
Fru-Spn-17	295	LL <b>L</b> <b>P</b> <b>K</b> -A <b>E</b> <b>L</b> <b>Q</b> <b>T</b> <b>L</b> <b>E</b> <b>G</b> <b>C</b> <b>L</b> <b>S</b> <b>S</b> <b>R</b> <b>R</b> <b>M</b> <b>S</b> <b>F</b> <b>W</b> <b>L</b> <b>S</b> <b>N</b> <b>L</b> <b>K</b> <b>P</b> <b>G</b> <b>R</b> <b>A</b> <b>E</b> <b>I</b> <b>L</b> <b>F</b> <b>P</b> <b>K</b> <b>F</b> <b>Q</b> <b>L</b> <b>R</b> <b>R</b> <b>S</b> <b>Y</b> <b>N</b> <b>V</b> <b>K</b> <b>N</b>	343
Tni-Spn-4	279	LL <b>L</b> <b>P</b> <b>R</b> -A <b>E</b> <b>L</b> <b>Q</b> <b>S</b> <b>L</b> <b>E</b> <b>G</b> <b>C</b> <b>L</b> <b>Y</b> <b>E</b> <b>G</b> <b>R</b> <b>M</b> <b>R</b> <b>F</b> <b>W</b> <b>L</b> <b>S</b> <b>N</b> <b>L</b> <b>K</b> <b>P</b> <b>G</b> <b>H</b> <b>A</b> <b>E</b> <b>I</b> <b>L</b> <b>F</b> <b>P</b> <b>K</b> <b>L</b> <b>Q</b> <b>L</b> <b>R</b> <b>R</b> <b>S</b> <b>Y</b> <b>N</b> <b>L</b> <b>E</b> <b>K</b>	327
		331c	
ALAT_HSA	302	V <b>L</b> <b>G</b> <b>Q</b> <b>L</b> <b>G</b> <b>I</b> <b>T</b> <b>K</b> <b>V</b> <b>S</b> <b>N</b> <b>G</b> <b>A</b> <b>D</b> <b>L</b> <b>S</b> <b>G</b> <b>V</b> <b>T</b> <b>E</b> <b>E</b> <b>A</b> <b>P</b> <b>L</b> <b>K</b> <b>L</b> <b>S</b> <b>K</b> <b>A</b> <b>V</b> <b>H</b> <b>K</b> <b>A</b> <b>V</b> <b>L</b> <b>T</b> <b>I</b> <b>D</b> <b>E</b> <b>K</b> <b>C</b> <b>T</b> <b>E</b> <b>A</b> <b>A</b> <b>G</b> <b>A</b> <b>M</b>	351
Fru-Spn-17	344	LL <b>K</b> <b>N</b> <b>S</b> <b>G</b> <b>A</b> <b>S</b> <b>S</b> <b>L</b> <b>F</b> <b>S</b> <b>D</b> <b>A</b> <b>P</b> <b>D</b> <b>F</b> <b>S</b> <b>G</b> <b>L</b> <b>S</b> <b>E</b> <b>K</b> <b>E</b> <b>T</b> <b>L</b> <b>R</b> <b>L</b> <b>V</b> <b>K</b> <b>A</b> <b>S</b> <b>Q</b> <b>E</b> <b>V</b> <b>M</b> <b>L</b> <b>E</b> <b>V</b> <b>E</b> <b>A</b> <b>K</b> <b>L</b> <b>E</b> <b>E</b> <b>G</b> <b>G</b> <b>G</b> -	392
Tni-Spn-4	328	LL <b>R</b> <b>S</b> <b>S</b> <b>G</b> <b>A</b> <b>S</b> <b>S</b> <b>L</b> <b>F</b> <b>S</b> <b>D</b> <b>L</b> <b>P</b> ---G <b>Q</b> <b>S</b> <b>E</b> <b>E</b> <b>K</b> <b>T</b> <b>P</b> <b>R</b> <b>L</b> <b>Q</b> <b>E</b> <b>A</b> <b>S</b> <b>H</b> <b>E</b> <b>V</b> <b>M</b> <b>L</b> <b>E</b> <b>V</b> <b>E</b> <b>A</b> <b>K</b> <b>L</b> <b>E</b> <b>E</b> <b>A</b> <b>G</b> <b>S</b> -	373
ALAT_HSA	352	<b>F</b> <b>L</b> <b>E</b> <b>A</b> <b>I</b> <b>P</b> <b>M</b> <b>S</b> <b>I</b> <b>P</b> <b>P</b> <b>E</b> <b>V</b> <b>K</b> <b>F</b> <b>N</b> <b>K</b> <b>P</b> <b>F</b> <b>V</b> <b>F</b> <b>L</b> <b>M</b> <b>I</b> <b>E</b> <b>Q</b> <b>N</b> <b>T</b> <b>K</b> <b>S</b> <b>P</b> <b>L</b> <b>E</b> <b>M</b> <b>G</b> <b>K</b> <b>V</b> <b>V</b> <b>N</b> <b>E</b> <b>T</b> <b>Q</b> <b>K</b> -	394
Fru-Spn-17	393	<b>Y</b> <b>D</b> <b>V</b> <b>P</b> <b>L</b> <b>D</b> <b>F</b> <b>S</b> <b>V</b> <b>P</b> <b>P</b> <b>R</b> <b>I</b> <b>T</b> <b>F</b> <b>N</b> <b>R</b> <b>P</b> <b>F</b> <b>V</b> <b>L</b> <b>L</b> <b>T</b> <b>Y</b> <b>D</b> <b>H</b> <b>V</b> <b>T</b> <b>G</b> <b>L</b> <b>V</b> <b>L</b> <b>L</b> <b>M</b> <b>G</b> <b>R</b> <b>I</b> <b>S</b> <b>D</b> <b>E</b> <b>A</b> <b>D</b> <b>V</b> -	435
Tni-Spn-4	374	<b>Y</b> <b>D</b> <b>V</b> <b>P</b> <b>L</b> <b>D</b> <b>F</b> <b>S</b> <b>M</b> <b>P</b> <b>P</b> <b>R</b> <b>I</b> <b>T</b> <b>F</b> <b>N</b> <b>R</b> <b>P</b> <b>F</b> <b>I</b> <b>L</b> <b>M</b> <b>V</b> <b>Y</b> <b>D</b> <b>R</b> <b>V</b> <b>T</b> <b>G</b> <b>L</b> <b>V</b> <b>L</b> <b>L</b> <b>L</b> <b>G</b> <b>R</b> <b>I</b> <b>S</b> <b>D</b> <b>E</b> <b>A</b> <b>H</b> <b>L</b> -	416



Appendix 8.3.17: Alignment of  $\alpha_1$ -antitrypsin like serpins from *Gallus gallus*. Gene specific features include an inhibitory RCL (red box). Conserved intron positions are indicated above the alignment.

ALAT_HSA	-	-----	-
ALAT_GGA	1	-----MPEQAVGNR-----	9
Gga-Spn-12	1	<b>MKYHLYLCTFLAGLCVVTQCH</b> -----RGKACHEVNNTN--	33
Gga-Spn-13	1	<b>MKPTFSLCFLLAGLYSVAQCH</b> -----QRPRYHNKQDNSKG	35
Gga-Spn-14	1	<b>MQVQDPSESNQLKSMKAILLLCLLL</b> AALYPAVPSMRQTDHHNEEPKAT-	49
Gga-Spn-15	1	<b>MKSALYLCLFLTGLQVQAL</b> -----PKPNHSNKHKEERP--	33
Gga-Spn-16	1	<b>MKTVFYICLLLAGLHAFAYGQ</b> -----LTASHHNGHNPN--	33
ZPI_GGA	1	<b>MKIRIYLLLLCELCFEISKAD</b> -----IKPKSPKKDKRLNF	35
ALAT_HSA	1	-----EDPQGDAAQKTDTSHHQDHPNFKITPNLAEF <b>AFSLYRQ</b>	40
ALAT_GGA	10	-----VCQFACCFYKE	20
Gga-Spn-12	34	-----LRDGTENLLKHN---CQKQGSTYTD <b>FAFRFYRQ</b>	63
Gga-Spn-13	36	A-----YYWGSSSHREGVFPNKNKTFVKVHNSAD <b>FALSFYKL</b>	73
Gga-Spn-14	50	-----HLHEQHPHEEDSLAFCQHIIP <b>SNTDFAFRFYRQ</b>	82
Gga-Spn-15	34	-----HSLGDHPHVEHKNLAHMKIAP <b>SNAEFAFRFYRQ</b>	66
Gga-Spn-16	34	-----EPKDHMHNAEAAACLKLV <b>PNNADFAFKFLNE</b>	65
ZPI_GGA	36	LGRNKNVSI <b>SEWHQHKNDHKPLEEQSFEELTLHN</b> FTEKTAN <b>GENLYRK</b>	85
ALAT_HSA	41	<b>LAHQSNSTN</b> IFFSPVSI <b>ATAFAMLSLGRKAD</b> THDEI <b>LEGLNFN</b> - <b>LTEIPE</b>	89
ALAT_GGA	21	<b>ISSHENS</b> GNIFFSP <b>LSISTAFAM</b> LT <b>LGARSDTLAQILRVLHFN</b> - <b>PRAISE</b>	69
Gga-Spn-12	64	<b>AISKEADKN</b> IFFSP <b>ISISTFAM</b> LAV <b>GA</b> KST <b>TLTQIF</b> EGLG <b>FDNLTETRI</b>	113
Gga-Spn-13	74	<b>VASEATDQ</b> NIFFSP <b>ISISTSLAM</b> LAL <b>GAKSVTLTQI</b> LEGL <b>LAFN</b> - <b>LKKTQD</b>	122
Gga-Spn-14	83	<b>ATVQAPGKN</b> IFFSPV <b>SVA</b> FALLAL <b>G</b> S <b>RAATQAQ</b> LL <b>LEGLAFN</b> - <b>LINNRE</b>	131
Gga-Spn-15	67	<b>VTEAGGNKN</b> IFFSP <b>LSLSTAFAM</b> LS <b>LGARSNTLSQ</b> LHK <b>CLTFN</b> - <b>LTEMEE</b>	115
Gga-Spn-16	66	<b>VAQEAPNKN</b> IFFSPV <b>SISA</b> FAM <b>LALGARSITKTQI</b> LEGL <b>LAFN</b> - <b>LTEIQE</b>	114
ZPI_GGA	86	<b>IAMKLDN</b> - <b>NIIT</b> SP <b>LSVTTLMATYLLA</b> EGET <b>HROIAKALNLHSILKDRDR</b>	134
ALAT_HSA	90	<b>AQTHEGFQ</b> ELL <b>R</b> TL <b>NQ</b> PD <b>SQ</b> L <b>QLTTG</b> N <b>GLFL</b> SE <b>GLKLV</b> DK <b>FLEDV</b> KK <b>LYH</b>	139
ALAT_GGA	70	<b>NEIHEGYR</b> QL <b>IQ</b> MIN <b>RKNEGLQ</b> LN <b>MGNVLFV</b> LD <b>R</b> L <b>KPQQR</b> FL <b>NSL</b> REF <b>YE</b>	119
Gga-Spn-12	114	<b>HDIHESF</b> YK <b>V</b> LAV <b>LNCTD</b> VN <b>ITL</b> N <b>IGNAF</b> FP <b>AI</b> GY <b>EPQ</b> ET <b>FLQ</b> N <b>VQ</b> Q <b>FYD</b>	163
Gga-Spn-13	123	<b>QEIHEGF</b> C <b>Q</b> LL <b>HMLNR</b> SD <b>SDLHLS</b> LG <b>N</b> TL <b>F</b> TE <b>ETL</b> K <b>P</b> L <b>Q</b> K <b>F</b> LD <b>DAK</b> S <b>FYQ</b>	172
Gga-Spn-14	132	<b>EEIHRG</b> F <b>HH</b> LL <b>LL</b> NR <b>PG</b> S <b>Q</b> VEL <b>SMGN</b> TL <b>F</b> MD <b>KHLK</b> P <b>LT</b> T <b>FL</b> K <b>DI</b> KK <b>LYK</b>	181
Gga-Spn-15	116	<b>QEIHEGF</b> Q <b>R</b> LL <b>Q</b> LL <b>ND</b> S <b>Q</b> RD <b>IQ</b> LN <b>MGN</b> TL <b>F</b> ID <b>ER</b> L <b>K</b> L <b>Q</b> Q <b>K</b> F <b>L</b> DD <b>V</b> T <b>N</b> F <b>Y</b>	165
Gga-Spn-16	115	<b>KEIHEGF</b> H <b>N</b> L <b>M</b> H <b>M</b> LS <b>HP</b> ES <b>GV</b> Q <b>LN</b> M <b>GNA</b> I <b>FL</b> T <b>K</b> KL <b>K</b> P <b>L</b> KK <b>F</b> LD <b>DAK</b> P <b>LYQ</b>	164
ZPI_GGA	135	<b>HYLPAL</b> E <b>K</b> Q <b>L</b> KN <b>IT</b> T <b>N</b> - <b>E</b> ELL <b>FV</b> Q <b>GL</b> LS <b>F</b> I <b>Q</b> K <b>D</b> F <b>T</b> V <b>R</b> E <b>A</b> FL <b>N</b> LS <b>K</b> Q <b>Y</b> F <b>D</b>	183

ALAT_HSA	140	SEFTVNEGDTEEAQKQINDYVEKGTQKIVDLVKEIDRDTVFALVNYIF	189
ALAT_GGA	120	GEIYPMNEKRSDDQAQTKINDYVAERTNGKIKDLINLDPLEILLISYTY	169
Gga-Spn-12	164	ADFFSTDFHKPEEAKLQINNYVKEKTQKIPELIARLDANTILVLVNYTY	213
Gga-Spn-13	173	SEVLSADFNNSGAENQINSYIEEKTNGKIVKLVENLDPLEAMVLVNYVF	222
Gga-Spn-14	182	GEIISSNFQNSTEAKKEINEHMKNKTHGKINQILKDLDPNSLMVLVNYTY	231
Gga-Spn-15	166	SEAVSMDFQNSEHAKKEINNYIKAKTHGKFLDLLDSIGKDVVMILVNYVY	215
Gga-Spn-16	165	LEVLATDFNNPTEAEKEINDYTEKKTQKITNLVKEIDPQTVMLLASFVF	214
ZPI_GGA	184	MEFLCVDFQNSTQAKFVINQNIKQRTKGISELFEEDRHSKLLLLDYIF	233

## 192a

ALAT_HSA	190	FKGKWERPFQVVDTEEEDEFHVDQVTTVKVPMMKRRLGMFNIQHCCKLSSWV	239
ALAT_GGA	170	FNAEWEKPFNPQYTKKEKFFVDGNKAVEVPMFMFGIGAFKHGYDEQLSSTV	219
Gga-Spn-12	214	FKAWEKPFDSSENTYEDDFVSANERVRVMMQOHENDYRYYDQDLSCQV	263
Gga-Spn-13	223	FKAHWEKPFSDSYTKKEDFFVDKKTQKIVDMMYRKYGYRNYFDEELSCWL	272
Gga-Spn-14	232	FKAYWENPFNTKGTTHKDYFYVNEKTLVEVKMMIRDSFYDIYSDDKLSCKV	281
Gga-Spn-15	216	FKGYWEEPFESYNTRDDDFVDAKHSVKVMMYKNTYNNIHRDEQLSCWV	265
Gga-Spn-16	215	FRGNWEKPFKPEPTEEREFFVDAETTQKIVPMFCRIGTFDLYFDKDLPCV	264
ZPI_GGA	234	FKGKWLYPENSEFTEIETFHINKYRSVQVPMFMFKSDKVNSTYDENLRQNV	283

## 282b

ALAT_HSA	240	LLMKYLGNAIAIFFLPDEGK-LQHLENELTHDIITKFLNEDRR-SASLH	287
ALAT_GGA	220	VQMDYKGGASAFFVLPDQGR-MRKLEKKLSCERMARWRTLVSKSNSVNLV	268
Gga-Spn-12	264	VELPYKGTAAQALLILPDDGK-MKQVENALSKETVCNWFQSKFETR-RLHLY	311
Gga-Spn-13	273	VQIPYNGNAALFVLPDEGK-MKQVEDALLKRTVSKTEKLLQHR-KIHLH	320
Gga-Spn-14	282	VRIPYKGNVSALFILPNEGK-LKWLEDGLKKTQVSKTEKSLERR-RMEVH	329
Gga-Spn-15	266	VEIPYRGNAAAFFVLPDEGS-MNQVEDALLQDTVSNWSQSLEGR-SIDLY	313
Gga-Spn-16	265	VRLHYNGSATAFLILPAKGG-MKQLEPTLDKERVKKWSDFHLFKS-KIQLY	312
ZPI_GGA	284	IKLPHYKGAAYMLIVIPEKGEDYVSLDHLTMELVESWLANMKSR-NMDIS	332

## 331c

ALAT_HSA	288	LPKLSITGT YDLKSVLGQLGITKVF SNGADLSGVTEEAPLKLKSKAVHKAV	337
ALAT_GGA	269	LPKFTLHG RYNLKNILYKMGIMDLFTDKADLSGITGQPQHRISQATHQAV	318
Gga-Spn-12	312	LPRISISGSYDVKDLFMBMGITDVFSSNADLSGISGSRTLQVVSQATHKAL	361
Gga-Spn-13	321	IPKLSISGT YDVKKIVREVG IIDLFTAQADLSGITEDPGLMVSKVIHRAV	370
Gga-Spn-14	330	IPKVSISGT YDLKMMAMNLGVTDVFSQADLSGITGKSDLKVSRAIHKGL	379
Gga-Spn-15	314	LPKFSISGSYDVKKLEFLKMGVTDVFSNADFSGVAKNTLLKVSRAIHKAK	363
Gga-Spn-16	313	FPKFSISGT YEITNLSKMGIVDVFTNQADLSGISGVEPKVSKVIHKAA	362
ZPI_GGA	333	FPKFKLEQKYKMKKLLYALG IKNLFA RTADLSHLTDQKHLTVSQVVKAV	382

---

```

ALAT_HSA 338 LTIDEKGTEAAGAMFLEAIPMSIPPE----VKFNKPFVFLMIEQNTKSPL 383
ALAT_GGA 319 VKVDETGTEAAAATGMEIVPMSVPPV----IRMNRPFLLVITLR--ENIL 362
Gga-Spn-12 362 LAVDETGTEAAGATAILLSKFSLPHI---TTKFNRPFIVLIFDKATSTTL 408
Gga-Spn-13 371 LNVHENGGTEAAGVTVTEITWRSGDFPRPFRVRFNRPFLLMILDKYAHTIL 420
Gga-Spn-14 380 LDIHENGGTEAAAVTGTEFAPHSVPPV----IKFNRPFLLLIVDQYTESIL 425
Gga-Spn-15 364 LNVNENGGTEAAAVTMVEMKVSAMID-PLEIKFNRPFVMMIFDKITNSIL 412
Gga-Spn-16 363 LDVDERGTEASATAATPKIMALSLAP---IIEFNRPFMLLIFDRDTNSTL 409
ZPI_GGA 383 IEVDEKGTEAAAATGSELLAFSVPPV----LKVDRPFLEFMIFEETFKILL 428

```

```

ALAT_HSA 384 FMGKVVNPTQK----- 394
ALAT_GGA 363 FMGKIVNPLEKD---- 374
Gga-Spn-12 409 FMGKIVDPTMK----- 419
Gga-Spn-13 421 FIGKIVNPLKNN---- 432
Gga-Spn-14 426 FIGKIVNPLKND---- 437
Gga-Spn-15 413 FMGKVVNPVAKED--- 425
Gga-Spn-16 410 FIGKIANPTTTSRTEI 425
ZPI_GGA 429 FIGRVVDPTET----- 439

```



**Appendix 8.3.18: Alignment of  $\alpha_1$ -antitrypsin like serpins from *Xenopus tropicalis*.** Gene specific features include an inhibitory RCL (red box). Conserved intron positions are indicated above the alignment.

AlAT_HSA	1	-----EDPQGDAAQKTDI	13
AlAT_XTR	1	<b>MRA-FLIVSLALLCAGVLAD</b> HDGQTKHGKDHNGHDHGDHDDHDDHHHGK	49
Xtr-Spn-9	1	MRAYFLFVSI <del>LLCAV</del> VFGDHE-----KHHESDHKEHGDHGNVHHHSD	44
Xtr-Spn-8	1	MRAYFLFVSI <del>LLCAV</del> VFGDHE-----KHHDSDHKEHGDHGNVHHHSDK	44
Xtr-Spn-10	1	---MRGLPYILFFITCIFAS-----HNDN	21
EP45_XTR	1	---MYLFVYLSLFIALTLAS-----VTDI	21
Xtr-Spn-12	1	---MRGICLVLLIAVIVESD-----NHDD	21
Xtr-Spn-13	-	-----	-
ZPI_XTR	1	-----MNSSLP	6
AlAT_HSA	14	SHHDQDHP-----TFN-----KTIIPNLAEIF	33
AlAT_XTR	50	GKHHDHKH-----HHHAGEDMACHKIAPSNISQF	77
Xtr-Spn-9	45	KHTDDSHG-----DHHHDESMPCCLKIAPYNANF	72
Xtr-Spn-8	45	THTDDSHG-----DHHHDESMPCCLKIAPYNANF	72
Xtr-Spn-10	22	SHQEHA-----NGHAGGKIAQEALGSANIDF	47
EP45_XTR	22	SLNKKQGNKQQHHHDNPKHCHQKDKQDQTWKAEGKLTkdkeVILSEENYDF	71
Xtr-Spn-12	22	GISEDN-----KELSE-EELKEILAQLNMHF	46
Xtr-Spn-13	1	-MKEEH-----HEHKI-REAQETIIQANRHF	24
ZPI_XTR	7	SGNEEQIP-----TVLSFANVVSQMSDF	29
AlAT_HSA	34	<b>AFSLYRQLAHQ</b> -----SNSTNIFFSPVSIATAFAMLSLGTKADTHDEIL	77
AlAT_XTR	78	<b>AFKLFROVVAD</b> -----HPSENIFFSPVSI <del>STALAMLSL</del> GARADTLN <del>QII</del>	121
Xtr-Spn-9	73	<b>GFSLYRQLAAD</b> -----HPTENIFFSPVSI <del>STVFAMLSL</del> GARSNTLN <del>QII</del>	116
Xtr-Spn-8	73	<b>GFSLYRQLAAD</b> -----HPTENIFFSPASISTVFAMLSL <del>GARSNTLN</del> <del>QII</del>	116
Xtr-Spn-10	48	<b>ALNLYKHLVTKTQAEKE</b> STQKNIVFSPLSILTAFSMLLLGAKSESHQ <del>QIL</del>	97
EP45_XTR	72	<b>TENLFNELSAECKR</b> ---SPKONIFFSPISISAAFYMLALGAKSKTHQ <del>QIL</del>	118
Xtr-Spn-12	47	<b>AVNIYKQVASSALKEKN</b> SEPKNIFFSPVSI <del>STSLAMLAL</del> GAKAETR <del>HQIL</del>	96
Xtr-Spn-13	25	<b>AINMFKHTAS</b> -----ESPKNIVFSPVSIYAAFAMLSIGARSKTERGIL	67
ZPI_XTR	30	<b>GFNLYRKTANK</b> -----HDNNIFFSPFSVSLGLSSILLGTRGNTYD <del>QLL</del>	72
AlAT_HSA	78	<b>EGLNFN-LTEIPEAQIHEGF</b> QELLRTLNQDPSQLQLTTGNGLFLSEGLKL	126
AlAT_XTR	122	<b>EGLNFN-NKITEEEEIHNGF</b> QHLLHMLNDPDRELQ <del>LN</del> SGNALFIDNNV <del>KL</del>	170
Xtr-Spn-9	117	<b>EGLKFN-RSELTEEEMHKG</b> FQHLLHMLNDPNSKVQ <del>LN</del> SGNALFIDKDL <del>QL</del>	165
Xtr-Spn-8	117	<b>EGLSFN-RSELTEEEMHKG</b> FQHLLHMLNDPNSKVQ <del>LN</del> SGNALFIDKDL <del>KL</del>	165
Xtr-Spn-10	98	<b>SGLSLN-QTQVPEEDMHEA</b> FEHLLQVLNRPKSDLOVKIGNAVFVEDTL <del>KI</del>	146
EP45_XTR	119	<b>QGLGFN-KKLNESQVHEA</b> FKGLIEDLNPMKDHQFTIGNALFVEQT <del>VNI</del>	167
Xtr-Spn-12	97	<b>NSLAPK-ETPTQEVKIHNA</b> FKHLLOTLNKP <del>KKDL</del> KTFVGNAAFVEEEL <del>KF</del>	145
Xtr-Spn-13	68	<b>ESLSFN-QTHYPD-QIHPG</b> EKDFILLALNKP <del>KN</del> LOVSTGNVLFVKDK <del>LET</del>	115
ZPI_XTR	73	<b>HGLNYPFKDQENPYLLPE</b> LLKTIKEKIAKNEELV <del>LN</del> IGSLSEFHET <del>F</del> SM	122

AlAT_HSA	127	VDKFLEDVVKLYHSEAFVNFVGDTEEAQKQINDYVEKGTQGGKIVDLVKEL	176
AlAT_XTR	171	IQQFIDDVKKYYESEAFSTDFNNAEEAKKQINSYVEKQTHGKIVDLLSSV	220
Xtr-Spn-9	166	IQKFVEDSKQFYEAEIFSTDFHNTTEEAQKQINTYABNKTKGKITDLLSSV	215
Xtr-Spn-8	166	IQKFVEDSKQFYKAEIIFSTDFHNTTEEAQKQINTYABNKTKGKITDLLSSV	215
Xtr-Spn-10	147	LDSFVQEIEHHYHAEIIFPSHFKNPAEAEKQINDFVNNKTEGRIQELVKDL	196
EP45_XTR	168	LRGFEENVKHYQAAIIFPINFRRDPNAKKQLNNYVKDKTHGAIQEMVRDL	217
Xtr-Spn-12	146	LKSFAHEAERYQADVFSTNFKNPKDAEQINNYVNNQINGKIKELVRGP	195
Xtr-Spn-13	116	LKSFLQKVKHHYQAEIIFTNFKNPKAEKQINDYVKNKTINGKIEELVQDL	165
ZPI_XTR	123	KDEFVNLTKKYFDMEYELDFHSS-KAKNEINAYVEKLTGKGLISNFYDFI	171

## 192a

AlAT_HSA	177	DRDITVFALVNYIFFKGGWERPFVVKDTEEDFHVDQVITVQVPMMKRLGM	226
AlAT_XTR	221	DKNVAVLYLLINYIFFRGGWEKPFEEKFTQDGIHVDENTNVTVPMMRRNGM	270
Xtr-Spn-9	216	DEKTIILVLLINYIIFYGWEKHFEEKWTKDGIHVDENTNVTVPMMHRNGM	265
Xtr-Spn-8	216	DEKTIILVLLINYIIFYGWEKHFEEKWTKDGIHVDENTNVTVPMMHRNGM	265
Xtr-Spn-10	197	SEATKLVVINFILFNAEQNPFSFFTHSRQFSVDENTTVEVQMMSKTDL	246
EP45_XTR	218	DANTEMLVLYNYVLEKGEWADTFNPSLTQKSIKSVDKNTKVTVQMMKRFGL	267
Xtr-Spn-12	196	SMDTKLLLVNYILEKGEWESPFSPDFIRLSVFSIDNRTKVEVKMMSRMGR	245
Xtr-Spn-13	166	DIETQLLVINYILEKGEWESPFSPSTHQSKFSIDNATVEVPMMSRTGI	215
ZPI_XTR	172	DPQTKLLLLDYIFFKGGQYFPNPALTEVDSEFIDKYNSVTVPMMYKTDK	221

