

瀬戸保彦, 池内義典, 川北栄継(1), 西川建(1), 金久實(2)  
(財)蛋白質研究奨励会 (1)蛋白工学研 (2)京大化研  
Seto, Y., Ikeuchi, Y., Kawakita, S. (1), Nishikawa, K. (1), and Kanehisa M. (2)  
Protein Res. Found., (1)Protein Eng. Res. Inst., and (2)Kyoto Univ.

We propose a semi-automatic method to classify protein sequences into superfamilies. The method is a combination of the global sequence homology search and the local sequence identifier. The former is performed by a pairwise homology search method developed by Nishikawa et al(2), which jointed five different measures to improve sensitivity for detecting sequence similarity. The latter is done by applying the fragment peptide library (FRAP) compiled by Seto et al(3). We evaluate our method by applying it to 89 superfamilies in the PIR database.

Out of 89 superfamilies, 65 superfamilies are classified by the combination method. We find that 20 superfamilies are multi-superfamilies: Different superfamilies belong to a same group by our criteria. We find difficulty to classify 4 superfamilies.

## 1. 序

我々はタンパク質の特定のグループに特徴的な断片配列、モチーフの収集について第2回ゲノム情報ワークショップで報告した。蓄積した断片配列は現在12198件に及ぶ。自然界におけるタンパク質の種類は数千種と見積もる事もできる。つまり殆どのタンパク質に固有のモチーフがそろそろ揃ってきたと考えられる。また長年の間に収集してきたタンパク質の配列は13000件に及ぶ。これらの配列の分類を始める時期が来たと考え、数年前よりその方法論の検討を始めた。

分類する事は同一性の規範、アイデンティティを決める事と言える。ふたつのタンパク質が同族か否かは物理的、化学的、生物学的な性質等から総合的に判断される。ゲノムプロジェクトの進行に伴い生産される膨大な情報はまず配列であり、タンパク質としてのその他の性質は不明である。

タンパク質の分類はDayhoff, Barkerらにより進化的視点から系統的になされてきた(1)。しかし、分類の難しさから、配列データの増加に比較して分類作業は順調には進んでいない。また長い配列や多機能性タンパク質の配列が増加してきたため、分類も複雑になっていく。本報告において、我々はできる限り配列の情報だけから自動的に分類する事を試みた。その方法は全体的な類似性と局所の同一性を組み合わせるものである。分類分けのされているPIRの配列データベースを利用して、多くの生物種から得られたタンパク質を含むスーパーファミリーをモデルとしてその方法論を評価した。

## 2. 方法

分類は2段階で行う。第1段階で西川らの開発した5パラメータ法(2)により、配列全体と比較して、タンパク質を類似した配列を持つグループに分ける。そのグループで弱い類似性しか示さず同族判定のむずかしい場合、第2段階として瀬戸-金久らの断片ペプチドライブラリー (FLAP) (3) のモチーフを共通に持つタンパク質を同族とする。

5パラメータ法は5個のスコアを利用してタンパク質の類似性を評価する。まず比較するタンパク質の配列をGoad-Kanehisa法により、整列表記する。進化的に変わり易いアミノ酸であるかどうかを考慮し、つまりDayhoff等のPAM行列数値を使い比較してスコア1を得る。窪田らの方法により、アミノ酸の物理化学的性質を考慮してスコア3を得る。ここに得たスコア1,3について配列の長さを考慮したスコア2,4をそれぞれ得る。その整列配列で、アミノ酸残基を直接比較して一致度の程度によりスコア5を得る。

西川はこれら5個のスコアについて同族のタンパク質になるための域値を立体構造データベースPDBのデータを利用して決めた。ふたつのタンパク質を比較して域値を満たすスコアがN個のものをランクNの類似度とする。ランク5の類似度のタンパク質は同族として大体間違いない。ランクを4にする事により、殆どの同族タンパク質を検索できる。しかし同族でないタンパク質(ノイズ)が類似配列として検索される。従って、検索で拾われたこれらのノイズが実際に同族でない事を別の手法で判断する必要がある。

FLAPにあるモチーフはタンパク質を特徴づける短い配列である。そこで、第2段階とし

てランク4,5の関係にあるタンパク質のグループの中で、そのグループに特徴的なモチーフを持たないタンパク質は同族でないとした。

配列データベースNBRF/PIRでは進化的に同じ起源から生成したと考えられるタンパク質を同じ族にまとめて、スーパーファミリーと呼んでいる。Rel.25のデータにおいて総エントリー約7000件、スーパーファミリー約2000件である。PIRのスーパーファミリー89個を選んで、この2つの方法の組み合わせによる再分類を試みた。これらのスーパーファミリーは次の15種に分けられた生物の中、どれか3種以上の生物から得たタンパク質を含む。従って配列の違いが比較的大きいタンパク質のグループと言える。

