

蛋白質工学支援システムにおける解析アルゴリズム

藤本哲知、高橋裕信、窪田 綏¹⁾、中島広志²⁾、大井龍夫³⁾
¹⁾三洋電機筑波研究所、²⁾金沢大学医療短大、³⁾京都女子大学

ワークステーション SUN3/260上に、蛋白質一次構造(アミノ酸配列)データベースと一次構造解析プログラムとを結合させた蛋白質工学支援システムを構築した。このシステムに従えば、データベースの任意の蛋白質を検索することができ、そのアミノ酸配列あるいは組成からその蛋白質の2次構造や folding typeの予測など様々な有用な情報を得ることができる。さらに、マウスの使用やマルチウインドウによる複数の解析結果の表示など、マンーマシンインターフェースの向上に努めており、その操作性は単純明快なものとなっている。本報告では、システムの概要並びに一次構造解析のアルゴリズムについて述べたい。

なお、本システムにおける一次構造データベース及び解析プログラムは我々独自によって構築且つ開発されてきたものである。

On the algorithm of numerical analysis of protein sequence based on the correlation function

Tetsunori FUJIMOTO, Hironobu TAKAHASHI, Yasushi KUBOTA¹⁾, Hiroshi NAKASHIMA²⁾, Tatsuo OOI³⁾

1) Laboratory of Intelligent Systems, Tsukuba Research Center, SANYO Electric Co., Ltd., 2-1 Koyadai, Tsukuba, Ibaraki, 305

2) The School of Allied Medical Professions, Kanazawa Univ., Kanazawa Ishikawa 920

3) Kyoto Women's University, Higashiyama, Kyoto 605

We are developing a system for protein sequence analysis, coupling the database of amino acid sequences of proteins with application programs on workstation SUN3/260. By the system we can retrieve any sequence data stored in the database and analyse them to get available information on folding type or secondary structure of proteins from their amino acid compositions and sequences. The fundamental algorithm of the analysis is based on the correlation function which is widely used for the analysis of random data.

Here, we will describe the algorithm of numerical analysis of protein used in the system.

1. 緒言

蛋白質は核酸と共に生命の最も基本的な担い手である。中でも特に蛋白質は生化学反応、病気に対する免疫、分子識別などその機能は生命の維持に本質的且つ不可欠な役割を演じている。この機能は立体構造と相関があり、そして立体構造はそのアミノ酸配列(一次構造)によって一意的に決定される。蛋白質研究の最も重要な課題は与えられたアミノ酸配列に従って生体内で形成される立体構造を予測することである。この問題の完全な解答は現在、得られていない。このように極めて明快でありながら解決が困難な問題では他にはフェルマーの大定理があげられるであろうか? さて、蛋白質はアミノ酸が線形に重合したものである。アミノ酸残基をある適当な物理化学的なパラメータで置き替えれば蛋白質の有限な数列表現が得られる。この線形性のために我々は種々の統計学や情報学の手法を駆使して様々な立体構造に関する情報を得ることが出来る。また、1970年代にはSangerやMaxam-Gilbert法によりDNA塩基配列決定に大きな進展を見るに至った。このため従来はアミノ酸分析によって蛋白質の一次構造の決定を行っていたものが、DNAの塩基配列から翻訳することにより行われるようになった。ここに至って、蛋白質の一次構造に関するデータは飛躍的に蓄積された。このように、今や蛋白質の一次構造の数学解析やデータ処理にコンピュータの利用は必要不可欠のものとなっている。

他方、microelectronics技術の急速な発展によりいまやその能力においてmainframeに匹敵するEWS(Engineering Workstation)が広く普及するに至った。UNIXをoperating systemとして採用されたEWSは計算能力のみならずbitmap displayやmultiwindow systemにサポートされた優れたman-machine interfaceを実現している。さらには、RPC(Remote Procedure Call)やNFS(Network File System)に代表される分散処理システムは、様々なタイプの分散されたソフトウェア、ハードウェア資源へのアクセスを可能にしている。こうした高度なソフトウェア環境のもとで、EWSは蛋白質構造の解析研究に新たな可能性を与えるであろう。このような理由のために我々はSUN 3/260上に蛋白質一次構造の解析システムとデータベースを結合した統合システムを構築しつつある。このような統合システムは、九州大学大型計算機センターにおいてmainframe(FACOM M-382 OS IV/F4)上で久原らによっても試みられている^{1,2)}。

