

Isofunctional Protein Subfamily Detection using Data Integration and Spectral Clustering

Elisa Boari de Lima^{1,2,*}, Wagner Meira Júnior², Raquel Cardoso de Melo-Minardi²

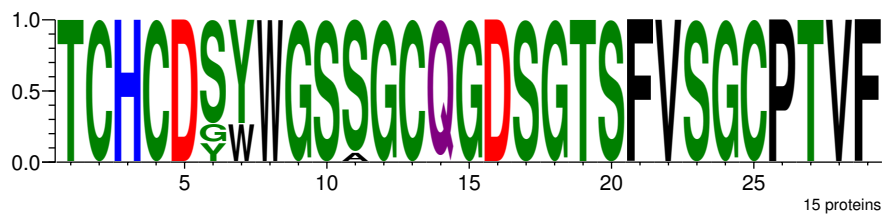
1 Department of Biochemistry and Immunology, Federal University of Minas Gerais, Belo Horizonte, MG, Brazil

2 Department of Computer Science, Federal University of Minas Gerais, Belo Horizonte, MG, Brazil

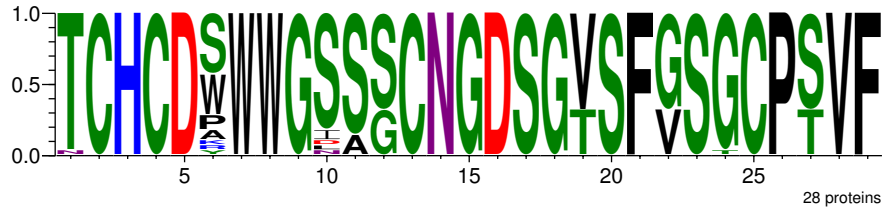
* eblima@dcc.ufmg.br

S9 Text: Dividing the serine proteases into eleven clusters

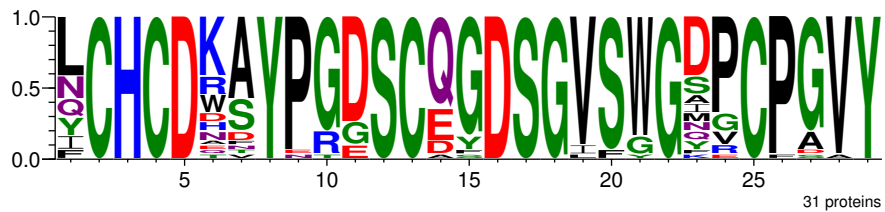
The second level of the hierarchical clustering produced by ASMC divided the serine proteases into eleven clusters, whose logos and compositions according to subfamily labels are presented in Fig. S9.1. This clustering has $MI = 10.59$.



