# Prediction of the Secondary Structures of Stefins and Cystatins, the Low-Molecular Mass Protein Inhibitors of Cysteine Proteinases

Andrej ŠALI and Vito TURK

Dept. of Biochemistry, J. Stefan Institute, Ljubljana, Yugoslavia

**Summary:** A procedure for classifying proteins of known sequence into structurally similar groups was developed on the basis of the Argos parametric approach. It is shown that stefins and cystatins constitute two structurally well resolved, but homologous groups of proteins. Furthermore, it is very probable that segments of secondary structures within each family are conserved, although significant differences between stefins and cystatins are indicated at the level of secondary structure.

Next, secondary structures of all sequenced stefins and cystatins were predicted and used in the construction of secondary structures of the "typical stefin" and the "typical cystatin". Results were interpreted in the light of evolution and inhibition mechanism:[1] Alignment of the "typical stefin" versus the "typical cystatin" secondary structure segments suggests that the divergence of stefin and cystatin families did not occur by a gene fusion event, but only by a mechanism of substitution, insertion and/or deletion.[2] The central region of low-molecular mass cystatins, which is assumed to interact with cysteine proteinases, is predicted to be in a β-sheet conformation. This resembles the β-sheet in the active site of "standard mechanism" serine proteinases inhibitors.

*Rückschluß von der Primärstruktur auf die Sekundärstrukturen von Stefinen und Cystatinen, den niedermolekularen Proteininhibitoren von Cystein-Proteinasen*

**Zusammenfassung:** Auf der Basis der Parameter-Näherung (parametric approach) von Argos wurde eine Methode zur Klassifizierung von Proteinen bekannter Sequenz in strukturell ähnliche Gruppen entwickelt. Es wird gezeigt, daß Stefine und Cystatine zwei strukturell gut getrennte, aber homologe Gruppen von Proteinen darstellen. Außerdem ist es sehr wahrscheinlich, daß innerhalb jeder der Proteinfamilien Abschnitte mit einer bestimmten Sekundärstruktur erhalten geblieben sind, obwohl in bezug auf die Sekundärstruktur zwischen Stefinen und Cystatinen signifikante Unterschiede erkannt wurden.

Außerdem wurden die Sekundärstrukturen aller sequenzierten Stefine und Cystatine vorausgesagt und die Sekundärstrukturen des „typischen Stefins" sowie des „typischen Cystatins" ermittelt. Die Ergebnisse werden in bezug auf die Evolution und den Hemm-Mechanismus interpretiert. Ein Homologie-Vergleich der Verteilung der verschiedenen Sekundärstrukturen des „typischen Stefins" mit denen des „typischen Cystatins" läßt vermuten, daß die Unterschiede zwischen den Familien der Stefine und der Cystatine nicht durch Genfusionen entstanden sind, sondern nur durch Substitutionen,

---

Insertionen und/oder Deletionen. Für den zentralen Bereich der niedermolekularen Cystatine, von dem man annimmt, daß er mit den Cystein-Proteinasen reagiert, wird eine β-Faltblattstruktur vorausgesagt. Das erinnert an die β-Faltblattstruktur im aktiven Zentrum der Serin-Proteinase-Inhibitoren, die nach dem „Standard-Mechanismus" arbeiten.

---

---

Inhibitors of cysteine proteinases have been found in tissues and body fluids of several mammalian species[1,2]. They are believed to be involved in the control mechanism of intracellular or extracellular protein breakdown.

Two recent works[3,4] demonstrate that all protein inhibitors of cysteine proteinases sequenced so far represent a homologous group, which can be clearly subdivided into three distinct families. It should be emphasized that the principle underlying classification in both cases was evolutionary relationship based on primary sequence homology and that no obvious indications about the similarity of three-dimensional structures could be obtained. At the First International Symposium on Cysteine proteinases and their Inhibitors (Portorož, Yugoslavia, 1985) the name "cystatins" was proposed for the whole superfamily. It was also suggested that families 1, 2 and 3 should be referred to as stefins, cystatins and kininogens, respectively[5]. Stefins and cystatins have molecular masses between 11 and 14 kDa and are termed low-molecular mass cystatins. On the other hand, kininogens comprise three low-molecular mass cystatin segments and have overall molecular masses of 50—120 kDa.

Neither the tertiary structure of any cystatin nor the exact mode of its binding to cysteine proteinases is yet known. However, X-ray analysis of chicken cystatin is in progress[6]. It would be of interest to predict some structural features of low-molecular mass cystatins from their primary structure as long as X-ray data are not available.