本報告では、我々が独自に構築し且つ開発して来たデータベースおよび蛋白質一次構造解析のアルゴリズムについて述べたい。

2. データベースと検索

蛋白質一次構造はその立体構造と機能に関する情報を秘めており、最も本質的な情報と考えられるので、我々はこの一次構造データを集め一つのデータベースを構築した。各データは蛋白質の名前とアミノ酸の個数が記述される" name "で始まり" // "で終わる。アミノ酸残基は一文字コードで表されている: 例えば m はメチオニン、e はグルタミン酸を表す等々。なお nbrf とあるのは 米国 N B R F (National Biomedical Research Foundation) のデータベースと同じ蛋白質の entry name である(図1)。全

```

Protein sequence database.
#name 759 p3 protein human influenza a virus updated 12/07/84
source human influenza a virus (strain a/pr/8/34)
nbrf p3iv34
reference 1 (sequence translated from the genomic rna sequence)
authors fields, s., and winter, g.;
journal cell 23, 303-313, 1982;
comment this protein is probably one of the three rna-dependent
rna polymerases.
sites from to description
matp 1 759 mature protein
sequence 759 aa
merikelrnlmsqsr treilktttvdhmaiikkytsgrqeknpalrmkwmmamkypitad 60
kritemipernegqgtlwskmndagsdrvmvplawtwnr ngpmtntvhyphkiyktye 120
rverlkhgtfgpvhfrnqvki rrryvdinpghadlsakeaqdvimevvpnevgariltse 180
sqtitkkekkeelqdcckisplmvaymlerelvrktrflpvaggtssvyievlhtqgtw 240
eqmlytpgsgevknndvdqsliaarnivrraavsadplasllemchstqiggirmvdiikq 300
npteeqgrratailrkatrrliqlivsgrdeqsi aeaiivamvfseqdcmikavrgdlmf 360
eeftmvgr ratalrkatrrliqlivsgrdeqsi aeaiivamvfseqdcmikavrgdlmf 420
vnr anqrlnqmqlrlrhfkdkakvlfnqwgvepidnvmgmigilpdmtpsiemsr gvri 480
skmgvdeysstervvvsidr flrvrdqrgnvlspvevsetgktekltitysssmmwein 540
gpesvlvntyqwli rnwetvkiqwsqntmlynk mefepfqsivpkairgqysgfvrtlf 600
qqmr dvlgtfdtaqliklplfaaappkqsrnqfssftvnrvgsgmriiivr gnspsvfnynk 660
atkr ltvlgkdagtl tedpdegtagvesavlrgfliigkedrrypalsinelsnlakge 720
kanvliqqgdvvlvmkrkr dssilt dsqtatkrirmain 759
//

```

図 1 アミノ酸配列データの format例

ての蛋白質はその生物学的観点から分類され41個のファイルに収められている。また、効果的な検索のためにアルファベット順に並べた蛋白質名とそれが記載されているファイル名からなるリスト・ファイルを用意した。キーボードから与えられる蛋白質名は、その蛋白質データが記載されているファイル名の情報を与えるこのリスト・ファイル上で探索される。このようにして、実際の蛋白質データへのアクセスはリスト・ファイル上で得られた情報をもとにしてなされる。なお、立体構造に関するデータは米国 Brookhaven 国立研究所のものを利用した³⁾。

物理化学的パラメータ

我々は疎水性や、 α らせん β 構造をとる傾向度などアミノ酸に固有な53個の物理化学的パラメータを用意している。

検索アルゴリズム

2つの文字列 $s_1s_2\cdots s_n$ 及び $t_1t_2\cdots t_n$ が与えられた時、2つの文字 s_i と t_j の間のズレの量として；

$$f(i, j) = \min\{f(i-1, j)+1, f(i, j-1)+1, f(i-1, j-1)+d(s_i, t_j)\}$$

を導入した⁴⁾。ここで、 $d(s_i, t_j)$ は $s_i=t_j$ の時 0、 $s_i\neq t_j$ の時 1の値を採る。この量をもとに、予め設定されたズレの許容範囲内の文字列を探し出すことで misspellingを許す検索を行う。なお、このアルゴリズムはC-言語で記述した。

3. 一次構造の数学解析

本システムには次のような解析プログラムが組み込まれている；

- ① 配列の繰り返し性の検出^{5, 6)}
- ② folding type (α , β , α/β , $\alpha+\beta$ などの折れたたみ方のタイプ)や細胞内外分泌予測^{7, 8)}
- ③ homology法による2次構造予測⁹⁾
- ④ 相関法による sequence homologyの検出^{6, 10)}