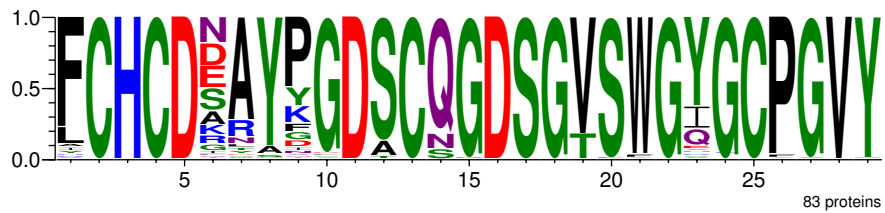
(a) Cluster I.A: 15 elastases



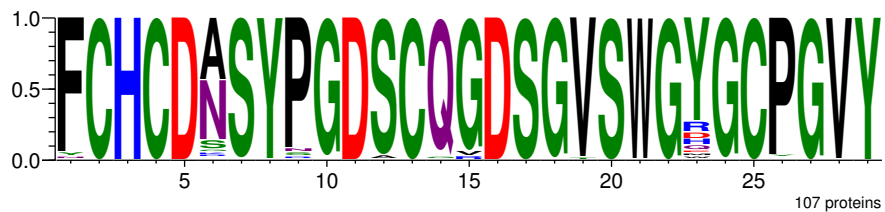
(b) Cluster I.B: 28 elastases



(c) Cluster II.A: 31 trypsins



(d) Cluster II.B: 83 trypsins



(e) Cluster II.C: 107 trypsins

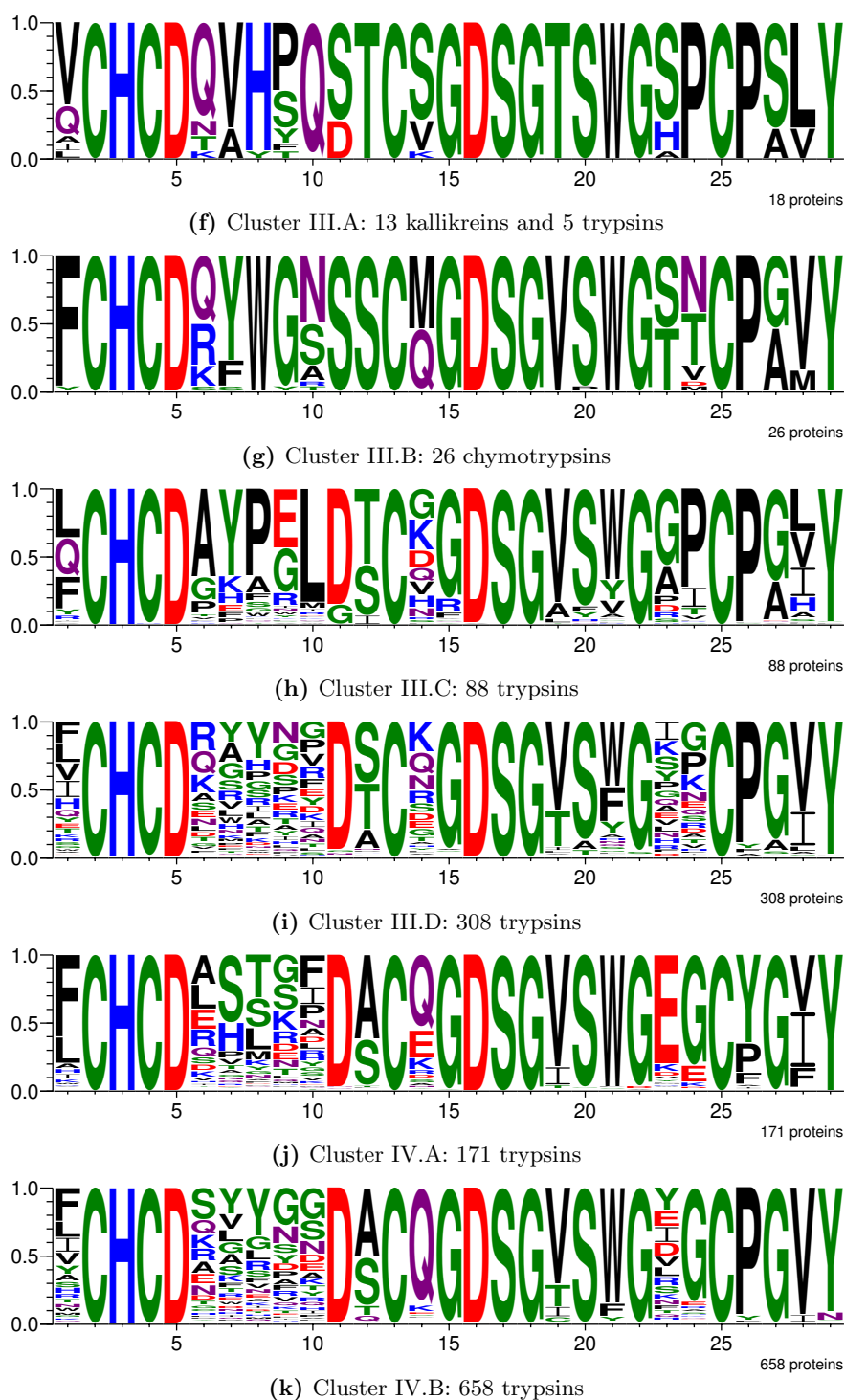


Figure S9.1. Serine protease division into eleven clusters in the second level of ASMC's hierarchical clustering.

