

A Protein Map and Its Application

Stephen S.-T. Yau,^{1,2} Chenglong Yu,³ and Rong He⁴

Graphical representation of gene sequences provides a simple way of viewing, sorting, and comparing various gene structures. Here we first report a two-dimensional graphical representation for protein sequences. With this method, we constructed the moment vectors for protein sequences, and mathematically proved that the correspondence between moment vectors and protein sequences is one-to-one. Therefore, each protein sequence can be represented as a point in a map, which we call protein map, and cluster analysis can be used for comparison between the points. Sixty-six proteins from five protein families were analyzed using this method. Our data showed that for proteins in the same family, their corresponding points in the map are close to each other. We also illustrate the efficiency of this approach by performing an extensive cluster analysis of the protein kinase C family. These results indicate that this protein map could be used to mathematically specify the similarity of two proteins and predict properties of an unknown protein based on its amino acid sequence.

Introduction

MANY METHODS HAVE BEEN REPORTED to analyze the huge amounts of gene data. One of them is the graphical representation of gene sequences, which is a very powerful tool for visual comparison of gene sequences. Hamori (1985) first used a three-dimensional H curve to represent a gene sequence. Gates (1985) later published a two-dimensional graphical representation that is simpler than the H curve. However, Gates' graphical representation has high degeneracy. We reported previously a new two-dimensional graphical representation of gene sequences (Yau *et al.*, 2003), which has no circuit or degeneracy, so the correspondence between gene sequences and gene graphs is one-to-one. Lately, many graphical representation methods for gene sequences have been proposed (Radic *et al.*, 2003a, 2003b); however, the method to make a protein sequence graph has never been shown. Unlike dealing with a gene or DNA sequence, from only four nucleotides, dealing with a protein sequence, from 20 amino acids, is more complicated. Here we report that a protein or amino acid sequence can be graphically represented and a universal protein map can be generated. This protein map can be used to predict the properties of proteins whose functions are not yet determined. We have analyzed 66 proteins from 5 protein families and an exhaustive set of 127 proteins from the protein kinase C (PKC) family using this method, and found it to be a useful predictive tool.

Protein Sequence Graphical Representation

Following our previous work (Yau *et al.*, 2003), we construct a protein sequence graph on two quadrants of the Cartesian coordinate system. The vectors corresponding to the 20 amino acids are lying in the line segment whose x -coordinate value is 1 and whose y -coordinate values are between -1 and 1 . The y -coordinate values of the 20 amino acid vectors are all distinct. The ordering of these y -coordinate values is based on amino acid hydrophobicity scale values (Fauchere and Pliska, 1983) because amino acid hydrophobicity plays an important role in protein folding. The 12 amino acids with positive hydrophobicity scale values were assigned in the first quadrant, and the difference of y -coordinate values between two amino acids next to each other is $1/13$. Because the hydrophobicity scale value of Gly is zero, its y -coordinate value was assigned to be zero. The other seven amino acids with negative hydrophobicity scale values were assigned in the fourth quadrant, and the difference of y -coordinate values between two amino acids next to each other is $1/8$. Thus, all y -coordinates of 20 amino acids are less than 1 and more than -1 . The y -coordinates for 20 amino acids are listed in Table 1, and 20 vectors are shown in Figure 1. The points in the graphical representation are obtained by the sum of vectors representing amino acids in the sequence. In Figure 2, we give the graphical representation of the first 10 vectors of human beta-globin amino acid sequence on the vector system shown in Figure 1, and the graphical representation of the whole

¹Department of Mathematics, Statistics, and Computer Science, University of Illinois at Chicago, Chicago, Illinois.

²Institute of Mathematics, East China Normal University, Shanghai, China.

³Department of Mathematics, The Chinese University of Hong Kong, Shatin, Hong Kong.

⁴Department of Pharmacology, University of Illinois at Chicago, Chicago, Illinois.

TABLE 1. HYDROPHOBICITY SCALE VALUES AND γ -COORDINATE OF 20 AMINO ACIDS

Amino acid	Hydrophobicity scale	γ -Coordinate
Trp (W)	+ 2.25	12/13
Ile (I)	+ 1.80	11/13
Phe (F)	+ 1.79	10/13
Leu (L)	+ 1.70	9/13
Cys (C)	+ 1.54	8/13
Met (M)	+ 1.23	7/13
Val (V)	+ 1.22	6/13
Tyr (Y)	+ 0.96	5/13
Pro (P)	+ 0.72	4/13
Ala (A)	+ 0.31	3/13
Thr (T)	+ 0.26	2/13
His (H)	+ 0.13	1/13
Gly (G)	0	0
Ser (S)	- 0.04	- 1/8
Gln (Q)	- 0.22	- 2/8
Asn (N)	- 0.60	- 3/8
Glu (E)	- 0.64	- 4/8
Asp (D)	- 0.77	- 5/8
Lys (K)	- 0.99	- 6/8
Arg (R)	- 1.01	- 7/8

human beta-globin amino acid sequence based on the same vector system is also shown in Figure 3. This protein sequence graphical representation has no circuits or degeneracy, and the correspondence between the sequence and the graphical curve can be mathematically proven to be one-to-one, as follows:

To prove there is no circuit or degeneracy in our two-dimensional graphical representation, we assume that (1) the