In the following presentation, low-molecular mass cystatins are grouped on a structural basis. The resulting homology groups have enabled the prediction of secondary structures for the "typical stefin" and the "typical cystatin". Implications of these structures for evolution and inhibition mechanism are considered.

## Materials and Methods

All the programs were executed on an Apple II+ microcomputer.

### Protein clustering on structural basis

The physical parametric approach to protein sequence comparison[7] was used to derive the maximal summed crosscorrelation coefficient with a lag value between −20 and 20 for all possible pairs of stefins and cystatins. Once all the homologues had been compared in this way, a difference matrix was constructed and systematic clustering according to the method of weighted pair-group with arithmetic averaging was undertaken[8].

### Prediction of secondary structure

#### A) The Chou & Fasman method[9]

Although a computerized version was used, the final predictions of α-helix and β-sheet regions are not exact, due to subjective decisions between overlapping α-helix and β-sheet preferences. Provisional α-helix and β-sheet regions were predicted as segments of at least three consecutive tetrapeptides with corresponding average potential greater than 1.0 and the exact lengths then obtained with the aid of the Chou & Fasman termination rule and boundary information. Finally, the occasional overlaps were resolved by inspection of boundary information and areas under the curve of plotted potentials. The threshold value used in β-turn predictions was one and a half times that of the average β-turn occurrence probability.

#### B) The method of Garnier et al.[10]

A computerized directional version was used to predict α-helix and β-sheet regions for all low-molecular mass cystatins. Furthermore, the prediction was improved by taking into account the sequences of homologous proteins: the information provided by each residue of the homologue was simply added, and the sum divided by the number of homologues with the residue at the treated position present. The alignment used was that of Salvesen et al.[4].

### Hydrophobicity profiles

Diagrams for each low-molecular mass cystatin were constructed according to the method of Kyte and Doolittle[11]. A moving segment of six amino acid residues was used. Averaged stefin and cystatin plots were obtained by averaging the original profiles, which were aligned according to Salvesen et al.[4]. The point plotted at $(i + 2.5)$ refers to the averaged hydrophobicity of the hexapeptide ranging from the i-th to $(i + 5)$-th amino-acid residue.

## Results and Discussion

### Clustering

The extended method of Argos et al.[7] was used to construct a dendrogram which results in a numerical index of the relationships between low-molecular mass cystatins with known amino-acid sequence. The comparison is based on the

Table 1. Correspondence between ASCC and secondary structure homology.

Alignment and structural assignments for globins, phospholipases and cysteine proteinases are from Lesk and Chothia[16], Dufton et al.[17] and Kamphius et al.[18], respectively.

| Compared proteins | ASCC | Secondary structure differences |
|---|---|---|
| Human deoxyhaemoglobin α<br>Horse deoxyhaemoglobin α | 5.24 | none |
| Bovine pancreas PLA$_2$<br>*Crotalus atrox* PLA$_2$ | 3.13 | 8 residue α-helix is inserted, several<br>1 residue insertions in non-helical/sheet regions |
| Papain<br>Actinidin | 2.89 | single 1, 2 and 4 residue deletions at helix boundaries<br>and 1 residue insertion in non-helical/sheet region. |
| Human deoxyhaemoglobin α<br>Human deoxyhaemoglobin β | 2.59 | single 1, 2 and 5 residue insertions,<br>length of aligned helical segments varies. |
| 1.4 < ASCC < 2.0: possible structural correspondence is indicated, albeit with significant insertions and/or deletions.<br>4.0 < ASCC < 6.0: suggests good structural equivalence with few insertions and/or deletions. | | |

physical properties thought to determine folding of a given polypeptide sequence, i.e. surrounding hydrophobicity, polarity, hydration potential and Chou & Fasman preferences for α-helix, β-sheet and β-turn. It therefore follows the conservation or change of conformation even where sequence homology is not detectable. This provides a more effective way of assessing the structural relation between different proteins than the methods based on evolutionary distance. However, interpretation of the numerical index provided by the Argos summed crosscorrelation coefficient (ASCC) in terms of the presence or absence of insertions, deletions, etc. needs to be established. To this end several pairs of reference proteins were examined for inserted segments and for differences in secondary structure judged by optimal alignment based on known tertiary structures. These differences with computed ASCCs and original conclusions of Argos are shown in Table 1. It can be predicted that the optimal alignment of two sequences with ASCC greater than 4.0 correlates with no significant insertions and/or deletions and that both proteins are likely to possess similar secondary structure segments which dif-