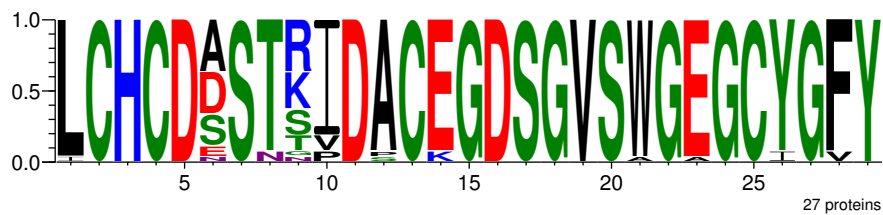
Given ASMC's criteria for considering as specificity determining positions (SDPs) those positions with p-values smaller than 0.0001 [1], the SDPs per cluster for this clustering are presented in Table S9.1, in which one may notice the presence of the known SDPs for this family, which were listed in the S8 Text.

Table S9.1. Cluster SDPs for the eleven serine protease clusters produced by ASMC.

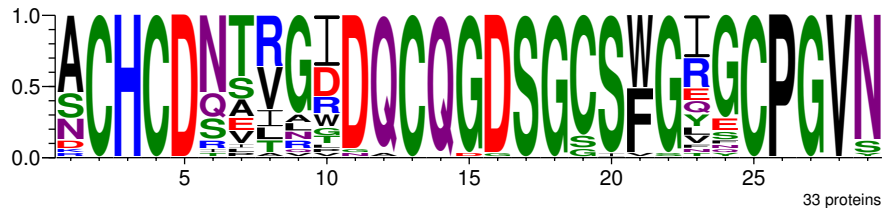
Cluster	Positions
I.A	11 ₁₈₉ , 22 ₂₁₆ , 27 ₂₂₆ , 29 ₁₇₃
I.B	8 ₁₇₂ , 11 ₁₈₉ , 14 ₁₉₂ , 27 ₂₂₆ , 29 ₁₇₃
II.A	-
II.B	-
II.C	7 ₁₇₁ , 9 ₁₇₃
III.A	11 ₁₈₉ , 27 ₂₂₆ , 28 ₂₂₇
III.B	11 ₁₈₉ , 14 ₁₉₂
III.C	8 ₁₇₂ , 10 ₁₇₄ , 24 ₂₁₉ , 28 ₂₂₇
III.D	21 ₂₁₅ , 24 ₂₁₉
IV.A	23 ₂₁₇ , 26 ₂₂₅ , 28 ₂₂₇
IV.B	12 ₁₉₀ , 14 ₁₉₂ , 28 ₂₂₇

Listed in order of active site position. Positions in bold correspond to known SDPs. Subscripted positions correspond to those in PDB structure 5PTP:A.

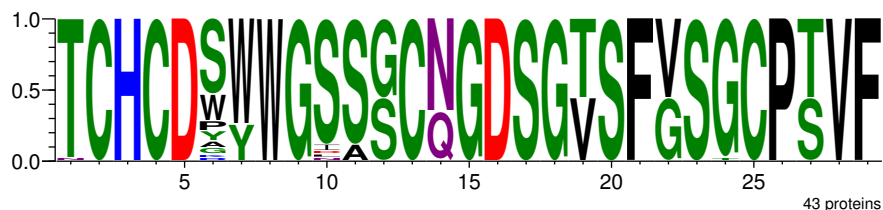
The best result obtained by the GP system for eleven clusters has $MI = 12.09$ and uses a combination of three data types, as shown in the main text. Cluster logos and compositions are presented in Fig. S9.2, while the residues which most distinguish each cluster are listed in S9.2, in which one can note the presence of the known SDPs for the cases in which the residues occurring in such positions distinguish the corresponding cluster from the others.