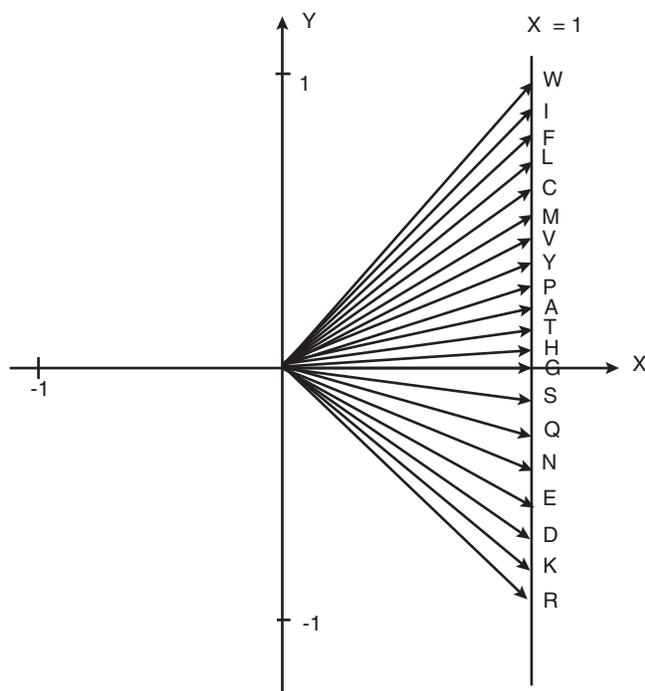


FIG. 1. Amino acid vector system based on Table 1.

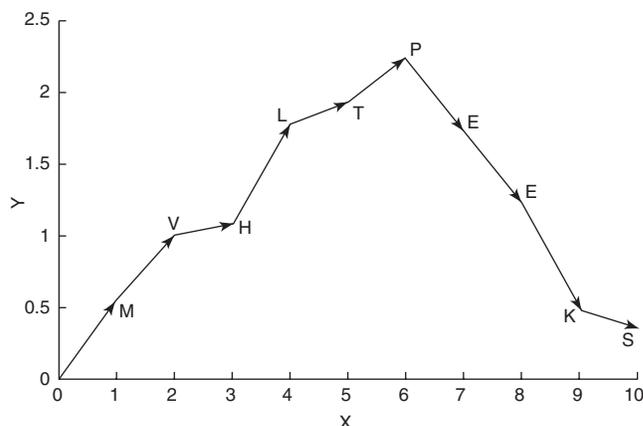


FIG. 2. Graphical representation of first 10 amino acids of human beta-globin sequence based on the vector system of Figure 1.

total number of amino acids is N ; (2) the number of amino acids R, K, D, E, N, Q, S, G, H, T, A, P, Y, V, M, C, L, F, I, and W is $r, k, d, e, n, q, s, g, h, t, a, p, y, v, m, c, l, f, i,$ and $w,$ respectively. Therefore, we have $r+k+d+e+n+q+s+g+h+t+a+p+y+v+m+c+l+f+i+w=N.$

If we further assume $rR, kK, dD, eE, nN, qQ, sS, gG, hH, tT, aA, pP, yY, vV, mM, cC, lL, fF, iI,$ and wW form a circuit, the following equation will hold:

$$\begin{aligned}
 & r\left(1, -\frac{7}{8}\right) + k\left(1, -\frac{6}{8}\right) + d\left(1, -\frac{5}{8}\right) + e\left(1, -\frac{4}{8}\right) \\
 & + n\left(1, -\frac{3}{8}\right) + q\left(1, -\frac{2}{8}\right) + s\left(1, -\frac{1}{8}\right) + g(1, 0) \\
 & + h\left(1, \frac{1}{13}\right) + t\left(1, \frac{2}{13}\right) + a\left(1, \frac{3}{13}\right) + p\left(1, \frac{4}{13}\right) \\
 & + y\left(1, \frac{5}{13}\right) + v\left(1, \frac{6}{13}\right) + m\left(1, \frac{7}{13}\right) + c\left(1, \frac{8}{13}\right) \\
 & + l\left(1, \frac{9}{13}\right) + f\left(1, \frac{10}{13}\right) + i\left(1, \frac{11}{13}\right) + w\left(1, \frac{12}{13}\right) = (0, 0)
 \end{aligned}$$

The sum of x -coordinates indicates that $r+k+d+e+n+q+s+g+h+t+a+p+y+v+m+c+l+f+i+w=0.$ It follows that $r=k=d=e=n=q=s=g=h=t=a=p=y=v=m=c=l=f=i=w=0$ as the number of amino acids is a nonnegative number. Therefore, no circuit exists in the graphical representation in a nontrivial case where $N > 0.$

A Moment Vector for Protein Sequences

Given the graphical curve of a protein sequence that can be represented by a sequence of points $(1, y_1), (2, y_2), \dots, (n, y_n),$ we can compute a sequence of numbers $1-y_1, 2-y_2, \dots, n-y_n.$ Conversely, if we know the sequence of numbers $1-y_1, 2-y_2, \dots, n-y_n,$ we can recover the graph $(1, y_1), (2, y_2), \dots, (n, y_n).$ Therefore, we would like to find a sequence of numbers, each of which uses the global information of the sequence of numbers $1-y_1, 2-y_2, \dots, n-y_n$ in such a way that this new sequence of numbers determines

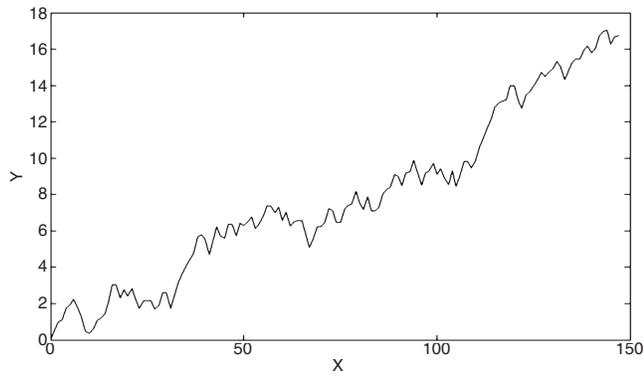


FIG. 3. Graphical representation of human beta-globin amino acid sequence based on the vector system of Figure 1.

