

## マルチドメインを持つ遠縁な相同タンパク質の検出手法

瀬下 真吾\* 松井 藤五郎† 大和田 勇人† 朽津 和幸‡

東京理科大学 理工学研究科 経営工学専攻\* 同 理工学部 経営工学科† 同 応用生物科学科‡

## 1 はじめに

タンパク質のドメインに注目してアミノ酸配列比較を行うことは、遠縁な相同タンパク質の発見に有効である。ドメインとはタンパク質ファミリー内で進化的に保存されており、機能的に重要なアミノ酸配列の領域である。複数のドメインによって機能を果たす時、それらをマルチドメインと呼ぶ。遠縁なタンパク質の配列間では、ドメインにおいても配列類似性は弱いと予測される。

マルチドメインを持つタンパク質を対象とした場合、既存の相同性検索ツール HMMER [1] では、各ドメインを考慮することができない。そのため、ドメイン以外の領域での類似配列を持つタンパク質が検出されてしまうという問題がある。

本論文では、ドメインごとに HMMER による検索を行い、その結果を統合することによりマルチドメインとの類似配列を持つタンパク質を検出する手法を提案する。これにより、遠縁な相同であるために配列類似性が弱いタンパク質であっても検出することができるようになる。

## 2 HMMER による相同タンパク質の検出

## 2.1 HMMER とは

HMMER は HMM (隠れマルコフモデル: Hidden Markov Model) を用いてドメイン配列群のモデル化・データベースへの検索等を行うためのツール群である。HMMER による検索では、まず、ファミリーに属するタンパク質からドメイン配列群 (ドメインにあたる配列を集めたもの) を取り出す。次に hmmbuild プログラムによりドメイン配列群の各位置でどのようなアミノ酸が出現するかをモデル化する。このモデルをプロファイル HMM と呼ぶ。ドメイン配列群には共通して出現するアミノ酸パターンが存在するため、構築されたモデルはそのドメイン配列群の特徴を表すことができる。このモデルをクエリーとしてタンパク質データベースへの検索を hmmsearch プログラムにより行い、モデルとの一致度が高い配列を持つものを検出する。ただし、hmmsearch はローカルアライメント (モデルと最も良く一致した部分領域) を求めるため、検出されたタンパク質中にモデル全長に渡っての一致

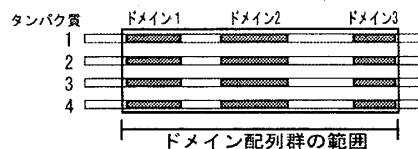


図1 マルチドメインを持つタンパク質

領域があるとは限らない。

検出されたタンパク質には Score と E-value の評価値が与えられる。

**Score** モデルとの適合を示すスコア。よく適合しているほど値が大きくなる。

**E-value** 検索に用いたデータベース中で、上記の Score 以上の類似領域を持つタンパク質数の期待値。E-value が低い値であるほど、そのタンパク質が検出されたことは偶然でない確かな一致であるとみなすことができる。

## 2.2 HMMER の問題点

図1に示すようなマルチドメインを持つタンパク質からプロファイル HMM を構築する際、HMMER による一般的な手法では、全てのドメインを含む範囲を1つのドメイン配列群として hmmbuild への入力とする。

こうして構築したプロファイル HMM を用いてデータベースへ検索を行った場合、マルチドメインを持たないタンパク質が検出されてしまうという問題がある。この原因は、hmmbuild ではマルチドメインを持つモデルを構築できないことにある。また、hmmsearch の検索アルゴリズムはモデルとの部分的な類似配列を持つものを検出してしまうため、モデル中にマルチドメインが存在しても、マルチドメインを含んだタンパク質が検出されるとは限らないのである。

## 3 マルチドメインに注目した検出手法

## 3.1 提案手法の概要

上で述べた問題点を解消するために、本研究では、マルチドメインを含むひとつの範囲からモデルの構築と検索を行う代わりに、ドメインごとにモデルの構築と検索をし、その結果を統合する手法を提案する。概要を図2に示す。

まず、Web 上で公開されているドメインデータベースからドメインごとのアミノ酸配列群が記述されたファイルを取得する。図中の  $D_i$  がこのファイルを表している。次に、ファイルごとに hmmbuild を用いてプロファイル HMM を構築する。続いてプロファイル HMM ごとに hmmsearch を用いてタンパク質データベースへの検索を行う。そして、ドメインごとに

The detection of the remote homologous protein which has a Multi-Domain

Shingo SEJIMO\*, Tohogoroh MATSUI†, Hayato OHWADA†, and Kazuyuki KUCHITSU‡

Department of Industrial Administration, Graduate school of Science and Technology, Tokyo University of Science\*, Department of Industrial Administration, Faculty of Science and Technology, Tokyo University of Science†, Department of Applied Biological Science, Faculty of Science and Technology, Tokyo University of Science‡