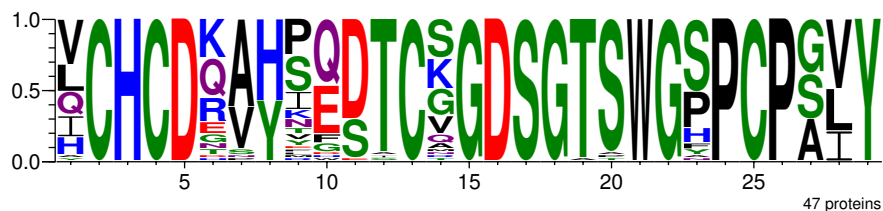
(a) Cluster I: 27 trypsins



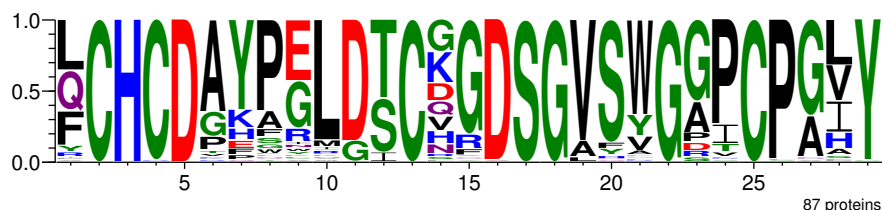
(b) Cluster II: 33 trypsins



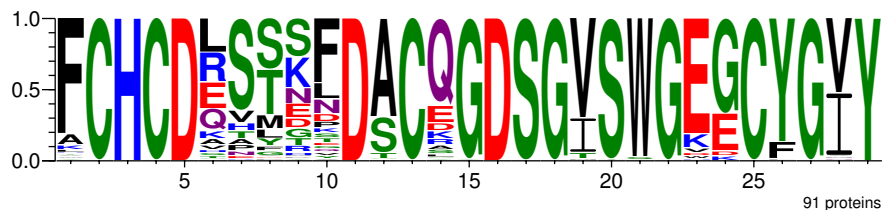
(c) Cluster III: 43 elastases



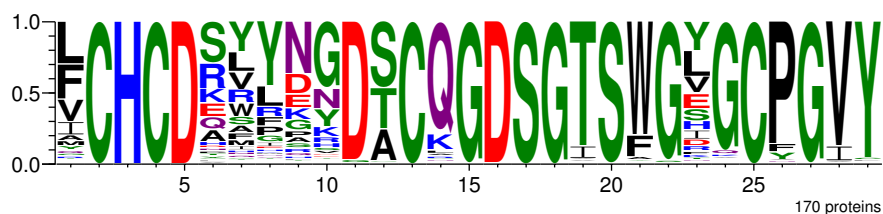
(d) Cluster IV: 34 trypsins and 13 kallikreins