AlAT_HSA	227	FNIQHCKKLSWVLLMKYLGNAATAIFFLDP-EGKLRQVLENEELTHDITKFL	275
AlAT_XTR	271	YNVAFDEKLGCTVVQIPYKGNATALEFILPD-EGKLRQVEEAELEKSTIMSW	319
Xtr-Spn-9	266	YNVAFDEKLGCTVVQIPYKGNATALEFILPD-EGKLRQVEEAELEKAVVKS	314
Xtr-Spn-8	266	YNVAFDEKLGCTVVQIPYKGNATALEFILPD-EGKLRQVEEAELEKAVVKS	314
Xtr-Spn-10	247	YQFYKDEKIPCSVLQLPYKNNASMLLIVPE-LGKIQEVEEALSVEITLKR	295
EP45_XTR	268	YKTYRDEDYCNCKIIEIPLKNDASMLLIVPQ-LGTIQEL--VLTPKLVTHW	314
Xtr-Spn-12	246	FDIFLDNELPCTVLKLPYIDDALMLLIMPE-LGKIQEVEEALSVDITLRW	294
Xtr-Spn-13	216	YNIYEDNKIPCIIVFQLPYKDNATMLLIVPK-LGKIQEVEEALSDETITRW	264
ZPI_XTR	222	VASVFDKDLSCITVFKLPIYRGNAMMLITKPEKEGDFGILEDHILTKEIN	271

## 282b

AlAT_HSA	276	LENEDRRSASLHLPKLSITGTIDIKSVLQGLGITKVFNSGADLSGVTEEA	325
AlAT_XTR	320	KKQFRYQSIELTIPKFSIMATLDLIEELKKFGVTDVFSQADLSGIVEGT	369
Xtr-Spn-9	315	KKLFRKRFVHLTLPKLSISATTDLVKELSKLGVTDVFSQNSDLGIVDVT	364
Xtr-Spn-8	315	KKLFRKQHRKAKPLWLNKSVIVEVGKKKRAFRAFKLAGTAETFIKYKEAN	364
Xtr-Spn-10	296	TSSAEKSFFELFLPKFSISSSLKLLKIDITDMGMGIIFTDAADFSGISENS	345
EP45_XTR	315	YESLTNSFVDLYMPIFSSISGKIVLKDITLHKMGISDIFTDKADLTGISQQT	364
Xtr-Spn-12	295	RNSTSKKFLRLYMPKFSISSSLKLLKEILSDVGMNSIFSDQADFSGIVEDG	344
Xtr-Spn-13	265	ESSGAKRRIEINMPKFSISSSLKLLKILSDMGMSIIFTDEADFSGIIVEDV	314
ZPI_XTR	272	QAKMQSRKTDIFFPKFKLDQKYKLLKSSILNELGIKELFTGKANLITDLTEER	321

		331c				
AlAT_HSA	326	PLKLSKAVHKAVLTID	---	EKGTEAAGAMF	LEAIPMSIPPEVKFNKPF	370
AlAT_XTR	370	PLKVS KAVHKAGLSVD	---	ETGTEAAAATA	FEIMPMMIPPHILLENRAF	414
Xtr-Spn-9	365	PLKVS KAVHKVLSID	---	ETGTEAAGVTV	MEIVPTMLPPRIEYNRPF	409
Xtr-Spn-8	365	-KACKKAINQAKIEMERDIAARR	---	RT EAAGATV	MEIMP NMAPLLITFNRPF	412
Xtr-Spn-10	346	RLKLSKVVHKAVLNVA	---	ENGTEAAAASA	VEGVLTSLMVQFVVVKPF	390
EP45_XTR	365	KLKVS MASHNAVLNVN	---	EFGTEAVGVTS	AQAIP TTSFPPFQIDSPF	409
Xtr-Spn-12	345	KLTL SKVIHKAVLDVD	---	EKGTEAAAATA	VNVIRYSLLKSQKVD RPF	389
Xtr-Spn-13	315	NLKVS KVIHKAILNVN	---	EEGTEAAAATAT	IETVPIMLYPSYSVDRPF	360
ZPI_XTR	322	NLMLTEITQQAMIEVD	---	ERGTEAAAVAG	AEI IAYS LPLTIRVNRPF	366
AlAT_HSA	371	VFLMIEQNTKSP	LFMGKVVNPTQK-	394		
AlAT_XTR	415	VVLIYDPIPKS	SILFVAKVVNPKN--	437		
Xtr-Spn-9	410	VLM IYEPTLRAN	LFMGRVMNPKE--	432		
Xtr-Spn-8	413	VLM IYEPTLRAN	LFMGRVMNP KK--	435		
Xtr-Spn-10	391	ITL ICSQEPYS	SILFMSRVIDPTEK-	414		
EP45_XTR	410	LVLIYSRTLGS	QLFMGKIMDPTNAK	434		
Xtr-Spn-12	390	LVV ICSKQTDI	ILFMGRIVNPTTEK-	413		
Xtr-Spn-13	361	LALLYCKDTKT	ILFMSRVVNPLEK-	384		
ZPI_XTR	367	LFMIFEEAYQS	LLFLGRVMDPTKL-	390		

**Appendix 8.3.19: Alignment of  $\alpha_1$ -antitrypsin like serpins from *Danio rerio*.** Gene specific features include an inhibitory RCL (red box). Conserved intron positions are indicated above the alignment.

AlAT_HSA	1	-----EDPQGDAAQKTDTSHHDDQDHP-	21
AlAT_DRE	1	MRGNIFCCAIAALLVATAWAAPHDGHEGHDHGSHTADHHHLLHHGKDEPH	50
Dre-Spn-8	1	MWGNIIYCCAIAALLVATAWAAPHDGHVGHHDHGSHTADHHHLLHHGKDEPH	50
Dre-Spn-9	-	-----	-
Dre-Spn-10	-	-----	-
Dre-Spn-12	-	-----	-
Dre-Spn-11	-	-----	-
ZPI1_DRE	1	-----MKMGFFTLIIASLLSVSV	19
ZPI2_DRE	1	-----MEFRLLLVIACFLCSAE	19
AlAT_HSA	22	-----TFNKITPNLAEFAFSLYRQLAHQSN--STNIFFSPPVS----IA	58
AlAT_DRE	51	P SHKGV DACHLLAPHNADF AFSLHKK LASNPD AQGKN IFFSPVG----IS	96
Dre-Spn-8	51	P SHKGV DACHLLAPHNADF AFSLYKK LASNPD GQGKN IFFSPVG----IS	96
Dre-Spn-9	1	-----MNNDFAFHLYKRLIESPDYQSKN IFFSPFS----VS	32
Dre-Spn-10	1	-----MNNDFAFHLYKRLIESPDYQSKN IFFSPFS----VS	32
Dre-Spn-12	1	-----MNNDFAFHLYKRVIELPDYQSKN IFFSPFS----VS	32
Dre-Spn-11	1	-----MNNDFAFHLYKRLIESPDYQSKN IFFSPFKMSNSVS	36
ZPI1_DRE	20	LG--QTTDVEELAIKNADFATRLYSKIAS SSD---DNVAVSTL G----AT	60
ZPI2_DRE	20	HEELRTPDISDLAFRNIDFAINLYRKISSLHD---RNVVFSPLS----VS	62
AlAT_HSA	59	TAFAMLSLGTKADITHDEILEGLNFNLTEIPEAQIHEGFQELLRTL NQPDS	108
AlAT_DRE	97	MALSLLAVGAKGSTLSQIYSGLGYS--ALTPEQVNEG YEHLLHMLGHSQD	144
Dre-Spn-8	97	MALSLLAVGAKASTLSQIYSGLGYS--ALTPEQVNEG YEHLLHMLGHSQD	144
Dre-Spn-9	33	MALSELSL GAGGDTKQQLSGIGYSSAIFSTEEMHQLFHS LLEDIGN-RT	81
Dre-Spn-10	33	MALSELSL GAGGDTKQQLSGIGYNSAIFSTEEMHQLFHS LLEEIGN-RT	81
Dre-Spn-12	33	MALSELSL GAGGDTKQQLSGIGYNSTIFSTEEMHQLFHS LLEDISN-RT	81
Dre-Spn-11	37	MALSELSL GAGGDTKQQLSGIGYNSTIFSTEEMHQLFHS LLEDIGN-RT	85
ZPI1_DRE	61	LALATLAAGAGGATQSELLQIGIGVD--SMVKDGEQERIQN ILQQLRE-DA	107
ZPI2_DRE	63	TCFSALLLAAQGSTRTIELKGLNLE--ALD-GGDSRRVPEL FQQLHQ-NI	108
AlAT_HSA	109	QLQLTTGNGLFLSEGLKLVDFLEDVKKLYHSEAFVNF GDTEEAKKQIN	158
AlAT_DRE	145	AMQLEAGAGVAIRDGFKVVDQFLKDAQHYYNSEAFGVDF SKPEIAAAEIN	194
Dre-Spn-8	145	AMQLEAGAGVAIRDGFKVVDQFLKDAQHYYNSEAFGVDF SKPEIAAAEIN	194
Dre-Spn-9	82	GVDIDVGTALYASDRFKPHSKFLEDMKEFYHSDGFTVDF SVKET-VDQIN	130
Dre-Spn-10	82	EVDIDVGTALYASDRFKPHSKFLEDMKEFYHSDGFTVDF SVKET-VDQIN	130
Dre-Spn-12	82	GVDIDVGTALYASDRFKPHSKFLEDMKEFYHSDGFTVDF RVKET-VDQIN	130
Dre-Spn-11	86	EVDIDVGTALYASDRFKPHSKFLEDMKEFYHSDGFTVDF RVKET-VDQIN	134
ZPI1_DRE	108	AQIPATG--LFIKQDVKADDSFSNQVKQYYNADVQNVNYANGQQAKGSIN	155
ZPI2_DRE	109	SLQMEQGTALFLDQHFHLQTNFSQQIQRFNAEVLRVDF SKPAVCRSLIN	158

		192a			
AlAT_HSA	159	DYVEKGTQGKIIVDLVKEILD	DRDITVFALVNYIFFKGGKWERPF	FEVKDIEEEDF	208
AlAT_DRE	195	KFIARKTHDKITNMVKDLDADT	VMMMLINMYFRGKWEKQF	DAKLTHKADF	244
Dre-Spn-8	195	KFIARKTHDKITNMVKDLDADT	VMMMLINMYFRGKWEKPF	DAKLTHKADF	244
Dre-Spn-9	131	KYVEEKTHGKINQAVDDLDADT	FMVLLTYTYFKGKWDKPF	NPKTISESTF	180
Dre-Spn-10	131	KYVEEKTHGKISQAVDNLEKDT	ILMFLTYTYFKGKWDKPF	KPETISESTF	180
Dre-Spn-12	131	NYAKKKTQGINQAVDNLEDDT	ILMFLTYTYFKGKWDKPF	KPEKIRESTF	180
Dre-Spn-11	135	KYVEEKTHGKINQAVDDLEEDT	ILMFLTYTYFKGKWDKPF	NPDTISESKF	184
ZPI1_DRE	156	DYVRGRTEGKVRDVVENVDPQ	SMAILISAFFTGQWLQPF	NATFTQEDRF	205
ZPI2_DRE	159	EFVSRKTGRKVLLEMLESVEPL	TQMMLNTIFYKGDWERPF	NPNNIEKSRF	208

AlAT_HSA	209	HVDQVITVQVPMMKRLGMFNIQ	HCKKISSWVLLMKVYLG	NATAIFFLDP--	256
AlAT_DRE	245	KVDQDITVQVDMMKRTGRYDI	YQDPVNQTTVLMVPYK	GNTSMLIVLDP--	292
Dre-Spn-8	245	KVDQDITVQVDMMKRTGRYDI	YQDPVNQTTVMMVPYK	GNTSMMIVLDP--	292
Dre-Spn-9	181	HIDDKITVQVQMMHQYERLKV	VYDAEELSTKVLCLDY	NDSFSMFLAVPDVH	230
Dre-Spn-10	181	YIDDKITVQVQMMHQYEHLKV	VYDVEELFTKVLCLDY	NDSFSMILAVPDVY	230
Dre-Spn-12	181	HIDDKITVQVQMMHQYERLKV	FYDAEELSTKVLCLDY	KDSFSMFLAVPDDK	230
Dre-Spn-11	185	NIDDKITVQVQMMHQYECLK	VYDVEELFSKVLCLDY	NDSFSMFLAVPDVH	234
ZPI1_DRE	206	YVNKYNTVQVPMMLRSGKY	YYLAYDPTFKVGI	LKLPCENGIAMLVLPDP--	253
ZPI2_DRE	209	YVDKYNIVQVPMMLLEEK	FSVVEDRDIRARVLR	LRPVRGGASMLLILDP--	256

		282b			
AlAT_HSA	257	--EGKLOHLENELTHDIT	TKFLENEDRRSASLHLP	PKLSITGTYDIKSVLG	304
AlAT_DRE	293	--DGKMKLEESICRHH	LKNTHDKLFRSSVDL	FMPKFSISATSKLDDITM	340
Dre-Spn-8	293	--DGKMKLEESICRHH	LKNTHDKLFRSSVDL	FMPKFSISATSKLDGILK	340
Dre-Spn-9	231	MGRKTIKOLEMTVSRQ	HVEKTRRSVSEKVDI	YVPKLSLKTSYSLKDILK	280
Dre-Spn-10	231	MKQKTIKOLEMTVSRQ	HIEKTRRSVSEKVDI	YVPKLSLKTSYSLKDILK	280
Dre-Spn-12	231	MEQKTIKOLEMTVSRQ	HVEKTRRSVAFKKT	VDIYVPKLSLKTSYSLKDILK	280
Dre-Spn-11	235	MGRKTIKOLEMTISRQ	HKKTRRGVSKDEVDI	YVPKLSLKTSYSLKDILK	284
ZPI1_DRE	254	-EDVDYTYVDESMT	GEVFRGWAKLKKTK	LEIQLPFSLKQSNLSVSLP	302
ZPI2_DRE	257	-ADADYTAIEDEISA	ERLHGWIKNMRR	IKTYTPHTETHNTHAHTCHVHTC	305

		331c				
AlAT_HSA	305	QLGITKVESNGADLSGV	TEEAPILKLSKAVHK	AVLTIDEK <b>GTEAAGAMFL</b>	353	
AlAT_DRE	341	DMGMTDAFDYKADF	SGMTEEVKVRVSRV	LHQAVMSVDEK <b>GTEAAAITTI</b>	389	
Dre-Spn-8	341	DMGMTDAFNDKADF	SGMTEEVKVKVSOV	LHQAVMSVDEK <b>GTEAAAITTI</b>	389	
Dre-Spn-9	281	GMGMADMFSDKAD	FTGVSEEKIVVSK	VLHKATLIDIDEK <b>GTTAAAVTTV</b>	328	
Dre-Spn-10	281	GMGMTDMFSDKAD	FTGVSEENIFVSK	VLHKATLIDIDEQ <b>GTTAAAVTGV</b>	328	
Dre-Spn-12	281	GMGMADMFSDKAN	FTGVSEEKIFVSK	VLHKATLIDIDEQ <b>GTTAAAVTGV</b>	328	
Dre-Spn-11	285	GMGMADMFSDKAD	FTGVSEE-----	-----	304	
ZPI1_DRE	303	SLGVKEIFG	STANLTGLISS	SEGLKLVVQKVA	VDVDES <b>GGSLAEASGN</b>	351
ZPI2_DRE	306	AIKHICFFCQVNES	VCGLAQALTLT	LSVSETKAVIEV	YEQ <b>GTSAASSTSV</b>	355

---

```

ALAT_HSA      354  EATPMSIPPEVKFN--KPFVFLMIEQNTKSPLEFMGKVVNPTQK-  394
ALAT_DRE      390  EIMPMSLPHTVILN--RPFLVLIVEDSTMSILFMGKITNPTA--  429
Dre-Spn-8     390  EIMPMSLPDTVILN--RPFLVLIVEDSTMSILFMGKITNPTA--  429
Dre-Spn-9     329  HLRFMSYSPMSDLSFDREPFMIFITDQTNDNILEFVGKVVNPNEKL  372
Dre-Spn-10    329  EIRPYSYDPLSDLKFDHPEPFMIFITDQTNDNILEFVGKVVNPNEKL  372
Dre-Spn-12    329  SMRVRLHNPLSILKFNRPEPFMIFITDQTNDNILEFVGKVVNPNEKL  372
Dre-Spn-11    -    -----
ZPI1_DRE      352  -LFMNP LPPRLTFN--RPFIFV VYHEVTKC ILYIGRVVDPTKN-  391
ZPI2_DRE      356  GITAYSLEDTFIIN--RPFFFFLYHEETASLLEFMGRVIDPTLS-  396

```



**Appendix 8.3.20: Alignment of  $\alpha_1$ -antitrypsin sequences from vertebrates.** Gene specific features include an inhibitory RCL (red box). Conserved intron positions are indicated above the alignment. In *Tetraodon*, A1AT\_TNI sequence has gap in genomic sequence at intron 331c position to RCL.

ALAT_HSA	-	-----	-
ALAT_MMU	1	MTPSISWGLLLL LAGLCCLVPS	21
ALAT_RNO	1	MAPSISRGLLLL LAAALCCLAPS	21
ALAT_GGA	-	-----	-
ALAT_XTR	1	MRAFLIVSLALLCAGVLAD	23
ALAT_FRU	1	MRDIIASCMLAALLAVASAD	22
ALAT_TNI	1	MSSVPAQAEKMHGMIAGCLLAALLAVASAD	33
ALAT_DRE	1	MRGNIFCCAIAALLVATAWAA	22
ALATL_PMA	1	MKLFILLLLAFCAALCSPCVGE	49
		-EDYDDRPYMQPFHLIPPLSVQATEQP	
ALAT_HSA	1	-----	19
ALAT_MMU	22	-----	37
ALAT_RNO	22	-----	37
ALAT_GGA	-	-----	-
ALAT_XTR	24	TKHGKDHDH-----	63
ALAT_FRU	23	HQHNN-----	39
ALAT_TNI	34	HDHHDHHDH-----	57
ALAT_DRE	23	HDGHEGHDHGSHTADHHHLLHHGKD	55
ALATL_PMA	50	LASNETWDYPEPLAPGQSPAASSEEGSSEEKGDERSHRGEGRRGRKDKY	99
ALAT_HSA	20	HPTFNKITPNLAEFAFSLYRQLAHQSNTS	67
ALAT_MMU	38	SPASHEIATNLGDF AISLYRELVHQSNTS	85
ALAT_RNO	38	SPTYRKISSNLADFAFSLYRELVHQSNTS	85
ALAT_GGA	1	-MPEQAVGNRVCQFAACCFYKEISSHENS	47
ALAT_XTR	64	DMACHKIAPSNSQFAFKLFQVADHPSE	111
ALAT_FRU	40	EDICHLVSNGNADEGFALYKQLNAKSDAG	88
ALAT_TNI	58	EHPCHLVSNNADEGFALYKHLKAKSDAK	106
ALAT_DRE	56	VDACHLLAPHNADEFAFSLYKKLASNPDAQ	105
ALATL_PMA	100	KSKTQRIASAVNGLGFRLYKQVLGGAGPA	148
ALAT_HSA	68	TKADTHDEILEGILNFNLTEIPE	115
ALAT_MMU	86	SKGDTHTQILEGILQFNLTQTSE	133
ALAT_RNO	86	SKGDTIRKQILEGILEFNLTQIPE	133
ALAT_GGA	48	ARSDTLAQILRVILHFNPRAISE	95
ALAT_XTR	112	ARADTLNQIIEGILNFNNTKITE	159
ALAT_FRU	89	ARGDTRSQLFSTLGY	134
ALAT_TNI	107	ARGDTRSQLFSSLGY	152
ALAT_DRE	106	AKGSTLSQIYSGLGY	151
ALATL_PMA	149	ANGSTRAELDTALGFKELLHGKKKAKSMKYFARLNSALYRRSAGFELMGK	198

ALAT_HSA	116	NGFLFSEGLKLVDFLEEDVKKLYHSEAFVNEFGDTEEAKKQINDYVEKGT	165
ALAT_MMU	134	NGLFVNNDLKLVEKFLLEEAKNHYQAEVFSVNEFAESEEAKKVINDFVEKGT	183
ALAT_RNO	134	NGLFVNKNLKLVEKFLLEEVKNHYHSEAFSVNEFADSEEAKKVINDYVEKGT	183
ALAT_GGA	96	NVLFVLDRLKPPQRFNLNLSREFYEGEITYPMNFKRSDQAQTKINDYVAERT	145
ALAT_XTR	160	NALFIDNNVKLIQQFLDDVKKYYESEAFSTDFENNAEEAKKQINSYVEKQT	209
ALAT_FRU	135	NAAAVDKTFNPLKAYMTDIKDYYSAEVLVDVDFKNPAEAAAEINKYIALNT	184
ALAT_TNI	153	NAAAVHQTTFNPLKTYMEDIKDYTAHVFDVDFETKPVDAAAEINKYIARNT	202
ALAT_DRE	152	AGVAIRDGFKVVDQFLKDAQHYNSEAFGVDFSKPEIAAAEINKFIARKT	201
ALATL_PMA	199	NVVFSSKGLWLYRQFTRTVAHLFKSNVRSVDFEGESKEAVELMAYIEKVT	248

## 192a

ALAT_HSA	166	QCKITVDLVKELDRDITVFALVNYIFFKGKWERPFEVKDTEEEDEFHVDQVTT	215
ALAT_MMU	184	QCKIAEAVKKLDQDITVFALANYILFKGKWKKEPDPENTEEAEFHVDESTT	233
ALAT_RNO	184	QCKITVDLMKQLDEDITVFALVNYIFFKGKWKRPENPEHTRDADFHVDKSTT	233
ALAT_GGA	146	NGKITKDLINNLDPLTELLLLISYIYFNAEWEKPFNPQYTKKEKEFFVDGNKA	195
ALAT_XTR	210	HCKITVDLLSSVDKNAVLYLILNYIFFRGKWEKPEEEKFTQDGFIEHVDENTN	259
ALAT_FRU	185	GDMIKDQVKDLDPDITAMVILINNYIFFKGEWERPENSILTKQKMDENVDESTK	234
ALAT_TNI	203	GDMIKDQVKDLDPDITVMMLINNYIFFKGEWEKPFENGLTRKMDFHVDESTN	252
ALAT_DRE	202	HDKITNMVKDLADITVMMLINNYIFFRGKWEKQFDAKLTHKADFKVDQDFTT	251
ALATL_PMA	249	SKKFTDVISDVDITATSLMIVNVYIFFKGSWANKEPDLTKNVREFWVNSSYS	298

ALAT_HSA	216	VKVPMMKRLGMFNIQHCKKLSWVLLMKYLGNATAIFFLPPD-EGKLQHLE	264
ALAT_MMU	234	VKVPMMTSLGMLHVHHCSTLSSWVLLMDYAGNATAVFLPPD-DGKMQHLE	282
ALAT_RNO	234	VKVPMMNRLGMFDMHYCSTLSSWVLLMDYLGNATAIFLPPD-DGKMQHLE	282
ALAT_GGA	196	VEVPMMFGIGAFKHGYDEQLSSTVVQMDYKGGASAFFVLPD-QGRMRKLE	244
ALAT_XTR	260	VTVPMMRRNGMYNVAFDEKLGCTVVQIPYKGNATALFILPPD-EGKLRQVE	308
ALAT_FRU	235	VQVDMRRTGRFDYYSDFDNHSSIIMLPYKGNATSMIILPN-EGKMKHVE	283
ALAT_TNI	253	VPVDMRRTGRFDYFDLDNHSSVIMLPYKGNATSMIILPS-EGKMEHVE	301
ALAT_DRE	252	VQVDMRRTGRYDIYQDPVNQTTVLMVPPYKGNATSMIIVLPPN-DGKMKLE	300
ALATL_PMA	299	MMVPTMHQRAKLSYAQDRKLRSTVIKLPYEGGASMLVIVPHRTEELPKVE	348

## 282b

ALAT_HSA	265	NELTHDIITKFLENE--DRRSASLHLPKLSITGTYDLKSVLGLGKITKVF	312
ALAT_MMU	283	QTLSEKELISKFLLNR--RRRLAQIHFPRLSISGEYNLKTLMSPLCITRIF	330
ALAT_RNO	283	QTLTKDLISRFLNR--QTRSAILYFPKLSISGTYNLKTLLSSLGICITRVF	330
ALAT_GGA	245	KKLSCERMARWRTLVS-KSNSVNLYLPKFTLHGRYNLKNILYKMGIMDLF	293
ALAT_XTR	309	EALEKSTIMSWKQF--RYQSIELTIPKFSIMATLDLIEELKKFGVTDVF	356
ALAT_FRU	284	NSISKEQILHWFNSL--FRMSVELMLPKFSISADASLNEVLQEMGVTVNF	331
ALAT_TNI	302	GSISKEQILHWHNSL--FRMSVELMLPKFSISADASLGEVLQEMGVTSVF	349
ALAT_DRE	301	ESICRHHLKNWHDKL--FRSSVDLFPKFSISATSKLDDILMDMGMTDAF	348
ALATL_PMA	349	ESVSQEQLEEWLSLLGPNHYVQLSLPKFKISVSYDLKAYLSAMGMPSMF	398



## 331c

ALAT_HSA	313	SNGADLSGVTEE-APL	KLSKAVH	KAVLTIDEK	GTEAAGAMFLEAIPMSIP	361
ALAT_MMU	331	NNGADLSGITEENAPL	KLSQAVH	KAVLTIDET	GTEAAAVTVLQMVPMSMP	380
ALAT_RNO	331	MNDADLSGITED-APL	KLSQAVH	KAVLTLDER	GTEAAGATVVEAVPMSLP	379
ALAT_GGA	294	TDKADLSGITGQ-PQ	HRIHQAIHQ	AVVKVDET	GTEAAAATGMEIVPMSVP	342
ALAT_XTR	357	SQNADLSGIVEG-TPL	KVSKAVH	KAGLSVDET	GTEAAAATAFEIMPMMIP	405
ALAT_FRU	332	SDAADLSGISQE-PK	LKVKVSK	SHRAVLDVDEK	GTTAAASTTIEIMPMSMP	380
ALAT_TNI	350	SDAADFSGISQE-PK	LKVKSK	-----	AAASTTIEIMPMSMP	383
ALAT_DRE	349	DYKADFSGMTEE-VK	VRVSRVLH	QAVMSVDEK	GTEAAAITTIEIMPMSLP	397
ALATL_PMA	399	SYGADLSRITGM-QK	LHVDKITH	KSVLHVNEE	GTEAKAETVVGIMPISMP	447

