

This information has not been peer-reviewed. Responsibility for the findings rests solely with the author(s).

Deposited research article

The rhomboids: a near ubiquitous family of intramembrane serine proteases evolved via multiple horizontal gene transfers

Eugene V Koonin¹, Kira S Makarova¹, Laetitia Davidovic² and Luca Pellegrini²

Addresses: ¹National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894, USA. ²Laboratory of molecular neurobiology, Centre de Recherche Université Laval Robert Giffard, Université Laval, Quebec, Canada.

Correspondence: Luca Pellegrini. E-mail: Luca.Pellegrini@crulrg.ulaval.ca

Posted: 3 October 2002

Received: 30 September 2002

Genome Biology 2002, **3(11)**:preprint0010.1-0010.26

The electronic version of this article is the complete one and can be found online at <http://genomebiology.com/2002/3/11/preprint/0010>

This is the first version of this article to be made available publicly. A peer-reviewed and modified version is now available in full at <http://genomebiology.com/2003/4/3/R19>

© BioMed Central Ltd (Print ISSN 1465-6906; Online ISSN 1465-6914)

comment

reviews

reports

deposited research

refereed research

interactions

information



deposited research

AS A SERVICE TO THE RESEARCH COMMUNITY, GENOME **BIOLOGY** PROVIDES A 'PREPRINT' DEPOSITORY TO WHICH ANY ORIGINAL RESEARCH CAN BE SUBMITTED AND WHICH ALL INDIVIDUALS CAN ACCESS FREE OF CHARGE. ANY ARTICLE CAN BE SUBMITTED BY AUTHORS, WHO HAVE SOLE RESPONSIBILITY FOR THE ARTICLE'S CONTENT. THE ONLY SCREENING IS TO ENSURE RELEVANCE OF THE PREPRINT TO GENOME **BIOLOGY**'S SCOPE AND TO AVOID ABUSIVE, LIBELLOUS OR INDECENT ARTICLES. ARTICLES IN THIS SECTION OF THE JOURNAL HAVE **NOT** BEEN PEER-REVIEWED. EACH PREPRINT HAS A PERMANENT URL, BY WHICH IT CAN BE CITED. RESEARCH SUBMITTED TO THE PREPRINT DEPOSITORY MAY BE SIMULTANEOUSLY OR SUBSEQUENTLY SUBMITTED TO GENOME **BIOLOGY** OR ANY OTHER PUBLICATION FOR PEER REVIEW; THE ONLY REQUIREMENT IS AN EXPLICIT CITATION OF, AND LINK TO, THE PREPRINT IN ANY VERSION OF THE ARTICLE THAT IS EVENTUALLY PUBLISHED. IF POSSIBLE, GENOME **BIOLOGY** WILL PROVIDE A RECIPROCAL LINK FROM THE PREPRINT TO THE PUBLISHED ARTICLE.



The rhomboids: a near ubiquitous family of intramembrane serine proteases evolved via multiple horizontal gene transfers

Eugene V. Koonin¹, Kira S. Makarova¹, Laetitia Davidovic², and Luca Pellegrini^{2*}

¹National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD 20894; ^{2,*}Laboratory of molecular neurobiology, Centre de Recherche Université Laval Robert Giffard, Université Laval, Quebec, Canada

*To whom correspondence should be addressed:

CRULRG, Rm F-5519, 2601 Chemin de la Canardiere, G1J 2G3 Quebec, QC, Canada
Tel. 1 418 663 5000 ext 6879

Fax 1 418 663 5775

Email: Luca.Pellegrini@crulrg.ulaval.ca

Abstract

Background. The rhomboid family consists of polytopic membrane proteins, which show a level of evolutionary conservation that is unique among membrane proteins. The rhomboids are present in nearly all sequenced genomes of archaea, bacteria and eukaryotes, with the exception of several species with small genomes. On the basis of experimental studies with the developmental regulator Rhomboid from *Drosophila* and the AarA protein from the bacterium *Providencia stuartii*, the rhomboids are thought to be intramembrane serine proteases whose signaling function is conserved in eukaryotes and prokaryotes.

Results. Phylogenetic tree analysis suggests that, despite the broad distribution in all three kingdoms of life, the rhomboid family was not present in the last universal common ancestor of the extant life forms, but instead evolved in bacteria and has been acquired by archaea and eukaryotes via several independent horizontal gene transfer events. In eukaryotes, two distinct, ancient horizontal acquisitions apparently gave rise to the two major subfamilies typified by Rhomboid and PARL (presenilin-associated Rhomboid-

like protein), respectively. The subsequent evolution of the rhomboid family in eukaryotes proceeded via multiple duplications and functional diversification through the addition of extra transmembrane helices and other domains in different orientations relative to the conserved core that harbors the protease activity.

Conclusions. Although the near universal presence of the rhomboid family in bacteria, archaea and eukaryotes appears to suggest that this protein is part of the heritage of the last universal common ancestor, phylogenetic tree analysis indicates bacterial origin with subsequent dissemination via horizontal gene transfer. This emphasizes the importance of explicit phylogenetic analysis for the reconstruction of ancestral life forms. A hypothetical scenario of origin of intracellular membrane proteases from membrane transporters is proposed.

Background

Polytopic transmembrane proteins are, in general, not particularly strongly conserved during evolution. Inspection of the database of Clusters of Orthologous Groups of proteins (COGs) [1] revealed only one family of such proteins that is represented in most of the sequenced bacterial, archaeal and eukaryotic genomes. The prototype of this family is the Rhomboid (RHO) protein from *Drosophila melanogaster*, a developmental regulator involved in epidermal growth factor (EGF)-dependent signaling pathways [2-4]. Not only were homologs of Rhomboid detected in prokaryotes and eukaryotes, but the pattern of sequence conservation in this family appeared uncharacteristic of non-enzymatic membrane proteins, such as transporters [5, 6]. Specifically, several polar amino acid residues are conserved in nearly all members of the Rhomboid family, suggesting the possibility of an enzymatic activity. Since three of these conserved residues were histidines, it has been hypothesized that rhomboid family proteins could function as metal-dependent membrane proteases [5, 6]. Recently, however, it has been shown that RHO cleaves a transmembrane helix (TMH) in the membrane-bound precursor of the TGF α -like growth factor Spitz, enabling the released Spitz to activate the EGF receptor, and that a conserved serine and a conserved histidine in RHO are essential for this cleavage [7, 8]. Thus, it appears that Rhomboid family proteins are a distinct group of intramembrane serine proteases. Altogether, the genome of *Drosophila* encodes 7 RHO paralogs (now designated RHO1-7, with the original Rhomboid becoming RHO-1), at least three of which are involved in distinct EGF-dependent pathways, apparently through proteolytic activation of diverse ligands of the EGF receptor [9, 10].

The newly discovered intramembrane proteolytic activity of RHO places the rhomboid family within the framework of regulated intramembrane proteolysis (RIP), a new paradigm of signal transduction, which appears to be prominent in all forms of life [11, 12]. Under RIP, signaling proteins undergo site-specific proteolysis within TMH, resulting in the release of active fragments, which are the actual effectors in signal transduction cascades. Until recently, the only characterized cases of RIP in eukaryotes involved presenilin, an aspartyl protease, which cleaves a transmembrane helix in type I membrane proteins, such as amyloid precursor protein (APP), Notch, and Ire1 [13], and S2P, a metalloprotease, which cleaves a TMH in a type 2 transmembrane protein, the sterol-dependent transcription factor SREBP [11]. Notably, S2P has highly conserved bacterial homologs, and the protease domain of presenilin also might be homologous to bacterial and archaeal type IV prepilin peptidases, although, in this case, the sequence similarity is very low [14, 15].

In the case of the rhomboid family, the existence of homologs of RHO in most prokaryotes is particularly remarkable because animal RHO proteins are involved in signaling pathways that are not found outside metazoa, which seems to make functional conservation in prokaryotes a remote possibility. The only prokaryotic protein of the rhomboid family that has been characterized experimentally in considerable detail is AarA from the bacterium *Providencia stuartii* [16, 17]. This protein is involved in the export of a quorum-sensing peptide, a function that, in physiological terms, resembles that of RHO, although the signaling molecules, other than RHO and AarA, are obviously unrelated [18]. In a striking recent development, two independent research groups have shown that several bacterial rhomboid family proteins, including AarA, were capable of cleaving the EGFR receptor ligands (Spitz, Keren and Gurken) that are normally cleaved

by RHO paralogs [19, 20]. The cleavage depended on the conserved serine and histidine residues paralogs [19] and, moreover, transgenic flies that expressed AarA developed a phenotype indistinguishable from that induced by overexpression of RHO, whereas RHO could substitute for AarA in *Providencia stuartii* [20]. These unexpected findings demonstrated the conservation of a RIP mechanism producing extracellular signals in eukaryotes and prokaryotes.

The near ubiquity of rhomboid family proteins among bacteria, archaea and eukaryotes, along with the remarkable functional conservation, suggest that a signaling mechanism mediated by rhomboids might have functioned already in the last common ancestor of all extant life forms, with subsequent loss in several lineages. To address this possibility, we performed a detailed phylogenetic analysis of the rhomboid family.