以下にこれらのアルゴリズムについて述べたい。

相関法による配列の繰り返し性^{5,6)}や sequence homology の検出^{6,10)}

20個のアミノ酸は疎水性のような物理化学的な量で表現できるので、与えられたアミノ酸配列はこのような値の数列に変換できる。従って、配列の繰り返し性の程度を計るために自己相関関数が定義できる。即ち、n 残基の長さの蛋白質 X の数列に対して、 τ 残基のラグの自己相関関数を：

$$\sum_{i=1}^{n-\tau} (x(i)-\bar{x})(x(i+\tau)-\bar{x})$$

$$A(\tau) = \frac{\sum_{i=1}^{n-\tau} (x(i)-\bar{x})(x(i+\tau)-\bar{x})}{\left[\left\{ \sum_{i=1}^{n-\tau} (x(i)-\bar{x})^2 \right\} \left\{ \sum_{i=1}^{n-\tau} (x(i+\tau)-\bar{x})^2 \right\} \right]^{1/2}}$$

但し、 $\bar{x} = \frac{1}{n-\tau} \left(\sum_{i=1}^{n-\tau} x(i) \right)$ 、 $\bar{x} = \frac{1}{n-\tau} \left(\sum_{i=1}^{n-\tau} x(i+\tau) \right)$

と定義する。ここで、 $\{x(i)\}$ はある適当な物理化学的な量で変換された数である。同様に、2つのアミノ酸配列、X (通常は X の部分列) と Y の間の homology の程度を計るために相互相関関数：

$$\sum_{i=1}^n (x(u+i-1)-\bar{x})(y(j+i-1)-\bar{y})$$

$$C(j) = \frac{\sum_{i=1}^n (x(u+i-1)-\bar{x})(y(j+i-1)-\bar{y})}{\left[\left\{ \sum_{i=1}^n (x(u+i-1)-\bar{x})^2 \right\} \left\{ \sum_{i=1}^n (y(j+i-1)-\bar{y})^2 \right\} \right]^{1/2}}$$

但し、 $\bar{x} = \frac{1}{n} \left(\sum_{i=1}^n x(u+i-1) \right)$ 、 $\bar{y} = \frac{1}{n} \left(\sum_{i=1}^n y(j+i-1) \right)$

を導入する。この式は、X のある固定した長さ n の部分列 (u 番目の残基で始まり、u+n-1 番目で終わる) を配列 Y の 1 番目の残基 (j=1) から順次走査しながら相関値を計算して行くことに相当する。更にこの手法は 2 次元的に拡張できる^{6,10)}。即ち、

$$\sum_{i=-k}^k (x_p(i+1)-\bar{c}_p)(y_p(j+1)-\bar{c}_p)$$

$$C_p(i, j) = \frac{\sum_{i=-k}^k (x_p(i+1)-\bar{c}_p)(y_p(j+1)-\bar{c}_p)}{\left[\left\{ \sum_{i=-k}^k (x_p(i+1)-\bar{c}_p)^2 \right\} \left\{ \sum_{i=-k}^k (y_p(j+1)-\bar{c}_p)^2 \right\} \right]^{1/2}}$$

ここで、(2k+1) は比較すべき segment (アミノ酸配列の断片) の長さ — (2k+1) の長さの "frame" 或いは、"window" — であり、 \bar{c}_p はパラメータ p の 20 個のアミノ酸上での平均値である。ここで、S/N 比を良くするため相関関数の平均を採る；

$$\bar{A}(\tau) = \frac{1}{n} \sum_{p=1}^n A_p(\tau), \quad \bar{C}(j) = \frac{1}{n} \sum_{p=1}^n C_p(j), \quad \bar{C}(i,j) = \frac{1}{n} \sum_{p=1}^n C_p(i,j)$$

ここで、 n はパラメターの個数である。上式は算術平均なので n 種のパラメターは互いに独立であることが望ましい。そうでなければ、たがいに相関するパラメターで求められた相関関数に weight がかかるからである。この問題を解決するために収集されたパラメター群に因子分析の手法を用いて以下のようなパラメター・セットを選び出した^{6, 10)};

- (1) 偏比容¹¹⁾、(2) 折り返しを形成する傾向度¹²⁾、(3) α -アミノグループの pK 値¹³⁾、(4) 極性度¹⁴⁾、(5) 相対的な突然変異度¹⁵⁾、(6) α -カルボキシルグループの pK 値¹³⁾

