



ACADEMIC
PRESS

Biochemical and Biophysical Research Communications 293 (2002) 23–29

BBRC

www.academicpress.com

Evolution of placentally expressed cathepsins[☆]

Katia Sol-Church,^a Gina N. Picerno,^a Deborah L. Stabley,^a Jennifer Frenck,^a Sixun Xing,^a
Greg P. Bertenshaw,^b and Robert W. Mason^{a,*}

^a Laboratory of Clinical Biochemistry, Alfred I duPont Hospital for Children, P.O. Box 269, Wilmington, DE 19899, USA

^b Department of Biochemistry and Molecular Biology, The Pennsylvania State University College of Medicine,
500 University Drive, Hershey, PA 17033, USA

Received 4 March 2002

Abstract

Species and strain variants of a family of placentally expressed cathepsins (PECs) were cloned and sequenced in order to identify evolutionary conserved structural characteristics of this large family of cysteine proteases. Cathepsins M, P, Q, and R, are conserved in mice and rats but homologs of these genes are not found in human or rabbit placenta, showing that this family of proteases are probably restricted to rodents. Species-specific gene duplications have given rise to variants of cathepsin M in mice, and cathepsin Q in rats. Although the PECs have diverged at a greater rate than the other lysosomal cathepsins, residues around the specificity sub-sites of the individual enzymes are conserved. Strain-specific polymorphisms show that the evolutionary rate of divergence of cathepsins M and 3, the most recently duplicated pair of mouse genes, is even higher than the other PECs. In human placenta, critical functions of the PECs are probably performed by broader specificity proteases such as cathepsins B and L. © 2002 Elsevier Science (USA). All rights reserved.

Inhibitor studies have implicated the mammalian cysteine proteases in placental function [1–4]. Activity and mRNA for cathepsins B and L can be demonstrated in placenta throughout gestation [5–7]. These observations led to the suggestion that cathepsins B and L are necessary for normal embryo development. However, cathepsin B and L-deficient mice are viable and fertile and appear normal at birth, indicating that there must be alternative targets for the protease inhibitors [8,9]. The first clue to indicate that the placenta contains novel cathepsins came from inhibitor labeling of cellular proteases [6]. A novel polypeptide, termed p14, was discovered in ectoplacental cone and placenta that specifically labeled with the cysteine protease inhibitor,

Fmoc-Tyr [¹²⁵I]–Ala–CHN₂. Mouse and rat placentae were subsequently screened by RT-PCR for novel cathepsin L-like proteases that might be a target of this inhibitor. Using this approach, we identified four novel genes, cathepsins M, P, Q, and R, expressed only in the placenta [10–13]. Four additional related placentally expressed cathepsins, cathepsins 1, 2, 3, and 6, have since been reported in mice [14–16]. The mouse genes are located in a tight cluster on chromosome 13 [15].

In this report, we show that similar enzymes are expressed in both rat and mouse placentae but are not expressed in rabbit or human placentae. Although the sequences of the PECs are diverging rapidly both within and between species, structural analysis indicates there is functional conservation of the enzymes. We have collectively called this family of proteases PECs, for placentally expressed cathepsins.

Materials and methods

Detection of placentally expressed cathepsins. Rat (*Rattus Norvegicus*) and mouse (Swiss Webster) placentae were kind gifts from Joan Pugarelli (duPont Hospital for Children) and Brian Malester (Nathan

[☆] Abbreviations: Fmoc, 9-fluorenylmethyloxycarbonyl; RT-PCR, reverse transcription-polymerase chain reaction; MCTSL, mouse cathepsin L; RCTSL, rat cathepsin L; MCTSM, mouse cathepsin M; RCTSM, rat cathepsin M; MCTSP, mouse cathepsin P; RCTSP, rat cathepsin P; MCTSQ, mouse cathepsin Q; RCTSQ, rat cathepsin Q; RCTSQ2, rat cathepsin Q2; MCTSR, mouse cathepsin R; RCTSR, rat cathepsin R.

* Corresponding author. Fax: +1-302-651-6767.

E-mail address: rmason@nemours.org (R.W. Mason).

Kline Institute). RNA was extracted from the placenta using a S.N.A.P. Total RNA Isolation kit from Invitrogen (Carlsbad, CA). RT-PCR were performed as previously described [13]. All primers used in this study were purchased from Integrated DNA Technologies (Coralville, IA). Specific primers were designed from the previously published sequences of mouse cathepsins M, P, and R to screen for their homologs in rat placenta. The mouse counterpart of rat cathepsin Q was detected by screening mice placenta with rat cathepsin Q-specific primers. The PCR products were ligated in pGEM-T easy (Promega, Madison, WI) and individual clones were sequenced using BigDye terminator cycle sequencing (PE, Foster City, CA). From these sequences, new species-specific primers were designed to complete the sequencing of the PEC homologs using a RACE kit from Gibco-BRL (Rockville, MD). Sequences obtained were analyzed using the BLAST program of the NCBI, and novel genes deposited in the GenBank database.