and is determined by the sequence of numbers $1 - y_1, 2 - y_2, \dots, n - y_n$. For this purpose, we decided to use moments to characterize a protein graphical curve. The moments are defined as follows:

$$M_j = \sum_{i=1}^n \frac{(x_i - y_i)^j}{n^j}, \quad j = 1, 2, \dots, n,$$

where n is the number of amino acids contained in a protein sequence, and (x_i, y_i) represents the position of the i th amino acid in the protein graphical curve. According to this definition, each protein sequence has an n -dimensional moment vector (M_1, M_2, \dots, M_n) associated with it. We should emphasize that two or three moments will already characterize the protein sufficiently well (as it is demonstrated later) and that this is a much simpler representation than the graph with hundreds of points (147 for human beta globin).

Theorem: Consider the set of protein sequences having the fixed number (n) of amino acids. Then the correspondence between a protein sequence and its n -dimensional moment vector (M_1, M_2, \dots, M_n) is one-to-one.

Proof: We have demonstrated that the correspondence between a protein sequence and its graphical curve is one-to-one (Yau *et al.*, 2003). To prove the theorem, we will need to prove that the correspondence between a protein graphical curve and its moment vector is one-to-one.

By the definition, one protein sequence graph has an n -dimensional moment vector (M_1, M_2, \dots, M_n) . Hence, we need to demonstrate that from any given protein moment vector, we can recover the protein curve, which means all (x_i, y_i) ($i = 1, 2, \dots, n$) can be recovered from any given protein moment vector.

x_i is the x -coordinate value of i th amino acid on a protein graph. Based on our assignment, x_i should be equal to i . y_i is the y -coordinate value of i th amino acid on a protein graph. The next step is to obtain y_i from moment vector. Let $z_i = x_i - y_i$, then the moments can be simplified as:

$$M_j = \sum_{i=1}^n \frac{z_i^j}{n^j}, \quad j = 1, 2, \dots, n.$$

To solve for z_i , let $\delta_j = M_j n^j$, then the δ_j can be obtained by M_j and n . δ_j and z_i have the relation given below:

$$\begin{cases} \delta_1 = z_1 + z_2 + \dots + z_n \\ \delta_2 = z_1^2 + z_2^2 + \dots + z_n^2 \\ \dots \\ \delta_n = z_1^n + z_2^n + \dots + z_n^n \end{cases}$$

The z_1, z_2, \dots, z_n can be the roots of a symmetric polynomial $a_0 + a_1 z + a_2 z^2 + \dots + a_n z^n = (z - z_1)(z - z_2) \dots (z - z_n)$. By using Newton's identities (Jacobson, 1974):

$$\delta_d - s_1 \delta_{d-1} + \dots + (-1)^{d-1} s_{d-1} \delta_1 + (-1)^d d s_d = 0,$$

where $d = 1, 2, \dots, n$; s_d is the elementary symmetric polynomials in z_1, z_2, \dots, z_n ; a_i can be obtained by δ_j as shown below:

$$\begin{cases} a_n = 1 \\ a_{n-1} = (-1) \delta_1 \\ a_{n-2} = \frac{1}{2} (\delta_1^2 - \delta_2) \\ a_{n-3} = (-1)^3 \frac{1}{6} (\delta_1^3 - 3 \delta_1 \delta_2 + 2 \delta_3) \\ a_{n-4} = \frac{1}{24} (\delta_1^4 - 6 \delta_1^2 \delta_2 + 3 \delta_2^2 + 8 \delta_1 \delta_3 - 6 \delta_4) \\ \vdots \\ \vdots \end{cases}$$

As a result, the coefficients of the symmetric polynomial $a_0 + a_1 z + a_2 z^2 + \dots + a_n z^n = (z - z_1)(z - z_2) \dots (z - z_n)$ can be confirmed, and the set of all roots can be obtained. Next we need to identify each root z_1, z_2, \dots, z_n .

Because the y -coordinate values of all 20 amino acids are between -1 and 1 , $x_i - y_i$ is greater than zero. As we have defined that the position of k th amino acid on a graph is (x_k, y_k) or (k, y_k) , the position of $(k + 1)$ th amino acid on a graph (x_{k+1}, y_{k+1}) can be represented as $(k + 1, y_k + u_{k+1})$, where u_{k+1} may be any of y -coordinate value of these 20 amino acids. Thus, $z_{k+1} = x_{k+1} - y_{k+1} = (k + 1) - (y_k + u_{k+1}) = (k - y_k) + (1 - u_{k+1}) > (k - y_k)$. Because $z_k = x_k - y_k = k - y_k$, z_{k+1} must be greater than z_k . As a consequence, z_i is strictly increasing and each root can be identified by this property, which means each value of y_i can be obtained. With all (x_i, y_i) , a protein graph can be recovered.

Therefore, we have successfully proven that the correspondence between a protein sequence and its moment vector obtained from its sequence graph is one-to-one.

Protein Map and Cluster Analysis

By using moments of our protein graphical curve, we change beta-globin sequences of human, gorilla, cod, duck, chicken, and tortoise (the length of these sequences are 147) into 147-dimensional moment vectors. In Table 2, we give the distances between human and the other five species for some different dimensional moment vectors. According to

TABLE 2. DISTANCES BETWEEN BETA-GLOBIN SEQUENCES OF HUMAN AND OTHER FIVE SPECIES FOR DIFFERENT DIMENSIONAL MOMENT VECTORS

	Gorilla	Cod	Duck	Chicken	Tortoise
2-dim	0.0667	1.9216	4.6743	4.4788	3.5836
3-dim	0.0927	2.3493	5.5714	5.2891	4.3612
4-dim	0.1141	2.6188	6.1769	5.8255	4.9223
5-dim	0.1314	2.7961	6.6031	6.1984	5.3410
10-dim	0.1761	3.1439	7.5224	6.9921	6.3494
100-dim	0.1947	3.2472	7.8132	7.2425	6.7329
147-dim	0.1947	3.2472	7.8132	7.2425	6.7329

this table, we find that the first two or three moments are most important because when higher moments are included the relationship of being close or further away remains unchanged. For example, distances between human and gorilla are always the smallest for any dimensional moment vector.