Results and Discussion

Sequence and structural features and phyletic distribution and of the rhomboid family

Although the sequence similarity between eukaryotic and prokaryotic rhomboid family proteins is relatively low (at the level of 15-10% identity in the conserved region), the entire superfamily could be retrieved from the protein sequence databases within three iterations of the PSI-BLAST program with a high statistical significance and without any false positives. The conserved core of the rhomboid family consists of six conserved TMH (Fig. 1). The predicted catalytic serine is located in TMH5, whereas the predicted catalytic histidine is in TMH7; TMH3 contains two additional histidines and an asparagine, which are conserved in the great majority of the rhomboid family proteins (Fig. 1). The roles of these conserved residues are not known, but, given the remarkable

evolutionary conservation, it seems likely that they also contribute to catalysis; indeed, it has been shown that the conserved asparagine is required for the cleavage of Spitz by RHO.

When examining the multiple alignment of the Rhomboid superfamily proteins, we noticed that several eukaryotic members appear to be inactivated proteases as indicated by the loss of the predicted catalytic serine or histidine (Fig. 1 and data not shown); these inactivated forms could be regulators of active rhomboid proteases. Several other proteins lack one or more of the conserved residues in TMH3; it remains unclear whether or not these are active proteases.

Bacterial and archaeal members of the Rhomboid superfamily contain 6 TMH, whereas the eukaryotic members typically have an additional, 7th TMH, which may be attached to the core either from the N-terminus or from the C-terminus as discussed below.

The phyletic distribution pattern of the rhomboid family shows that this intramembrane protease is extremely common in all three kingdoms of life, but is not necessarily essential for cell function. Rhomboids are missing in the microsporidium *Encephalitozoon cuniculi*, a eukaryotic intracellular parasite with a highly degraded genome, the archaea *Methanothermobacter thermoautotrophicus* and *Thermoplasma volcanium*, and several bacterial species, primarily parasites with small genomes but also species with moderate-size genomes, such as *Xylella fastidiosum* (see COG0705 at <http://www.ncbi.nlm.nih.gov/COG/>). On two occasions, a representative of the rhomboid family is present in only one of a pair of relatively close genomes (present in *T. acidophilum* but missing in *T. volcanium*; present in the spirochete *Treponema pallidum*

but missing in the related bacterium *Borrelia burgdorferi*), which suggests relatively recent, repeated losses of this gene. Most of the prokaryotic species encode a single gene coding for a rhomboid family protein, although some have two-three paralogs (see COG0705 at <http://www.ncbi.nlm.nih.gov/COG/>); in contrast, eukaryotes show expansion of the rhomboid family, with 7 members in *Drosophila*, and as many as 13 in *Arabidopsis*.

Phylogeny and evolutionary history of the rhomboid family

The multiple alignment of the 6-TMH core of the rhomboid family (Fig. 1) was employed to construct a phylogenetic tree using the least-square algorithm with subsequent optimization using the maximum likelihood method (see Materials and Methods). Only the conserved regions including the TMH and short adjacent stretches shown in Figure 1 were used as the input for tree building, whereas the poorly conserved intervening regions were omitted to avoid the noise from potentially misaligned residues. The resulting phylogenetic tree of the rhomboid family presents a complex and unexpected picture (Fig. 2). Neither the eukaryotic nor the archaeal subsets of the family appear to form monophyletic clades. Instead, the eukaryotic rhomboids are split between two major subfamilies, which are positioned in the midst of different prokaryotic branches (Fig. 2). The first subfamily, which includes 6 of the 7 *Drosophila* rhomboids, clusters with a distinct prokaryotic assemblage, which consists primarily of Gram-positive bacteria as well as a subset of archaeal rhomboids; this clade is strongly supported by bootstrap analysis (Fig. 2). The proteins in this group of eukaryotic rhomboids, which we designated the RHO-subfamily, typically have an extra TMH

added C-terminally of the 6-TMH core; some of these proteins also contain EF-hand Ca-binding domains N-terminally of the core (Fig. 2).

The second eukaryotic subfamily, which we designated the PARL-subfamily, after presenilin-associated rhomboid-like protein (PARL), the human ortholog of *Drosophila* Rhomboid 7 [6], resides within a large, heterogeneous prokaryotic cluster (Fig. 2). Within this subfamily, PARL and its orthologs from other animals and fungi, have a distinct domain architecture, with an extra TMH added to the N-terminus of the core, whereas the rest have only the core (a C-terminal TMH and a ubiquitin-associated domain are appended in one *Arabidopsis* protein; Fig. 2). Thus, the existence of two distinct subfamilies of eukaryotic rhomboids is supported by features of domain architectures that appear to comprise shared derived characters. Within these two major eukaryotic subfamilies, evolution apparently proceeded via both ancient and more recent duplications. Several lineage-specific expansions of paralogs [21] are noticeable, in insects, mammals and plants (Fig. 2).

Archaeal rhomboids are scattered over the phylogenetic tree, with two major clusters and three more isolated proteins joining different bacterial branches (Fig. 2). There is no indication of an affinity between any of the archaeal and eukaryotic rhomboids. Although many of the bacterial rhomboids form phylogenetically coherent clusters corresponding to the established bacterial lineages, there are also several clusters that have odd composition, such as grouping of proteobacterial and Gram-positive species; some of these clusters are well supported by bootstrap (see clusters 1-4 in Fig. 2).

The phylogenetic tree of the rhomboid family tree shown in Figure 2 clearly follows neither the “standard model” scenario [22, 23], with the major split between the archaeo-eukaryotic and bacterial lineages nor the “mitochondrial” scenario, which postulates acquisition of a gene by eukaryotes from the pro-mitochondrial endosymbiont. Neither can this tree be explained by postulating a small number of lineage-specific gene losses. The parsimonious interpretation of the rhomboid family seems to be that the evolutionary history of this family had been replete with horizontal gene transfer (HGT) and lineage-specific gene loss events. In particular, in spite of the presence of rhomboids in the majority of modern life forms from all three primary kingdoms, phylogenetic analysis suggests that this family had not been inherited from LUCA. Instead, the tree topology seems to indicate that this family emerged in some bacterial lineage and afterwards had been widely disseminated via HGT, and then lost in some lineages. Both archaea and eukaryotes seems to have acquired rhomboids on several independent occasions. In particular, at least two HGT events seem to have contributed to the origin of eukaryotic rhomboids, one of them yielding the RHO-subfamily and the other one the PARL-subfamily, with a possible additional HGT in plants (Fig. 2). Given the broad phyletic representation of both subfamilies of eukaryotic rhomboids, both the RHO-subfamily and the PARL-subfamily must have been acquired via HGT at an early stage of eukaryotic evolution, definitely before the divergence of the major crown-group lineages. This early epoch in eukaryotic evolution is thought to have been dominated by HGT from multiple bacterial symbionts [24, 25].

Two alternatives to this multiple-HGT scenario may be considered. One of them would postulate that LUCA already had multiple, paralogous rhomboids, which evolved

via a series of ancient gene duplications, and the odd topology of the phylogenetic tree is due primarily to differential loss of these ancient paralogs. Although this cannot be ruled formally, this hypothesis implies the existence of extremely elaborate signaling system in LUCA, which is hardly compatible with the existing general notions regarding this primitive life form. The second possibility is that the topology of the tree in Figure 2 is simply random. However, the strong bootstrap support for many nodes and the presence of several phylogenetically coherent clusters (above all, the RHO and PARL subfamilies in eukaryotes, but also some of the archaeal and bacterial clusters) seem to argue against this explanation.

The multiple-HGT interpretation of the evolutionary history of the rhomboid family is, at least at first glance, distinctly counter-intuitive, given that this family is nearly ubiquitous among the extant life forms. Indeed, when attempts are made to construct parsimonious evolutionary scenarios on the basis of phyletic patterns [25, 26], there is no chance that such a widespread family is not assigned to LUCA. It should be realized, however, that these approaches are inherently probabilistic and extensive HGT can fool them. For the rhomboid family, this mode of evolution seems to be particularly plausible (Fig. 3). It seems likely that the ultimate ancestor of the rhomboid family evolved from a non-enzymatic integral membrane protein, probably a transporter that might have been involved in an early, primitive form of export of signaling peptides in bacteria. The protease active center might have evolved in such a transporter by chance emergence of the suitable catalytic amino acid within two or three of the TMH (Fig. 3). This would enable the transition from simple transport to the RIP mode of controlled export of signaling molecules. Emergence of RIP could have conferred a major selective

advantage on the respective bacteria and might have resulted in an evolutionary sweep whereby the gene carrying this trait had been repeatedly fixed, rather than eliminated, after HGT. In terms of the evolution of sequence itself, the requirements for the conservation of the protease activity apparently “locked” the rhomboid family in the regime of relatively slow evolution, which ensures the significant sequence similarity among all family members (Fig. 1). The scenario of origin from non-catalytic transporters might potentially apply to other integral membrane enzymes, including intramembrane proteases involved in RIP, such as presenilins and their homologs [14, 15] and the archaeo-eukaryotic signal peptide peptidase [27].