従って、 $n=6$ である。なお、 $C_p(i,j)$ を計算するための "window" の長さは 11 が適切である、即ち、 $(2k+1)=11$ 。

これら 6 つのパラメターからなるセットを使えば良好な構造 homology が得られる (即ち、高い相関——0.4 以上——が得られる領域に高い立体構造上の homology が認められる)。

folding type や細胞内外分泌予測^{7, 8)}

蛋白質のアミノ酸組成は 20 次元空間の 1 点として表される。従って、蛋白質の集団はこの 20 次元の組成空間にある分布を呈するであろう。そこで細胞の内外など生体内での位置が既知のものや、folding type の既知の蛋白質を収集しその組成空間内においていくつかのグループに分けることを試みる。即ち、主アークのノミノ酸組成の平均を原点にとり、各座標軸を適当に規格化しておく。まず、第一ステップとしてある半径内の蛋白質を中心グループとして集める。次に、この中心グループ外の蛋白質を reference として一つ採り、この reference point の方向から 60 度の立体角に入る蛋白質を求め、そして最も高い密度の方向を、reference とする蛋白質をシフトして行くことで探し求める。この手順を集められる蛋白質が 10 個以下になるまで次々で行う。こうして 13 個に分類されたグループにはそれぞれ生体内の位置 (細胞の内外) や生物学的機能 (enzyme や nonenzyme) や folding type に強い相関があることが判明した。

Homology 法による 2 次構造の予測⁹⁾

この手法の基本的アイデアは、配列の似た (homologous な) 蛋白質は類似のコンフォメーションを採るであろうと言う単純な発想から出ている。配列が既知の或る蛋白質を立体構造が既知のいくつかの蛋白質と比較し、これと homologous な配列の 2 次構造 (α 、 β 、Coil) を調べ、その個数 (n_α 、 n_β 、 n_c) を数え上げる。即ち $\bar{C}(i,j)$ が 0.3 以上のものが 8 残基以上続いた時、その残基のペア (i, j) を homologous とする。ここで、2 種類の weighting factor ν と W を導入する。即ち、 $0.3 \leq \bar{C} < 0.4$ の時 $\nu = 1$ 、 $0.4 \leq \bar{C} < 0.5$ の時 $\nu = 2$ 、等々。 W は 3 つのコンフォメーション状態 α 、 β 、coil の存在比率を考慮した factor である。こうして予測に使われるパラメター q として

$$q_k = W_k \sum_{j=1}^{n_k} \nu_j \quad (k = \alpha, \beta \text{ or coil})$$

但し、 $W_n=1.3$ 、 $W_s=1.4$ 、 $W_c=1.0$ が採用される。

この予測法による精度は約60%程度であり⁹⁾Chou-Fasman¹⁷⁾やRobsonら¹⁸⁾による方法よりも良い値となっている¹⁶⁾。

4. 解析例

ここでは、本システムによる3つの代表的な蛋白質の解析例を示したい。まず、自己相関関数を羊のケラチンB2A¹⁹⁾に適用し配列の繰返し性の有無を見てみる。次に、ヒトのインターロイキンil-2²⁰⁾に対してホモロジー法による2次構造の予測を行う。最後に、ペニシロペアシン²¹⁾、エンドチアペアシン²²⁾及び血圧調節に関与すると言われているヒトのレニン²³⁾についてホモロジーの有無を検出する。レニンを除いてこの2つの立体構造は知られている。図2はケラチンが5残基ごとのはっきりした周期性を有していることがわかる(左下のウインドウ)。50残基のラグ付近から不規則な領域が続き、73残基から再び周期性が表れる。この蛋白質のパワースペクトルを右下のウインドウに示す。