fer slightly only in their boundaries. On the other hand, it is very probable that in order to align two sequences with ASCC less than 3.0, some insertions and/or deletions – which might also change their pattern of secondary structure segments – are necessary.
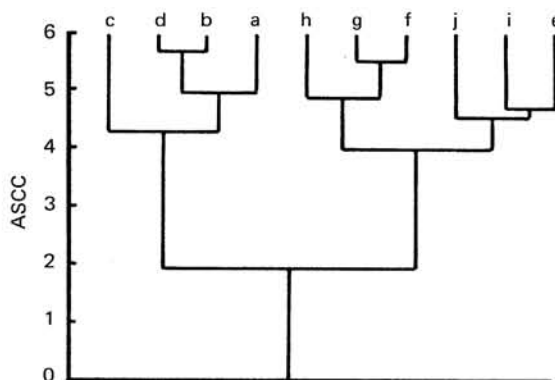


Fig. 1. Dendrogram constructed from the similarity matrix.

The cophenetic correlation coefficient is 0.963. For sequence key (small letters) see Fig. 4.

Table 2. Similarity matrix based on the method of Argos et al.

For sequence key (small letters) see Fig. 4.

|   | a | b | c | d | e | f | g | h | i | j |
|---|---|---|---|---|---|---|---|---|---|---|
| a | 6.00 | 5.05 | 4.45 | 4.92 | 1.45 | 1.16 | 1.34 | 2.06 | 1.91 | 1.19 |
| b | 5.05 | 6.00 | 4.07 | 5.69 | 1.94 | 1.29 | 1.64 | 2.12 | 2.58 | 1.46 |
| c | 4.45 | 4.07 | 6.00 | 4.30 | 2.76 | 1.96 | 1.74 | 2.25 | 2.28 | 2.02 |
| d | 4.92 | 5.69 | 4.30 | 6.00 | 2.02 | 1.41 | 1.62 | 1.98 | 2.74 | 2.06 |
| e | 1.45 | 1.94 | 2.76 | 2.02 | 6.00 | 4.12 | 4.01 | 3.83 | 4.73 | 4.50 |
| f | 1.16 | 1.29 | 1.96 | 1.41 | 4.12 | 6.00 | 5.54 | 5.09 | 3.45 | 4.47 |
| g | 1.34 | 1.64 | 1.74 | 1.62 | 4.01 | 5.54 | 6.00 | 4.66 | 3.76 | 4.62 |
| h | 2.06 | 2.12 | 2.25 | 1.98 | 3.83 | 5.09 | 4.66 | 6.00 | 3.58 | 3.96 |
| i | 1.91 | 2.58 | 2.28 | 2.74 | 4.73 | 3.45 | 3.76 | 3.58 | 6.00 | 4.61 |
| j | 1.19 | 1.46 | 2.02 | 2.06 | 4.50 | 4.47 | 4.62 | 3.96 | 4.61 | 6.00 |

The ASCCs for all possible pairs of low-molecular mass cystatins are shown in Table 2. This matrix was used to construct the dendrogram in Fig. 1. The very high cophenetic correlation coefficient indicates that pairwise relations are not considerably distorted in this particular dendrogram. Clustering clearly reveals two structurally distinct groups, namely stefins and cystatins, and therefore further supports the already suggested scheme of evolution[3,4].

Bearing in mind the results shown in Table 1 it is concluded that inhibitors within both families are practically identical at the level of secondary structure. It is also anticipated that considerable gaps have to be introduced in the optimal alignment of stefins and cystatins.