Conclusions

The rhomboid family may be the most widespread and conserved group of integral membrane proteins. In and by itself, this would suggest that this family is part of the gene repertoire of LUCA. However, phylogenetic analysis strongly suggests a different scenario, one of emergence in a bacterial lineage with subsequent multiple independent HGT events and gene losses. In particular, eukaryotes probably acquired their two major rhomboid subfamilies, RHO and PARL, as the result of two independent, early HGT events. These events introducing RIP as a means of intercellular communication might have been pivotal in the evolution of eukaryotic multicellularity along the lines discussed previously with regard to the apparent bacterial origin of key components of eukaryotic programmed cell death machinery [28]. Subsequent evolution of rhomboids in eukaryotes proceeded via lineage-specific expansion of paralogs [21], followed by diversification through the addition of an extra TMH in different positions

relative to the catalytic core, some limited domain accretion (Fig. 2), and sequence divergence.

Phylogenetic analysis of the rhomboid family described here carries a general message for studies aimed at the reconstruction of ancestral life forms, particularly LUCA. Although most of the (nearly) ubiquitous protein families probably do derive from LUCA, explicit phylogenetic analysis is required to ascertain this in each individual case.

Material and Methods

The non-redundant (NR) protein sequence database at the National Center for Biotechnology Information (NIH, Bethesda) was searched iteratively using the PSI-BLAST program with multiple starting queries [29]. PSI-BLAST was normally run with expectation (E) value of 0.01 as the cut-off for inclusion of sequences into the position-specific scoring matrix. Multiple alignments of protein sequences were constructed using the ClustalW program [30] and manually adjusted on the basis of the examination of PSI-BLAST search outputs and the superposition of the predicted transmembrane helices. Transmembrane helices were predicted using the programs TMpred[31] and TMAP[32].

Phylogenetic trees were built using the least-square method [33] implemented in the FITCH program of the PHYLIP package [34], with subsequent local rearrangement using the PROTML program of the MOLPHY package to obtain the maximum likelihood tree [35]. The reliability of the tree topology was assessed using the RELB bootstrap method of MOLPHY, with 10000 replications [36].

Acknowledgments

L.P. is supported by a grant from NSERC Canada.

References

1. Tatusov RL, Natale DA, Garkavtsev IV, Tatusova TA, Shankavaram UT, Rao BS, Kiryutin B, Galperin MY, Fedorova ND, Koonin EV: **The COG database: new developments in phylogenetic classification of proteins from complete genomes.** *Nucleic Acids Res* 2001, **29**:22-28.
2. Sturtevant MA, Roark M, Bier E: **The Drosophila rhomboid gene mediates the localized formation of wing veins and interacts genetically with components of the EGF-R signaling pathway.** *Genes Dev* 1993, **7**:961-973.
3. Sturtevant MA, Roark M, O'Neill JW, Biehs B, Colley N, Bier E: **The Drosophila rhomboid protein is concentrated in patches at the apical cell surface.** *Dev Biol* 1996, **174**:298-309.
4. Guichard A, Biehs B, Sturtevant MA, Wickline L, Chacko J, Howard K, Bier E: **rhomboid and Star interact synergistically to promote EGFR/MAPK signaling during Drosophila wing vein development.** *Development* 1999, **126**:2663-2676.
5. Mushegian AR, Koonin EV: **Sequence analysis of eukaryotic developmental proteins: ancient and novel domains.** *Genetics* 1996, **144**:817-828.
6. Pellegrini L, Passer BJ, Canelles M, Lefterov I, Ganjei JK, Fowlkes BJ, Koonin EV, D'Adamio L: **PAMP and PARL, two novel putative metalloproteases interacting with the COOH-terminus of Presenilin-1 and -2.** *J Alzheimers Dis* 2001, **3**:181-190.

7. Urban S, Lee JR, Freeman M: **Drosophila rhomboid-1 defines a family of putative intramembrane serine proteases.** *Cell* 2001, **107**:173-182.
8. Klambt C: **EGF receptor signalling: roles of star and rhomboid revealed.** *Curr Biol* 2002, **12**:R21-23.
9. Guichard A, Roark M, Ronshaugen M, Bier E: **brother of rhomboid, a rhomboid-related gene expressed during early Drosophila oogenesis, promotes EGF-R/MAPK signaling.** *Dev Biol* 2000, **226**:255-266.
10. Wasserman JD, Urban S, Freeman M: **A family of rhomboid-like genes: Drosophila rhomboid-1 and roughoid/rhomboid-3 cooperate to activate EGF receptor signaling.** *Genes Dev* 2000, **14**:1651-1663.
11. Brown MS, Ye J, Rawson RB, Goldstein JL: **Regulated intramembrane proteolysis: a control mechanism conserved from bacteria to humans.** *Cell* 2000, **100**:391-398.
12. Urban S, Freeman M: **Intramembrane proteolysis controls diverse signalling pathways throughout evolution.** *Curr Opin Genet Dev* 2002, **12**:512.
13. Wolfe MS, Xia W, Ostaszewski BL, Diehl TS, Kimberly WT, Selkoe DJ: **Two transmembrane aspartates in presenilin-1 required for presenilin endoproteolysis and gamma-secretase activity.** *Nature* 1999, **398**:513-517.
14. Steiner H, Kostka M, Romig H, Basset G, Pesold B, Hardy J, Capell A, Meyn L, Grim ML, Baumeister R, Fichteler K, Haass C: **Glycine 384 is required for presenilin-1 function and is conserved in bacterial polytopic aspartyl proteases.** *Nat Cell Biol* 2000, **2**:848-851.

15. Sreekumar KR, Aravind L, Koonin EV: **Computational analysis of human disease-associated genes and their protein products.** *Curr Opin Genet Dev* 2001, **11**:247-257.
16. Rather PN, Orosz E: **Characterization of aarA, a pleiotropic negative regulator of the 2'-N-acetyltransferase in *Providencia stuartii*.** *J Bacteriol* 1994, **176**:5140-5144.
17. Rather PN, Ding X, Baca-DeLancey RR, Siddiqui S: ***Providencia stuartii* genes activated by cell-to-cell signaling and identification of a gene required for production or activity of an extracellular factor.** *J Bacteriol* 1999, **181**:7185-7191.
18. Gallio M, Kylsten P: ***Providencia* may help find a function for a novel, widespread protein family.** *Curr Biol* 2000, **10**:R693-694.
19. Urban S, Schlieper D, Freeman M: **Conservation of intramembrane proteolytic activity and substrate specificity in prokaryotic and eukaryotic rhomboids.** *Curr Biol* 2002, **12**:1507.
20. Gallio M, Sturgill G, Rather P, Kylsten P: **A conserved mechanism for extracellular signaling in eukaryotes and prokaryotes.** *Proc Natl Acad Sci U S A* 2002, **99**:12208-12213.
21. Lespinet O, Wolf YI, Koonin EV, Aravind L: **The role of lineage-specific gene family expansion in the evolution of eukaryotes.** *Genome Res* 2002, **12**:1048-1059.
22. Brown JR, Doolittle WF: **Archaea and the prokaryote-to-eukaryote transition.** *Microbiol Mol Biol Rev* 1997, **61**:456-502.

23. Woese CR, Kandler O, Wheelis ML: **Towards a natural system of organisms: proposal for the domains Archaea, Bacteria, and Eucarya.** *Proc Natl Acad Sci U S A* 1990, **87**:4576-4579.
24. Doolittle WF: **You are what you eat: a gene transfer ratchet could account for bacterial genes in eukaryotic nuclear genomes [In Process Citation].** *Trends Genet* 1998, **14**:307-311.
25. Koonin EV, Galperin MY: *Sequence - Evolution- Function. Computational Approaches in Comparative Genomics.* New York: Kluwer Acad Publ; 2002.
26. Snel B, Bork P, Huynen MA: **Genomes in flux: the evolution of archaeal and proteobacterial gene content.** *Genome Res* 2002, **12**:17-25.
27. Weihofen A, Binns K, Lemberg MK, Ashman K, Martoglio B: **Identification of signal peptide peptidase, a presenilin-type aspartic protease.** *Science* 2002, **296**:2215-2218.
28. Koonin EV, Aravind L: **Origin and evolution of eukaryotic apoptosis: the bacterial connection.** *Cell Death Differ* 2002, **9**:394-404.
29. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ: **Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.** *Nucleic Acids Res* 1997, **25**:3389-3402.
30. Thompson JD, Higgins DG, Gibson TJ: **CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice.** *Nucleic Acids Res* 1994, **22**:4673-4680.

31. Hofmann K, Stoffel W: **TMbase - A database of membrane spanning proteins segments.** *Biol. Chem. Hoppe-Seyler* 1993, **374**:166.
32. Persson B, Argos P: **Prediction of membrane protein topology utilizing multiple sequence alignments.** *J Protein Chem* 1997, **16**:453-457.
33. Fitch WM, Margoliash E: **Construction of phylogenetic trees.** *Science* 1967, **155**:279-284.
34. Felsenstein J: **Inferring phylogenies from protein sequences by parsimony, distance, and likelihood methods.** *Methods Enzymol* 1996, **266**:418-427.
35. Adachi J, Hasegawa M: *MOLPHY: Programs for Molecular Phylogenetics.* Tokyo: Institute of Statistical Mathematics; 1992.
36. Kishino H, Miyata T, Hasegawa M: **Maximum likelihood inference of protein phylogeny and the origin of chloroplasts.** *J. Mol. Evol.* 1990, **31**:151-160.