Pr:prokaryota Fu:fungi Li:lichenes Pl:plantae Pz:protozoa  
Ne:nematoda Mo:mollusca Ar:arthropoda Ag:agnatha Ch:chondrichthyes  
Os:osteichthyes Am:amphibia Re:reptilia Av:aves Ma:mammalia

各スーパーファミリーについて上のいづれかの生物から得たタンパク質を一つずつ選びそれをプローブとしてPIR全データに対して類似配列検索を行った。例えばSUPERFAMILY NUMBER=958のACh receptorはそのスーパーファミリーにENTRY=10件のタンパク質がある。ソースとなった生物ごとのタンパク質のエントリー数はaves=2, chondrichthyes=4, mammalia=4で、このなかからACCHG1(Av), ACRYA1(Ch), ACHUA1(Ma)のコードを持つタンパク質がACh receptor族のプローブとして検索、分類の出発点となった。個々のスーパーファミリーのプローブの記載は省く。

### 3. 結果と考察

スーパーファミリー89個のタンパク質を対象にして5パラメータ法で検索した結果を表に示す。PROTEIN欄はスーパーファミリーの名前、ENTRY欄はそのスーパーファミリーに属するタンパク質の総数、MISSED欄はプローブタンパク質によって検索できなかったそのスーパーファミリーのタンパク質数、NOISE欄はプローブで検索されたそのスーパーファミリー以外のタンパク質総数を示す。MISSED欄が0のスーパーファミリーはENTRYの全てのタンパク質が検索されている。NOISE欄のmanyはスーパーファミリー以外のタンパク質が20以上数多く検索された事を意味する。MOTIF欄の配列の間の数字はふたつのモチーフの間にどれだけのアミノ酸残基があるかを特定の配列で計算したもので、同族であってもタンパク質により異なる事が多い。モチーフ間の距離の目安として記載してある。

対象となったスーパーファミリーは4グループに分ける事ができる。

グループI:同族タンパク質を5パラメータ法で全て検索できる。MISSED欄が0でないホルモン類とhistone H1は検索漏れがある。しかしタンパク質として考えた場合、それらが短い部分配列であることから分類の対象からははずす方が適当である。NOISE欄に示すようにスーパーファミリー以外のタンパク質が拾われるが、グループIにおいてはMOTIF欄に示すモチーフ配列がこれらのタンパク質にない事からノイズを落とす事ができる。

グループII:プローブ配列で検索されるノイズは少ないが、検索できない同族タンパク質がある。IIのノイズはモチーフの有無で除く事ができる。しかし、検索漏れの同族タンパク質を拾うためには5パラメータ法でランクを落とし更に弱い類似検索をする必要がある。検索にもれたものもその多くはモチーフを持っている。しかしこのIIのタンパク質には配列だけから同族とするのは困難なものがPIRにある事を示す。

グループIII:PIRでは別のスーパーファミリーに属しているタンパク質が我々の規範からは同族となるものを含むグループである。類似配列の検索で整理したときの同一アミノ酸が10から20パーセントという弱い類似度の場合、局所に同一のモチーフを持つか否かで同族にするかどうかを決める。従って、特定のグループのタンパク質に特徴的であるべきモチーフの質が分類の信頼度を決める。

グループIV:このグループの配列の分散度は大きく、従って弱い類似性を持った多くのタンパク質が検索される。同族性の判断が困難なグループである。これらのスーパーファミリーのモチーフがFRAPライブラリーにあるが、いづれもスーパーファミリー内の特定な族のモチーフであると考えられる。モチーフをスーパーファミリー全体に一般化するにはより洗練されたパターンにする事が必要である。

第一段階ではプローブから出発してペアワイズな配列の整理表記で一つ一つのタンパク質の同族性を判断した。グループIIのcytochrome c, globinのように多くの生物種の配列がある場合、プローブの選び方によって類似配列検索で拾えないタンパク質が多くなる。つまり、ある族内の配列の分散が大きいとABC問題が生ずる。A, B, Cがあったとき、 $A=B, B=C, A=C$ となったとき $A=C$ に変更するのか。A, Cがあり、AはCと同族でないがBが現れて、A, CをBがつないでA, B, Cを同族にする事はないのか。現在この問題は人間が判断している。cytochrome c, globinの場合、NOISEが少ない事から分かるように類似配列の検索にミスは少ない。プローブを増す事で検索漏れはなくなる。