ALAT_HSA	362	PEVKFNKPFVFLMIEQNTKSPLEMGKVVNPTQK-	394
ALAT_MMU	381	PILRFDHPFLFIIFEHTQSPIFLGKVVDPHTK-	413
ALAT_RNO	380	PQVKFDHPFIFMIVESETQSPLEVGKVIDPTR--	411
ALAT_GGA	343	VVIRMNRPFLLVITLR--ENILEMGKIVNPLEKD	374
ALAT_XTR	406	PHILFNRAFVVIIYDIPKSIKLVAKVVNPKN--	437
ALAT_FRU	381	GTMKVDRPFLVLILERSTRSILEMGKINNPTAQ-	413
ALAT_TNI	384	ETMVVNRPFLLVLIHSTRSILEMGKVNNP----	413
ALAT_DRE	398	HTVILNRPFLVLIVEDSTMSILEMGKITNPTA--	429
ALATL_PMA	448	PTVTVDRPFVVLIIYDEKTRAVIEMGRVADEKQ--	479

Appendix 8.3.21: Alignment of THBG (serpinA7) protein sequences from vertebrates. Common features incorporated (as section 8.3)

ALAT_HSA	1	-----EDPQGDAAQKTDTSHHQDQDPTFNKITPNLA	31
THBG_HSA	1	MSPFLYLVLVLLVGLHATIHCA---SPEGKVTACHSSQPNATLYKMSSINA	47
THBG_MMU	1	MSVFFYLFVLFVGLQATIHCAPHNSSEGKVTTCMLPQQNATLYKMP SINA	50
THBG_RNO	1	MSMFFYLFLLVGLQATIHCAPHNSSEGKVTTCMLPQQNATLYKMP SINA	50
ALAT_HSA	32	EFAFSLYRQLAHQSNSTNIFFSPVSIATAFAMLSLGTKADITHDEILEGIDN	81
THBG_HSA	48	DFAFNLYRRFTVETPDKNIFFSPVSIISAAVLVLSFGACCSTQTEIVETLG	97
THBG_MMU	51	DFAFSLYRRLSVENPDLNIFFSPVSIISVALAMLSFGSGSSTQTQILEVLG	100
THBG_RNO	51	DFAFRLYRKL SVENPDLNIFFSPVSIISAAVLVLSFGSGSSTQTQILEVLG	100
ALAT_HSA	82	FNLTEIPEAQIHEGFQELLRTL NQPDSQLQLTTGNGLFLSEGLKLVDKFL	131
THBG_HSA	98	FNLTDTPMVEIQHGFQHLICSLNFPKKELELQIGNALFIGKHLKPLAKEL	147
THBG_MMU	101	FNLTDTPVTELQQGFQHLICSLNFPKNELELQMGNAVF IGQQLKPLAKEL	150
THBG_RNO	101	FNLTDTPVKELQQGFQHLICSLNFPNNELELQMGNAVF IGQQLKPLAKEL	150
ALAT_HSA	132	EDVKKLYHSEAFVNFVGDTEEAKQINDYVEKGIQKTIIVDLVKELDRDITV	181
THBG_HSA	148	NDVKTLYETEVEFSTDFSNISAAKQEIINSHVEMQIKKGVVGLIQDLKPNII	197
THBG_MMU	151	DDVKTLYETEVEFSTDFSNVSAAQHKIINSYVEKQIKKTIIVGLIQGLKLNII	200
THBG_RNO	151	DDVKTLYETEVEFSTDFSNVSAAQHEIINSYVEKQIKKTIIVGLIQDLKLNII	200
		192a	
ALAT_HSA	182	FALVNYIFFKCKWERPFEVKDIEEE-DEHVDQVTTVKVPMMKRLGMFNIQ	230
THBG_HSA	198	MVLVNYIHFKAQWANPEFDP SKIEDSSSELIDKTTTVQVPMMHQMEQYYHL	247
THBG_MMU	201	MILVNYIHFRAQWANPEFRVSKIEESSNESVDKSTTVQVPMMHQLEQYYHY	250
THBG_RNO	201	MILVNYIHFKAQWANPEFRVSKIEESSNESVDKSTTVQVPMMHQLEQYYHY	250

---

```

ALAT_HSA 231 HCKKLSWVLLMKYLGNATAIFFLPDEGKLOHLENELTHDIITKFLNED 280
THBG_HSA 248 VDMELNCTVLQMDYSKNALALFVLPKEGQMSVEAAMS SKTLKKWNRLQ 297
THBG_MMU 251 VDMELNCTVLQMDYSENALALFVLPKEGHMEWVEAAMS SKTLKKWNYLLQ 300
THBG_RNO 251 VDVELNCTVLQMDYSANALALFVLPKEGHMEWVEAAMS SKTLKKWNHLLQ 300

```

## 282b

```

|
ALAT_HSA 281 RRSASLHLEPKLSITGTYDLKSVLGQLGITKVF SNGADLSGVTEEAPDKLS 330
THBG_HSA 298 KGWVDLFVPEKFSISATYDLGATLLKMGIQHAYSENADF SGLTEDNGLKLS 347
THBG_MMU 301 KGWVELFVPEKFSISATYDLGSTLQKMGMRDAFAESADFP GITEDSGLKLS 350
THBG_RNO 301 KGWVELFVPEKFSISATYDLGSTLQKMGMRDAFAESADFP GITKDNGLKLS 350

```

## 331c

```

|
ALAT_HSA 331 KAVHKAVLTIDEKGTEAAAGAMFLEATPMS----TPPEVKFNKDFVFLMIE 376
THBG_HSA 348 NAAHKAVLHIGEEKGTEAAAVPEVELSDQPENTFLHPIIQIDRSFLLILE 397
THBG_MMU 351 YAFHKAVLHIGEEGTTKEGASPEVGSLDQQEVPPLHPVIRLDRAFLMILE 400
THBG_RNO 351 YAFHKAVLHIGEEGTTKEGASPEAGSLDQQEVAPLHAVIRLDRTFLLMILE 400

```

```

ALAT_HSA 377 QNTKSPLEMGKVVNPTQK 394
THBG_HSA 398 RSTRSILFLGKVVNPTQA 415
THBG_MMU 401 KRTRSVLFLGKLVNPTKQ 418
THBG_RNO 401 KRTRSVLFLGKVVDPTKE 418

```

**Appendix 8.3.22: Alignment of PAI1 sequences from vertebrates.** Gene specific features include inhibitory RCL (red box) and group V3 specific discriminating amino acid indels (indicated by \*\* or \*) and an intron indel (indicated by \$). Conserved intron positions are indicated above the alignment. A predicted low complexity region (GENSCAN) within the intron at position 290b for PAI1\_FRU was deleted manually.

```

A1AT_HSA 1  -----EDPQGDAAQKTDTSHHDDHPFTFNKITPNLAEEFAFSLYRQLAHQSNSTNIFFSPV 55
PAI1_HSA 1  MQMSPALTCVLVGLALVFGEGSAVHHPPSYVAHLASDFGVRVFOQVAQASKDRNVVFSFY 60
PAI1_MMU 1  MQMSSALACLILGLVLVSGKGF TLPLRESHTAHQATDFGVKVFOQVVQASKDRNVVFSFY 60
PAI1_RNO 1  MQMSSALTCLTGLGLVLVFGKGFASPLPESHTAQQATNFGVKVFOHVQVQASKDRNVVFSFY 60
PAI1_XTR 1  --MIVVLLSLASVTSACNR-----VSRVAQKGTSGFLRRLFQEVLDQWGKMLGFSFY 50
PAI1_FRU 1  MLFA Y T L L L L A L S R A A L S S -----LQDKQTD FGLKVFSQLSQSSVDKNVAMSPY 49
PAI1_TNI 1  MLLTY L L L L A L N H A G L G L G S -----LQDKQTD FGLKLFSQLSQSLADKNLAMSPY 51
PAI1_DRE 1  MQSLSVLLIFALCASSLCN-----LIQDKQTD FGLQVFAEAVQSA PDRNLALSPY 50

```

86a

```

A1AT_HSA 56  SIATAFAMLSLGTAKADTHDEILEGLNFNFLTEIPEAQIHEGFQELLRTL NQP-DSQLQLTT 114
PAI1_HSA 61  GVASVLA MLQ LTTGGETQQQIQAAAMGFKIDDKGMAPALRHL YKELMGPW----NKDEIST 116
PAI1_MMU 61  GVSSVLA MLQ LTTAGKTRRQIQDAMGFKVNEKGTAHALRQLSKELMGPW----NKNEIST 116
PAI1_RNO 61  GVSSVLA MLQ LTTAGKTRQIQDAMGFNISERGTAPALRKL SKELMGSW----NKNEIST 116
PAI1_XTR 51  GVTSA LSVLQSGAAGTLDQIRKALNYGHKEWAVALALNKLREQISGQOKSAEDPKPVHI 110
PAI1_FRU 50  GAVSVLA MAQ LCAAGKTLRALNSAMGFSLLARGMSRQORLLHRDLS-----SEDGVET 102
PAI1_TNI 52  GAVSVIA MAQ LCAAGKTLRALDSAMGYSLLARGMSRQORLLQRDLS-----SEEGVET 104
PAI1_DRE 51  GIASVLA GMAQMGAYGATL KLLASKMGYSLQERGM PKLQORLLQRDLA-----SEDGVEV 103

```

167a\$

```

A1AT_HSA 115  GNGLFLSEGLKLVDRFLEDVKKLYHSEAF TVMFGDTEEAKKQINDYVEKGTGGKIVD--L 172
PAI1_HSA 117  TDAIFVQRDLKLVQGFMPHFFR LFRSTVKQVDFSEVERARFIINDWVKTHTKGMISNLLG 176
PAI1_MMU 117  ADAIFVQRDLKLVQGFMPHFFKLFRTTVKQVDFSEVERARFIINDWVERHTKGMINDLLA 176
PAI1_RNO 117  ADAIFVQRDLKLVQGFMPHFFKLFRTTVKQVDFSEVERARFIINDWVERHTKGMISDLLA 176
PAI1_XTR 111  ADGLFVQRDLSTPGFLQRFQATFHRHLSQVNF TDVAQAKDIINQWVENKTDGMIKDLVG 170
PAI1_FRU 103  ASAVMVERKMSLEKGYRRALVKAFQTHPHQVDFTRPEQAVGVINEWVSDHTAGAI P DFLQ 162
PAI1_TNI 105  ASAAMVERKMSLEKGYRRALVKAFQTHPHQVDFTKPEQAVNIINEWVSDHTAGAI P DFLA 164
PAI1_DRE 104  ASGVMVDRKIILEKVFRRSLSKAFQSVPHQIDFSQPEMARQVIMS WTS DHTGGMISEFLP 163

```

230a

```

A1AT_HSA 173  VKELDRDTVFA LVMYIFFKCKMERPF EVKDTTEEEDFHVDQVTTVKVPMMKRLGMFNIQHC 232
PAI1_HSA 177  KGAVDQLTRLV LVMALYFMCQKTPFPDSS THRRLFHKSDGSTVSVPMMAQTNKFNYTEF 236
PAI1_MMU 177  KGAVDELTRLV LVMALYFSCQKTPFLEASTHQRLFHKSDGSTVSVPMMAQSNKFNYTEF 236
PAI1_RNO 177  KGAVNELTRLV LVMALYFMCQKTPFLEASTHQRLFHKSDGSTISVPMMAQNNKFNYTEF 236
PAI1_XTR 171  SNNIPPLTRLV LLSAVHFSCRWTVPFLEKATHQRPFYRSDGSHVQVQMMANTGKYNCSEF 230
PAI1_FRU 163  SGSLTDETRLV LLMALS FQAPWKVPFDPKRTAERMFHCANGSTVVPVHMMTLTNHYHYGEF 222
PAI1_TNI 165  SGSLTDETRLV LLMALS FQALWKVPFDPKQTAERMFHCANGSVVPVHMMTLTNYFHYGEF 224
PAI1_DRE 164  SGVLSLTRLV LLMALHFHCWVKTPFDP RNTREQLFHTVNGSAVSVPMMTTTTQKFNYTEF 223

```

		*						
A1AT_HSA	233	KKLSS---	WVLLMKYLGN-ATAIFFLP--	DEGKLQHLENELTHDIITKFL	ENEDRRSASL 286			
PAI1_HSA	237	TTPDGHYYD	ILELPYHGD	TLSMFI	AAPYEKEVPLSALTNILSAQLISHW	KGNMTRLPRLL 296		
PAI1_MMU	237	TTPDGL	EYDVVELPYQ	RD	TLSMFI	AAPFEKDVHLSALTNILDAELIRQW	KGNMTRLPRLL 296	
PAI1_RNO	237	TTPDGH	EYDILEL	PYHGET	TLSMFI	AAPFEKDVPLSAITNILDAELIRQW	KSNMTRLPRLL 296	
PAI1_XTR	231	TTPDGD	FYDVIEL	PYEGEELS	M	LI	AAPYEKNVPLSAITNILTPELIAQW	KAQMKKVTRLL 290
PAI1_FRU	223	VTTEGID	YDVIEVP	YEGD	SLSMLLVSP	IEREVPLSALIGDLSSQ	RIRQWRQELRRVKRQL 282	
PAI1_TNI	225	VTTEGID	YSVIEVP	YDGD	TLSMLLASP	IESDVPLDKVIADLSSKRIHQWRQELRRVKRQL 284		
PAI1_DRE	224	VSKDGD	VDVIEIP	YEGE	SISMLLVTP	FEKDVPLSALNKELSSSRIHQWRQEMRKISKQL 283		
		290b		323a				
A1AT_HSA	287	HLPKLSITGTYD	LKSVL	GQLGITKVF	SNG-ADLSGV	TEEA	PLKLSKAVH	KAVLTIDEKGT 345
PAI1_HSA	297	VLPKFSLETEVD	LRKPLE	NLCMTDM	FRQFQADFT	SLSDOEPL	HVAQALQKVKIEV	NESGT 356
PAI1_MMU	297	ILPKFSLETEVD	LRGPLE	EKLGMPDM	F	SATLADFT	SLSDOEQLSVAQALQKVR	IEVNESGT 356
PAI1_RNO	297	ILPKFSLETEVD	LRGPLE	EKLGMTD	IF	SSTQADFT	SLSDOEQLSVAQALQKVKIEV	NESGT 356
PAI1_XTR	291	VLPKFSLLSEVD	LKKPLER	LGITDM	F	TOETADFS	SRLSSEKPLYVSEAFQKIKVEV	TEKGT 350
PAI1_FRU	283	SMPRFTLNSEVNF	KSALLNM	CLGDVFN	L	ATADFT	TRITTEERLCVSKIMQKIKIEV	NEHGT 342
PAI1_TNI	285	SMPRFTFNSEVDF	KSALLK	MCLGDVFN	M	ATADFT	TRITTEERLCVSKIMQKIKIEV	NEHGT 344
PAI1_DRE	284	SIPRFSMDTEID	LKSTLS	SRMCLGD	I	F	SQSRADFSRITTEEPLCVSKVLQ	RVKLEVNEEGT 343
		352a		380a				
A1AT_HSA	342	EAAGAMFLEAIPMS	IPPE	VKFNKPFV	FLMIEQNTKSPLFMGKV	VNPTQK----	394	
PAI1_HSA	357	VASSSTAVIVSARM	APEE	IIMDRPFL	FVVRHNP	TGTVLFMGQVMEP	-----	402
PAI1_MMU	357	VASSSTAFVISARM	APTE	MVIDRSFL	FVVRHNP	TETILFMGQVMEP	-----	402
PAI1_RNO	357	VASSSTAILVSARM	APTE	MVLDRSFL	FVVRHNP	TETILFMGQVMEP	-----	402
PAI1_XTR	351	RASAAT	-----	-----	-----	-----	-----	356
PAI1_FRU	343	KAAAATAAVMF	SRMAVEE	IALDRPFL	FLIQHKPTG	TLLFMGQFNHP	QQQ----	391
PAI1_TNI	345	KASAASA	AVMF	SRMAVEE	IALDRPFL	FLIQHKPTG	-----	379
PAI1_DRE	344	KGSSATAA	VIYSRMAVEE	ITLDRPFL	FLIQHKPTG	ALLFSGQLTQP	QEY----	392

**Appendix 8.3.23: Alignment of GDN sequences from vertebrates.** Gene specific features include an inhibitory RCL (red box), helix-D (yellow box), a conserved N-glycosylation site (cyan box) and group V3 specific discriminating amino acid indels (indicated by \*\* or \* ) and an intron indel (indicated by \$). Conserved intron positions are indicated above the alignment. A predicted low complexity regions (GENSCAN) within the intron at position 290b for GDN\_FRU was deleted manually.

A1AT_HSA	1	--EDPQGDAAQKTDTSHHDDQHPFTFNKITPNLAEEF	AFSLYRQLAHQSNSTN	IFFSPV	58		
GDN_HSA	1	--MNWHLPLFLLASVTLP	SI	CSHFNPLSLEELGSNTGIQVFNQIVKSRPHDMIVISPHGIA	59		
GDN_MMU	1	--MNWHFPFFILTTVTLYSVHSCFNLSLSLEELGSNTGIQVFNQIIKSRPHENVVVS	PHGIA		59		
GDN_RNO	1	--MNWHFPFFILTTVTLLSSVYSQLNLSLSLEELGSDTGIQVFNQIIKSPHENVVI	SPHGIA		59		
GDN_GGA	1	--MNWHFSLFLG--TLASVCSQFNFYPLEELSSDVGIQVFNQIVKAKPQDMVVV	SPHGIA		58		
GDN_XTR	1	--MRLVIFPFLVAF	FLASVQPELDPLSLEELGSDIGIQVFNQVARTRPHENIVM	SPHGIS	58		
GDN_FRU	1	MKTL	SFLCLFLGLVVLRGHGALSQAPS	YGERGSDLGIQVFQOEVRSRPLDMIVL	SPHGVA		
GDN_TNI	1	MNHLVFLCLLGLATF	HSHDGAHSLASSYGERGSDLGIQVFQREVHSRPLDMIVL	SPHGVA	60		
GDN_DRE	1	--MCVLFRCGV--LFL	CSVSVS	QSSSYGARGSDLGLQVFMQVLQDRAQENVLL	SPHGVA		
					58		
					86a		
A1AT_HSA	59	TAFAMLSL	GTKADTHDEILEGLNFNLTEIPEAQIHEGFQELLRTL	NQPDSQLQLTTGNGL	118		
GDN_HSA	60	SVL	GMLQLGADGR	TKKQLAMVMRYGVNGV----	GKILKKINKAIVSKKNKDIVTVANAV		
GDN_MMU	60	SIL	GMLQLGADGR	TKKQLSTVMRYNVNGV----	GKVLKKINKAIVSKKNKDIVTVANAV		
GDN_RNO	60	SIL	GMLQLGADGR	TKKQLSTVMRYNVNGV----	GKVLKKINKAIVSKKNKDIVTVANAV		
GDN_GGA	59	SVL	GVLQLGADGR	TKKQLTMMRYSVNGV----	GKALKKINRLIVSKKNKDIVTIANAV		
GDN_XTR	59	SVL	GMLQLGADGR	TKKQLTMVMRYKINEV----	AKSLKKTNRRAIVAKKNKDIVTTANGV		
GDN_FRU	61	SIL	GMLLP	GAHGETRQVLTALRYKKNP----	YKMLKKLHKTLTAKANQDSLLIANAM		
GDN_TNI	61	SIL	GMLLP	GAHGETRQVLTALRYKKNP----	YKMLRKLHKTLTAKANQDSVLIANAM		
GDN_DRE	59	SVL	GMLLP	GAHGDTRRQLLNGLKYKKNP----	YKMLRKLHKSLTTKSNADIVTIANAL		
					113		
					167a\$		
A1AT_HSA	119	FLSEGLKLVDR	FLEDVKKLYHSEAF	TVMFGDTEEAKKQ	INDYVEKGTQGGKIVD--LVKEL	176	
GDN_HSA	115	FVK	NASEIEVPPFVTRNKDVFQCEVRNVMF	EDPASACDS	IMAWVKNETRDMIDNLLSPDLI	174	
GDN_MMU	115	FLR	NGFKMEVPPFAVRNKDVFQCEVQVMVF	QDPASASES	IMFWVKNETRGMIDNLLSPNLI	174	
GDN_RNO	115	FVR	NGFKMEVPPFAARNKEVFQCEVQSVVF	QDPASACDA	IMFWVKNETRGMIDNLLSPNLI	174	
GDN_GGA	114	FAK	SQFKMEVPPFVTRNKEVFQCSVKSVDF	EDPNTACDS	IMQWVKNETRGMIDQVVAPDDI	173	
GDN_XTR	114	FAS	AFKVEGSFVYKNKDIFHSDVRSVDF	QEKNTAASI	IMQWVKNETKGMIEGLISPELL	173	
GDN_FRU	116	FTK	EGFPMKEAFVATNKANFQCESRSLDF	FRHPSKAADD	INEWVSNKTKGHIPSLVKADML	175	
GDN_TNI	116	FTK	DGFPMEETFRATNKANFQCESRSLDF	FRHPQTAAD	INEWVSNKTKGHIPSLIKADML	175	
GDN_DRE	114	FPN	EFGSMKEDFLSANRENFLCESHSVDYSDPEAAAQS	INDWVKNSTKGCIPSVVTADMF		173	
						173	
						230a	
A1AT_HSA	177	DRD--TVF	ALVNIYIFFKCKMERPF	EVKDTTEEEDFHVDQVTTVKVPMMKRLGMFNIQHCKKL		235	
GDN_HSA	175	DGVL	TRLVLMNAVYFRCLWKS	RFQPE	TKKRTFVAADGKSYQVPMLAQLSVFRCGSTSAP	234	
GDN_MMU	175	DGAL	TRLVLMNAVYFRCLWKS	RFQPE	STKKRTFVAGDGKSYQVPMLAQLSVFRSGSTRTP	234	
GDN_RNO	175	DSAL	TKLVLMNAVYFRCLWKS	RFQPE	TKKRTFVAGDGKSYQVPMLAQLSVFRSGSTKTP	234	
GDN_GGA	174	D--SL	TRLVLMNAVYFRCLWKS	RFQPE	TKKRPFFYAGDGKTYQVPMLSQLSIFRCGTTSTP	232	
GDN_XTR	174	DSSV	TRLVLMNAVYFRCLWKS	RFHP	ENTTKKRTFHGPDGKDRQVPMLAQLSLFRSGSASTP	233	
GDN_FRU	176	DSAL	TRLVAVNSIYFRCLWKS	RFQ	AEDTKMRPFTSGDGTVHKVPMMSQLSVFNIGMVTTTP	235	
GDN_TNI	176	DSAL	TRLVAVNSIYFRCLWKS	RFQ	PE	TKLRHFTGGDGNVSKVPMMSQLSIFNISMATTP	235
GDN_DRE	174	DTAL	TRLVAVNSIFFRCLWKS	RFQ	PQSTKPRSF	TAGDGNTYKVPMMSQLSVFNMGOASTP	233

		*			
A1AT_HSA	236	SS---WVLLMKYLGN-ATAIFFLP--DEGKLOHLENELTHDIITKFLNEDRRSASLHLP			289
GDN_HSA	235	NDLWYNFIELPYHGESISMLIALP TESSTPLSAIIPHISTKTIDSWMSIMVPKRVQVILP			294
GDN_MMU	235	NGLWYNFIELPYHGESISMLIALP TESSTPLSAIIPHITTKTIDSWMNTMVPKRMQLVLP			294
GDN_RNO	235	NGLWYNFIELPYHGESISMLIALP TESSTPLSAIIPHISTKTINSWMNTMVPKRMQLVLP			294
GDN_GGA	233	NELWYNIIELPYHGEMISMLIALP TENTTPLSAIIPHISTKTIGSWMTTMVAKRVQVILP			292
GDN_XTR	234	NGLWYNVIELPYHGGSISMLVALP TEKSTPLSAIIPHISTKTLQSWM-TMSPKRVQLILP			292
GDN_FRU	236	QGLKYKVIELPYHGNTVSMIALP SEENTPLSHIIP TISTASVQNWTKLMHMMKIRLLIP			295
GDN_TNI	236	QGLKYKVIELPYHGNTVSMIALP SEEDTPLSHIIPHISTATVQSWTQLMHRKIRLLIP			295
GDN_DRE	234	DGQKYIVIELPYHGNSMSMFIALP TEDSTPLSSILPHISTNTIQSWTKLMNPRRMRLLP			293
		290b		323a	
A1AT_HSA	290	KLSITGTYDLKSVLGQLCITKVFNSG-ADLSGVTE-EAPLKLSKAVHKAVLTIDEKGT <b>EA</b>			347
GDN_HSA	295	KFTAVAQTDLKEPLKVLGITDMFDSSKANFAKITTGSENLHVSHILOKAKIEVSE <b>DGTKA</b>			354
GDN_MMU	295	KFTAVAQTDLKEPLKALGITEMFEPSKANFTKITR-SESLHVSHILOKAKIEVSE <b>DGTKA</b>			353
GDN_RNO	295	KFTALAQTDLKEPLKALGITEMFEPSKANFAKITR-SESLHVSHILOKAKIEVSE <b>DGTKA</b>			353
GDN_GGA	293	KFTAVAETDLKDPLKALGITDMFDESNSNFAKITR-TEGLHVSHVLQKTIEVSE <b>DGTKA</b>			351
GDN_XTR	293	KFSVEAEADLKEPLRNLGITEMFDVSKANFAKISR-SESLHVSHLLQKAKIEVNE <b>EGTKA</b>			351
GDN_FRU	296	KFTADAEVDLKGSLSALCLTDMFSSERADFRHLSA--EPLYVSTALQKAKIEVNE <b>DGTKA</b>			353
GDN_TNI	296	KFTADAEVDLKESLSALCLTDMFSVERADFRHLSA--EPVYVSKALQKAKIEVNE <b>DGTKA</b>			353
GDN_DRE	294	KFTVEQELDLETPLKALCIKIDFQNKADFRHLS--ESIYVSKALQKAKIEVNE <b>DGTKA</b>			351
		352a		380a	
A1AT_HSA	348	<b>AGAMFLEAIPMSIPPE</b> VKFNKPFVFLMIEQNTKSPLFMGKVVNP <b>TQK</b>			394
GDN_HSA	355	<b>SAATTAILIARSSPPWF</b> IVDRPFLFFIRHNPTGAVLFMGQIN <b>KP---</b>			398
GDN_MMU	354	<b>SAATTAILIARSSPPWF</b> IVDRPFLFSIRHNPTGAILFLGQV <b>NKP---</b>			397
GDN_RNO	354	<b>AVVTTAILIARSSPPWF</b> IVDRPFLFCIRHNPTGAILFLGQV <b>NKP---</b>			397
GDN_GGA	352	<b>SAATTAILIARSSPPWF</b> IVDRPFVFFIRHNPTGTILFMGQIN <b>KP---</b>			395
GDN_XTR	352	<b>SGATTAVLIARSSPRWF</b> TVDRPFLFFIRHNPTGAVLFTGQIN <b>KP---</b>			395
GDN_FRU	354	<b>SAATTAILIARSSPPWF</b> VAVDRPFLFLIRHNPTGTILFMGQIN <b>QP---</b>			397
GDN_TNI	354	<b>SAATTAILLARSSPPWF</b> TVDRPFLFLIRHNPTGTILFMGQIN <b>QP---</b>			397
GDN_DRE	352	<b>SATTSVILHARSSPPWF</b> TVDRPFLFLIRHNSSGTILF <b>AGQINKP---</b>			395