Figure legends

Figure 1. **Multiple alignment of the conserved core of the rhomboid family proteins.**

The alignment includes the majority of the detected rhomboid family proteins; some closely related sequences were omitted. Only the six conserved (predicted) transmembrane helices and short surrounding regions are shown. The boundaries of the predicted TMH are indicated by shading and overline and the TMH are numbered 1-6. The number of amino acid residues in the omitted terminal and internal regions are indicated. The consensus shows amino acid residues present in at least 90% of the aligned sequences; h stands for hydrophobic residues (A,C,I,L,V,M,F,Y,W) and s for small residues (G,A,S,D,N,V). The proposed catalytic serine (TMH4) and histidine (TMH6) as well as conserved residues in TMH3 with possible ancillary roles in catalysis are highlighted with color. The proteins are identified with the gene identification (GI) number from the non-redundant database and an abbreviated species name. Bacterial species are color-coded green, eukaryotic species blue and archaeal species orange.

Species name abbreviations: Aerpe, *Aeropyrum pernix*, Agrtu, *Agrobacterium tumefaciens*, Anoga, *Anopheles gambiae*, Arath, *Arabidopsis thaliana*; Arcfu, *Archaeoglobus fulgidus*, Bacsu, *Bacillus subtilis*, Brume, *Brucella melitensis*, Caeel, *Caenorhabditis elegans*, Caucr, *Caulobacter crescentus*, Chlte, *Chlorobium tepidum*, Cloac, *Clostridium acetobutlicum*, Corgl, *Corynebacterium glutamicum*, Deira, *Deinococcus radiodurans*, Dicdi, *Dictyostelium discoideum*, Drome, *Drosophila melanogaster*, Escco, *Escherichia coli*, Haein, *Haemophilus influenzae*, Halsp, *Halobacterium sp.*, Homsa, *Homo sapiens*, Lacla, *Lactococcus lactis*, Lisin, *Listeria*

innocua, Metja, *Methanococcus jannaschii*, Metka, *Methanopyrus kandleri*, Metma, *Methanosarcina mazei*, Meslo, *Mesorhisobium loti*, Mycle, *Mycobacterium leprae*, Myctu, *Mycobacterium tuberculosis*, Neucr, *Neurospora crassa*, Nosp, *Nostoc sp.*, Prost, *Providencia stuartii*, Pyrab, *Pyrococcus abyssi*, Pyrae, *Pyrobaculum aerophilum*, Ralso, *Ralstonia solanaraceum*, Sacce, *Saccharomyces cerevisiae*, Schpo, *Schizosaccharomyces pombe*, Sinme, *Sinorhisobium meliloti*, Strco, *Streptomyces coelicolor*, Strpn, *Streptococcus pneumoniae*, Sulso, *Sulfolobus solfataricus*, Sulto, *Sulfolobus tokodaii*, Synsp, *Synechocystis sp.*, Theac, *Thermoplasma acidophilum*, Thema, *Thermotoga maritima*, Thete, *Thermus thermophilus*, Vibch, *Vibrio cholerae*, Xanca, *Xanthomonas campestris*, Xylfa, *Xylella fastidiosa*.

Figure 2. Phylogenetic tree of the rhomboid family.

The sequences and their regions used for the construction of the tree are exactly those shown in Fig. 1. The color coding and abbreviations are as in Fig. 1. The two major eukaryotic subfamilies are denoted as RHO and PARL (see text) and four clusters containing unexpected, from a phylogenetic viewpoint, sets of species are denoted 1-4. Although the tree is shown in a rooted form for convenience, this is an unrooted tree; in particular, the placement of the “root” in the midst of the PARL subfamily is arbitrary. Internal nodes with at least 70% RELL bootstrap supported are denoted by circles. Domain architectures are connected to the respective proteins by brackets or by lines. The domain key is shown in the bottom of the figure.

Figure 3

A hypothetical scenario for the origin and dissemination of the rhomboid family proteases.

The figure schematically shows the proposed three stages of evolution of the rhomboid family:

I – the progenitor of the rhomboid family functions as a transporter for a regulatory peptide in some bacterial lineage

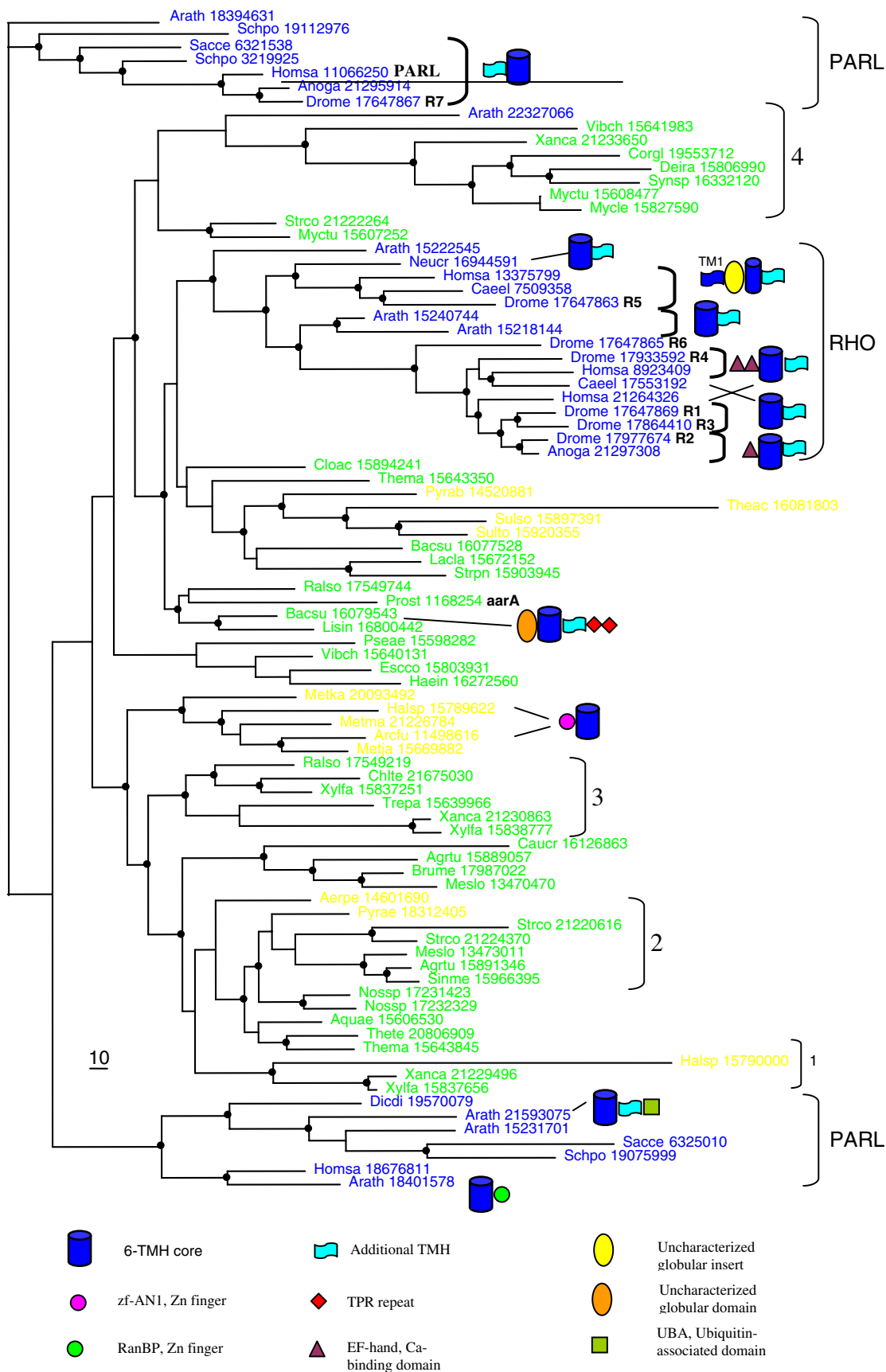
II – the catalytic site of the intramembrane protease evolves allowing the switch to RIP as the mechanism of the regulatory peptide release

III – the emergence of RIP is followed by a burst of HGT

R, regulatory peptide; the transmembrane helices of rhomboid are designated as in Fig. 1, their topology in the membrane is based on that proposed in Ref. 7; the catalytic histidine and serine are shown and connected by a dotted line to indicate the proposed charge-relay system of the protease; possible ancillary catalytic residues are not shown.

