

Evolution in the Structure and Function of Aspartic Proteases

Jordan Tang and Ricky N.S. Wong

*Laboratory of Protein Studies, Oklahoma Medical Research Foundation and
The Department of Biochemistry and Molecular Biology, University of Oklahoma
Health Sciences Center, Oklahoma City, Oklahoma 73104*

Aspartic proteases (EC3.4.23) are a group of proteolytic enzymes of the pepsin family that share the same catalytic apparatus and usually function in acid solutions. This latter aspect limits the function of aspartic proteases to some specific locations in different organisms; thus the occurrence of aspartic proteases is less abundant than other groups of proteases, such as serine proteases. The best known sources of aspartic proteases are stomach (for pepsin, gastricsin, and chymosin), lysosomes (for cathepsins D and E), kidney (for renin), yeast granules, and fungi (for secreted proteases such as rhizopuspepsin, penicillopepsin, and endothiasepsin). These aspartic proteases have been extensively studied for their structure and function relationships and have been the topics of several reviews or monographs (Tang: *Acid Proteases, Structure, Function and Biology*. New York: Plenum Press, 1977; Tang: *J Mol Cell Biochem* 26:93-109, 1979; Kostka: *Aspartic Proteinases and Their Inhibitors*. Berlin: Walter de Gruyter, 1985). All mammalian aspartic proteases are synthesized as zymogens and are subsequently activated to active proteases. Although a zymogen for a fungal aspartic protease has not been found, the cDNA structure of rhizopuspepsin suggests the presence of a "pro" enzyme (Wong et al: *Fed Proc* 44:2725, 1985). It is probable that other fungal aspartic proteases are also synthesized as zymogens.

It is the aim of this article to summarize the major models of structure-function relationships of aspartic proteases and their zymogens with emphasis on more recent findings. Attempts will also be made to relate these models to other aspartic proteases.

Key words: aspartic proteases, pepsin, gastric enzymes, lysosomal protease processing, zymogen activation, cDNA

The amino acid sequences of a number of aspartic proteases have now been determined. By comparing the completed sequences available so far, it is clear that regardless of the biological sources, the aspartic proteases are homologous in se-

Received June 10, 1986; accepted August 19, 1986.

© 1987 Alan R. Liss, Inc.

Code 190 200 210 220 230 240 250 260
PF YWQITLDSITMDEGTI--ACSGGQQAIVDTGTSLLTGPSTSAIA-NIQSDIGA--SENSDGDGMVVISCSISLSDLPDIVF
HP YWQITVDSITMNGEAI--ACAEGCQQAIVDTGTSLLTGPSTSAIA-NIQSDIGA--SENSDGDGMVVISCSISLSDLPDIVF
YFQITVDSITMNGEAI--ACAEGCQQAIVDTGTSLLTGPSTSAIA-NIQSDIGA--SENSDGDGMVVISCSISLSDLPDIVF
CF YWQITVDRWTVGKVV--ACFPFCQAIVDTGTSLLVMVQGAAYN-RIKDLDGV--S--SDGE--ISCDISKLPDIVF
HG YWQITGVEFLICGQASGWCSEGCQAIVDTGTSLLTVPQQYMS--ALLQATGA--QEDDEYGGFLVNCMSIQLMLTLTF
MG YWQITGVEFLICGQASGWCSEGCQAIVDTGTSLLTVPQQYMS--ALLQATGA--QEDDEYGGFLVNCMSIQLMLTLTF
KC YWQITGVEFLICGQASGWCSEGCQAIVDTGTSLLTVPQQYMS--ALLQATGA--QEDDEYGGFLVNCMSIQLMLTLTF
HR YWQIQMKGVSVGSSIL-LCEEGDCLAVDTGTSYFSCSSISSIE-KLMEALGA--K-KRLLFDYVVKVCKMGGTLPDIF
OR SWQITMKGVSVGSSIL-LCEEGDCLAVDTGTSYFSCSSISSIE-KLMEALGA--K-KRLLFDYVVKVCKMGGTLPDIF
SR SWQITMKGVSVGSSIL-LCEEGDCLAVDTGTSYFSCSSISSIE-KLMEALGA--K-KRLLFDYVVKVCKMGGTLPDIF
FD YWQIHNNQVAVGSLL-LCKXGCEAIVDTGTSLLVGVQPEVR-ELKXAIQA-VPLIQEYMYPCCKVSTLPAITL
HD YWQVHLDQVEVAGSLL-LCKXGCEAIVDTGTSLLVGVQPEVR-ELKXAIQA-VPLIQEYMYPCCKVSTLPAITL
RZ WCGITVDRATVGTST--VAAS-FDGLDGTLLLLLPNNVAASVARY-GA--SDNDGDTYITISDTRFKPK-LVF
PW FWSFVWDSYTAGSGG--DGF--SCIDATGTLTLLNDSVVSQYVSGAQQSSNSAGCYVFDCT--NLPDFSV
EN FWEVTSITAVGSGV--KSTS-IDGIDATGTLTLLYLFATVVSAYWAQVSGAQQSSNSAGCYVFDCT--NLPDFSV
PA YWEVRFEGIGLDETV--AELESHGAIDTGTSLITLPSGLAE-MINAEIGA-KKGWIGQYITLDCNTRNLDPLTF