		230a							
			*						
A1AT_HSA	214	TTV	QVPMKRLGMFNIQHCK---	KLSSWVLLMKYLG	ATAIFFLDEGK--LQHLENEL	267			
E3_HSA	218	LVLQVPMH	HQTTEVNYGQFQDTAGHQVGVLEL	PYLGSAVSLFLVLE	PRDKDTPLSHIEPHL	277			
E3_MMU	215	LVLQVPMH	HQVAEVSYGQFQDAAGHEI	AVLELLYLGRVASLLLVLE	PQDKGTPLDHIEPHL	274			
E3_RNO	216	LVLQVPMH	HQVAEVSYGQFQDAAGHKVDVLE	LELLYLGRVASLLLVLE	PQDKGTPLDHIEPHL	275			
E3_GGA	212	STLKVP	TMHHTAEVNYGQFQTATQDAF	SVIELPYLGEKLSMFIVLE	SHKRTPLSHIESHL	271			
E3_XTR	218	STLKVP	TMHHTAEVNYGQFETPSL	KRFTVVELPYIGNTVSMFVVR	SDRNTPLSCIEANL	277			
E3_FRU	28	SSIKVPM	MYQATEVSFGQFRTSADQRY	TVLELPPFLGRTL	SLQVVLESERKTPSSLESQ	87			
E3_TNI	220	GAIKVP	MYLVDTVAFGQFR	TAAEQRYTVLELPPFLGRTL	SLQVVLESERKAPLASLEAQL	279			
E3_DRE	223	STVKVPM	MYQSSEVNI	HGFRLPSEQEYTVLELPPFL	DHSLRLLVALPSDRKTPSSQLEKQI	282			
		290b		323a					
A1AT_HSA	268	THDII	TKFLENEDRRSASLHLPKLSIT	GTYYDKSVLGQLGITKVES	-NGADLSGVTEEAP	326			
E3_HSA	278	TASTI	HLWTTSLRRARMDVFLPRFRI	QNFNLSILNSWCVTDLE	DPLKANLKGISGQDG	337			
E3_MMU	275	TARV	LHLWTTTLKRARMDVFLPRFKI	QNFVKSILRSWGITDLE	DPLKANLKGISGQDG	334			
E3_RNO	276	TARV	IHLWTTTLKRARMDVFLPRFRI	QNFNLSILRSWGITDLE	DPLKANLKGISGRDG	335			
E3_GGA	272	SAKT	IALWSSSLKRMKMDIFLPRFSI	QSLFDLKTVSALGIRDA	EDPITANFKGISEQAG	331			
E3_XTR	278	TSKS	MAQWANSMKRMKMDVFLPRFRL	QSHSNLRNVLPALGATDLE	DPWKA	NFKGISEQSG	337		
E3_FRU	88	TARQ	VASWDFGLRRTKMDIFLPRFKI	QNFNLSVLPALGITDA	ENPTTADFSGISAEEER	147			
E3_TNI	280	TAGQ	VASWESGLRRTKMDVFLPRFKI	QNFNLSVLPAMGISDA	ENPTTADFSGISAEEK	339			
E3_DRE	283	TARAV	GLWDTGLRRTKMDIFLPRFKM	QSKINLKPVLQSLGVSDIE	SPSADDFRGISD	TDG	342		
		352a		380a					
A1AT_HSA	327	LKL	SKAVTKAVLTIDEKGT	EAAGAMFLEAIPMSIPPE	VKFNKPPVFLMIEQNTKSPLE	385			
E3_HSA	338	FYV	SEATPKAKIEVLEEG	TRSSGATALLLKR	SRIPIFKADRPFIYFLREPNTGITV	FD	397		
E3_MMU	335	FYV	SQLTKAKMELSEEG	TRSSAA	TAVLLLR	SRTPSAFKADRPFIYFLREHSTGFV	-ESI	393	
E3_RNO	336	FYV	SEVTKAKMELSEEG	TKSCAA	TAVLLLR	SRTPSAFKADRPFIYFLREHNTGFV	-ESI	394	
E3_GGA	332	LYI	SEATPKAKIEVTEDG	TKASGATAMVLLKR	SRTPIFKADRPFTFFL	RQANTGSVLE-I	390		
E3_XTR	338	LYI	SQATPKAEIEVAEGG	TRASGVTAMVLLKR	SRMPVFKADRPFFFL	RQASSGSILE-I	396		
E3_FRU	148	LYV	SDAFTEVRIEVTEDG	TKAAAA	SMVLLKR	SRAPVFKADRPFL	FLLRQ	TSTG-----	201
E3_TNI	340	LYV	SDAFTEARIEVTEDG	TKAAAA	T-----	-----	-----	364	
E3_DRE	343	IFV	SEAFTEARIEVTEAG	TKAASA	TAMVLLKR	SRSAVFKADRPFL	ILRQISTGSL	LE-I	401
A1AT_HSA	386	GK	-----	VVNPTQK	-----	394			
E3_HSA	398	RIQ	IYQCLSSNKG	SFVHYPLK	NKHSF	424			
E3_MMU	394	GR	-----	VSNPLD	-----	401			
E3_RNO	395	GR	-----	VSNPLD	-----	402			
E3_GGA	391	GR	-----	VTNP	-----	396			
E3_XTR	397	GR	-----	VTNPLE	-----	404			
E3_FRU	-		-----		-----	-			
E3_TNI	-		-----		-----	-			
E3_DRE	402	GR	-----	VVNP	-----	407			

**Appendix 8.3.25: Alignment of pancpin sequences from vertebrates.** Gene specific features include an inhibitory RCL (red box), a conserved C-terminal extension (blue box), group V3 specific discriminating amino acid indels (indicated by \*\* or \* ) and intron indel (indicated by \$). Conserved intron positions are indicated above the alignment. Additional introns at positions 205b and 217a in PANC\_XTR are indicated by !.

A1AT_HSA	1	-EDPQGDAAQKTDTSHHDDQDHTFKNITPNLAEEFAFSLYRQLAHQSNSTNIFFSPVSIAT	59
PANC_HSA	1	MD-----TIFLWSLLLLFFGSQASRCSACKNTEFAVDLYQEVSLSH-KDNIIFSPLGITL	54
PANC_MMU	1	MN-----KTILWSFLLFFSGSQTSRATDQKIADFAVDLYKAISLSH-KNNIIFSPLGTTM	54
PANC_RNO	1	MK-----KTILWSFLLFLSGSQTSRRTMDQKNAEFAVDLYKAISLSN-KNNVIFSPLGTTV	54
PANC_XTR	1	MKRLYPGVIVAVTLCGVWITCNASRLWGDAITELAVDLSRAIHSSCTEENIIFSPLGTSL	60
90a			
A1AT_HSA	60	AFAMLSLGTRKADTHDEILEGLNFNMLTEIPEAQIHEGFQELLRTLNPQDSQLQLTTGNGLF	119
PANC_HSA	55	VLEMVQLGAKGKAQQQIRQTLKQ--QETSAGEEFFVLKSFSSAISEKKQEFFTNLANALY	112
PANC_MMU	55	LLGMVQLGAKGKAQQQILKTLRM--RGTPAGEEFSVLKSLFSAISKKKQEFFTNLASALY	112
PANC_RNO	55	LLGMVQLGAKGKAQQQIMQTLRM--QKTSTGEEFSVLKSLFSAISKKKQEFFTNLASALY	112
PANC_XTR	61	ILGMIKLARGAALSQIQQAKL--QGNQDSEEFSELKTL LAVISEENKEFTFNLANALY	118
167a\$			
A1AT_HSA	120	LSEGLKLVDRFLEDVKKLYHSEAFVTFVFGDTEEAKKQINDYVEKGTQGKIIVD--LVKELD	177
PANC_HSA	113	LQEGFTVKEQYLHGNKEFFQSAIKLVDFQDAKACAEMISTWVERKTDGKIKDMFSGEEFG	172
PANC_MMU	113	LQEGFIVKETYLHNSKEFFQSATKLVDFELDAKTSAQAISTWVESKTDGKIKNMFSEEEFG	172
PANC_RNO	113	LQEGFIVKESYLHNSKEFFQSATKLVDFELDAKTSAQAISTWVESKTDGKIKNMFSEEDFG	172
PANC_XTR	119	LQEGFQVKEQYLHNSNRDVFMSAIKLVDFQDVKASAEITISEWVQRQTHVEVK-LQDNSKLF	177
205b!                          217a!                          230a			
A1AT_HSA	178	RDIVFALVNYIFFKGRMERPFVVKDTEEEEDFHV-DQVTTVKVPMMKRLGMFNIQ-HCKKL	235
PANC_HSA	173	PLTRLVLVMAIYFKGDUKQRFKEDTQLINFTK-KNGSTVKIPMMKALLRTKYGYFSESS	231
PANC_MMU	173	PLTRLVLVMAIYFKGDUKQRFKEDTEMDFTK-KDGSTVKVPMMKALLRAQYGYFSQSS	231
PANC_RNO	173	PLTRLVLVMAIYFKGDUKQRFKEDTEMDFSK-KDGSTVKIPMMKALLRAKYGYFSESS	231
PANC_XTR	178	IQDKQLLTKLMYKVLQRGRRGIHQDAGPFTYSRTGRMLKVPMMHLQTTTKLGYFSVKN	237
290b			
A1AT_HSA	236	SSWVLLMK-VLGN-ATAIFFLPDEGKLOHLENELTHDIIITKFLNEDRRSASLH-LPKLS	292
PANC_HSA	232	LNYQVLELSYKGFDFSLIIILPAEGMDIEEVEKLITAOQILKWLSEMQEEVEISLPRFK	291
PANC_MMU	232	MTCQVLELPYKADEFSLVILLPTEDTSEEEVENQVTAPHVRRWFSELHEEEVEVSLPRFK	291
PANC_RNO	232	MTYQVLELPYKADEFSLVILLPTEDVMIEEVEKQVTARHVQKWFSELHEEEVEVSLPRFK	291
PANC_XTR	238	ASYKVLELPYKGDKFSLLLTLPAEDVEIGELEKIVTATMIKTWFADMKEEVVEISLPRFK	297

```

                                     323a
                                     |
A1AT_HSA 293  ITGTYD[KSV]GQLG[ITK]V[SNG]D[LSG]VTEEA[PLK]L[KSKAVH]KAVLTIDEK[GTE]AAGAMF 352
PANC_HSA 292  VEQKVD[K]FDVLYSLNITEIFSGGCDLSGITDSSEVYVSQVTQKVFFEINEDGSEAAATSTG 351
PANC_MMU 292  IEQKLD[K]KEALYSLNVTEIFSGGCDLSGITDSSEVYVSRVMQKVFFEINEDGSEAAASTG 351
PANC_RNO 292  IEQKLD[K]KEALYSLNVTEIFSGGCDLSGITDSSELVYVSRAMQKVFFEINEDGSEAAASTG 351
PANC_XTR 298  VEHKID[K]KKSFLNLMNITDIFMEGCDLSGITESP[NLY]ISKVFKQVFLEINEEGSEAAASTG 357

                                     352a                                     380a
                                     |                                     |
A1AT_HSA 353  LEAIPMSIPPEVKF--NRPFVFLMIEQNTKSPLFMGKVVNPTQK----- 394
PANC_HSA 352  IHIPVIMSLACSQFIANHPFLFIMKHNPTESILFMGRVTNPDTOEIKGRDLDL 405
PANC_MMU 352  INIPAIMSLTCTQFLANHPFLFILKHIRTESILFMGKVTDPDIIQTTKGRDLDL 405
PANC_RNO 352  INIPAIMSLTCTQFLANHPFLFIMKHIQTESILFMGKVTDPDIIHTVKGRDLDL 405
PANC_XTR 358  MQVSAMSMSEH-F AANRPFLFFIRHIQSGMILFMGKVMNPDFYDALGRDVESL 410

```

**Appendix 8.3.26: Alignment of neuroserpin sequences from vertebrates.** Gene specific features include an inhibitory RCL (red box), a conserved C-terminal extension (blue box), an N-glycosylation site (cyan box), group V3 specific discriminating amino acid indels (indicated by \*\* or \*) and an intron indel (indicated by \$). Conserved intron positions are indicated above the alignment. Gaps in the genomic region within NEUS\_TNI gene are responsible for gaps in coding region and the absence of an intron at position 90a (indicated by ?).

A1AT_HSA	1	EDPQGDAAQKTDTS----HHDQDHTFNKITPNLAEF	AFSLYRQLAHQSNSTNIFFSPVS	56												
NEUS_HSA	1	MAFLGLFSLLVLC-----SMATGATFP	EEAIADLSVNMYNRLRATGEDEMIIFSPLS	52												
NEUS_MMU	1	MTYLELLALLALC-----SVVTGATFP	DETITEWSVNMYNHLRGTGEDEMIIFSPLS	52												
NEUS_RNO	1	MAYLGLLSLVALC-----SLVTGATFP	DETIAEWSVNVYNHLRATGEDEMIIFSPLS	52												
NEUS_GGA	1	MYFLGLLSLLVLP-----SKAFKTNFP	DETIAELSVNVYNQLRAAREDEMIIFCPLS	52												
NEUS_XTR	1	MHHLSSLALIVMC-----ALVFGTSV	HDETVNEFSIKVYHKLRTIEDIIFIIFSPLS	52												
NEUS_FRU	1	MLSDDTASSQ-----DVYPECNVP	EDPTAEFSVRLYHLLQAGGDQDNIIFSPLS	49												
NEUS_TNI	1	MSTLDDLPSLLLLLLTVLLRCH	HCRETDPVEDALADFSVRLYQQLQAGGEQDNLVFSPLS	60												
NEUS_DRE	1	MLLLVVLPLLLLLRC-----CFCCASD	VPEDVTAEF SVRLYHQLQISSGEENIIFSPLS	54												
			90a?													
A1AT_HSA	57	IATAFAMLSLGTKADTHDEI	LEGLNFNLTEIPEAQIHEGFQELLRTLNQPD	SQLQLTTGN 116												
NEUS_HSA	53	IALAMGMELG	AQGSTRKEIRHSMGY--DSLKNGE	EFSFLKEFSNMVTAKESQYVMKIAN 110												
NEUS_MMU	53	IALAMGMELG	AQGSTRKEIRHSMGY--EGLKNGE	EFSFLRDFSMMASAEENQYVMKLAN 110												
NEUS_RNO	53	IALAMGMELG	AQGSTLKEIRHSMGY--ESLKS	GEEFSFLRDFSSMVASAEEGQYVMKIAN 110												
NEUS_GGA	53	IAIAMGMIELG	AHGTTLKEIRHSLDGF--DSLKNGE	EFTFLKDLSDMATTEESHYVLNMAN 110												
NEUS_XTR	53	TAIALGMVELG	ARGSSLKEIRHVLGYS--DKLNGE	EFSLLKDLNMLTAQEKHYVLSIAN 110												
NEUS_FRU	50	VAVALGMVGLG	ARGVSL	EQIRKVGAF--SHLVSGGEFS	LLOQLTAPLADKEAHHVVFAN 107											
NEUS_TNI	61	VAVALGMVR	SCSR-----	TPDDELQTSAAAPLFP 89												
NEUS_DRE	55	VALALGMVELG	ARGSSLQ	EIRQAVGY--SHFREDEEFS	LLRNLSQALSTDEEQYVVRLAN 112											
			167a\$													
A1AT_HSA	117	GLFLSEGLKLVDRFL	EDVKKLYHSEAF	TVNFGDTEEA	AKKQIMDYVEKGT	IQGKIIVD--LVK 174										
NEUS_HSA	111	SLFVQNGFHVNEE	FLQM	MKKYFNAAVN	HVDFSONVAVANY	IMKWVENNTN	NLVKDLVSPR 170									
NEUS_MMU	111	SLFVQNGFHVNEE	FLQ	LKMYFNAAVN	HVDFSONVAVANS	IMKWVENY	TNSLLKDLVSP	170								
NEUS_RNO	111	SLFVQNGFHINEE	FLQM	MKKYFNAAVN	HVDFSENVAVANY	IMKWVENY	TNSLLKDLVSP	170								
NEUS_GGA	111	SLYVQNGFHVSEK	FLQL	VKKYFKA	EVENIDFSQSA	AAVATHINKW	VENHTN	NMIKDFVSSR 170								
NEUS_XTR	111	SLYLQNGFHISDK	FIQL	MKKYFKA	EVENVDFSQS	TVANHIM	WVENHT	NNRIRDLVTAD 170								
NEUS_FRU	108	ILFLQQGVT	FNPEFLH	LMKKYFKA	HVMVD	FSQSA	AAVAQIMTW	VENHTESMIRELMSAE 167								
NEUS_TNI	90	SLFLQQGVT	FNPEFLR	LMRKYFKA	EVETVDFSQ	PAVAQ	IMSW	VENRTGKIGELLAAE 149								
NEUS_DRE	113	SLFLQSGV	HFNEDFL	QLMKKYF	RAE	TVTDFSQ	S	TAVAEIRINSW	VLNHTESKIQNLVSAE 172							
			230a													
A1AT_HSA	175	ELDRDITVF	ALVMYIFFKGR	MERPF	EVKDT	TEEDFH	HVDQVT	TVKVPM	MMKRLGMFNIQHCKK 234							
NEUS_HSA	171	DFDAATYL	ALINAVYF	KGN	WKSQFRP	ENTRTFS	FTKDD	EVQIP	PMYQQGEFYG-EFS 229							
NEUS_MMU	171	DFDGVNL	ALINAVYF	KGN	WKSQFRP	ENTRTFS	FTKDD	EVQIP	PMYQQGEFYG-EFS 229							
NEUS_RNO	171	DFDAVTHL	ALINAVYF	KGN	WKSQFRP	ENTRTFS	FTKDD	EVQIP	PMYQQGEFYG-EFS 229							
NEUS_GGA	171	DFSALTHL	VLINAVYF	KGN	WKSQFRP	ENTRTFS	FTKDD	ETEVQIP	PMYQQGEFYG-EFS 229							
NEUS_XTR	171	DFTNLT	KLVLN	AVYF	KGN	WKSQFRP	ENTRTFS	FTKDD	EVQIPPMYQKGEFYVGEFT 230							
NEUS_FRU	168	DVSGIT	RLLV	N	AVYF	RGS	W	KIQFRP	ENTRTFSFSKDDG	EVQTO	PMYQQGDF	FYG-EFS 226				
NEUS_TNI	150	DLSTI	RLLV	N	AVYF	RGS	W	KIQFRP	ENTRTFSFSR	DDGSE	VHT	PMYQQGDF	FYG-EFS 208			
NEUS_DRE	173	DFSSS	I	M	L	V	N	AVYF	RGS	W	KIQFRP	ENTRTFS	TRDDG	SEVQTL	PMYQQGDF	FYG-EFS 231

\*

A1AT_HSA	235	LSS-----WVLLMKYLGN-ATAIFFLPD-EGKLOHLENELTHDIIITKFLNEDRRSAS	285
NEUS_HSA	230	DGSNEAGGIYQVLEIPYEGDEISMMLVLSRQEVPLATLEPLVKAQLVEEWANSVKKQKVE	289
NEUS_MMU	230	DGSNEAGGIYQVLEIPYEGDEISMMLALSQRQEVPLATLEPLLKAQLIEEWANSVKKQKVE	289
NEUS_RNO	230	DGSNEAGGIYQVLEIPYEGDEISMMLVLSRQEVPLATLEPLLKPQLIEEWANSVKKQKVE	289
NEUS_GGA	230	DGSNEAGGIYQVLEIPYEGDEISMMLVLSRQEVPLVTLEPLVKASLINEWANSVKKQKVE	289
NEUS_XTR	231	DGSNEAGGVYQVLELPYEGDEISLIIVLSRQEVPLATLEPLLKAPLIEEWANSVKKQKVE	290
NEUS_FRU	227	DGSQEAGGMYQVLEMPYEGEDLSMMIVLPRQEVPLSSLEPIIKAPLLEEWANNVKLQKVE	286
NEUS_TNI	209	DGSQEAGGVYQVLEMPYEGEDMSMMIVLPRQ-----AQDWANNVKLQKVE	253
NEUS_DRE	232	DGTTEAGGVYQVLEMLYEGEDMSMMIVLPRQEVPLASLEPIIKAPLLEEWANNVKRQKVE	291

		290b		323a	
A1AT_HSA	286	LHLPKLSITGTYDLKSVLGQLGITKVFNSGADLSGVTEEAPLKLSKAVHKAVLTIDEKGT	345		
NEUS_HSA	290	VYLPRFTVEQEIDLKDVLLKALGITEIFIKDANLTGLSDNKEIFLSKAIHKSFLEVNEEGS	349		
NEUS_MMU	290	VYLPRFTVEQEIDLKDLKALGVTEIFIKDANLTAMSDKKELFLSKAVHKSCIEVNEEGS	349		
NEUS_RNO	290	VYLPRFTVEQEIDLKDLKALGVTEIFIKDANLTAMSDKKELFLSKAVHKSFIEVNEEGS	349		
NEUS_GGA	290	VYLPRFTVEQEIDLKDVLLKGLGITEVFSRSADLTAMSDNKELYLAKAFHKAFLEVNEEGS	349		
NEUS_XTR	291	VYLPRFKVEEVVNLKDVLMQLGITKIFSGEADLSAVSDSKDLFVAKAVHKSFLEVNEEGS	350		
NEUS_FRU	287	VYLPRFKMEQKIDLRKTLQELGIKSVFSTEADLSSMIAGKDLYIGKAVQKAYLEVTEEGL	346		
NEUS_TNI	254	VYLPRFKVEQKMDLRKTLQELGIKSIFSTEADLSAMTDGKDLYIGKAVQKAYLEVTEEGS	313		
NEUS_DRE	292	VYLPRFKVEQKIDLRRESLQQLCIRSFISKDADLSAMTDGQDLFIGKAVQKAYLEVTEEGA	351		

		352a		380a	
A1AT_HSA	346	EAAGAMFLEAIPMSIPPE--VKFNKPFVFLMIEQNTKSPLFMGKVVMPPTQK-----	394		
NEUS_HSA	350	EAAAVSGMIAISRMAVLYPQVIVDHPFFFLIRNRRTGTILFMGRVMHPETMNTSGHDFEE	409		
NEUS_MMU	350	EAAAASGMIAISRMAVLYPQVIVDHPFLYLIRNRKSGIILFMGRVMMPETMNTSGHDFEE	409		
NEUS_RNO	350	EAAVASGMIAISRMAVLYPQVIVDHPFLFLIKNRKTGTILFMGRVMHPETMNTSGHDFEE	409		
NEUS_GGA	350	EAAAASGMIAISRMAVLYPQVIVDHPFFFLVRNRRTGTVLFMGRVMHPEAMNTSGHDFEE	409		
NEUS_XTR	351	EAAAASGMIAISRMAVLYPQVIVDHPFFVIRNRKTGSVLFMGRVMHPETLHTIGHDFEE	410		
NEUS_FRU	347	EGAVGSGLVALTRTLVLYPQVMADHPFFFVIRDRRTGSILFMGRVMTDPVIDATGPDFDS	406		
NEUS_TNI	314	EGAAGSGMMALTRTLVLYPQVMADHPFFFVVRERRTGSILFMGRVTTPEVIDAGDGDFDS	373		
NEUS_DRE	352	EGAAGSGMIALTRTLVLYPQVMADHPFFFIRNRKTGSILFMGRVMNPELIDPFDMNFDI	411		

A1AT_HSA	-	-	-
NEUS_HSA	410	L	410
NEUS_MMU	410	L	410
NEUS_RNO	410	L	410
NEUS_GGA	410	L	410
NEUS_XTR	411	L	411
NEUS_FRU	407	L	407
NEUS_TNI	374	L	374
NEUS_DRE	412	M	412

**Appendix 8.3.27: Alignment of PEDF sequences from vertebrates.** Gene specific features include non-inhibitory RCL (red box) and a nuclear localization signal (brown box). Conserved intron positions are indicated above the alignment.