	TMH1	TMH2	TMH3
6325010 Sacce	17 LTTGLVPLTATYLLSFFIFA	14 LQMSRLSLYPLHLHLSLPHLFLNVLAIWAPLNLFEET	4 YTGVLNLSALFAGIYLLCLLGLKLLY
19075999 Schpo	10 ILKLPITWQIITYIALLVYA	21 RQLEVEIITYVTLHLSMLHLVFNVSFLPAMSQFEKK	5 CILVTVIPYTLFPGIMHLLVYHFFFL
21593075 Arath	25 LTSSVVVVCVGIYILCLLGT	17 FQVYRFYTAIFHGSLHLVLFNMALVPMGSELERI	6 LYLTVLLATTNAVLHLLIASLAGYN
19570079 Diced	39 ATKVYIICSILFALSFLVAP	19 LDNRLLIILSNFASLSIYHIVYVNMITFDLAK-LERL	1 FGTLLKYLLLFFLFGITITNLIICLFYI
18676811 Homsa	28 PVPVTLATLNLWFLNPNQ	15 KDWQRLLISPLHADDVHLVFNMMASMLWKGINLERR	0 LGSRFWAFYVITLFSVLTGVYVLLLYQ
18401578 Arath	33 PVPVTLASLLAANTLVYLRPAF	21 KDLKRLFLSAFYHVNPHLVYVNMSSLLWKGIKLETS	0 MGSSEFASMVFTLIGMSQGVTLTLLA
11498616 Arcfuc	133 ANNTVLIICTILFFISIVAP	17 AMPVQLITSMFLHVEFWHFFVNMVFLFFGTLELRR	0 LGDRKYLEIFFVSGLAGNVGYIAYS
6321538 Sacce	143 KNLVYALLGINVAVFGLWQ	18 TSKISIIIGSAFSSHQEFWHLGMMNLALWSFGTSLATM	0 LGASNFVSLYMSGAIVAGLSLWYMP
11066250 Homsa	166 QRTVTGIIAANVLVFLWRV	18 VLCSPMLLSTFSHFLSFLHMAANMYVLSFSSSIVNI	0 LGQEQFMAVYLSAGVISNFVSYLKG
17647867 Drome	145 DKMFAPILLCNLAVAFAMVRV	18 VVCWMPFLSTFSHYSAMHLFANMYVMHSFANAAAVS	0 LGKEQFLVAVLSAGVFSGLMSVLYK
18394631 Arath	133 RDVVLGLVLIANAGVFMWVRV	19 GRHLTLITSAFSDIDIGHIVSNMIGLYFFGTSIARN	0 FGFQFLKLYLAGALGGSVFLYLIHH
19112976 Schpo	117 IMVAVIVCLVNGVVFVHWDL	30 GRWWTLVVSIFFSHQNLALHLLVNCVAIYSFSLIVYK	0 FGVWKALSVMYLVGAGVFGNYVALQRM
21295914 Anoga	163 ERIFAPICALNVIVYGLWRI	18 AVCWMPFLSTFSHYSFLHILANMYVLSHFSHAAVAT	0 LGREQFLVYLSAGVIFASPVASHVFK
22327066 Arath	81 ANGIFWITLILNIGIYADHF	15 PAWYQVITATFCHANWNHLSNLLFFLYIFGKLVEEE	0 EGNFGLWLSYLTFGVGANVSLWLV
7509358 Caeel	392 PWFYTYWITTIQIFVCLLSLL	257 NQFYRLFTSLFVHAGVILHALSLLFQYVVMKDLENL	0 IASKRMALYFASGIGGNLASAIFV
13375799 Homsa	165 PFYTYWLVFVHVIITLLVIC	230 DQFYRLWLSLFLHAGVVHCLVSVVVFQMTILRDLKLL	0 AGWHRITAIIFILSGITGNLASAIFV
17647863 Drome	1246 PFFIISISLAELAVFIYYAV	236 DQLYRLLTSLCMLHAGILHLAIFLTLFQHLFLADLERL	0 IGTVTRAVIYFIMSGFAGNLTSAIFL
15240744 Arath	55 SWLVPMFVVANVAVFVAMF	57 KEGWRLLTCTIWLHAGVILHGANMLSLVFIGIRLEQQ	0 FGFVIRIGVYLLSGIGGSVLSLFI
16944591 Neucr	161 PFFVYVFTTQIAVFAIELV	56 NQWWRFITPMLHAGVILHIFGNMMLQMTIGKEMERS	0 IGSIRFFIVYVSAGIFGFVGMGNFA
8923409 Homsa	61 PFFIISISLAELAVFIYYAV	26 EAWRFISYMLVHAGVILHIFGNMMLQMTIGKEMERS	0 HKGLRVLVYLSAGVIFASPVASHVFK
17647865 Drome	72 PWFILLMSFVQISLHWHIASE	13 VEYWRLLTYMMLHSDYWHLSLNICFCQCFIGICLEVE	0 QGHWRVAVVYMGVAVGAGSLANAWLQ
17647869 Drome	102 PWFILVISIIEIAIFAYDRY	26 LQVWRFFSYMFLHANWFHLGFNIVIQVFFGIPLEVM	0 HGTARIGVIYAGVAVGAGSLGTSVVD
17864410 Drome	98 PFFIISISLAELAVFIYYAV	15 LQLWRFLSYALLHSWHLHIFGNMMLQMTIGKEMERS	0 HGSLRTGVYLSAGVIFASPVASHVFK
21264326 Homsa	163 PWFMITVTLLEVAFFLYNGV	26 AQVWRLLTYIFMAGIEHLGLNVVQLLVGVPLEVM	0 HGATRIGLVYVAVGAGSLAVSVAD
17933592 Drome	179 PLTMVLFVSIIEIIMFLVDVI	31 YEGWRVFSYMFVHVGIMHLMNLIQIFLGLIALELV	0 HHWRVGLVYLSAGVIFASPVASHVFK
17977674 Drome	168 PFFIISISLAELAVFIYYAV	24 HEIWRFLTYMFLHAGVILHIFGNMMLQMTIGKEMERS	0 HGSTRITACVYLSAGVIFASPVASHVFK
17553192 Caeel	174 PIFMILLITIQVIGIFFYWE	33 GEAWRFFSYMFLHAGLHNLGNVLIQVLLVGPLEVA	0 HKIWRIGVYLSAGVIFASPVASHVFK
21297308 Anoga	157 PLFVILVTFVGLGFFVYHSL	24 QEVWRFLFYMVLHAGVILHIFGNMMLQMTIGKEMERS	0 HGSTRIGVYLSAGVIFASPVASHVFK
3219925 Schpo	77 RSLVLSIIGINVGVPALWRA	20 INMPSMIVSAFSSHQEFWHLGMMNLALWSFGTSLATM	0 FGNMQVAVYLSAGVIFASPVASHVFK
15218144 Arath	48 TWLVSVFLLQVIMLFAVMTG	52 HEIWRILTSPLWHSGLPHLFINLGLSIFVGYIMEQQ	0 FGPLRIAVIYLSAGVIFASPVASHVFK
15222545 Arath	153 RRWTNVLLAINVIMYIAQIA	18 GQLWRLATASVLANPMLHMLNCYSLNSIGPTAESL	0 GGPKRFLAVYLSAGVIFASPVASHVFK
15231701 Arath	14 ATSCVITLCSVIVFVLIQSLI	15 GHYWRMLTSAHLSVHLHIFGNMMLQMTIGKEMERS	8 YLHYTLVLLVSVSGVAVIYGLYHLLI
18312405 Pyrae	15 PFTVKALVFINVAVFIYELL	16 SEPVRWVTHMFLHGGGLHIVGNMIVLWVFGDNVEDH	0 YGHFRFLALYLMWGLAAAFVHYWAV
15789622 Halsp	94 AFLFLGMVMTVFAVIQYIAP	22 EYVWTVTSVFAGGGFSHIVLNSIVLYFFGPIVEDR	0 IGSKKFVALFLGAGILAGLAQVAGS
20093492 Metka	1 MSLTMLMFLNVLAVIYVSVG	21 VHPRECLITSMFLHAGVILHIFGNMMLQMTIGKEMERS	0 LSTSEFLVYLSAGVIFASPVASHVFK
21226784 Metma	24 APSMMAIFLCLVCSFFLEMV	19 TRPWLTVTYIFLHAGLGHLLFNMIIVLYFFGTALERK	0 VGNKQLLGLFFFTAGILSAIGYFTFL
14520881 Pyrab	28 TFSMLIITAVFIYIEVIGF	16 QQWRLTLTAIFLHMGFVHALNAFWLFLYGLDLEGI	0 VGTKRFLIVFASLAGVIFASPVASHVFK
14601690 Aerpe	19 PIVNMSIIALNFAAFIVGLT	29 ERLYTVFTSMFLHAGVILHIFGNMMLQMTIGKEMERS	0 LGRARYIILYLSAGVIFASPVASHVFK
15669882 Metja	1 -MINLIVIGICIAMPIISVF	16 NMPQVLTSTIFMAGITHLVNLMLVLFIFGTYLENI	0 VGSKKYLIIFLPSGIIINLAIYAYA
15799000 Halsp	96 GVPWGTLLVAGIVAGFYTLV	18 APYLGVTSTPIAGHANGLHIVGNMIVLWVFGDNVEDH	7 RTAARFAGVITRPFYVAVGVPVAGV
15897391 Sulso	35 TFFLMLFVTLGFMVGLLATF	18 GYSELFTSIFITNSFDVIFNFISLVVYIYIFGSR	0 AGKHEY-GIFILAGILGNLLTVIFY
15920355 Sulto	28 TVVLTILITIGYIIGQILSL	18 GFYVQLVTSIFVTPNFPDWAFTIYAMFYIYWLKGE	0 AGKLEY-IIFLIVAGIIVNLSLYLY
16081803 Theac	2 FLFALFLLGGLLISYYPGA	7 RTPWGLTSTIPIYDGSNGVYFLIFALLFSANISH	6 KRTAVALLASVLSGIIANLLDLALY
15598282 Pseae	85 SPMTAAVLLLTFVVAAVTYL	33 QQWRLFTPLMLHFGWHLHAMNAMWFELGRRIEFR	0 QGRPMLLGLTLLFGLVSNVQYAVS