Screening of human and rabbit placenta. Rabbit (New Zealand White) placenta were purchased from Pel Freeze (Rogers, AR). Rabbit RNA extraction was performed as above. Human placenta RNA was purchased from Ambion (Austin, TX). First strand cDNA synthesis was performed as previously described [13]. Degenerate primers were designed based upon the conserved regions around the cysteine and asparagine active site residues of mouse cathepsins L, P, Q, M, and R. F-CYS25 (5'-TCTTGTTGGTCTTTY-3') and R-ASN (5'-KYACCCMWGWSWGTTYT-3') are expected to amplify a region approximately 510 bp in length. An additional primer, F-PEC (5'-GCTGTRGCAWSTAHWGG-3'), designed from a region that is conserved among the PEC genes, was used with R-ASN to amplify a region 200 bp in length. The PCR products were cloned into pGEM-T easy as above. Clones were picked at random and screened by PCR using vector-specific primers flanking the inserts (T7 and SP6). All candidates containing an insert of the expected size range were sequenced using T7 and SP6 using cycle sequencing.

Sequence analysis tools. Sequences generated during this study were assembled and analyzed using the MacVector Sequence Analysis package. The mouse and rat PECs homologs were aligned using the ClustalW program. Phylogenetic analysis of the mature portions of the proteins was performed using the PHYLIP programs available on the server of the Pasteur Institute (<http://bioweb.pasteur.fr/seqanal/phylogeny/phytip-uk.html>). A distance matrix for protein sequences was computed using the ProtDist program and evolutionary trees calculated using the Neighbor Joining and UPGMA methods.

Structural modeling. The sequences of the PECs protein were submitted to SWISS-Model for homology modeling in the optimized mode. Modeling was performed using the coordinates deposited for the crystal structures of human procathepsin L (1CJL.pdb) and human procathepsin K (1BY8.pdb) as templates. All images were made using WebLab ViewerLite 4.0 (Molecular Simulations, San Diego, CA).

Results and discussion

Screening of mammalian placenta for PECs

In this study, we have cloned and sequenced mouse cathepsin Q and rat cathepsins M, P, and R (AF456461, AF456462, AF456458, AF456459, respectively). The deduced amino acid sequences of these rodent PECs were aligned with published sequences of other PECs using mouse cathepsin L as a reference (Fig. 1). The novel rat genes are 57–64% identical to each other at the protein level, and are 77–84% identical to their mouse homologs. Mouse cathepsin Q is 86% and 82% identical

to the cDNA and deduced protein sequences of rat cathepsin Q, respectively [10]. We have also discovered a variant of cathepsin Q in the rat genome that we have called rat cathepsin Q2 (AF456460). This novel gene is 88% and 86% identical at the nucleotide level to rat and mouse cathepsin Q, respectively, and probably arose from duplication of the rat ancestral cathepsin Q gene.

For human placental tissue, RT-PCR using degenerate primers did not reveal the presence of sequences related to the PECs although cDNAs from related enzymes cathepsins L and F were detected. No PECs could be identified by screening of the human genome database at the NCBI, leading us to conclude that there are no PEC homologs in the human genome. The human genome does contain some pseudogenes related to cathepsins L that have premature stop codons, but these are not closely related to the PECs [17]. The human genome also contains cathepsin V, which is primarily expressed in thymus and testis [18]. This protein is closely related to cathepsin L but is not a human homolog of the PECs.

RT-PCR of rabbit placental cDNA using degenerate primers only yielded cathepsin clones that are closely related to cathepsin L (AF458078). We conclude that rabbit probably also does not express any PEC-like genes, and that PECs have arisen after divergence of rodents and lagomorphs.

Analysis of primary protein sequence

Each of the PEC sequences contains a putative signal peptide, pro-region, and mature region (Fig. 1). In all of the PECs, the essential catalytic residues (Cys138, His276, Asn300, Gln132) are conserved, indicating evolutionary pressure to conserve function. The cysteines that form disulfide bonds to stabilize the structures of cathepsins L, H, and S (Cys 135, 169, 178, 211, 269, 322) are also conserved. Tryptophan is encoded by a single codon but is conserved in all of the PECs at seven sites. Hydrophobic amino acids that have been shown to be embedded in the core of the cathepsin L molecule are conserved as hydrophobic amino acids in the PECs. The conservation of amino acid type is typically required to retain structural integrity of proteins and shows that the PECs are unlikely to be non-functional products.

Despite the overall similarities between the PECs, there are regions of greater sequence divergence at both the DNA and amino acid level. Significant divergence occurs in bases that code for amino acids in the N-terminus of the pro-peptide, close to the site of signal peptide cleavage. The crystal structure of cathepsin L shows that this region of the pro-peptide is unstructured, indicating flexibility and hence fewer structural constraints. Alterations in this region are unlikely to have a major impact on the activity or structural integrity of the PECs, but are probably important for