A1AT_HSA	1	-----	EDPQGDAAQKTDTSHHDDHPTFNKITPNLA	31
PEDF_HSA	1	<b>MQALVLLL</b> CIGALLGHSSC	NPASPEEGSPDPDSTGALVEEED-PFFKVPVNKLAAAVS	59
PEDF_MMU	1	<b>MQALVLLL</b> WTGALLGHGSSC	NVPSSSE-GSPVPDSTGEPVEEED-PFFKVPVNKLAAAVS	58
PEDF_RNO	1	<b>MQTLVLLL</b> WTGALLGHGSSC	NVPDSSQ-DSPAPDSTGEPVVEEDPFFKAPVNKLAAAVS	59
PEDF_GGA	1	<b>MQIPAVLLL</b> GLLTIPSKSC	NSPAGQN--SPTTDGTGVEVEEED-PFYKTPINKLAAAVS	57
PEDF_XTR	1	<b>MKIYLALL</b> FTGSFLSYTSAQN	-----AADEVPTVEEED-PFYKSPINRLASSAS	49
PEDF2_FRU	1	<b>MKGTTFL</b> LVIGVILRFCAQS	-----ETEAEESVAEEHVELFTTAQTKMGAATS	50
PEDF2_TNI	-	-----	-----	-
PEDF2_DRE	1	<b>MKKIVLLV</b> GLWSLLSLSHAQ	-----LADTTDAEGEEEAVDLFTTPRTKLAAATS	49
67a				
A1AT_HSA	32	EFAFSLYRQLAHQSNSTNIFFSPVSIATAFAMLSL	GTKADTHDEILEGLNFNLTEIPEAQ	91
PEDF_HSA	60	NFGYDLYRVRSSMSPTTNVLLSPLSVATALSALS	SLGAEQRTESIHRALYYDL--ISSPD	117
PEDF_MMU	59	NFGYDLYRLRSSASPTGNVLLSPLSVATALSALS	SLGAEHRTESVIHRALYYDL--ITNPD	116
PEDF_RNO	60	NFGYDLYRLRSGAVSTGNILLSPLSVATALSALS	SLGAEQRTESVIHRALYYDL--INNPD	117
PEDF_GGA	58	NFGYDLYRQSSRTATANVLLSPFSLATALSGLS	LGAGERTEDVISRALFYDL--LNKAE	115
PEDF_XTR	50	NFGYDLYRMQANKNPNSNIIISPLSIATSLSSL	SLGGQRTESLIQRSLYYDL--LNDPE	107
PEDF2_FRU	51	DFGYNLFRALASQEAGNWFAPISVSAVLTQLS	MGSEHAQSOLFRAIRYHT--LHDPQ	108
PEDF2_TNI	-	-----	-----	-
PEDF2_DRE	50	DFGYNLFRQLASRDTKASVFLSPMSISA	AFTQLSMGASERAEKQIYRALRYHT--LQDSQ	107
123a				
A1AT_HSA	92	IHEGFQELLRTLNPDSQLQLTTGNGLFLSEGL	KLVDKFLLEDVKLYHSEFTVNF	151
PEDF_HSA	118	IHGTYKELLDTVTAPQK-- <b>NLKSASRIVFEK</b>	<b>KLRIKSSFWAPLEKSYG</b> -TRPRVLTGNPR	174
PEDF_MMU	117	IHSTYKELLASVTAPEK-- <b>NLKSASRIVFER</b>	<b>KLRVKSSFWAPLEKSYG</b> -TRPRILTGNPR	173
PEDF_RNO	118	IHSTYKELLASVTAPEK-- <b>NFKSASRIVFER</b>	<b>KLRVKSSFWAPLEKSYG</b> -TRPRILTGNPR	174
PEDF_GGA	116	VHNTYKDLLASVTGPEK-- <b>SLKSASRIIVEK</b>	<b>RLRVKSTFHSQLEKSYR</b> -MRLRALSGNTQ	172
PEDF_XTR	108	VHATYKDLLASFTSQAS-- <b>GLKSTWRIMLER</b>	<b>RRLRLRMDFVTQVEKFG</b> -NPKKVLTGSTR	164
PEDF2_FRU	109	LHDTLKNILATVKAPGK-- <b>GLSTAARLYSR</b>	<b>RRLRLKQEF</b> LALVENQYN-VRPKAVLG---	162
PEDF2_TNI	-	-----	-----	-
PEDF2_DRE	108	LHDTLRDLSSLRASAK-- <b>GFKSAERILLAR</b>	<b>KRLRLLEYLNSVEKQYG</b> -ERPQILAGGA-	163
192a				
A1AT_HSA	152	EAKKQINDYWEKGTQGGKIVDLVK-ELDRD	IVFALVNIYIFK-----GKMERPFV	205
PEDF_HSA	175	LDLQEIINNWWQAQMKGKLARSTK-EIPDE	ISILLGVAYFK-----GQWTKFDSR	228
PEDF_MMU	174	VDLQEIINNWWQAQMKGKLARSTR-EMPS	ALSILLGVAYFK-----GQWTKFDSR	227
PEDF_RNO	175	IDLQEIINNWWQAQMKGKLARSTR-EMPS	ALSILLGVAYFK-----GQWTKFDSR	228
PEDF_GGA	173	LDLQEIINNWRQOTRGRILRFMK-DMPTD	VSILLAGAAYFKAFKKTGTGKTKFDTKR	231
PEDF_XTR	165	LDLQEAANDFIQKQTQGGKVVVFK-EIPT	SVSILLGTTYLK-----GQWTKFNP	218
PEDF2_FRU	163	KDIKEVNDWVSSQITGRKVOGFLASNF	PRNSGANAVSAAYFK-----GKMTRE	216
PEDF2_TNI	-	-----	-----	-
PEDF2_DRE	164	RDLKTVNDWFKQQTGGKVDQVVPSP	LPRTALLPVGSAAYFK-----GKMITRE	217



238c

A1AT_HSA	206	EDFHVDQVTTVKVPMMKR-LGMFNIQHCKKLSSWVLLMKYLGNATAIFFLPDEG--KLQH	262
PEDF_HSA	229	EDFYLDEERTVRVPMMSDPKAVLRYGLSDLSCKIAQLPLTGSMSIIFFLPLKVTQNLTL	288
PEDF_MMU	228	QDFHLDEDRTVRVPMMSDPKAILRYGLSDLNCKIAQLPLTGSMSIIFFLPLTVTQNLTM	287
PEDF_RNO	229	QDFHLDEDRTVRVPMMSDPKAILRYGLSDLNCKIAQLPLTGSMSIIFFLPLTVTQNLTM	288
PEDF_GGA	232	KDFHLDEDRTVQVSMMSDPKAILRYGFDSELNCKIAQLPLTEGVSAMFFLPKTKVTQNMML	291
PEDF_XTR	219	REFHLDEQTSVTVPMMSSKNIPVRYGLSDFNCKIVQLPLTGGVSIIMFFLPLNTVTQNLTM	278
PEDF2_FRU	217	DTFQVADGAPVSIIPMMKQDNYPVKMGVSDLKCTIAQIPMQDDVSMFLPLPDDLSSNMTQ	276
PEDF2_TNI	-	-----	-
PEDF2_DRE	218	ETFRRDGQAPAVIPMMEQENYPVKMGIDSDLGCTIAQVPMEDGVSMYFFLPDEVTVQNLTL	277

307a

A1AT_HSA	263	LENELTHDIIITKFLNEDRRSASLHLPRLSITGTYDLKSVLGQLGITKVFSSNGADLSGVT	322
PEDF_HSA	289	IEESLTSEFIHDIDRELKTVQAVLTVPKLKLSEGEVTKSLQEMKQSLFDS-PDFSKIT	347
PEDF_MMU	288	IEESLTSEFIHDIDRELKTIQAVLTVPKLKLSEGEVTKSLQDMKQSLFES-PDFSKIT	346
PEDF_RNO	289	IEESLTSEFVHDIDRELKTIQAVLTVPKLKLSEGEVTVNSLQDMKQSLFES-PDFSKIT	347
PEDF_GGA	292	IEESLTSEFVHDVRELKTVHAVLSLPRKLKLNVEEALGNTVKETRLQSLFES-PDFTKIS	350
PEDF_XTR	279	IEEGLTSEFVHDIDQALQPINLVLSVPRKLKLNVEEALKEALQESKQSLFES-PDFSKIS	337
PEDF2_FRU	277	LEESLTAEFVQDLSMTLLPAQVSLTLPVLRLSYSKDLLPLLGDLGLSDWLLN-TELQKIS	335
PEDF2_TNI	-	-----	-
PEDF2_DRE	278	IEEALTAEFVQDLSNSLHTVKVLLTLPVIKLSYKTNLLPSLSDLGLSEWLAE-TDLTKIT	336

A1AT_HSA	323	EEAPLKLKSKAVHKAVLTIDEKGTAAAGAMFLEAIPMSIPPEVKFNKPFVFLMIEQNTKSP	382
PEDF_HSA	348	GK-PIKLTQVEHRAAFEWNEEGAGTTPSPGLQPAHLTFPLDYHLNQPFIFVLRDQDTGAL	406
PEDF_MMU	347	GK-PVKLTQVEHRAAFEWNEEGAGSSPSPGLQPVRLTFPLDYHLNQPFIFVLRDQDTGAL	405
PEDF_RNO	348	GK-PVKLTQVEHRAAFEWNEEGAGTSSNPDLQPVRLTFPLDYHLNRPFFIFVLRDQDTGAL	406
PEDF_GGA	351	AK-PIKLSHVQHKAVLELNEDGEKSTPNPGVNAARLTFPIEYHVDKPFLLVLRDQDTGTL	409
PEDF_XTR	338	SK-PLKLSYVVKATLELNEDGAETAPKP-EDSHRNYFPLEYHLDHPFLFVLRANDNGAL	395
PEDF2_FRU	336	PQ-PVKLTSVNHKVVMEVMAPEGNQYPPSS-SAPTHLS----YRADRPFLYLIRDETSGAL	389
PEDF2_TNI	1	-----LSSVRHKVVMEVMAPEGNQYPPSS-SAPSHLS----YRVDRPFLYLIRDETSGAL	49
PEDF2_DRE	337	SQ-PVKLNAVHVKVLETAPEGAEYASTT-PSATGQSLGLSYRVDRPFLFVLRDEPSGAL	394

A1AT_HSA	383	LFMGKVVNPTQK--	394
PEDF_HSA	407	LFIGKILDPRGP--	418
PEDF_MMU	406	LFIGRILDPSST--	417
PEDF_RNO	407	LFIGRILDPSST--	418
PEDF_GGA	410	LFIGKILDPRTHEF	423
PEDF_XTR	396	LFIGKVMDPKGFSF	409
PEDF2_FRU	390	LFIGRVVNPTGLTI	403
PEDF2_TNI	50	LFIGRVVNPTG---	60
PEDF2_DRE	395	LFIGKVLNPSDL--	406

**Appendix 8.3.28: Alignment of  $\alpha_2$ -antiplasmin from vertebrates.** Gene specific features include inhibitory RCL with putative overlapping two reactive sites P2-P1 and P1-P1' (red boxes), cysteine residues [blue box] and N- and C-terminal extensions. Conserved intron positions are indicated above the alignment. Predicted (by GENSCAN and FGENESH) low complexity regions at the intron at position 67a for A2AP2\_FRU and A2AP2\_TNI and at the intron at position 123a for A2AP1\_FRU were deleted manually.

A1AT_HSA	1	-----EDPQGDAAQKTDTSHHDDQHP-----	21
A2AP_HSA	1	MALLWGLLVLSWSCLQGPCSVFSPVSAEPLGRQLTSGPN-----	40
A2AP_MMU	1	MALLRGLLVLSLSCLQGPCFTFSPVSAVDLPGQQPVSEQA-----	40
A2AP_RNO	1	MALLRGLLVLSLSCLQGPSSMFPVSAVDLPGQQPVSEQA-----	40
A2AP_GGA	1	MVLLWGLLLLLLSLSALHSHPRPLAHAVE-----	28
A2AP_XTR	-	-----	-
A2AP1_FRU	1	MALWFRRLLLNRRK-----	15
A2AP2_FRU	1	MNLHLLLLLLLCLCCPGLTE-----PTPGVTNA-SVPVSD-----	33
A2AP2_TNI	1	MKLPHLLLLLLLCLCCPGLTE-----	19
A2AP2_DRE	1	MNLCFLAFLLLCYSKQGWTD--DVTDPEDGVDPVIPLIPLTPSKPISDLKATLDPTVTEELTV	61
A2APL2_PMA	-	-----	-
A2APL1_PMA	1	MASLSPLFVSLLVLTILSLGFADHHGHKTPGAPPVSAATAIS-----	40
A1AT_HSA	-	-----	-
A2AP_HSA	41	-----QEQVSPLTLLKLG-----	54
A2AP_MMU	41	-----QQKLPLPALFKLDN-----	54
A2AP_RNO	41	-----QQKLPLALLKLG-----	54
A2AP_GGA	29	-----QQHSSDKAVDLKNLKSGGDEESALPEAIPTLLDAKLADTWE	70
A2AP_XTR	1	-----LVLDDNNEVEQES-----	13
A2AP1_FRU	16	-----PENNSATKVPAAANT-SQPDS-----	35
A2AP2_FRU	-	-----	-
A2AP2_TNI	20	-----GLTEPTPGATDT--QVPVSD-----	38
A2AP2_DRE	62	NPDGLEADPPTPGPSGG---QKEGS-----	83
A2APL2_PMA	-	-----	-
A2APL1_PMA	-	-----	-
A1AT_HSA	22	-----TFNK	25
A2AP_HSA	55	-----QEPGGQTALKSPPGVCS-RDP-----TPEQTHR	81
A2AP_MMU	55	-----QDFGDHATLKRSPGHCK-SVP-----TAEETRR	81
A2AP_RNO	55	-----QDLGDHATLKRSPGDCK-SAP-----TTEETRR	81
A2AP_GGA	71	TYGTTPSISTSAETEEEEESPGDKATAGAVSCHEQEPSGKTLSSSEEEGEGEEESCDITWKKSQK	133
A2AP_XTR	14	-----CDENAS-----LEEMRK	25
A2AP1_FRU	36	-----SEDGRNEDYCLIGRSL-----SRE--A	56
A2AP2_FRU	34	-----EEDTTEAHNCR-TQLVS-----TEEQRS	55
A2AP2_TNI	39	-----EEDSKKADSCG-GQLFS-----SEERRS	60
A2AP2_DRE	84	-----SEEELDTLCDGDMT-----GKQIKRT	104
A2APL2_PMA	1	-----MTR	3
A2APL1_PMA	41	-----PFVVSR	46



67a  
|

A1AT_HSA	26	ITPNLAEFAFSLYRQLAHQSNSTWIFFSPVSIATAFAMLSLGTKADTHDEILEGLNFNLTEIP	88
A2AP_HSA	82	LARAMMAFTADLFLSLVAQTSTCPNLLSPLSVALALSHLALGAQNHTLQRLQQVLAHAGSG---	141
A2AP_MMU	82	LAQAMMAFTTDLFLSLVAQTSTSSMLVLSPLSVALALSHLALGAQNQTLHSLHRVLAHMNTG---	141
A2AP_RNO	82	LSQAMMAFTTDLFLSLVAQTSTSSMLVLSPLSVALALSHLALGARNQTLLENLQRVLAHMNMG---	141
A2AP_GGA	134	LANGLMRFSTDLLREVQESNGNMVILSPLSIALALSNLALGAANQTEKRLEAMHLESV---	193
A2AP_XTR	26	FSQAITFESIDLLKEIDPESKKPSVVMSPFSIALGLLQSLGAGKEMQNKLME TLHVESL---	85
A2AP1_FRU	57	IAAAIQKLGVLQONLEATPEQPMIIISPLSISLALSQLALGAVNETRELLMHHLHERAL---	116
A2AP2_FRU	56	LGGAEQLGLQLENLPIVVSQPMVILSPLSVALALAHLTGGAHNETENLLKALHAHNL---	115
A2AP2_TNI	61	LGGTIERLGLQLENLPIVQPMIILSPLSVALALAHLTGGAHNETEQLLKTLHAHNL---	120
A2AP2_DRE	105	IGNGIMKLGLLFLENLKPSPDQPMVIFSPLSLSVALSQLALGATNDTEELLHHLHADAL---	164
A2APL2_PMA	4	LAMSQANFGFDLYRAVAQESPGEMIFMSPLTTSLVLA MLTAGAHGATEQALARALYFTHLRN-	65
A2APL1_PMA	47	LAGSQGDFGFGQFFHKLGEASPGQNVLFSPLTTSAA LMMLLAGSGDKTETQLTNALRLQFLRD-	108

123a!  
|

A1AT_HSA	89	EAQIHEGFQELLRTLNPQDSQLQLTTGNGFLFSEGLKLVDFLEEDVKKLYHSEAFVNF GDTE	151
A2AP_HSA	142	-PCLPHLLSRLCQDLGPG--AFRLAAR--MYLQKGFPIKEDFLEQSEQLFGAKPV-SLTGKQE	198
A2AP_MMU	142	-SCLPHLLSHFYQNLGPG--TIRLAAR--IYLQKGFPIKDDFLEQSERLFGAKPV-KLTGKQE	198
A2AP_RNO	142	-SCIPHLLSHFCQNLNPG--TIRLAAR--IYLQKGFPIKDDFLEQSEKLFAGKPV-KLTGRQE	198
A2AP_GGA	194	-PCLHMLSSLRRLRAGA--TLSLASR--VYLQKGYEVKEEFLEESEKFGYAKPV-TLSGSSE	250
A2AP_XTR	86	-HCLHNKLTVRKELSKS--ILRTATR--IYLKKGFIKDSFLKSSEKWYGSKPL-HLGGSKK	142
A2AP1_FRU	117	-PCYHESLHMLAGLRKN--DLQIATQ--IFLRQ FQPKQDFVNKSRHLYGSEPA-ELKS---	170
A2AP2_FRU	116	-PCYHHILGGLLAHFKNT--SLEVATR--MYLRPGSEIKRSFVEESLARYQSRPV-PLVS---	169
A2AP2_TNI	121	-PCYHHILGHLLPHFKNT--SLEVATR--MYLRPGFEVKLSFVEESLARYQSRPV-PLVS---	174
A2AP2_DRE	165	-PCYHTALSSLLRNFRKR--SMP IASR--IYLKTGFKAQSDFMEDSQLYDSEPA-TLTD---	218
A2APL2_PMA	66	-PNLHGTFRDLIQKITSGKSVSKMAAR--IFAARNIKIKKDFLDVVEQNYHAKPE-NLNGPEE	124
A2APL1_PMA	109	-PNPQASFQALVSKLHHGRDSTNIAAR--IFTAKHATIKQCFLDAVEKYYKAKPQ-KLIGNMK	167

192a  
|

A1AT_HSA	152	EAKKQINDYWEKGTQGGKIVDLVKELDRD TVFALVNYIFFKGKWERPFEVKDTEEDFHVDQVT	214
A2AP_HSA	199	DDLANINQWVKEATEGKIQEFLSGLPEDTVLLLLNAIHFGGFURNKFDPSLTQRDSFHLDEQF	261
A2AP_MMU	199	EDLANINQWVKEATEGKIEDFLSELPDSTVLLLLNAIHFGGFURTKFDPSLTQKDSFHLDERF	261
A2AP_RNO	199	EDLMNINKWVKEATEGKIEDFLSELPDNTVLLLLNAIHFGGFURTKFDPSLTQKDSFHLDEQF	261
A2AP_GGA	251	DDLTAINKWVKEATNGQIPTFLQQLPGD TVMLLLNAIHFGGFURNKFDASFTAPDAFHLSDDF	313
A2AP_XTR	143	KNLESINKWVKDITEGQIPHFLSDLPQDVLLILLNAMHFKGVWKNFTDPSLTSEDSFYINDDM	205
A2AP1_FRU	171	--LQQINDWVYNATNGKMPQFLSALPLNVLMMLINAVHFKGDVVARFDPRFTSRGAFYLDNND	231
A2AP2_FRU	170	--VEEVNQWVENATNGHISNFLERIPHGVVLMMLINAVYFKGEWQTRFDPLETSKGVFYLDNKN	230
A2AP2_TNI	175	--VEEVNHWIENATNGHISNFLSIPHNVLMMLINAVYFKGEWQTRFDPTMTFKGVFYLDNKN	235
A2AP2_DRE	219	--VNDVNEWVKKVTNGHISEFLSSLPSSAVMMLINAMHYKGEWLTFRFDPHFTSTENFYIDENQ	279
A2APL2_PMA	125	KDLKRINSWEEKTDGKIKDFLKPENLRMLLLSAIYFK-AMINPLLDEGICKPLPSLMVLQ	186
A2APL1_PMA	168	EDVALINKWVAEKTEGHIPDFVKELPEELQLFIVSAIFFKGKULKPFQVESTSPRPFHLSPSN	230

238c

A1AT_HSA	215	TVKVPMMKRLG-MFNIOHCKKLLSSWVLLMKYLGNATAIFFFLPDEG--KLQHLENELTHDIITK	274
A2AP_HSA	262	TVPVEMMQARTYPLRWFLLEQPEIQVAHFPPKNNMSFVVVLPVTHFEWNVSQVLANLS---WDT	321
A2AP_MMU	262	TVSVDMMHAVSYPLRWFLLEQPEIQVAHFPPKNNMSFVVVMPVTFYFENNVSEVLANLT---WDT	321
A2AP_RNO	262	TVPVAMMHAQSYPLRWFLLEQPEIQVAHFPPQNNMSFVVIMPTFYFGWNVSEVLANLT---WDT	321
A2AP_GGA	314	VVSVEMMKAQRYPLSWFTLESQDIQVAKFPFKGNMSFVVIVPMQYTWNTSHVLENFP---YGQ	373
A2AP_XTR	206	SVPVEMMSAQKYPFSWFFLESIESQVAKFQFKGNMSFVVVLMVYSSTWNLKLLANFS---QSD	265
A2AP1_FRU	232	MIDVEVMEDAKHPLSLFIDNEMDAQVARFRFRKLMSELLVVMPTSSQVSVASLLPKLN---VSK	291
A2AP2_FRU	231	SVSVDMMKSFQYPPFRLLDPELKSQVASFQFKGNTSFLVVMPLPGSGNVSSLLPKLN---ISD	290
A2AP2_TNI	236	SVSVDMMSSQYPPFRLLDPELKAQVASFQFKGNTSFLIVMPPVPGIGNVSSVLPKLN---ISD	295
A2AP2_DRE	280	IVNVDMMLGPKYPLSVFTHHELDAQVARFPFKGDRSLLVVMPTSGHVNVSAIAAKLN---ISD	339
A2APL2_PMA	187	CNDPPTILHR----RRFNLSISWVDVAQLEFQGDKNMLIFLPLDEVTTNLTALEQSLSSDLLLN	245
A2APL1_PMA	231	ETQVPTMFAAGYPIKKGRHPSLPVTVAKIQFQGNLSLLLFVPEDAVSTNLSALESSLSSQLVTT	293

307a

A1AT_HSA	275	FLENEDR-RSASLHLPRLSITGTYDLKSVLGQLGITKVFSGADLSGVTEEAPLKLKSAVHKA	336
A2AP_HSA	322	LHPLVWERPTKVRLPKLYLKHQMDLVATLSQLGLQELFQA-PDLRGISEQS-LVVSGVQHQS	382
A2AP_MMU	322	LYHPSLQERPTKVWLPKHLHQQLDLVATLSQLGLQELFQG-PDLRGISEQN-LVVSSVQHQS	382
A2AP_RNO	322	LYQPSMREKPTKVRLPKLHLEQHLDLVATLSKLGQLQELFQS-PDLRGISDQS-LVVSSVQHQS	382
A2AP_GGA	374	LCRLFPKEVPTTVKIPKITLDYQLELNSVLSHMGLQELFIS-PNLQKISDEP-LFVSSIQHQS	434
A2AP_XTR	266	LYSRFPREKNTNLKMPKLNLDYKLELRNPLTNLGLGQLFTN-PDLSGITNEA-LVVSSIQHQS	326
A2AP1_FRU	292	LYSRLPKERAVQVKVPKFKLEYSQELQEVFTKIGLGEIFSR-PNLAEIADGP-LLVSSVMHKS	352
A2AP2_FRU	291	LYRRLPQEKIMHVSLPKVKLQYRQELQEALTSMGLCSLFSG-PNLSGISDYP-LRVGSVRHAS	351
A2AP2_TNI	296	LYRRLPQERIMHVNLPKVKLQYRQELQEALTSMGLGSLFSG-PDLSGITDHP-LRVGSVRHAS	356
A2AP2_DRE	340	LYSRLPRERNMQVKLPKFKLDFNQDLQEAMTSMGLGKLFSSH-PKLDRITEVP-LFVSSVQHMS	400
A2APL2_PMA	246	LTESELKSGNRIVYLPRLRLKMKKDLSTALNHLGLDDLQFMA-PDFNKISEEP-LLVSAVTHVA	306
A2APL1_PMA	294	LVEETLVQKKIDLYLPLISLDVESNIEQKLTIDIGLDLTKT-PDLSKISDIP-LRVSKVIHRA	354

A1AT_HSA	337	VLTIDERGTEAAGAMFLEAIPMSIPPEVK---F-NKPFVFLMIEQNTKSPLFMCKVVMNPTQK-	394
A2AP_HSA	383	TLELSEVGVAAAAATSIAMSRMSLSSF-----SVNRPFLFFIFEDTTGLPLFVGSVRNPNPSA	440
A2AP_MMU	383	TMELSEAGVAAAAATSVAMNRMSSLSSF-----TVNRPFLFFIMEDTIGVPLFVGSVRNPNPSA	440
A2AP_RNO	383	TMELSEAGVAAAAATSTAMTRMSLSSF-----FLNRPFLFFIMEETIGIPLFVGSVRNPNPSA	440
A2AP_GGA	435	TMELKEDGVESAATGVMISRSLSAFSDID-----RPFIFILFEEEMGIPLFVGSVKNPNPSA	491
A2AP_XTR	327	SLELNEEGVESAAVTAVITSRSHSVYRIN-----RPFLLFFLFDFTMGIPLFMCHVRNPNPGF	383
A2AP1_FRU	353	TMEINEEGAEAAAATTVVISRASSPVFHMT-----QPFFFAVMDDTTEVPFIMCVVMNPNPGA	410
A2AP2_FRU	352	TVELNEEGVKASAATVVTTLSISMFSVNS-----PFLFALVDDASLVPLFMGIVTNPAPDN	408
A2AP2_TNI	357	TVELSEKGVESAATVVTTMRSISMFSVNS-----PFLFAIVDDASLVPLFMGIVTNPAP--	411
A2AP2_DRE	401	SVEINEEGAEAAVAATSVVISRNSPFTVNN-----QPFFFALMDDLSTPLFLGVISNPNPGA	457
A2APL2_PMA	307	TMDLTEEGAEAAAATGVFLSRTPNPIYPVFK---VDRPFLFLIRDSTGTGLFLGRVMDPTDAA	366
A2APL1_PMA	355	TMTLNEEGVKATAATGIMISLMSVQHS--EELKVDKPFVFLIRDDDETGALLFVGRVTSPPPVP	415

---

```

A1AT_HSA      -  -----
A2AP_HSA     442  PRELK-----EQQDSPGNKDFLQSLKGFPRGDKLFGPDLKLVPPMEEDYPQFGSPK 491
A2AP_MMU     442  LPQLQ-----EQRDSPDNRLIGQNDKADFHGKTFGPDLKLAPRMEEDYPQFSSPK 491
A2AP_RNO     442  QPQPQ-----EQQDSPDNRRLDQNDKADIPGGKTFAPDLKLVPRLEEDYPQFSSPK 491
A2AP_GGA     493  APQVK-----ELQDLPD-----ATDDNEYTMPK 514
A2AP_XTR     385  QKTGK-----DPKNFDKGFLPK----- 400
A2AP1_FRU    412  PVMQT-----GDKVGFPIDKSMTRFEGPPK----- 435
A2AP2_FRU    410  DRMSN-----DDPLVNGTMSDQPVSVDNKNMNLFTERAACSAPTGENMMMD---- 455
A2AP2_TNI     -  -----
A2AP2_DRE    459  STMIT-----NPGNADKTDD--KPFVHPK----- 480
A2APL2_PMA   368  VEEEGISSVHYWRFLKSTGDGQHHVQALCQSSDYRGFGLRIWDPWCNITTNVIDESPVANGI 429
A2APL1_PMA   417  EKKKK-----EKHGDS-----SSSEEHGEGKGKGGKGGKHHH----- 447

```

```

A1AT_HSA      -  -----
A2AP_HSA     -  -----
A2AP_MMU     -  -----
A2AP_RNO     -  -----
A2AP_GGA     -  -----
A2AP_XTR     -  -----
A2AP1_FRU    -  -----
A2AP2_FRU    -  -----
A2AP2_TNI    -  -----
A2AP2_DRE    -  -----
A2APL2_PMA   431  PNSAVVIMASPKRT 443
A2APL1_PMA   -  -----

```

**Appendix 8.3.29: Alignment of protein sequences of C1 inhibitor and fish specific group V4 (FSG4).** Gene specific features include two pairs of conserved disulfide bridges (blue boxes and marked as pair C1 and C2), inhibitory RCL (red box). Additionally, FSG4 possess two extra Ig domains (labeled as Ig1 and Ig2 domains) with invariant cysteine residues (pink boxes). Conserved intron positions are indicated above the alignment. Predicted introns in Ig domains of FSG4 are indicated by #.