Thus, we can use the first two components of the moment vector (M_1, M_2) of a protein sequence graph to represent a protein as a point in a two-dimensional space and generate a universal protein map. Using the distance between two points as an index for comparison, we can perform cluster analysis for protein sequences on this protein map. If two sequences are similar, the distance between two corresponding points should be small. Therefore, we may use this universal protein

map to predict properties of newly found proteins by performing clustering analysis.

Fifty beta-globin sequences of different species were extracted from Swiss-Prot (<http://au.expasy.org/>). Using the amino acid vectors shown in Figure 1, the sequence graphs of these beta-globins were obtained. We used our two-dimensional moment vector (M_1, M_2) system to characterize these 50 protein graphical curves and calculated the 50 points shown in Figure 4. From this figure, we note that these 50 beta-globins are separated into two main clusters. One cluster contains mammalian beta-globins, and the other contains beta-globins from avian, fish, and reptilian species. Because the chimpanzee beta-globin sequence is the same as the human beta-globin sequence, these two proteins have the same two-dimensional moment vector. For the same reason, dog and coyote beta-globin sequences have the same moment vector, as do black bear and polar bear beta-globin sequences. Figure 4 also shows that the distances between beta-globin sequences from several primatal species (human, grivet, gorilla, langur, gibbon, and chimpanzee) are very small, and they form a subcluster.

We have also developed a three-dimensional moment vector (M_1, M_2, M_3) system to characterize graphical curves of these 50 proteins using the first three moments of their protein graph. As a result, these 50 beta-globins can be represented as 50 points in a three-dimensional map. By computing distances between these points, we did cluster analysis. The data are shown in a hierarchical tree in Figure 5, which indicates the evolutionary relationships between these proteins.

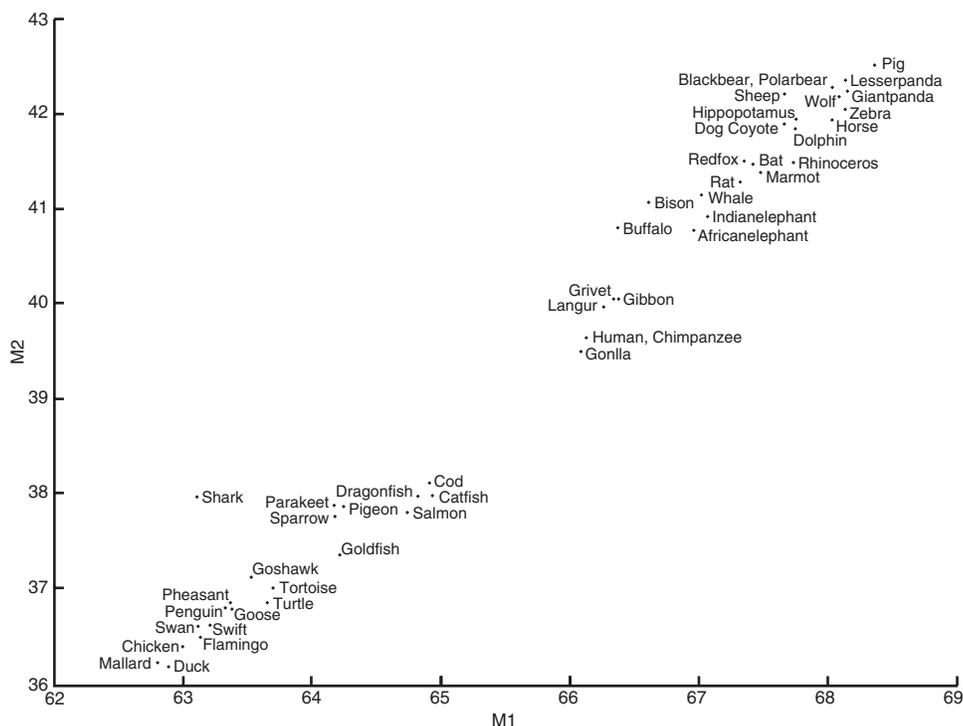


FIG. 4. Two-dimensional moment vector points of beta-globin sequences of 50 species.