17549219 Raliso	1 -MISILLANVIVVAELF	24 FSPWQLLTYAFHLHAGVILHIFGNMMLQMTIGKEMERS	0 LGRVRTGVYLSAGVIFASPVASHVFK
17549744 Raliso	205 PHLTHALIALNVLAWLATLV	26 GEWWRLTSATFLHAGVILHIFGNMMLQMTIGKEMERS	0 YGVPVYLLIYLAGLGLSALSLSFA
17987022 Brume	17 VIALIGLCVAVVYQNYILS	27 AVIFTFISYSFMSHGSFAHIAVNMIVLWAAFGSPLAGR	0 IGAVRMLFVWVTSVAVGLTHYALH
19553712 Corgl	45 VRTGLTIAIGYVVVIAVHLL	23 SALWGIFTSPLHGSFSLHIGNTVPGFIFSPILIGMS	3 VFWEVTIAGLIGGLGTWIFGGIT
20806909 Thete	14 PVITLSLIIINSLIFFTLSS	32 SNLYPFTSMFLHAGVILHIFGNMMLQMTIGKEMERS	0 MGHIRFLIFPYSLSAGVIFASPVASHVFK
21220616 Strco	39 LCCLLFLISPAAGLNPVYGT	27 GSALTPTALFVHGSWVHLLGNMFLVYVFGAMTEER	0 MGRLOFALFYLGCYALVGYAGAN
21222264 Strco	84 HLVTKILIGINAVAFIAVQA	24 SPEVLSLVTMFTHEEIHIFGNMMLQMTIGKEMERS	0 LGRARYLALYLSAGVIFASPVASHVFK
21224370 Strco	135 ANLVVFLTPFGMAGIASGDG	58 GSWLRLFTALFLHADWSHLLGNLFLVFLFGLPAERI	0 MGHVPPFLYVYGVCGYAAATYFALLD
21229496 Xanca	13 PRWAVPLLFAAVWLAFLWSI	33 GSVLRLFTALFLHADWSHLLGNLFLVFLFGLPAERI	0 LGPWRLLLFLGGAASNLAAIFAI
21230863 Xanca	1 -MITLILITAITGIVSWMAFN	18 QYDLRILTYGFIHADLGHVFNMTLFFFGRYIEDV	0 MTRLRTGSVLYPLFYLGLALVLSILP
21233650 Xanca	140 SRVLRAFNLISLAVLWLVAV	19 KDGILGTAPLLHAGVILHIFGNMMLQMTIGKEMERS	3 ATAMALPMLVTLGGLVWLLGDDPS
21675030 Chlte	17 PPAIKAIITINIVFLFQNS	24 FHLWQPITYLFLHGSFAHIFNMALWFMFVGEIENY	0 WGRTRNFVSFYFCIGGALINLLAT
1168254 Prost	21 IALTTLTLLNVAIVFYQIV	25 GDWWRYPITSMMLHSGNTHLAFNCLALFVIGICERA	0 YGKFKLLAIYIISGIGAAFLSAYWQ
13470470 Meslo	16 LAVLIGICAAVFLVQYVNL	26 FLFTRPPTYAFMHHGGFAHIAINMVMWLAAGSPLANR	0 LGLRFLRAFVAVTGLASVAFWYALL
13473011 Meslo	17 QVVTIGLIVNALVYCATAL	33 PESLSYLYTSFLHADIFHLGGNMLFLVWFGDNVEDA	0 LGHIRYLIYFYLCAIAGAAFGQLVA
15606530 Aquae	14 PIVNLSIIVACSLIWLWYEW	31 QKPYTLTSMFLHAGVILHIFGNMMLQMTIGKEMERS	0 LGKFRYIIFYLICGLGALITQTFIS
15607252 Myctu	37 PVTYTLISLNLAVFVFMQVT	17 GQTYRLLTSMFLHAGVILHIFGNMMLQMTIGKEMERS	0 LGLRFRFALYVAVGAGVILVYLIA
15608477 Myctu	37 VVGGTITLTFVALLYLVELI	18 DGLWGVIFAPLLHANWHLMANTIPLLVGLPMTLA	3 RFWWATAIWIWGLGLTGLIGNVGS
15639966 Trepn	13 TNVTLVSLVANGAVFVITSL	18 RMYQVIFTYGFIHSGVWHLLFNMLGLVFFGQTIIEKK	0 MGSSEMLLFLYVGLTLCGAGACAAV
15640131 Vibch	97 GVFTLFLIMALCIIIFTLQTF	19 WQIWRVYSHALLHFSVWHLAFNLLWWWQFGDLEQR	0 LGSVRLIKLFLYVAVCAIAGALFHGVA
15641983 Vibch	32 LGTITGHVDNLYLLLAISL	23 QQWWRILTGNFAHTNFAHWANMLAALWISFVFKPT	0 ARQLLIPLLIISLAVGVMLILASDMQ
15643350 Thema	3 KRAVYFILLNFNAFIVMMTF	29 GDWFRILTALFVHGGILHIFNLSYALYFFGLIVEDI	0 YGTEKFLVGYFFGTIGVGNLATHVFI
15643845 Thema	14 PYVTIALILINVVVVFVYELM	30 FSLPLFITMFLHGGFWHLGNMMLQMTIGKEMERS	0 MGHVGYTLFVLSAGIIFAALTQGFVT
15672152 Lacla	15 ATYILSITILLVWLWQFFTY	25 SQMWRLLTALFIIHIGWAVLLNVATLFFIGRQIENV	0 FGWLRFTLTYLLSGIFGNAMVFLLT
15803931 Escco	94 GPVTVWMMIACVVVFIAMQI	19 FEFWRVYFTHALMFLSMLHIFNLLWWWYLGGAVEKR	0 LGSGKLIVITLISALLSGYVQKFS
15806990 Deira	50 VKAAAGVTAGLIALWQQEVE	27 FTFWHVFTAPFLHAGVILHIFGNMMLQMTIGKEMERS	3 RFLVATFLIALISGGLVWLLGRSRS
15827590 Mycle	36 MVGGVTITLTFMALLYLVELI	18 DVLWGISFAPVLANWQHLVANTIPLLVGLFGLIALA	3 RFIWWTAMVWVIFGGSATWLIIGNMGS
15837251 Xylfa	10 PTVTKGLLLTNVVVFLFQMM	27 FMPWQLLTYGFLHGGFQHLFFNMLAVFMFGAALEHT	0 WGEKRFLLYLVAVAGVAVGCVQLLVS
15837656 Xylfa	19 WLNAVPLFFFAVLIADFWSI	33 GSALRLFTALFLHADWHLHIFGNMMLQMTIGKEMERS	0 LGSWRLLLFLGGLANLANAAVLTGI
15838777 Xylfa	4 LMITLILITAMNAVSWLSFN	18 RQYDRLLTYGFIHANISHLLFNMTLYFFGSMIEAV	0 MGELTGSLLTYPLFYLGLLVSILP
15889057 Agrtu	32 LVGILAAALAIYVVPAYLLS	27 EWLTVPTVYSFLHGGIEHLIFNGLWMAFGAPVLR	0 IGTVRFVLLWCISAASVAFGHAAALN
15891346 Agrtu	36 QVVTIGLIVNALVYCATAL	34 PDDLTVVTYAFPLHDFVHMLHIFGNMMLQMTIGKEMERS	0 LGHFRFLIFLAVCAIAGALFHGVA
15894241 Cloac	141 MRVTFWLLVIVNFIVYGISAW	26 GQYRLLTSMFLHAGVILHIFGNMMLQMTIGKEMERS	0 YGKLRVTAIYFISGITASFFSYIFS
15903945 Strpn	12 VTSFLLVLTALVFLMLVTA	25 EQVWRLLSAIFVHIGWEHFIVNMLSLYLRQVVEEI	0 FGSKQFFFLYLLSGMGNLGVVVFVS
15966395 Sinme	17 QVVTIGLIVNALVYCATAL	34 PDEFVTVYSFLHGGFWHLGNMMLQMTIGKEMERS	0 LGHFRFLIFLAVCAIAGALFHGVA
16077528 Bacsu	15 YPVVTFILALQAVLWLFSSL	21 GEWWRLTTPILLHAGFTHLFNSMSIFLPAALERM	0 LGKARFLLYVYAGSGIIGNIATVYTE
16079543 Bacsu	177 PTFYTLFIALQILMFSLEI	23 GEWWRLTTPILLHAGFTHLFNSMSIFLPAALERM	0 YGSGRFLLYLVAAGITVCSIASFVS
16126863 Caucr	12 NAPWPAALVAAVAPIPHLLL	20 GRWWTLVVSIFFSHQNLALHLLVNCVAIYSFSLIVYK	0 LGLNVRGGGIFCLFYLVCVIGAVG
16272560 Haein	9 GKITLILITLALCVIYLAQQL	19 SEVWRVISHVTLHLSNLHIFNLSWFFIFGGMIEERT	0 FGSVLLMLYVAVSAGITGVYQNVYS
16332120 Synsp	13 LQSQFSIIVSFLAIFLWLEI	20 EGLRGRVYFAPFLHADFGHLIANSVPMVFLVAVLWMLQ	3 DFVWVTITMVMVAVGGLTGLIAPNT
16800442 Lisin	182 PIVTYFSIIGLIVAVFLWVTF	23 GEWWRFISPIFLHSGLHIFNLSAVMLYVIGAWAERI	0 YGKWRVILILLGGICGNIASPALN
17231423 Nosspp	14 PYFTYGLIGMNVLVFLYHEAS	25 GEWPTLTSQFLHGGFWHLGNMMLQMTIGKEMERS	0 LGHFRFLIFLAVCAIAGALFHGVA
17232329 Nosspp	14 PYVTYGLIANSVFLYHEAS	33 PEWATLTSQFLHGGFWHLGNMMLQMTIGKEMERS	0 LGHFRFLIFLAVCAIAGALFHGVA
consensus/90%h.....hhh..h....hh...h.H.sh.HhhhN.h..h.hs..ht..hhh...shhs.hh..h..