A1AT_HSA	-	-----		#	-
C1IN_HSA	1	MASRLTLLTLLLLLAGDRASSNPNATSSSSQDPESLQDRGEGKVATTVISKMLFVEPIL			60
C1IN_MMU	1	MASRLTPLTLLLLLAGDRAFSDEATSHSTQDPLEAQAKSRESFPERDDSWSP-PEPTV			59
C1IN_RNO	1	MASKLTPLTLLLLLAGDRAFSDEVTSHSSQDPLVVQEGSRDVSUPERDGRSP-IEHTG			59
C1IN_GGA	1	M---KLCLLMVCLVAVMTATV-----TPALEPLGP-----VEFVA			32
FSG4_FRU	1	MRFQATLCFLLQLIFEHSL--CTH-LQVTSGSSLELPCLPA-VTEFITAPASWTFNGVNL			56
FSG4_TNI	1	MKLQATFCFLLQLFEHSL--CTR-LTVTTGVLSLELPCFPARLPELVTA PITWTFNGMNL			57
FSG4_DRE	1	MYR-----WLLLLCVGLSLSSCDITVLLYSSIS--LPCVDPNAPALAGSTYIWNF-T---			49
					Ig1 Dom-
				#	
A1AT_HSA	-	-----			-
C1IN_HSA	61	EVSSLPTTNSTTNSATKITANTTDEPTTQPTTE-----			93
C1IN_MMU	60	LPSTWPTTSVAITITNDTMGKVANESFSQHSQP-----			92
C1IN_RNO	60	QSSTWPTTSGSTKISNDTMDQVANESFIQHVQP-----			92
C1IN_GGA	33	AP-----EESGTAANGSVATPEPE-----			51
FSG4_FRU	57	SAAVSDSVRIKRDGLYLSISPITAAHQGYACLVK----YINMDIVRITYDIAVINDS----			108
FSG4_TNI	58	GAEASDSVRIKENGSYLSIFPVTVAHQGYVCLVN----QTNMNILRAYTITVINDS----			109
FSG4_DRE	50	-APQTEPHTLSEKGIKILTLKNVNSSYSYGQYKCVQEGYRDEARVRRSRTFSLQVEEPPLLQ			108
		ain			
					Ig2 Domain
A1AT_HSA	-	-----			-
C1IN_HSA	94	-----PTTQPTIQPTQP---TTQLP---TDSPTQ-----			116
C1IN_MMU	93	-----AAQLPTDSPGQPPLNSSSQPSTASDLPTQ-----			121
C1IN_RNO	93	-----AAQLPEDSPSQSPVNSSSPPSTASAPPTQ-----			121
C1IN_GGA	52	-----TTTGPPSTPPGT-----INPP-----			67
FSG4_FRU	109	FAYDINAAQGSTIYLPYFPHSSQ-ILANALWYKETGIG-----QKTLHFSKEPLDK			160
FSG4_TNI	110	LFYEINALQGSTIYLPYIVPQSNQ-IPAKALWYKETDAG-----QRS--DFREGASGE			159
FSG4_DRE	109	EWQVIRVEAGYDVILPKVSFSNETISPSVVWKQVGTGQAVLLNPDKNTDVEEKKEKDKE			168
					Ig2 Domain
A1AT_HSA	1	-----EDPQGD			6
C1IN_HSA	117	-----PTTGSFCPG			125
C1IN_MMU	122	-----ATTEPFCPE			130
C1IN_RNO	122	-----APTEPLCPE			130
C1IN_GGA	68	-----CPG			70
FSG4_FRU	161	LER-VQQIYPLDQDQSVKFTMVLMEDSGTYRCEESPEGEELSVVRVTVKVAPTPVPYSCGE			219
FSG4_TNI	160	MQR-LQR-YPLAHDQSVIITMVVSEDSGTYRCEYYPVGVKLSLILVTVK---APVPYSCGE			214
FSG4_DRE	169	PQRVFWDISPEEKDWAIKISQTRWKDAGMYQCVINTNQTLLELEMEG----PPPHCEG			224
					C1

---

A1AT_HSA	7	AAQKTDTS	HH	QD	HPTF	NK	I	T	P	N	L	A	E	F	A	F	S	L	Y	R	Q	L	A	H	Q	-	S	N	S	T	N	I	F	F	S	P	V	S	I	A	T	A	F	A	M	L	S	65													
C1IN_HSA	126	---	P	V	T	L	C	S	D	L	E	S	H	S	T	E	A	V	L	G	D	A	L	V	D	F	S	L	K	L	Y	H	A	F	S	A	M	K	K	V	E	T	N	M	A	F	S	P	F	S	I	A	S	L	L	T	Q	V	L	182	
C1IN_MMU	131	---	P	L	A	C	C	S	D	S	D	R	D	S	S	E	A	K	L	S	E	A	L	T	D	F	S	V	K	L	Y	H	A	F	S	A	T	K	M	A	K	T	N	M	A	F	S	P	F	S	I	A	S	L	L	T	Q	V	L	187	
C1IN_RNO	131	---	P	L	A	W	C	S	D	S	D	R	D	S	S	E	A	T	L	S	E	A	L	T	D	F	S	V	K	L	Y	H	A	F	S	A	T	K	K	A	E	T	N	M	A	F	S	P	F	S	I	A	S	L	L	T	Q	V	L	187	
C1IN_GGA	71	D	E	E	P	A	E	T	C	W	A	P	T	R	E	Q	K	E	-	V	A	M	A	L	G	T	F	A	L	R	F	Y	Q	H	M	A	E	S	A	K	P	D	T	N	L	L	F	S	P	V	N	V	A	L	G	L	S	H	L	129	
FSG4_FRU	220	W	T	A	A	W	N	P	C	H	D	Q	E	D	H	T	G	E	A	V	L	Q	E	S	M	A	E	F	S	M	K	L	Y	S	F	V	R	E	S	-	Q	L	S	N	N	L	L	S	P	L	S	I	S	A	L	L	S	H	L	278	
FSG4_TNI	215	F	A	T	E	W	K	P	C	L	D	Q	S	R	S	D	E	A	V	L	Q	E	S	L	A	E	F	S	M	K	L	Y	S	F	L	R	E	S	-	L	P	S	D	M	I	L	V	S	P	L	S	I	S	T	L	L	S	H	L	273	
FSG4_DRE	225	Y	T	D	P	W	E	S	C	N	D	P	D	S	R	S	S	K	A	I	L	Q	E	S	L	G	D	F	S	T	S	V	Y	S	R	L	K	G	S	-	K	A	K	A	M	L	I	F	S	P	I	S	I	A	A	A	L	S	N	L	283

C1

67a

123a

A1AT_HSA	66	L	G	T	K	A	D	I	H	D	E	I	L	E	G	L	N	F	N	L	T	E	I	P	E	A	Q	I	H	E	G	F	Q	E	L	L	R	T	L	N	Q	P	D	S	Q	L	T	T	G	N	G	L	F	L	S	E	G	L	125
C1IN_HSA	183	L	G	A	G	E	N	T	K	T	N	L	E	S	I	L	S	Y	---	P	K	D	F	T	C	V	H	Q	A	L	K	G	F	T	---	---	---	T	K	G	V	T	S	V	S	Q	I	F	H	S	P	D	L	A	231				
C1IN_MMU	188	L	G	A	G	D	S	T	K	S	N	L	E	S	I	L	S	Y	---	P	K	D	F	A	C	V	H	Q	A	L	K	G	F	S	---	---	---	S	K	G	V	T	S	V	S	Q	I	F	H	S	P	D	L	A	236				
C1IN_RNO	188	L	G	A	G	D	S	T	K	S	N	L	E	D	I	L	S	Y	---	P	K	D	F	A	C	V	H	Q	T	L	K	A	F	S	---	---	---	S	K	G	V	T	S	V	S	Q	I	F	H	S	P	D	L	A	236				
C1IN_GGA	130	L	G	A	R	G	E	T	Q	Q	R	L	A	A	I	L	G	Y	---	Q	P	G	L	A	C	V	H	S	A	L	Q	L	V	N	---	---	---	V	S	G	L	L	S	A	T	V	I	F	H	H	P	D	L	H	179				
FSG4_FRU	279	L	G	A	R	D	I	T	Q	R	A	I	E	R	A	L	A	V	---	P	H	D	F	S	C	V	H	F	Q	M	A	K	L	R	E	K	L	---	---	---	A	S	S	L	Q	M	S	S	Q	I	Y	H	P	K	M	I	330		
FSG4_TNI	274	L	G	A	R	G	K	T	Q	R	A	I	E	S	A	L	S	V	---	P	H	E	F	S	C	I	H	F	H	I	K	K	L	E	K	L	---	---	---	D	T	S	L	Q	M	A	S	Q	I	Y	H	P	D	I	N	325			
FSG4_DRE	284	L	G	A	R	G	K	T	R	M	H	L	E	G	A	L	G	L	---	P	L	G	F	S	C	L	H	T	E	L	K	K	L	R	G	V	M	---	---	---	K	D	T	L	K	M	A	S	A	I	F	Y	N	P	E	Q	335		

C2

A1AT_HSA	126	L	V	D	R	F	L	E	D	V	K	K	L	Y	H	S	E	A	F	T	V	M	F	G	D	T	E	E	A	K	K	Q	I	N	D	Y	V	E	K	G	T	Q	G	K	I	V	D	L	V	K	E	L	D	R	D	T	V	F	A	L	V	185
C1IN_HSA	232	I	R	D	T	F	V	N	A	S	R	T	L	Y	S	S	S	P	R	V	L	S	N	-	N	S	D	A	N	L	E	L	I	N	T	W	V	A	K	N	T	N	N	K	I	S	R	L	L	D	S	L	P	S	D	T	R	L	V	L	290	
C1IN_MMU	237	I	R	D	T	Y	V	N	A	S	Q	S	L	Y	G	S	S	P	R	V	L	G	P	-	D	S	A	A	N	L	E	L	I	N	T	W	V	A	E	N	T	N	H	K	I	R	K	L	L	D	S	L	P	S	D	T	R	L	V	L	295	
C1IN_RNO	237	I	R	D	T	Y	V	N	A	S	L	S	L	Y	G	S	S	P	R	V	L	G	P	-	D	G	D	A	N	L	K	L	I	N	T	W	V	A	E	N	T	N	H	K	I	N	E	L	L	D	S	L	P	S	D	T	R	L	V	L	295	
C1IN_GGA	180	L	R	P	R	F	L	N	E	S	W	R	F	Y	K	A	R	P	R	E	L	S	G	-	N	G	S	L	D	L	Q	R	I	N	E	W	R	K	A	T	H	G	L	V	P	Q	L	L	S	Q	L	P	D	E	P	R	L	V	L	238		
FSG4_FRU	331	L	S	E	S	F	T	N	S	I	Q	F	Y	E	S	V	P	T	R	L	L	E	-	T	S	E	E	N	T	N	M	I	S	W	V	A	N	K	T	N	N	K	I	Q	H	L	V	D	S	V	S	P	S	T	Q	L	M	F	L	389		
FSG4_TNI	326	L	S	E	S	F	T	N	H	S	I	Q	F	Y	E	A	V	P	T	R	L	L	E	-	T	S	E	E	N	T	N	M	I	S	W	V	A	N	K	T	N	N	K	I	Q	R	L	V	D	S	V	S	S	T	Q	L	M	L	384			
FSG4_DRE	336	L	A	E	A	F	I	N	Q	S	K	E	F	Y	E	F	V	P	Q	L	T	N	-	D	S	T	R	N	V	A	L	I	N	K	W	V	E	N	K	T	N	K	I	T	Q	L	I	D	D	V	D	P	S	T	T	F	V	L	394			

192a

238c

A1AT_HSA	186	N	Y	I	F	F	K	G	K	D	E	R	P	F	E	V	K	D	T	E	E	E	D	F	H	V	D	Q	V	T	T	V	K	V	P	M	M	K	R	L	G	-	M	F	N	I	Q	H	C	K	L	S	S	W	L	L	M	K	Y	244	
C1IN_HSA	291	N	A	I	Y	L	S	A	K	K	T	T	F	D	P	K	K	T	R	M	E	P	F	H	F	K	-	N	S	V	I	K	V	P	M	M	S	S	K	K	P	V	A	H	F	I	D	Q	T	L	K	A	K	V	G	Q	L	Q	L	349	
C1IN_MMU	296	N	A	V	Y	L	S	A	K	K	I	T	F	E	P	K	K	-	M	M	A	P	F	F	Y	K	-	N	S	M	I	K	V	P	M	L	S	S	K	K	Y	P	L	A	L	F	N	D	Q	T	L	K	A	K	V	G	Q	L	Q	L	353
C1IN_RNO	296	N	A	V	Y	L	S	A	K	K	I	T	F	E	Q	K	K	-	M	M	A	S	F	L	Y	K	-	N	S	M	I	K	V	P	M	L	S	S	K	K	Y	P	L	A	L	F	N	D	Q	T	L	K	A	K	V	G	Q	L	Q	L	353
C1IN_GGA	239	S	A	V	H	F	Q	A	R	W	Q	K	P	F	K	T	K	H	T	V	L	L	P	F	M	R	H	G	H	R	P	V	D	V	L	T	M	T	S	K	K	P	V	A	S	F	T	D	P	R	L	Q	V	Q	V	R	L	E	L	298	
FSG4_FRU	390	N	A	V	S	F	K	G	Q	E	L	K	F	D	L	N	P	-	R	D	A	L	F	S	K	L	N	G	D	L	V	S	V	P	V	F	H	H	P	D	Y	L	L	A	T	M	I	D	N	A	L	K	A	Q	V	G	R	F	A	L	448
FSG4_TNI	385	N	A	V	S	F	K	G	Q	E	L	K	F	D	S	K	P	-	K	K	R	H	F	T	K	P	N	G	D	L	V	S	V	W	V	L	Y	H	Q	S	P	L	A	M	T	L	D	T	D	L	K	A	M	V	A	R	F	A	L	443	
FSG4_DRE	395	N	A	V	Y	N	G	K	K	T	V	E	S	T	N	-	N	K	E	R	F	T	M	F	S	G	E	T	K	D	V	K	T	L	Y	S	S	N	Y	I	L	Q	M	G	Y	N	K	Q	L	K	A	D	V	G	K	F	P	L	453		

A1AT_HSA	245	LGNATAIFFL	PDE---	GK	LQ	HLENEL	THD	I	IT	KF	LE	NEDR---	RS	AS	LHL	PK	LS	IT	G	TYD	298																																	
C1IN_HSA	350	SHNLSLVIL	V	P	Q	N	-	L	K	H	R	L	E	D	M	E	Q	A	L	S	P	S	V	F	K	A	I	M	E	K	L	E	M	S	K	F	Q	P	T	L	L	T	L	P	R	I	K	V	T	T	S	Q	408	
C1IN_MMU	354	SHNLSFVIV	V	P	V	F	-	P	K	H	Q	L	K	D	V	E	K	A	L	N	P	T	V	F	K	A	I	M	K	K	L	E	L	S	K	F	L	P	T	Y	L	T	M	P	H	I	K	V	K	S	S	Q	412	
C1IN_RNO	354	SHNLSFVIM	V	P	Q	S	-	P	T	H	Q	L	E	D	M	E	K	A	L	N	P	T	V	F	K	A	I	L	K	K	L	E	L	S	K	F	O	P	T	Y	V	M	M	P	R	I	K	V	K	S	S	Q	412	
C1IN_GGA	299	SGGLSLVVL	V	P	R	G	-	P	P	E	A	L	E	A	V	E	R	A	L	D	P	P	T	F	L	A	L	L	Q	R	A	A	N	T	P	V	R	P	T	T	V	A	L	P	R	L	H	L	D	L	A	V	D	357
FSG4_FRU	449	TGDSSLYIL	L	P	N	--	T	V	D	L	Q	L	V	E	S	M	T	Y	D	T	L	R	Q	L	M	D	K	M	K	T	V	V	P	Q	K	T	E	V	T	L	P	K	I	K	L	D	V	E	P	D	506			
FSG4_TNI	444	TGETSLYVL	L	P	A	S	H	T	M	A	D	L	Q	V	E	R	M	T	D	T	A	L	L	R	M	I	H	N	M	K	T	I	V	P	Q	K	A	E	V	I	L	P	R	I	K	L	D	V	K	P	D	503		
FSG4_DRE	454	TGQNSLYIL	V	P	R	T	L	S	E	S	F	L	L	M	E	N	N	I	N	R	N	T	L	E	E	M	V	S	E	M	N	Q	T	P	A	Q	S	A	E	V	T	L	P	A	I	K	L	T	M	T	T	Q	513	

307a

A1AT_HSA	299	L	K	S	V	L	G	Q	L	G	I	T	K	V	F	S	N	G	A	D	L	S	G	V	T	E	E	A	P	L	K	L	S	K	-	A	V	H	K	A	V	L	T	I	D	E	K	G	T	E	A	A	G	A	M	F	L	E	A	I	P	357	
C1IN_HSA	409	M	L	S	I	M	E	K	L	E	F	F	D	-	F	S	Y	D	L	N	L	C	G	L	T	E	D	P	D	L	Q	V	S	A	-	M	Q	H	O	T	V	L	E	L	T	T	E	T	G	V	E	A	A	A	A	S	A	I	S	V	-	A	465
C1IN_MMU	413	M	L	S	V	M	E	K	L	E	F	F	D	-	F	T	Y	D	L	N	L	C	G	L	T	E	D	P	D	L	Q	V	S	A	-	M	K	H	E	T	V	L	E	L	T	E	S	G	V	E	A	A	A	A	S	A	I	S	F	-	G	469	
C1IN_RNO	413	M	L	S	I	M	E	K	L	E	F	F	D	-	F	T	Y	D	L	N	L	C	G	L	T	E	D	P	D	L	Q	V	S	S	-	M	K	H	E	T	V	L	E	L	T	E	T	G	V	E	A	A	A	A	S	T	I	S	V	-	A	469	
C1IN_GGA	358	V	V	A	K	V	H	D	M	D	F	G	-	L	F	L	-	D	A	E	L	C	G	L	A	O	G	P	E	V	A	V	-	D	A	A	Q	H	R	A	V	L	T	L	D	E	K	G	V	E	A	A	G	A	M	A	T	S	L	-	A	413	
FSG4_FRU	507	M	F	M	L	M	K	K	L	G	L	S	S	L	F	E	-	D	A	N	L	C	G	L	Y	S	E	D	R	L	V	L	-	D	D	V	R	H	R	G	L	L	A	L	T	E	H	G	V	E	A	V	A	V	T	S	T	T	F	-	S	563	
FSG4_TNI	504	M	F	M	L	M	K	K	L	G	L	S	S	L	F	E	-	D	A	N	L	C	G	L	Y	S	E	D	R	L	V	L	-	D	E	V	R	H	R	G	F	L	A	L	T	E	Q	G	V	E	A	V	A	V	T	S	V	S	F	-	S	560	
FSG4_DRE	514	V	D	D	L	R	N	M	G	L	S	D	L	F	M	N	-	P	N	L	C	G	M	F	P	G	E	P	E	S	F	I	S	D	V	R	H	R	A	F	L	S	L	T	E	K	G	V	E	A	A	A	T	S	I	S	F	-	S	571			

C2

A1AT_HSA	358	M	S	I	P	P	E	V	K	F	N	K	P	F	V	L	M	I	E	Q	N	T	K	S	P	L	F	M	G	K	V	V	M	P	T	Q	K	394	
C1IN_HSA	466	R	T	L	L	V	-	F	E	V	Q	P	P	L	F	V	L	W	D	Q	Q	H	K	F	P	V	F	M	G	R	V	Y	D	P	R	A	-	500	
C1IN_MMU	470	R	S	L	P	I	-	F	E	V	Q	R	P	P	L	F	L	W	D	Q	Q	H	R	F	P	V	F	M	G	R	V	Y	D	P	R	G	-	504	
C1IN_RNO	470	R	N	L	L	I	-	F	E	V	Q	P	P	L	F	L	W	D	Q	Q	H	K	F	P	V	F	M	G	R	V	Y	D	P	R	A	-	504		
C1IN_GGA	414	R	I	A	L	C	-	L	E	A	L	Q	P	P	L	F	V	L	W	D	E	G	N	A	I	P	L	F	M	G	R	L	S	D	P	Q	A	-	448
FSG4_FRU	564	R	T	Y	N	S	-	F	S	A	L	H	P	P	I	F	L	L	W	S	D	R	A	N	V	P	L	F	V	G	R	V	E	P	-	-	-	596	
FSG4_TNI	561	R	T	Y	I	S	-	F	S	A	L	Q	P	P	V	F	L	L	W	S	D	Q	A	N	V	P	L	F	V	G	R	V	I	D	P	-	-	593	
FSG4_DRE	572	R	S	F	S	S	-	F	S	A	L	Q	P	P	V	L	I	L	W	S	D	E	A	A	V	P	L	F	M	G	R	I	N	P	-	-	-	604	

**Appendix 8.3.30: Alignment of ATIII sequences from vertebrates.** The ATIII gene specific features are conserved helix-D (yellow box), eight residues involved in heparin binding (orange box), six cysteine residues (blue box) and N-glycosylation sites (cyan boxes).

A1AT_HSA	-	-----	-----	-----	-----	-----	-----	-----	-----	-	
ATIII_HSA	1	MYSNVIGTVTSGKRKVVLLSLLLIGFWDCVTC	HGSPV-DIC	TAKPRD	IPMNP	CIYRSP-				58	
ATIII_MMU	1	MYSPGAGSGAAGERKLCLLSLLLIGALGCAIC	HGNPVDDIC	IAKPRD	IPVNPL	CIYRSP-				59	
ATIII_RNO	1	MYSPGIGSAVAGERKLCLLSLLLIGALGCAVCH	HGNPVDDIC	IAKPRD	IPVNPM	CIYRSP-				59	
ATIII_GGA	1	MHLFMVCLFGLWGMASPAYAVE	-----	DIC	TAKPRD	IPVNP	CIYRNP-			44	
ATIII_XTR	1	MYLLSLLLLSLLGSAYLQ	-----	PQH-ADIC	LAKPKD	IPLTPM	CVYRKPL			44	
ATIII_FRU	1	MPASDWLLLLASLHVVSAD	-----	VLDIC	GAKPRD	DLALEPR	CIYRSPD			43	
ATIII_TNI	1	MFYLLRTQASYWLLLLAALPAVAD	-----	VLDIC	SAKPRD	DLALEPR	CIYRSPD			48	
ATIII_DRE	1	MKLLACMWALWAFALCSI	-----	HATKDIC	NAKPRD	LDLPLEPM	CIYRNPD			44	
				1C	2C						
A1AT_HSA	1	-----	EDPQGDAAQKTD	TSHHDQD	HPTFNKI	TPNLAEF	AFSLYRQLAHQ-SNSTN	IFF		52	
ATIII_HSA	59	----	EKKATEDEGSEQ----	KIPEATN	RE	VWELSKANSR	FATTFYQHLADSKNDNDN	IFL		110	
ATIII_MMU	60	----	GKKATEEDGSEQ----	KVPEATN	RE	VWELSKANSR	FATNFYQHLADSKNDNDN	IFL		111	
ATIII_RNO	60	----	AKKATEEDVLEQ----	KVPEATN	RE	VWELSKANSR	FATNFYQHLADSKNDNDN	IFL		111	
ATIII_GGA	45	----	EKKPQESKVLDPGKGRIPDF	TNP	RE	VWELSRANSR	FAVVFYKHL PNSKDKEE	IFL		99	
ATIII_XTR	45	----	EVVETEKEKEPTTQEQ----	KVPESTN	RE	VYELSQANAR	FAIAFYKNLADSKRDKE	IFM		101	
ATIII_FRU	44	----	PEAPEPLTTHP----	VPGSTN	RE	VWELSKANAR	FAMSLYKQVASSRGPES	IFM		93	
ATIII_TNI	49	----	PDGAQPPTTPP----	VPESTN	RE	VWELSKANGR	FALALFKQVASS-RPED	NVFM		97	
ATIII_DRE	45	----	EIQPNKEPEN----	IPVGTN	RE	VWELSKANSR	FALSFLKQLAEGKSND	ENIFL		93	
				78c							
A1AT_HSA	53	SPVSI	IATAFAMLSL	GTKAD	THDEILEG	LNFN-LTEI	PEAQIHEGFQELLR	TLN-QPDSQL		110	
ATIII_HSA	111	SPLS	ISTAFAMTKL	GACND	TLQQLMEVF	KFDTISEK	TSDQIHFFF	AKLNCRL	LYRKANKSS	170	
ATIII_MMU	112	SPLS	ISTAFAMTKL	GACND	TLKQLMEVF	KFDTISEK	TSDQIHFFF	AKLNCRL	LYRKANKSS	171	
ATIII_RNO	112	SPLS	ISTAFAMTKL	GACNNT	TLKQLMEVF	KFDTISEK	TSDQIHFFF	AKLNCRL	LYRKANKSS	171	
ATIII_GGA	100	SPLS	ISTAFAMTKL	GACGD	TLQQLMEVF	QFDTISEK	TSDQVHFFF	AKLNCRL	LYKKANKSS	159	
ATIII_XTR	102	SPLS	ISQAF	TMAKL	GACNNT	TLKQLMEVF	HFDTVSERASDQI	HYFFF	AKLNCRL	FRKANKSS	161
ATIII_FRU	94	SPIS	ISTAFAMTKL	GACN	QTLQ	LMRVF	QFDTIKEK	TSDQVHFFF	AKLNCRL	YRKKDKSN	153
ATIII_TNI	98	SPIS	ISTAFAMTKL	GACN	QTLQ	LMKVF	QFDTIKEK	TSDQVHFFF	AKLNCRL	YRKKDATT	157
ATIII_DRE	94	SPIS	ISTAFAMTKL	GACN	TTLEQ	LMKVF	QFDTIKEK	TSDQVHFFF	AKLNCRL	YRKKHETT	153
				2C				1C			
A1AT_HSA	111	QLTTG	NGLFLSEGL	KLVDK	FLEDV	VKLYHSEAF	TVNFGD	TEE-AKKQ	INDYVEKGTQ	GKI	169
ATIII_HSA	171	KLVS	ANRLFGDKSL	TFNETY	QDISEL	VYGAKL	QPLDF	KENAEQSRAA	INKWVSNK	TEGRI	230
ATIII_MMU	172	DLVS	ANRLFGDKSL	TFNESY	QDVSEV	VYGAKL	QPLDF	KENPEQSRVT	INNWWANK	TEGRI	231
ATIII_RNO	172	NLVS	ANRLFGDKSL	TFNESY	QDVSEI	VYGAKL	QPLDF	KENPEQSRVT	INNWWANK	TEGRI	231
ATIII_GGA	160	ELIS	ANRLFGEKSL	VFNETY	QNI	SEIVYGAKL	WPLM	FKEKPELSR	KIIMEWVANK	TERRI	219
ATIII_XTR	162	ELVSV	NRLFGEKSL	TFNETY	QDISEI	VYGAKL	WPLM	FRDKPELSREI	INNWVSNK	TEKRI	221
ATIII_FRU	154	ELVS	ANRLFGDKSL	AFDQTY	QNI	SETVYGAKL	LP	DFKDDPEKARVT	INNWISNK	TENLI	213
ATIII_TNI	158	ELVS	ANRLFGDQSD	LQSQTY	QNI	SETVYGAKL	LP	DFKDFRRARLT	IMSWISNK	TRNLI	217
ATIII_DRE	154	ELIS	ANRLFGDKST	TFNETY	QHI	SETVYGAKL	MPLD	FKEKPEASRIT	IMEWIANK	TENRI	213
								148c			