*Secondary structure predictions*

The prediction of α-helix and β-sheet segments according to the methods of Chou and Fasman[9]
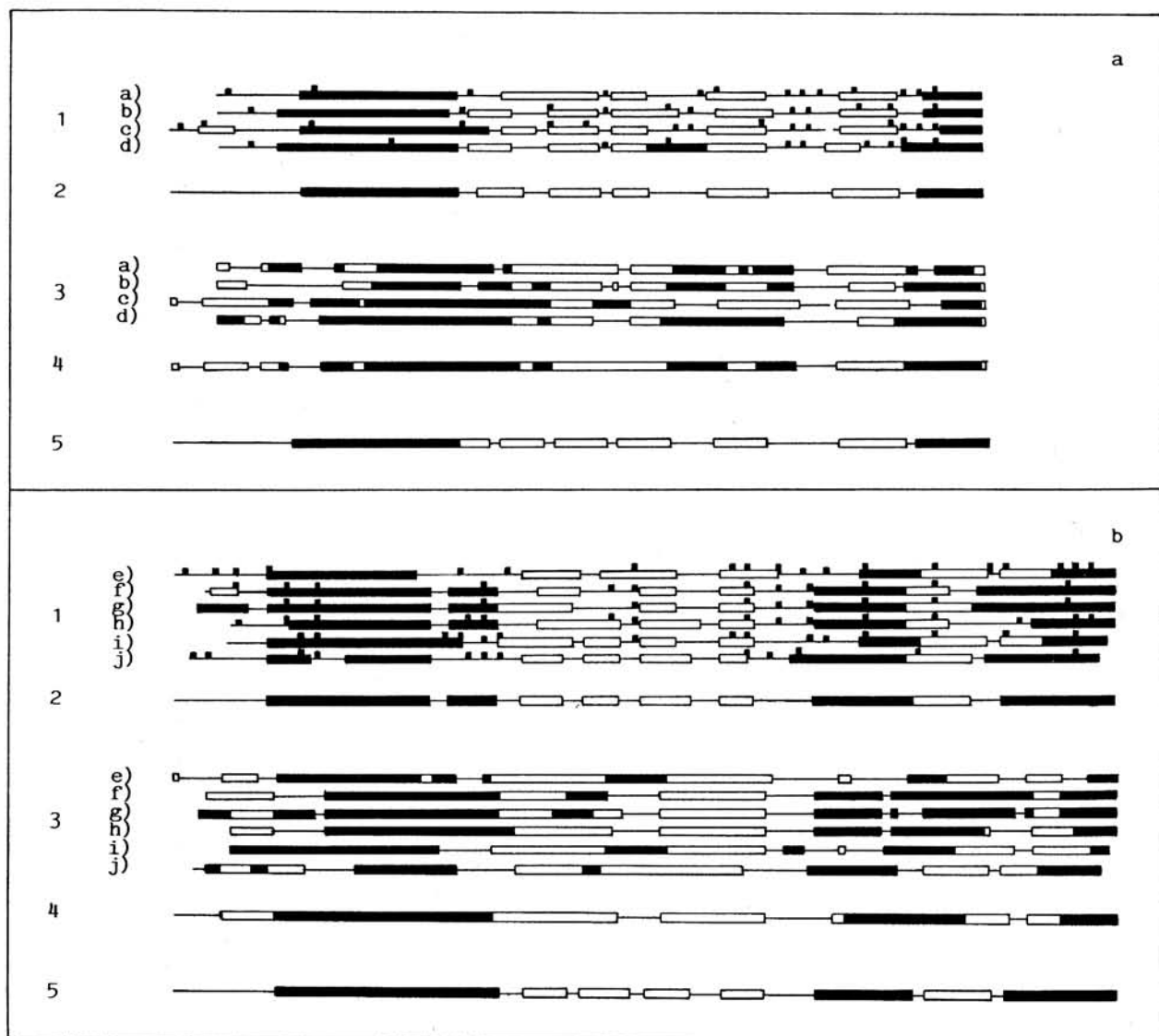


Fig. 2. Prediction of secondary structure.

1) Chou & Fasman predictions, 2) averaged Chou & Fasman prediction, 3) Garnier et al. predictions, 4) family Garnier et al. prediction, 5) typical family member, final combined prediction. Filled horizontal bars indicate helical regions and white horizontal bars code for β-strands. β-Turns predicted according to Chou and Fasman are designated by small squares at their first residue. For sequence key (small letters) see Fig. 4. *a) Stefins:* In the Garnier et al. prediction the decision constants were − 40, − 88 and 0 for α-helix, β-sheet and β-turn, respectively. *b) Cystatins:* In the Garnier et al. prediction the decision constants were − 75, − 88 and 0 for α-helix, β-sheet and β-turn, respectively. These decision constants were chosen because they minimize the difference between Garnier et al. and Chou & Fasman predictions.

and Garnier et al.[10] is shown for all the low-molecular mass cystatins (Fig. 2). In the next step the two methods were still considered separately and "averaged" structures for stefin and cystatin family were obtained (Fig. 2). For the Chou and Fasman method the "typical" secondary structure was constructed visually, whereas for the Garnier et al. method averaging was performed as suggested. Finally, secondary structures of the "typical stefin" and the "typical cystatin" were constructed from the averaged predictions of both methods (Fig. 2).