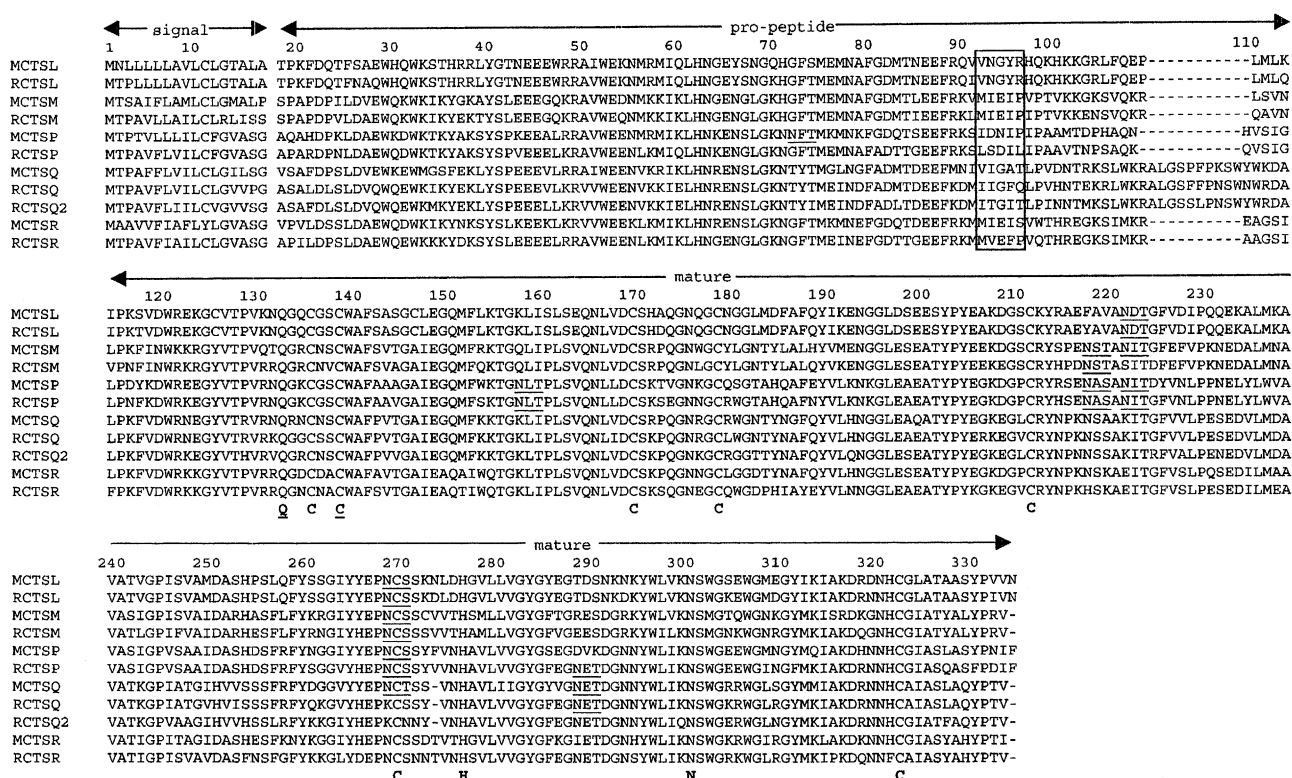


Fig. 1. Amino acid sequences of PECs. Full length amino acid sequences of both rat and mouse cathepsins L, M, P, Q, and R, and rat cathepsin Q2 are aligned. Conserved active site residues and cysteines are shown in bold below the sequences. Potential N-glycosylation sites are underlined. Putative signal, pro-, and mature peptides are marked above the amino acid sequences. Numbering is for cathepsins L, M, and P. Residues of the pro-peptide that align with those shown to bind in the S'₂ to S₃ sub-sites of cathepsin L are boxed.

sorting and signal peptide cleavage. A large variable region spans the inhibitory binding loop and the linker region between the pro-peptide and the mature enzyme. In the cathepsin Q and Q2 sequences there is an insertion of 33 bases to give an additional peptide loop in the linker region. Variability in the linker region might be expected because this region is exposed at the surface of the molecule and its primary role seems to be to link the pro-peptide to the mature protease. Processing within this region results in activation of the enzymes. The N-terminal portion of this variable region binds across the substrate binding cleft in reverse orientation to substrates, effectively inhibiting the mature enzyme (boxed in Fig. 1). The residues of the pro-peptides of cathepsins K, L, and S occupying the substrate specificity pockets show some correspondence to the known substrate specificity of these enzymes, particularly with respect to the residue bound in the S₂ pocket [19]. The sequences of the corresponding peptide in each of the PECs are quite variable, implying that the active PECs probably evolved to bind a different range of peptides. Additional smaller regions of variability at the DNA level have more subtle changes in amino acid sequences. One region of DNA variability codes for amino acids 178–182 in cathepsin L (CNGGL). The cysteine and other glycine are conserved in all of the cathepsins but the

other amino acids are varied in the PECs. In cathepsin L, these amino acids contribute to the S₁ and S₂ binding sites in the substrate binding cleft [20]. Changes in this region are likely to alter the substrate binding properties of the PECs. Another variable region codes for amino acids 251–258 in cathepsin L (ASHPSLQF) and these residues are involved in S'₁ substrate binding pocket for cathepsin L.

Overall, the major differences in the primary amino acid sequences of the PECs are either in regions of protein flexibility or in regions that define enzyme specificity, while essential structural components are conserved. This comparison of primary amino acid sequence with known proteases indicates that the PECs will be active enzymes with unique specificities.

Polymorphisms in PECs

The published mouse PEC sequences originated from the placenta of C57BL/6 mice [12]. We have now cloned and sequenced cathepsins M, P, Q, and R, from Swiss Webster strains of mice. For cathepsins P, Q, and R, only one polymorphism was detected for each gene product and none resulted in a change in the amino acid sequence of translated products. By contrast, for cathepsin M 11 of the 999 bp in the coding sequence

differed between Swiss Webster and C57BL/6 (Fig. 2, AF456463). Six of these differences result in amino acid changes (E54D, R130Q, P157Q, E198D, S272C, A286T; Swiss Webster to C57BL/6), although the changes appear to be conservative.