Code 270 280 290 300 310 320
PF IDLGGVYPLSPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
HP IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
YF IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
CF IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
HG IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
MG IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
KC IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
HR IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
OR IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
SR IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
FD IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
HD IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
RZ IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
PW IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
EN IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA
PA IINGVYVFPASAYILQDDSDS---CTSGFEGMDVPTSSGE-L--WILGDVFIKQYTYVFDRAANNKVGGLAPVA

Code 60 70 80 90 100 110 120
PF DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
HP DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
YF DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
CF DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
HG DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
MG DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
KC DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
HR DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
OR DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
SR DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
FD DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
HD DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
RZ DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
PW DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
EN DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL
PA DSST-YQTSSTGSLSTIYGTGSM--TGILGYDVTQV---GGLSDTNQIFGLSETEPGSLIYVAFDGL

Code 130 140 150 160 170 180
PF GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
HP GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
YF GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
CF GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
HG GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
MG GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
KC GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
HR GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
OR GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
SR GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
FD GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
HD GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
RZ GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
PW GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
EN GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG
PA GLAYPSISASGAT--PVFDNLDQGLVSDQDIFSVYLSN---DQ--GSVVLGLGIDSSYVGGSLNWPV-SVGG

Code P1 P10 P20 P30
Pepsinogen, Pig (77,78) PF LVKVFVLRKKSRLRQNLKDGKLLK-DFLKTKHKNPASKYV-----
Human (52) HP IYKVFVLRKKSRLRQNLKDGKLLK-DFLKTKHKNPASKYV-----
Monkey (33) MF IYKVFVLRKKSRLRQNLKDGKLLK-DFLKTKHKNPASKYV-----
Chicken (79) CP SIHRRVFLKKGKSLRQNLKDGKLLK-DFLKTKHKNPASKYV-----
Progastatin, Human (33) HG AVKVFVLRKKSRLRQNLKDGKLLK-DFLKTKHKNPASKYV-----
Monkey (30) MC AVKVFVLRKKSRLRQNLKDGKLLK-DFLKTKHKNPASKYV-----
Prochymosin, Bovine (34,35) BC AETIRFLYKKSRLRQNLKDGKLLK-DFLKTKHKNPASKYV-----
Renin, Human (80) HR TFLGPDITFPRIFLKRKMPISRESLKERGV-----DMARLGPWSPQPKR-----
Mouse kidney (81) HP TFLGPDITFPRIFLKRKMPISRESLKERGV-----DMARLGPWSPQPKR-----
Mouse Submaxillary Gland (54,82) SR TFLGPDITFPRIFLKRKMPISRESLKERGV-----DMARLGPWSPQPKR-----
Cathepsin D, Pig (40) PU LVRIPLHKFISIRRTMSEVGGSV-EDLIAXG--PVSKYSQAV-----
Prochymosin, Human (41) RD QLTLPLETRKSA-IPLAKPNMNY-PSAKNAIQRAIKYKHKINTSTGG
Rhinopapsinogen (55) RZ
Penicillopepsin (7) PW
Endothiapepsin (83) EN
Pro-proteinase A, yeast (50,84) PA AKV--HKAIKYKHELSDENKEVIFQHLAHL--GQKYLTPQEKANPEV

Code P40 1 10 20 30 40 50
PF --P--P--GAALIGDE-PLNLYDTEVY--GTIGGTPAQDFTVIFDTGSSNLWVPSVYCS--LACGDHNRQ-FNPD
HP --P--P--PQAEATLWDEGPLENYLDMEVY--GTIGGTPAQDFTVIFDTGSSNLWVPSVYCS--LACTHNRH-FNPF
YF --P--P--PQAEATLWDEGPLENYLDMEVY--GTIGGTPAQDFTVIFDTGSSNLWVPSVYCS--LACTHNRH-FNPF
CF --P--P--VYATISYPMYNDASVY--GTIGGTPAQDFTVIFDTGSSNLWVPSVYCS--LACTHNRH-FNPF
HG --P--P--SVTYEPM-AYMDAAVY--GSEISGTPQNLVLFVDTGSSNLWVPSVYCS--QACTSHSR-FNPS
MG --P--P--SVTYEPM-AYMDAAVY--GSEISGTPQNLVLFVDTGSSNLWVPSVYCS--QACTSHSR-FNPS
KC --P--P--GVAASVYPLNYLDSQYV--GKILYLGTFPQEFVLFVDTGSSDFV-PSVYCS--NACKNHRQ-FDPR
HR --P--P--LIGNTISSVILNMYDIOVY--GEIGCTPPQTFKVVFDTGSSVWVPSKSRRLVACGVHKL-PDAS
OR --P--P--PSSLNLSVYVLLNYDIOVY--GEIGCTPPQTFKVVFDTGSSVWVPSKSRRLVACGVHKL-PDAS
SR --P--P--PSSLNLSVYVLLNYDIOVY--GEIGCTPPQTFKVVFDTGSSVWVPSKSRRLVACGVHKL-PDAS
FD --P--P--*GPIPEVLRKMYDAQYV--GEIGCTPPQTFVVFDTGSSNLWVPSHCKLLDIAQVIRHK-YNSG
HD --P--P--*GPIPEVLRKMYDAQYV--GEIGCTPPQTFVVFDTGSSNLWVPSHCKLLDIAQVIRHK-YNSG
RZ IVPD--*AGVGVYVPHDYG-NDEYVGVYVIG--TPGKKNLDFDTGSSDLWVSTLCT--NCGSKRKYVDPK
PW *AASGVYATIFIAN-DEYVYVYVIG--GT--LWLRFDGADLVVFTSLPA-SQSGSRYVPS
EN *STGSAITTFDSDLDAVYVYVYVIG--TPAQTLNDFDTGSSDLWVSTLCT--NCGSKRKYVDPK
PA VFSREHPPFT*GCHDVPVLTNYLNAQYV--DITLGTTPQNFKVLIDTGSSNLWVPSNEGCS--LACFLHRSK-YDHE