SUPERFAMILIES CLASSIFIED BY NEAR HOMOLOGY SEARCH AND NOTIPS FOR EACH SUPERFAMILY  
I SUPERFAMILIES EASILY CLASSIFIED

SUPERFAMILY	PROTEIN	ENTRY	NOISE	NOTIP
HM ER			MISSED	
356:ACH receptor:		10	0	many
27:alcohol dehydrogenase:		12	0	4
585:antithrombin:		20	0	4
535:argininosuccinate lyase:		5	0	24
434:ATPase alpha/beta:		15	0	23
401:ATPase F:		12	0	many
400:ATPase lipid binding:		11	0	2
476:citrate synthase:		3	0	5
643:corticosteroids:		5	0	0
461:corticosteroids:		17	0	2
238:creatine kinase M:		5	0	3
905:crystallin alpha:		55	0	7
906:crystallin beta,gamma:		19	0	4
14:cytochrome b5:		9	0	0
93:cytochrome c oxidase 2:		12	0	many
100:cytochrome c oxidase 3:		3	0	23
122:elastin factor:		1	0	1
443:gastrin:		12	0	4
678:gastrin releasing peptide:		4	0	2
68:Gln dehydrogenase:		5	0	11
676:glucagon:		25	0	1
51:glycerol-3-phosphate dehydrogenase:		13	0	0
172:glycogen phosphorylase:		4	0	10
645:guadonin:		5	0	0
773:histone H1:		19	0	12
772:histone H2A:		25	0	1
773:histone H2B:		24	0	1
774:histone H3:		35	0	1
775:histone H4:		14	0	1
776:histone H5:		5	0	4
1015:lysozyme:		5	0	24
1016:lysozyme:		2	0	20
492:insulin:		43	0	0
31:lactate dehydrogenase:		13	0	10
959:lectin:		20	0	10
354:lysozyme:		26	0	0
1084:metallothionein:		18	0	19
43:NAD dehydrogenase 1:		1	0	many
45:NAD dehydrogenase 2:		7	0	many
46:NAD dehydrogenase 4:		7	0	many
47:NAD dehydrogenase 5:		7	0	many
679:pancreatic hormone:		13	0	0
399:papain:		9	0	1
493:pepsin:		18	0	5
235:phosphoglycerate kinase:		5	0	4
219:pyruvate kinase:		4	0	20
854:ribosomal protein L14:		14	0	3
804:ribosomal protein S8:		3	0	1
913:ribosomal protein S11:		10	0	1
130:serine/threonine kinase Ca/Zn:		19	0	0
131:serine/threonine kinase Fe/Mn:		4	0	4
209:thymidine kinase:		3	0	1
612:transforming protein myb:		10	0	4
521:triosephosphate isomerase:		0	0	0
544:trypsin inhibitor (PST1):		14	0	0
515:ubiquitin:		15	0	3
II SUPERFAMILIES WITH MISSING COMPONENTS				
- 443:ATPase F:		5	1	0
- 1:cytochrome c:		123	39	0
- 15:cytochrome P450:		22	1	13
- 705:globin:		429	88	2
- 983:lactoglobulin:		21	8	1
- 1017:prostatein:		3	1	0
- 232:phospholipase A2:		44	2	7
- 616:transforming protein ras:		14	1	5
- 354:trypsin:		64	4	3
III MULTIPLE SUPERFAMILIES ARE CLASSIFIED INTO A SAME GROUP				
* 317:actin:		17	0	5
* 194:asp aminotransferase:		7	0	5
* 921:calmodulin:		72	0	2
* 571:carbamoylphosphate synthetase:		7	0	3
* 10:cytochrome b:		10	0	many
* 38:cytochrome c oxidase 1:		11	0	many
* 73:dihydrofolate reductase:		16	0	3
* 747:DNA dependent RNA polymerase:		0	0	many
* 14:ferrodoxin:		70	47	1
* 1019:koeoo proteins:		4	0	3
* 84:NAD dehydrogenase 2:		5	0	many
* 87:NAD dehydrogenase 4:		10	0	many
* 89:NAD dehydrogenase 5:		10	0	many
* 313:protein kinase:		31	4	many
* 843:ribosomal protein L12:		13	0	2
* 1020:RNA dependent RNA polymerase:		43	0	23
* 142:thymidylate synthase:		13	3	7
* 614:transforming protein myc:		23	0	10
* 581:trypsin inhibitor (PPI):		21	0	1
* 706:ubiquitin:		13	0	0
IV SUPERFAMILIES DIFFICULTY TO BE CLASSIFIED				
# 752:lg:		207	—	1
# 757:lg:		59	—	many
# 908:keratin:		26	1	many
# 317:agrosin B:		5	0	many