Several clones giving a novel gene product, related to cathepsin M (76% amino acid identity), were also discovered in Swiss Webster mouse placental cDNA. The sequences of these clones are closely related to a gene product from C57BL/6 mouse placenta (BAB23995) with 16 bp differences (Fig. 2, AF456464). The sequence of the C57BL/6 cDNA has been confirmed (apart from a single base change) and designated cathepsin 3 (AAK58450, 15). The translated gene products predict 11 different amino acids between the proteins from Swiss Webster and C57BL/6 strains of mice (A20S, V129A, K157Q, A183S, L184F, N234S, I245V, G272A, N307Y, L328V, Y329F). These results show that both cathepsins M and 3 are considerably more polymorphic than the

other PECs, implying that these two genes are still evolving rapidly. There are presumably fewer constraints preventing genetic variation of the cathepsin M and 3 genes. The number of polymorphisms is over 10 times greater than would be expected for a neutral rate of evolution, indicating that these genes may be significantly less stable than the genes for the other PECs. Nevertheless, the mutations are clearly not random. For both the cathepsin M and cathepsin 3 variants, most of the amino acid changes are either conservative or code for amino acids on the protein surface and the differences would probably not significantly alter the structure or specificity of the individual enzymes. The translated products have retained essential catalytic components and stop codons or frame-shift deletions or insertions have not occurred. The reason for the highly polymorphic nature of cathepsins M and 3 is not known, but it may be that these new proteases have not yet evolved into their ideal structures.

SWM M T S A I F L A M L C L G M A L P S P A P D P I L D V E W Q K K W K I K Y G K A Y S L E
 SWM ATGACTTCTGCTATCTTCTCGGCCATGCTCTGCTTGGGAATGGCTTTACCATCACCAGCACCTGATCCCATTTTGGATGTTGAATGGCAGAAATGGAAATAAAATATGGAAAAGCATACAGTCTGGAG
 C57MC.....G.....CT.G.G...T...C.C.T...T...C...G.....A.....
 SW3C.....G.....C.....G.....A.....
 SW3 S P A A P D P I L D A E W Q K K W K I K Y G K T Y S L E

 SWM E E G Q K R A V W E E E N M K K I K L H N G E N G L G K H G F T M E M N A F G D M T L E
 SWM GAAGAAGGACAGAAGAGCAGTATGGGAAGAAATATGAAAAGATCAAACCTGCACAACTGGGGAGAATGGCGTGGGGAAGCATGGTTTCCACATGGAAATGAATGCCTTTGGTGACATGACACTTGAA
 C57MC.....T.....
 C573T.....
 SW3T.....
 SW3 E E G Q K R A V W E E N M K K I K L H N G E N G L G K H G F T M E M N A F G D M T L E

 SWM E F R K V M I E I P V P T V K K G K S V Q K R L S V N L P K F I N W K K R G Y V T P V
 SWM GAATTCCAGGAAAGTGATGATTGAATCCAGTCCCACTGTCAAGAAAGGAAAAAGTGTCCAGAAACGTCTGTCTGTTAACCTGCCCAAGTTTATAAATCGAAAAAGAGAGGCTATGTTTACTCTCTGTG
 C57MT.....C.....
 C573A.....
 SW3A.....
 SW3 E F R K E M I E I P V P T V K K G K S V Q K R L S V N L P K F I N W K K R G Y V T P V

 SWM R Q T Q G R C N S C W A F S V T G A I E G Q M F R K T G P Q L I P L S V Q N L V D C S R P
 SWM CCGACACAGGGCAGATTAATCTTGTGGGCTTTTCTGTGACTGGTGCCATAGAGGACAGATGTCCGGAAGAACAGCTGCTAGCTGATCCCTCTGAGTGACAGAACCTTGTGGACTGTTCTAGGCGCT
 C57MA.....
 C573ATTGC.....A.....T.....C/A.....A.....GT.GAC--
 SW3ATTGC.....A.....CAA.....A.....GT.GAC--
 SW3 R T Q I A C N S C W A I S V T G A I E G Q M F R K T G K Q L I P L S V Q N L V D C V D

 SWM G G N W G C Y L L W N T Y L A L H Y V M E N G G L E S E A T Y P Y E E K E D G S C R Y S P
 SWM CAAGGCACTGGGCTGTTATTGGCAATACATCTTGCATTGCACATATGTGATGGAAATGGAGGCTAGAGTCTGAGGCAACCTATCCTTATGAAGAAAAGGAAGGATCGTCAGGTCAGTCCCA
 C57MA.GT---...C...GCA...G.GTG.TAGA.T.T...T...G...T...G...A.....C...C...A.....
 C573A.GT---...C...GCA...G.GTG.TAGA.T.T...A...G...T...G...A.....C...C...A.....
 SW3A.GT---...C...GCA...G.GTG.TAGA.T.T...A...G...T...G...A.....C...C...A.....
 SW3 G S G C H A G S V L D A S L E K Y L M E K G G L E S E A T Y P Y E D K Q G S C R Y N P