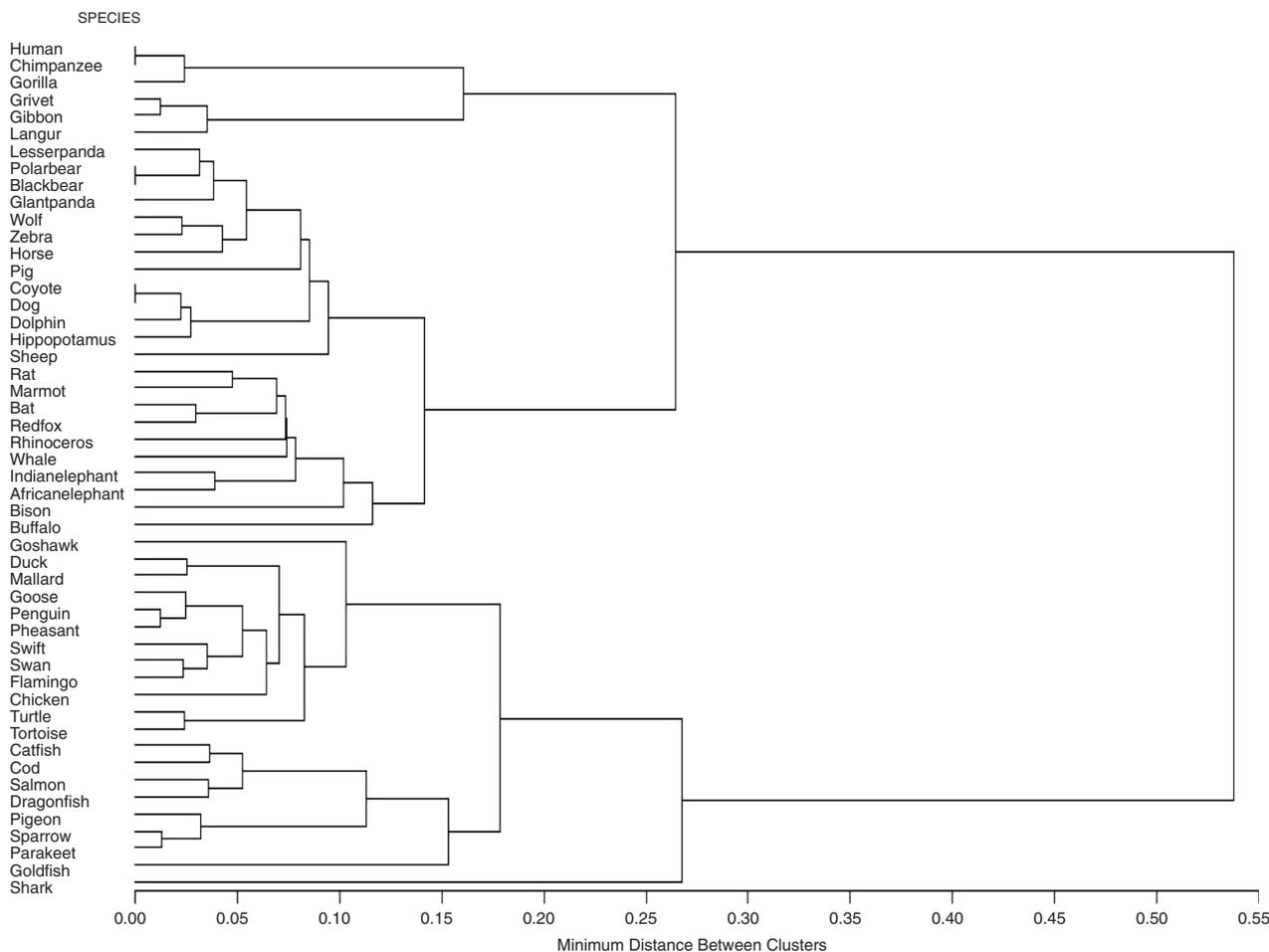


FIG. 5. Clustering hierarchical tree of beta-globin sequences of 50 species by three-dimensional moment vector.

To demonstrate the effectiveness of our clustering approach, we extracted another 16 proteins in 4 families (insulin family, catalase family, CCN family, and thiolase family) from Swiss-Prot. By using our two-dimensional moment vector system, these 16 proteins and the 50 beta-globins are represented as 66 points shown in Figure 6. The clusters of five families are identified in the figure.

We also applied our protein map to a set of 127 proteins from the PKC family (Table 3). PKCs can be divided into two functional parts, the regulatory and catalytic domains. The regulatory domain is the principal determinant of classification, as the catalytic domain tends to be highly conserved. Since PKCs were originally identified in mammals, their classification system is strongly characterized by varieties of PKCs found in mammals. In particular, mammalian PKCs are generally divided into three subfamilies: conventional PKCs (cPKCs: α , β I, β II, and γ), novel PKCs (nPKCs: θ , ϵ , δ , and η), and atypical PKCs (aPKCs: λ / ι and ζ). A controversial group of potential PKCs including PKC ν and PKC μ /PKD (protein kinase D) (mouse) (Webb *et al.*, 2000)

has similar regulatory domains to PKCs, but catalytic domains are more similar to the myosin light-chain kinase of *Dictostelium* (Hurley *et al.*, 1997). Fungi have PKC homologs that characteristically contain more residues than mammalian PKCs with significantly different regulatory domains but similar catalytic domains. There are also PKC-related kinases (PRKs) that are found in many animals and have features similar to fungal PKCs. Thus, we identified six categories of structural architecture for PKCs and PKC-related molecules: cPKC, nPKC, aPKC, PKC μ (ν , μ , and D2 types), PKC1 (fungal PKCs), and PRK. We used our two-dimensional protein map to characterize the regulatory domains of these 127 proteins from PKC family and calculated the 127 points in Figure 7. Four clusters, PKC, PKC μ , PKC1, and PRK, can be seen in this figure. To study the PKC cluster clearly, we expanded it and obtained Figure 8. In this figure, we classified these points into three categories (cPKC, nPKC, and aPKC). In Table 3, we list our clustering results and made a comparison with the categories established in the literature.