*Improvement of joint secondary structure predictions*

We used reverse turn prediction to improve rough secondary structure predictions of the typical stefin and typical cystatin. Turns were predicted applying the following considerations. First, independent Chou & Fasman $\beta$-turn prediction was used as an indicator of plausible chain reversals (Fig. 2). Second, we tended to locate reverse turns at the local minima of an averaged hydrophobicity profile (Fig. 3) and between already predicted secondary structure segments[12]. Third, we assume the globularity of a protein and consequently limited the maximal allowed length of $\alpha$-helices and $\beta$-segments to approximately 3 nm[12].

The only major impact of the refinement on the secondary structure predictions is a shortening of the first helix in both the typical stefin and typical cystatin. The final results in Fig. 3 indicate that

1) N- and C-terminals of the typical stefin and typical cystatin are helical,
2) the middle parts adopt the $\beta$-sheet conformation and
3) the typical cystatin has an additional helical segment preceding the last $\beta$-strand of the chain.

Unfortunately, individual predictions differ appreciably even within each family. In particular, boundaries are subject to the greatest uncertainties. Nevertheless, we can be rather confident that the regions of $\alpha$-helix and extended conformation are predicted with reasonable accuracy due to consideration of information from the whole family. The idea to improve secondary structure prediction using related proteins has been already discussed by Garnier et al.[10].

*Alignment*

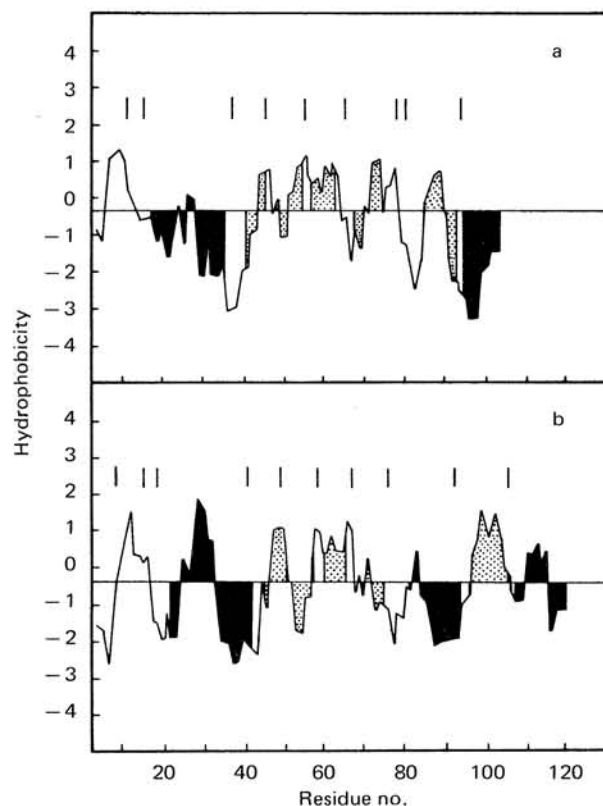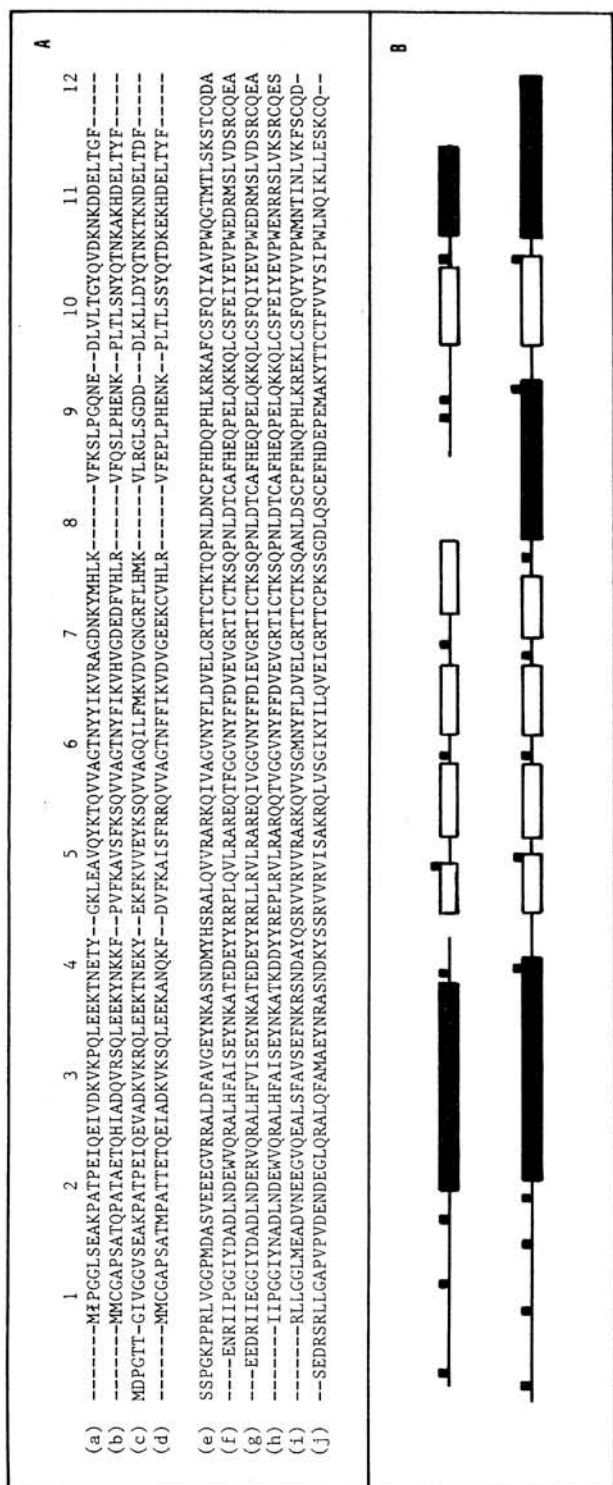The aligned segments of secondary structure of the "typical stefin" and the "typical cystatin"