	TMH4	TMH5	TMH6			
6325010	Sacce	4 VAGASGWCPTLFLFAYFVSESQI	9 DYSIPTLYTLPLVLLVAIAVVI	2 SSFWG	FFGLCVGYAIGYKESWF	196
19075999	Schpo	6 IAGLSGWAFPAFISASCVHSPQR	6 LFSIPACYCFPIIYLIMTITLV	2 ASFIG	ASGAVMGYCTPFMLGSI	196
21593075	Arath	12 AIGFSGILFSMIVIETSLSGVT	6 LFNVPKALYPWILLIVFQLLM	2 VSLLG	LCGILSGFSYSYGLFNF	214
19570079	Dicdi	8 HLGFSGVLFPALYIEIENSNGRD	5 AVKIPSKLYPWAMLLIAHVVF	2 SSFIG	FSGIVVGLIFKIGYLDI	219
18676811	Homsa	16 AVGFSGVLFPALKVLNHYCPGG	5 GFPVPPGRFACWVELVAIHLFS	2 TSFAG	LAGLVGLVGLMYTQGPLK	212
18401578	Arath	17 AVGFSGVLFPAMKVVLSQAEDY	4 GILVPTKYAAWAEILLVQMFV	2 ASFLG	LGGILAGIYYLKLKGSY	223
11498616	Arcfu	8 ALGASAAIPGVMGCLAIAIPEI	8 IPINIRNTALLFAAYDFWMMV	10 VANIA	LAGLAVGLYYGKRLGRR	322
6321538	Sacce	4 FGGLSGVLVGLLGHWCWIFQYLA	3 AYRLPRGVVAMMLIWLVLVCL	10 TANGA	VGGLLVGCLSGLLGLL	265
11066250	Homsa	25 TLGASGAIAMLMGVLTLNPLGL	7 IPMLWLATGLFAAYSIFVSG	8 VAQLA	LAGLIGLGLYAKLKR	302
17647867	Drome	5 VVGASAAIPGLLGCCLTMLRPMS	6 IPMPLALFVLYAALALFVIQ	6 VAHAG	LVGMIVGGVLLALYRPS	184
18394631	Arath	10 SLGASGALFGVLGCLSYLFPFA	5 VFPVPGGAVWAFSLASVAVNAA	8 FDYAA	LGGSMGVLYGWIYSKA	330
19112976	Schpo	8 SLGASGAIMTVLAAVCTKIPEG	6 LPMFTFTAGNALKAIAMDTA	8 FDHAA	LGGALFGIYVVTYGH	352
21295914	Anoga	8 SLGASGAIMTLLAYVCTQYPTD	6 LPALTFSAGAGIKVLMGIDFA	8 FDHAA	LGGAMFGIFWATYGAQ	330
22327066	Arath	22 GLGASGAVNAIMLLDLFLHPRA	6 FIPVPAMLLGIFLIGKDIIRI	6 ISGSA	LGGAAVAA-IAWARIRK	331
7509358	Caeel	60 LLGASGAVYATAAIFACLFPYT	4 FFVYVVKAGIFMPLDFIAEYV	11 VAFDA	VSGTFFGVVSSFLLLPA	368
13375799	Homsa	8 SLGASGAIMGILAYVCSQYPTD	6 LPMYTFSSAGAAIKVIMGIDLA	8 FDHAA	LGGALFGLFWCHFGSQN	349
17647863	Drome	5 SVGASGAVPGLFAISVLVKMSW	8 LILGLQFVIERVMEAAQASAGL	12 VNHA	LQSGALVGVVLLVWLSKF	267
15240744	Arath	4 AVGFSQAQCGILAAVIVCCDN	8 WALVQHILVTLVLVLCIGFIPW	0 -VDNWA	LFGTIFGLLTTIIFPY	807
16944591	Neucr	4 EVGPAGSQFLLACLFLVFLQS	8 KAFNLNSAIVLFLFCIGLLPW	0 -IDNIA	IFGFLSGLLLAFAFLPY	553
8923409	Homsa	4 EVGFSASLGGVVALVLRGIFLN	9 IALPKLILLCSVLVIGITLTPY	1 LNFGLLGLAGVICGLLTMSPVF	1642	
17647865	Drome	4 SVGASGALFGLLGSMSELFTN	8 AALLTLFVILINLAIGLILPH	0 -VDNFA	VGGFVTFGLLGFILLAR	270
17647869	Drome	5 TTGASGALFGIALLLDDLLYS	8 KDLFLIGLDIVISFVLGLLPG	0 -LDNFA	IGGFLAGLALGICVQLS	418
17864410	Drome	4 LVGASGGVYALMGGYFMNVLVN	10 FRLLIITLIVLDMGFALYRR	9 VSFPA	IAGGFAGMSIGYTVFSC	256
21264326	Homsa	4 LMGASGVYAMLGSHVPHLVLN	8 ARIASLILLSDVGFTTYHF	9 TSLEA	IGGGVAGILCGFIYVRR	252
17933592	Drome	4 LVGASGGVYALLAAHLANITLN	8 TQLGSSVIVFVSCDLGYALYQ	12 VSYIA	LTGALAGLTIGFLVLKN	298
17977674	Drome	4 LVGASGGVYALLAAHLANVLLN	8 IQLMVLVLFVFCDLGYALYSR	12 VSYIA	MTGALAGISVGLLLLRQ	283
17553192	Caeel	4 VVGSAGVYALVSAHLANIVMN	10 LRMAVALICMSMEFRAVWLR	11 PSFVA	LGGVAVGITLGVVLRN	360
21297308	Anoga	4 LAGASGGVYALTAHATIIMN	8 VQLLAFVFCFDLGTSVYRH	7 IGYVA	LSGAVAGLLVIGIVLRN	375
3219925	Schpo	4 LVGASGGVYALLAAHLANVLLN	8 IKLLHLVLFVFSDFGFAIYAR	25 VSYVA	LAGAIAGLTIGLLVLKS	375
15218144	Arath	4 LVGASAGVYALIFAHVANVLLN	8 IRVLVLFVFIPLDFGGAIHRR	8 VSHLA	IAGAVTGLFFGYVVLVN	373
15222545	Arath	9 IVGASGGVYALLAYAVLFPFR	9 PMPAWLFAVYALVELTLGIS	5 IAHFA	LGGMAGSGVLLWRW-LR	191
15231701	Arath	5 VVGASGAVPGVAGALVAVIRQY	8 SKRLLTQIGLFLVLSLVQGLT	3 VDNAA	LIGGLGCLLACLIPAR	393
18312405	Pyrae	6 LVGASGAISSGMMGAAARYGFRR	21 LKPVLIIPVGVVFLINIVTGLY	9 IAWEA	IGGFVAGFVGLPMDRP	226
15789622	Halsp	20 AVGASGAISSGVLGAYMVLPHFA	15 IP-AWAYIGVFWFIYQLFYGAL	9 VAYFA	IGGFVAGLALYIYRR	220
20093492	Metka	1 HIGASGLIYGLWGLYLVIRGFIN	3 KQPLGLVLAFAIYISGLFWGLL	5 VSWQG	LFGLAGGIGAGAFIASD	226
21226784	Metma	6 VVGASGAIAGIMGAYVFLPESA	16 PIPAVVYVFLVFLVQLYSGMV	11 IAWWA	IGGFVAGVLLNRFFLRD	225
14520881	Pyrab	6 LVGASGAISSAVLGAFLFLPRA	14 RFPAWVPLFPVWSLQWLAAGR	6 VAYLA	LVGFLGFAFAVAVRFR	238
14601690	Aerpe	5 TLGASGAIPLFGATLALVLR--	1 LNADMRPVVILLVLSLIFTFT	3 ISWQA	VGGLVAGAVIGYAMLHA	265
15669882	Metja	6 LIGASGAIAGVLYGAYLVLYPRA	14 RLPAPLWLGWFGVQLQAVYSSG	8 VAYVA	VVGFVVGMLIAWPLRRG	363
15790000	Halsp	10 MVGASGAIYGVFAALVLEPNL	6 VPMRLKHALLLFAVDFDLMVN	4 IAHTA	LSGLFVGLYMGYRIRKM	209
15897391	Sulso	6 IIGASGAVSALIGTYLALFPFA	15 RVPAPLILGAWAVLQVVFAYI	6 VAWSA	IAGVVFVIGVYGLYVRA	219
15920355	Sulto	12 SLGASGAVSAMPLFAFLFLKPT	7 PAPAIIYAVFYVGYSLWMDRR	4 INHSA	LAGAAFGVMMFLMIEPR	187
16081803	Theac	1 HLGASGVTHGLMFLVFLVGLLR	3 PAIATSMIFLIFYGMLMRTL	5 VSWQG	LGGAVAGLIAALLRLR	303
15598282	Pseae	4 LVGASGGVYALLAAHLANVMLN	8 LRLLAIFLFAVSCDVGFAYSR	11 VSYVA	LTGALAGLTIGLLVLKN	350