TABLE 3. CLUSTERING RESULTS OF 127 PROTEINS FROM PKC FAMILY BY USING PROTEIN MAP

Number	Species	Sequence ID	Category	Architecture
1	<i>Gallus gallus</i>	NP001006133	delta	nPKC
2	<i>Canis familiaris</i>	NP001008716	delta	nPKC
3	<i>Xenopus tropicalis</i>	NP001012707	iota	aPKC
4	<i>Hydra vulgaris</i>	O01715	cPKC	cPKC
5	<i>Oryctolagus cuniculus</i>	O19111	zeta	aPKC
6	<i>Cochliobolus heterostrophus</i>	O42632	PKC1	PKC1
7	<i>Sycon raphanus</i>	O61224	cPKC	cPKC
8	<i>Sycon raphanus</i>	O61225	nPKC	nPKC
9	<i>Suberites domuncula</i>	O62567	cPKC	cPKC
10	<i>Suberites domuncula</i>	O62569	nPKC	nPKC
11	<i>Calliphora vicina</i>	O76850	cPKC	cPKC
12	<i>Homo sapiens</i>	O94806	mu	PKCmu
13	<i>Rhabdocalypus dawsoni</i>	O96942	cPKC	cPKC
14	<i>Geodia cydonium</i>	O96997	cPKC	cPKC
15	<i>Bos taurus</i>	P04409	alpha	cPKC
16	<i>Bos taurus</i>	P05126	beta	cPKC
17	<i>Homo sapiens</i>	P05129	gamma	cPKC
18	<i>Drosophila melanogaster</i>	P05130	cPKC	cPKC
19	<i>Rattus norvegicus</i>	P05696	alpha	cPKC
20	<i>Homo sapiens</i>	P05771	beta	cPKC
21	<i>Oryctolagus cuniculus</i>	P05772	beta	cPKC
22	<i>Rattus norvegicus</i>	P09215	delta	nPKC
23	<i>Rattus norvegicus</i>	P09216	epsilon	nPKC
24	<i>Rattus norvegicus</i>	P09217	zeta	aPKC
25	<i>Oryctolagus cuniculus</i>	P10102	alpha	cPKC
26	<i>Oryctolagus cuniculus</i>	P10829	gamma	cPKC
27	<i>Oryctolagus cuniculus</i>	P10830	epsilon	nPKC
28	<i>Drosophila melanogaster</i>	P13677	cPKC	cPKC
29	<i>Mus musculus</i>	P16054	epsilon	nPKC
30	<i>Homo sapiens</i>	P17252	alpha	cPKC
31	<i>Mus musculus</i>	P20444	alpha	cPKC
32	<i>Mus musculus</i>	P23298	eta	nPKC
33	<i>Saccharomyces cerevisiae</i>	P24583	PKC1	PKC1
34	<i>Homo sapiens</i>	P24723	eta	nPKC
35	<i>Mus musculus</i>	P28867	delta	nPKC
36	<i>Caenorhabditis elegans</i>	P34885	nPKC	nPKC
37	<i>Schizosaccharomyces pombe</i>	P36582	PKC1	PKC1
38	<i>Schizosaccharomyces pombe</i>	P36583	PKC1	PKC1
39	<i>Homo sapiens</i>	P41743	iota	aPKC
40	<i>Candida albicans</i>	P43057	PKC1	PKC1
41	<i>Mus musculus</i>	P63318	gamma	cPKC
42	<i>Rattus norvegicus</i>	P68403	beta	cPKC
43	<i>Mus musculus</i>	P68404	beta	cPKC
44	<i>Neurospora crassa</i>	P87253	PKC1	PKC1
45	<i>Caenorhabditis elegans</i>	P90980	cPKC	cPKC
46	<i>Aspergillus niger</i>	Q00078	PKC1	PKC1
47	<i>Mus musculus</i>	Q02111	theta	nPKC
48	<i>Homo sapiens</i>	Q02156	epsilon	nPKC
49	<i>Mus musculus</i>	Q02956	zeta	aPKC
50	<i>Homo sapiens</i>	Q04759	theta	nPKC
51	<i>Homo sapiens</i>	Q05513	zeta	aPKC
52	<i>Homo sapiens</i>	Q05655	delta	nPKC
53	<i>Homo sapiens</i>	Q15139	mu	PKCmu
54	<i>Aplysia californica</i>	Q16974	cPKC	cPKC
55	<i>Aplysia californica</i>	Q16975	nPKC	nPKC
56	<i>Caenorhabditis elegans</i>	Q19266	aPKC	aPKC
57	<i>Lytechinus pictus</i>	Q25378	cPKC	cPKC
58	<i>Xenopus tropicalis</i>	Q28EN9	iota	aPKC
59	<i>Mus musculus</i>	Q2NKI4	cPKC	cPKC
60	<i>Aspergillus oryzae</i>	Q2U6A7	PKC1	PKC1
61	<i>Mus musculus</i>	Q3UEA6	PRK	PRK
62	<i>Xenopus laevis</i>	Q498G7	nPKC	nPKC
63	<i>Bombyx mori</i>	Q4AED5	aPKC	aPKC

(continued)

TABLE 3. (CONTINUED)