			191c				
A1AT_HSA	170	VDLVKELDRD----	TVFALVNYIFFKGRD	ERPEVVKDTEEDF	FHVDQVTTVKVPMMKRLG	225	
ATIII_HSA	231	TDVIPSEAINEL--	TVLVLVNTIYFKGL	DKSKFSPENTR	KELFYKADGESCSAS	MMYQEG 288	
ATIII_MMU	232	KDVIPOGAINEL--	TALVLVNTIYFKGL	DKSKFSPENTR	KEPFFYKVDGQSC	PVPMMYQEG 289	
ATIII_RNO	232	KDVIPOGAIDEL--	TALVLVNTIYFKGL	DKSKFSPENTR	KEPFFHKVDGQSC	LVPMMYQEG 289	
ATIII_GGA	220	TEVIPEKGIDDL--	TVLVLVNTIYFKGH	DKSKFPAPNTRL	DLFHKANGETCN	VPIIMYQES 277	
ATIII_XTR	222	TDVIPKDAITPD--	TVLVLINAIYFKGL	DKSKFENSEM	TKMDQFHPAKNSN	CLTATMYQEG 279	
ATIII_FRU	214	QDTLP-PGVLD	SN-TVLVLVNTIYFKGH	DKSKFDKDNVYV	SEFHSSQTRS	CSVMMYQER 271	
ATIII_TNI	218	QDTLPDTGVLGLQD	TVLVLVNTIYFK--	SERNKFDKDNVYV	SDFHVS	PARTCSARMMYQEA 276	
ATIII_DRE	214	KDTLP-EGSIDTN-	TILVLVNAIYFKGQ	DKNKF	DKQNVMKLDFHVS	SPTHKCPVPMYQEK 271	
					3C		
				262c!			
A1AT_HSA	226	MFNIQHCKKLSSWVLLMKYLG-NATAIFL	LPDEGK-LQHLE	NELTHDII	TKFLE	NEDRRS 283	
ATIII_HSA	289	KFRYRRVAEG-TQVLEL	PFKGGDITMVLIL	PKPEKSLAK	VEKELTPEVL	QEWLDELEEMM 347	
ATIII_MMU	290	KFKYRRVAEG-TQVLEL	PFKGGDITMVLIL	PKPEKSLAK	VEQELTPELL	QEWLDELSETM 348	
ATIII_RNO	290	KFKYRRVGEGETQVLE	MPFKGGDITMVLIL	PKPEKSLAK	VEQELTPELL	QEWLDELSEVM 348	
ATIII_GGA	278	RFRYAFIQEDKVQVLEL	PFYKGGDITMVLV	LPKAGTPLVE	VERDLTSDKLQD	WIDSMMEVS 337	
ATIII_XTR	280	TFRYGSFKDDGVQVLEL	PFYKGGDITMVLV	LPSSQETPL	TTVEQNL	TLEKLGWNLQKSRELO 339	
ATIII_FRU	272	RFRYKHFPEQVQVLE	MPYRGDDITMVIIL	PSQGTALS	VEEVLDL	LKKLSAWLDQMKETT 331	
ATIII_TNI	277	RFRYRHVPEDHVQVLE	MPYRGDITMVIIL	PSRG	TALRVEEVLDL	LKKVSAWLDQMKETM 336	
ATIII_DRE	272	KFYAKIPEDKVKILEL	PFYNGG	ITMVLILP	IEGATLS	VVANMNLKKLVGWLHAMKETT 331	
				320a	339c		
A1AT_HSA	284	ASLHLPKLSITGTYDLKSVL	QGLG	ITKVF	NSGDL	SLGVTEEA---PLKLSKAVHKAVLT 339	
ATIII_HSA	348	LVVHMPFRFRIEDGFS	LKEQLQDMGLVDL	FSPEKSL	PGIVAEG-RDD	LYVSDAFHKAFLE 406	
ATIII_MMU	349	LVVHMPFRFRTEDGFS	LKEQLQDMGLIDL	FSPEKSL	QPGIVAAG-RDD	LYVSDAFHKAFLE 407	
ATIII_RNO	349	LVVHMPFRFRIEDSFS	LKEQLQDMGLVDL	FSPEKSL	PGIIAEG-RDD	LYVSDAFHKAFLE 407	
ATIII_GGA	338	LTVSFPFRFVEKGF	SVKEKLRKMGLEDL	FSPEKSL	PGIVAAG-RTD	LYVSDAFHKAFLE 396	
ATIII_XTR	340	LSVYLPFRFRVEDS	FSVKEKLEMG	LVDFDPNSAK	LPGIIAAG-RTD	LYVSDAFHKAFLE 398	
ATIII_FRU	332	VSVHVPFRFRVEDS	FSLKEKQLQLGL	TDLFDPNKAS	LPGMLEDG-VEGL	HISDAYHKAFLE 390	
ATIII_TNI	337	VSVHVPFRFRMEDS	FRLKEKLOVGL	SLDFSPDRAS	LPGMLEDG-GEGL	HISDAYHKAFLE 395	
ATIII_DRE	332	VAVQIPFRFRVEDS	FSLKEQLTKMGLEDL	FSPANAS	LPGMVADAE	GNLFI	SDAYHKAFLE 391
A1AT_HSA	340	IDEKGTAAAGAMFLEAIPMSIPPE---	VKFNKPFVFLMIE	QNTKSPLFMG	KVVNPTQK-	394	
ATIII_HSA	407	VNEEGSEAAASTAVVI	AGRSLNPNRV	TFKANRPFLV	FIREVPLNTIIF	MGRVANPCVK- 464	
ATIII_MMU	408	VNEEGSEAAASTSVVITGRSLNPNRV	TFKANRPFLV	LIREVALNTIIF	MGRVANPCVN-	465	
ATIII_RNO	408	VNEEGSEAAASTSVVITGRSLNPNRV	TFKANRPFLV	LIREVALNTIIF	MGRVSNPCVN-	465	
ATIII_GGA	397	VNEEGSEASAATAVVISGRSFP	MRIIFEANRP	FLFIREATLNTIIF	MGRISDPCS--	453	
ATIII_XTR	399	VNEEGSEAAASTAVILTGRSLN	LNRIIFRANRP	FLVIREVA	INAILFMGRVANPC	TE- 456	
ATIII_FRU	391	VNEEGSEAAAATAAVATGRS	INLNREIFQANRP	FLLLIREAS	INTLLFIARVAEPC	DR- 448	
ATIII_TNI	396	VNEEGSEAGAATAVAVGRS	IHFSREVFQANRP	FLLLIREAS	INTLLFVARVAQPC	SP- 453	
ATIII_DRE	392	VNEEGSEASAATAVAVATGRS	LNIFREQVADRP	FLFIRESS	INALIFTGRVANPC	RSS 450	
					3C		



**Appendix 8.3.31: Alignment of HSP47 homologs from vertebrates.** Group V6 gene specific features are a non-inhibitory RCL (red boxes) and an endoplasmic reticulum (ER) retention signal (light blue boxes). *Fugu* gene HSP47\_FRU1 has two additional unique introns at positions 36b and 102c. In HSP47\_TNI, the intron at position 192a is tentatively assigned (indicated by ?) and an additional low complexity region at the intron at position 300c (predicted by GENSCAN and FGENESH) was deleted manually.

A1AT_HSA	1	-----EDPQGDAAQKTD TSHHDQDHP TFNKI TPNLAE	34
HSP47_HSA	1	MRSLLL GTL CLLAVALAA-----EVKKPVEAAAPGTAEKLSSKATTLAEPSTGLA	50
HSP47_MMU	1	MRSLLL GTL CLLAVALAA-----EVKKPLEAAAPGTAEKLSSKATTLAERSTGLA	50
HSP47_RNO	1	MRSLLL GTL CLLAVALAA-----EVKKPVEAAAPGTAEKLSSKATTLAERSTGLA	50
HSP47_GGA	1	MQIFLVLALCGLAAAV-----PSEDRKLSDKATTLADRSTTLA	38
HSP47_XTR	1	MWMIKLLALSILLVVDAAVDKPKVVADKPKVAPVVE--PPVEKKISQHANVLADKSAGLA	58
HSP47_1_FRU	1	MRAANTVVL SLLALLASAE D-----KKL SNHATTLADNSANLA	38
HSP47_2_FRU	1	MLPRLPVYILLFLPLAPVQRSTADSSSEKSSASS-----PHLPPPPLG---DPSWALG	49
HSP47_TNI	1	MERNLD-----VIGERGTE--PP-----	16
HSP47_1_DRE	1	MWVSSLIALCLLAVAVSGE-----DKKLS THATSMADTSANLA	38
HSP47_2_DRE	1	MLASNVLLLCLLATVSAN-----KTLSSSIATTLADNSATLA	36
HSP47_3_DRE	1	MQPIFPVPLFLLLAQQSVWSS-----TPQEPKVQGGSPPEISSLHHP TWSLG	47
HSP47_PMA	1	MLLLEALASGALAAAAAD--GKKATVSKADAAAANNA TAPPKNLSEHAKKVGEGNWA	58

## 36b!

A1AT_HSA	35	FSLYRQLAHQS-NSTRIFFSVVS IATAFAMLSLGTAKADTHDEILEGLNFNL TEIPEAQIH	93
HSP47_HSA	51	FSLYQAMAKDQ-AVENILLSPVVVASSLGLVSLGGKATIASQAKAVL--SAEQLRDEEVH	107
HSP47_MMU	51	FSLYQAMAKDQ-AVENILLSPLVVASSLGLVSLGGKATIASQAKAVL--SAEKL RDEEVH	107
HSP47_RNO	51	FSLYQAMAKDQ-AVENILLSPLVVASSLGLVSLGGKATIASQAKAVL--SAEKL RDEEVH	107
HSP47_GGA	39	FNLYHAMAKDK-NMENILLSPVVVASSLGLVSLGGKATIASQAKAVL--SADKLND DYVH	95
HSP47_XTR	59	FNLYQTMAKDK-NVENILLSPVVVASSLGLVSLGGQASIAAQAKAVL--SADKLSDEH IH	115
HSP47_1_FRU	39	FSLYHNMAKDK-NVENILLSPVV LASSLGMVALGGKASIASQVKTVL--SADKLKDEHLH	95
HSP47_2_FRU	50	LRLYQALRSDS-RSVNTLFSPLLAASSL GALGGGSAGASASQFQDLL--KASSS-AKAGA	105
HSP47_TNI	17	-----GAME-GESVTEG-----RQDST-GRWNA	37
HSP47_1_DRE	39	FNLYHNVAKEK-GLENILLSPVVVASSLGMVAMGSKSSIASQVKSVL--KADALKDEHLH	95
HSP47_2_DRE	37	FNLYQNMAKDK-DIENILLSPVVVASSLGLVALGGKSNIASQVKTVL--SAASVKDEQLH	93
HSP47_3_DRE	48	LQLYRSLRTNG-SQTNTFIPPLLLANSL LALGGGAKGSIVSQFHDLL--RITKN-ENVVG	103
HSP47_PMA	59	IDLYQSVAKAVPAMERVVLSVVLVASALGAAQLGASSG IASRLKAT--NPSGLPGE GFH	116

## 102c!

A1AT_HSA	94	EGFQELLRTL NQPDS-----QLQLTTGNGLFLSEGLKLVDKLEL DVKKLYHSEAF	143
HSP47_HSA	108	AGLGELLRSLSNSTA-----RNVTWKLGSRLYGPSSVSFADDFVRS SKQHYNCEHS	158
HSP47_MMU	108	TGLGELLRSLSNSTA-----RNVTWKLGSRLYGPSSVSFADDFVRS SKQHYNCEHS	158
HSP47_RNO	108	TGLGELVRSLSNSTA-----RNVTWKLGSRLYGPSSVSFADDFVRS SKQHYNCEHS	158
HSP47_GGA	96	SGLSELLNEVSNSTA-----RNVTWKIGNRLYGPASINFADDFVKNSKKHYNYEHS	146
HSP47_XTR	116	SGLAELLNEVSNSTA-----RNVTWKIGNRLYGPSSISF TDDEVKNSKKHYNYEHS	166
HSP47_1_FRU	96	AGLSELLT LSDADK-----RNTTWKINNRLYGPSSVSFSDDFVKS SKKHYYDHS	146
HSP47_2_FRU	106	ELLSESLKSLGKSNG-----TSFHAAHASTALFSKEAPQV SQAVKDSQARFGLQH Q	156
HSP47_TNI	38	PPFGVAKDSIK-----CDMF SRELAKPERVLLNKTYP----ERQ	72
HSP47_1_DRE	96	TGLSELLTEVSDPQT-----RNVTWKISNRLYGPSSVSFAEDFVKNSKKHYNYEHS	146
HSP47_2_DRE	94	SGLSELLTEVSNPKA-----RNVTWKISNRFYGPSSVSFVDDFLKSKKHYYNDHS	144
HSP47_3_DRE	104	ETLTTAQKAVHESNG-----TSYILRSSALFSKQAPELEKSEKLQTHFGMQHV	154
HSP47_PMA	117	SGLAEVLDGLASQEEEEAAAAAATWRNHTWKAASRVYAPSGVTF SQGFVSSSKARYGLQHD	176

		192a!!	
A1AT_HSA	144	TVNFGDTEEAKKQINDYVEKGTGGKIVDLVKELDRD <sup>I</sup> IV--FAIVNYIFFKGFWERPF <sup>I</sup> EVK	201
HSP47_HSA	159	KINFPDKRSALQSINEWAAQT <sup>I</sup> TDGKLEVTKDVERTD <sup>I</sup> G--ALLV <sup>I</sup> RAMFFKPHWDEK <sup>I</sup> FHHK	216
HSP47_MMU	159	KINFRDKRSALQSINEWASQT <sup>I</sup> TDGKLEVTKDVERTD <sup>I</sup> G--ALLV <sup>I</sup> RAMFFKPHWDER <sup>I</sup> FHHR	216
HSP47_RNO	159	KINFRDKRSALQSINEWASQT <sup>I</sup> TDGKLEVTKDVERTD <sup>I</sup> G--ALLV <sup>I</sup> RAMFFKPHWDEK <sup>I</sup> FHHK	216
HSP47_GGA	147	KINFRDKRSALKSINEWAAQT <sup>I</sup> TDGKLEVTKDVEK <sup>I</sup> TDG--ALIV <sup>I</sup> RAMFFKPHWDEK <sup>I</sup> FHHK	204
HSP47_XTR	167	KINFRDKRSTLRSINEWASQA <sup>I</sup> TDGKLEVTSDMERT <sup>I</sup> DG--ALIV <sup>I</sup> RAMFFKPHWDER <sup>I</sup> FHHQ	224
HSP47_1_FRU	147	KINFRDKRSAVNSINEWAAKA <sup>I</sup> TDGKLEI <sup>I</sup> TKDVQ <sup>I</sup> NADG--AMIV <sup>I</sup> RAMFFKPHWDER <sup>I</sup> FHDK	204
HSP47_2_FRU	157	PLGKGD <sup>I</sup> SKAD-----L <sup>I</sup> KRLWEREF <sup>I</sup> GEG	178
HSP47_TNI	73	-----HARVP <sup>I</sup> LVKREFT <sup>I</sup> QG	87
HSP47_1_DRE	147	KINFRDKRSAIN <sup>I</sup> SINEWAAKT <sup>I</sup> TDGKLEI <sup>I</sup> TKDVK <sup>I</sup> NTD <sup>I</sup> G--AMIV <sup>I</sup> RAMFFKPHWDEK <sup>I</sup> FHHK	204
HSP47_2_DRE	145	KINFRDKRSAVKAI <sup>I</sup> NDWASK <sup>I</sup> STDGKLEVTKDVEK <sup>I</sup> TDG--AMI <sup>I</sup> RAMFFKPHWDEK <sup>I</sup> FHHK	202
HSP47_3_DRE	155	ALEDAQKQSDMEK <sup>I</sup> LQYWA <sup>I</sup> KS <sup>I</sup> GMD <sup>I</sup> CEETAAL <sup>I</sup> KTAL <sup>I</sup> E <sup>I</sup> FKAGAM <sup>I</sup> TANAL <sup>I</sup> HEK <sup>I</sup> GL <sup>I</sup> DRG <sup>I</sup> EYHE	214
HSP47_PMA	177	KVNLKDKRGAL <sup>I</sup> KAL <sup>I</sup> NEWAAQ <sup>I</sup> NT <sup>I</sup> GCKVKEVAKEL <sup>I</sup> DGADG--AVFV <sup>I</sup> RAL <sup>I</sup> FFKGRWNEK <sup>I</sup> FHHQ	234

		225a	
A1AT_HSA	202	D <sup>I</sup> TEED <sup>I</sup> EHVD <sup>I</sup> QVT <sup>I</sup> TVK <sup>I</sup> VM <sup>I</sup> KK <sup>I</sup> RL <sup>I</sup> GMFNI <sup>I</sup> QHCK <sup>I</sup> KL <sup>I</sup> SSVLL <sup>I</sup> MK <sup>I</sup> LGN-ATA <sup>I</sup> IF <sup>I</sup> LE <sup>I</sup> D <sup>I</sup> E <sup>I</sup> GK-	259
HSP47_HSA	217	MVDNRG <sup>I</sup> GMV <sup>I</sup> TRS <sup>I</sup> YTV <sup>I</sup> GV <sup>I</sup> VM <sup>I</sup> HRT <sup>I</sup> GLYNY <sup>I</sup> YD <sup>I</sup> DEKE <sup>I</sup> KL <sup>I</sup> QLVEM <sup>I</sup> PLA <sup>I</sup> HKL <sup>I</sup> SSL <sup>I</sup> I <sup>I</sup> IL <sup>I</sup> MP <sup>I</sup> H <sup>I</sup> VEP	276
HSP47_MMU	217	MVDNRG <sup>I</sup> GMV <sup>I</sup> TRS <sup>I</sup> YTV <sup>I</sup> GV <sup>I</sup> VM <sup>I</sup> HRT <sup>I</sup> GLYNY <sup>I</sup> YD <sup>I</sup> DEKE <sup>I</sup> KL <sup>I</sup> QVEM <sup>I</sup> PLA <sup>I</sup> HKL <sup>I</sup> SSL <sup>I</sup> I <sup>I</sup> IL <sup>I</sup> MP <sup>I</sup> H <sup>I</sup> VEP	276
HSP47_RNO	217	MVDNRG <sup>I</sup> GMV <sup>I</sup> TRS <sup>I</sup> YTV <sup>I</sup> GV <sup>I</sup> VM <sup>I</sup> HRT <sup>I</sup> GLYNY <sup>I</sup> YD <sup>I</sup> DEKE <sup>I</sup> KL <sup>I</sup> QLVEM <sup>I</sup> PLA <sup>I</sup> HKL <sup>I</sup> SSL <sup>I</sup> I <sup>I</sup> IL <sup>I</sup> MP <sup>I</sup> H <sup>I</sup> VEP	276
HSP47_GGA	205	MVDNRG <sup>I</sup> GMV <sup>I</sup> TRS <sup>I</sup> YTV <sup>I</sup> GV <sup>I</sup> VM <sup>I</sup> HRT <sup>I</sup> GLYNY <sup>I</sup> YD <sup>I</sup> EAE <sup>I</sup> KL <sup>I</sup> QVEM <sup>I</sup> PLA <sup>I</sup> HKL <sup>I</sup> SSM <sup>I</sup> IF <sup>I</sup> IM <sup>I</sup> PH <sup>I</sup> VEP	264
HSP47_XTR	225	MVDNRG <sup>I</sup> GMV <sup>I</sup> TR <sup>I</sup> SFT <sup>I</sup> VS <sup>I</sup> VM <sup>I</sup> HRT <sup>I</sup> GLYNY <sup>I</sup> LE <sup>I</sup> DEK <sup>I</sup> NGL <sup>I</sup> QILEM <sup>I</sup> PLA <sup>I</sup> HKL <sup>I</sup> SSM <sup>I</sup> L <sup>I</sup> F <sup>I</sup> IM <sup>I</sup> PH <sup>I</sup> VEP	284
HSP47_1_FRU	205	MVDTRG <sup>I</sup> ELV <sup>I</sup> TR <sup>I</sup> SHT <sup>I</sup> IGI <sup>I</sup> SM <sup>I</sup> HRT <sup>I</sup> GLYDF <sup>I</sup> YD <sup>I</sup> EVN <sup>I</sup> RI <sup>I</sup> YVL <sup>I</sup> NM <sup>I</sup> PL <sup>I</sup> GQK <sup>I</sup> QAS <sup>I</sup> MIL <sup>I</sup> IM <sup>I</sup> PH <sup>I</sup> LEP	264
HSP47_2_FRU	179	SSDL <sup>I</sup> RTEL <sup>I</sup> GKK--YTKI <sup>I</sup> MM <sup>I</sup> HRA <sup>I</sup> GLY <sup>I</sup> RF <sup>I</sup> HED <sup>I</sup> I <sup>I</sup> QNM <sup>I</sup> V <sup>I</sup> QV <sup>I</sup> LE <sup>I</sup> EAP <sup>I</sup> LW <sup>I</sup> GGK <sup>I</sup> AS <sup>I</sup> VLL <sup>I</sup> LP <sup>I</sup> FH <sup>I</sup> VED	236
HSP47_TNI	88	SGDL <sup>I</sup> RTEL <sup>I</sup> GKK--YTKI <sup>I</sup> MM <sup>I</sup> HRA <sup>I</sup> GLY <sup>I</sup> RF <sup>I</sup> YED <sup>I</sup> MM <sup>I</sup> NV <sup>I</sup> QV <sup>I</sup> LE <sup>I</sup> EAP <sup>I</sup> LW <sup>I</sup> GGK <sup>I</sup> AS <sup>I</sup> VLL <sup>I</sup> LP <sup>I</sup> FH <sup>I</sup> VES	145
HSP47_1_DRE	205	MVDNRG <sup>I</sup> ELV <sup>I</sup> TR <sup>I</sup> SHT <sup>I</sup> VS <sup>I</sup> VM <sup>I</sup> HRT <sup>I</sup> GIY <sup>I</sup> GF <sup>I</sup> YED <sup>I</sup> TEN <sup>I</sup> RFL <sup>I</sup> IV <sup>I</sup> SM <sup>I</sup> PLA <sup>I</sup> HKK <sup>I</sup> SSM <sup>I</sup> IF <sup>I</sup> IM <sup>I</sup> PH <sup>I</sup> VEP	264
HSP47_2_DRE	203	MVDNRG <sup>I</sup> ELV <sup>I</sup> HR <sup>I</sup> SFT <sup>I</sup> VS <sup>I</sup> VM <sup>I</sup> HRT <sup>I</sup> GIY <sup>I</sup> GF <sup>I</sup> LDD <sup>I</sup> T <sup>I</sup> NKLL <sup>I</sup> V <sup>I</sup> EM <sup>I</sup> PLA <sup>I</sup> HKK <sup>I</sup> SSL <sup>I</sup> V <sup>I</sup> IM <sup>I</sup> PH <sup>I</sup> VES	262
HSP47_3_DRE	215	NQD <sup>I</sup> VRSEL <sup>I</sup> GTK--YTKI <sup>I</sup> VM <sup>I</sup> HRS <sup>I</sup> GY <sup>I</sup> RF <sup>I</sup> YED <sup>I</sup> MM <sup>I</sup> NV <sup>I</sup> QV <sup>I</sup> LE <sup>I</sup> GL <sup>I</sup> WE <sup>I</sup> GK <sup>I</sup> AS <sup>I</sup> VLL <sup>I</sup> LP <sup>I</sup> FH <sup>I</sup> VES	272
HSP47_PMA	235	MVDTRG <sup>I</sup> EL <sup>I</sup> TR <sup>I</sup> SHT <sup>I</sup> ISI <sup>I</sup> QMM <sup>I</sup> HRT <sup>I</sup> GFY <sup>I</sup> NFY <sup>I</sup> H <sup>I</sup> DEKA <sup>I</sup> QV <sup>I</sup> QLLE <sup>I</sup> M <sup>I</sup> QL <sup>I</sup> KGN <sup>I</sup> LES <sup>I</sup> LL <sup>I</sup> IAL <sup>I</sup> PL <sup>I</sup> H <sup>I</sup> TES	294