Fig. 1

quence. Figure 1 illustrates that different enzymes from stomach, lysosomes, kidney, and fungi are of similar molecular size (about 330 residues), and they share, in addition to identical residues, important structural features such as positions of disulfide pairs and location of active sites. Of the sequences listed in Figure 1, four aspartic protease crystal structures have been solved at high resolution. These are porcine pepsin [1-4], rhizopuspepsin [5,6], penicillopepsin [7-9], and endothiasepsin [5,10,11]. All four crystal structures are closely related in overall shapes and the tracing of chain foldings. (The atomic coordinates of the fungal protease crystal structures are available in the Brookhaven Data Bank.) Also, the structures at the active sites are nearly identical among the four crystal structures. These comparisons, together with the homology in primary structures (Fig. 1), predict that the tertiary structures of all the aspartic proteases are similar and that they are derived from the same ancestral protein by a divergent evolutionary process. Because of the similarity in the primary and tertiary structures, it has been possible to model the tertiary structures of aspartic proteases by fitting the amino acid sequences onto the existing, presumably homologous crystal structures [6,7].

EVOLUTION OF THE ACTIVE CENTER

The catalytic apparatus of aspartic proteases consists primarily of two aspartic acid residues (positions 32 and 215, Fig. 1), which were originally identified by active-site-directed reagents: diazoacetyl norleucine methyl ester [12] and 1,2-epoxy-3-(p-nitrophenoxy) propane [13]. These two aspartyl side chains, which are located in the center of the apparent substrate binding cleft, are within the hydrogen bond distance of each other [4,6,7,11] and are in essentially the same relationship in all four crystals. This high conformational similarity among the four extends also to the polypeptide chains near the active-site aspartyls. These observations support the notion that the catalytic apparatus in all the aspartic proteases is virtually the same, and the differences among these enzymes are due mainly to the differences in specificities resulting from the structural evolution of the sites for substrate side chain bindings. The hypothesis that the aspartic proteases share the same catalytic apparatus is supported also by the fact that they are universally inhibited by the two active-site reagents mentioned above and by pepstatin, a transition-state analogue inhibitor [14,15].

In contrast to the good agreement in the structure of catalytic apparatus, little agreement exists in the translation of this apparatus into a catalytic mechanism. It appears clear now that no stable covalent acyl or amide intermediate is present during catalysis [16,17] and also that, owing to conformational reasons, neither of the active-site carboxyls can serve directly as a nucleophile in the peptide bond hydrolysis.

Fig. 1. Homology alignment in the available amino acid sequences of aspartic proteases and zymogens. The residue numbering is based on porcine pepsinogen and pepsin. The activation peptide residue numbers carry a prefix P. The numbering for the enzyme part starts at the N-terminus of porcine pepsin. The sequences are grouped according to their sources. The first group of seven sequences are gastric enzymes. The second group of three sequences are renins. The third group of two sequences are fungal enzymes. The alignment positions are based on both the maximum homology and the structurally or functionally important residues (such as KY at P36-P37, see text). The N-terminal positions of active proteases are marked by solid diamonds. References are in brackets.

Thus, it is likely that a concerted mechanism exists that produces, at the transition state, two tetrahedral atoms at the carboxyl and amide positions of the scissile bond. This is supported by the structure of the transition state inhibitor statine [15,18,19] as well as the ^{13}C nuclear magnetic resonance data of statine derivatives [20].

The precise catalytic roles of the aspartyls are unclear. Since a water molecule is found equidistant from these carboxyls [9,21], it may have an enhanced nucleophilicity in the attack on the carbon atom of the carbonyl group of the scissile peptide bond [22]. Alternatively, the aspartyl groups may serve to polarize the substrate carboxyl [23] in a manner similar to the "oxyanion hole" of the serine proteases.

The amino acid sequences around the two active-site aspartyl residues are highly similar. In addition, there is a near symmetry of the conformation surrounding those residues. This twofold symmetry in conformation can be extended to the entire molecule in each of the four available aspartic protease crystal structures. This is the result of the separate but nearly identical folding of the two halves of the aspartic protease molecules [24], as illustrated in Figure 2. The substrate binding clefts are located between the N- and C-terminal lobes and extend nearly the entire width of the enzyme molecules. Within each of the lobes, another twofold symmetry in polypeptide foldings can be seen [10,25], although these similarities are less obvious. Therefore, the genes for the aspartic proteases must have originated from an ancestral gene of about one-fourth the size of the present gene. The evolution of the ancestral gene to the genes from the aspartic proteases must have involved at least two separate gene duplications and fusions [24]. From this evolutionary scheme, it can perhaps be speculated that the emergence of a primordial aspartic protease took place after the first gene duplication and fusion to produce an enzyme with two identical subunits, each contributing an active-site aspartic acid.