 SWM E N S T A N I T G F E F V P K N E D A L M N A V A S I G P I S V A I D A R H A S F L F
 SWM GAAAAATCTACTGCTAATAACAGCCTTGTGAATTTGGTTCCAAAGATGAGGATGCCCATATGAATGCTGTGGCAAGTATAGGCGCCATTTCTGTTGCAATCGATGCAAGCATGCACTTCTCTGTTC
 C57MG...C...T.....A.C...C...T.A...G...C.T.....T.....A.....
 C573G...C...T.....A.C...C...A.T.A...G...C.T.....A.T.....T.....A.....
 SW3G...C...T.....A.C...C...A.T.A...G...C.T.....A.T.....T.....A.....
 SW3 E N S T A S I T G F E F I P N N E V D L M N S A V A S L G P I S V I I D A W H E S F L F

 SWM Y K R G I Y Y E P N C S S S C V - - V T H S M L L V G Y G F A T G R E S D G R K Y W L V K
 SWM TATAAAAGAGGTATATATTACGAGCCAACTGCAGTAGCTGTGTT-----GTCACACATTCTATGTTACTAGTTGGCTATGGCTTTGCAGGAAGAGAATCGGATGGCAGGAATACTGGCTGGTCAAG
 C57MT.....G.....A.....A.....
 C573T.....A.A.G.T.ATTGGTCT...G...G.G.C...G...T...AT...A.A.....A.CA...
 SW3T.....A.A.G.T.ATTGGTCT...G...G.G.C...G...T...AT...A.A.....A.CA...
 SW3 Y K R G I Y Y E P N C N N S L F G A L R H A V L L V G Y G F I G R E S E G R K Y W I I K

 SWM N S M G T O W G N K G Y M K I S R D K G N H C G I A T Y A L Y P R V
 SWM AACAGCATGGGTACACAATGGGGCAATAAAGGCTATATGAAGATTTCAGAGACAGGGAAACCATCTGTGGAATTGCTACATATGCCCTTTATCCCCAGAGTG
 C57MC.....A.....T.....G.A...C.....T.C.T.C.G...T.....
 C573C.....A.....G.A...C.....T.C.T.C...
 SW3 N S L G T K W G N V K G Y M K I A K D O G N H C G I A S L P L V Y P R V

Fig. 2. Polymorphisms in mouse cathepsins M and 3. cDNA Sequences of cathepsins M and 3 from C57BL/6 mice (marked C57M and C573, respectively) were obtained from the NCBI database and the homologous sequences from Swiss Webster mice (SWM and SW3) were obtained in this study. Only bases that differ from Swiss Webster cathepsin M are noted. The amino acid sequences shown are for the proteins in Swiss Webster mice with polymorphic differences in the C57BL/6 mice shown in lower case. Sites of polymorphism between the two different strains of mice are shaded.

Structural models of PECs

Structural models offer a better view of the contribution of amino acid changes to enzyme structure and permit comparisons of features of the species variants of the enzymes. Conserved features are more likely to be important for enzyme function. The close similarity between the sequences of cathepsins K and L and the PECs allowed the generation of molecular models of each enzyme (Fig. 3). All of the models predict stable structures for the PECs (Fig. 3), with overall energies similar to those of cathepsin L. Significant differences are noted in residues aligning the specificity pockets of the PECs. In general, the active site clefts of the PECs are predicted to be narrower than that of cathepsin L, and this may restrict access of protein substrates. The amino acids surrounding the specificity sites of the PECs are generally conserved between the mouse PECs and their rat counterparts, implying evolutionary pressure to maintain substrate specificity of the individual PECs.

Conservation of residues around the active site of the two species variants of cathepsin R is particularly striking. Two of the aspartic acid residues in mouse cathepsin R that line the S' specificity pockets are replaced by asparagine residues in rat cathepsin R by alteration of all three bases that code for each amino acid. The retention of the carbonyl group of the ASX residues indicates evolutionary pressure to retain the potential substrate binding characteristics of these residues. A third aspartic acid residue near the S₂ pocket is conserved, and this might support aminopeptidase activity of these enzymes. A lysine residue is also conserved near to the S'₂ pocket and could support carboxypeptidase activity.

For cathepsin Q, a tryptophan appears to block the S₁ and S₂ substrate binding pockets. In models of the pro-enzyme, this tryptophan rotates to accommodate the pro-peptide, and a similar rotation would be required for mature cathepsin Q to bind substrates. Interestingly, in rat cathepsin Q2 glycine replaces this tryptophan, which might provide a different proteolytic function for this enzyme. Rat cathepsin P has a tryptophan residue in this same position whereas the mouse homolog does not. The two cathepsin P structures look similar if this residue is rotated to expose the active-site cleft. Conserved lysine residues near the S'₂ pocket might support carboxypeptidase activity of cathepsin P.

Closer examination of the catalytic triad for each of the models show that the active-site asparagine can form a weak hydrogen bond to the active-site histidine to orientate the latter residue correctly for catalysis with the active-site cysteine. The predicted active-site residue geometries are similar to those of the structures of cathepsin L. These models show that each of the PECs are probably catalytically functional.