		300c	
A1AT_HSA	260	LQHLENEL <sup>I</sup> THD <sup>I</sup> IIT <sup>I</sup> K <sup>I</sup> FLEN <sup>I</sup> EDRR <sup>I</sup> SAS <sup>I</sup> LHL <sup>I</sup> PK <sup>I</sup> LSIT <sup>I</sup> GT <sup>I</sup> YD <sup>I</sup> LK <sup>I</sup> SVL <sup>I</sup> GQL <sup>I</sup> GIT <sup>I</sup> KV <sup>I</sup> SNG-ADL	318
HSP47_HSA	277	LERLEKLL <sup>I</sup> TKE <sup>I</sup> QL <sup>I</sup> K <sup>I</sup> TW <sup>I</sup> M <sup>I</sup> GK <sup>I</sup> M <sup>I</sup> QK <sup>I</sup> KAVA <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> GV <sup>I</sup> VEV <sup>I</sup> THD <sup>I</sup> LQK <sup>I</sup> HLAG <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> DK <sup>I</sup> NK <sup>I</sup> ADL	336
HSP47_MMU	277	LERLEKLL <sup>I</sup> TKE <sup>I</sup> QL <sup>I</sup> KAW <sup>I</sup> M <sup>I</sup> GK <sup>I</sup> M <sup>I</sup> QK <sup>I</sup> KAVA <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> GV <sup>I</sup> VEV <sup>I</sup> THD <sup>I</sup> LQK <sup>I</sup> HLAG <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> DK <sup>I</sup> NK <sup>I</sup> ADL	336
HSP47_RNO	277	LERLEKLL <sup>I</sup> TKE <sup>I</sup> QL <sup>I</sup> KTW <sup>I</sup> M <sup>I</sup> GK <sup>I</sup> M <sup>I</sup> QK <sup>I</sup> KAVA <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> GV <sup>I</sup> VEV <sup>I</sup> THD <sup>I</sup> LQK <sup>I</sup> HLAG <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> DK <sup>I</sup> NK <sup>I</sup> ADL	336
HSP47_GGA	265	LERVEKLL <sup>I</sup> NRE <sup>I</sup> QL <sup>I</sup> K <sup>I</sup> TW <sup>I</sup> ASK <sup>I</sup> M <sup>I</sup> KK <sup>I</sup> RSVA <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> V <sup>I</sup> V <sup>I</sup> LEV <sup>I</sup> SHD <sup>I</sup> LQK <sup>I</sup> H <sup>I</sup> ADL <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> DK <sup>I</sup> TK <sup>I</sup> ADL	324
HSP47_XTR	285	LERVEKLL <sup>I</sup> TRE <sup>I</sup> QV <sup>I</sup> KT <sup>I</sup> W <sup>I</sup> VG <sup>I</sup> M <sup>I</sup> TK <sup>I</sup> KAV <sup>I</sup> AV <sup>I</sup> SL <sup>I</sup> PK <sup>I</sup> V <sup>I</sup> S <sup>I</sup> LEV <sup>I</sup> SHD <sup>I</sup> LQK <sup>I</sup> H <sup>I</sup> LDL <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> DK <sup>I</sup> TK <sup>I</sup> ADL	344
HSP47_1_FRU	265	LERLEKLL <sup>I</sup> SKK <sup>I</sup> QVD <sup>I</sup> TW <sup>I</sup> ISK <sup>I</sup> MT <sup>I</sup> NKAVA <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> IS <sup>I</sup> VD <sup>I</sup> V <sup>I</sup> SHN <sup>I</sup> I <sup>I</sup> QK <sup>I</sup> Y <sup>I</sup> SEL <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> VD <sup>I</sup> KAK <sup>I</sup> ADL	324
HSP47_2_FRU	237	LARLDKLL <sup>I</sup> TV <sup>I</sup> QL <sup>I</sup> V <sup>I</sup> SK <sup>I</sup> W <sup>I</sup> LEK <sup>I</sup> SS <sup>I</sup> SS <sup>I</sup> SS <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> AN <sup>I</sup> ISS <sup>I</sup> AL <sup>I</sup> SL <sup>I</sup> QK <sup>I</sup> PL <sup>I</sup> SAL <sup>I</sup> GL <sup>I</sup> VD <sup>I</sup> AWD <sup>I</sup> QK <sup>I</sup> VAD <sup>I</sup> F	296
HSP47_TNI	146	LARLDRL <sup>I</sup> SL <sup>I</sup> QL <sup>I</sup> MS <sup>I</sup> K <sup>I</sup> W <sup>I</sup> LEK <sup>I</sup> SS <sup>I</sup> SS <sup>I</sup> SS <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> AN <sup>I</sup> ISS <sup>I</sup> TL <sup>I</sup> SL <sup>I</sup> QK <sup>I</sup> PL <sup>I</sup> SAL <sup>I</sup> GL <sup>I</sup> VD <sup>I</sup> AWD <sup>I</sup> QK <sup>I</sup> VAD <sup>I</sup> F	205
HSP47_1_DRE	265	LDRL <sup>I</sup> ENLL <sup>I</sup> TR <sup>I</sup> Q <sup>I</sup> LD <sup>I</sup> TW <sup>I</sup> ISK <sup>I</sup> LE <sup>I</sup> RAVA <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> V <sup>I</sup> S <sup>I</sup> MEV <sup>I</sup> SHD <sup>I</sup> LQK <sup>I</sup> H <sup>I</sup> L <sup>I</sup> GEL <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> VD <sup>I</sup> KSK <sup>I</sup> ADL	324
HSP47_2_DRE	263	LERVEKLL <sup>I</sup> TR <sup>I</sup> Q <sup>I</sup> L <sup>I</sup> NT <sup>I</sup> V <sup>I</sup> VSAME <sup>I</sup> QKAVA <sup>I</sup> ISL <sup>I</sup> PK <sup>I</sup> V <sup>I</sup> S <sup>I</sup> MEV <sup>I</sup> SHN <sup>I</sup> LQK <sup>I</sup> H <sup>I</sup> LAEL <sup>I</sup> LGL <sup>I</sup> TEA <sup>I</sup> VD <sup>I</sup> KAK <sup>I</sup> ADL	322
HSP47_3_DRE	273	LARLDRL <sup>I</sup> TL <sup>I</sup> DRLE <sup>I</sup> KWF <sup>I</sup> GK <sup>I</sup> L <sup>I</sup> NST <sup>I</sup> SMAL <sup>I</sup> SL <sup>I</sup> PK <sup>I</sup> RT <sup>I</sup> K <sup>I</sup> MS <sup>I</sup> SAVN <sup>I</sup> LQK <sup>I</sup> QLAAM <sup>I</sup> GL <sup>I</sup> VD <sup>I</sup> AWNET <sup>I</sup> SAD <sup>I</sup> F	332
HSP47_PMA	295	LERLEKLL <sup>I</sup> TK <sup>I</sup> Q <sup>I</sup> LEEW <sup>I</sup> T <sup>I</sup> SK <sup>I</sup> L <sup>I</sup> QK <sup>I</sup> T <sup>I</sup> IAV <sup>I</sup> SM <sup>I</sup> PK <sup>I</sup> GL <sup>I</sup> LQGSAD <sup>I</sup> IKNS <sup>I</sup> LADL <sup>I</sup> GLA <sup>I</sup> EV <sup>I</sup> GDKAK <sup>I</sup> AD <sup>I</sup> F	354

A1AT_HSA	319	SGVTEE--APTKL SKAVIKRAVLTIDEK-GTEAAGAMFLEAIPMSIPPEVKFNKQEVFLMI	375
HSP47_HSA	337	SRMSGK--KDIYLASVFTATAFELDDT-GNPFQDIYGREE-LRSPKLFYADHPFI FLVR	392
HSP47_MMU	337	SRMSGK--KDIYLASVFTATAFEWDTG-GNPFQDIYGREE-LRSPKLFYADHPFI FLVR	392
HSP47_RNO	337	SRMSGK--KDIYLASVFTATAFEWDTG-GNPFQDIYGREE-LRSPKLFYADHPFI FLVR	392
HSP47_GGA	325	SKISGK--KDIYLSNVFTAAALEWDDT-GNPDYADIYGREE-MRNPKLFYADHPFI FMIK	380
HSP47_XTR	345	SKISGK--KDIYLASVFTAAALEWDDT-GNPFSDIYSREE-LRAPKLFYVDHPFVFLIK	400
HSP47_1_FRU	325	SNISGK--KDIYLSNVFTASAVELDVD-GNPDYTSIFGTEK-LKNPKLFYVDHPFI FLVK	380
HSP47_2_FRU	297	SGVSGKAKGKQHL SAVLQWT SLELAAQAG-P-GEDQLEEEI-IEKPKLFYADHPFVFLVR	353
HSP47_TNI	206	SGVSGKSEGKQHL GAVLQWT SLELAAQAG-P-GEEELEEEK-IEAPKLFYADHPFVFLVR	262
HSP47_1_DRE	325	SNISGK--KDIYLSNVFTASALEWDTG-GNPFDPSTIFGSEK-MRNPKLFYADHPFI FLVK	380
HSP47_2_DRE	323	SNISGK--KDIYLSNVFTASAMEWDTG-GNPPDTSIFGTDQ-LKNPKLFYADHPFVFLVK	378
HSP47_3_DRE	333	STLSSLGQKQHL GAVLQWT TLELAPESG-S-KDDVLEDED-VKKPKLFYADHSEFI ILVR	389
HSP47_PMA	355	SGMTGG--REIHLGSLLTAALEFDTE-GEEYDMSVHGHPD-MRNPHLFYIDHPFFLVR	410

A1AT_HSA	376	EQNTKSPLEIMGKVVNP TQK-----	394
HSP47_HSA	393	DTQSGSLLEFIGRLVRLKGDKMRDEL	417
HSP47_MMU	393	DNQSGSLLEFIGRLVRLKGDKMRDEL	417
HSP47_RNO	393	DNQSGSLLEFIGRLVRLKGDKMRDEL	417
HSP47_GGA	381	DSKTNSILEFIGRLVRLKGDKMRDEL	405
HSP47_XTR	401	DEKTDSTLEFIGRLVRLKGDKIRDEL	425
HSP47_1_FRU	381	DNKTNSIMYIGRVVVKPKGDKMRDEL	405
HSP47_2_FRU	354	DNATGALLLMGALDHVEGEAVHDEL	378
HSP47_TNI	263	DNATGALLLMGALDHVEGEAVHDEL	287
HSP47_1_DRE	381	DNKTNSILEFIGRLVRLKGDKMRDEL	405
HSP47_2_DRE	379	DNKTNSILEMGRLLIRPKGDKMRDEL	403
HSP47_3_DRE	390	DNSTGALLMTGALDHTDGPATHDEL	414
HSP47_PMA	411	DARSGATLLIGRCMRMMSGRHDEL	435

## Appendix 8.4: List of marker genes flanking serpin genes.

### Appendix 8.4.1 Marker genes flanking *Ciona* serpins.

Name	Accession id	Brief description
bZIP	ci0100130316	Basic-leucine zipper (bZIP) transcription factor
Pleckstrin-like	ci0100130317	proteins involved in intracellular signaling or as constituents of the cytoskeleton
Lactase	ci0100150830	Carbohydrate transport and metabolism
SEC25	ci0100151625	Vesicle coat complex COPII, subunit SEC23
fCRD	ci0100134603	Secreted frizzled-related protein working as a Wnt antagonist.
CGI-69	ci0100134650	Mitochondrial carrier protein CGI-69
ZPR-1	ci0100134707	Zn-finger, ZPR1 type
PreRYK	ci0100134759	RYK receptor-like tyrosine kinase precursor
DZIP1	ci0100141107	Zn-finger, C2H2 type
GKReg	ci0100135373	Glucokinase regulatory protein
Kv+	ci0100147129	potassium voltage gated Kv channel
NGAP like	ci0100147243	NGAP like protein
SOH1	ci0100147293	Transcriptional regulator SOH1
VTRS	ci0100152710	Valyl-tRNA synthetase
STXBP1	ci0100147806	syntaxin 5
ChMH	ci0100147808	Chondromodulin-1 precursor
SUR4	ci0100147829	Surfeit locus protein 4
NET-7	ci0100147944	tetraspanin 15
SPARC	ci0100148004	Calcium-binding EF-hand
Fbox like	ci0100131254	Cyclin like F-box
WD40	ci0100147754	G-protein beta WD-40 protein

### Appendix 8.4.2 Marker genes flanking group V1 serpin genes in vertebrates.

Gene	Brief description
GMD	GDP-mannose 4,6-dehydratase
WHIP	Werner helicase interacting protein 1
RIPK1	Receptor (TNFRSF)-interacting serine-threonine kinase 1
BPHL	Biphenyl hydrolase-like (serine hydrolase; breast epithelial mucin-associated antigen)
TUBB2A	Tubulin, beta 2A
TUBB2B	Tubulin, beta 2B
SEC5	EXOC2, Exocyst complex component 2
FOXQ1	Forkhead box Q
FOXF2	Forkhead box F2
DUSP15	Dual specificity protein phosphatase 15
IRF4	Interferon regulatory factor 4
PECI	Peroxisomal D3,D2-enoyl-CoA isomerase
RPP40	Ribonuclease P/MRP 40kDa subunit
CDYL	Chromodomain protein, Y-like
GPS	GDP-fucose synthetase
P5CR	Pyrraline-5-carboxylate reductase
EF1D	Elongation factor 1 delta
FVT1	Follicular lymphoma variant translocation 1
VPS4B	Vacuolar protein sorting 4 homolog B ( <i>S. cerevisiae</i> )
SNX-16	Sorting nexin 16

ZFPH	Zink finger protein homolog
MT11	Mitochondrial topoisomerase I
ACVR2B	Activin receptor IIb
FAM82B	Family with sequence similarity 82, member B
SPSB3	SplA/ryanodine receptor domain and SOCS box containing 3
WWP1	WW domain containing E3 ubiquitin protein ligase 1
CKI	Type I cytokeratin
STARTD3	StAR-related lipid transfer (START) domain containing 3
SCF9	Solute carrier family 9
FAM110C	Similar to FAM110C
SHEPB	Similar to sodium-hydrogen exchange protein-beta

#### Appendix 8.4.3: Marker genes flanking group V2 serpin cluster.

Gene	Brief description
DICER	Dicer1, Dcr-1 homolog
GSC	Goosecoid
HEATL	KIAA1622 or HEAT-like repeat-containing protein
DEADB	DEAD (Asp-Glu-Ala-Asp) box polypeptide 24
ITPK1	Inositol 1,3,4-triphosphate 5/6 kinase
GLRX5	Glutaredoxin 5 homolog ( <i>S. cerevisiae</i> )

#### Appendix 8.4.4: Marker genes flanking serpinA7 in mammals.

Gene	Brief description
IL1RAPL2	Interleukin 1 receptor accessory protein-like 2
NRK	Nik related kinase
MUM1L1	Melanoma associated antigen (mutated) 1-like 1

#### Appendix 8.4.5: Marker genes flanking group V3 serpins in vertebrates.

Gene	Brief description
Flanking PAI1 gene	
AP1S1	Adaptor-related protein complex 1, sigma 1 subunit
MUC3B	mucin 3B
MUC12	mucin 12
MUC17	mucin 17
Flanking neuroserpin-pancpin genes	
PDCD10	Programmed cell death 10
GOLPH4	Golgi phosphoprotein 4
Flanking GDN gene	
CUL3	cullin 3
AP1S3	Adaptor-related protein complex 1, sigma 3 subunit
WDFY	WD repeat and FYVE domain containing 1
S28	Serine carboxypeptidase S28
Flanking serpinE3 gene	
GUCY1B2	Guanylate cyclase 1, soluble, beta
ARL11	ADP-ribosylation factor-like 11
WDFY2	WD repeat and FYVE domain containing 2
INTS6	Integrator complex subunit 6/DEADH box 26

**Appendix 8.4.6: Marker genes flanking group V4 serpins in vertebrates.**

Gene	Brief description
OATP	Organic anion transporter polypeptide
SCF	Solute carrier family
SCFL	Solute carrier family like.
NUP98	98 KDa Nucleoporin
RTN4R	Reticulin 4 like receptor 1
DPH1	Region containing DPH1-OVCA2
RPA1	Replication protein A1 70 KDa
WDRD	WD repeat domain
SMYD4	SET and MYND domain containing 4
FMO	Flavin-containing monooxygenase
ZDHHC5	Zinc finger, DHHC domain containing 5
ABP	ATP/GTP-binding protein
5HTAR	5-Hydroxytryptamine 4 receptor
TAR1	Trace amine receptor 1
DOC2B	Double C2-like domain-containing protein beta (Doc2-beta) similar to mouse rabphilin 3A homolog

**Appendix 8.4.7: Marker genes flanking group V5 serpins in vertebrates.**

Gene	Brief description
RC3H1	Ring finger and CCCH-type zinc finger domains 1
STIL	TAL1 (SCL) interruptin
ZBTB37	Zinc finger and BTB domain containing 37

**Appendix 8.4.8: Marker genes flanking group V6 serpins in vertebrates.**

Gene	Brief description
ARRB1	Arrestin, beta 1
GDPD5	Glycerophosphodiester phosphodiesterase domain containing 5
MAP6	Microtubule-associated protein 6
MOGAT2	Monoacylglycerol O-acyltransferase 2
DGAT2	Diacylglycerol O-acyltransferase homolog 2 (mouse)
RPS3	Ribosomal protein S3
EFNB3	Ephrin B3
TRAP	Tudor repeat associator with PCTAIRE 2
GUCY2F	Guanylate cyclase 2F, retinal
LRR	Leucine rich region

## Appendix 8.5: List of Figures.

Figures	Caption to figures	Page No.
1	Three-dimensional structure of uncleaved $\alpha$ 1-antitrypsin (PDB code 1HP7), a member of the serpin family.	4
2	Branched pathway model.	5
3	Physiological functions of selected serpins in vertebrates (A) or invertebrate model organisms (B)	6
4	Gene structure-based phylogenetic classification of vertebrate serpins.	7
5	The phylogenetic tree of animal evolution.	8
6	Models and examples of intron gain.	10
7	Models (A-B) and example (C) of intron loss (taken from Roy and Gilbert, 2006).	11
8	Fate of duplicated single genes (A-C) and duplicated gene families (D-E).	12
9	Exonization of intron sequences to generate novel gene products.	13
10	Exonization of an Alu element.	14
11	Exonization through RNA editing.	14
12	Intron early (E) hypothesis and intron first (F) hypothesis during evolution of life.	15
13	Intron late (L) hypothesis during evolution of life.	16
14	BLAST algorithm.	21
15	Overview of FASTA algorithm.	22
16	An example of the hidden Markov model for protein sequence alignment based on SAM (sequence alignment and modeling) (Krogh et al., 1994).	24
17	Summary of multiple sequence alignment algorithms (Thompson et al., 1999).	24
18	Steps in CLUSTAL algorithm.	25
19	Algorithm of MUSCLE.	27
20	Algorithm of T-COFFEE.	27
21	The UCSC genome browser as seen in Firefox 2.0 web browser.	30
22	The Ensembl browser as seen in Firefox 2.0 web browser.	31
23	The NCBI mapviewer as seen in Firefox 2.0 web browser.	31
24	Protocol for phylogenetic study of vertebrate serpins.	33
25	Generalized scheme of intron position mapping.	35
26	<i>Xenopus tropicalis</i> .	42
27	<i>Danio rerio</i> .	43
28.	<i>Tetraodon nigroviridis</i> .	45
29.	<i>Fugu rubripes</i> .	46
30	<i>Petromyzon marinus</i> .	47
31	<i>Ciona intestinalis</i> .	48
32.	Comparison of gene structure of vertebrate serpins and <i>Ciona</i> serpins.	50
33	Genomic organization of <i>Ciona</i> serpins	51
34	Comparison of genomic organization of <i>Ciona</i> serpin Ci-Spn-5 and Human ATNII.	51
35	Comparison of genomic localization of <i>Ciona</i> serpin Ci-Spn-9, Ci-Spn-10, and human HSP47.	52
36	Phylogenetic tree of <i>Ciona</i> serpins	52
37	<i>Branchiostoma floridae</i>	53
38	<i>Strongylocentrotus purpuratus</i>	54
39	<i>Nematostella vectensis</i> .	55
40	Synteny organization of group V1 serpins in vertebrates.	59
41	Genomic localization of serpinB6-like genes in fishes.	60
42	Genomic localization of Dre-Spn-4 from <i>Danio rerio</i> .	60
43	Genomic localization of four group V1 serpins from <i>Danio rerio</i> namely Dre-Spn-6, Dre-Spn-29, Dre-Spn-30, and Dre-Spn-31.	61
44	Evolutionary tree of group V1 serpins from different vertebrates.	62
45	Synteny of group V2 ( $\alpha$ 1-antitrypsin like) serpin genes in vertebrates	66
46	Spn_94a orthologs unraveled by chromosomal gene order from selected fishes.	67
47	Synteny organization of serpinA7 gene in mammalian genomes.	68
48	Synteny of the angiotensinogen (AGT) genes in vertebrate genomes.	68
49	Synteny analysis of the heparin cofactor II (HCII) gene in vertebrates.	69
50	Genomic organization of the fish specific group V2 serpins – Fru-Spn-7 and Tni-Spn-3.	70
51	Genomic organization of the <i>Danio</i> specific group V2 serpin genes – Dre-Spn-9, Dre-Spn-10, Dre- Spn-11 and Dre -Spn-12.	70
52	Phylogenetic tree of group V2 serpins from different vertebrates.	72
53	Genomic localization of PAI1 genes in vertebrates.	75

---

54	Genomic localization of GDN genes in vertebrates.	76
55	Genomic localization of serpinE3 genes.	77
56	Genomic localization of neuroserpin and pancpin genes in vertebrates and comparative analysis of micro-synteny with serpins of higher invertebrates – Bfl-spn-1 ( <i>B. floridae</i> ) and Spu-spn-1 ( <i>S. purpuratus</i> ).	78
57	Comparison of discriminating amino acid indels among selected human serpins and invertebrate serpins.	80
58	Sequence comparisons among selected group V3 serpins and invertebrate serpins.	81
59	Evolutionary tree of group V3 serpins and related serpins from lancelets and sea urchins.	82
60	Synteny of the group V4 genes, $\alpha$ 2-AP and PEDF.	85
61	Genomic localization of the C1 Inhibitor gene and a fish specific group V4 (FSG4) gene.	86
62	Domain architecture comparisons of C1 inhibitor and a fish specific group V4 (FSG4) serpin.	87
63	Phylogenetic tree of group V4 serpins based the Maximum Parsimony method.	88
64	Synteny comparison of ATIII genes in different vertebrate genomes.	90
65	Sequence logo of RCL region of ATIII from different vertebrates.	91
66	Genomic localization of HSP47 homologs in different vertebrate genomes.	93
67	Evolutionary tree of HSP47 homologs from lamprey to human created with the UPGMA method using MEGA4.	95
68	Distribution of vertebrate serpins based on their intron-coded classification into six groups (V1-V6).	97
69	Comparison of reactive center loops (RCL) of selected group V2 serpins in (A) chicken, (B) <i>X. tropicalis</i> and (C) zebrafish.	99
70	Intron positions of PDCD10 genes in metazoans.	102
71	Exon-intron organisation of the neuroserpin gene lineage.	103
72	Summary of evolutionary history of HSP47-related serpins from fishes.	105
73	Gene structure comparisons between vertebrate (V1-V6) and lancelet (L1-L3) serpin groups.	106
74	Phylogenetic tree of vertebrates emphasizing timescale and lineages displaying intron gain in serpin genes.	108



## Appendix 8.6: List of tables.

Tables	Captions to tables	Page No.
1	Classification of vertebrate serpins.	3
2	Overview of rare genomic change (RGC) markers for phylogenetic purposes	9
3	Genomes analyzed.	18
4	Major databases used.	18
5	Variants of BLAST suite.	22
6	Variants of FASTA.	23
7	Tools for multiple sequence alignment.	25
8	Tools for multiple sequence alignment editing and representation	28
9	Major phylogenetic tools	28
10	Major genome browsers	30
11	Locations of signature sequences in a typical serpin.	34
12	List of genomes and corresponding genome browsers used in building maps synteny of different serpins genes.	37
13	Summary of the phylogenetic methods.	39
14	List of serpins of <i>Gallus gallus</i> .	41
15	List of serpins from <i>Xenopus tropicalis</i> genome	42
16	List of serpins from <i>Danio rerio</i> genome	44
17	List of serpins from <i>Tetraodon nigroviridis</i> .	45
18	List of serpins from <i>Fugu rubripes</i>	46
19	List of serpins from <i>Petromyzon marinus</i>	47
20	List of serpins from <i>Ciona</i> genome draft version v1.95.	49
21	List of serpins from <i>Branchiostoma floridae</i>	53
22	List of serpins from sea urchin - <i>Strongylocentrotus purpuratus</i> genome	54
23	List of serpins from <i>Nematostella vectensis</i>	55
24	Physiological roles of group V1 serpins and associated diseases/syndromes.	56
25	Intron positions of group V1 genes in different vertebrates.	57
26	Physiological roles of group V2 serpins and associated diseases/syndromes.	63
27	Intron positions of group V2 genes in different vertebrates.	64
28	Physiological roles of group V3 serpins and associated diseases/syndromes.	73
29	Intron positions of group V3 genes in vertebrates.	74
30	Intron positions of group V4 genes.	84
31	Intron positions of the ATIII gene in different vertebrates.	89
32	Intron positions of group V6 genes in different vertebrates.	92
33	Sequence comparisons of HSP47 homologs in vertebrates.	94
34	Sequences flanking the insertion points of novel introns in vertebrate <i>serpin</i> genes.	107
35	Genome size of selected animals.	110

## Appendix 8.7 Abbreviations

$\alpha_1$ -AT	$\alpha_1$ -antitrypsin
$\alpha_2$ -AP	$\alpha_2$ -antiplasmin
aa	Amino acid
ATIII	Antithrombin III
AGT	Angiotensinogen
Bfl	<i>Branchiostoma floridae</i>
cDNA	Complementary DNA
Chr.	Chromosome(s)
Ci/CIN	<i>Ciona intestinalis</i>
cpDNA	Chloroplast DNA
Dme/DME	<i>Drosophila melanogaster</i>
Dru/DRE	<i>Danio rerio</i>
EST	Expressed sequence tag
Fru/FRU	<i>Fugu rubripes</i>
Gga/GGA	<i>Gallus gallus</i>
HCII	Heparin cofactor II
HSA	<i>Homo sapiens</i>
HSP47	Heat Shock Protein 47kDa
Id.	Identity
Indel	Insertions/deletions
kb	Kilobase(s)
kD or kDa	Kilodalton
Mb	Megabase (10 <sup>6</sup> )
mRNA	Messenger RNA
Mmu/MMU	<i>Mus musculus</i>
mtDNA	Mitochondrial DNA
MY	Million year(s)
MYA	Million years ago
Nt	Nucleotide(s)
NEURO	Neuroserpin
Nve / NVE	<i>Nematostella vectensis</i>
P1-P1'	Cleavage site in the reactive center loop
PAI1	Plasminogen activator inhibitor-1
PAI2	Plasminogen activator inhibitor-2
PANC	Pancpin
PEDF	Pigment epithelium derived factor
PI	Protease Inhibitor
PMA	<i>Petromyzon marinus</i>
RCL	Reactive center loop
RGC	Rare genomic changes
Rno/RNO	<i>Rattus norvegicus</i>
Scaf.	Scaffold
Spn	Serpins
Spu/SPU	<i>Strongylocentrotus purpuratus</i>
Tni/TNI	<i>Tetraodon nigroviridis</i>
WGD	Whole Genome duplication(s)
Xtr/XTR	<i>Xenopus tropicalis</i>
ZGC	Zebrafish gene collection
ZPI	Protein Z-dependent protease inhibitor

---

## 9. Acknowledgements

I would like to express my gratitude to my supervisor, Prof. Dr. Hermann Ragg, whose expertise, understanding, and patience, added considerably to my graduate experience.

Special thank goes to Prof. Dr. Robert Giegerich for his support as a second supervisor and for his novel ideas and guidelines under scheme of graduate training via graduate college (GK) and the graduate school (GS) specialized in Bioinformatics at the University of Bielefeld.

I would like to thank Prof. Dr. Erwin Flaschel for taking time out from his busy schedule to serve as head of my thesis committee. I would like to thank the other members of my thesis committee for the assistance they provided during thesis submission process.

My sincere thanks also goes to Dr. Heino Büntemeyer for his supportive role in thesis corrections and his support related to computer administration during my work at the Cellular Genetics.

I would like to thank graduate college Bioinformatics “GK635” for their generous support and fellowship in order to complete my PhD.

I am thankful to the University of Bielefeld for its support during my work and stay at Bielefeld. I thank many of known and unknown people in Bielefeld, who somehow or other helped me. I thank Bielefeld as a city and to its inhabitants. It was a remarkable part of my life and I have learned a lot in every sphere of life from this lovely city.

I would like to thanks all my teachers and mentors at different levels of education from school to courses I took in University of Bielefeld.

I would like to thank everyone in AG Zelluläre Genetik in the University of Bielefeld with special thanks to Dr. Olaf Krüger and Dr. Ulf Tödttmann for their support at the beginning of my research work in this group.

I would like to thank many of my friends without naming anyone without their imminent supports; I would be not what I am today.

Last but not least, I would like to thank my parents, my brothers and others in my family for their constant supports.

# Declaration

I, Abhishek Kumar declare that this thesis and the work presented in it are my own and has been generated by me as the result of my own original research.

I confirm that:

1. This work was done wholly or mainly while in candidature for a Ph.D. degree at this University.
2. Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
3. Where I have consulted the published work of others, this is always clearly attributed.
4. Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
5. I have acknowledged all main sources of help.

---

**Abhishek Kumar**

**Bielefeld, 26.02.2010**