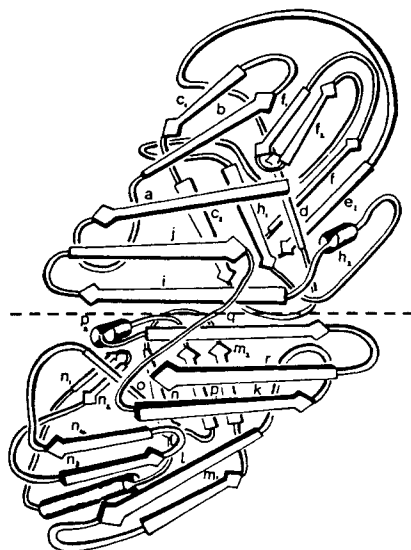


Fig. 2. Schematic presentation of two-fold symmetry in the conformation of aspartic proteases. The view is from the opposite side of the substrate-binding cleft. The symmetry is on the two sides of the dotted line.

The substrate binding clefts of the aspartic proteases are located essentially between the N- and C-terminal lobes and run across the entire width of the molecules. The cleft can accommodate the side chains of eight residues in a polypeptide substrate, equally divided on both sides of the catalytic aspartyls. The binding positions of major substrate side chains, those near the catalytic apparatus, have been identified in some of the crystal structures [6,23]. There is a lack of agreement in the assigned substrate binding residues among the aspartic proteases. This is not surprising since the specificities of the aspartic protease are drastically different. They range from a highly stringent structural requirement for renin to a fit-all design of pepsin and the fungal enzymes.

The discussions above suggest that during evolution the catalytic apparatus in the aspartic proteases was essentially unaltered, and the functional selection of the specificities was accomplished by changes of the substrate binding cleft. In this respect, the evolution of aspartic proteases is similar to that of other proteases, such as the difference between the active center of trypsin and chymotrypsin in serine protease. However, unlike the case of serine proteases, for which the increase of specificity is often accomplished by the addition of protein components to the enzyme (see Neurath, this UCLA Symposia volume), the enhanced specificity in aspartic proteases, such as renin, is achieved by the restricted steric requirements of the extended binding cleft.

GASTRIC PROTEASES

The best studied gastric aspartic proteases are those from the stomachs of high mammals. There are three different gastric aspartic proteases. Pepsin (or pepsin A, EC3.4.23.1) from several species has been sequenced (see Fig. 1), but the pig enzyme is the only gastric protease for which a high resolution crystal structure exists [4]. Gastricsin (or pepsin C, EC3.4.23.3) differs from pepsin in pH optimum, in some specificity, and considerably in primary structure (Fig. 1). Gastricsin is present in significant quantity in the stomachs of man [26,27] and monkey [28]. It is apparently the major protease in rat stomach [29]. Chymosin (EC3.4.23.4), which is present in the stomach of newborn ruminants, has low proteolytic activity but can effectively partially hydrolyze casein and clot milk.

Recent sequence determinations of gastricsin from monkey [30] and man [31] enable comparisons of the structural relatedness of all three gastric proteases (see Fig. 1). It is interesting that the sequences of pepsin from humans [32] and monkeys [33] are closer to the bovine chymosin sequence [34,35] than the gastricsin sequence from the same species. The closer structural relationships of pepsin and chymosin are achieved in spite of the cross-species comparison and specificity differences. These observations suggest that the divergence of pepsin and chymosin is the most recent event in the evolution of three gastric aspartic proteases.

An interesting structural relationship may exist between stomach gastricsin and the aspartic protease of seminal plasma. The seminal enzyme cross-reacts with antibody against gastricsin [36,37], and the two enzymes are close in amino acid compositions [38]. The seminal protease is synthesized in the prostate and secreted as a zymogen [39]. Although the physiological function of this protease is still not known, existing evidence seems to point to a close structural relationship with the gastric enzymes.

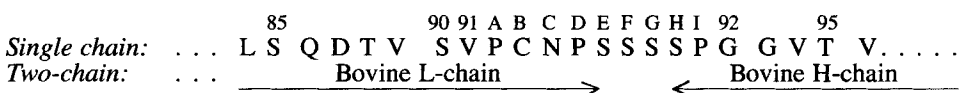
Another aspect of interest is the close relationship of the structures of gastric and fungal aspartic proteases. Considering that these are all "general purpose" proteases secreted to function outside of cells, the strong structural conservation is not surprising.

MAMMALIAN LYSOSOMAL CATHEPSIN D

Cathepsin D is a major endopeptidase present in the lysosomes. The physiological function of this enzyme is to degrade proteins. Cathepsin D has a broad peptide bond specificity similar to pepsin and gastricsin. However, unlike the gastric proteases, cathepsin D functions inside the cells. Thus it is interesting to compare the structure-function of cathepsin D with other aspartic proteases.