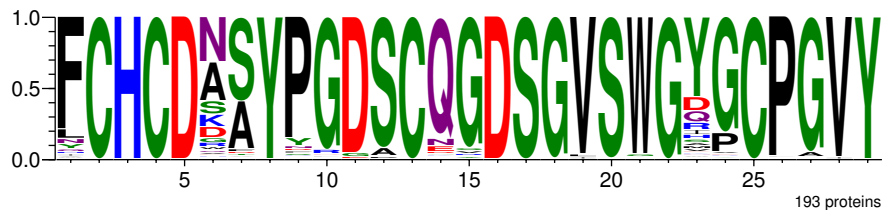
(e) Cluster V: 87 trypsins



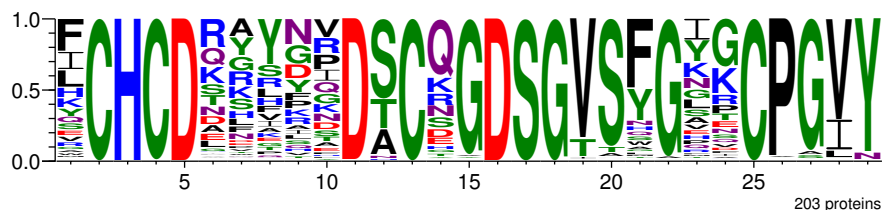
(f) Cluster VI: 91 trypsins



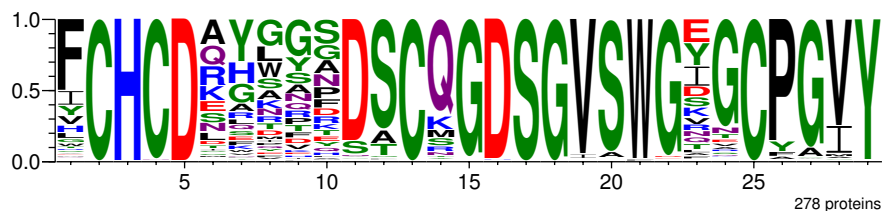
(g) Cluster VII: 170 trypsins



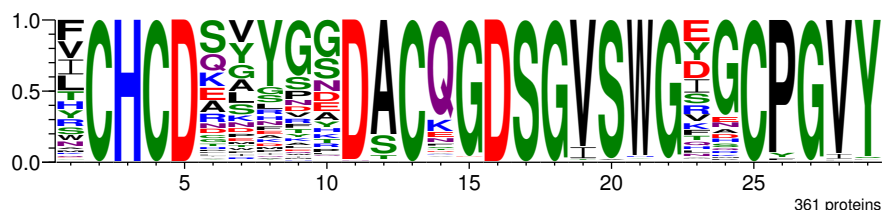
(h) Cluster VIII: 193 trypsins



(i) Cluster IX: 203 trypsins



(j) Cluster X: 252 trypsins and 26 chymotrypsins



(k) Cluster XI: 361 trypsins

Figure S9.2. Serine protease division into eleven clusters by the GP system.

Table S9.2. Most important residues for the eleven serine protease clusters produced by the GP system.

Cluster	Residues
I	F28 ₂₂₇ , E14 ₁₉₂ , T8₁₇₂ , I10 ₁₇₄ , Y26 ₂₂₅ , L1 ₄₁ , E23 ₂₁₇ , S8₁₇₂
II	Q12 ₁₉₀ , C19 ₂₁₃ , N29 ₂₂₉ , A1 ₄₁
III	F29 ₂₂₉ , W8₁₇₂ , T1 ₄₁ , S11₁₈₉ , F21 ₂₁₅ , W7 ₁₇₁ , V22₂₁₆ , S10 ₁₇₄ , T27₂₂₆
IV	P24 ₂₁₉ , T19 ₂₁₃ , T12 ₁₉₀ , H8₁₇₂
V	L10 ₁₇₄ , P8₁₇₂ , P24 ₂₁₉ , G23 ₂₁₇ , A6 ₁₇₀ , E9 ₁₇₃ , Y7 ₁₇₁
VI	Y26 ₂₂₅ , E23 ₂₁₇ , F10 ₁₇₄ , S8 ₁₇₁
VII	T19 ₂₁₃ , N9 ₁₇₃
VIII	P9 ₁₇₃ , G10 ₁₇₄ , Y8₁₇₂ , S8 ₁₇₁ , F1 ₄₁ , S12 ₁₉₀ , Y23 ₂₁₇
IX	F21 ₂₁₅ , Y21 ₂₁₅
X	S12 ₁₉₀ , V19 ₂₁₃ , W21 ₂₁₅
XI	A12 ₁₉₀ , W21 ₂₁₅ , G9 ₁₇₃

Listed in decreasing order of partial MI value. Residues in bold correspond to known SDPs. Subscripted positions correspond to those in PDB structure 5PTP:A.

References

1. Melo-Minardi RC, Bastard K, Artiguenave F. Identification of subfamily-specific sites based on active sites modeling and clustering. *Bioinformatics*. 2010 Dec;26(24):3075–3082.