Fig. 3. Averaged hydrophobicity profile and improved secondary structure prediction.

Filled areas code for $\alpha$-helix segments, dotted regions indicate $\beta$-sheet segments and horizontal bars code for reverse turns. a) Typical stefin, b) typical cystatin.

are shown in Fig. 4b. This kind of alignment may reveal the relation between corresponding structures and also shed light on the evolution of stefin and cystatin families from their common ancestor. It is evident that, by the criteria of secondary structure, stefins and cystatins are very similar at the N-terminal and central region, whereas there are greater differences at the C-terminus. This is in agreement with amenability of that part of the molecule to accept evolutionary mutations as a consequence of its exposure (which can be deduced from the hydrophobicity profile). It should be mentioned that the alignment cannot be clearly determined. Nevertheless, the transformation between $\alpha$-helix and $\beta$-sheet has not yet been observed within a homologous group of proteins[13]. The simplest way to take this into account is an alignment with some minor insertions, deletions and/or substitutions at the C-terminus. It is interesting to note that the sequence alignment[4] based on the program "Align" (Dayhoff) is almost identical with the one presented here (Fig. 4), although different principles are involved.

**A**

```
                    1            2             3              4               5                6                7                 8              9            10           11           12
(a)  -------M?PGGLSEAKPATPEIQEIVDKVKPQLEEKTNETY--GKLEAVQYKTQVVAGTNYYIKVRAGDNKYMHLK-------VFKSLPGQNE--DLVLTGYQVDKNKDDELTGF------
(b)  -------MMCGAPSATQPATAETQHIADQVRSQLEEKYNKKF--PVFKAVSFKSQVVAGTNYFIKVHVGDEDFVHLR-------VFQSLPHENK--PLTLSNYQTNKAKHDELTYF------
(c)  MDPGTT-GIVGGVSEAKPATPEIQEVADKVKRQLEEKTNEKY--EKFKVVEYKSQVVAGQLLFMKVDVGNGRFLHMK-------VLRGLSGDD---DLKLLDYQTNKTKNDELTDF------
(d)  -------MMCGAPSATMPAITTETQEIADKVKSQLEEKANQKF--DVFKAISFRRQVVAGTNFFIKVDVGEEKCVHLR-------VFEPLPHENK--PLTLSSYQTDKEKHDELTYF------

(e)  SSPGKPPRLVGGPMDASVEEEGVRRALDFAVGEYNKASNDMYHSRALQVVRARKQIVAGVNYFLDVELGRTTCTKTQPNLDNCPFHDQPHLKRKAFCSFQIYAVPWQGTMTLSKSTCQDA
(f)  -------ENRIIPGGIVDADLNDEWVQRALHFAISEYNKATEDEYYRRPLQVLRAREQTFGGVNYFFDVEVGRTICTKSQPNLDTCAFHEQPELQKKQLCSFEIYEVPWEDRMSLVDSRCQEA
(g)  -------EEDRIEGGIYDADLNDERVQRALHFVISEYNKATEDEYYRRLLRVLRAREQIVGGVNYFFDIEVGRTICTKSQPNLDTCAFHEQPELQKKQLCSFQIYEVPWEDRMSLVDSRCQEA
(h)  -------IIPGGIYNADLNDEWVQRALHFAISEYNKATKDDYYREPLRVLRARQQTVGGVNYFFDVEVGRTICTKSQPNLDTCAFHEQPELQKKQLCSFEIYEVPWENRRSLVKSRCQES
(i)  -------RLLGGLMEADVNEEGVQEALSFAVSEFNKRSNDAYQSRVVRVVRARKQVVSGMNYFLDVELGRTTCTKSQANLDSCPFHNQPHLKREKLCSFQVYVVPWMNTINLVKFSCQD-
(j)  --SEDRSRLLGAPVPVDENDEGLQRALQFAMAEYNRASNDKYSSRVVRVISAKRQLVSGIKYILQVEIGRTTCPKSSGDLQSCEFHDEPEMAKYTTCTFVVYSIPWLNQLLLESKCQ--
```