The amino acid sequences of porcine cathepsin D [40] and human cathepsin D [41] are very homologous to other gastric and fungal aspartic proteases. In the alignment of Figure 1, cathepsin D shares 49% identical residues with mouse submaxillary renin, 46% with human renin, 48% with porcine pepsin, and 26% with penicillopepsin. This similarity includes the most conserved regions around the active-site aspartyl residues. Thus, it can be predicted that the three-dimensional structure of cathepsin D is very closely related to the crystal structures of other aspartic proteases discussed above.

Two aspects of cathepsin D structure are different from other aspartic proteases. First, porcine cathepsin D is glycosylated at Asn residues at positions 67 and 183 (Fig. 1). Five oligosaccharide units were found at position 67, and three oligosaccharides were found at position 183 [42]. Most of these carbohydrate structures are variants of N-linked high-mannose oligosaccharides, which are known to contain the targeting signal, mannose-6-phosphate, for the lysosomal hydrolases [43]. Like other lysosomal enzymes, cathepsin D appears to contain a site on its polypeptide structure that is recognized for the enzymic phosphorylation of its mannose by a phosphotransferase (N-acetylglucosaminylphosphotransferase) [44]. The function of this site seems to depend on the conformation of the native enzyme; thus the recognition by phosphotransferase diminishes when the light and heavy chains of the porcine cathepsin D are separated [41]. For this reason, it has been difficult to identify the location of the site. The second major structural difference in cathepsin D is the sequence involved in the processing of single chain to two-chain enzyme *in vivo*. The structures of the proteolytically processed sites in cathepsin D are of considerable interest since it is known that all the lysosomal hydrolases are proteolytically processed after biosynthesis [45-47]. The predominant form of the porcine cathepsin D is a two-chain enzyme that has been cleaved between residues 97 and 98 (Fig. 1). Bovine cathepsin D, however, exists in equal amounts of single and two-chain species [48]. The amino acid sequence at the processing site in these enzymes have been determined recently as follows [49]:



These results indicate that the processing of bovine single-chain cathepsin D takes

place at two Ser-Ser bonds between residues 91E-91F and 91G-91H. Another interesting aspect is the conformation of the processing region. Alignment of the bovine sequence above against the corresponding region of porcine pepsin shows that nine residues from 91A to 91I are insertions in cathepsin D. The insertion position in pepsin (residues 91-92) is located at the end of an antiparallel β -hairpin structure, which slightly protrudes from the surface of pepsin crystal structure. Thus, the additional region of nine inserted residues probably loops outside the cathepsin D molecular surface. Alignment of the human cathepsin D sequence with the pepsin sequence indicates the presence of a similar insertion at the same region (Fig. 1). This unique conformation should be readily hydrolyzed by another protease of appropriate specificity. However, it is not certain at the present whether the purpose of this conformation serves as recognition for the processing of cathepsin D. To summarize, cathepsin D is structurally very similar to other "general purpose" aspartic proteases. The special structural features appear to function in lysosome targeting and proteolytic processing.

It is appropriate to compare the structure and function relationships of cathepsin D and yeast proteinase A. Both enzymes occur in granules and hydrolyze proteins intracellularly. In amino acid sequence, cathepsin D is more closely related to proteinase A than other mammalian aspartic proteases [50] (see also Fig. 1). Like cathepsin D, proteinase A contains two Asn-linked oligosaccharide units [51,52]. It is interesting to note that for the two potential glycosylation sites seen in the proteinase sequence, one is identical with that in cathepsin D (residue 68, Fig. 1) The other glycosylation site in proteinase A (residue 267, Fig. 1) is quite far from that of cathepsin D (residue 183, Fig. 1). However, based on the crystal structures of aspartic proteases, these two glycosylation sites are located on two adjacent β -strands quite near each other. Although another yeast granule enzyme, carboxypeptidase Y, is known to be phosphorylated on the oligosaccharides [53], it is not known whether this is the case for proteinase A. The above discussion nevertheless suggests a close relationship in the structure and function of these two enzymes.

ACTIVATION OF ASPARTIC PROTEASE ZYMOGENS

It is likely that all aspartic proteases are biosynthesized as larger proenzymes, which are subsequently converted to mature enzymes. In addition to the gastric zymogens, such as pepsinogen, the existence of the zymogens has been established for renin [54], cathepsin D [41,46], seminal plasma acid protease [38], yeast proteinase A [50], and fungal rhizopuspepsin [55]. The amino acid sequences of most of these zymogens are shown in Figure 1. The conversion of zymogens to enzymes involves the removal of the N-terminal region, which is about 45 residues. The conformation of the activation peptide of porcine pepsinogen is known from the solution of the crystal structure of this protein [56]. As shown in Figure 1, the sequences of activation peptides are of similar length and are apparently homologous to one another. Thus, the conformations of the activation peptides in aspartic protease zymogens are likely quite homologous.