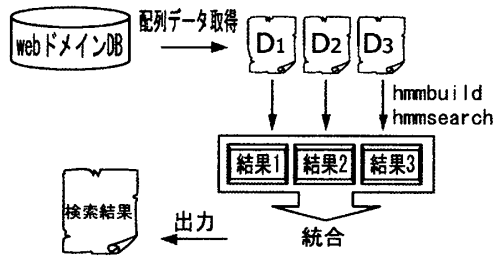


図2 提案手法の概要

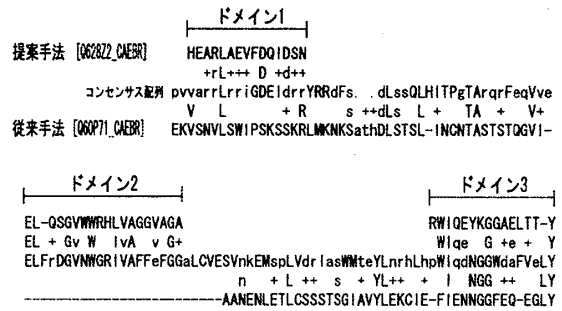


図3 最上位に検出されたタンパク質の比較

得られた検索結果を統合する。

検索結果の統合にはマルチドメインを持つタンパク質を絞り込む過程と、絞り込まれたタンパク質の確率的類似度を計算し再評価する過程がある。

### 3.2 マルチドメインを持つタンパク質の抽出

絞り込みの第一段階として、ドメインごとの検索結果からタンパク質に付けられた固有の識別番号を取り出し、共通要素を求めることで全てのドメインで検出されたタンパク質の集合を得る。

第二段階では、集合に含まれる個々のタンパク質に対し、モデルと類似している領域の先頭位値と最後尾位置の情報を用いてドメインが正しい順序で並んでいるかを判定する。これにより、マルチドメインが正しい並び順で全て存在するタンパク質を得ることができる。

### 3.3 結合 E-value を用いた確率的類似度の再評価

複数のドメインが同時に一致する確率を個々の確率的類似度から求めるために本論文では**結合 E-value** という新しい評価値を提案する。

結合 E-value を求めるためには Bailey [2] らの提案した結合 P-value を用いる。タンパク質  $s$  のドメイン  $d$  における Score が  $x$  であり、ローカルアライメントの長さが  $l$  であるときに、長さ  $l$  の部分配列の Score が  $x$  以上になる確率を  $P_d(s)$  とする。ドメイン数が  $n$  のマルチドメインタンパク質においてドメインごとの  $P_d(s)$  が独立の時、その同時確率  $Z_n(s)$  を次のように定義する。

$$Z_n(s) = \prod_{i=1}^n P_i(s) \quad (1)$$

ここで  $P_i(s)$  はドメイン  $d_i$  における  $P_{d_i}(s)$  を短く表したものである。

この時、ドメイン数  $n$  の同時確率  $Z_n$  がとりうる値の中で、 $Z_n(s)$  以下の値をとる確率を結合 P-value と呼び、次式で計算する。ただし、 $Z_n(s)$  は  $p$  として表示する。

$$F_n(p) = p \sum_{i=0}^{n-1} \frac{(-\ln p)^i}{i!} \quad (2)$$

本研究では  $P_i(s)$  を求める方法として (2) 式の代わりに E-value から計算する手法を用いる。E-value は P-value と  $DBsize$  (検索に用いたデータベースの大きさ) の積によって計算されているので、式変換を行い次式を得る。

$$P_i(s) = \frac{Evalue_i(s)}{DBsize} \quad (3)$$

これにより  $P_i(s)$  を容易に計算することができる。

また、式 (4) によって得られた  $F_n(p)$  に  $DBsize$  を掛けることで結合 E-value :  $E(s)$  を求める。

$$E(s) = F_n(p) \times DBsize \quad (4)$$

HMMER における E-value と同じく  $E(s)$  は数値が小さいほど統計的に良い。そこで、検出されたタンパク質の  $E(s)$  を昇順に並び替えて検索結果を出力する。

## 4 検出されたタンパク質の比較

本研究の提案手法に基づいて相同性検索を行った。図3は線虫のタンパク質データベースに対して3つのドメインを持つ Bcl-2 ファミリーを検索した結果、最上位に検出されたタンパク質の比較である。最上段と最下段の配列は、それぞれ提案手法と従来手法によって検出されたタンパク質である。コンセンサス配列と一致している箇所にはアミノ酸配列を表す文字が、類似している箇所には + のマークが示されている。

従来手法ではドメイン2との類似配列が存在していないタンパク質が検出されてしまったが、提案手法ではマルチドメインの各ドメインとの類似配列を持つタンパク質が検出できた。

## 5 まとめ

本論文ではマルチドメインを持つ遠縁な相同タンパク質の検出を行うために、ドメインごとの検索結果を統合し、結合 E-value により確率的類似度を評価する手法を提案した。

提案手法によって遠縁であってもマルチドメインを持つタンパク質の検出が可能となった。しかし、本手法ではドメイン部分にしか着目していないため、検出されたタンパク質の長さが同じファミリーに属する他のタンパク質と大きく異なることがある。ドメイン間の長さを考慮したタンパク質の評価については今後の課題である。

## 参考文献

- [1] S.R. Eddy. Multiple alignment using hidden markov models. *Ismb*, Vol. 3, pp. 114–120, 1995.
- [2] T.L. Bailey, and M. Gribskov. Methods and statistics for combining motif match scores. *J. Comput. Biol.*, Vol. 5, pp. 211–221, 1998.