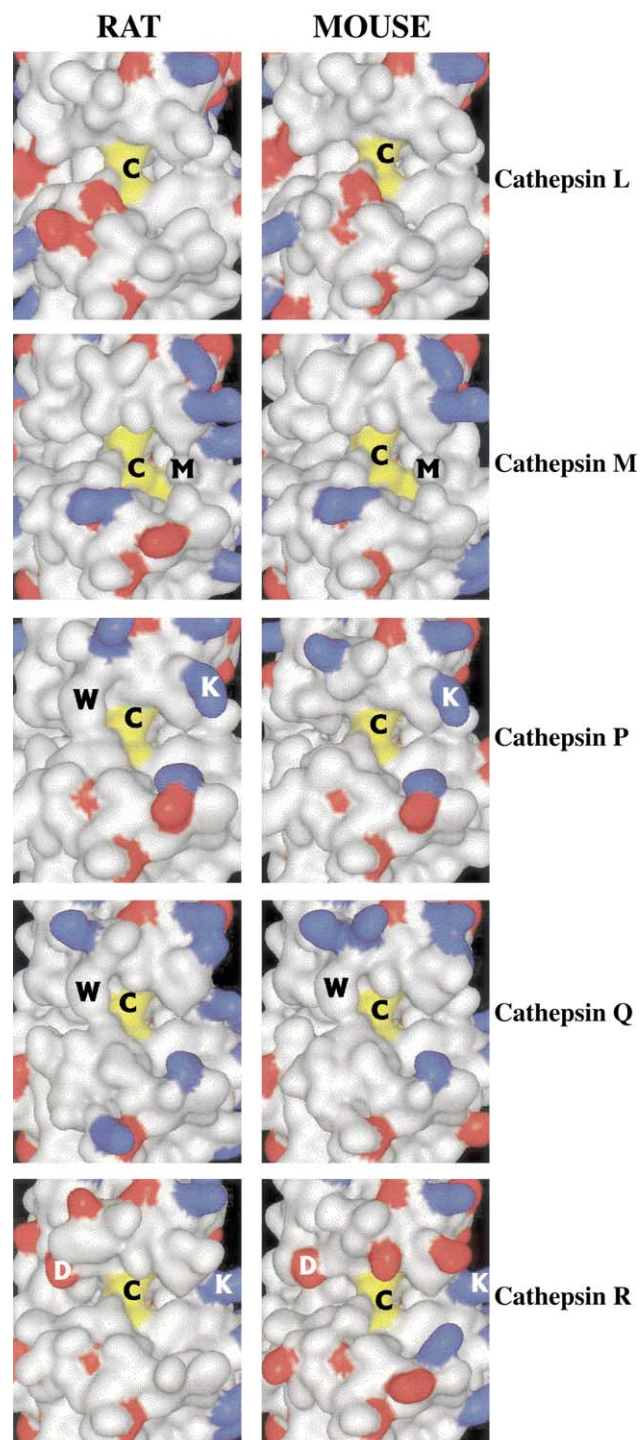


Fig. 3. Molecular models of rat and mouse cathepsins L, M, P, Q, and R. Molecular models of the mature portions of rat and mouse cathepsins L, M, P, Q, and R were prepared using SWISS-Model as described in Materials and methods. The images displayed focus on the active-site region of the proteases, and substrates bind from left to right in an N–C direction in the active-site groove. The active-site residues (cysteine, histidine, and asparagine) are colored yellow and positive and negative charges are colored blue and red, respectively. Single letter codes are used to identify important residues. Final energies of the models were –4315, –5633, –4013, –5233, –3988, –4542, –5427, –4638, –2418, –2029 kJ/mol for rat and mouse cathepsins L, M, P, Q, and R, respectively.

The cathepsin M models predict a more open active site than the other PECs, although a methionine residue partially occludes the S' specificity pockets. Structural models of cathepsin 3 from C57BL/6 and Swiss Webster mice are similar to each other, but quite distinct from cathepsin M (not shown). The methionine that partially occludes the S' pockets of cathepsin M is not conserved in cathepsin 3, making the S' pockets more open. Although the primary sequences of cathepsins M and 3 are closely related, the distinct active-site clefts indicate and each enzyme has probably evolved a unique catalytic property.

Evolutionary relationships

The sequences of the PEC family members were compared to other cathepsins by phylogenetic analysis and found to be most closely related to cathepsin L (Fig. 4). The upper portion of Fig. 4 shows the relationship between published cathepsin L sequences derived from 14 different species, including mammals, birds, fish, and insects. Although this phylogenetic analysis implies divergence of the PECs before species divergence of the cathepsin L, PECs are only found in rodent placenta. The rate of divergence of the PEC DNA sequences must be significantly greater than that for cathepsin L to account for the large differences in the PEC sequences.

Within the PEC sequences there are at least four groups of enzymes common to both rats and mice, i.e., cathepsins M, P, Q, and R. The sequences of the mature portions of rat cathepsins M, P, Q, and R are 83%, 86%, 88%, and 78% identical to their mouse counterparts, respectively. By contrast, mouse and rat counterparts of cathepsins L are less divergent, being 94% identical to each other.

The sequences of cathepsins M and 3 are very closely related and cDNA differences only result in two amino acid changes in the 96 amino acid pro-peptide (Fig. 2). Most of the differences between cathepsins M and 3 are clustered around the two sites where it is necessary to insert gaps to align the sequences of cathepsin 3 with the other cathepsins. Deletions and insertions usually require mutations in adjacent regions to permit correct folding of proteins. If the multiple base-changes at these sites are considered as single mutational events, phylogenetic analysis indicates that cathepsins M and 3 are even more closely related and probably arose by gene duplication after divergence of rats and mice. Differences outside these regions are spread throughout the mature polypeptide. At these sites, rat cathepsin M shares amino acid identity equally with either mouse cathepsin M or cathepsin 3, consistent with our proposal that gene duplication to generate the two mouse cathepsins probably occurred after divergence of rats and mice.