Goad-Kanehisa法はローカル類似配列、弱い類似配列も強力に検索する。このことは一方で、ローカルの程度を個々に人間が判断する必要があった。またアミノ酸組成や並びに片寄りがある場合や繰り返し配列を持つ場合にはこの類似性のスコアは適当な同族性の判断にはならない。例えばLeuが多い場合それがLeu rich regionにおいてどのようなパターンで並んでいようと類似性は高くなる。GlyGlyProの繰り返しがあると、GlyGlyLysの繰り返し配列と整列させると66パーセントは同じになる。このような場合も自動化を目指す場合には判断の論理をプログラムする必要があり、今回は人間が判断した。

タンパク質の同一性を決める上でモチーフが鍵になると我々は考えている。あるタンパク質がモチーフを持つかどうかの判断は今回は人間が行った。この判断をコンピュータにさせる事が次の段階であるが、いくつかの解決すべき問題がある。モチーフは通常複数個ある。同族タンパク質がその全てを持つとは限らない。次の図に示すProtein kinaseのスーパーファミリーでは、VALYDY或いはFLVRESのモチーフしかもたない配列が検索され、それらはdomain SH3, domain SH2として知られる。IHRDL, DFGLAR, KWTAPEはkinase domainに特徴的なモチーフであるがこれらがその配列のままいづれのkinaseにもあるわけではない。これらのモチーフの有無の判断は実際には整列表記したときに複数モチーフ全体として10残基程度が一致する場合にモチーフありとした。

自然界ではあるモチーフに起こったミスをはかのモチーフが残基を変化させて補うといったrevertantもある。モチーフ自身の中である残基は実際には他の残基に変わり得るものもある。実験データが数多くたまってモチーフを確定的な配列に洗練する事は原理的に難しい。

(A) PROTEIN KINASE: motif

VALYDY(77)FLVRES(102)VAVKTLK(90)IHRDL(16)DFGLAR(17)KWTAPE(10)SDVWA  
SH3 SH2 IV Vib VII VIII IX

(B) PROTEIN KINASE: domain structure

22---19---15---15---13---11---23---18---16---16---7---16---7---20---7---26---4  
I II III IV ---V--- VIa Vb VII VIII ---IX--- X ---XI---

(C) PROTEIN KINASE: gap position

10-17/4-17/2-9/ 1-6/ 1/ 3-20/ 5-9/7-11/2-19/3-11/ 1/ 1/ 7/ 44/ 1/ 1

上の図(A)はモチーフがある一つのタンパク質でどれほどのアミノ酸残基数の間隔で並んでいるか、そして(B)のサブドメインのどこに存在しているかを示す。図(B)はHanks, S.K.(4)らが分けているkinase catalytic domainのサブドメイン構造を示し、数字は各サブドメインの長さを残基数で示す。図(C)はサブドメインの間のギャップの長さを示す。整列表記した配列により1から44残基の様々な長さのギャップがある事が分かる。図(A), (B), (C)を比較すると分かるように整列表記におけるギャップの長さに変化があるためモチーフの存在場所の判断が難しい。モチーフの変化と並びをどの様に評価するかが分類の自動化のための課題である。本報告では人間が判断した。

結論

タンパク質の分類を二つの手法の組み合わせで行った。弱い類似性の検索でまず配列をグループ化し、更にFRAPライブラリーのモチーフの有無で分類した。この方法を多くの生物から得られたタンパク質を含むスーパーファミリー89個についてPIRのデータを利用して評価した。この方法で56個は容易に分類できた。9個は更に弱い類似性を検索し、モチーフを利用する事で分類が可能である。20個については、我々の分類規範では複数スーパーファミリーが同族となった。これはタンパク質がドメイン構造を持っている事とも関係があり、Doolittle(5)が指摘している分類規範の問題である。残り4個は本方法では分類が困難で同一性を決める方法を開発する事が必要であった。

参考文献

- 1) Dayhoff, M.O., Barker, W.C.: *Methods Enzymol.*, 91, 524-545 (1983)
- 2) Nishikawa, K., Nakashima, H., Kanehisa, M., Ooi, T.: *Protein Seq Data Anal.*, 1, 107-116 (1987)
- 3) Seto, Y., Ikeuchi, Y., Kanehisa, M.: *PROTEINS*, 8, 341-351 (1990)
- 4) Hanks, S.K., Quinn, A.M.: *Methods Enzymol.*, 200, 38-62 (1991)
- 5) Doolittle, R.F.: *Searching through Sequence Databases: Methods Enzymol.*, 183, 99-110 (1990)