Number	Species	Sequence ID	Category	Architecture
64	<i>Bombyx mori</i>	Q4AED6	cPKC	cPKC
65	<i>Macaca fascicularis</i>	Q4R4U2	cPKC	cPKC
66	<i>Pongo pygmaeus</i>	Q5R4K9	aPKC	aPKC
67	<i>Danio rerio</i>	Q5TZD4	nPKC	nPKC
68	<i>Mus musculus</i>	Q62074	iota	aPKC
69	<i>Mus musculus</i>	Q62101	mu	PKCmu
70	<i>Rattus norvegicus</i>	Q64617	eta	nPKC
71	<i>Schistosoma mansoni</i>	Q69G16	cPKC	cPKC
72	<i>Xenopus laevis</i>	Q6AZF7	cPKC	cPKC
73	<i>Debaryomyces hansenii</i> CBS767	Q6BI27	PKC1	PKC1
74	<i>Yarrowia lipolytica</i>	Q6C292	PKC1	PKC1
75	<i>Xenopus laevis</i>	Q6DCJ8	delta1	nPKC
76	<i>Rattus norvegicus</i>	Q6DUV1	epsilon	nPKC
77	<i>Candida glabrata</i> CBS138	Q6FJ43	PKC1	PKC1
78	<i>Xenopus laevis</i>	Q6GNZ7	nPKC	nPKC
79	<i>Homo sapiens</i>	Q6P5Z2	PRK	PRK
80	<i>Rattus norvegicus</i>	Q6P748	PRK	PRK
81	<i>Cryptococcus neoformans</i> var.	Q6UB96	PKC1	PKC1
82	<i>Cryptococcus neoformans</i> var.	Q6UB97	PKC1	PKC1
83	<i>Eremothecium gossypii</i>	Q75BT0	PKC1	PKC1
84	<i>Aspergillus nidulans</i>	Q76G54	PKC1	PKC1
85	<i>Xenopus laevis</i>	Q7LZQ8	cPKC	cPKC
86	<i>Xenopus laevis</i>	Q7LZQ9	cPKC	cPKC
87	<i>Anopheles gambiae</i>	Q7QCP8	nPKC	nPKC
88	<i>Danio rerio</i>	Q7SY24	cPKC	cPKC
89	<i>Xenopus laevis</i>	Q7SZH7	delta2	nPKC
90	<i>Xenopus laevis</i>	Q7SZH8	delta1	nPKC
91	<i>Danio rerio</i>	Q7T2C5	cPKC	cPKC
92	<i>Homo sapiens</i>	Q86XJ6	delta	nPKC
93	<i>Pichia pastoris</i>	Q86ZV2	PKC1	PKC1
94	<i>Leptosphaeria maculans</i>	Q873Y9	PKC1	PKC1
95	<i>Homo sapiens</i>	Q8IUU5	PRK	PRK
96	<i>Kluyveromyces lactis</i>	Q8J213	PKC1	PKC1
97	<i>Takifugu rubripes</i>	Q8JFZ9	cPKC	cPKC
98	<i>Mus musculus</i>	Q8K1Y2	mu	PKCmu
99	<i>Mus musculus</i>	Q8K2K8	eta	nPKC
100	<i>Limulus polyphemus</i>	Q8MXB6	nPKC	nPKC
101	<i>Homo sapiens</i>	Q8NE03	eta	nPKC
102	<i>Danio rerio</i>	Q90XF2	iota	aPKC
103	<i>Xenopus laevis</i>	Q91569	iota	aPKC
104	<i>Xenopus</i> sp.	Q91948	PRK	PRK
105	<i>Hypocrea jecorina</i>	Q99014	PKC1	PKC1
106	<i>Homo sapiens</i>	Q9BZL6	D2	PKCmu
107	<i>Drosophila melanogaster</i>	Q9GSZ3	aPKC	aPKC
108	<i>Blumeria graminis</i>	Q9HF10	PKC1	PKC1
109	<i>Tuber borchii</i>	Q9HGK8	PKC1	PKC1
110	<i>Botryotinia fuckeliana</i>	Q9UVJ5	PKC1	PKC1
111	<i>Sporothrix schenckii</i>	Q9Y792	PKC1	PKC1
112	<i>Magnaporthe grisea</i>	Q9Y7C1	PKC1	PKC1
113	<i>Apis mellifera</i>	XM391874	cPKC	cPKC
114	<i>Rattus norvegicus</i>	XP001066028	theta	nPKC
115	<i>Macaca mulatta</i>	XP001116804	gamma	cPKC
116	<i>Pan troglodytes</i>	XP001147999	theta	nPKC
117	<i>Bos taurus</i>	XP001250401	delta	nPKC
118	<i>Rattus norvegicus</i>	XP234108	PKCmu	PKCmu
119	<i>Gallus gallus</i>	XP421417	eta	nPKC
120	<i>Canis familiaris</i>	XP540151	PKCmu	PKCmu
121	<i>Canis familiaris</i>	XP541432	gamma	cPKC
122	<i>Bos taurus</i>	XP583587	epsilon	nPKC
123	<i>Bos taurus</i>	XP602125	gamma	cPKC
124	<i>Danio rerio</i>	XP683138	nPKC	cPKC
125	<i>Canis familiaris</i>	XP849292	theta	nPKC
126	<i>Canis familiaris</i>	XP851386	PKCmu	PKCmu
127	<i>Canis familiaris</i>	XP851861	epsilon	nPKC

The number column provides identifying numbers used in Figures 7 and 8. Sequence IDs refer to NCBI or SwissProt accession numbers. Categories are those established in the literature, whereas Architectures are our clustering results.