There appear to be two ways the aspartic protease zymogens can be activated under the physiological conditions. Most of these zymogens are activated upon acidification. This route of activation is known for pepsinogen [57], other gastric zymogens [58,59], seminal protease zymogen [38], and procathepsin D [60]. The

mechanism of acid activation has been well studied for porcine pepsinogen and documented in reviews [61,62]. Porcine pepsinogen activates by an intramolecular mechanism under the conditions (pH 1–2) similar to those expected in the stomach [63–65]. The intramolecular mechanism involves first the “local denaturation” of activation peptide to reveal the pepsin active site [66,67]. The activation peptide is then bound to the substrate cleft, which results in the cleavage of Leu-Ile bond at residues 16–17 and the release of the N-terminal peptide [68]. The remainder of the activation peptide (residues 17–44) is then removed, probably by another pepsin molecule. The intramolecular activation mechanism of porcine pepsinogen has been studied kinetically in some detail [69,70]. In the activation of pepsinogen and progastricsin from several species of mammals, the first cleavage sites in the activation peptides appear to differ [71]. This is probably related to the sequence changes in the activation peptides and their suitability as substrates of their respective active sites in the intramolecular hydrolysis. When the activation peptide sequence is not suited for the enzyme specificity, the activation peptide may be removed in its entirety, as in the case of monkey pepsinogen [72]. Whether this is an intramolecular process is not yet known. Judging from the pepsinogen crystal structure [56], the N-terminal 12 residues of pepsin molecule must also be displaced with the activation peptide during the “local denaturation” step at the beginning of the acid activation. Thus, there seems no conformational reason why the activation peptide can not be removed en bloc during the intramolecular catalysis. Alternatively, a pepsin-catalyzed pepsinogen activation in acid has been demonstrated [65] and can serve to remove the entire activation peptide. This bimolecular mechanism, which is predominant at high zymogen concentration and pH range of 2.5 to 3.5 [65], is probably a minor route under physiological conditions.

Procathepsin D is known to activate in acid [60]; thus it may have an intramolecular activation mechanism similar to that of pepsinogen. The pulse-chase experiments suggested that this activation takes place in the acidic environment inside of the lysosomes. Although the mode of activation is not known for rhizopuspepsinogen, it may also be activated in the secretory granule by an intramolecular mechanism.

The activation mechanism for the aspartic protease zymogen prorenin is unique. Since prorenin cannot be activated by acidification and it contains at the end of activation peptide a Lys-Arg sequence that is characteristic for prohormone activation, the involvement of a second protease that recognizes two basic residues is assumed. Perhaps owing to this difference in mode of activation, the activation peptides of prorenins contain fewer basic residues (from 6 to 8) than those in other aspartic protease zymogens (usually from 12 to 14). Many of these basic residues are proposed to take part in charge interaction in the native porcine pepsinogen molecule and thus contribute to the “local denaturation” of the activation peptide upon acidification [56]. Notably absent in the renin group is residue Lys-P36, which is involved in the interaction with active-site aspartyls [56] and is otherwise conserved in other zymogen sequences (see Fig. 1).

Comparing the two activation systems, the intramolecular mechanism is probably selected in evolution for its speed and for its cell economy of not requiring another activating enzyme. It seems an advantageous mechanism for digestive proteases for which regulation is not essential. Another conceivable advantage for the intramolecular mechanism is that the resulting aspartic protease can initiate the activation of other proteases in a cascade. Supporting evidence exists for the possible function of

proteinase A in the activation of other proteases in the intracellular granules [73]. When regulation is obviously required, as in the case of prorenin, the activation by another protease then becomes necessary.

CONCLUDING REMARKS

The aspartic proteases discussed above are all derived from the same evolutionary origin. Although their functions are diverse, their amino acid sequences appear to be highly conserved. One of the reasons perhaps is the large substrate binding cleft, which allows the broad specificity alteration without having to provide additional protein domains, as in the case of serine proteases evolution. However, some aspartic proteases are reported to have larger molecular weights, such as cathepsin E [74] and a pituitary enzyme that is involved in the prohormone processing [75]. It remains to be seen whether additional domains exist in these enzymes.

Although the aspartic proteases are apparently derived in evolution from gene duplications and fusions, no clear primordial enzyme has so far been observed. A possible candidate for this is the pepstatin-insensitive acid protease B₁ from *Scytalidium lignicolum* [76]. This enzyme contains 204 residues and appears to have an active-site aspartic acid that can react with the epoxide inhibitor. There is no apparent sequence homology between this enzyme and other aspartic proteases. A distant relationship of this enzyme to the others might be revealed by determination of its X-ray structure.

ACKNOWLEDGMENTS

The authors wish to thank Dr. Tom Stevens for making his manuscript available to us before publication and to Dr. Jean Hartsuck for help in manuscript preparation. This work was supported by NIH research grants AM01107 and GM35424.