Rat counterparts of cathepsins 6, 1, and 2 have not yet been cloned, but genomic sequences related to

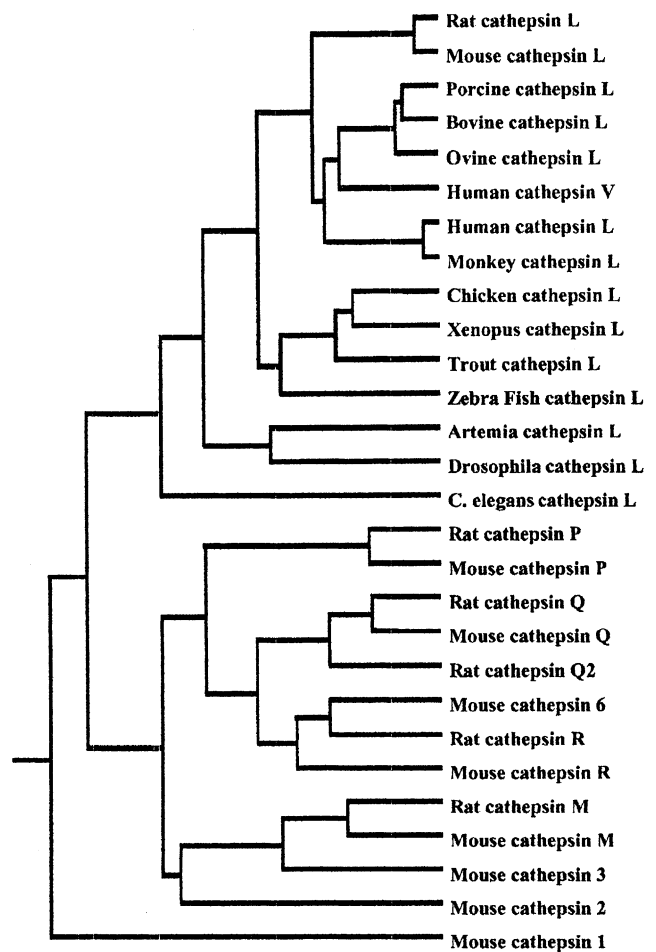


Fig. 4. Phylogenetic analysis of cathepsin L and the PECs. Sequences of each of the cathepsins were aligned using the Phylip program as described in Materials and methods. Cathepsin L sequences are shown at the top. Lengths of lines indicate differences between protein sequences. Accession numbers used to obtain amino acid sequence data were: mouse cathepsin L (P06797), mouse cathepsin M (AF202528), mouse cathepsin P (AF158182), mouse cathepsin Q (AF456461), mouse cathepsin R (AF245399), mouse cathepsin 3 (AF456464), mouse cathepsin 6 (AAK00510), mouse cathepsin 1 (AF250837), mouse cathepsin 2 (AF250840), rat cathepsin L (P07154), rat cathepsin M (AF456462), rat cathepsin P (AF456458), rat cathepsin Q (AF187323), rat cathepsin Q2 (AF456460), rat cathepsin R (AF456459), human cathepsin L (CAA30981), human cathepsin V (XP_005571), monkey cathepsin L (AF201700), porcine cathepsin L (Q28944), bovine cathepsin L (P25975), ovine cathepsin L (Q10991), chicken cathepsin L (KHCHL), xenopus cathepsin L (CAA75862), trout cathepsin L (AF358668), zebrafish cathepsin L (CAA69623), artemia cathepsin L (AF147207), drosophila cathepsin L (CATL_DROME), c. elegans cathepsin L (CAB07275).

cathepsins 1 and 2 can be found in the NCBI rat genome trace sequence database. Cathepsin 6 is closely related to cathepsin R and, like cathepsin 3, may be mouse-specific.

Rapid rates of divergence have been seen for other species-specific gene duplications and it has been proposed that this may function to generate diversity of specificity. Of particular note is a family of aspartic

proteases called pregnancy associated glycoproteins (PAGs) expressed in the placenta of ruminants [21]. The appearance of two different families of proteases in placenta of two different sub-sets of mammals indicates that proteases probably play important roles in placentation. Duplication of ancestral proteases that have multiple proteolytic functions would allow individual enzymes to evolve more unique specificities. The gene duplications and molecular models of PECs indicate that the rapid divergence of this family of proteases has given rise to proteins that have unique active site clefts and consequently unique catalytic or binding properties. The biological function of the PECs are not yet known, but the gene duplications and subsequent mutations conserve structural motifs and give rise to similar enzymes in both rats and mice, indicating that these enzymes have evolved to perform a range of functions in rodent placenta. Mammals that do not express PECs are presumably dependent on other proteases such as cathepsins B and L for placenta function.

Acknowledgments

The work was supported in part by the Nemours Foundation and NIGMS (GM59183).