**B**

◁

Fig. 4. a) Amino-acid sequences of low-molecular mass cystatins aligned according to Salvesen et al.[4].

Key and references: a) human stefin A[19], b) human stefin B[20,21], c) rat cystatin α[22], d) rat cystatin β[23], e) human cystatin C[24], f) human cystatin S7[25], g) human cystatin S5[26], h) human cystatin SN[27], i) beef colostrum cystatin[28], j) chicken cystatin[29]. b) Alignment of secondary structure segments of the typical stefin and the typical cystatin.
For secondary structure key see Fig. 2.

## Implication for inhibition mechanism

The second observation evident from Fig. 4 is that the central part of both types of low-molecular mass cystatins is in a β-sheet conformation. It has already been suggested[14], although without firm evidence, that the active-site residues are located in the conserved central region of the sequence and that the inhibition mechanism might be similar to that of the "classical mechanism" of protein serine proteinases inhibitors[15]. It is worth noting that the interaction of a "classical inhibitor" with the serine proteinase is mainly through β-sheet hydrogen bonds. Our prediction of a β-sheet in the central region indicates that this region might be involved in the interaction with cysteine proteinases similar to that of the "classical mechanism" inhibitors.

## Conclusions

A general method was devised to assess the relation between secondary structure of proteins with known amino-acid sequences. It was used to compare low-molecular mass cystatins and to group them into two structurally homologous families, cystatins and stefins. The members of each group were found to be similar enough to allow construction of a secondary structure representing a typical protein of the family. By using the extensive information available from the whole homology group, it was possible to rationally bypass differences due to choosing different prediction methods and/or different proteins within the same family. We believe that the secondary structure resulting from such pooling of information from closely related proteins is more accurate than that relying on only one sequence and one prediction method.

This work and new sequences may also represent starting point for further theoretical work, such as prediction of supersecondary structure. Application of the emerging tertiary structure of any cystatin or stefin to other members

## The evolutionary implications

If we do accept that the alignment in Fig. 4 accounts for evolutionary changes, it follows that the divergence of stefin and cystatin families probably did not occur by a gene fusion event[3] but by simple substitution, deletion and/or insertion[4].

of the superfamily may now also be more straightforward.

Above all, we are looking forward to compare our predictions to the first X-ray data.

*Literature*

1  Turk, V., Brzin, J., Kopitar, M., Kotnik, M., Lenarčič, B., Popovič, T., Ritonja, A., Trstenjak, M., Rozman, B. & Machleidt, W. (1986) in *Proteinases in Inflammation and Tumor Invasion* (Tschesche, H., ed.) pp. 77–92, Walter de Gruyter, Berlin.

2  Barrett, A.J., Rawlings, N.D., Davies, M.E., Machleidt, W., Salvesen, G. & Turk, V. (1986) in *Proteinase Inhibitors* (Barrett, A.J. & Salvesen, G., eds.) pp. 515–569, Elsevier, Amsterdam.

3  Müller-Esterl, W., Iwanaga, S. & Nakanishi, S. (1986) *Trends Biochem. Sci.* **11**, 336–339.

4  Salvesen, G., Parkes, C., Rawlings, N.D., Brown, M.A., Barrett, A.J., Abrahamson, M. & Grubb, A. (1986) in *Cysteine Proteinases and Their Inhibitors* (Turk, V., ed.) pp. 413–428, Walter de Gruyter, Berlin.