REFERENCES

1. Andreeva NS, Gustchina AE, Federov AA, Volnova TV, Shutzkerver NE: *Adv Exp Med Biol* 95:23-31, 1977.
2. Andreeva NS, Federov AA, Gustchina AE, Risculov RR, Sufro MG, Shutzkever NE: *Mol Biol (Mosk)* 12:704-717, 1978.
3. Andreeva NS, Gustchina AE, Federov AA, Shutzkever NE, Volnova TV: In Tang J (ed): "Acid Proteases, Structure, Function and Biology," New York: Plenum Press, 1977, pp 23-31.
4. Andreeva NS, Zdanov AS, Gustchina AE, Federov AA: *J Biol Chem* 259:11353-11365, 1984.
5. Subramanian E, Swan IDA, Liu M, Davies DR, Jenkins JA, Tickle IJ, Blundell TL: *Proc Natl Acad Sci USA* 74:556-559, 1977.
6. Bott R, Subramanian E, Davies DR: *Biochemistry* 21:6956-6962, 1982.
7. Hsu I-N, Delbaere LTJ, James MNG, Hofmann T: *Nature* 266:140-145, 1977.
8. Hsu: In Tang J (ed): "Acid Proteases, Structure, Function and Biology," New York: Plenum Press, 1977, pp 61-81.
9. James MNG, Sielecki AR: *J Mol Biol* 163:299-361, 1983.
10. Blundell TL, Sewell BT, McLachlin AD: *Biochim Biophys Acta* 580:24-31, 1979.
11. Blundell T, Jenkins J, Pearl L, Sewell T: In Kostka V (ed): "Aspartic Proteinases and Their Inhibitors," Berlin: Walter de Gruyter, 1985, pp 151-161.
12. Ragagopalan TG, Stein WH, Moore S: *J Biol Chem* 241:4295-4297, 1966.
13. Chen KCS, Tang J: *J Biol Chem* 247:2566, 1972.
14. Umezawa H, Aoyagi T, Morishima H, Matzusaki M, Hamada H, Tekeuchi T: *J Antibiot (Tokyo)* 23:259, 1970.

62:JCB Tang and Wong

15. Marciniszyn J, Jr, Hartsuck JA, Tang J: *J Biol Chem* 251:7088–7094, 1976.
16. Dunn BM, Fink AL: *Biochemistry* 23:5241–5247, 1984.
17. Hofmann T, Fink AL: *Biochemistry* 23:5247–5256, 1984.
18. Marshall GR: *Fed Proc Am Soc Exp Biol* 35:2494, 1976.
19. Rich DH, Bernatowicz MS, Agarwal NS, Kawai M, Salituro FG, Schmidt PG: *Biochemistry* 24:3165–3173, 1985.
20. Rich DH, Bernatowicz MS, Schmidt PG: *J Am Chem Soc* 104:3535–3536, 1982.
21. Pearl LH, Blundell TL: *FEBS Lett* 174:96–101, 1984.
22. Pearl LH: In Kostka V (ed): “Aspartic Proteinases and Their Inhibitors,” Berlin: Walter de Gruyter, 1985, pp 189–195.
23. James MNG, Sielecki AR, Hofmann I: In Kostka V (ed): “Aspartic Proteinases and Their Inhibitors,” Berlin: Walter de Gruyter, 1985, pp 163–177.
24. Tang J, James M, Hsu IN, Jenkins JA, Blundell TL: *Nature* 271:618, 1977.
25. Andreeva NS, Gustchina AE: *Biochem Biophys Res Commun* 87:32–42, 1979.
26. Tang J, Wolf S, Caputto R, Trucco RE: *J Biol Chem* 234:1174, 1959.
27. Chiang L, de Chiang L, Wolf S, Tang J: *Proc Soc Exp Biol Med* 122:700, 1966.
28. Kageyama T, Takahashi K: *J Biochem (Tokyo)* 80:983–992, 1976.
29. Arai KM, Muto N, Tani S, Akahane K: *Biochim Biophys Acta* 788:256–261, 1984.
30. Kageyama T, Takahashi K: *J Biol Chem* 261:4406–4419, 1986.
31. Wong RNS, Tang J: *Fed Proc* 45:105, 1986.
32. Sogawa K, Fujii-Kuriyama Y, Mizukami Y, Ichihara Y, Takahashi K: *J Biol Chem* 258:5306–5311, 1983.
33. Kageyama T, Takahashi K: *J Biol Chem* 261:4395–4405, 1986.
34. Foltmann B, Pedersen VB, Kaufman D, Wybrandt G: *J Biol Chem* 254:8447–8456, 1979.
35. Harris TJR, Lowe PA, Lyons A, Thomas PG, Eaton MAW, Millican TA, Patel TP, Bose CC, Carey NH, Doel MT: *Nucleic Acids Res* 10:2177–2187, 1982.
36. Hirsch-Marie H, Conte M: *Bull Soc Chim Biol* 49:147–155, 1967.
37. Samloff IM, Liebman WM: *Clin Exp Immunol* 11:405–414, 1972.
38. Ruenwongsa P, Chulavatnatol M: *J Biol Chem* 250:7574–7578, 1975.
39. Chiang L, Contreras L, Chiang J, Ward PH: *Arch Biochem Biophys* 210:14–20, 1981.
40. Shewale JG, Tank J: *Proc Natl Acad Sci USA* 80:3703–3707, 1984.
41. Faust PL, Kornfeld S, Chirgwin JM: *Proc Natl Acad Sci USA* 82:4910–4914, 1985.
42. Takahashi T, Tang J: *J Biol Chem* 258:6435–6443, 1983.
43. Sly WS, Fischer HD: *J Cell Biochem* 18:67–85, 1982.
44. Lang L, Reitman M, Tang J, Roberts RM, Kornfeld S: *J Biol Chem* 259:14663–14671, 1984.
45. Hasilik A, Neufeld EF: *J Biol Chem* 255:4946–4950, 1980.
46. Erickson AH, Conner GE, Blobel G: *J Biol Chem* 256:11224–11231, 1981.
47. Waheed A, Hasilik A, von Figura K: *Eur J Biochem* 123:317–321, 1982.
48. Huang JS, Huang SS, Tang J: In Mildner P, Ries B (eds): “Enzyme Regulation and Mechanism of Action,” Oxford: Pergamon Press, 1980, pp 289–306.
49. Yonezawa S, Takahashi T, Dehdarani MA, Tang J: manuscript in preparation.
50. Ammerer G, Hunter CP, Rothman JH, Saari GC, Valls LA, Stevens TH: *Mol Cell Biol* 7:2490–2499, 1986.
51. Mechler B, Muller M, Muller H, Meussdoerffer F, Wolf DH: *J Biol Chem* 257:11203–11206, 1982.
52. Meussdoerffer F, Tortora P, Holzer H: *J Biol Chem* 255:12087–12093, 1980.
53. Hashimoto C, Cohen RF, Zhang W-J, Ballou CE: *Proc Natl Acad Sci USA* 78:2244–2248, 1981.
54. Panthier JJ, Foote S, Chambraud B, Strosberg AD, Corvol P, Rougeon F: *Nature* 298:90–92, 1982.
55. Wong NS, Delaney R, Tang J: *Fed Proc* 44:2725, 1985, and *J. Biol. Chem.* in press, 1987.
56. James MNG, Sielecki AR: *Nature* 319:33–38, 1986.
57. Herriott RM: *J Gen Physiol* 21:501–540, 1938.
58. Foltmann B, Jensen AL: *Eur J Biochem* 128:63–70, 1982.
59. Foltmann B: *Methods Enzymol* 19:421–435, 1970.
60. Hasilik A, von Figura K, Conzelmann E, Nehr Korn H, Sandhoff K: *Eur J Biochem* 125:317–321, 1982.
61. Hartsuck JA, Marciniszyn J, Jr, Huang J-S, Tang J: In Tang J (ed): “Acid Proteases, Structure, Function and Biology,” Plenum Press, 1977, pp 85–102.
62. Hartsuck JA, Tang J: In Magnusson (ed): “Regulatory Proteolytic Enzymes,” Pergamon Press, 1978, p 35.