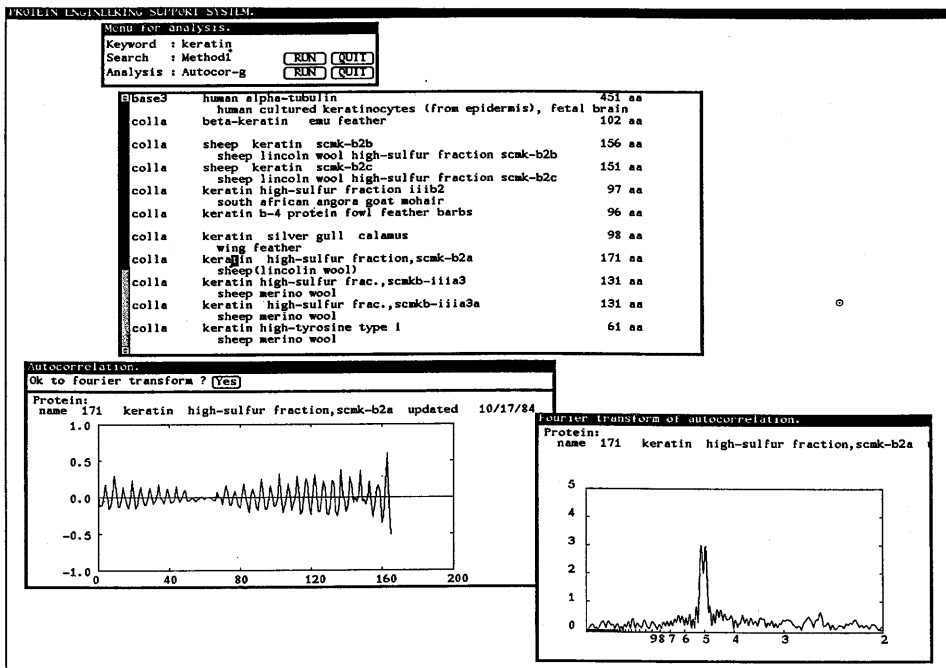


図2 自己相関関数によるケラチンの解析

図3はヒトのインターロイキンに対する結果である;この蛋白質のfolding typeは α タイプであることが予測されている(左下のウインドウ)。この上のウインドウにはこの予測の信頼率が表にされている。左上のウインドウはホモロジー法による2次構造予測の結果をシンボライズして表現されている。参考のために右側2枚のウインドウにRobson及びChou-Fasman法による予測も表示した。

PROTEIN ENGINEERING SUPPORT SYSTEM

Prediction of secondary structure
Homology method (by K.Nishikawa)
name 153 human interleukin 2 (il-2)
updated 01/07/85

Page No. 1/2 next previous

1 13 17 32 48
myrmqllscialslalvtinsaptssstkkqlqlehlldlqmlnginn
51 55 61 65 69 81 84
yknpklt r m t f k f y m p k k a t e l k h l q e l e e l k p l e e v i n l a q s k n f h l
101 104 111 116 121 125 131
r p r d i l s n i n v i l e l k g s e t t f m e y a d e t a t i v e f l n r w i t f o q s i i s

Prediction of secondary structure
Robson's method
name 153 human interleukin 2 (il-2)
updated 01/07/85

Page No. 1/2 next previous

1 9 18 27 30 42 45 50
myrmqllscialslalvtinsaptssstkkqlqlehlldlqmlnginn
51 54 59 65 69
yknpklt r m t f k f y m p k k a t e l k h l q e l e e l k p l e e v i n l a q s k n f h l
101 104 109 115 118 130 147
r p r d i l s n i n v i l e l k g s e t t f m e y a d e t a t i v e f l n r w i t f o q s i i s

Reliability.
The three-dimensional structures are known.
Average ratio (%) of fitness: 69.63%

	predicted					
	a	b	a/b	a+b	irr	ratio (%)
a	27	1	2	1	0	37.10
b	3	22	5	4	0	64.71
a/b	2	3	33	1	0	84.61
a+b	4	3	10	10	0	37.04
irr	0	0	0	2	2	50.00

Prediction of folding type.
Protein: name 153 human interleukin 2 (il-2) updated 01/07/85
Predicted as: alpha

Prediction of secondary structure
Chou-Fasman's method
name 153 human interleukin 2 (il-2)
updated 01/07/85

Page No. 1/2 next previous

1 18 22 25 41 48
myrmqllscialslalvtinsaptssstkkqlqlehlldlqmlnginn
51 54 59 65 69
yknpklt r m t f k f y m p k k a t e l k h l q e l e e l k p l e e v i n l a q s k n f h l
101 105 109 114 117 120 125 130 141
r p r d i l s n i n v i l e l k g s e t t f m e y a d e t a t i v e f l n r w i t f o q s i i s

図 3 インターロイキン il-2 の folding type と 2 次構造の予測

PROTEIN ENGINEERING SUPPORT SYSTEM

Menu for analysis.
Keyword : penicillopepsin
Search : Method1
Analysis : Homology

ec34 penicillopepsin ec 3.4.23.7
penicillium janthinellum

Homology.
abscissa: name 318 acid proteinase (endothiaepsin)
ordinate: name 323 penicillopepsin ec 3.4.23.7