References

- [1] S. Afonso, L. Romagnano, B. Babiaryz, The expression and function of cystatin C and cathepsin B and cathepsin L during mouse embryo implantation and placentation, *Development* 124 (1997) 3415–3425.
- [2] S. Afonso, L. Romagnano, B. Babiaryz, Expression of cathepsin proteinases by mouse trophoblast in vivo and in vitro, *Dev. Dyn.* 216 (1999) 374–384.
- [3] S.J. Freeman, J.B. Lloyd, Inhibition of proteolysis in rat yolk sac as a cause of teratogenesis. Effects of leupeptin in vitro and in vivo., *J. Embryol. Exp. Morph.* 78 (1983) 183–193.
- [4] J.L. Ambroso, C. Harris, In vitro embryotoxicity of the cysteine proteinase inhibitors benzyloxycarbonyl-phenylalanine-alanine-diazomethane (Z-Phe-Ala-CHN₂) and benzyloxycarbonyl-phenylalanine-phenylalanine-diazomethane (Z-Phe-Phe-CHN₂), *Teratology* 50 (1994) 214–228.
- [5] J.D. Grubb, T.R. Koszalka, J.J. Drabick, R.M. Mettrione, The activities of thiol proteases in the rat visceral yolk sac increase during late gestation, *Placenta* 12 (1991) 143–151.
- [6] K. Sol-Church, J. Shipley, D.A. Beckman, R.W. Mason, Expression of cysteine proteases in extraembryonic tissues during mouse embryogenesis, *Arch. Biochem. Biophys.* 372 (1999) 375–381.
- [7] M. Nilsen-Hamilton, Y.J. Jang, M. Delgado, J.K. Shim, K. Bruns, C.P. Chiang, Y. Fang, C.L. Parfett, D.T. Denhardt, R.T. Hamilton, Regulation of the expression of mitogen-regulated protein (MRP; proliferin) and cathepsin L in cultured cells and in the murine placenta, *Mol. Cell Endocrinol.* 77 (1991) 115–122.
- [8] W. Halangk, M.M. Lerch, B. Brandt-Nedele, W. Roth, M. Ruthenburger, T. Reinheckel, W. Domschke, H. Lippert, C. Peters, J. Deussing, Role of cathepsin B in intracellular trypsinogen activation and the onset of acute pancreatitis, *J. Clin. Invest* 106 (2000) 773–781.
- [9] T. Nakagawa, W. Roth, P. Wong, A. Nelson, A. Farr, J. Deussing, J.A. Villadangos, H. Ploegh, C. Peters, A.Y. Rudensky, Cathepsin L: critical role in *Ii* degradation and CD4 T cell selection in the thymus, *Science* 280 (1998) 450–453.
- [10] K. Sol-Church, J. Frenck, R.W. Mason, Cathepsin Q, a novel lysosomal cysteine protease highly expressed in placenta, *Biochem. Biophys. Res. Commun.* 267 (2000) 791–795.
- [11] K. Sol-Church, J. Frenck, D. Troeber, R.W. Mason, Cathepsin P, a novel protease in mouse placenta, *Biochem. J.* 343 (1999) 307–309.
- [12] K. Sol-Church, J. Frenck, R.W. Mason, Mouse cathepsin M, a placenta-specific lysosomal cysteine protease related to cathepsins L and P, *Biochim. Biophys. Acta* 1491 (2000) 289–294.
- [13] K. Sol-Church, J. Frenck, G. Bertenshaw, R.W. Mason, Characterization of mouse cathepsin R, a new member of a family of placentally expressed cysteine proteases, *Biochim. Biophys. Acta* 1492 (2000) 488–492.
- [14] A. Nakajima, K. Kataoka, Y. Takata, N.H. Huh, Cathepsin-6, a novel cysteine proteinase showing homology with and co-localized expression with cathepsin J/P in the labyrinthine layer of mouse placenta, *Biochem. J.* 349 (2000) 689–692.
- [15] J. Deussing, M. Kouadio, S. Rehman, I. Werber, A. Schwinde, C. Peters, Identification and characterization of a dense cluster of placenta-specific cysteine peptidase genes and related genes on mouse chromosome 13, *Genomics* 79 (2002) 225–240.
- [16] M. Hemberger, H. Himmelbauer, J. Ruschmann, C. Zeitz, R. Fundele, cDNA subtraction cloning reveals novel genes whose temporal and spatial expression indicates association with trophoblast invasion, *Dev. Biol.* 222 (2000) 158–169.
- [17] S.D. Bryce, S. Lindsay, A.J. Gladstone, K. Braithwaite, C. Chapman, N.K. Spurr, J. Lunec, A novel family of cathepsin L-like (CTSLL) sequences on human chromosome 10q and related transcripts, *Genomics* 24 (1994) 568–576.
- [18] D. Bromme, Z.Q. Li, M. Barnes, E. Mehler, Human cathepsin V functional expression, tissue distribution, electrostatic surface potential, enzymatic characterization, and chromosomal localization, *Biochemistry* 38 (1999) 2377–2385.
- [19] J. Guay, J.P. Falgout, A. Ducret, M.D. Percival, J.A. Mancini, Potency and selectivity of inhibition of cathepsin K, L and S by their respective propeptides, *Eur. J. Biochem.* 267 (2000) 6311–6318.
- [20] R. Coulombe, J.S. Mort, Structure of human procathepsin L reveals the molecular basis of inhibition by the prosegment, *EMBO J.* 15 (1996) 5492–5503.
- [21] A.L. Hughes, J.A. Green, J.M. Garbayo, R.M. Roberts, Adaptive diversification within a large family of recently duplicated, placentally expressed genes, *Proc. Natl. Acad. Sci. USA* 97 (2000) 3319–3323.