17549219	Ralso	6 LIGASGAIPLGVLAFVMMFPDR	8 PIKTKYVAGYALIEFIMGLG	9 IAYFA	LGGMLFYIYIYVIRRN	210
17549744	Ralso	19 VVGASGAIMGVAIAAIVYLKIV	14 QKYQLYNLIAMIALTLINLQ	2 VDNAA	IGGAIAGLISIAIYILV	227
17987022	Brume	11 SLGASGAIYAAATSYYFFPNA	6 LPFPIKIGVALLGLMAFDW	15 IDHAA	LGGGIFGWLKAKYGYST	275
19553712	Corgl	6 LVGASGAISSGMMVAAARFGFRT	21 SRGVVPLVAVVMIINLATGLL	9 IAWEA	IGGFVAGFGLRWFDRR	224
20806909	Thete	6 LIGASGAIAGVYVAYLLYPERV	12 RIPAFIPILVWLVFQVFMFAA	5 ISWAG	IGGIAGAVLVLVLR	219
21220616	Strco	5 SGGASGGLFVIGALLSIEGVL	3 IQKALINALALFLINSIF--	2 VNIFA	FGGLVTLGLVLYFYGIW	197
21222264	Strco	22 AVGASGAISSGVLGAYALLIPFS	15 SVPASIPFGWVYVQLVMGLA	9 IAFWA	VGGLTGVVALAPLVDK	240
21244370	Strco	4 SISSGAAFFGLIGAMLSALAKN	8 SALAIIFTITFVNFILGFLF	0 -IDNFA	IGGFVAGLGLVLLFK	258
21229496	Xanca	10 SVGASGAIPLVGSVAVFVIRH	8 EDLMQIAQIIALNMAMGLMSR	1 IDNWG	IGGLGGTAMTWLLGPQ	336
21230863	Xanca	12 AVGYSQVVPWMTILSVKQPS	6 LLSLIPISFAPFESLIFTSIV	2 ASFLG	LSGLLVGYAISWGLIG	202
21233650	Xanca	9 MVGASGAISSGVLGAYMMPFPHA	15 ELPAVPIPLGWLFFFQIINGII	9 VAWYA	IGGFVAGLGLVLYVFRKR	224
21675030	Chlte	5 TAGASGAVFLFGATFMVAR--	1 LHLVDRVVALIVINLAFTFL	3 ISWQG	VGGLVTLGALVAATYVYA	207
1168254	Prost	6 HIGASGLIFGWLAFVFLVGLFV	3 WDIVIGLVVLFVYGGILLGAM	8 VSWQG	LSGAVAGVVAAYLLSAP	221
13470470	Meslo	9 LIGASGAIPLGILGLFAAGFRK	8 PIPAPLIVGYLIFEIFDLFF	4 VSHLT	LLGVLFAGVYIRIRFGI	198
13473011	Meslo	3 FGGLSGVVYALAGYLWILQRA	3 GLSIPRSLMGFMLIWLVLGYV	5 IANTA	LAGLISGVVLAWFDSQR	273
15606530	Aquae	1 YVGLGTLHGLFAYYALNEALN	5 WLLVGLVIGKVAWEQWFGASA	9 VATEA	LAGLVGGLLLAAGHCLF	216
15607252	Myctu	4 SVGASGAIPLGILGLFAAGFRK	3 FFMKPVTVGWSLPLIILNVVY	7 INNAA	IGGFLVGLVLLGYTSPF	192
15608477	Myctu	6 MVGASGAVSVMGAYVFLFPYS	15 EIPAFYIYLMIFWFIQVNLGLV	4 IAWWA	IGGFVYGMIVGYILRMR	215
15639966	Trepa	8 SVGASGAIPLGALALAPLH	8 IPVNIIRVAVIIFALIDLILL	6 IAHTI	LAGLITGLIFGKLLYRK	184
15640131	Vibch	4 SAGASISPLGFAAVVGLAFFT	4 LQIQIRMTVTLIVANLVMLNF	4 VSIWA	IGGAIAGLGLLSAIFAPK	198
15641983	Vibch	14 AVGFSGVVPAFAGFALLKYPLA	4 VAARDAISVFWRTLLEPVTEA	13 VAVQG	LFGLLGLAALAVAVLVH	298
15643350	Thema	3 FGGLSGVVYALMGYVWVLRGERD	3 GIYLRGLIIFALIVIVAGWF	6 MANGA	IAGLAVGLAMAFVDSL	271
15643845	Thema	1 HLGASLVPFVGLAYLLGVGWKE	3 LSVVAVIAFALYGGVWLVGL	5 ISWEA	ILFGVIGGLVAAALLHRK	228
15672152	Lacla	6 HIGVSLIFGWLAFVFLVGLFV	3 W-DIIGCMVLFAYGGVLLGVM	8 VSWQG	LCCAISGVVVAAYLLSAP	219
15803931	Escco	8 VLGASGGVFLMLMAYGMLFPNE	9 PMKARTFVILYGVIELLMGIT	5 VAHFT	LGGMLFVGLLIRYWRGQ	205
15806990	Deira	6 IIGASGAVSALIGSYLALPFGA	15 RVPAPFLIGFWALLQVVFAYT	6 VAWSA	LAGFVSGVYVGSVCRAT	225
15827590	Mycte	12 SLGASGAVSAMPLFAAVLLQPWA	7 PAPAIFYAVFVYVGSIMMERR	4 INHSA	LSGAAFGVVMFLCMEQP	191
15837251	Xylfa	6 LIGASGVVYALMGAACRFAPVP	22 NRTVLIPTLMLVFGNVLIAIG	10 IAWDA	VFGFLGLFVFLSFLDRP	243
15837656	Xylfa	6 LIGASGAVSVAAYFLLHPRV	12 PLPAPLIPALWIGQQFLMLAL	5 VSWAG	VGGILAGAIMVIFMRRP	239
15838777	Xylfa	4 SVGASGAIPLGLGAAVFGFKL	4 GKAFANMVGVFALNIFISFT	3 IDIFA	FGGLVGGVVSIVLIGRT	324
15889057	Agрту	5 SSGASGGIFGLLSYTYFDLKL	4 GYVGLVFLVSVFVGSDDLIF--	2 VNVVA	IGGILGGIMYAVVYLLI	207
15891346	Agрту	4 AAGASVSLVGLFYAIIVLRAT	4 IQQLGQSYLTLFVNIIGSVL	3 ISLAG	IGGAVGAFVAVIFVPR	194
15894241	Cloac	5 SAGASGGIFGLFYAYTVTDYKL	4 NQISILLVSVFVILSDTLFP	2 VDIWA	TGGLLTGLLSLFFKFI	201
15903945	Strpn	6 LIGASGAISSVVAAYFLLHPKV	12 PLPAAIPLAFVWIGQQFMMFLA	5 VSWSA	VGGIVAGLVLVLLVLRP	220
15966395	Sinme	5 HVGASGAIPLGILGLFAVFLFR	4 QAEHSMKITLALFAVLMSPFI	3 INMMA	LFGLVGGFLLSFLCVQK	194
16077528	Bacsu	3 SAGASGAIPLGILGLFAVFLFR	4 LRTIGTNIIVIIINLGFGFA	3 IDNSG	IGGLVGGFFAAALGLP	356
16079543	Bacsu	8 SYGQSGVVYALMGSAASMLLD	34 SLALIFLLTFVFMALDIKAFY	5 IDSFV	AMAFGSSAIFIIISYT	208
16126863	Caucr	11 VVGASGAIAGLMSAAARTMDSA	8 GPRVIVSLGLVWVNLVLAVT	10 VAWEA	LIGFVAGVLLIGPFARW	207
16272560	Haain	3 FGLSGVVYAVLGVYVIRDKLN	2 LFDLPEGFFTMLVGIAGLFI	7 MGNAA	ISGLIVGLVIGFIDSKL	186
16332120	Synsp	1 TVGASLIFGLVGLLFRGWFO	3 ASIVLSIVVLYVGLSALWGL	5 VSWQG	ILFGVIGGAAIAAWLIARE	191
16800442	Lisin	3 SVGASVAVFVAGLGLVYVLLK	4 AKTIGTISIALVAINLLIDVF	3 IDIAG	IGGLVGGFLLAGALSPL	361
17231423	Noss	6 SLGASGAISSVLYGAYLIRFPA	15 SVPALVLIIGFVQNVISGLV	14 VAYWA	IGGFVAGFVILAPLGLF	220
17232329	Noss	6 SLGASGAIAGVMGAYLIRFPNA	15 RVPAYFVILGFWFLQSFYGLA	14 IAYWA	AGGFVAGLGLLGLLGL	217
consensus/90%		.hGhSssh.uhhh.....hh.hh.h.h....	hs..sHh.Ghh.Ghhh.....		



Koonin et al., Fig. 2