144:PBCB

63. Bustin M, Conway-Jacobs A: *J Biol Chem* 246:615-620, 1971.
64. McPhie P: *J Biol Chem* 247:4277-4281, 1972.
65. Al-Janabi J, Hartsuck JA, Tang J: *J Biol Chem* 247:4628-4632, 1972.
66. Marciniszyn J, Jr, Huang J-S, Hartsuck JA, Tang J: *J Biol Chem* 251:7095-7102, 1976.
67. Glick DM, Auer HE, Rich DH, Kawai M, Kamath A: *Biochemistry* 25:1858-1864, 1986.
68. Dykes CW, Kay J: *Biochem J* 153:141-144, 1976.
69. Twining SS, Sealy RC, Glick DM: *Biochemistry* 20:2375-2379, 1981.
70. Auer HE, Glick DM: *Biochemistry* 23:2735-2739, 1984
71. Takahashi K, Kageyama T: In Kostka V (ed): "Aspartic Proteinases and Their Inhibitors," Berlin: Walter deGruyter, 1985, pp 265-282.
72. Kageyama T, Takahashi K: *J Biochem (Tokyo)* 88:9-16, 1980.
73. Zubenko GS, Park FJ, Jones EW: *Proc Natl Acad Sci USA* 80:510-514, 1983.
74. Lapresle C, Webb T: *Biochem J* 84:455-462, 1962.
75. Loh YP, Parish DC, Tuteja R: *J Biol Chem* 260:7194-7205, 1985.
76. Maita T, Nagata S, Matsuda G, Maruta S, Oda K, Murao S, Tsuru D: *J Biochem (Tokyo)* 95:465-475, 1984.
77. Sepulveda P, Marciniszyn J, Jr, Liu D, Tang J: *J Biol Chem* 250:5082-5088, 1975.
78. Ong EB, Perlmann GE: *J Biol Chem* 243:6104-6109, 1968.
79. Baudys M, Kostka V: *Eur J Biochem* 136:89-99, 1983.
80. Imai T, Miyazaki H, Hirose S, Hori H, Hayashi T, Kageyama R, Ohkubo H, Nakanishi S, Murakami K: *Proc Natl Acad Sci USA* 80:7405-7409, 1983.
81. Holm I, Ollo R, Panthier J-J, Rougeon F: *EMBO J* 3:557-562, 1984.
82. Misono KS, Chang J-J, Inagami T: *Proc Natl Acad Sci USA* 79:4858-4862, 1982.
83. Pedersen VB, Foltmann B: (unpublished results).
84. Dreyer T, Halkier B, Svendsen I, Ottesen M: *Carlsberg Res Commun* 51:27-41, 1986.