5  Barrett, A.J., Fritz, H., Grubb, A., Isemura, S., Jarvinen, M., Katunuma, N., Machleidt, W., Müller-Esterl, W., Sasaki, M. & Turk, V. (1986) *Biochem. J.* **236**, 312.

6  Bode, W., Brzin, J. & Turk, V. (1985) *J. Mol. Biol.* **181**, 331–332.

7  Argos, P., Hanei, M., Wilson, J.M. & Kelley, W.N. (1983) *J. Biol. Chem.* **258**, 6450–6457.

8  Davis, J.C. & Sampson, R.J. (1973) *Statistics and Data Analysis in Geology*, pp. 456–473, J. Willey, New York.

9  Chou, P.Y. & Fasman, G.D. (1978) *Adv. Enzymol.* **47**, 45–148.

10  Garnier, J., Osguthorpe & Robson, B. (1977) *J. Mol. Biol.* **120**, 97–120.

11  Kyte, J. & Doolittle, R.F. (1982) *J. Mol. Biol.* **157**, 105–132.

12  Cohen, F.E., Abarbanel, R.M., Kuntz, I.D. & Fletterick, R.J. (1986) *Biochemistry* **25**, 266–275.

13  Creighton, T.E. (1984) in *Proteins*, p. 255, W.H. Freeman & Co., New York.

14  Barrett, A.J. (1985) in *Intracellular Protein Catabolism* (E.A. Khairallah, Bond, J.S., Bird, J.W.C., eds.) pp. 105–116, A.R. Liss, New York.

15  Laskowski, M., Jr. & Kato, J. (1980) *Annu. Rev. Biochem.* **49**, 593–626.

16  Lesk, A. & Chothia, C. (1980) *J. Mol. Biol.* **136**, 227–270.

17  Dufton, M.J., Eaker, D. & Hider, R.C. (1983) *Eur. J. Biochem.* **137**, 537–544.

18  Kamphius, I.G., Drenth, J. & Baker, N.E. (1985) *J. Mol. Biol.* **182**, 317–329.

19  Machleidt, W., Borchart, U., Fritz, H., Brzin, J., Ritonja, A. & Turk, V. (1983) *Hoppe-Seyler's Z. Physiol. Chem.* **364**, 1481–1486.

20  Lenarčič, B., Ritonja, A., Šali, A., Kotnik, M., Turk, V. & Machleidt, W. (1986) in *Cysteine Proteinases and Their Inhibitors* (Turk, V., ed.) pp. 473–487, Walter de Gruyter, Berlin.

21  Ritonja, A., Machleidt, W. & Barrett, A.J. (1985) *Biochem. Biophys. Res. Commun.* **131**, 1187–1192.

22  Takio, K., Kominami, E., Bando, Y., Katunuma, N. & Titani, K. (1984) *Biochem. Biophys. Res. Commun.* **121**, 149–154.

23  Takio, K., Kominami, E., Wakamatsu, N., Katunuma, N. & Titani, K. (1983) *Biochem. Biophys. Res. Commun.* **115**, 902–908.

24  Grubb, A. & Lofberg, H. (1982) *Scand. J. Clin. Lab. Invest.* **45**, suppl. 177, 7–13.

25  Isemura, S., Saitoh, E. & Sanada, K. (1984) *J. Biochem.* **96**, 489–498.

26  Isemura, S., Saitoh, E., Sanada, K., Isemura, M. & Ito, S. (1986) in *Cysteine Proteinases and Their Inhibitors* (Turk, V., ed.) pp. 497–505, Walter de Gruyter, Berlin.

27  Isemura, S., Saitoh, E. & Sanada, K. (1986) *FEBS Lett.* **198**, 145–149.

28  Hirado, M., Tsunasawa, S., Sakiyama, F., Niinobe, M. & Fujii, S. (1985) *FEBS Lett.* **186**, 41–45.

29  Turk, V., Brzin, J., Longer, M., Ritonja, A., Eropkin, M., Borchart, U. & Machleidt, W. (1983) *Hoppe-Seyler's Z. Physiol. Chem.* **364**, 1487–1496.

Andrej Šali and Dr. Vito Turk, Dept. of Biochemistry, J. Stefan Institute, Jamova 39, 61000 Ljubljana, Yugoslavia.