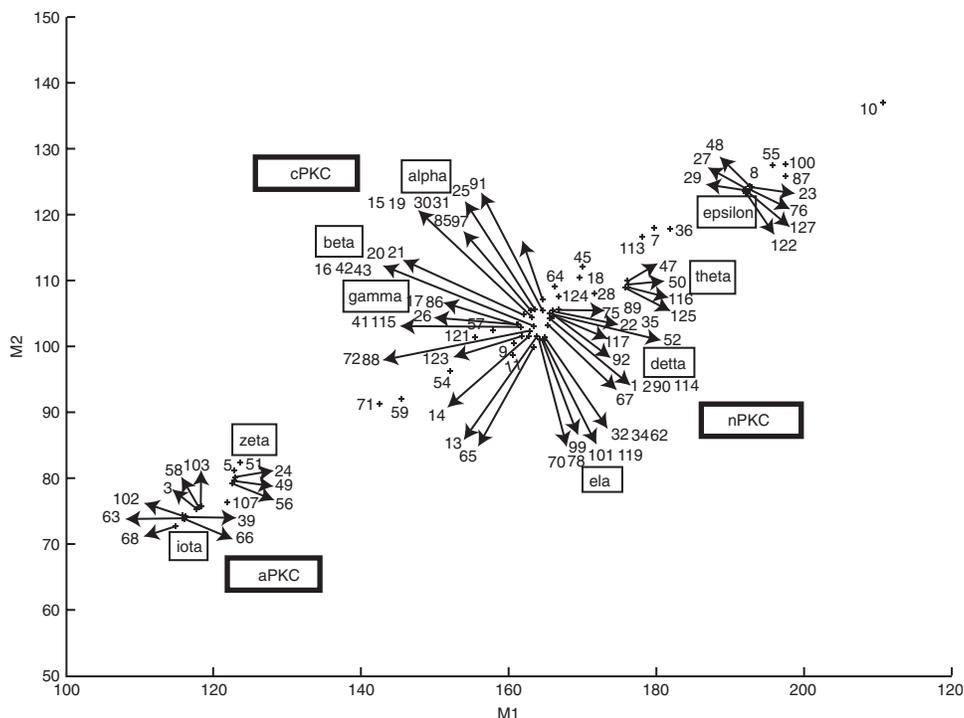


FIG. 8. Two-dimensional moment vector points of regulatory domains of proteins from PKC cluster in Figure 7.

specific yeast-like growth patterns and filamentous fungi that display mycelial growth patterns. This distinction is supported by the evidence of fungal PKCs involvement in control of growth patterns (Aquino-Piñero and Valle, 2002). We believe that our protein map provides a rapid and accurate way of identifying the likely category of a PKC family member.

Discussion and Conclusion

In this paper, we report a two-dimensional graphical representation for protein sequences. A moment vector system to represent a protein sequence is introduced, and the correspondence between a protein sequence and its moment vector is mathematically proven to be one-to-one. With this moment vector system, a protein can be represented as a point on a map, which is called protein map. Proteins with similar properties plot close together. Thus, the protein map allows us to analyze protein sequences using clustering methods and to predict properties of newly found proteins. This method will provide a new tool for protein functional studies.

To represent a protein sequence as a two-dimensional graph, we assigned 20 vectors to 20 amino acids. The *x*-coordinate value of a vector was chosen greater than zero to avoid circuit or degeneracy. The *y*-coordinate value of a vector was assigned based on amino acid hydrophobicity scale values. Because many other amino acid physicochemical properties, such as polarity (Grantham, 1974) and refractivity (Jones, 1975), should be considered, further studies will be needed to decide which combination of amino acid prop-

erties is most biologically meaningful for determining its *y*-coordinate value.

With our moment vector system, a protein sequence graph can be represented as a point in a two-dimensional or three-dimensional map depending on whether first two or first three moments of this sequence graph are used as its moment vector. By applying clustering methods to either of these maps, protein sequences, protein domains, and even arbitrary amino acid sequences can be efficiently analyzed.

Acknowledgment

We thank Dr. Max Benson for critically reading and editing the manuscript. We also thank Kareem Carr for providing us the data of PKC family.

References

Aquino-Piñero, E., and Valle, N.R.-D. (2002). Characterization of a protein kinase C gene in *Sporothrix schenckii* and its expression during the yeast-to-mycelium transition. *Med Mycol* **40**, 185–199.
 Fauchere, J., and Pliska, V. (1983). Hydrophobic parameters of amino-acid side-chains from the partitioning of N-acetyl-amino-acid amides. *Eur J Med Chem* **18**, 369–375.
 Gates, M.A. (1985). Simpler DNA sequence representations. *Nature* **316**, 219.
 Grantham, R. (1974). Amino acid difference formula to help explain protein evolution. *Science* **185**, 862–864.
 Hamori, E. (1985). Novel DNA sequence representation. *Nature* **314**, 585–586.
 Hurley, J.H., Newton, A.C., Parker, P.J., Blumberg, P.M., and Nishizuka, Y. (1997). Taxonomy and function of C1 protein kinase C homology domains. *Protein Sci* **6**, 477–480.

- Jacobson, N. (1974). *Basic Algebra*, vol. 1. India: Hindustan Publishing Corporation, p. 135.
- Jones, D.D. (1975). Amino acid properties and side-chain orientation in proteins: a cross correlation approach. *J Theor Biol* **50**, 167–184.
- Randic, M., Vracko, M., Lers, N., and Plavsic, D. (2003a). Novel 2-D graphical representation of DNA sequences and their numerical characterization. *Chem Phys Lett* **368**, 1–6.
- Randic, M., Vracko, M., Lers, N., and Plavsic, D. (2003b). Analysis of similarity/dissimilarity of DNA sequences based on novel 2-D graphical representation. *Chem Phys Lett* **371**, 202–207.
- Webb, B.L.J., Hirst, S.J., and Giembycz, M.A. (2000). Protein kinase C isoenzymes: a review of their structure, regulation and role in regulating airways smooth muscle tone and mitogenesis. *Br J Pharmacol* **130**, 1433–1452.
- Yau, S.S.-T., Wang, J., Niknejad, A., Lu, C., Jin, N., and Ho, Y. (2003). DNA sequence representation without degeneracy. *Nucleic Acids Res.* **31**, 3078–3080.

Address reprint requests to:
Stephen S.-T. Yau, Ph.D.

Distinguished Professor
Department of Mathematics, Statistics, and Computer Science
University of Illinois at Chicago
851 S. Morgan St.
Chicago, IL 60607-7045

Director
Institute of Mathematics
East China Normal University
Shanghai 200062
China

E-mail: yau@uic.edu

Received for publication August 23, 2007; received in revised form December 7, 2007; accepted December 24, 2007.

Copyright of DNA & Cell Biology is the property of Mary Ann Liebert, Inc. and its content may not be copied or emailed to multiple sites or posted to a listserv without the copyright holder's express written permission. However, users may print, download, or email articles for individual use.