PROTEIN NAME:
acid proteinase [e.c.3.4.23.7], penicillopepsin

CHAIN NUMBER: 1
ALPHA HELIX: 6
BETA SHEET: 32
INTERACTION:

VIEW alpha
 all atom
 wire
 ribbon
ROTATION reset
 x: 0 reset
 y: 0 reset
 z: 0 reset
 all reset
PARK VIEW
 alpha helix
 beta sheet
 from:
to:
 reset
ZOOM on off
 PICK UP RESIDUE

Homology.
abscissa: name 406 human renin precursor ec 3.4.23.15
ordinate: name 323 penicillopepsin ec 3.4.23.7

図 4 penicillopepsin の、endothiaepsin と renin に対する homology

図 4 ではペニシロペアシンに対して、エンドチアペアシンとレニンとのホモロジーを見た。エンドチアペアシンに対しては対角線上に 0.6 以上の $\bar{C}(i, j)$

j)が連続して見られ(左上のウィンドウ)これらの蛋白質の立体構造が似ていることが推測される。一方レニンに対してはそれ程のホモロジーは認められない(右下のウィンドウ)。右上のウィンドウにはペニシロペプシンの立体構造が3Dグラフィクスにより示されている。

5. 結語

本システムによる一次構造解析の基本アルゴリズムは相関関数にもとずいている。この関数はいくつかの物理化学的パラメータで計算された相関値の算術平均をとることでS/N比を高めることに成功している。

更に、本システムは良好なマン-マシン インターフェイスが得られるように、マルチ ウィンドウによる複数の解析結果を画面に表示するなど蛋白質工学研究者の新たな発想を促すのに役立てられると思われる。またマウスによる操作性の向上もあわせて計られている。

なお、本システムは蛋白質の立体構造予測に向けて引き続き発展中である。

参考文献

- (1) S.Kuhara, et al: *Nucleic Acids Research*, 12, 89 (1984)
- (2) 久原 哲、他: 情報学基礎研資料 86-3、情報処理学会 (1986)
- (3) F.C.Bernstein, et al: *J.Mol.Biol.*, 112, 535 (1977)
- (4) 電子通信学会編: "パターン情報処理", pp134-135, コロナ社 (1983)
- (5) Y.Kubota, et al: *J.Theor.Biol.*, 91, 347 (1981)
- (6) Y.Kubota: *Bull.Inst.Chem.Res., Kyoto Univ.*, 60, 309 (1982)
- (7) K.Nishikawa, Y.Kubota and T.Ooi: *J.Biochem.*, 94, 981 (1983)
- (8) K.Nishikawa, Y.Kubota and T.Ooi: *J.Biochem.*, 94, 997 (1983)
- (9) K.Nishikawa and T.Ooi: *Biochim.Biophys.Acta*, 871, 45 (1986)
- (10) Y.Kubota, et al: *Biochim.Biophys.Acta*, 701, 242 (1982)
- (11) E.J.Cohn and J.T.Edsall: "Proteins, Amino Acids, and Peptides" Van Nostrand-Reinhold, Princeton, New Jersey (1943)
- (12) M.Levitt: *Biochemistry*, 17, 4277 (1978)
- (13) H.A.Sober Ed.: "Handbook of Biochemistry, Selected Data for Molecular Biology" 2nd ed., The Chemical Robber Co., Cleveland, Ohio (1970)
- (14) R.Grantham: *Science*, 185, 862 (1974)
- (15) M.O.Dayhoff Ed.: "Atlas of Protein Sequence and Structure" Vol.5, Suppl.3, National Biomedical Research Foundation, Washington, D.C. (1978)
- (16) K.Nishikawa: *Biochim.Biophys.Acta*, 748, 285 (1983)
- (17) P.Y.Chou and G.D.Fasman: *Adv.Enzymol.*, 47, 45 (1978)
- (18) J.Garnier, D.J.Osguthorpe and B.Robson: *J.Mol.Biol.*, 120, 97 (1978)
- (19) T.C.Elleman: *Biochem.J.*, 130, 833 (1972)
- (20) T.Taniguti, et al: *Nature*, 302, 305 (1983)
- (21) J.Tang, et al: *Nature*, 271, 618 (1978)
- (22) T.L.Blundell, B.L.Sibanda and L.Pearl: *Nature*, 304, 273 (1983)
- (23) T.Imai, et al: *Proc.Natl.Acad.Sci.USA*, 80, 7405 (1983)